

# Chapter 16

## Basic Concepts

### 16.1 Languages, grammars and automata

At one level of description, a natural language is simply a set of strings—finite sequences of words, morpheme, phonemes, or whatever. Not every possible sequence is in the language: we distinguish the *grammatical* strings from those that are *ungrammatical*. A *grammar*, then, is some explicit device for making this distinction; it is, in other words, a means for selecting a subset of strings, those that are grammatical, from the set of all possible strings formed from an initially given alphabet or vocabulary.

In this chapter we will consider two classes of formal devices which can function as grammars in this very general sense: (1) automata, which are abstract computing machines, and (2) string rewriting systems, which generally bear the name “grammar” or “formal grammar”. The latter will be familiar to linguists inasmuch as grammars in this sense have formed the basis of much of the work in generative transformational theory.

We begin by considering certain properties of strings and sets of strings. Given a finite set  $A$ , a *string on* (or *over*)  $A$  is a finite sequence of occurrences of elements from  $A$ . For example, if  $A = \{a, b, c\}$ , then  $acbaab$  is a string on  $A$ . Strings are by definition finite in length. (Infinite sequences of symbols are also perfectly reasonable objects of study, but they are not suitable as models for natural language strings.) The set from which strings are formed is often called the *vocabulary* or *alphabet*, and this too is always assumed to be finite. The length of a string is, of course, the number of occurrences of symbols in it (i.e., the number of tokens, not the number of types). The string  $acbaab$  thus is of length 6.

Because we are dealing with tokens of an alphabet, there is an important difference between the linearly ordered sequences we call strings and a linearly ordered set. If the set  $A = \{a, b, c\}$  were linearly ordered, say,  $b \rightarrow a \rightarrow c$ , each element of  $A$  would occupy a unique place in the ordering. In a string, e.g.,  $acbaab$ , tokens of  $a$ , occur in the first, fourth, and fifth positions.

To be formal, one could define a string of length  $n$  over the alphabet  $A$  to be a function mapping the first  $n$  positive integers into  $A$ . For example,  $acbaab$  would be the function  $\{\langle 1, a \rangle, \langle 2, c \rangle, \langle 3, b \rangle, \langle 4, a \rangle, \langle 5, a \rangle, \langle 6, b \rangle\}$ . There is little to be gained in this case by the reduction to the primitives of set theory, however, so we will continue to think of strings simply as finite sequences of symbols. A string may be of length 1, and so we distinguish the string  $b$  of length 1 from the symbol  $b$  itself. We also recognize the (unique) string of length 0, the *empty string*, which we will denote  $e$  (some authors use  $\Lambda$ ). Two strings are identical if they have the same symbol occurrences in the same order; thus,  $acb$  is distinct from  $abc$ , and strings of different length are always distinct.

An important binary operation on strings is concatenation, which amounts simply to juxtaposition. For example, the strings  $abca$  and  $bac$  can be concatenated, in the order mentioned, to give the string  $abcabac$ . Sometimes concatenation is denoted with the symbol " $\wedge$ " thus,  $abca \wedge bac$ . Concatenation is associative since for any strings  $\alpha, \beta, \gamma$ ,  $(\alpha \wedge \beta) \wedge \gamma = \alpha \wedge (\beta \wedge \gamma)$ , but it is not commutative, since in general  $\alpha \wedge \beta \neq \beta \wedge \alpha$ . The empty string is the identity element for concatenation; i.e., for any string  $\alpha$ ,  $\alpha \wedge e = e \wedge \alpha = \alpha$ .

Given a finite set  $A$ , the set of all strings over  $A$ , denoted  $A^*$ , together with the operation of concatenation constitutes a monoid. Concatenation is well-defined for any pair of strings in  $A^*$  and the result is a string in  $A^*$ ; the operation is associative; and there is an identity element  $\langle A^*, \wedge \rangle$  fails to be a group since no element other than  $e$  has an inverse: no string concatenated with a non-empty string  $x$  will yield the empty string. Since concatenation is not commutative,  $\langle A^*, \wedge \rangle$  is not an Abelian monoid.

A frequently encountered unary operation on strings is reversal. The reversal of a string  $x$ , denoted  $x^R$ , is simply the string formed by writing the symbols of  $x$  in the reverse order. Thus  $(acbab)^R = babca$ . The reversal of  $e$  is just  $e$  itself. To be formal, we could define reversal by induction on the length of a string:

DEFINITION 16.1 Given an alphabet  $A$ :

- (1) If  $x$  is a string of length 0, then  $x^R = x$  (i.e.,  $e^R = e$ )
- (2) If  $x$  is a string of length  $k + 1$ , then it is of the form  $wa$ , where  $a \in A$  and  $w \in A^*$ ; then  $x^R = (wa)^R = aw^R$ .

■

Concatenation and reversal are connected in the following way: For all strings  $x$  and  $y$ ,  $(x \frown y)^R = y^R \frown x^R$ . For example,

$$(16-1) \quad (bca \frown ca)^R = (ca)^R \frown (bca)^R = ac \frown acb = acacb$$

Given a string  $x$ , a *substring* of  $x$  is any string formed from contiguous occurrences of symbols in  $x$  taken in the same order in which they occur in  $x$ . For example,  $bac$  is a substring of  $abacca$ , but neither  $bcc$  nor  $cb$  is a substring. Formally,  $y$  is a substring of  $x$  iff there exist strings  $z$  and  $w$  such that  $x = z \frown y \frown w$ . In general,  $z$  or  $w$  (or both) may be empty, so every string is trivially a substring of itself. (Non-identical substrings can be called *proper* substrings.) The empty string is a substring of every string; i.e., given  $x$  we can choose  $z$  in the definition as  $e$  and  $w$  as  $x$  so that  $x = e \frown e \frown x$ .

An initial substring is called a *prefix*, and a final substring, a *suffix*. Thus,  $ab$  is a (proper) prefix of  $abacca$ , and  $cca$  is a (proper) suffix of this string.

We may now define a *language* (over a vocabulary  $A$ ) as any subset of  $A^*$ . Since  $A^*$  is a denumerably infinite set, it has cardinality  $\aleph_0$ ; its power set, i.e., the set of all languages over  $A$ , has cardinality  $2^{\aleph_0}$  and is thus non-denumerably infinite. Since the devices for characterizing languages which we will consider, *viz.*, formal grammars and automata, form denumerably infinite classes, it follows that there are infinitely many languages—in fact, non-denumerably infinitely many—which have no grammar. What this means in intuitive terms is that there are languages which are such motley collections of strings that they cannot be completely characterized by any finite device. The languages which *are* so characterizable exhibit a certain amount of order or pattern in their strings which allows these strings to be distinguished from others in  $A^*$  by a grammar or automaton with finite resources. The study of formal languages is essentially the investigation of a scale of

complexity in this patterning in strings. For example, we might define a language over the alphabet  $\{a, b\}$  in the following way:

$$(16-2) \quad L = \{x \mid x \text{ contains equal numbers of } a\text{'s and } b\text{'s (in any order)}\}$$

We might then compare this language with the following:

$$(16-3) \quad L_1 = \{x \in \{a, b\}^* \mid x = a^n b^n (n \geq 0)\}, \text{ i.e., strings consisting of some number of } a\text{'s followed by the same number of } b\text{'s}$$

$$L_2 = \{x \in \{a, b\}^* \mid x \text{ contains a number of } a\text{'s which is the square of the number of } b\text{'s}\}$$

Is  $L_1$  or  $L_2$  in some intuitive sense more complex than  $L$ ? Most would probably agree that  $L_2$  is a more complex language than  $L$  in that greater effort would be required to determine that the members of  $a$ 's and  $b$ 's stood the "square" relation than to determine merely that they were equal. In other words, a device which could discriminate strings from non-strings of  $L_2$  would have to be more powerful or more "intelligent" than a device for making the comparable discrimination for  $L$ .

What of  $L_1$  and  $L$ ? Here our intuitions are much less clear. Some might think that it would require a less powerful device to recognize strings in  $L$  reliably than to recognize strings in  $L_1$ ; others might think it is the other way around or see no difference. As it happens, the particular scale of complexity we will investigate (the so-called Chomsky Hierarchy) does regard  $L_2$  as more complex than  $L$  but puts  $L_1$  and  $L$  in the same complexity class. At least this is so for the overall complexity measure. Finer divisions could be established which might distinguish  $L_1$  from  $L$ .

One linguistic application of these investigations is to try to locate natural languages on this complexity scale. This is part of the overall task of linguistics to characterize as precisely as possible the class of (potential and actual) natural languages and to distinguish this class from the class of all language-like systems which could not be natural languages. One must keep clearly in mind the limitations of this enterprise, however, the principal one being that languages are regarded here simply as string sets. It is clear that sentences of any natural language have a great deal more structure than simply the concatenation of one element with another. Thus, to establish a complexity scale for string sets and to place natural languages on this scale may, because of the neglect of other important structural properties, be to classify natural language along an ultimately irrelevant dimension. Extend-

ing results from the study of formal languages into linguistic theory must therefore be done with great caution.

## 16.2 Grammars

A formal grammar (or simply, grammar) is essentially a deductive system of axioms and rules of inference (see Chapter 8), which generates the sentences of a language as its theorems. By the usual definitions, a grammar contains just one axiom, the string consisting of the *initial symbol* (usually  $S$ ), and a finite number of rules of the form  $\psi \rightarrow \omega$ , where  $\psi$  and  $\omega$  are strings, and the interpretation of a rule is the following: whenever  $\psi$  occurs as a substring of any given string, that occurrence may be replaced by  $\omega$  to yield a new string. Thus if a grammar contained the rule  $AB \rightarrow CDA$ , we could derive from the string  $EBABCC$  the string  $EBCDACC$ .

Grammars use two alphabets: a *terminal alphabet* and a *non-terminal alphabet*, which are assumed to be disjoint. The strings we are interested in deriving, i.e., the sentences of the language, are strings over the terminal alphabet, but intermediate strings in derivations (proofs) by the grammar may contain symbols from both alphabets. We also require in the rules of the grammar that the string on the left side not consist entirely of terminal symbols. Here is an example of a grammar meeting these requirements:

$$\begin{aligned}
 (16-4) \quad V_T \text{ (the terminal alphabet)} &= \{a, b\} \\
 V_N \text{ (the non-terminal alphabet)} &= \{S, A, B\} \\
 S &\text{ (the initial symbol—a member of } V_N) \\
 R \text{ (the set of rules)} &= \left\{ \begin{array}{l} S \rightarrow ABS \\ S \rightarrow e \\ AB \rightarrow BA \\ BA \rightarrow AB \\ A \rightarrow a \\ B \rightarrow b \end{array} \right\}
 \end{aligned}$$

A common notational convention is to use lower case letters for the terminal alphabet and upper case letters for the non-terminal alphabet.

A derivation of the string  $abba$  by this grammar could proceed as follows:

$$\begin{aligned}
 (16-5) \quad S &\Rightarrow ABS \Rightarrow ABABS \Rightarrow ABAB \Rightarrow ABBA \Rightarrow ABbA \Rightarrow \\
 &aBbA \Rightarrow abbA \Rightarrow abba
 \end{aligned}$$

Here we have used the symbol " $\Rightarrow$ " to mean "yields in one rule application." Note that  $abba$  is not subject to further rewriting inasmuch as it consists entirely of terminal symbols and no rule licenses rewriting strings of terminals. The sequence (16-5) is said to be a *derivation (of  $abba$  from  $S$ )*, and the string  $abba$  is said to be *generated by* the grammar. The *language generated by* the grammar is the set of all strings generated. Here are the formal definitions:

**DEFINITION 16.2** Let  $\Sigma = V_T \cup V_N$ . A (formal) grammar  $G$  is a quadruple  $\langle V_T, V_N, S, R \rangle$ , where  $V_T$  and  $V_N$  are finite disjoint sets,  $S$  is a distinguished member of  $V_N$ , and  $R$  is a finite set of ordered pairs in  $\Sigma^* V_N \Sigma^* \times \Sigma^*$ . ■

We have written  $\psi \rightarrow \omega$  above for clarity instead of  $\langle \psi, \omega \rangle$ . The last condition simply says that a rule rewrites a string containing at least one non-terminal as some (possibly empty) string.

**DEFINITION 16.3** Given a grammar  $G = \langle V_T, V_N, S, R \rangle$ , a derivation is a sequence of strings  $x_1, x_2, \dots, x_n$  ( $n \geq 1$ ) such that  $x_1 = S$  and for each  $x_i$  ( $2 \leq i \leq n$ ),  $x_i$  is obtained from  $x_{i-1}$  by one application of some rule in  $R$ . ■

To be completely formal, we would spell out in detail what it means to apply a rule of  $R$  to a string. The reader may want to do this as an exercise.

**DEFINITION 16.4** A grammar  $G$  generates a string  $x \in V_T^*$  if there is a derivation  $x_1, \dots, x_n$  by  $G$  such that  $x_n = x$ . ■

Note that by this definition only strings of terminal symbols are said to be generated.

**DEFINITION 16.5** The language generated by a grammar  $G$ , denoted  $L(G)$ , is the set of all strings generated by  $G$ . ■

The language generated by the grammar in the example of (16-4) is  $\{x \in \{a, b\}^* \mid x \text{ contains equal numbers of } a\text{'s and } b\text{'s}\}$ .

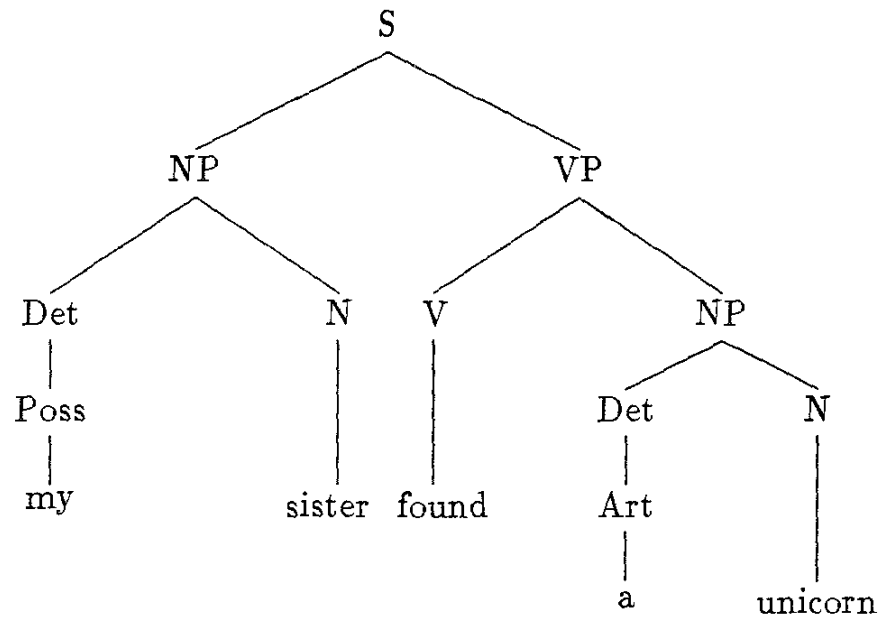


Figure 16-1: A typical constituent structure tree

### 16.3 Trees

When the rules of a grammar are restricted to rewriting only a single non-terminal symbol, it is possible to contrive grammars as generating *constituent structure trees* rather than simply strings. An example of such a tree is shown in Fig. 16-1.

Such diagrams represent three sorts of information about the syntactic structure of a sentence:

1. The hierarchical grouping of the parts of the sentence into constituents
2. The grammatical type of each constituent
3. The left-to-right order of the constituents

For example, Fig. 16-1 indicates that the largest constituent, which is labeled by S (for Sentence), is made up of a constituent which is a N(oun) P(hrase) and one which is a V(erb) P(hrase) and that the noun phrase is composed of two constituents: a Det(erminer) and a N(oun), etc. Further,

in the sentence constituent the noun phrase precedes the verb phrase, the determiner precedes the noun in the noun phrase constituents, and so on. The tree diagram itself is said to be composed of *nodes*, or points, some of which are connected by lines called *branches*. Each node has associated with it a *label* chosen from a specified finite set of grammatical categories (S, NP, VP, etc.) and formatives (*my*, *sister*, etc.). As they are customarily drawn, a tree diagram has a vertical orientation on the page with the nodes labeled by the formatives at the bottom. Because a branch always connects a higher node to a lower one, it is an inherently directional connection. This directionality is ordinarily not indicated by an arrow, as in the usual diagrams of relations, but only by the vertical orientation of the tree taken together with the convention that a branch extends *from* a higher node *to* a lower node.

### 16.3.1 Dominance

We say that a node  $x$  *dominates* a node  $y$  if there is a connected sequence of branches in the tree extending from  $x$  to  $y$ . This is the case when all the branches in the sequence have the same orientation away from  $x$  and toward  $y$ . For example, in Fig. 16-1 the node labeled VP dominates the node labeled Art, since the sequence of branches connecting them is uniformly descending from the higher node VP to the lower node Art. The node labeled VP does not dominate the node labeled Poss, since the path by which they are joined first ascends from VP to S and then descends through NP and Det.

Given a tree diagram, we represent the fact that  $x$  dominates  $y$  by the ordered pair  $\langle x, y \rangle$ . The set of all such ordered pairs for a given tree is said to constitute the *dominance relation* for that tree. Dominance is clearly a transitive relation. If  $x$  is connected to  $y$  by a sequence of descending branches and  $y$  is similarly connected to  $z$ , then  $x$  dominates  $z$  because they are also connected by a sequence of descending branches, specifically, by the sequence passing through  $y$ . As a technical convenience, it is usually assumed that every node dominates itself, i.e., that the dominance relation is reflexive. Further, if  $x$  dominates  $y$ , then  $y$  can dominate  $x$  only if  $x = y$ ; or in other words, dominance is antisymmetric. Thus, the relation of dominance is a weak partial ordering of the nodes of a tree.

If  $x$  and  $y$  are distinct,  $x$  dominates  $y$ , and there is no distinct node between  $x$  and  $y$ , then  $x$  *immediately dominates*  $y$ . In Fig. 16-1, the node labeled VP immediately dominates the node labeled V but not the node labeled *found*. A node is said to be the *daughter* of the node immediately



dominating it, and distinct nodes immediately dominated by the same node are called *sisters*. In Fig. 16-1, the node labeled VP has two daughters, viz., the node labeled V and the rightmost node labeled NP. The latter two nodes are sisters. A node which is minimal in the dominance relation, i.e., which is not dominated by any other node, is called a *root*. In Fig. 16-1 there is one root, the node labeled S. Maximal elements are called *leaves*, and in Fig. 16-1 these are the nodes labeled by the formatives, *my*, *sister*, etc. Note that a tree diagram is ordinarily drawn upside down since the root is at the top and the leaves are at the bottom.

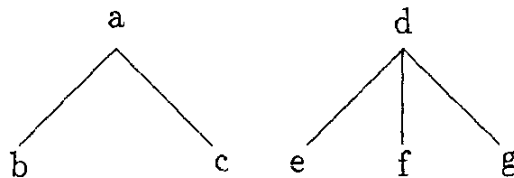


Figure 16-2: A multiply rooted "tree"

Mathematicians sometimes use the term *tree* for a configuration with more than one root, e.g., that shown in Fig. 16-2. For linguists, however, a tree is invariably singly rooted, the configuration in Fig. 16-2 being considered a "forest" of trees. We shall adhere to linguistic usage and accordingly we have the following condition:

**The Single Root Condition:** In every well-formed constituent structure tree there is exactly one node that dominates every node.

The root node is, therefore, a least element (and necessarily also a minimal element) in the dominance relation. We note, incidentally, that the Single Root Condition is met in the trivial case of a tree that has only one node, which is simultaneously root and leaf. The condition would not be met by an "empty" tree with no nodes at all, since it asserts that a node with the specified property exists in the tree.

### 16.3.2 Precedence

Two nodes are ordered in the left-to-right direction just in case they are not ordered by dominance. In Fig. 16-1 the node labeled V precedes (i.e., is to

the left of) its sister node labeled NP and all the nodes dominated by this NP node; it neither precedes nor follows the nodes labeled S, VP, V, and *found*, i.e., the nodes that either dominate or are dominated by the V node. It is not logically necessary that the relations of dominance and left-to-right precedence be mutually exclusive, but this accords with the way in which tree diagrams are usually interpreted.

Given a tree, the set of all ordered pairs  $\langle x, y \rangle$  such that  $x$  precedes  $y$  is said to define the *precedence relation* for that tree. To ensure that the precedence and dominance relations have no ordered pairs in common, we add the Exclusivity Condition:

**The Exclusivity Condition:** In any well-formed constituent structure tree, for any nodes  $x$  and  $y$ ,  $x$  and  $y$  stand in the precedence relation  $P$ , i.e., either  $\langle x, y \rangle \in P$  or  $\langle y, x \rangle \in P$ , if and only if  $x$  and  $y$  do not stand in the dominance relation  $D$ , i.e., neither  $\langle x, y \rangle \in D$  nor  $\langle y, x \rangle \in D$ .

Like dominance, precedence is a transitive relation, but precedence is irreflexive rather than reflexive. The latter follows from the Exclusivity Condition, since for every node  $x$ ,  $\langle x, x \rangle \in D$  and therefore  $\langle x, x \rangle \notin P$ . If  $x$  precedes  $y$ , then  $y$  cannot precede  $x$ , and thus the relation is asymmetric. Precedence, therefore, defines a strict partial order on the nodes of the tree.

One other condition on the dominance and precedence relations is needed to exclude certain configurations from the class of well-formed trees. An essential characteristic of a tree that distinguishes it from a partially ordered set in general is that no node can have more than one branch entering it; i.e., every node has at most one node immediately dominating it. The structure shown in Fig. 16-3(a) has a node  $d$  with two immediate predecessors,  $b$  and  $c$ , and therefore it is not a tree. Another defining property of trees is that branches are not allowed to cross. Figure 16-3(b) illustrates the sort of structure that is forbidden. Both types of ill-formedness can be ruled out by adding the Nontangling Condition:

**The Nontangling Condition:** In any well-formed constituent structure tree, for any nodes  $x$  and  $y$ , if  $x$  precedes  $y$ , then all nodes dominated by  $x$  precede all nodes dominated by  $y$ .

The configuration in Fig. 16-3(a) fails to meet this condition because  $b$  precedes  $c$ ,  $b$  dominates  $d$ , and  $c$  dominates  $d$ , and therefore  $d$  ought to precede  $d$ . This is impossible, however, since precedence is irreflexive. In Fig. 16-3(b),  $b$  precedes  $c$ ,  $b$  dominates  $d$ , and  $c$  dominates  $e$ . Thus, by the Nontangling Condition,  $d$  should precede  $e$ , but in fact the reverse is true.

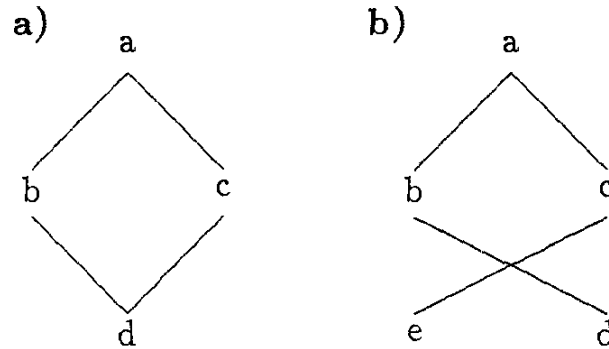


Figure 16-3: Structures excluded as trees by the Nontangling Condition

### 16.3.3 Labeling

To complete the characterization of trees we must consider the labeling of the nodes. It is apparent from Fig. 16-1 that distinct nodes can have identical labels attached to them, e.g., the two nodes labeled NP. Since each node has exactly one label, the pairing of nodes and labels can be represented by a *labeling function*  $L$ , whose domain is the set of nodes in the tree and whose range is a set (in syntactic trees, a set of grammatical categories and formatives). The mapping is, in general, an *into* function. In summary, we have the following definition:

**DEFINITION 16.6** A (constituent structure) tree is a mathematical configuration  $\langle N, Q, D, P, L \rangle$ , where

$N$  is a finite set, the set of nodes

$Q$  is a finite set, the set of labels

$D$  is a weak partial order in  $N \times N$ , the dominance relation

$P$  is a strict partial order in  $N \times N$ , the precedence relation

$L$  is a function from  $N$  into  $Q$ , the labeling function

and such that the following conditions hold:

(1)  $(\exists x \in N)(\forall y \in N)\langle x, y \rangle \in D$  (Single Root Condition)

(2)  $(\forall x, y \in N)((\langle x, y \rangle \in P \vee \langle y, x \rangle \in P) \leftrightarrow (\langle x, y \rangle \notin D \ \& \ \langle y, x \rangle \notin D))$   
(Exclusivity Condition)

- (3)  $(\forall w, x, y, z \in N)((\langle w, x \rangle \in P \ \& \ \langle w, y \rangle \in D \ \& \ \langle x, z \rangle \in D) \rightarrow \langle y, z \rangle \in P)$   
 (Nontangling Condition) ■

Given this definition, one can prove theorems of the following sort:

**THEOREM 16.1** *Given a tree  $T = \langle N, Q, D, P, L \rangle$ , every pair of sister nodes is ordered by  $P$ .* ■

*Proof:* Take  $x$  and  $y$  as sisters immediately dominated by some node  $z$ . By the definitions of ‘sister’ and ‘immediate domination,’  $x, y$ , and  $z$  must all be distinct. As an assumption to be proved false, let  $x$  dominate  $y$ . Therefore,  $x$  must dominate  $z$ , since  $z$  immediately dominates  $y$ . But  $z$  also dominates  $x$ , and  $x$  and  $z$  are distinct, so this violates the condition that dominance is antisymmetric. Therefore,  $x$  cannot dominate  $y$ . By a symmetrical argument, we can show that  $y$  does not dominate  $x$ . Thus,  $\langle x, y \rangle \notin D$  and  $\langle y, x \rangle \notin D$ , and by the Exclusivity Condition it follows that  $\langle x, y \rangle \in P \vee \langle y, x \rangle \in P$ ; i.e.,  $x$  and  $y$  are ordered by  $P$ . ■

**THEOREM 16.2** *Given a tree  $T = \langle N, Q, D, P, L \rangle$ , the leaves are totally ordered by  $P$ .* ■

*Proof:* Let  $M$  be the set of leaves, and let  $R$  be the restriction of the relation  $P$  to the set  $M$ ; i.e.,  $R = \{\langle x, y \rangle \in M \times M \mid \langle x, y \rangle \in P\}$ .  $R$  is a strict partial order, since if there were any ordered pairs violating the conditions of irreflexivity, asymmetry, and transitivity in  $R$ , then because  $R \subseteq P$ , these pairs would also appear in  $P$ , and  $P$  would not be a strict partial order. By definition, a leaf dominates no node except itself, and therefore for every pair of distinct leaves  $x$  and  $y$ ,  $\langle x, y \rangle \notin D$  and  $\langle y, x \rangle \notin D$ . Thus, by the Exclusivity Condition  $\langle x, y \rangle \in P \vee \langle y, x \rangle \in P$ . Since  $x$  and  $y$  are leaves,  $\langle x, y \rangle \in R \vee \langle y, x \rangle \in R$ , by the definition of  $R$ , and thus  $R$  is connex. Therefore,  $R$  is a strict total order. ■

Every statement about the formal properties of a constituent structure tree can be formulated in terms of the dominance and precedence relations and the labeling function. For example, one useful predicate on trees is that of *belonging to*. A node will be said to belong to the next highest  $S$  node that dominates it. Formally, the definition is as follows:

DEFINITION 16.7 Given a tree  $T = \langle N, Q, D, P, L \rangle$ , node  $x$  belongs to node  $y$  iff

- (1)  $x \neq y$
- (2)  $\langle y, x \rangle \in D$
- (3)  $\langle y, S \rangle \in L$
- (4)  $\sim (\exists w \in N)(\langle w, S \rangle \in L \ \& \ w \neq y \ \& \ w \neq x \ \& \ \langle y, w \rangle \in D \ \& \ \langle w, x \rangle \in D)$ .

■

Parts 2 and 3 of this definition specify that the node to which  $x$  belongs is labeled  $S$  and dominates  $x$ . Part 4 prohibits any  $S$  node from standing between  $x$  and  $y$  in the dominance relation, and part 1 excludes the case of an  $S$  node belonging to itself. To illustrate, let us consider the tree in Fig 16-4.

The node *Prn* belongs to the circled  $S$  node since this is the next highest  $S$  node dominating it. *Prn* does not belong to the highest  $S$  (i.e., the root) of the tree because the circled  $S$  node is between the root and *Prn* in the dominance relation.

With this definition we can easily define some other predicates. Two nodes are called *clause mates* iff neither dominates the other and both belong to the same node. In Fig 16-4 the nodes labeled *John* and *him* are clause mates since neither dominates the other and both belong to the circled  $S$  node. *Fred* and *him* are not clause mates since they do not belong to the same node, and *Prn* and *him* are not clause mates since *Prn* dominates *him*.

If we let  $B\langle x, y \rangle$  denote ' $x$  belongs to  $y$ ,' we can state the definition of clause mates as follows:

DEFINITION 16.8 Given a tree  $T = \langle N, Q, D, P, L \rangle$ , nodes  $x$  and  $y$  are clause mates iff  $\langle x, y \rangle \notin D \ \& \ \langle y, x \rangle \notin D \ \& \ (\exists z \in N)(\langle x, z \rangle \in B \ \& \ \langle y, z \rangle \in B)$ . ■

A node  $x$  is said to *command* a node  $y$  iff neither dominates the other and  $x$  belongs to a node  $z$  that dominates  $y$  (Langacker, 1969). In Fig 16-4 the node labeled *Fred* commands the node labeled *him* since neither dominates the other and *Fred* belongs to the root node  $S$ , which also dominates *him*. The node *him* does not command *Fred*, however, since the node to which *him* belongs—the circled  $S$  node—does not dominate *Fred*. Note, further, that *John* commands *him* and vice versa. Formally, the definition is as follows

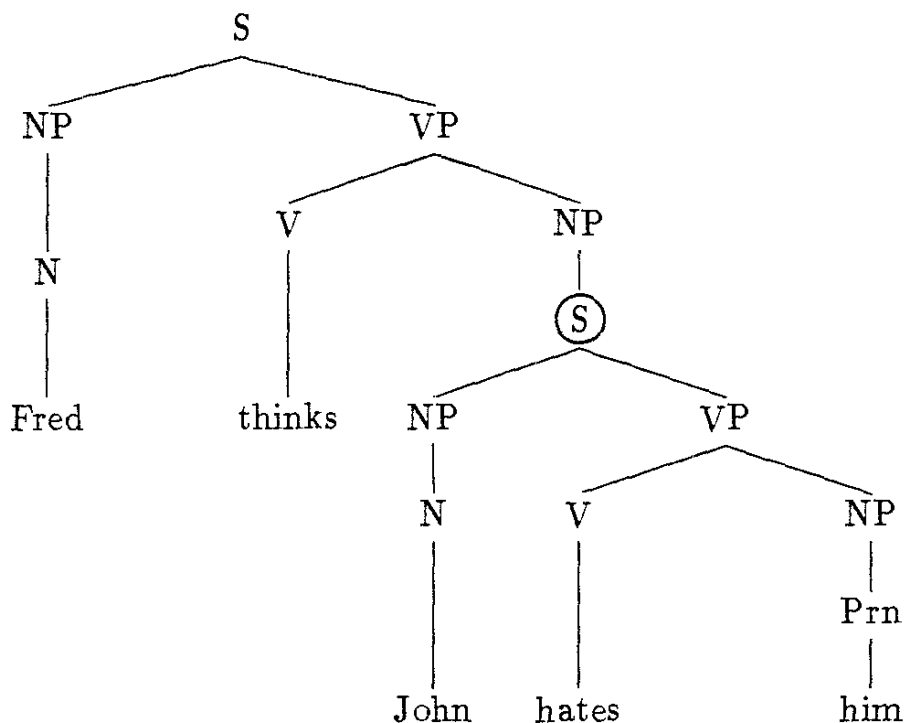


Figure 16-4: Tree illustrating the definitions of 'belonging to' and 'command'

**DEFINITION 16.9** Given a tree  $T = \langle N, Q, D, P, L \rangle$ , node  $x$  commands node  $y$  iff  $\langle x, y \rangle \notin D$  &  $\langle y, x \rangle \notin D$  &  $(\exists z \in N)(\langle x, z \rangle \in B$  &  $\langle z, y \rangle \in D)$ . ■

*Problem:* Prove that two nodes are clause mates iff each commands the other.

## 16.4 Grammars and trees

As we have said, if a grammar has only rules of the form  $A \rightarrow \psi$ , where  $A$  is a nonterminal symbol, there is a natural way to associate applications of such rules with the generation of a tree. For example, if the grammar contains the rule  $A \rightarrow aBc$ , we can associate this with the (sub)tree in Fig. 16-5.

in which  $A$  immediately dominates  $a, B$ , and  $c$ , and the latter three elements stand in the precedence relation in the order given. Further, if the grammar

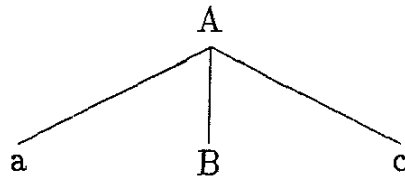


Figure 16-5.

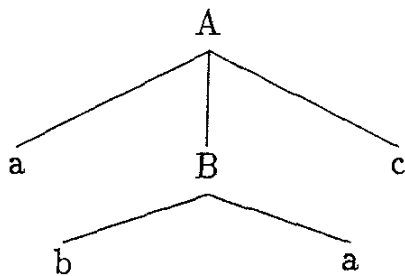


Figure 16-6

also contains the rule  $B \rightarrow ba$ , we can apply this rule at the node labelled  $B$  in the preceding tree to produce the tree shown in Fig. 16-6.

Let us define the *yield* of a tree as the string formed by its leaves ordered according to the precedence relation. The yield of the tree in Fig. 16-6, for example, is  $abac$ ; that of Fig. 16-5 is  $aBc$ . We can now say:

**DEFINITION 16 10** *A grammar (having all rules of the form  $A \rightarrow \psi$ ) generates a tree iff all the following hold:*

- (i) *the root is labelled with the initial symbol of the grammar*
- (ii) *the yield is a string of terminal symbols*

- (iii) *for each subtree of the form  $\begin{array}{c} A \\ \triangle \\ \alpha_1 \dots \alpha_n \end{array}$  in the tree, where  $A$  immediately dominates  $\alpha_1 \dots \alpha_n$ , there is a rule in the grammar  $A \rightarrow \alpha_1 \dots \alpha_n$ .*

■

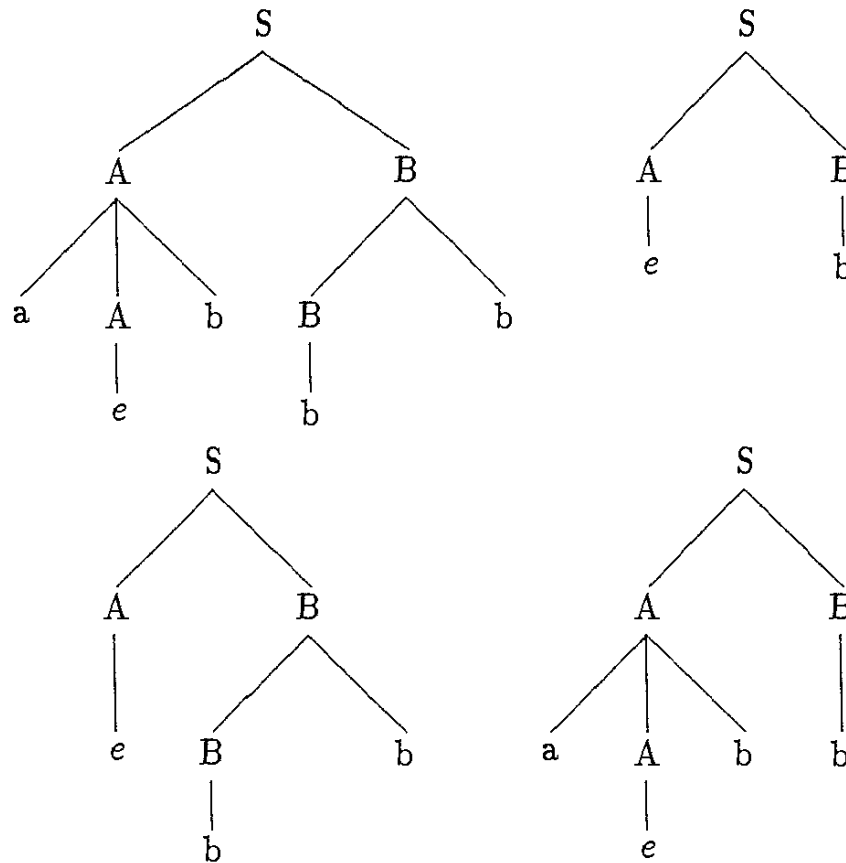


Figure 16-7

Thus the grammar  $G = \langle \{a, b\}, \{S, A, B\}, S, R \rangle$  where

$$R = \left\{ \begin{array}{ll} S \rightarrow AB & B \rightarrow Bb \\ A \rightarrow aAb & B \rightarrow b \\ A \rightarrow e & \end{array} \right\}$$

generates trees such as those in Fig. 16-7. We can further say that a string is generated by such a grammar iff it is the yield of some tree which is generated. The language generated is, as usual, the set of all strings generated. For grammars in which there is only a single symbol on the left side of each rule, this definition and the earlier definition of generation of a string turn out to be equivalent: a string is generated (by the earlier definition) iff it is the yield of some generated tree.

*Problem:* What language is generated by the above grammar?

Such grammars have interested linguists precisely because of the possi-

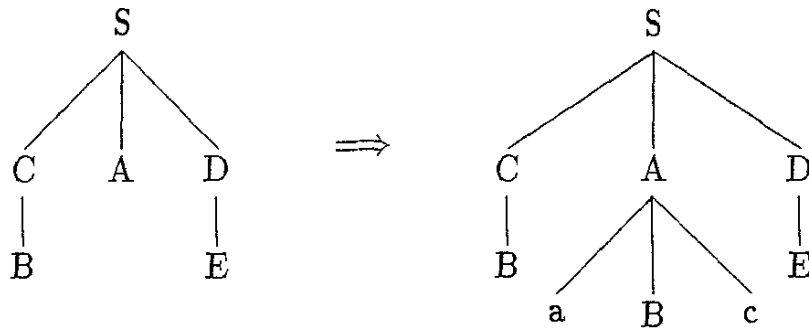


bility of specifying a constituent structure tree for each string generated. In attempting to write such grammars for natural languages, however, linguists have noted that often such rules are not universally applicable but may be allowed only in certain contexts. For example, a rule rewriting  $\text{Det}(\text{erminer})$  as *many* might be applied only if the following noun were a plural form. Such considerations led to the investigation of formal grammar rules of the form  $A \rightarrow \psi/\alpha\_ \beta$ , where the “/” is read “in the context”, and where “\_” marks the position of the  $A$ . The interpretation of such a rule is that the symbol  $A$  can be replaced by the string  $\psi$  in a derivation only when the string  $\alpha$  immediately precedes  $A$  and the string  $\beta$  immediately follows  $A$ . The context specifications are not necessarily exhaustive: additional symbols may occur to the left of the  $\alpha$  and to the right of  $\beta$ . For example, if the rule were  $A \rightarrow aBc/C\_Dc$ , then the string  $BECADcbA$  could be rewritten as  $BECaBcDcbA$ .

Such rules are called *context sensitive* in contrast to rules of the form  $A \rightarrow \psi$ , which are called *context free*. A context free rule, thus, is a context sensitive rule in which the context is null.

A context sensitive rule  $A \rightarrow \psi/\alpha\_ \beta$  can also be written as  $\alpha A \beta \rightarrow \alpha \psi \beta$  in conformity with the schema for grammar rules generally. So long as we regard these grammars as string rewriting systems the notations are interchangeable: in either case we may replace  $A$  by  $\psi$  when we find the substring  $\alpha A \beta$ . However, if we want to think of context sensitive rules as generating trees, the two representations may not be equivalent. For example, the rule  $CABD \rightarrow CAaBD$  could be construed either as  $A \rightarrow Aa/C\_BD$  or as  $B \rightarrow aB/CA\_D$ , and the associated trees would obviously differ depending on whether an  $A$  node or a  $B$  node was expanded.

Another problem which arises is how the context restriction is to be satisfied by the tree. If we think of the rules as specifying how one tree is to be converted into the next in a derivation, then does a rule such as  $A \rightarrow aBc/C\_D$  mean that the  $C$  and  $D$  must be *leaves* immediately to the left and right, respectively, of  $A$  when the rule is applied, or is it sufficient that the  $C$  *immediately precede* the  $A$  and the  $D$  *immediately follow*, without necessarily being leaves along with  $A$ ? Under the latter interpretation, the following derivational step would be allowed, but by the former it would not.



Note also that in the definition of tree derivation by means of context free rules in Def 16-10 above, we essentially thought of the trees being somehow given in advance and then checked for well-formedness by the grammar rules. That is, the rules served as so-called “node admissibility conditions” rather than as directions for converting one tree into another. In the context free case, the two points of view are equivalent, but this is not the case for context sensitive rules. For example, the grammar

$$\begin{aligned}
 (16-6) \quad & S \rightarrow AB \\
 & A \rightarrow a/_b \\
 & B \rightarrow b/a_
 \end{aligned}$$

will generate the tree

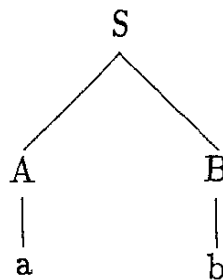


Figure 16-8.

if the rules are interpreted as node admissibility conditions but not if they are interpreted as tree generating rules (the problem being that the  $A$  cannot be rewritten until the  $B$  has, and vice versa.

## 16.5 The Chomsky Hierarchy

By putting increasingly stringent restrictions on the allowed forms of rules we can establish a series of grammars of decreasing generative power. Many such series are imaginable, but the one which has received the most attention is due to Chomsky and has come to be known as the Chomsky Hierarchy. At the top are the most general grammars of the sort we defined above in Section 16.2. There are no restrictions on the form of the rules except that the left side must contain at least one non-terminal symbol (Actually, even this restriction could be eliminated in favor of one which says simply that the left side cannot be the empty string. The formulation we have chosen is essentially a technical convenience) Chomsky dubbed such grammars ‘Type 0,’ and they are also sometimes called unrestricted rewriting systems (urs). The succeeding three types are as follows:

**Type 1:** each rule is of the form  $\alpha A \beta \rightarrow \alpha \psi \beta$ , where  $\psi \neq \epsilon$ .

**Type 2:** each rule is of the form  $A \rightarrow \psi$ .

**Type 3:** each rule is of the form  $A \rightarrow xB$  or  $A \rightarrow x$

In the above  $\alpha$ ,  $\beta$ , and  $\psi$  are arbitrary strings (possibly empty unless otherwise specified) over the union of the terminal and non-terminal alphabets;  $A$  and  $B$  are non-terminals, and  $x$  is a string of terminal symbols

Type 1 grammars are also called *context sensitive*; an equivalent formulation is to say that each rule is of the form  $\psi \rightarrow \omega$ , where  $\omega$  is at least as long as  $\psi$  (i.e., the rules are “non-shrinking”). Type 2 grammars are called *context free*, and Type 3 grammars are called *regular* or *right linear* for reasons which will become apparent in the next section.

Note that these classes of grammars do not form a strict hierarchy in the sense that each type is a subclass of the one with the next lower number. Every Type 1 grammar is also a Type 0 grammar, but because rules of the form  $A \rightarrow \epsilon$  are allowed in Type 2 grammars, these are not properly contained in Type 1. Type 3 grammars, however, are properly contained in the Type 2 grammars. It is nonetheless apparent, technical details concerning the empty string aside, that the hierarchy represents a series of generally increasing restrictions on the allowed form of rules.

The question then arises whether the languages generated by such grammars stand in an analogous relationship. We say that a language is of Type

$n$  ( $n = 0, 1, 2,$  or  $3$ ) iff it is generated by some grammar of Type  $n$ . For example, we saw in Section 16.2 that  $L = \{x \in \{a, b\}^* \mid x \text{ contains equal numbers of } a\text{'s and } b\text{'s}\}$  is of Type 0 inasmuch as it is generated by the grammar given in 16-4. But one might wonder whether it could also be generated by a grammar of some other type—say of Type 2. This is indeed the case; this language is generated by the following Type 2 grammar:

(16-7)  $G = \langle \{a, b\}, \{S, A, B\}, S, R \rangle$  where

$$R = \left\{ \begin{array}{ll} S \rightarrow e & A \rightarrow a \\ S \rightarrow aB & A \rightarrow aS \\ S \rightarrow bA & A \rightarrow bAA \\ B \rightarrow b & B \rightarrow aBB \\ B \rightarrow bS & \end{array} \right\}$$

This fact immediately establishes this language as Type 0 also, since every Type 2 grammar is perforce a Type 0 grammar. (It does not at the same time establish it as a Type 1 language since the given grammar is not Type 1, because of the rule  $S \rightarrow e$ . In fact, this language could not be Type 1 since Type 1 languages can never contain  $e$ .)

Is this language also Type 3? It turns out that it is not, but to prove this is not a simple matter. One must show somehow that *no* Type 3 grammar, however elaborate, can generate this language. We will consider techniques for proving such results in later sections.

Note that if one has two classes of grammars  $G_i$  and  $G_j$  such that  $G_i$  is properly contained in  $G_j$ , it does not necessarily follow that the corresponding classes of languages stand in the proper subset relation. Because every Type  $i$  grammar is also a Type  $i+1$  grammar it *does* follow that every Type  $i$  language is also a Type  $i+1$  language, i.e.,  $L_i \subseteq L_{i+1}$ . But it might also be the case that every Type  $i+1$  language happens to have some Type  $i$  grammar which generates it. In such a case  $L_i$  is a subset of  $L_{i+1}$  but not a proper subset. Among the earliest results achieved in the study of formal grammars and languages were proofs that the inclusions among the languages of the Chomsky hierarchy are in fact proper inclusions. Specifically,

- (i) the Type 3 languages are properly included in the Type 2 languages;
- (ii) the Type 2 languages not containing the empty string are properly included in the Type 1 languages;
- (iii) the Type 1 languages are properly included in the Type 0 languages.

Some of the proofs will be sketched in the following chapters

## 16.6 Languages and automata

As we mentioned at the beginning of this section, languages can also be characterized by abstract computing devices called automata. Ultimately we will define a hierarchy of automata and establish correspondences between them and the grammars of the Chomsky Hierarchy. This gives us yet another point of view from which to examine the notion of ‘complexity of a language’ which we hope eventually to put to use in characterizing natural language.

Before turning to the detailed study of the various classes of automata, it would be well to make a few general remarks about these devices.

An automaton is an idealized abstract computing machine—that is, it is a mathematical object rather than a physical one. An automaton is characterized by the manner in which it performs computations: for any automaton there is a class of *inputs* to which it reacts, and a class of *outputs* which it produces, the relation between these being determined by the *structure*, or internal organization of the automaton. We will consider only automata whose inputs and outputs are discrete (e.g., strings over an alphabet) rather than continuous (e.g., readings on a dial), and we will not deal with automata whose behavior is probabilistic.

Central to the notion of the structure of an automaton is the concept of a *state*. A state of an automaton is analogous to the arrangement of bits in the memory banks and registers of an actual computer, but since we are abstracting away from physical realizations here, we can think of a state as a characteristic of an automaton which in general changes during the course of a computation and which serves to determine the relationship between inputs and outputs. We will consider only automata which have a finite number of states (cf. a computer whose internal hardware at any given moment can be in only one of a finite number of different arrangements of 1's and 0's.)

An automaton may also have a *memory*. For the simplest automata, the memory consists simply of the states themselves. More powerful automata may be outfitted with additional devices, generally “tapes” on which the machine can read and write symbols and do “scratch work.” Since the amount of memory available on such tapes is potentially unlimited, these machines can in effect overcome the limitations inherent in having only a

finite number of states. We will see that the most powerful automata, Turing machines, are capable in principle of performing any computation for which an explicit set of instructions can be given

Automata may be regarded as devices for computing functions, i.e., for pairing inputs with outputs, but we will normally view them as *acceptors*, i.e., devices which, when given an input, either accept or reject it after some finite amount of computation. In particular, if the input is a string over some alphabet  $A$ , then an automaton can be thought of as the acceptor of some language over  $A$  and the rejector of its complement. As we will see, it is also possible to regard automata as *generators* of strings and languages in a manner similar to grammars