

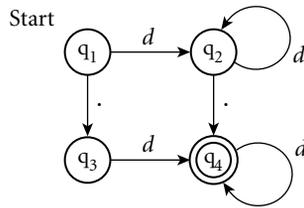
# Programming Language Syntax

## 2.4 Theoretical Foundations

As noted in the main text, scanners and parsers are based on the finite automata and pushdown automata that form the bottom two levels of the Chomsky language hierarchy. At each level of the hierarchy, machines can be either *deterministic* or *nondeterministic*. A deterministic automaton always performs the same operation in a given situation. A nondeterministic automaton can perform any of a *set* of operations. A nondeterministic machine is said to accept a string if there exists a choice of operation in each situation that will eventually lead to the machine saying “yes.” As it turns out, nondeterministic and deterministic finite automata are equally powerful: every DFA is, by definition, a degenerate NFA, and the construction in Example 2.14 (page 56) demonstrates that for any NFA we can create a DFA that accepts the same language. The same is not true of push-down automata: there are context-free languages that are accepted by an NPDA but not by any DPDA. Fortunately, DPDAs suffice in practice to accept the syntax of real programming languages. Practical scanners and parsers are always deterministic.

### 2.4.1 Finite Automata

Precisely defined, a deterministic finite automaton (DFA)  $M$  consists of (1) a finite set  $Q$  of *states*, (2) a finite alphabet  $\Sigma$  of input symbols, (3) a distinguished *initial* state  $q_1 \in Q$ , (4) a set of distinguished *final* states  $F \subseteq Q$ , and (5) a *transition function*  $\delta : Q \times \Sigma \rightarrow Q$  that chooses a new state for  $M$  based on the current state and the current input symbol.  $M$  begins in state  $q_1$ . One by one it consumes its input symbols, using  $\delta$  to move from state to state. When the final symbol has been consumed,  $M$  is interpreted as saying “yes” if it is in a state in  $F$ ; otherwise it is interpreted as saying “no.” We can extend  $\delta$  in the obvious way to take strings, rather than symbols, as inputs, allowing us to say that  $M$  accepts string  $x$  if  $\delta(q_1, x) \in F$ . We can then define  $L(M)$ , the language accepted by  $M$ , to be the set  $\{x \mid \delta(q_1, x) \in F\}$ . In a nondeterministic finite automaton (NFA),



**Figure 2.32** Minimal DFA for the language consisting of all strings of decimal digits containing a single decimal point. Adapted from Figure 2.10 in the main text. The symbol *d* here is short for “0, 1, 2, 3, 4, 5, 6, 7, 8, 9”.

the transition function,  $\delta$ , is multivalued: the automaton can move to any of a set of possible states from a given state on a given input. In addition, it may move from one state to another “spontaneously”; such transitions are said to take input symbol  $\epsilon$ .

**EXAMPLE 2.56**

Formal DFA for  $d^*(.d|d.)d^*$

We can illustrate these definitions with an example. Consider the circles-and-arrows automaton of Figure ©2.32 (adapted from Figure 2.10 in the main text). This is the minimal DFA accepting strings of decimal digits containing a single decimal point.  $\Sigma = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, .\}$  is the machine’s input alphabet.  $Q = \{q_1, q_2, q_3, q_4\}$  is the set of states;  $q_1$  is the initial state;  $F = \{q_4\}$  (a singleton in this case) is the set of final states. The transition function can be represented by a set of triples  $\delta = \{(q_1, 0, q_2), \dots, (q_1, 9, q_2), (q_1, ., q_3), (q_2, 0, q_2), \dots, (q_2, 9, q_2), (q_2, ., q_4), (q_3, 0, q_4), \dots, (q_3, 9, q_4), (q_4, 0, q_4), \dots, (q_4, 9, q_4)\}$ . In each triple  $(q_i, a, q_j)$ ,  $\delta(q_i, a) = q_j$ . ■

Given the constructions of Examples 2.12 and 2.14, we know that there exists an NFA that accepts the language generated by any given regular expression, and a DFA equivalent to any given NFA. To show that regular expressions and finite automata are of equivalent expressive power, all that remains is to demonstrate that there exists a regular expression that generates the language accepted by any given DFA. We illustrate the required construction below for our decimal strings example (Figure ©2.32). More formal and general treatment of all the regular language constructions can be found in standard automata theory texts [HMU01, Sip97].

**From a DFA to a Regular Expression**

To construct a regular expression equivalent to a given DFA, we employ a dynamic programming algorithm that builds solutions to successively more complicated subproblems from a table of solutions to simpler subproblems. We begin with a set of simple regular expressions that describe the transition function,  $\delta$ . For all states  $i$ , we define

$$r_{ii}^0 = a_1 | a_2 | \dots | a_m | \epsilon$$

where  $\{a_1 \mid a_2 \mid \dots \mid a_m\} = \{a \mid \delta(q_i, a) = q_i\}$  is the set of characters labeling the “self-loop” from state  $q_i$  back to itself. If there is no such self-loop,  $r_{ij}^0 = \epsilon$ . Similarly, for  $i \neq j$ , we define

$$r_{ij}^0 = a_1 \mid a_2 \mid \dots \mid a_m$$

where  $\{a_1 \mid a_2 \mid \dots \mid a_m\} = \{a \mid \delta(q_i, a) = q_j\}$  is the set of characters labeling the arc from  $q_i$  to  $q_j$ . If there is no such arc,  $r_{ij}^0$  is the empty regular expression. (Note the difference here: we can stay in state  $q_i$  by not accepting any input, so  $\epsilon$  is always one of the alternatives in  $r_{ii}^0$ , but not in  $r_{ij}^0$  when  $i \neq j$ .)

Given these  $r^0$  expressions, the dynamic programming algorithm inductively calculates expressions  $r_{ij}^k$  with larger superscripts. In each,  $k$  names the highest-numbered state through which control may pass on the way from  $q_i$  to  $q_j$ . We assume that states are numbered starting with  $q_1$ , so when  $k = 0$  we must transition directly from  $q_i$  to  $q_j$ , with no intervening states.

#### EXAMPLE 2.57

Reconstructing a regular expression for the decimal string DFA

In our small example DFA,  $r_{11}^0 = r_{33}^0 = \epsilon$ , and  $r_{22}^0 = r_{44}^0 = 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \mid \epsilon$ , which we will abbreviate  $d \mid \epsilon$ . Similarly,  $r_{13}^0 = r_{24}^0 = \cdot$ , and  $r_{12}^0 = r_{34}^0 = d$ . Expressions  $r_{14}^0, r_{21}^0, r_{23}^0, r_{31}^0, r_{32}^0, r_{41}^0, r_{42}^0$ , and  $r_{43}^0$  are all empty.

For  $k > 0$ , the  $r_{ij}^k$  expressions will generally generate multicharacter strings. At each step of the dynamic programming algorithm, we let

$$r_{ij}^k = r_{ij}^{k-1} \mid r_{ik}^{k-1} r_{kk}^{k-1} r_{kj}^{k-1}$$

That is, to get from  $q_i$  to  $q_j$  without going through any states numbered higher than  $k$ , we can either go from  $q_i$  to  $q_j$  without going through any state numbered higher than  $k - 1$  (which we already know how to do), or else we can go from  $q_i$  to  $q_k$  (without going through any state numbered higher than  $k - 1$ ), travel out from  $q_k$  and back again an arbitrary number of times (never visiting a state numbered higher than  $k - 1$  in between), and finally go from  $q_k$  to  $q_j$  (again without visiting a state numbered higher than  $k - 1$ ). If any of the constituent regular expressions is empty, we omit its term of the outermost alternation. At the end, our overall answer is  $r_{1f_1}^n \mid r_{1f_2}^n \mid \dots \mid r_{1f_l}^n$ , where  $n = |Q|$  is the total number of states and  $F = \{q_{f_1}, q_{f_2}, \dots, q_{f_l}\}$  is the set of final states.

Because  $r_{11}^0 = \epsilon$  and there are no transitions from States 2, 3, or 4 to State 1, nothing changes in the first inductive step in our example; that is,  $\forall i [r_{ii}^1 = r_{ii}^0]$ . The second step is a bit more interesting. Since we are now allowed to go through State 2, we have  $r_{22}^2 = r_{22}^1 r_{22}^1 = (d \mid \epsilon) \mid (d \mid \epsilon) \mid (d \mid \epsilon)^* \mid (d \mid \epsilon) \mid = d^+$ . Because  $r_{21}^1, r_{23}^1, r_{32}^1$ , and  $r_{42}^1$  are empty, however,  $r_{11}^2, r_{33}^2$ , and  $r_{44}^2$  remain the same as  $r_{11}^1, r_{33}^1$ , and  $r_{44}^1$ . In a similar vein, we have

$$r_{12}^2 = d \mid d(d \mid \epsilon)^*(d \mid \epsilon) = d^+$$

$$r_{14}^2 = d(d \mid \epsilon)^* = d^+$$

$$r_{24}^2 = \cdot \mid (d \mid \epsilon)(d \mid \epsilon)^* = d^*$$

Missing transitions and empty expressions from the previous step leave  $r_{13}^2 = r_{13}^1 = \cdot$  and  $r_{34}^2 = r_{34}^1 = d$ . Expressions  $r_{21}^2, r_{23}^2, r_{31}^2, r_{32}^2, r_{41}^2, r_{42}^2$ , and  $r_{43}^2$  remain empty.

In the third inductive step, we have

$$\begin{aligned} r_{13}^3 &= \cdot \mid \cdot \epsilon^* \epsilon = \cdot \\ r_{14}^3 &= d^+ \cdot \mid \cdot \epsilon^* d = d^+ \cdot \mid \cdot d \\ r_{34}^3 &= d \mid \epsilon \epsilon^* d = d \end{aligned}$$

All other expressions remain unchanged from the previous step.

Finally, we have

$$\begin{aligned} r_{14}^4 &= (d^+ \cdot \mid \cdot d) \mid (d^+ \cdot \mid \cdot d) (d \mid \epsilon)^* (d \mid \epsilon) \\ &= (d^+ \cdot \mid \cdot d) \mid (d^+ \cdot \mid \cdot d) d^* \\ &= (d^+ \cdot \mid \cdot d) d^* \\ &= d^+ \cdot d^* \mid \cdot d^+ \end{aligned}$$

Since  $F$  has a single member ( $q_4$ ), this expression is our final answer. ■

### Space Requirements

In Section 2.2.1 we noted without proof that the conversion from an NFA to a DFA may lead to exponential blow-up in the number of states. Certainly this did not happen in our decimal string example: the NFA of Figure 2.8 has 14 states, while the equivalent DFA of Figure 2.9 has only 7, and the minimal DFA of Figures 2.10 and ©2.32 has only 4.

#### EXAMPLE 2.58

A regular language with a large minimal DFA

Consider, however, the subset of  $(a \mid b \mid c)^*$  in which some letter appears at least three times. The minimal DFA for this language has 28 states. As shown in Figure ©2.33, 27 of these are states in which we have seen  $i, j$ , and  $k$  as, bs, and cs, respectively. The 28th (and only final) state is reached once we have seen at least three of some specific character.

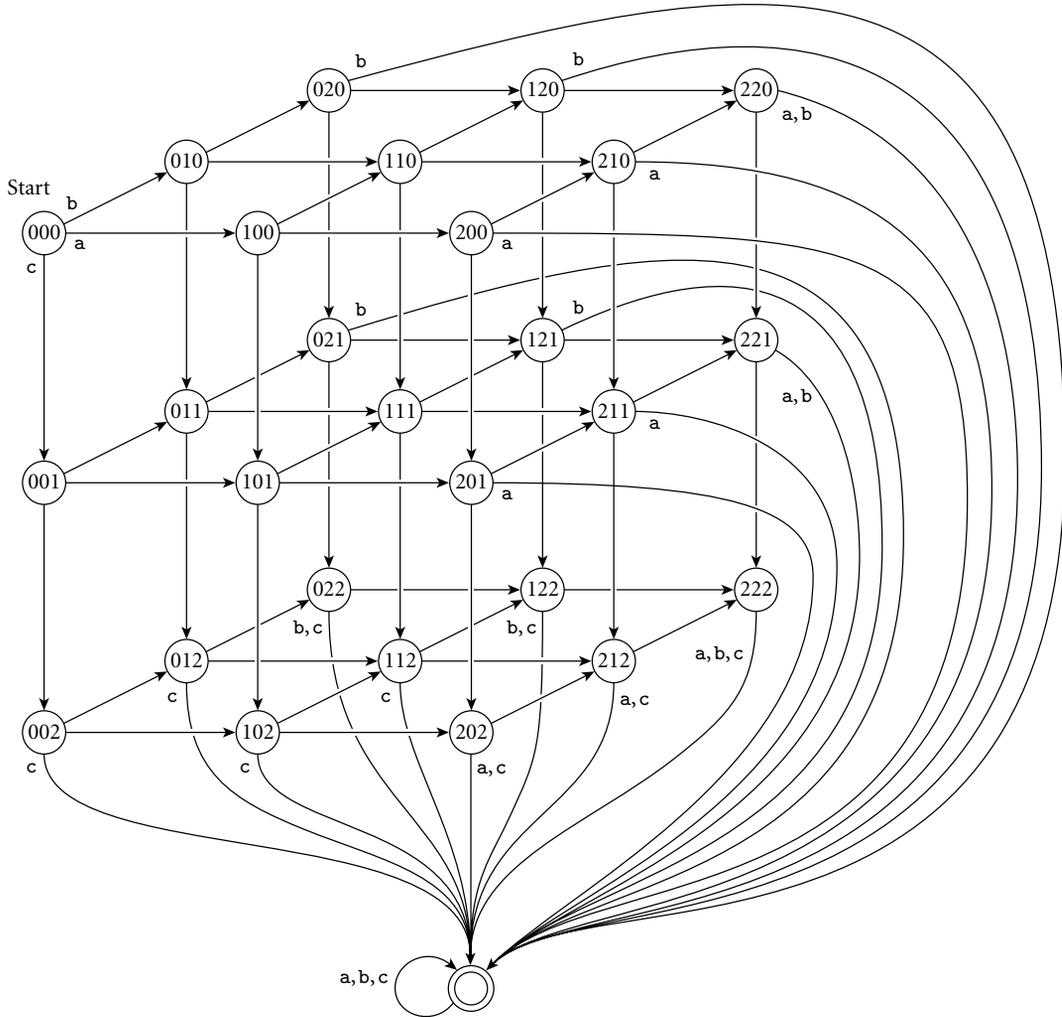
By contrast, there exists an NFA for this language with only eight states, as shown in Figure ©2.34. It requires that we “guess,” at the outset, whether we will see three as, three bs, or three cs. It mirrors the structure of the natural regular expression  $(a \mid b \mid c)^* a (a \mid b \mid c)^* a (a \mid b \mid c)^* \mid (a \mid b \mid c)^* b (a \mid b \mid c)^* b (a \mid b \mid c)^* \mid (a \mid b \mid c)^* c (a \mid b \mid c)^* c (a \mid b \mid c)^*$ . ■

Of course, the eight-state NFA does not emerge directly from the construction of Section 2.2.1; that construction produces a 52-state machine with a certain amount of redundancy, and with many extraneous states and epsilon productions.

#### EXAMPLE 2.59

Exponential DFA blow-up

But consider the similar subset of  $(0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9)^*$  in which some digit appears at least ten times. The minimal DFA for this language has 10,000,000,001 states: a non-final state for each combination of zeros through nines with less than ten of each, and a single final state reached once any digit has appeared at least ten times. One possible regular expression for this language is



**Figure 2.33** DFA for the language consisting of all strings in  $(a | b | c)^*$  in which some letter appears at least three times. State name  $ijk$  indicates that  $i$  as,  $j$  bs, and  $k$  cs have been seen so far. Within the cubic portion of the figure, most edge labels are elided: **a** transitions move to the right, **b** transitions go back into the page, and **c** transitions move down.

$$\begin{aligned}
 & ((0 | 1 | \dots | 9)^* 0 (0 | 1 | \dots | 9)^* 0 (0 | 1 | \dots | 9)^* 0 (0 | 1 | \dots | 9)^* 0 \\
 & (0 | 1 | \dots | 9)^* 0 \\
 & (0 | 1 | \dots | 9)^* 0 (0 | 1 | \dots | 9)^* 0 (0 | 1 | \dots | 9)^* ) \\
 & | ((0 | 1 | \dots | 9)^* 1 (0 | 1 | \dots | 9)^* 1 (0 | 1 | \dots | 9)^* 1 (0 | 1 | \dots | 9)^* 1 \\
 & (0 | 1 | \dots | 9)^* 1 \\
 & (0 | 1 | \dots | 9)^* 1 (0 | 1 | \dots | 9)^* 1 (0 | 1 | \dots | 9)^* ) \\
 & | \dots \\
 & | ((0 | 1 | \dots | 9)^* 9 (0 | 1 | \dots | 9)^* 9 (0 | 1 | \dots | 9)^* 9 (0 | 1 | \dots | 9)^* 9 \\
 & (0 | 1 | \dots | 9)^* 9 \\
 & (0 | 1 | \dots | 9)^* 9 (0 | 1 | \dots | 9)^* 9 (0 | 1 | \dots | 9)^* )
 \end{aligned}$$

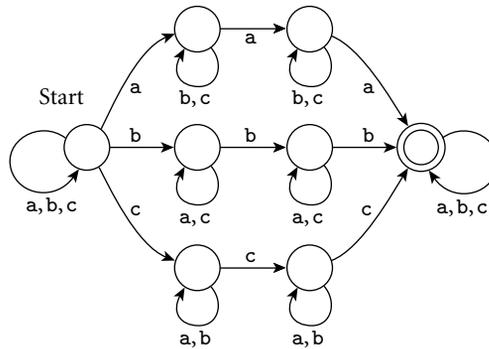


Figure 2.34 NFA corresponding to the DFA of Figure 2.33.

Our construction would yield a very large NFA for this expression, but clearly many orders of magnitude smaller than ten billion states! ■

### 2.4.2 Push-Down Automata

A deterministic push-down automaton (DPDA)  $N$  consists of (1)  $Q$ , (2)  $\Sigma$ , (3)  $q_1$ , and (4)  $F$ , as in a DFA, plus (6) a finite alphabet  $\Gamma$  of stack symbols, (7) a distinguished initial stack symbol  $Z_1 \in \Gamma$ , and (5') a transition function  $\delta : Q \times \Gamma \times \{\Sigma \cup \{\epsilon\}\} \rightarrow Q \times \Gamma^*$ , where  $\Gamma^*$  is the set of strings of zero or more symbols from  $\Gamma$ .  $N$  begins in state  $q_1$ , with symbol  $Z_1$  in an otherwise empty stack. It repeatedly examines the current state  $q$  and top-of-stack symbol  $Z$ . If  $\delta(q, \epsilon, Z)$  is defined,  $N$  moves to state  $r$  and replaces  $Z$  with  $\alpha$  in the stack, where  $(r, \alpha) = \delta(q, \epsilon, Z)$ . In this case  $N$  does not consume its input symbol. If  $\delta(q, \epsilon, Z)$  is undefined,  $N$  examines and consumes the current input symbol  $a$ . It then moves to state  $s$  and replaces  $Z$  with  $\beta$ , where  $(s, \beta) = \delta(q, a, Z)$ .  $N$  is interpreted as accepting a string of input symbols if and only if it consumes the symbols and ends in a state in  $F$ .

As with finite automata, a nondeterministic push-down automaton (NPDA) is distinguished by a multivalued transition function: an NPDA can choose any of a set of new states and stack symbol replacements when faced with a given state, input, and top-of-stack symbol. If  $\delta(q, \epsilon, Z)$  is nonempty,  $N$  can also choose a new state and stack symbol replacement without inspecting or consuming its current input symbol. While we have seen that nondeterministic and deterministic finite automata are equally powerful, this correspondence does not carry over to push-down automata: there are context-free languages that are accepted by an NPDA but not by any DPDA.

The proof that CFGs and NPDAs are equivalent in expressive power is more complex than the corresponding proof for regular expressions and finite automata. The proof is also of limited practical importance for compiler construction; we do

not present it here. While it is possible to create an NPDA for any CFL, that NPDA may in some cases require exponential time to recognize strings in the language. (The  $O(n^3)$  algorithms mentioned in Section 2.3 do not take the form of PDAs.) Practical programming languages can all be expressed with LL or LR grammars, which can be parsed with a (deterministic) PDA in linear time.

An LL(1) PDA is very simple. Because it makes decisions solely on the basis of the current input token and top-of-stack symbol, its state diagram is trivial. All but one of the transitions is a self-loop from the initial state to itself. A final transition moves from the initial state to a second, final state when it sees \$\$ on the input and the stack. As we noted in Section 2.3.3 (page 91), the state diagram for an SLR(1) or LALR(1) parser is substantially more interesting: it's the characteristic finite-state machine (CFSM). Full LR(1) parsers have similar machines, but usually with many more states, due to the need for path-specific look-ahead.

A little study reveals that if we define every state to be accepting, then the CFSM, without its stack, is a DFA that recognizes the grammar's *viable prefixes*. These are all the strings of grammar symbols that can begin a sentential form in the canonical (right-most) derivation of some string in the language, and that do not extend beyond the end of the handle. The algorithms to construct LL(1) and SLR(1) PDAs from suitable grammars were given in Sections 2.3.2 and 2.3.3.

### 2.4.3 Grammar and Language Classes

#### EXAMPLE 2.60

$0^n 1^n$  is not a regular language

As we noted in Section 2.1.2, a scanner is incapable of recognizing arbitrarily nested constructs. The key to the proof is to realize that we cannot count an arbitrary number of left-bracketing symbols with a finite number of states. Consider, for example, the problem of accepting the language  $0^n 1^n$ . Suppose there is a DFA  $M$  that accepts this language. Suppose further that  $M$  has  $m$  states. Now suppose we feed  $M$  a string of  $m + 1$  zeros. By the *pigeonhole principle* (you can't distribute  $m$  objects among  $p < m$  pigeonholes without putting at least two objects in some pigeonhole),  $M$  must enter some state  $q_i$  twice while scanning this string. Without loss of generality, let us assume it does so after seeing  $j$  zeros and again after seeing  $k$  zeros, for  $j \neq k$ . Since we know that  $M$  accepts the string  $0^j 1^j$  and the string  $0^k 1^k$ , and since it is in precisely the same state after reading  $0^j$  and  $0^k$ , we can deduce that  $M$  must also accept the strings  $0^j 1^k$  and  $0^k 1^j$ . Since these strings are not in the language, we have a contradiction:  $M$  cannot exist. ■

Within the family of context-free languages, one can prove similar theorems about the constructs that can and cannot be recognized using various parsing algorithms. Though almost all real parsers get by with a single token of look-ahead, it is possible in principle to use more than one, thereby expanding the set of grammars that can be parsed in linear time. In the top-down case we can redefine FIRST and FOLLOWsets to contain pairs of tokens in a more or less straightforward fashion. If we do this, however, we encounter a more serious

version of the immediate error detection problem described in Section ©2.3.4. There we saw that the use of context-independent FOLLOW sets could cause us to overlook a syntax error until after we had needlessly predicted one or more epsilon productions. Context-specific FOLLOW sets solved the problem, but did not change the set of *valid* programs that could be parsed with one token of look-ahead. If we define  $LL(k)$  to be the set of all grammars that can be parsed predictively using the top-of-stack symbol and  $k$  tokens of look-ahead, then it turns out that for  $k > 1$  we must adopt a context-specific notion of FOLLOW sets in order to parse correctly. The algorithm of Section 2.3.2, which is based on context-independent FOLLOW sets, is actually known as SLL (simple LL), rather than true LL. For  $k = 1$ , the  $LL(1)$  and  $SLL(1)$  algorithms can parse the same set of grammars. For  $k > 1$ , LL is strictly more powerful. Among the bottom-up parsers, the relationships among  $SLR(k)$ ,  $LALR(k)$ , and  $LR(k)$  are somewhat more complicated, but extra look-ahead always helps.

**EXAMPLE 2.61**

Separation of grammar classes

Containment relationships among the classes of grammars accepted by popular linear-time algorithms appear in Figure ©2.35. The LR class (no suffix) contains every grammar  $G$  for which there exists a  $k$  such that  $G \in LR(k)$ ; LL, SLL, SLR, and LALR are similarly defined. Grammars can be found in every region of the figure. Examples appear in Figure ©2.36. Proofs that they lie in the regions claimed are deferred to Exercise ©2.30. ■

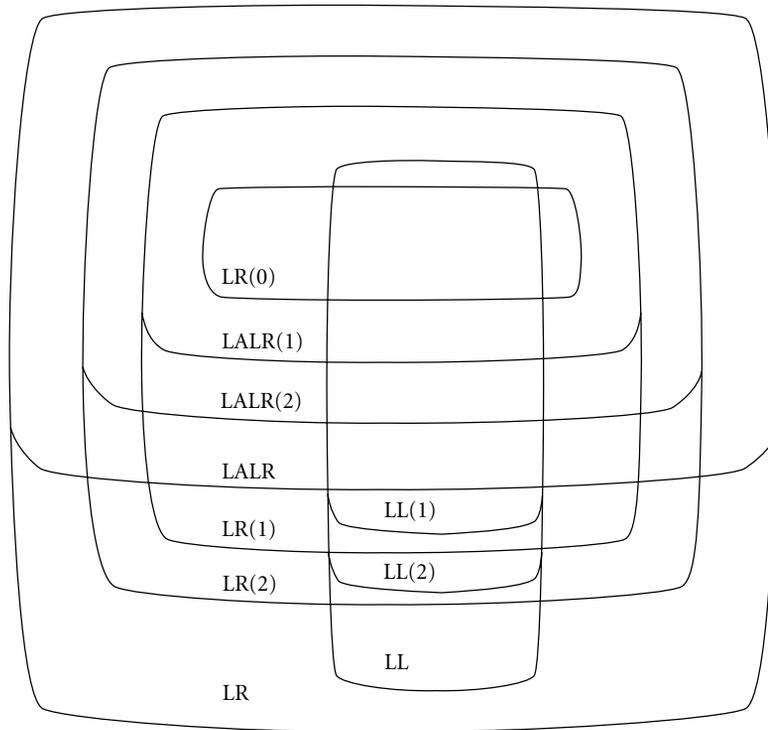
For any context-free grammar  $G$  and parsing algorithm  $P$ , we say that  $G$  is a  $P$  grammar (e.g., an  $LL(1)$  grammar) if it can be parsed using that algorithm. By extension, for any context-free language  $L$ , we say that  $L$  is a  $P$  language if there exists a  $P$  grammar for  $L$  (this may not be the grammar we were given). Containment relationships among the classes of languages accepted by the popular parsing algorithms appear in Figure ©2.37. Again, languages can be found in every region. Examples appear in Figure ©2.38; proofs are deferred to Exercise ©2.31. ■

**EXAMPLE 2.62**

Separation of language classes

Note that every context-free language that can be parsed deterministically has an  $SLR(1)$  grammar. Moreover, any language that can be parsed deterministically and in which no valid string can be extended to create another valid string (this is called the *prefix property*) has an  $LR(0)$  grammar. If we restrict our attention to languages with an explicit  $$$$  marker at end-of-file, then they all have the prefix property, and therefore  $LR(0)$  grammars.

The relationships among language classes are not as rich as the relationships among grammar classes. Most real programming languages can be parsed by any of the popular parsing algorithms, though the grammars are not always pretty, and special-purpose “hacks” may sometimes be required when a language is almost, but not quite, in a given class. The principal advantage of the more powerful parsing algorithms (e.g., full LR) is that they can parse a wider variety of grammars for a given language. In practice this flexibility makes it easier for the compiler writer to find a grammar that is intuitive and readable, and that facilitates the creation of semantic action routines.



**Figure 2.35** Containment relationships among popular grammar classes. In addition to the containments shown,  $SLL(k)$  is just inside  $LL(k)$ , for  $k \geq 2$ , but has the same relationship to everything else, and  $SLR(k)$  is just inside  $LALR(k)$ , for  $k \geq 1$ , but has the same relationship to everything else.

$LL(2)$  but not  $SLL$ :

$$\begin{aligned} S &\rightarrow a A a \mid b A b a \\ A &\rightarrow b \mid \epsilon \end{aligned}$$

$SLL(k)$  but not  $LL(k-1)$ :

$$S \rightarrow a^{k-1} b \mid a^k$$

$LR(0)$  but not  $LL$ :

$$\begin{aligned} S &\rightarrow A b \\ A &\rightarrow A a \mid a \end{aligned}$$

$SLL(1)$  but not  $LALR$ :

$$\begin{aligned} S &\rightarrow A a \mid B b \mid c C \\ C &\rightarrow A b \mid B a \\ A &\rightarrow D \\ B &\rightarrow D \\ D &\rightarrow \epsilon \end{aligned}$$

$SLL(k)$  and  $SLR(k)$  but not  $LR(k-1)$ :

$$\begin{aligned} S &\rightarrow A a^{k-1} b \mid B a^{k-1} c \\ A &\rightarrow \epsilon \\ B &\rightarrow \epsilon \end{aligned}$$

$LALR(1)$  but not  $SLR$ :

$$\begin{aligned} S &\rightarrow b A b \mid A c \mid a b \\ A &\rightarrow a \end{aligned}$$

$LR(1)$  but not  $LALR$ :

$$\begin{aligned} S &\rightarrow a C a \mid b C b \mid a D b \mid b D a \\ C &\rightarrow c \\ D &\rightarrow c \end{aligned}$$

Unambiguous but not  $LR$ :

$$S \rightarrow a S a \mid \epsilon$$

**Figure 2.36** Examples of grammars in various regions of Figure ©2.35.

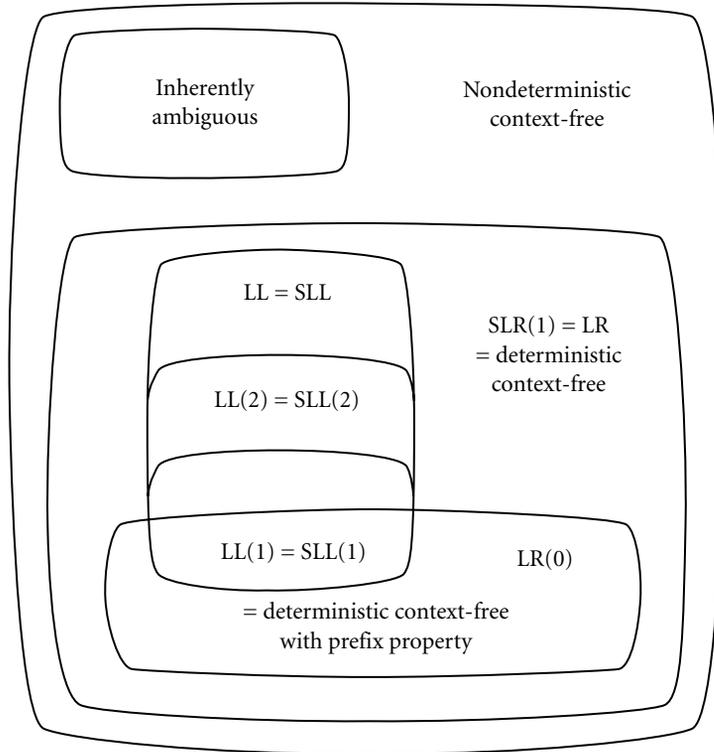


Figure 2.37 Containment relationships among popular language classes.

Nondeterministic language:

$$\{a^n b^n c : n \geq 1\} \cup \{a^n b^{2n} d : n \geq 1\}$$

Inherently ambiguous language:

$$\{a^i b^j c^k : i = j \text{ or } j = k; i, j, k \geq 1\}$$

Language with LL(k) grammar but no LL(k-1) grammar:

$$\{a^n (b | c | b^k d)^n : n \geq 1\}$$

Language with LR(0) grammar but no LL grammar:

$$\{a^n b^n : n \geq 1\} \cup \{a^n c^n : n \geq 1\}$$

Figure 2.38 Examples of languages in various regions of Figure 2.37.

**✓ CHECK YOUR UNDERSTANDING**

- 56. What formal machine captures the behavior of a scanner? A parser?
- 57. State three ways in which a real scanner differs from the formal machine.

58. What are the formal components of a DFA?
  59. Outline the algorithm used to construct a regular expression equivalent to a given DFA.
  60. What is the inherent “big-O” complexity of parsing with an NPDA? Why is this worse than the  $O(n^3)$  time mentioned in Section 2.3?
  61. How many states are there in an LL(1) PDA? An SLR(1) PDA? Explain.
  62. What are the *viable prefixes* of a CFG?
  63. Summarize the proof that a DFA cannot recognize arbitrarily nested constructs.
  64. Explain the difference between LL and SLL parsing.
  65. Is every LL(1) grammar also LR(1)? Is it LALR(1)?
  66. Does every LR language have an SLR(1) grammar?
  67. Why are the containment relationships among grammar classes more complex than those among language classes?
-

