# Chapter 3

# Lexical Analysis

In this chapter we show how to construct a lexical analyzer. To implement a lexical analyzer by hand, it helps to start with a diagram or other description for the lexemes of each token. We can then write code to identify each occurrence of each lexeme on the input and to return information about the token identified.

We can also produce a lexical analyzer automatically by specifying the lexeme patterns to a *lexical-analyzer generator* and compiling those patterns into code that functions as a lexical analyzer. This approach makes it easier to modify a lexical analyzer, since we have only to rewrite the affected patterns, not the entire program. It also speeds up the process of implementing the lexical analyzer, since the programmer specifies the software at the very high level of patterns and relies on the generator to produce the detailed code. We shall introduce in Section 3.5 a lexical-analyzer generator called *Lex* (or *Flex* in a more recent embodiment).

We begin the study of lexical-analyzer generators by introducing regular expressions, a convenient notation for specifying lexeme patterns. We show how this notation can be transformed, first into nondeterministic automata and then into deterministic automata. The latter two notations can be used as input to a "driver," that is, code which simulates these automata and uses them as a guide to determining the next token. This driver and the specification of the automaton form the nucleus of the lexical analyzer.

## 3.1 The Role of the Lexical Analyzer

As the first phase of a compiler, the main task of the lexical analyzer is to read the input characters of the source program, group them into lexemes, and produce as output a sequence of tokens for each lexeme in the source program. The stream of tokens is sent to the parser for syntax analysis. It is common for the lexical analyzer to interact with the symbol table as well. When the lexical analyzer discovers a lexeme constituting an identifier, it needs to enter that lexeme into the symbol table. In some cases, information regarding the

kind of identifier may be read from the symbol table by the lexical analyzer to assist it in determining the proper token it must pass to the parser.

These interactions are suggested in Fig. 3.1. Commonly, the interaction is implemented by having the parser call the lexical analyzer. The call, suggested by the *getNextToken* command, causes the lexical analyzer to read characters from its input until it can identify the next lexeme and produce for it the next token, which it returns to the parser.
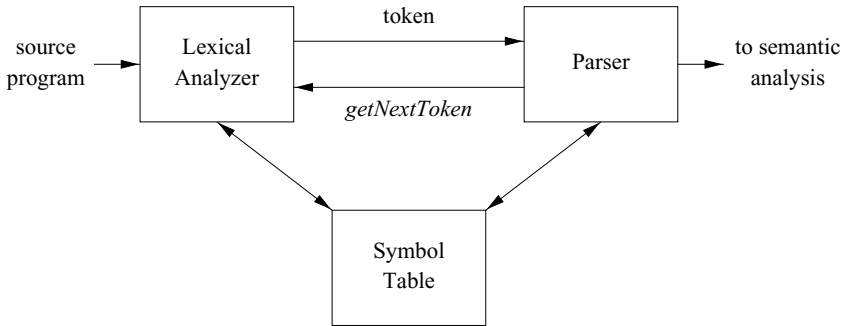
Figure 3.1: Interactions between the lexical analyzer and the parser

Since the lexical analyzer is the part of the compiler that reads the source text, it may perform certain other tasks besides identification of lexemes. One such task is stripping out comments and *whitespace* (blank, newline, tab, and perhaps other characters that are used to separate tokens in the input). Another task is correlating error messages generated by the compiler with the source program. For instance, the lexical analyzer may keep track of the number of newline characters seen, so it can associate a line number with each error message. In some compilers, the lexical analyzer makes a copy of the source program with the error messages inserted at the appropriate positions. If the source program uses a macro-preprocessor, the expansion of macros may also be performed by the lexical analyzer.

Sometimes, lexical analyzers are divided into a cascade of two processes:

a) *Scanning* consists of the simple processes that do not require tokenization of the input, such as deletion of comments and compaction of consecutive whitespace characters into one.

b) *Lexical analysis* proper is the more complex portion, which produces tokens from the output of the scanner.

## 3.1.1   Lexical Analysis Versus Parsing

There are a number of reasons why the analysis portion of a compiler is normally separated into lexical analysis and parsing (syntax analysis) phases.

1. Simplicity of design is the most important consideration. The separation of lexical and syntactic analysis often allows us to simplify at least one of these tasks. For example, a parser that had to deal with comments and whitespace as syntactic units would be considerably more complex than one that can assume comments and whitespace have already been removed by the lexical analyzer. If we are designing a new language, separating lexical and syntactic concerns can lead to a cleaner overall language design.

2. Compiler efficiency is improved. A separate lexical analyzer allows us to apply specialized techniques that serve only the lexical task, not the job of parsing. In addition, specialized buffering techniques for reading input characters can speed up the compiler significantly.

3. Compiler portability is enhanced. Input-device-specific peculiarities can be restricted to the lexical analyzer.

## 3.1.2 Tokens, Patterns, and Lexemes

When discussing lexical analysis, we use three related but distinct terms:

- A *token* is a pair consisting of a token name and an optional attribute value. The token name is an abstract symbol representing a kind of lexical unit, e.g., a particular keyword, or a sequence of input characters denoting an identifier. The token names are the input symbols that the parser processes. In what follows, we shall generally write the name of a token in boldface. We will often refer to a token by its token name.

- A *pattern* is a description of the form that the lexemes of a token may take. In the case of a keyword as a token, the pattern is just the sequence of characters that form the keyword. For identifiers and some other tokens, the pattern is a more complex structure that is *matched* by many strings.

- A *lexeme* is a sequence of characters in the source program that matches the pattern for a token and is identified by the lexical analyzer as an instance of that token.

**Example 3.1 :** Figure 3.2 gives some typical tokens, their informally described patterns, and some sample lexemes. To see how these concepts are used in practice, in the C statement

```
printf("Total = %d\n", score);
```

both `printf` and `score` are lexemes matching the pattern for token **id**, and `"Total = %d\n"` is a lexeme matching **literal**.  □

In many programming languages, the following classes cover most or all of the tokens:

| TOKEN | INFORMAL DESCRIPTION | SAMPLE LEXEMES |
|---|---|---|
| **if** | characters `i`, `f` | `if` |
| **else** | characters `e`, `l`, `s`, `e` | `else` |
| **comparison** | `<` or `>` or `<=` or `>=` or `==` or `!=` | `<=`, `!=` |
| **id** | letter followed by letters and digits | `pi`, `score`, `D2` |
| **number** | any numeric constant | `3.14159`, `0`, `6.02e23` |
| **literal** | anything but `"`, surrounded by `"`'s | `"core dumped"` |

Figure 3.2: Examples of tokens

1. One token for each keyword. The pattern for a keyword is the same as the keyword itself.

2. Tokens for the operators, either individually or in classes such as the token `comparison` mentioned in Fig. 3.2.

3. One token representing all identifiers.

4. One or more tokens representing constants, such as numbers and literal strings.

5. Tokens for each punctuation symbol, such as left and right parentheses, comma, and semicolon.

### 3.1.3   Attributes for Tokens

When more than one lexeme can match a pattern, the lexical analyzer must provide the subsequent compiler phases additional information about the particular lexeme that matched. For example, the pattern for token **number** matches both 0 and 1, but it is extremely important for the code generator to know which lexeme was found in the source program. Thus, in many cases the lexical analyzer returns to the parser not only a token name, but an attribute value that describes the lexeme represented by the token; the token name influences parsing decisions, while the attribute value influences translation of tokens after the parse.

We shall assume that tokens have at most one associated attribute, although this attribute may have a structure that combines several pieces of information. The most important example is the token **id**, where we need to associate with the token a great deal of information. Normally, information about an identifier — e.g., its lexeme, its type, and the location at which it is first found (in case an error message about that identifier must be issued) — is kept in the symbol table. Thus, the appropriate attribute value for an identifier is a pointer to the symbol-table entry for that identifier.

---

**Tricky Problems When Recognizing Tokens**

Usually, given the pattern describing the lexemes of a token, it is relatively simple to recognize matching lexemes when they occur on the input. However, in some languages it is not immediately apparent when we have seen an instance of a lexeme corresponding to a token. The following example is taken from Fortran, in the fixed-format still allowed in Fortran 90. In the statement

```
DO 5 I = 1.25
```

it is not apparent that the first lexeme is DO5I, an instance of the identifier token, until we see the dot following the 1. Note that blanks in fixed-format Fortran are ignored (an archaic convention). Had we seen a comma instead of the dot, we would have had a do-statement

```
DO 5 I = 1,25
```

in which the first lexeme is the keyword DO.

---

**Example 3.2:** The token names and associated attribute values for the Fortran statement

```
E = M * C ** 2
```

are written below as a sequence of pairs.

> $<$**id**, pointer to symbol-table entry for E$>$
> $<$**assign_op**$>$
> $<$**id**, pointer to symbol-table entry for M$>$
> $<$**mult_op**$>$
> $<$**id**, pointer to symbol-table entry for C$>$
> $<$**exp_op**$>$
> $<$**number**, integer value 2$>$

Note that in certain pairs, especially operators, punctuation, and keywords, there is no need for an attribute value. In this example, the token **number** has been given an integer-valued attribute. In practice, a typical compiler would instead store a character string representing the constant and use as an attribute value for **number** a pointer to that string. □

## 3.1.4 Lexical Errors

It is hard for a lexical analyzer to tell, without the aid of other components, that there is a source-code error. For instance, if the string fi is encountered for the first time in a C program in the context:

```
fi ( a == f(x)) ...
```

a lexical analyzer cannot tell whether `fi` is a misspelling of the keyword `if` or an undeclared function identifier. Since `fi` is a valid lexeme for the token **id**, the lexical analyzer must return the token **id** to the parser and let some other phase of the compiler — probably the parser in this case — handle an error due to transposition of the letters.

However, suppose a situation arises in which the lexical analyzer is unable to proceed because none of the patterns for tokens matches any prefix of the remaining input. The simplest recovery strategy is "panic mode" recovery. We delete successive characters from the remaining input, until the lexical analyzer can find a well-formed token at the beginning of what input is left. This recovery technique may confuse the parser, but in an interactive computing environment it may be quite adequate.

Other possible error-recovery actions are:

1. Delete one character from the remaining input.

2. Insert a missing character into the remaining input.

3. Replace a character by another character.

4. Transpose two adjacent characters.

Transformations like these may be tried in an attempt to repair the input. The simplest such strategy is to see whether a prefix of the remaining input can be transformed into a valid lexeme by a single transformation. This strategy makes sense, since in practice most lexical errors involve a single character. A more general correction strategy is to find the smallest number of transformations needed to convert the source program into one that consists only of valid lexemes, but this approach is considered too expensive in practice to be worth the effort.

### 3.1.5  Exercises for Section 3.1

**Exercise 3.1.1:** Divide the following C++ program:

```
float limitedSquare(x) float x; {
    /* returns x-squared, but never more than 100 */
    return (x<=-10.0||x>=10.0)?100:x*x;
}
```

into appropriate lexemes, using the discussion of Section 3.1.2 as a guide. Which lexemes should get associated lexical values? What should those values be?

! **Exercise 3.1.2:** Tagged languages like HTML or XML are different from conventional programming languages in that the punctuation (tags) are either very numerous (as in HTML) or a user-definable set (as in XML). Further, tags can often have parameters. Suggest how to divide the following HTML document:

```
Here is a photo of <B>my house</B>:
<P><IMG SRC = "house.gif"><BR>
See <A HREF = "morePix.html">More Pictures</A> if you
liked that one.<P>
```

into appropriate lexemes. Which lexemes should get associated lexical values, and what should those values be?

## 3.2 Input Buffering

Before discussing the problem of recognizing lexemes in the input, let us examine some ways that the simple but important task of reading the source program can be speeded. This task is made difficult by the fact that we often have to look one or more characters beyond the next lexeme before we can be sure we have the right lexeme. The box on "Tricky Problems When Recognizing Tokens" in Section 3.1 gave an extreme example, but there are many situations where we need to look at least one additional character ahead. For instance, we cannot be sure we've seen the end of an identifier until we see a character that is not a letter or digit, and therefore is not part of the lexeme for **id**. In C, single-character operators like -, =, or < could also be the beginning of a two-character operator like ->, ==, or <=. Thus, we shall introduce a two-buffer scheme that handles large lookaheads safely. We then consider an improvement involving "sentinels" that saves time checking for the ends of buffers.

### 3.2.1 Buffer Pairs

Because of the amount of time taken to process characters and the large number of characters that must be processed during the compilation of a large source program, specialized buffering techniques have been developed to reduce the amount of overhead required to process a single input character. An important scheme involves two buffers that are alternately reloaded, as suggested in Fig. 3.3.
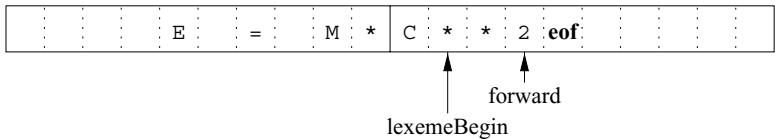


Figure 3.3: Using a pair of input buffers

Each buffer is of the same size $N$, and $N$ is usually the size of a disk block, e.g., 4096 bytes. Using one system read command we can read $N$ characters into a buffer, rather than using one system call per character. If fewer than $N$ characters remain in the input file, then a special character, represented by **eof**,

marks the end of the source file and is different from any possible character of the source program.

Two pointers to the input are maintained:

1. Pointer `lexemeBegin`, marks the beginning of the current lexeme, whose extent we are attempting to determine.

2. Pointer `forward` scans ahead until a pattern match is found; the exact strategy whereby this determination is made will be covered in the balance of this chapter.

Once the next lexeme is determined, `forward` is set to the character at its right end. Then, after the lexeme is recorded as an attribute value of a token returned to the parser, `lexemeBegin` is set to the character immediately after the lexeme just found. In Fig. 3.3, we see `forward` has passed the end of the next lexeme, `**` (the Fortran exponentiation operator), and must be retracted one position to its left.

Advancing `forward` requires that we first test whether we have reached the end of one of the buffers, and if so, we must reload the other buffer from the input, and move `forward` to the beginning of the newly loaded buffer. As long as we never need to look so far ahead of the actual lexeme that the sum of the lexeme's length plus the distance we look ahead is greater than $N$, we shall never overwrite the lexeme in its buffer before determining it.

### 3.2.2  Sentinels

If we use the scheme of Section 3.2.1 as described, we must check, each time we advance `forward`, that we have not moved off one of the buffers; if we do, then we must also reload the other buffer. Thus, for each character read, we make two tests: one for the end of the buffer, and one to determine what character is read (the latter may be a multiway branch). We can combine the buffer-end test with the test for the current character if we extend each buffer to hold a *sentinel* character at the end. The sentinel is a special character that cannot be part of the source program, and a natural choice is the character **eof**.

Figure 3.4 shows the same arrangement as Fig. 3.3, but with the sentinels added. Note that **eof** retains its use as a marker for the end of the entire input. Any **eof** that appears other than at the end of a buffer means that the input is at an end. Figure 3.5 summarizes the algorithm for advancing `forward`. Notice how the first test, which can be part of a multiway branch based on the character pointed to by `forward`, is the only test we make, except in the case where we actually are at the end of a buffer or the end of the input.

## 3.3  Specification of Tokens

Regular expressions are an important notation for specifying lexeme patterns. While they cannot express all possible patterns, they are very effective in spec-

---

### Can We Run Out of Buffer Space?

In most modern languages, lexemes are short, and one or two characters of lookahead is sufficient. Thus a buffer size $N$ in the thousands is ample, and the double-buffer scheme of Section 3.2.1 works without problem. However, there are some risks. For example, if character strings can be very long, extending over many lines, then we could face the possibility that a lexeme is longer than $N$. To avoid problems with long character strings, we can treat them as a concatenation of components, one from each line over which the string is written. For instance, in Java it is conventional to represent long strings by writing a piece on each line and concatenating pieces with a + operator at the end of each piece.

A more difficult problem occurs when arbitrarily long lookahead may be needed. For example, some languages like PL/I do not treat keywords as *reserved*; that is, you can use identifiers with the same name as a keyword like DECLARE. If the lexical analyzer is presented with text of a PL/I program that begins DECLARE ( ARG1, ARG2,... it cannot be sure whether DECLARE is a keyword, and ARG1 and so on are variables being declared, or whether DECLARE is a procedure name with its arguments. For this reason, modern languages tend to reserve their keywords. However, if not, one can treat a keyword like DECLARE as an ambiguous identifier, and let the parser resolve the issue, perhaps in conjunction with symbol-table lookup.

---

ifying those types of patterns that we actually need for tokens. In this section we shall study the formal notation for regular expressions, and in Section 3.5 we shall see how these expressions are used in a lexical-analyzer generator. Then, Section 3.7 shows how to build the lexical analyzer by converting regular expressions to automata that perform the recognition of the specified tokens.

## 3.3.1 Strings and Languages

An *alphabet* is any finite set of symbols. Typical examples of symbols are letters, digits, and punctuation. The set $\{0,1\}$ is the *binary alphabet*. ASCII is an important example of an alphabet; it is used in many software systems. Uni-
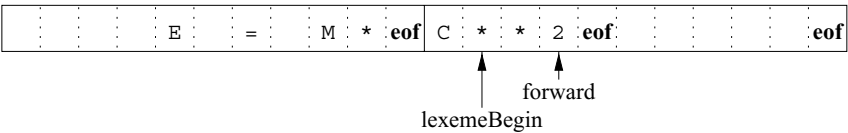


Figure 3.4: Sentinels at the end of each buffer

```
switch ( *forward++ ) {
      case eof:
            if (forward is at end of first buffer ) {
                  reload second buffer;
                  forward = beginning of second buffer;
            }
            else if (forward is at end of second buffer ) {
                  reload first buffer;
                  forward = beginning of first buffer;
            }
            else /* eof within a buffer marks the end of input */
                  terminate lexical analysis;
            break;
      Cases for the other characters
}
```

Figure 3.5: Lookahead code with sentinels

---

### Implementing Multiway Branches

We might imagine that the switch in Fig. 3.5 requires many steps to ex-
ecute, and that placing the case **eof** first is not a wise choice. Actually,
it doesn't matter in what order we list the cases for each character. In
practice, a multiway branch depending on the input character is made in
one step by jumping to an address found in an array of addresses, indexed
by characters.

---

code, which includes approximately 100,000 characters from alphabets around
the world, is another important example of an alphabet.

A *string* over an alphabet is a finite sequence of symbols drawn from that
alphabet. In language theory, the terms "sentence" and "word" are often used
as synonyms for "string." The length of a string $s$, usually written $|s|$, is the
number of occurrences of symbols in $s$. For example, banana is a string of
length six. The *empty string*, denoted $\epsilon$, is the string of length zero.

A *language* is any countable set of strings over some fixed alphabet. This
definition is very broad. Abstract languages like $\emptyset$, the *empty set*, or $\{\epsilon\}$, the
set containing only the empty string, are languages under this definition. So
too are the set of all syntactically well-formed C programs and the set of all
grammatically correct English sentences, although the latter two languages are
difficult to specify exactly. Note that the definition of "language" does not
require that any meaning be ascribed to the strings in the language. Methods
for defining the "meaning" of strings are discussed in Chapter 5.

---

### Terms for Parts of Strings

The following string-related terms are commonly used:

1. A *prefix* of string $s$ is any string obtained by removing zero or more symbols from the end of $s$. For example, `ban`, `banana`, and $\epsilon$ are prefixes of `banana`.

2. A *suffix* of string $s$ is any string obtained by removing zero or more symbols from the beginning of $s$. For example, `nana`, `banana`, and $\epsilon$ are suffixes of `banana`.

3. A *substring* of $s$ is obtained by deleting any prefix and any suffix from $s$. For instance, `banana`, `nan`, and $\epsilon$ are substrings of `banana`.

4. The *proper* prefixes, suffixes, and substrings of a string $s$ are those, prefixes, suffixes, and substrings, respectively, of $s$ that are not $\epsilon$ or not equal to $s$ itself.

5. A *subsequence* of $s$ is any string formed by deleting zero or more not necessarily consecutive positions of $s$. For example, `baan` is a subsequence of `banana`.

---

If $x$ and $y$ are strings, then the *concatenation* of $x$ and $y$, denoted $xy$, is the string formed by appending $y$ to $x$. For example, if $x = $ `dog` and $y = $ `house`, then $xy = $ `doghouse`. The empty string is the identity under concatenation; that is, for any string $s$, $\epsilon s = s\epsilon = s$.

If we think of concatenation as a product, we can define the "exponentiation" of strings as follows. Define $s^0$ to be $\epsilon$, and for all $i > 0$, define $s^i$ to be $s^{i-1}s$. Since $\epsilon s = s$, it follows that $s^1 = s$. Then $s^2 = ss$, $s^3 = sss$, and so on.

## 3.3.2 Operations on Languages

In lexical analysis, the most important operations on languages are union, concatenation, and closure, which are defined formally in Fig. 3.6. Union is the familiar operation on sets. The concatenation of languages is all strings formed by taking a string from the first language and a string from the second language, in all possible ways, and concatenating them. The (*Kleene*) *closure* of a language $L$, denoted $L^*$, is the set of strings you get by concatenating $L$ zero or more times. Note that $L^0$, the "concatenation of $L$ zero times," is defined to be $\{\epsilon\}$, and inductively, $L^i$ is $L^{i-1}L$. Finally, the positive closure, denoted $L^+$, is the same as the Kleene closure, but without the term $L^0$. That is, $\epsilon$ will not be in $L^+$ unless it is in $L$ itself.

| OPERATION | DEFINITION AND NOTATION |
|---|---|
| *Union* of $L$ and $M$ | $L \cup M = \{s \mid s \text{ is in } L \text{ or } s \text{ is in } M\}$ |
| *Concatenation* of $L$ and $M$ | $LM = \{st \mid s \text{ is in } L \text{ and } t \text{ is in } M\}$ |
| *Kleene closure* of $L$ | $L^* = \cup_{i=0}^{\infty} L^i$ |
| *Positive closure* of $L$ | $L^+ = \cup_{i=1}^{\infty} L^i$ |

Figure 3.6: Definitions of operations on languages

**Example 3.3 :** Let $L$ be the set of letters $\{A, B, \ldots, Z, a, b, \ldots, z\}$ and let $D$ be the set of digits $\{0, 1, \ldots 9\}$. We may think of $L$ and $D$ in two, essentially equivalent, ways. One way is that $L$ and $D$ are, respectively, the alphabets of uppercase and lowercase letters and of digits. The second way is that $L$ and $D$ are languages, all of whose strings happen to be of length one. Here are some other languages that can be constructed from languages $L$ and $D$, using the operators of Fig. 3.6:

1. $L \cup D$ is the set of letters and digits — strictly speaking the language with 62 strings of length one, each of which strings is either one letter or one digit.

2. $LD$ is the set of 520 strings of length two, each consisting of one letter followed by one digit.

3. $L^4$ is the set of all 4-letter strings.

4. $L^*$ is the set of all strings of letters, including $\epsilon$, the empty string.

5. $L(L \cup D)^*$ is the set of all strings of letters and digits beginning with a letter.

6. $D^+$ is the set of all strings of one or more digits.

☐

### 3.3.3 Regular Expressions

Suppose we wanted to describe the set of valid C identifiers. It is almost exactly the language described in item (5) above; the only difference is that the underscore is included among the letters.

In Example 3.3, we were able to describe identifiers by giving names to sets of letters and digits and using the language operators union, concatenation, and closure. This process is so useful that a notation called *regular expressions* has come into common use for describing all the languages that can be built from these operators applied to the symbols of some alphabet. In this notation, if *letter_* is established to stand for any letter or the underscore, and *digit* is

established to stand for any digit, then we could describe the language of C identifiers by:

$$letter\_ \ ( \ letter\_ \ | \ digit \ )^*$$

The vertical bar above means union, the parentheses are used to group subexpressions, the star means "zero or more occurrences of," and the juxtaposition of *letter_* with the remainder of the expression signifies concatenation.

The regular expressions are built recursively out of smaller regular expressions, using the rules described below. Each regular expression $r$ denotes a language $L(r)$, which is also defined recursively from the languages denoted by $r$'s subexpressions. Here are the rules that define the regular expressions over some alphabet $\Sigma$ and the languages that those expressions denote.

**BASIS**: There are two rules that form the basis:

1. $\epsilon$ is a regular expression, and $L(\epsilon)$ is $\{\epsilon\}$, that is, the language whose sole member is the empty string.

2. If $a$ is a symbol in $\Sigma$, then **a** is a regular expression, and $L(\mathbf{a}) = \{a\}$, that is, the language with one string, of length one, with $a$ in its one position. Note that by convention, we use italics for symbols, and boldface for their corresponding regular expression.[1]

**INDUCTION**: There are four parts to the induction whereby larger regular expressions are built from smaller ones. Suppose $r$ and $s$ are regular expressions denoting languages $L(r)$ and $L(s)$, respectively.

1. $(r)|(s)$ is a regular expression denoting the language $L(r) \cup L(s)$.

2. $(r)(s)$ is a regular expression denoting the language $L(r)L(s)$.

3. $(r)^*$ is a regular expression denoting $\big(L(r)\big)^*$.

4. $(r)$ is a regular expression denoting $L(r)$. This last rule says that we can add additional pairs of parentheses around expressions without changing the language they denote.

As defined, regular expressions often contain unnecessary pairs of parentheses. We may drop certain pairs of parentheses if we adopt the conventions that:

a) The unary operator $*$ has highest precedence and is left associative.

b) Concatenation has second highest precedence and is left associative.

---

[1] However, when talking about specific characters from the ASCII character set, we shall generally use teletype font for both the character and its regular expression.

c) | has lowest precedence and is left associative.

Under these conventions, for example, we may replace the regular expression $(\mathbf{a})|((\mathbf{b})^*(\mathbf{c}))$ by $\mathbf{a}|\mathbf{b}^*\mathbf{c}$. Both expressions denote the set of strings that are either a single $a$ or are zero or more $b$'s followed by one $c$.

**Example 3.4:** Let $\Sigma = \{a, b\}$.

1. The regular expression $\mathbf{a}|\mathbf{b}$ denotes the language $\{a, b\}$.

2. $(\mathbf{a}|\mathbf{b})(\mathbf{a}|\mathbf{b})$ denotes $\{aa, ab, ba, bb\}$, the language of all strings of length two over the alphabet $\Sigma$. Another regular expression for the same language is $\mathbf{aa}|\mathbf{ab}|\mathbf{ba}|\mathbf{bb}$.

3. $\mathbf{a}^*$ denotes the language consisting of all strings of zero or more $a$'s, that is, $\{\epsilon, a, aa, aaa, \dots\}$.

4. $(\mathbf{a}|\mathbf{b})^*$ denotes the set of all strings consisting of zero or more instances of $a$ or $b$, that is, all strings of $a$'s and $b$'s: $\{\epsilon, a, b, aa, ab, ba, bb, aaa, \dots\}$. Another regular expression for the same language is $(\mathbf{a}^*\mathbf{b}^*)^*$.

5. $\mathbf{a}|\mathbf{a}^*\mathbf{b}$ denotes the language $\{a, b, ab, aab, aaab, \dots\}$, that is, the string $a$ and all strings consisting of zero or more $a$'s and ending in $b$.

□

A language that can be defined by a regular expression is called a *regular set*. If two regular expressions $r$ and $s$ denote the same regular set, we say they are *equivalent* and write $r = s$. For instance, $(\mathbf{a}|\mathbf{b}) = (\mathbf{b}|\mathbf{a})$. There are a number of algebraic laws for regular expressions; each law asserts that expressions of two different forms are equivalent. Figure 3.7 shows some of the algebraic laws that hold for arbitrary regular expressions $r$, $s$, and $t$.

| LAW | DESCRIPTION |
|---|---|
| $r|s = s|r$ | | is commutative |
| $r|(s|t) = (r|s)|t$ | | is associative |
| $r(st) = (rs)t$ | Concatenation is associative |
| $r(s|t) = rs|rt;\ (s|t)r = sr|tr$ | Concatenation distributes over | |
| $\epsilon r = r\epsilon = r$ | $\epsilon$ is the identity for concatenation |
| $r^* = (r|\epsilon)^*$ | $\epsilon$ is guaranteed in a closure |
| $r^{**} = r^*$ | * is idempotent |

Figure 3.7: Algebraic laws for regular expressions

### 3.3.4 Regular Definitions

For notational convenience, we may wish to give names to certain regular expressions and use those names in subsequent expressions, as if the names were themselves symbols. If $\Sigma$ is an alphabet of basic symbols, then a *regular definition* is a sequence of definitions of the form:

$$
\begin{array}{ccc}
d_1 & \rightarrow & r_1 \\
d_2 & \rightarrow & r_2 \\
& \cdots & \\
d_n & \rightarrow & r_n
\end{array}
$$

where:

1. Each $d_i$ is a new symbol, not in $\Sigma$ and not the same as any other of the $d$'s, and

2. Each $r_i$ is a regular expression over the alphabet $\Sigma \cup \{d_1, d_2, \ldots, d_{i-1}\}$.

By restricting $r_i$ to $\Sigma$ and the previously defined $d$'s, we avoid recursive definitions, and we can construct a regular expression over $\Sigma$ alone, for each $r_i$. We do so by first replacing uses of $d_1$ in $r_2$ (which cannot use any of the $d$'s except for $d_1$), then replacing uses of $d_1$ and $d_2$ in $r_3$ by $r_1$ and (the substituted) $r_2$, and so on. Finally, in $r_n$ we replace each $d_i$, for $i = 1, 2, \ldots, n - 1$, by the substituted version of $r_i$, each of which has only symbols of $\Sigma$.

**Example 3.5 :** C identifiers are strings of letters, digits, and underscores. Here is a regular definition for the language of C identifiers. We shall conventionally use italics for the symbols defined in regular definitions.

$$
\begin{array}{rcl}
letter\_ & \rightarrow & \texttt{A} \mid \texttt{B} \mid \cdots \mid \texttt{Z} \mid \texttt{a} \mid \texttt{b} \mid \cdots \mid \texttt{z} \mid \texttt{\_} \\
digit & \rightarrow & \texttt{0} \mid \texttt{1} \mid \cdots \mid \texttt{9} \\
id & \rightarrow & letter\_ \; (\; letter\_ \mid digit \;)^*
\end{array}
$$

□

**Example 3.6 :** Unsigned numbers (integer or floating point) are strings such as `5280`, `0.01234`, `6.336E4`, or `1.89E-4`. The regular definition

$$
\begin{array}{rcl}
digit & \rightarrow & \texttt{0} \mid \texttt{1} \mid \cdots \mid \texttt{9} \\
digits & \rightarrow & digit \; digit^* \\
optionalFraction & \rightarrow & \texttt{.} \; digits \mid \epsilon \\
optionalExponent & \rightarrow & (\; \texttt{E} \; (\; \texttt{+} \mid \texttt{-} \mid \epsilon \;) \; digits \;) \mid \epsilon \\
number & \rightarrow & digits \; optionalFraction \; optionalExponent
\end{array}
$$

is a precise specification for this set of strings. That is, an *optionalFraction* is either a decimal point (dot) followed by one or more digits, or it is missing (the empty string). An *optionalExponent*, if not missing, is the letter E followed by an optional + or − sign, followed by one or more digits. Note that at least one digit must follow the dot, so *number* does not match `1.`, but does match `1.0`. □

### 3.3.5  Extensions of Regular Expressions

Since Kleene introduced regular expressions with the basic operators for union, concatenation, and Kleene closure in the 1950s, many extensions have been added to regular expressions to enhance their ability to specify string patterns. Here we mention a few notational extensions that were first incorporated into Unix utilities such as Lex that are particularly useful in the specification lexical analyzers. The references to this chapter contain a discussion of some regular-expression variants in use today.

1. *One or more instances.*  The unary, postfix operator $^+$ represents the positive closure of a regular expression and its language. That is, if $r$ is a regular expression, then $(r)^+$ denotes the language $\big(L(r)\big)^+$. The operator $^+$ has the same precedence and associativity as the operator $*$. Two useful algebraic laws, $r^* = r^+|\epsilon$ and $r^+ = rr^* = r^*r$ relate the Kleene closure and positive closure.

2. *Zero or one instance.* The unary postfix operator ? means "zero or one occurrence." That is, $r$? is equivalent to $r|\epsilon$, or put another way, $L(r?) = L(r) \cup \{\epsilon\}$. The ? operator has the same precedence and associativity as $*$ and $+$.

3. *Character classes.*  A regular expression $a_1|a_2|\cdots|a_n$, where the $a_i$'s are each symbols of the alphabet, can be replaced by the shorthand $[a_1a_2\cdots a_n]$.  More importantly, when $a_1, a_2, \ldots, a_n$ form a logical sequence, e.g., consecutive uppercase letters, lowercase letters, or digits, we can replace them by $a_1$-$a_n$, that is, just the first and last separated by a hyphen.  Thus, [**abc**] is shorthand for **a**|**b**|**c**, and [**a-z**] is shorthand for **a**|**b**|$\cdots$|**z**.

**Example 3.7:** Using these shorthands, we can rewrite the regular definition of Example 3.5 as:

$$
\begin{array}{rcl}
letter\_ & \to & [\texttt{A-Za-z\_}] \\
digit & \to & [\texttt{0-9}] \\
id & \to & letter\_\ (\ letter\_ \mid digit\ )^*
\end{array}
$$

The regular definition of Example 3.6 can also be simplified:

$$
\begin{array}{rcl}
digit & \to & [\texttt{0-9}] \\
digits & \to & digit^+ \\
number & \to & digits\ (\texttt{.}\ digits)?\ (\ \texttt{E}\ [\texttt{+-}]?\ digits\ )?
\end{array}
$$

□

### 3.3.6 Exercises for Section 3.3

**Exercise 3.3.1:** Consult the language reference manuals to determine (*i*) the sets of characters that form the input alphabet (excluding those that may only appear in character strings or comments), (*ii*) the lexical form of numerical constants, and (*iii*) the lexical form of identifiers, for each of the following languages: (a) C (b) C++ (c) C# (d) Fortran (e) Java (f) Lisp (g) SQL.

**! Exercise 3.3.2:** Describe the languages denoted by the following regular expressions:

   a) $\mathbf{a(a|b)^*a}$.

   b) $\mathbf{((\epsilon|a)b^*)^*}$.

   c) $\mathbf{(a|b)^*a(a|b)(a|b)}$.

   d) $\mathbf{a^*ba^*ba^*ba^*}$.

   **!!** e) $\mathbf{(aa|bb)^*((ab|ba)(aa|bb)^*(ab|ba)(aa|bb)^*)^*}$.

**Exercise 3.3.3:** In a string of length $n$, how many of the following are there?

   a) Prefixes.

   b) Suffixes.

   c) Proper prefixes.

   **!** d) Substrings.

   **!** e) Subsequences.

**Exercise 3.3.4:** Most languages are *case sensitive*, so keywords can be written only one way, and the regular expressions describing their lexemes are very simple. However, some languages, like SQL, are *case insensitive*, so a keyword can be written either in lowercase or in uppercase, or in any mixture of cases. Thus, the SQL keyword `SELECT` can also be written `select`, `Select`, or `sElEcT`, for instance. Show how to write a regular expression for a keyword in a case-insensitive language. Illustrate the idea by writing the expression for "select" in SQL.

**! Exercise 3.3.5:** Write regular definitions for the following languages:

   a) All strings of lowercase letters that contain the five vowels in order.

   b) All strings of lowercase letters in which the letters are in ascending lexicographic order.

   c) Comments, consisting of a string surrounded by `/*` and `*/`, without an intervening `*/`, unless it is inside double-quotes (`"`).

!! d) All strings of digits with no repeated digits. *Hint*: Try this problem first with a few digits, such as $\{0, 1, 2\}$.

!! e) All strings of digits with at most one repeated digit.

!! f) All strings of $a$'s and $b$'s with an even number of $a$'s and an odd number of $b$'s.

g) The set of Chess moves, in the informal notation, such as p-k4 or kbp×qn.

!! h) All strings of $a$'s and $b$'s that do not contain the substring *abb*.

i) All strings of $a$'s and $b$'s that do not contain the subsequence *abb*.

**Exercise 3.3.6:** Write character classes for the following sets of characters:

a) The first ten letters (up to "j") in either upper or lower case.

b) The lowercase consonants.

c) The "digits" in a hexadecimal number (choose either upper or lower case for the "digits" above 9).

d) The characters that can appear at the end of a legitimate English sentence (e.g., exclamation point).

The following exercises, up to and including Exercise 3.3.10, discuss the extended regular-expression notation from Lex (the lexical-analyzer generator that we shall discuss extensively in Section 3.5). The extended notation is listed in Fig. 3.8.

**Exercise 3.3.7:** Note that these regular expressions give all of the following symbols (*operator characters*) a special meaning:

    \ " . ^ $ [ ] * + ? { } | /

Their special meaning must be turned off if they are needed to represent themselves in a character string. We can do so by quoting the character within a string of length one or more; e.g., the regular expression "**" matches the string **. We can also get the literal meaning of an operator character by preceding it by a backslash. Thus, the regular expression \*\* also matches the string **. Write a regular expression that matches the string "\.

**Exercise 3.3.8:** In Lex, a *complemented character class* represents any character except the ones listed in the character class. We denote a complemented class by using ^ as the first character; this symbol (caret) is not itself part of the class being complemented, unless it is listed within the class itself. Thus, [^A-Za-z] matches any character that is not an uppercase or lowercase letter, and [^\^] represents any character but the caret (or newline, since newline cannot be in any character class). Show that for every regular expression with complemented character classes, there is an equivalent regular expression without complemented character classes.

| EXPRESSION | MATCHES | EXAMPLE |
|---|---|---|
| $c$ | the one non-operator character $c$ | `a` |
| $\backslash c$ | character $c$ literally | `\*` |
| $"s"$ | string $s$ literally | `"**"` |
| $.$ | any character but newline | `a.*b` |
| $\hat{\ }$ | beginning of a line | `^abc` |
| $\$$ | end of a line | `abc$` |
| $[s]$ | any one of the characters in string $s$ | `[abc]` |
| $[\hat{\ }s]$ | any one character not in string $s$ | `[^abc]` |
| $r*$ | zero or more strings matching $r$ | `a*` |
| $r+$ | one or more strings matching $r$ | `a+` |
| $r?$ | zero or one $r$ | `a?` |
| $r\{m,n\}$ | between $m$ and $n$ occurrences of $r$ | `a{1,5}` |
| $r_1 r_2$ | an $r_1$ followed by an $r_2$ | `ab` |
| $r_1 \mid r_2$ | an $r_1$ or an $r_2$ | `a\|b` |
| $(r)$ | same as $r$ | `(a\|b)` |
| $r_1/r_2$ | $r_1$ when followed by $r_2$ | `abc/123` |

Figure 3.8: `Lex` regular expressions

**! Exercise 3.3.9 :** The regular expression $r\{m,n\}$ matches from $m$ to $n$ occurrences of the pattern $r$. For example, `a{1,5}` matches a string of one to five $a$'s. Show that for every regular expression containing repetition operators of this form, there is an equivalent regular expression without repetition operators.

**! Exercise 3.3.10 :** The operator ^ matches the left end of a line, and $ matches the right end of a line. The operator ^ is also used to introduce complemented character classes, but the context always makes it clear which meaning is intended. For example, `^[^aeiou]*$` matches any complete line that does not contain a lowercase vowel.

   a) How do you tell which meaning of ^ is intended?

   b) Can you always replace a regular expression using the ^ and $ operators by an equivalent expression that does not use either of these operators?

**! Exercise 3.3.11 :** The UNIX shell command `sh` uses the operators in Fig. 3.9 in filename expressions to describe sets of file names. For example, the filename expression `*.o` matches all file names ending in `.o`; `sort1.?` matches all filenames of the form `sort1.c`, where $c$ is any character. Show how `sh` filename

| EXPRESSION | MATCHES | EXAMPLE |
|---|---|---|
| $'s'$ | string $s$ literally | '\' |
| $\backslash c$ | character $c$ literally | \' |
| $*$ | any string | *.o |
| ? | any character | sort1.? |
| $[s]$ | any character in $s$ | sort1.[cso] |

Figure 3.9: Filename expressions used by the shell command sh

expressions can be replaced by equivalent regular expressions using only the basic union, concatenation, and closure operators.

**! Exercise 3.3.12 :** SQL allows a rudimentary form of patterns in which two characters have special meaning: underscore (_) stands for any one character and percent-sign (%) stands for any string of 0 or more characters. In addition, the programmer may define any character, say $e$, to be the escape character, so an $e$ preceding _, %, or another $e$ gives the character that follows its literal meaning. Show how to express any SQL pattern as a regular expression, given that we know which character is the escape character.

## 3.4    Recognition of Tokens

In the previous section we learned how to express patterns using regular expressions. Now, we must study how to take the patterns for all the needed tokens and build a piece of code that examines the input string and finds a prefix that is a lexeme matching one of the patterns. Our discussion will make use of the following running example.

$$
\begin{aligned}
stmt &\rightarrow \quad \textbf{if } expr \textbf{ then } stmt \\
&\mid \quad \textbf{if } expr \textbf{ then } stmt \textbf{ else } stmt \\
&\mid \quad \epsilon \\
expr &\rightarrow \quad term \textbf{ relop } term \\
&\mid \quad term \\
term &\rightarrow \quad \textbf{id} \\
&\mid \quad \textbf{number}
\end{aligned}
$$

Figure 3.10: A grammar for branching statements

**Example 3.8 :** The grammar fragment of Fig. 3.10 describes a simple form of branching statements and conditional expressions. This syntax is similar to that of the language Pascal, in that **then** appears explicitly after conditions.

For **relop**, we use the comparison operators of languages like Pascal or SQL, where = is "equals" and <> is "not equals," because it presents an interesting structure of lexemes.

The terminals of the grammar, which are **if**, **then**, **else**, **relop**, **id**, and **number**, are the names of tokens as far as the lexical analyzer is concerned. The patterns for these tokens are described using regular definitions, as in Fig. 3.11. The patterns for *id* and *number* are similar to what we saw in Example 3.7.

$$
\begin{aligned}
digit &\rightarrow & &\texttt{[0-9]} \\
digits &\rightarrow & &digit^{+} \\
number &\rightarrow & &digits\ (\texttt{.}\ digits)?\ (\ \texttt{E}\ \texttt{[+-]?}\ digits\ )? \\
letter &\rightarrow & &\texttt{[A-Za-z]} \\
id &\rightarrow & &letter\ (\ letter\ |\ digit\ )^{*} \\
if &\rightarrow & &\texttt{if} \\
then &\rightarrow & &\texttt{then} \\
else &\rightarrow & &\texttt{else} \\
relop &\rightarrow & &\texttt{<}\ |\ \texttt{>}\ |\ \texttt{<=}\ |\ \texttt{>=}\ |\ \texttt{=}\ |\ \texttt{<>}
\end{aligned}
$$

Figure 3.11: Patterns for tokens of Example 3.8

For this language, the lexical analyzer will recognize the keywords `if`, `then`, and `else`, as well as lexemes that match the patterns for *relop*, *id*, and *number*. To simplify matters, we make the common assumption that keywords are also *reserved words*: that is, they are not identifiers, even though their lexemes match the pattern for identifiers.

In addition, we assign the lexical analyzer the job of stripping out white-space, by recognizing the "token" *ws* defined by:

$$ws \rightarrow (\ \textbf{blank}\ |\ \textbf{tab}\ |\ \textbf{newline}\ )^{+}$$

Here, **blank**, **tab**, and **newline** are abstract symbols that we use to express the ASCII characters of the same names. Token *ws* is different from the other tokens in that, when we recognize it, we do not return it to the parser, but rather restart the lexical analysis from the character that follows the whitespace. It is the following token that gets returned to the parser.

Our goal for the lexical analyzer is summarized in Fig. 3.12. That table shows, for each lexeme or family of lexemes, which token name is returned to the parser and what attribute value, as discussed in Section 3.1.3, is returned. Note that for the six relational operators, symbolic constants LT, LE, and so on are used as the attribute value, in order to indicate which instance of the token **relop** we have found. The particular operator found will influence the code that is output from the compiler. □

| LEXEMES | TOKEN NAME | ATTRIBUTE VALUE |
|---|---|---|
| Any *ws* | − | − |
| if | **if** | − |
| then | **then** | − |
| else | **else** | − |
| Any *id* | **id** | Pointer to table entry |
| Any *number* | **number** | Pointer to table entry |
| < | **relop** | LT |
| <= | **relop** | LE |
| = | **relop** | EQ |
| <> | **relop** | NE |
| > | **relop** | GT |
| >= | **relop** | GE |

Figure 3.12: Tokens, their patterns, and attribute values

### 3.4.1   Transition Diagrams

As an intermediate step in the construction of a lexical analyzer, we first convert patterns into stylized flowcharts, called "transition diagrams." In this section, we perform the conversion from regular-expression patterns to transition diagrams by hand, but in Section 3.6, we shall see that there is a mechanical way to construct these diagrams from collections of regular expressions.

*Transition diagrams* have a collection of nodes or circles, called *states*. Each state represents a condition that could occur during the process of scanning the input looking for a lexeme that matches one of several patterns. We may think of a state as summarizing all we need to know about what characters we have seen between the *lexemeBegin* pointer and the *forward* pointer (as in the situation of Fig. 3.3).

*Edges* are directed from one state of the transition diagram to another. Each edge is *labeled* by a symbol or set of symbols. If we are in some state $s$, and the next input symbol is $a$, we look for an edge out of state $s$ labeled by $a$ (and perhaps by other symbols, as well). If we find such an edge, we advance the *forward* pointer and enter the state of the transition diagram to which that edge leads. We shall assume that all our transition diagrams are *deterministic*, meaning that there is never more than one edge out of a given state with a given symbol among its labels. Starting in Section 3.5, we shall relax the condition of determinism, making life much easier for the designer of a lexical analyzer, although trickier for the implementer. Some important conventions about transition diagrams are:

1. Certain states are said to be *accepting*, or *final*. These states indicate that a lexeme has been found, although the actual lexeme may not consist of all positions between the *lexemeBegin* and *forward* pointers. We always

indicate an accepting state by a double circle, and if there is an action
to be taken — typically returning a token and an attribute value to the
parser — we shall attach that action to the accepting state.

2. In addition, if it is necessary to retract the *forward* pointer one position
   (i.e., the lexeme does not include the symbol that got us to the accepting
   state), then we shall additionally place a * near that accepting state. In
   our example, it is never necessary to retract *forward* by more than one
   position, but if it were, we could attach any number of *'s to the accepting
   state.

3. One state is designated the *start state*, or *initial state*; it is indicated by
   an edge, labeled "start," entering from nowhere. The transition diagram
   always begins in the start state before any input symbols have been read.

**Example 3.9 :** Figure 3.13 is a transition diagram that recognizes the lexemes
matching the token **relop**. We begin in state 0, the start state. If we see < as the
first input symbol, then among the lexemes that match the pattern for **relop**
we can only be looking at <, <>, or <=. We therefore go to state 1, and look at
the next character. If it is =, then we recognize lexeme <=, enter state 2, and
return the token **relop** with attribute LE, the symbolic constant representing
this particular comparison operator. If in state 1 the next character is >, then
instead we have lexeme <>, and enter state 3 to return an indication that the
not-equals operator has been found. On any other character, the lexeme is <,
and we enter state 4 to return that information. Note, however, that state 4
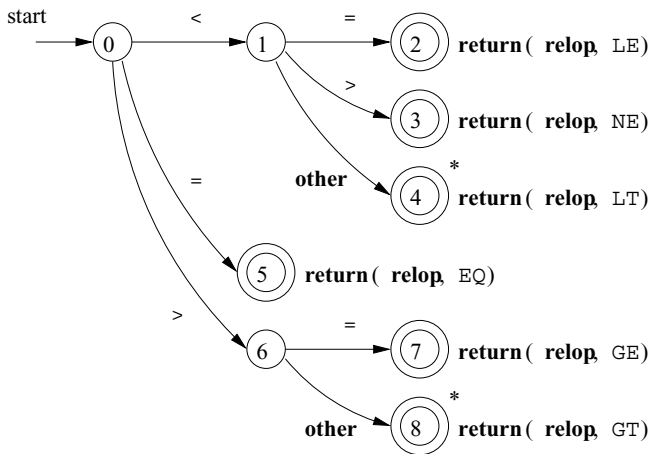has a * to indicate that we must retract the input one position.



Figure 3.13: Transition diagram for **relop**

On the other hand, if in state 0 the first character we see is =, then this one
character must be the lexeme. We immediately return that fact from state 5.

The remaining possibility is that the first character is >. Then, we must enter state 6 and decide, on the basis of the next character, whether the lexeme is >= (if we next see the = sign), or just > (on any other character). Note that if, in state 0, we see any character besides <, =, or >, we can not possibly be seeing a relop lexeme, so this transition diagram will not be used.   □

## 3.4.2   Recognition of Reserved Words and Identifiers

Recognizing keywords and identifiers presents a problem. Usually, keywords like if or then are reserved (as they are in our running example), so they are not identifiers even though they *look* like identifiers. Thus, although we typically use a transition diagram like that of Fig. 3.14 to search for identifier lexemes, this diagram will also recognize the keywords if, then, and else of our running example.



Figure 3.14: A transition diagram for **id**'s and keywords

There are two ways that we can handle reserved words that look like identifiers:

1. Install the reserved words in the symbol table initially. A field of the symbol-table entry indicates that these strings are never ordinary identifiers, and tells which token they represent. We have supposed that this method is in use in Fig. 3.14. When we find an identifier, a call to *installID* places it in the symbol table if it is not already there and returns a pointer to the symbol-table entry for the lexeme found. Of course, any identifier not in the symbol table during lexical analysis cannot be a reserved word, so its token is **id**. The function *getToken* examines the symbol table entry for the lexeme found, and returns whatever token name the symbol table says this lexeme represents — either **id** or one of the keyword tokens that was initially installed in the table.

2. Create separate transition diagrams for each keyword; an example for the keyword then is shown in Fig. 3.15. Note that such a transition diagram consists of states representing the situation after each successive letter of the keyword is seen, followed by a test for a "nonletter-or-digit," i.e., any character that cannot be the continuation of an identifier. It is necessary to check that the identifier has ended, or else we would return token **then** in situations where the correct token was **id**, with a lexeme like thenextvalue that has then as a proper prefix. If we adopt this approach, then we must prioritize the tokens so that the reserved-word

tokens are recognized in preference to **id**, when the lexeme matches both patterns. We *do not* use this approach in our example, which is why the states in Fig. 3.15 are unnumbered.
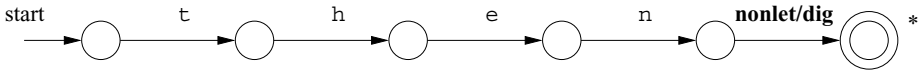


Figure 3.15: Hypothetical transition diagram for the keyword `then`

## 3.4.3 Completion of the Running Example

The transition diagram for **id**'s that we saw in Fig. 3.14 has a simple structure. Starting in state 9, it checks that the lexeme begins with a letter and goes to state 10 if so. We stay in state 10 as long as the input contains letters and digits. When we first encounter anything but a letter or digit, we go to state 11 and accept the lexeme found. Since the last character is not part of the identifier, we must retract the input one position, and as discussed in Section 3.4.2, we enter what we have found in the symbol table and determine whether we have a keyword or a true identifier.

The transition diagram for token **number** is shown in Fig. 3.16, and is so far the most complex diagram we have seen. Beginning in state 12, if we see a digit, we go to state 13. In that state, we can read any number of additional digits. However, if we see anything but a digit, dot, or E, we have seen a number in the form of an integer; 123 is an example. That case is handled by entering state 20, where we return token **number** and a pointer to a table of constants where the found lexeme is entered. These mechanics are not shown on the diagram but are analogous to the way we handled identifiers.
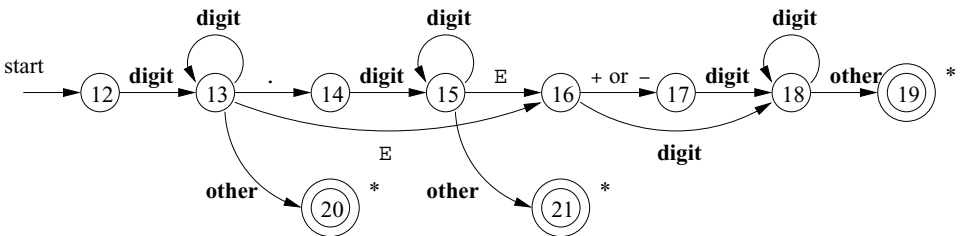


Figure 3.16: A transition diagram for unsigned numbers

If we instead see a dot in state 13, then we have an "optional fraction." State 14 is entered, and we look for one or more additional digits; state 15 is used for that purpose. If we see an E, then we have an "optional exponent," whose recognition is the job of states 16 through 19. Should we, in state 15, instead see anything but E or a digit, then we have come to the end of the fraction, there is no exponent, and we return the lexeme found, via state 21.

The final transition diagram, shown in Fig. 3.17, is for whitespace. In that diagram, we look for one or more "whitespace" characters, represented by **delim** in that diagram — typically these characters would be blank, tab, newline, and perhaps other characters that are not considered by the language design to be part of any token.
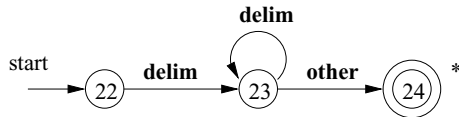


Figure 3.17: A transition diagram for whitespace

Note that in state 24, we have found a block of consecutive whitespace characters, followed by a nonwhitespace character. We retract the input to begin at the nonwhitespace, but we do not return to the parser. Rather, we must restart the process of lexical analysis after the whitespace.

### 3.4.4 Architecture of a Transition-Diagram-Based Lexical Analyzer

There are several ways that a collection of transition diagrams can be used to build a lexical analyzer. Regardless of the overall strategy, each state is represented by a piece of code. We may imagine a variable `state` holding the number of the current state for a transition diagram. A switch based on the value of `state` takes us to code for each of the possible states, where we find the action of that state. Often, the code for a state is itself a switch statement or multiway branch that determines the next state by reading and examining the next input character.

**Example 3.10:** In Fig. 3.18 we see a sketch of `getRelop()`, a C++ function whose job is to simulate the transition diagram of Fig. 3.13 and return an object of type `TOKEN`, that is, a pair consisting of the token name (which must be **relop** in this case) and an attribute value (the code for one of the six comparison operators in this case). `getRelop()` first creates a new object `retToken` and initializes its first component to `RELOP`, the symbolic code for token **relop**.

We see the typical behavior of a state in case 0, the case where the current state is 0. A function `nextChar()` obtains the next character from the input and assigns it to local variable $c$. We then check $c$ for the three characters we expect to find, making the state transition dictated by the transition diagram of Fig. 3.13 in each case. For example, if the next input character is =, we go to state 5.

If the next input character is not one that can begin a comparison operator, then a function `fail()` is called. What `fail()` does depends on the global error-recovery strategy of the lexical analyzer. It should reset the `forward` pointer to `lexemeBegin`, in order to allow another transition diagram to be applied to

```
TOKEN getRelop()
{
    TOKEN retToken = new(RELOP);
    while(1) { /* repeat character processing until a return
                  or failure occurs */
        switch(state) {
            case 0: c = nextChar();
                    if ( c == '<' ) state = 1;
                    else if ( c == '=' ) state = 5;
                    else if ( c == '>' ) state = 6;
                    else fail(); /* lexeme is not a relop */
                    break;
            case 1: ...
            ...
            case 8: retract();
                    retToken.attribute = GT;
                    return(retToken);
        }
    }
}
```

Figure 3.18: Sketch of implementation of **relop** transition diagram

the true beginning of the unprocessed input. It might then change the value of `state` to be the start state for another transition diagram, which will search for another token. Alternatively, if there is no other transition diagram that remains unused, `fail()` could initiate an error-correction phase that will try to repair the input and find a lexeme, as discussed in Section 3.1.4.

We also show the action for state 8 in Fig. 3.18. Because state 8 bears a *, we must retract the input pointer one position (i.e., put $c$ back on the input stream). That task is accomplished by the function `retract()`. Since state 8 represents the recognition of lexeme >, we set the second component of the returned object, which we suppose is named `attribute`, to `GT`, the code for this operator. □

To place the simulation of one transition diagram in perspective, let us consider the ways code like Fig. 3.18 could fit into the entire lexical analyzer.

1. We could arrange for the transition diagrams for each token to be tried sequentially. Then, the function `fail()` of Example 3.10 resets the pointer `forward` and starts the next transition diagram, each time it is called. This method allows us to use transition diagrams for the individual keywords, like the one suggested in Fig. 3.15. We have only to use these before we use the diagram for **id**, in order for the keywords to be reserved words.

2. We could run the various transition diagrams "in parallel," feeding the next input character to all of them and allowing each one to make whatever transitions it required. If we use this strategy, we must be careful to resolve the case where one diagram finds a lexeme that matches its pattern, while one or more other diagrams are still able to process input. The normal strategy is to take the longest prefix of the input that matches any pattern. That rule allows us to prefer identifier `thenext` to keyword `then`, or the operator `->` to `-`, for example.

3. The preferred approach, and the one we shall take up in the following sections, is to combine all the transition diagrams into one. We allow the transition diagram to read input until there is no possible next state, and then take the longest lexeme that matched any pattern, as we discussed in item (2) above. In our running example, this combination is easy, because no two tokens can start with the same character; i.e., the first character immediately tells us which token we are looking for. Thus, we could simply combine states 0, 9, 12, and 22 into one start state, leaving other transitions intact. However, in general, the problem of combining transition diagrams for several tokens is more complex, as we shall see shortly.
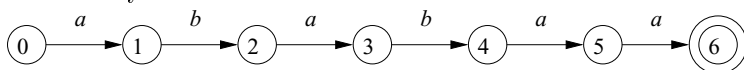
### 3.4.5   Exercises for Section 3.4

**Exercise 3.4.1:** Provide transition diagrams to recognize the same languages as each of the regular expressions in Exercise 3.3.2.

**Exercise 3.4.2:** Provide transition diagrams to recognize the same languages as each of the regular expressions in Exercise 3.3.5.

The following exercises, up to Exercise 3.4.12, introduce the Aho-Corasick algorithm for recognizing a collection of keywords in a text string in time proportional to the length of the text and the sum of the length of the keywords. This algorithm uses a special form of transition diagram called a *trie*. A trie is a tree-structured transition diagram with distinct labels on the edges leading from a node to its children. Leaves of the trie represent recognized keywords.

Knuth, Morris, and Pratt presented an algorithm for recognizing a single keyword $b_1b_2\cdots b_n$ in a text string. Here the trie is a transition diagram with $n + 1$ states, 0 through $n$. State 0 is the initial state, and state $n$ represents acceptance, that is, discovery of the keyword. From each state $s$ from 0 through $n-1$, there is a transition to state $s+1$, labeled by symbol $b_{s+1}$. For example, the trie for the keyword `ababaa` is:



In order to process text strings rapidly and search those strings for a keyword, it is useful to define, for keyword $b_1b_2\cdots b_n$ and position $s$ in that keyword (corresponding to state $s$ of its trie), a *failure function*, $f(s)$, computed as in

Fig. 3.19. The objective is that $b_1 b_2 \cdots b_{f(s)}$ is the longest proper prefix of $b_1 b_2 \cdots b_s$ that is also a suffix of $b_1 b_2 \cdots b_s$. The reason $f(s)$ is important is that if we are trying to match a text string for $b_1 b_2 \cdots b_n$, and we have matched the first $s$ positions, but we then fail (i.e., the next position of the text string does not hold $b_{s+1}$), then $f(s)$ is the longest prefix of $b_1 b_2 \cdots b_n$ that could possibly match the text string up to the point we are at. Of course, the next character of the text string must be $b_{f(s)+1}$, or else we still have problems and must consider a yet shorter prefix, which will be $b_{f(f(s))}$.

```
1)    t = 0;
2)    f(1) = 0;
3)    for (s = 1; s < n; s + +) {
4)            while (t > 0 && b_{s+1} ! =  b_{t+1}) t = f(t);
5)            if (b_{s+1} == b_{t+1}) {
6)                    t = t + 1;
7)                    f(s + 1) = t;
              }
8)            else f(s + 1) = 0;
      }
```

Figure 3.19: Algorithm to compute the failure function for keyword $b_1 b_2 \cdots b_n$

As an example, the failure function for the trie constructed above for `ababaa` is:

| $s$    | 1 | 2 | 3 | 4 | 5 | 6 |
|--------|---|---|---|---|---|---|
| $f(s)$ | 0 | 0 | 1 | 2 | 3 | 1 |

For instance, states 3 and 1 represent prefixes `aba` and `a`, respectively. $f(3) = 1$ because `a` is the longest proper prefix of `aba` that is also a suffix of `aba`. Also, $f(2) = 0$, because the longest proper prefix of `ab` that is also a suffix is the empty string.

**Exercise 3.4.3 :** Construct the failure function for the strings:

a) `abababaab`.

b) `aaaaaa`.

c) `abbaabb`.

! **Exercise 3.4.4 :** Prove, by induction on $s$, that the algorithm of Fig. 3.19 correctly computes the failure function.

!! **Exercise 3.4.5 :** Show that the assignment $t = f(t)$ in line (4) of Fig. 3.19 is executed at most $n$ times. Show that therefore, the entire algorithm takes only $O(n)$ time on a keyword of length $n$.

Having computed the failure function for a keyword $b_1 b_2 \cdots b_n$, we can scan a string $a_1 a_2 \cdots a_m$ in time $O(m)$ to tell whether the keyword occurs in the string. The algorithm, shown in Fig. 3.20, slides the keyword along the string, trying to make progress by matching the next character of the keyword with the next character of the string. If it cannot do so after matching $s$ characters, then it "slides" the keyword right $s - f(s)$ positions, so only the first $f(s)$ characters of the keyword are considered matched with the string.

```
1)    s = 0;
2)    for (i = 1; i ≤ m; i++) {
3)            while (s > 0 && aᵢ ! =  b_{s+1}) s = f(s);
4)            if (aᵢ == b_{s+1}) s = s + 1;
5)            if (s == n) return "yes";
      }
6)    return "no";
```

Figure 3.20: The KMP algorithm tests whether string $a_1 a_2 \cdots a_m$ contains a single keyword $b_1 b_2 \cdots b_n$ as a substring in $O(m + n)$ time

**Exercise 3.4.6:** Apply Algorithm KMP to test whether keyword `ababaa` is a substring of:

a) `abababaab`.

b) `abababbaa`.

!! **Exercise 3.4.7:** Show that the algorithm of Fig. 3.20 correctly tells whether the keyword is a substring of the given string. *Hint*: proceed by induction on $i$. Show that for all $i$, the value of $s$ after line (4) is the length of the longest prefix of the keyword that is a suffix of $a_1 a_2 \cdots a_i$.

!! **Exercise 3.4.8:** Show that the algorithm of Fig. 3.20 runs in time $O(m + n)$, assuming that function $f$ is already computed and its values stored in an array indexed by $s$.

**Exercise 3.4.9:** The *Fibonacci strings* are defined as follows:

1. $s_1 = $ b.

2. $s_2 = $ a.

3. $s_k = s_{k-1} s_{k-2}$ for $k > 2$.

For example, $s_3 = $ ab, $s_4 = $ aba, and $s_5 = $ abaab.

a) What is the length of $s_n$?

b) Construct the failure function for $s_6$.

c) Construct the failure function for $s_7$.

!! d) Show that the failure function for any $s_n$ can be expressed by $f(1) = f(2) = 0$, and for $2 < j \le |s_n|$, $f(j)$ is $j - |s_{k-1}|$, where $k$ is the largest integer such that $|s_k| \le j + 1$.

!! e) In the KMP algorithm, what is the largest number of consecutive applications of the failure function, when we try to determine whether keyword $s_k$ appears in text string $s_{k+1}$?

Aho and Corasick generalized the KMP algorithm to recognize any of a set of keywords in a text string. In this case, the trie is a true tree, with branching from the root. There is one state for every string that is a prefix (not necessarily proper) of any keyword. The parent of a state corresponding to string $b_1 b_2 \cdots b_k$ is the state that corresponds to $b_1 b_2 \cdots b_{k-1}$. A state is accepting if it corresponds to a complete keyword. For example, Fig. 3.21 shows the trie for the keywords he, she, his, and hers.



Figure 3.21: Trie for keywords he, she, his, hers

The failure function for the general trie is defined as follows. Suppose $s$ is the state that corresponds to string $b_1 b_2 \cdots b_n$. Then $f(s)$ is the state that corresponds to the longest proper suffix of $b_1 b_2 \cdots b_n$ that is also a prefix of *some* keyword. For example, the failure function for the trie of Fig. 3.21 is:

| $s$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| $f(s)$ | 0 | 0 | 0 | 1 | 2 | 0 | 3 | 0 | 3 |

! **Exercise 3.4.10:** Modify the algorithm of Fig. 3.19 to compute the failure function for general tries. *Hint*: The major difference is that we cannot simply test for equality or inequality of $b_{s+1}$ and $b_{t+1}$ in lines (4) and (5) of Fig. 3.19. Rather, from any state there may be several transitions out on several characters, as there are transitions on both e and i from state 1 in Fig. 3.21. Any of

those transitions could lead to a state that represents the longest suffix that is also a prefix.

**Exercise 3.4.11 :** Construct the tries and compute the failure function for the following sets of keywords:

    a) `aaa`, `abaaa`, and `ababaaa`.

    b) `all`, `fall`, `fatal`, `llama`, and `lame`.

    c) `pipe`, `pet`, `item`, `temper`, and `perpetual`.

! **Exercise 3.4.12 :** Show that your algorithm from Exercise 3.4.10 still runs in time that is linear in the sum of the lengths of the keywords.

## 3.5     The Lexical-Analyzer Generator `Lex`

In this section, we introduce a tool called `Lex`, or in a more recent implementation `Flex`, that allows one to specify a lexical analyzer by specifying regular expressions to describe patterns for tokens. The input notation for the `Lex` tool is referred to as the *Lex language* and the tool itself is the *Lex compiler*. Behind the scenes, the Lex compiler transforms the input patterns into a transition diagram and generates code, in a file called `lex.yy.c`, that simulates this transition diagram. The mechanics of how this translation from regular expressions to transition diagrams occurs is the subject of the next sections; here we only learn the Lex language.

### 3.5.1    Use of `Lex`

Figure 3.22 suggests how `Lex` is used. An input file, which we call `lex.l`, is written in the Lex language and describes the lexical analyzer to be generated. The Lex compiler transforms `lex.l` to a C program, in a file that is always named `lex.yy.c`. The latter file is compiled by the C compiler into a file called `a.out`, as always. The C-compiler output is a working lexical analyzer that can take a stream of input characters and produce a stream of tokens.

The normal use of the compiled C program, referred to as `a.out` in Fig. 3.22, is as a subroutine of the parser. It is a C function that returns an integer, which is a code for one of the possible token names. The attribute value, whether it be another numeric code, a pointer to the symbol table, or nothing, is placed in a global variable `yylval`,[2] which is shared between the lexical analyzer and parser, thereby making it simple to return both the name and an attribute value of a token.

---

[2] Incidentally, the `yy` that appears in `yylval` and `lex.yy.c` refers to the `Yacc` parser-generator, which we shall describe in Section 4.9, and which is commonly used in conjunction with `Lex`.
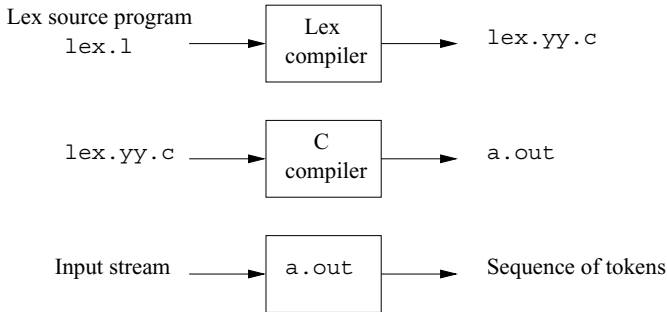
Figure 3.22: Creating a lexical analyzer with Lex

## 3.5.2   Structure of Lex Programs

A Lex program has the following form:

> declarations
> %%
> translation rules
> %%
> auxiliary functions

The declarations section includes declarations of variables, *manifest constants* (identifiers declared to stand for a constant, e.g., the name of a token), and regular definitions, in the style of Section 3.3.4.

The translation rules each have the form

> Pattern   { Action }

Each pattern is a regular expression, which may use the regular definitions of the declaration section. The actions are fragments of code, typically written in C, although many variants of Lex using other languages have been created.

The third section holds whatever additional functions are used in the actions. Alternatively, these functions can be compiled separately and loaded with the lexical analyzer.

The lexical analyzer created by Lex behaves in concert with the parser as follows. When called by the parser, the lexical analyzer begins reading its remaining input, one character at a time, until it finds the longest prefix of the input that matches one of the patterns $P_i$. It then executes the associated action $A_i$. Typically, $A_i$ will return to the parser, but if it does not (e.g., because $P_i$ describes whitespace or comments), then the lexical analyzer proceeds to find additional lexemes, until one of the corresponding actions causes a return to the parser. The lexical analyzer returns a single value, the token name, to the parser, but uses the shared, integer variable yylval to pass additional information about the lexeme found, if needed.

**Example 3.11 :** Figure 3.23 is a `Lex` program that recognizes the tokens of Fig. 3.12 and returns the token found. A few observations about this code will introduce us to many of the important features of `Lex`.

In the declarations section we see a pair of special brackets, `%{` and `%}`. Anything within these brackets is copied directly to the file `lex.yy.c`, and is not treated as a regular definition. It is common to place there the definitions of the manifest constants, using C `#define` statements to associate unique integer codes with each of the manifest constants. In our example, we have listed in a comment the names of the manifest constants, `LT`, `IF`, and so on, but have not shown them defined to be particular integers.[3]

Also in the declarations section is a sequence of regular definitions. These use the extended notation for regular expressions described in Section 3.3.5. Regular definitions that are used in later definitions or in the patterns of the translation rules are surrounded by curly braces. Thus, for instance, *delim* is defined to be a shorthand for the character class consisting of the blank, the tab, and the newline; the latter two are represented, as in all UNIX commands, by backslash followed by `t` or `n`, respectively. Then, *ws* is defined to be one or more delimiters, by the regular expression `{delim}+`.

Notice that in the definition of *id* and *number*, parentheses are used as grouping metasymbols and do not stand for themselves. In contrast, `E` in the definition of *number* stands for itself. If we wish to use one of the `Lex` meta-symbols, such as any of the parentheses, `+`, `*`, or `?`, to stand for themselves, we may precede them with a backslash. For instance, we see `\.` in the definition of *number*, to represent the dot, since that character is a metasymbol representing "any character," as usual in UNIX regular expressions.

In the auxiliary-function section, we see two such functions, `installID()` and `installNum()`. Like the portion of the declaration section that appears between `%{...%}`, everything in the auxiliary section is copied directly to file `lex.yy.c`, but may be used in the actions.

Finally, let us examine some of the patterns and rules in the middle section of Fig. 3.23. First, *ws*, an identifier declared in the first section, has an associated empty action. If we find whitespace, we do not return to the parser, but look for another lexeme. The second token has the simple regular expression pattern `if`. Should we see the two letters `if` on the input, and they are not followed by another letter or digit (which would cause the lexical analyzer to find a longer prefix of the input matching the pattern for **id**), then the lexical analyzer consumes these two letters from the input and returns the token name `IF`, that is, the integer for which the manifest constant `IF` stands. Keywords `then` and `else` are treated similarly.

The fifth token has the pattern defined by *id*. Note that, although keywords like `if` match this pattern as well as an earlier pattern, `Lex` chooses whichever

---

[3] If `Lex` is used along with `Yacc`, then it would be normal to define the manifest constants in the `Yacc` program and use them without definition in the `Lex` program. Since `lex.yy.c` is compiled with the `Yacc` output, the constants thus will be available to the actions in the `Lex` program.

```
%{
    /* definitions of manifest constants
    LT, LE, EQ, NE, GT, GE,
    IF, THEN, ELSE, ID, NUMBER, RELOP */
%}

/* regular definitions */
delim      [ \t\n]
ws         {delim}+
letter     [A-Za-z]
digit      [0-9]
id         {letter}({letter}|{digit})*
number     {digit}+(\.{digit}+)?(E[+-]?{digit}+)?

%%

{ws}       {/* no action and no return */}
if         {return(IF);}
then       {return(THEN);}
else       {return(ELSE);}
{id}       {yylval = (int) installID(); return(ID);}
{number}   {yylval = (int) installNum(); return(NUMBER);}
"<"        {yylval = LT; return(RELOP);}
"<="       {yylval = LE; return(RELOP);}
"="        {yylval = EQ; return(RELOP);}
"<>"       {yylval = NE; return(RELOP);}
">"        {yylval = GT; return(RELOP);}
">="       {yylval = GE; return(RELOP);}

%%

int installID() {/* function to install the lexeme, whose
                    first character is pointed to by yytext,
                    and whose length is yyleng, into the
                    symbol table and return a pointer
                    thereto */
}

int installNum() {/* similar to installID, but puts numer-
                     ical constants into a separate table */
}
```

Figure 3.23: Lex program for the tokens of Fig. 3.12

pattern is listed first in situations where the longest matching prefix matches two or more patterns. The action taken when *id* is matched is threefold:

1. Function `installID()` is called to place the lexeme found in the symbol table.

2. This function returns a pointer to the symbol table, which is placed in global variable `yylval`, where it can be used by the parser or a later component of the compiler. Note that `installID()` has available to it two variables that are set automatically by the lexical analyzer that `Lex` generates:

    (a) `yytext` is a pointer to the beginning of the lexeme, analogous to `lexemeBegin` in Fig. 3.3.

    (b) `yyleng` is the length of the lexeme found.

3. The token name ID is returned to the parser.

The action taken when a lexeme matching the pattern *number* is similar, using the auxiliary function `installNum()`.   □

### 3.5.3   Conflict Resolution in `Lex`

We have alluded to the two rules that `Lex` uses to decide on the proper lexeme to select, when several prefixes of the input match one or more patterns:

1. Always prefer a longer prefix to a shorter prefix.

2. If the longest possible prefix matches two or more patterns, prefer the pattern listed first in the `Lex` program.

**Example 3.12 :** The first rule tells us to continue reading letters and digits to find the longest prefix of these characters to group as an identifier. It also tells us to treat `<=` as a single lexeme, rather than selecting `<` as one lexeme and `=` as the next lexeme. The second rule makes keywords reserved, if we list the keywords before **id** in the program. For instance, if `then` is determined to be the longest prefix of the input that matches any pattern, and the pattern `then` precedes `{id}`, as it does in Fig. 3.23, then the token THEN is returned, rather than ID.   □

### 3.5.4   The Lookahead Operator

`Lex` automatically reads one character ahead of the last character that forms the selected lexeme, and then retracts the input so only the lexeme itself is consumed from the input. However, sometimes, we want a certain pattern to be matched to the input only when it is followed by a certain other characters. If so, we may use the slash in a pattern to indicate the end of the part of the

pattern that matches the lexeme. What follows / is additional pattern that must be matched before we can decide that the token in question was seen, but what matches this second pattern is not part of the lexeme.

**Example 3.13 :** In Fortran and some other languages, keywords are not reserved. That situation creates problems, such as a statement

```
IF(I,J) = 3
```

where IF is the name of an array, not a keyword. This statement contrasts with statements of the form

```
IF( condition ) THEN ...
```

where IF is a keyword. Fortunately, we can be sure that the keyword IF is always followed by a left parenthesis, some text — the condition — that may contain parentheses, a right parenthesis and a letter. Thus, we could write a Lex rule for the keyword IF like:

```
IF / \( .* \) {letter}
```

This rule says that the pattern the lexeme matches is just the two letters IF. The slash says that additional pattern follows but does not match the lexeme. In this pattern, the first character is the left parentheses. Since that character is a Lex metasymbol, it must be preceded by a backslash to indicate that it has its literal meaning. The dot and star match "any string without a newline." Note that the dot is a Lex metasymbol meaning "any character except newline." It is followed by a right parenthesis, again with a backslash to give that character its literal meaning. The additional pattern is followed by the symbol *letter*, which is a regular definition representing the character class of all letters.

Note that in order for this pattern to be foolproof, we must preprocess the input to delete whitespace. We have in the pattern neither provision for whitespace, nor can we deal with the possibility that the condition extends over lines, since the dot will not match a newline character.

For instance, suppose this pattern is asked to match a prefix of input:

```
IF(A<(B+C)*D)THEN...
```

the first two characters match IF, the next character matches \(, the next nine characters match .*, and the next two match \) and *letter*. Note the fact that the first right parenthesis (after C) is not followed by a letter is irrelevant; we only need to find some way of matching the input to the pattern. We conclude that the letters IF constitute the lexeme, and they are an instance of token **if**. □

### 3.5.5  Exercises for Section 3.5

**Exercise 3.5.1:** Describe how to make the following modifications to the Lex program of Fig. 3.23:

- a) Add the keyword `while`.

- b) Change the comparison operators to be the C operators of that kind.

- c) Allow the underscore (_) as an additional letter.

! d) Add a new pattern with token STRING. The pattern consists of a double-quote ("), any string of characters and a final double-quote. However, if a double-quote appears in the string, it must be escaped by preceding it with a backslash (\), and therefore a backslash in the string must be represented by two backslashes. The lexical value, which is the string without the surrounding double-quotes, and with backslashes used to escape a character removed. Strings are to be installed in a table of strings.

**Exercise 3.5.2:** Write a Lex program that copies a file, replacing each non-empty sequence of white space by a single blank.

**Exercise 3.5.3:** Write a Lex program that copies a C program, replacing each instance of the keyword `float` by `double`.

! **Exercise 3.5.4:** Write a Lex program that converts a file to "Pig latin." Specifically, assume the file is a sequence of words (groups of letters) separated by whitespace. Every time you encounter a word:

1. If the first letter is a consonant, move it to the end of the word and then add `ay`.

2. If the first letter is a vowel, just add `ay` to the end of the word.

All nonletters are copied intact to the output.

! **Exercise 3.5.5:** In SQL, keywords and identifiers are case-insensitive. Write a Lex program that recognizes the keywords SELECT, FROM, and WHERE (in any combination of capital and lower-case letters), and token ID, which for the purposes of this exercise you may take to be any sequence of letters and digits, beginning with a letter. You need not install identifiers in a symbol table, but tell how the "install" function would differ from that described for case-sensitive identifiers as in Fig. 3.23.

# 3.6    Finite Automata

We shall now discover how `Lex` turns its input program into a lexical analyzer. At the heart of the transition is the formalism known as *finite automata*. These are essentially graphs, like transition diagrams, with a few differences:

1. Finite automata are *recognizers*; they simply say "yes" or "no" about each possible input string.

2. Finite automata come in two flavors:

    (a) *Nondeterministic finite automata* (NFA) have no restrictions on the labels of their edges. A symbol can label several edges out of the same state, and $\epsilon$, the empty string, is a possible label.

    (b) *Deterministic finite automata* (DFA) have, for each state, and for each symbol of its input alphabet exactly one edge with that symbol leaving that state.

Both deterministic and nondeterministic finite automata are capable of recognizing the same languages. In fact these languages are exactly the same languages, called the *regular languages*, that regular expressions can describe.[4]

## 3.6.1    Nondeterministic Finite Automata

A *nondeterministic finite automaton* (NFA) consists of:

1. A finite set of states $S$.

2. A set of input symbols $\Sigma$, the *input alphabet*. We assume that $\epsilon$, which stands for the empty string, is never a member of $\Sigma$.

3. A *transition function* that gives, for each state, and for each symbol in $\Sigma \cup \{\epsilon\}$ a set of *next states*.

4. A state $s_0$ from $S$ that is distinguished as the *start state* (or *initial state*).

5. A set of states $F$, a subset of $S$, that is distinguished as the *accepting states* (or *final states*).

We can represent either an NFA or DFA by a *transition graph*, where the nodes are states and the labeled edges represent the transition function. There is an edge labeled $a$ from state $s$ to state $t$ if and only if $t$ is one of the next states for state $s$ and input $a$. This graph is very much like a transition diagram, except:

---

[4]There is a small lacuna: as we defined them, regular expressions cannot describe the empty language, since we would never want to use this pattern in practice. However, finite automata *can* define the empty language. In the theory, $\emptyset$ is treated as an additional regular expression for the sole purpose of defining the empty language.

a) The same symbol can label edges from one state to several different states, and

b) An edge may be labeled by $\epsilon$, the empty string, instead of, or in addition to, symbols from the input alphabet.

**Example 3.14:** The transition graph for an NFA recognizing the language of regular expression $(\mathbf{a}|\mathbf{b})^*\mathbf{abb}$ is shown in Fig. 3.24. This abstract example, describing all strings of $a$'s and $b$'s ending in the particular string $abb$, will be used throughout this section. It is similar to regular expressions that describe languages of real interest, however. For instance, an expression describing all files whose name ends in .o is $\mathbf{any}^*.\mathbf{o}$, where $\mathbf{any}$ stands for any printable character.



Figure 3.24: A nondeterministic finite automaton

Following our convention for transition diagrams, the double circle around state 3 indicates that this state is accepting. Notice that the only ways to get from the start state 0 to the accepting state is to follow some path that stays in state 0 for a while, then goes to states 1, 2, and 3 by reading $abb$ from the input. Thus, the only strings getting to the accepting state are those that end in $abb$.   □

## 3.6.2   Transition Tables

We can also represent an NFA by a *transition table*, whose rows correspond to states, and whose columns correspond to the input symbols and $\epsilon$. The entry for a given state and input is the value of the transition function applied to those arguments. If the transition function has no information about that state-input pair, we put $\emptyset$ in the table for the pair.

**Example 3.15:** The transition table for the NFA of Fig. 3.24 is shown in Fig. 3.25.   □

The transition table has the advantage that we can easily find the transitions on a given state and input. Its disadvantage is that it takes a lot of space, when the input alphabet is large, yet most states do not have any moves on most of the input symbols.

| STATE | $a$ | $b$ | $\epsilon$ |
|-------|-----|-----|-----------|
| 0 | $\{0,1\}$ | $\{0\}$ | $\emptyset$ |
| 1 | $\emptyset$ | $\{2\}$ | $\emptyset$ |
| 2 | $\emptyset$ | $\{3\}$ | $\emptyset$ |
| 3 | $\emptyset$ | $\emptyset$ | $\emptyset$ |

Figure 3.25: Transition table for the NFA of Fig. 3.24

### 3.6.3 Acceptance of Input Strings by Automata

An NFA *accepts* input string $x$ if and only if there is some path in the transition graph from the start state to one of the accepting states, such that the symbols along the path spell out $x$. Note that $\epsilon$ labels along the path are effectively ignored, since the empty string does not contribute to the string constructed along the path.

**Example 3.16 :** The string *aabb* is accepted by the NFA of Fig. 3.24. The path labeled by *aabb* from state 0 to state 3 demonstrating this fact is:

$$0 \xrightarrow{\ a\ } 0 \xrightarrow{\ a\ } 1 \xrightarrow{\ b\ } 2 \xrightarrow{\ b\ } 3$$

Note that several paths labeled by the same string may lead to different states. For instance, path

$$0 \xrightarrow{\ a\ } 0 \xrightarrow{\ a\ } 0 \xrightarrow{\ b\ } 0 \xrightarrow{\ b\ } 0$$

is another path from state 0 labeled by the string *aabb*. This path leads to state 0, which is not accepting. However, remember that an NFA accepts a string as long as *some* path labeled by that string leads from the start state to an accepting state. The existence of other paths leading to a nonaccepting state is irrelevant. □

The *language defined* (or *accepted*) by an NFA is the set of strings labeling some path from the start to an accepting state. As was mentioned, the NFA of Fig. 3.24 defines the same language as does the regular expression $(\mathbf{a}|\mathbf{b})^*\mathbf{abb}$, that is, all strings from the alphabet $\{a, b\}$ that end in *abb*. We may use $L(A)$ to stand for the language accepted by automaton $A$.

**Example 3.17 :** Figure 3.26 is an NFA accepting $L(\mathbf{aa}^*|\mathbf{bb}^*)$. String *aaa* is accepted because of the path

$$0 \xrightarrow{\ \varepsilon\ } 1 \xrightarrow{\ a\ } 2 \xrightarrow{\ a\ } 2 \xrightarrow{\ a\ } 2$$

Note that $\epsilon$'s "disappear" in a concatenation, so the label of the path is *aaa*. □

### 3.6.4 Deterministic Finite Automata

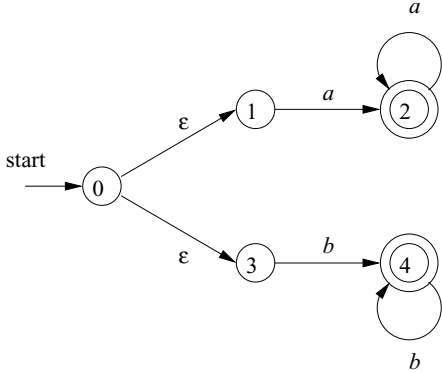A *deterministic finite automaton* (DFA) is a special case of an NFA where:

Figure 3.26: NFA accepting **aa**$^*$|**bb**$^*$

1. There are no moves on input $\epsilon$, and

2. For each state $s$ and input symbol $a$, there is exactly one edge out of $s$ labeled $a$.

If we are using a transition table to represent a DFA, then each entry is a single state. We may therefore represent this state without the curly braces that we use to form sets.

While the NFA is an abstract representation of an algorithm to recognize the strings of a certain language, the DFA is a simple, concrete algorithm for recognizing strings. It is fortunate indeed that every regular expression and every NFA can be converted to a DFA accepting the same language, because it is the DFA that we really implement or simulate when building lexical analyzers. The following algorithm shows how to apply a DFA to a string.

**Algorithm 3.18 :** Simulating a DFA.

**INPUT**: An input string $x$ terminated by an end-of-file character **eof**. A DFA $D$ with start state $s_0$, accepting states $F$, and transition function *move*.

**OUTPUT**: Answer "yes" if $D$ accepts $x$; "no" otherwise.

**METHOD**: Apply the algorithm in Fig. 3.27 to the input string $x$. The function $move(s, c)$ gives the state to which there is an edge from state $s$ on input $c$. The function *nextChar* returns the next character of the input string $x$.    □


**Example 3.19 :** In Fig. 3.28 we see the transition graph of a DFA accepting the language $(\mathbf{a}|\mathbf{b})^*\mathbf{abb}$, the same as that accepted by the NFA of Fig. 3.24. Given the input string $ababb$, this DFA enters the sequence of states $0, 1, 2, 1, 2, 3$ and returns "yes."    □

$$s = s_0;$$
$$c = nextChar();$$
**while** ( $c$ != **eof** ) {
$\quad s = move(s, c);$
$\quad c = nextChar();$
}
**if** ( $s$ is in $F$ ) **return** "yes";
**else return** "no";
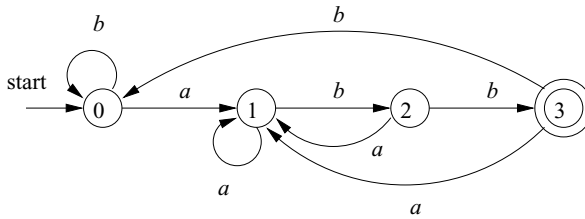
Figure 3.27: Simulating a DFA



Figure 3.28: DFA accepting **(a|b)\*abb**

## 3.6.5  Exercises for Section 3.6

! **Exercise 3.6.1 :** Figure 3.19 in the exercises of Section 3.4 computes the failure function for the KMP algorithm. Show how, given that failure function, we can construct, from a keyword $b_1b_2 \cdots b_n$ an $n + 1$-state DFA that recognizes $.^*b_1b_2 \cdots b_n$, where the dot stands for "any character." Moreover, this DFA can be constructed in $O(n)$ time.

**Exercise 3.6.2 :** Design finite automata (deterministic or nondeterministic) for each of the languages of Exercise 3.3.5.

**Exercise 3.6.3 :** For the NFA of Fig. 3.29, indicate all the paths labeled *aabb*. Does the NFA accept *aabb*?



Figure 3.29: NFA for Exercise 3.6.3

Figure 3.30: NFA for Exercise 3.6.4

**Exercise 3.6.4:** Repeat Exercise 3.6.3 for the NFA of Fig. 3.30.

**Exercise 3.6.5:** Give the transition tables for the NFA of:

a) Exercise 3.6.3.

b) Exercise 3.6.4.

c) Figure 3.26.

## 3.7    From Regular Expressions to Automata

The regular expression is the notation of choice for describing lexical analyzers and other pattern-processing software, as was reflected in Section 3.5.  However, implementation of that software requires the simulation of a DFA, as in Algorithm 3.18, or perhaps simulation of an NFA. Because an NFA often has a choice of move on an input symbol (as Fig. 3.24 does on input $a$ from state 0) or on $\epsilon$ (as Fig. 3.26 does from state 0), or even a choice of making a transition on $\epsilon$ or on a real input symbol, its simulation is less straightforward than for a DFA. Thus often it is important to convert an NFA to a DFA that accepts the same language.

In this section we shall first show how to convert NFA's to DFA's. Then, we use this technique, known as "the subset construction," to give a useful algorithm for simulating NFA's directly, in situations (other than lexical analysis) where the NFA-to-DFA conversion takes more time than the direct simulation. Next, we show how to convert regular expressions to NFA's, from which a DFA can be constructed if desired. We conclude with a discussion of the time-space tradeoffs inherent in the various methods for implementing regular expressions, and see how to choose the appropriate method for your application.

### 3.7.1    Conversion of an NFA to a DFA

The general idea behind the subset construction is that each state of the constructed DFA corresponds to a set of NFA states.  After reading input

$a_1 a_2 \cdots a_n$, the DFA is in that state which corresponds to the set of states that the NFA can reach, from its start state, following paths labeled $a_1 a_2 \cdots a_n$.

It is possible that the number of DFA states is exponential in the number of NFA states, which could lead to difficulties when we try to implement this DFA. However, part of the power of the automaton-based approach to lexical analysis is that for real languages, the NFA and DFA have approximately the same number of states, and the exponential behavior is not seen.

**Algorithm 3.20:** The *subset construction* of a DFA from an NFA.

**INPUT**: An NFA $N$.

**OUTPUT**: A DFA $D$ accepting the same language as $N$.

**METHOD**: Our algorithm constructs a transition table *Dtran* for $D$. Each state of $D$ is a set of NFA states, and we construct *Dtran* so $D$ will simulate "in parallel" all possible moves $N$ can make on a given input string. Our first problem is to deal with $\epsilon$-transitions of $N$ properly. In Fig. 3.31 we see the definitions of several functions that describe basic computations on the states of $N$ that are needed in the algorithm. Note that $s$ is a single state of $N$, while $T$ is a set of states of $N$.

| OPERATION | DESCRIPTION |
|-----------|-------------|
| $\epsilon\text{-}closure(s)$ | Set of NFA states reachable from NFA state $s$ on $\epsilon$-transitions alone. |
| $\epsilon\text{-}closure(T)$ | Set of NFA states reachable from some NFA state $s$ in set $T$ on $\epsilon$-transitions alone; $= \cup_{s \text{ in } T} \epsilon\text{-}closure(s)$. |
| $move(T, a)$ | Set of NFA states to which there is a transition on input symbol $a$ from some state $s$ in $T$. |

Figure 3.31: Operations on NFA states

We must explore those sets of states that $N$ can be in after seeing some input string. As a basis, before reading the first input symbol, $N$ can be in any of the states of $\epsilon\text{-}closure(s_0)$, where $s_0$ is its start state. For the induction, suppose that $N$ can be in set of states $T$ after reading input string $x$. If it next reads input $a$, then $N$ can immediately go to any of the states in $move(T, a)$. However, after reading $a$, it may also make several $\epsilon$-transitions; thus $N$ could be in any state of $\epsilon\text{-}closure\big(move(T, a)\big)$ after reading input $xa$. Following these ideas, the construction of the set of $D$'s states, *Dstates*, and its transition function *Dtran*, is shown in Fig. 3.32.

The start state of $D$ is $\epsilon\text{-}closure(s_0)$, and the accepting states of $D$ are all those sets of $N$'s states that include at least one accepting state of $N$. To complete our description of the subset construction, we need only to show how

      initially, $\epsilon\text{-}closure(s_0)$ is the only state in *Dstates*, and it is unmarked;
    **while** ( there is an unmarked state $T$ in *Dstates* ) {
        mark $T$;
        **for** ( each input symbol $a$ ) {
            $U = \epsilon\text{-}closure\big(move(T, a)\big)$;
            **if** ( $U$ is not in *Dstates* )
                add $U$ as an unmarked state to *Dstates*;
            $Dtran[T, a] = U$;
        }
    }

Figure 3.32: The subset construction

$\epsilon\text{-}closure(T)$ is computed for any set of NFA states $T$. This process, shown in Fig. 3.33, is a straightforward search in a graph from a set of states. In this case, imagine that only the $\epsilon$-labeled edges are available in the graph. □

      push all states of $T$ onto *stack*;
      initialize $\epsilon\text{-}closure(T)$ to $T$;
      **while** ( *stack* is not empty ) {
          pop $t$, the top element, off *stack*;
          **for** ( each state $u$ with an edge from $t$ to $u$ labeled $\epsilon$ )
             **if** ( $u$ is not in $\epsilon\text{-}closure(T)$ ) {
                add $u$ to $\epsilon\text{-}closure(T)$;
                push $u$ onto *stack*;
             }
      }

Figure 3.33: Computing $\epsilon\text{-}closure(T)$

**Example 3.21 :** Figure 3.34 shows another NFA accepting $(\mathbf{a}|\mathbf{b})^*\mathbf{abb}$; it happens to be the one we shall construct directly from this regular expression in Section 3.7. Let us apply Algorithm 3.20 to Fig. 3.34.

    The start state $A$ of the equivalent DFA is $\epsilon\text{-}closure(0)$, or $A = \{0, 1, 2, 4, 7\}$, since these are exactly the states reachable from state 0 via a path all of whose edges have label $\epsilon$. Note that a path can have zero edges, so state 0 is reachable from itself by an $\epsilon$-labeled path.

    The input alphabet is $\{a, b\}$. Thus, our first step is to mark $A$ and compute $Dtran[A, a] = \epsilon\text{-}closure\big(move(A, a)\big)$ and $Dtran[A, b] = \epsilon\text{-}closure\big(move(A, b)\big)$. Among the states 0, 1, 2, 4, and 7, only 2 and 7 have transitions on $a$, to 3 and 8, respectively. Thus, $move(A, a) = \{3, 8\}$. Also, $\epsilon\text{-}closure(\{3, 8\}) = \{1, 2, 3, 4, 6, 7, 8\}$, so we conclude
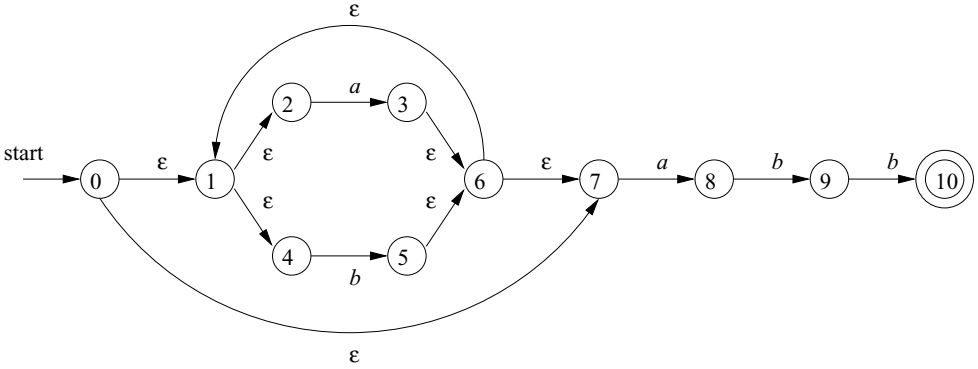
Figure 3.34: NFA $N$ for $(\mathbf{a}|\mathbf{b})^*\mathbf{abb}$

$$Dtran[A, a] = \epsilon\text{-}closure\big(move(A, a)\big) = \epsilon\text{-}closure(\{3, 8\}) = \{1, 2, 3, 4, 6, 7, 8\}$$

Let us call this set $B$, so $Dtran[A, a] = B$.

Now, we must compute $Dtran[A, b]$. Among the states in $A$, only 4 has a transition on $b$, and it goes to 5. Thus,

$$Dtran[A, b] = \epsilon\text{-}closure(\{5\}) = \{1, 2, 4, 5, 6, 7\}$$

Let us call the above set $C$, so $Dtran[A, b] = C$.

| NFA State | DFA State | $a$ | $b$ |
|:---:|:---:|:---:|:---:|
| $\{0, 1, 2, 4, 7\}$ | $A$ | $B$ | $C$ |
| $\{1, 2, 3, 4, 6, 7, 8\}$ | $B$ | $B$ | $D$ |
| $\{1, 2, 4, 5, 6, 7\}$ | $C$ | $B$ | $C$ |
| $\{1, 2, 4, 5, 6, 7, 9\}$ | $D$ | $B$ | $E$ |
| $\{1, 2, 4, 5, 6, 7, 10\}$ | $E$ | $B$ | $C$ |

Figure 3.35: Transition table $Dtran$ for DFA $D$

If we continue this process with the unmarked sets $B$ and $C$, we eventually reach a point where all the states of the DFA are marked. This conclusion is guaranteed, since there are "only" $2^{11}$ different subsets of a set of eleven NFA states. The five different DFA states we actually construct, their corresponding sets of NFA states, and the transition table for the DFA $D$ are shown in Fig. 3.35, and the transition graph for $D$ is in Fig. 3.36. State $A$ is the start state, and state $E$, which contains state 10 of the NFA, is the only accepting state.

Note that $D$ has one more state than the DFA of Fig. 3.28 for the same language. States $A$ and $C$ have the same move function, and so can be merged. We discuss the matter of minimizing the number of states of a DFA in Section 3.9.6.
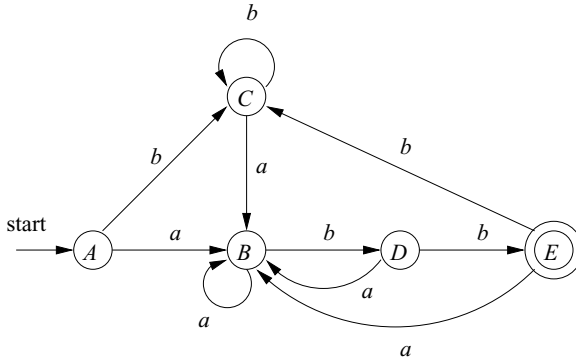□

Figure 3.36: Result of applying the subset construction to Fig. 3.34

## 3.7.2   Simulation of an NFA

A strategy that has been used in a number of text-editing programs is to construct an NFA from a regular expression and then simulate the NFA using something like an on-the-fly subset construction. The simulation is outlined below.

**Algorithm 3.22 :** Simulating an NFA.

**INPUT**: An input string $x$ terminated by an end-of-file character **eof**. An NFA $N$ with start state $s_0$, accepting states $F$, and transition function *move*.

**OUTPUT**: Answer "yes" if $N$ accepts $x$; "no" otherwise.

**METHOD**: The algorithm keeps a set of current states $S$, those that are reached from $s_0$ following a path labeled by the inputs read so far. If $c$ is the next input character, read by the function *nextChar*(), then we first compute *move*$(S, c)$ and then close that set using $\epsilon$-*closure*(). The algorithm is sketched in Fig. 3.37. □

```
1)    S = ε-closure(s₀);
2)    c = nextChar();
3)    while ( c != eof ) {
4)          S = ε-closure(move(S, c));
5)          c = nextChar();
6)    }
7)    if ( S ∩ F != ∅ ) return "yes";
8)    else return "no";
```

Figure 3.37: Simulating an NFA

### 3.7.3 Efficiency of NFA Simulation

If carefully implemented, Algorithm 3.22 can be quite efficient. As the ideas involved are useful in a number of similar algorithms involving search of graphs, we shall look at this implementation in additional detail. The data structures we need are:

1. Two stacks, each of which holds a set of NFA states. One of these stacks, *oldStates*, holds the "current" set of states, i.e., the value of $S$ on the right side of line (4) in Fig. 3.37. The second, *newStates*, holds the "next" set of states — $S$ on the left side of line (4). Unseen is a step where, as we go around the loop of lines (3) through (6), *newStates* is transferred to *oldStates*.

2. A boolean array *alreadyOn*, indexed by the NFA states, to indicate which states are in *newStates*. While the array and stack hold the same information, it is much faster to interrogate *alreadyOn*[$s$] than to search for state $s$ on the stack *newStates*. It is for this efficiency that we maintain both representations.

3. A two-dimensional array *move*[$s, a$] holding the transition table of the NFA. The entries in this table, which are sets of states, are represented by linked lists.

To implement line (1) of Fig. 3.37, we need to set each entry in array *alreadyOn* to FALSE, then for each state $s$ in $\epsilon$-*closure*($s_0$), push $s$ onto *oldStates* and set *alreadyOn*[$s$] to TRUE. This operation on state $s$, and the implementation of line (4) as well, are facilitated by a function we shall call *addState*($s$). This function pushes state $s$ onto *newStates*, sets *alreadyOn*[$s$] to TRUE, and calls itself recursively on the states in *move*[$s, \epsilon$] in order to further the computation of $\epsilon$-*closure*($s$). However, to avoid duplicating work, we must be careful never to call *addState* on a state that is already on the stack *newStates*. Figure 3.38 sketches this function.

```
9)    addState(s) {
10)        push s onto newStates;
11)        alreadyOn[s] = TRUE;
12)        for ( t on move[s, ε] )
13)            if ( !alreadyOn[t] )
14)                addState(t);
15)    }
```

Figure 3.38: Adding a new state $s$, which is known not to be on *newStates*

We implement line (4) of Fig. 3.37 by looking at each state $s$ on *oldStates*. We first find the set of states *move*[$s, c$], where $c$ is the next input, and for each

of those states that is not already on *newStates*, we apply *addState* to it. Note that *addState* has the effect of computing the $\epsilon$-*closure* and adding all those states to *newStates* as well, if they were not already on. This sequence of steps is summarized in Fig. 3.39.

```
16)     for ( s on oldStates ) {
17)         for ( t on move[s, c] )
18)             if ( !alreadyOn[t] )
19)                 addState(t);
20)         pop s from oldStates;
21)     }

22)     for ( s on newStates ) {
23)         pop s from newStates;
24)         push s onto oldStates;
25)         alreadyOn[s] = FALSE;
26)     }
```

Figure 3.39: Implementation of step (4) of Fig. 3.37

Now, suppose that the NFA $N$ has $n$ states and $m$ transitions; i.e., $m$ is the sum over all states of the number of symbols (or $\epsilon$) on which the state has a transition out. Not counting the call to *addState* at line (19) of Fig. 3.39, the time spent in the loop of lines (16) through (21) is $O(n)$. That is, we can go around the loop at most $n$ times, and each step of the loop requires constant work, except for the time spent in *addState*. The same is true of the loop of lines (22) through (26).

During one execution of Fig. 3.39, i.e., of step (4) of Fig. 3.37, it is only possible to call *addState* on a given state once. The reason is that whenever we call *addState*(*s*), we set *alreadyOn*[*s*] to TRUE at line (11) of Fig. 3.38. Once *alreadyOn*[*s*] is TRUE, the tests at line (13) of Fig. 3.38 and line (18) of Fig. 3.39 prevent another call.

The time spent in one call to *addState*, exclusive of the time spent in recursive calls at line (14), is $O(1)$ for lines (10) and (11). For lines (12) and (13), the time depends on how many $\epsilon$-transitions there are out of state $s$. We do not know this number for a given state, but we know that there are at most $m$ transitions in total, out of all states. As a result, the aggregate time spent in lines (12) and (13) over all calls to *addState* during one execution of the code of Fig. 3.39 is $O(m)$. The aggregate for the rest of the steps of *addState* is $O(n)$, since it is a constant per call, and there are at most $n$ calls.

We conclude that, implemented properly, the time to execute line (4) of Fig. 3.37 is $O(n + m)$. The rest of the while-loop of lines (3) through (6) takes $O(1)$ time per iteration. If the input $x$ is of length $k$, then the total work in that loop is $O(k(n + m))$. Line (1) of Fig. 3.37 can be executed in $O(n + m)$ time, since it is essentially the steps of Fig. 3.39 with *oldStates* containing only

---

### Big-Oh Notation

An expression like $O(n)$ is a shorthand for "at most some constant times $n$." Technically, we say a function $f(n)$, perhaps the running time of some step of an algorithm, is $O(g(n))$ if there are constants $c$ and $n_0$, such that whenever $n \geq n_0$, it is true that $f(n) \leq cg(n)$. A useful idiom is "$O(1)$," which means "some constant." The use of this *big-oh notation* enables us to avoid getting too far into the details of what we count as a unit of execution time, yet lets us express the rate at which the running time of an algorithm grows.

---

the state $s_0$. Lines (2), (7), and (8) each take $O(1)$ time. Thus, the running time of Algorithm 3.22, properly implemented, is $O(k(n + m))$. That is, the time taken is proportional to the length of the input times the size (nodes plus edges) of the transition graph.

### 3.7.4 Construction of an NFA from a Regular Expression

We now give an algorithm for converting any regular expression to an NFA that defines the same language. The algorithm is syntax-directed, in the sense that it works recursively up the parse tree for the regular expression. For each subexpression the algorithm constructs an NFA with a single accepting state.

**Algorithm 3.23 :** The McNaughton-Yamada-Thompson algorithm to convert a regular expression to an NFA.

**INPUT**: A regular expression $r$ over alphabet $\Sigma$.

**OUTPUT**: An NFA $N$ accepting $L(r)$.

**METHOD**: Begin by parsing $r$ into its constituent subexpressions. The rules for constructing an NFA consist of basis rules for handling subexpressions with no operators, and inductive rules for constructing larger NFA's from the NFA's for the immediate subexpressions of a given expression.

**BASIS**: For expression $\epsilon$ construct the NFA



Here, $i$ is a new state, the start state of this NFA, and $f$ is another new state, the accepting state for the NFA.

For any subexpression $a$ in $\Sigma$, construct the NFA

where again $i$ and $f$ are new states, the start and accepting states, respectively. Note that in both of the basis constructions, we construct a distinct NFA, with new states, for every occurrence of $\epsilon$ or some $a$ as a subexpression of $r$.

**INDUCTION**: Suppose $N(s)$ and $N(t)$ are NFA's for regular expressions $s$ and $t$, respectively.

a) Suppose $r = s|t$. Then $N(r)$, the NFA for $r$, is constructed as in Fig. 3.40. Here, $i$ and $f$ are new states, the start and accepting states of $N(r)$, respectively. There are $\epsilon$-transitions from $i$ to the start states of $N(s)$ and $N(t)$, and each of their accepting states have $\epsilon$-transitions to the accepting state $f$. Note that the accepting states of $N(s)$ and $N(t)$ are not accepting in $N(r)$. Since any path from $i$ to $f$ must pass through either $N(s)$ or $N(t)$ exclusively, and since the label of that path is not changed by the $\epsilon$'s leaving $i$ or entering $f$, we conclude that $N(r)$ accepts $L(s) \cup L(t)$, which is the same as $L(r)$. That is, Fig. 3.40 is a correct construction for the union operator.



Figure 3.40: NFA for the union of two regular expressions

b) Suppose $r = st$. Then construct $N(r)$ as in Fig. 3.41. The start state of $N(s)$ becomes the start state of $N(r)$, and the accepting state of $N(t)$ is the only accepting state of $N(r)$. The accepting state of $N(s)$ and the start state of $N(t)$ are merged into a single state, with all the transitions in or out of either state. A path from $i$ to $f$ in Fig. 3.41 must go first through $N(s)$, and therefore its label will begin with some string in $L(s)$. The path then continues through $N(t)$, so the path's label finishes with a string in $L(t)$. As we shall soon argue, accepting states never have edges out and start states never have edges in, so it is not possible for a path to re-enter $N(s)$ after leaving it. Thus, $N(r)$ accepts exactly $L(s)L(t)$, and is a correct NFA for $r = st$.
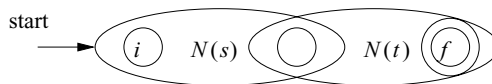


Figure 3.41: NFA for the concatenation of two regular expressions

c) Suppose $r = s^*$. Then for $r$ we construct the NFA $N(r)$ shown in Fig. 3.42. Here, $i$ and $f$ are new states, the start state and lone accepting state of $N(r)$. To get from $i$ to $f$, we can either follow the introduced path labeled $\epsilon$, which takes care of the one string in $L(s)^0$, or we can go to the start state of $N(s)$, through that NFA, then from its accepting state back to its start state zero or more times. These options allow $N(r)$ to accept all the strings in $L(s)^1$, $L(s)^2$, and so on, so the entire set of strings accepted by $N(r)$ is $L(s^*)$.



Figure 3.42: NFA for the closure of a regular expression

d) Finally, suppose $r = (s)$. Then $L(r) = L(s)$, and we can use the NFA $N(s)$ as $N(r)$.

□

The method description in Algorithm 3.23 contains hints as to why the inductive construction works as it should. We shall not give a formal correctness proof, but we shall list several properties of the constructed NFA's, in addition to the all-important fact that $N(r)$ accepts language $L(r)$. These properties are interesting in their own right, and helpful in making a formal proof.

1. $N(r)$ has at most twice as many states as there are operators and operands in $r$. This bound follows from the fact that each step of the algorithm creates at most two new states.

2. $N(r)$ has one start state and one accepting state. The accepting state has no outgoing transitions, and the start state has no incoming transitions.

3. Each state of $N(r)$ other than the accepting state has either one outgoing transition on a symbol in $\Sigma$ or two outgoing transitions, both on $\epsilon$.

**Example 3.24 :** Let us use Algorithm 3.23 to construct an NFA for $r = (\mathbf{a}|\mathbf{b})^*\mathbf{abb}$. Figure 3.43 shows a parse tree for $r$ that is analogous to the parse trees constructed for arithmetic expressions in Section 2.2.3. For subexpression $r_1$, the first $\mathbf{a}$, we construct the NFA:
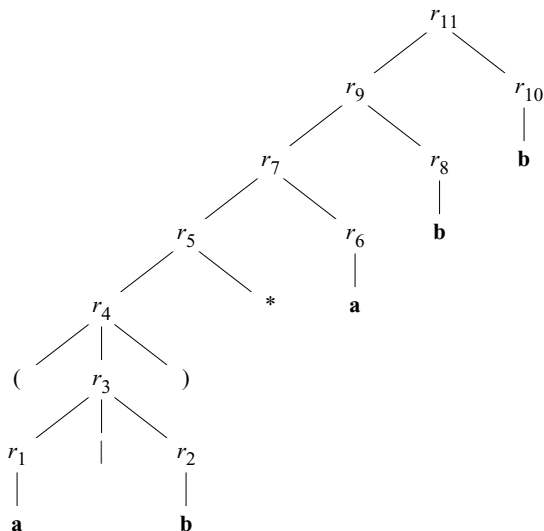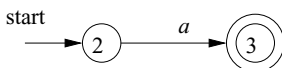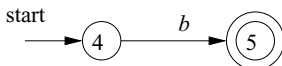
Figure 3.43: Parse tree for $(\mathbf{a}|\mathbf{b})^*\mathbf{abb}$



State numbers have been chosen for consistency with what follows. For $r_2$ we construct:



We can now combine $N(r_1)$ and $N(r_2)$, using the construction of Fig. 3.40 to obtain the NFA for $r_3 = r_1|r_2$; this NFA is shown in Fig. 3.44.
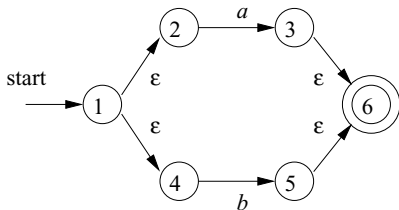


Figure 3.44: NFA for $r_3$

The NFA for $r_4 = (r_3)$ is the same as that for $r_3$. The NFA for $r_5 = (r_3)^*$ is then as shown in Fig. 3.45. We have used the construction in Fig. 3.42 to build this NFA from the NFA in Fig. 3.44.

Now, consider subexpression $r_6$, which is another **a**. We use the basis construction for $a$ again, but we must use new states. It is not permissible to reuse
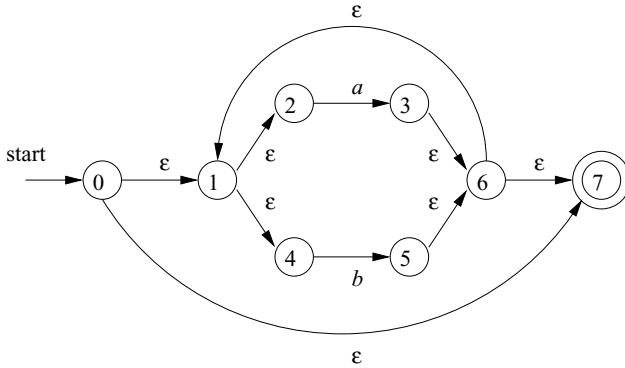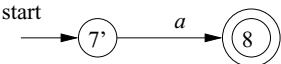
Figure 3.45: NFA for $r_5$

the NFA we constructed for $r_1$, even though $r_1$ and $r_6$ are the same expression. The NFA for $r_6$ is:



To obtain the NFA for $r_7 = r_5 r_6$, we apply the construction of Fig. 3.41. We merge states 7 and 7′, yielding the NFA of Fig. 3.46. Continuing in this fashion with new NFA's for the two subexpressions **b** called $r_8$ and $r_{10}$, we eventually construct the NFA for $(\mathbf{a}|\mathbf{b})^*\mathbf{abb}$ that we first met in Fig. 3.34. □
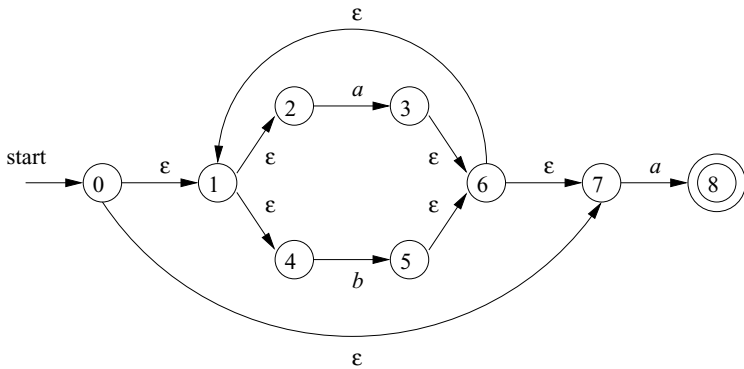


Figure 3.46: NFA for $r_7$

## 3.7.5 Efficiency of String-Processing Algorithms

We observed that Algorithm 3.18 processes a string $x$ in time $O(|x|)$, while in Section 3.7.3 we concluded that we could simulate an NFA in time proportional to the product of $|x|$ and the size of the NFA's transition graph. Obviously, it

is faster to have a DFA to simulate than an NFA, so we might wonder whether it ever makes sense to simulate an NFA.

One issue that may favor an NFA is that the subset construction can, in the worst case, exponentiate the number of states. While in principle, the number of DFA states does not influence the running time of Algorithm 3.18, should the number of states become so large that the transition table does not fit in main memory, then the true running time would have to include disk I/O and therefore rise noticeably.

**Example 3.25 :** Consider the family of languages described by regular expressions of the form $L_n = (\mathbf{a}|\mathbf{b})^*\mathbf{a}(\mathbf{a}|\mathbf{b})^{n-1}$, that is, each language $L_n$ consists of strings of $a$'s and $b$'s such that the $n$th character to the left of the right end holds $a$. An $n + 1$-state NFA is easy to construct. It stays in its initial state under any input, but also has the option, on input $a$, of going to state 1. From state 1, it goes to state 2 on any input, and so on, until in state $n$ it accepts. Figure 3.47 suggests this NFA.
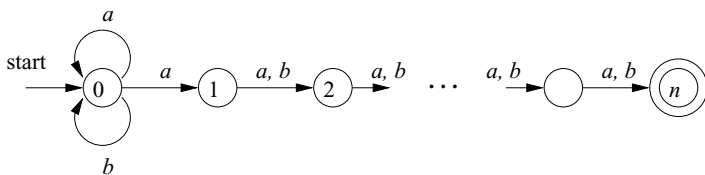


Figure 3.47: An NFA that has many fewer states than the smallest equivalent DFA

However, any DFA for the language $L_n$ must have at least $2^n$ states. We shall not prove this fact, but the idea is that if two strings of length $n$ can get the DFA to the same state, then we can exploit the last position where the strings differ (and therefore one must have $a$, the other $b$) to continue the strings identically, until they are the same in the last $n-1$ positions. The DFA will then be in a state where it must both accept and not accept. Fortunately, as we mentioned, it is rare for lexical analysis to involve patterns of this type, and we do not expect to encounter DFA's with outlandish numbers of states in practice.   □

However, lexical-analyzer generators and other string-processing systems often start with a regular expression. We are faced with a choice of converting the regular expression to an NFA or DFA. The additional cost of going to a DFA is thus the cost of executing Algorithm 3.20 on the NFA (one could go directly from a regular expression to a DFA, but the work is essentially the same). If the string-processor is one that will be executed many times, as is the case for lexical analysis, then any cost of converting to a DFA is worthwhile. However, in other string-processing applications, such as `grep`, where the user specifies one regular expression and one or several files to be searched for the pattern

of that expression, it may be more efficient to skip the step of constructing a DFA, and simulate the NFA directly.

Let us consider the cost of converting a regular expression $r$ to an NFA by Algorithm 3.23. A key step is constructing the parse tree for $r$. In Chapter 4 we shall see several methods that are capable of constructing this parse tree in linear time, that is, in time $O(|r|)$, where $|r|$ stands for the *size* of $r$ — the sum of the number of operators and operands in $r$. It is also easy to check that each of the basis and inductive constructions of Algorithm 3.23 takes constant time, so the entire time spent by the conversion to an NFA is $O(|r|)$.

Moreover, as we observed in Section 3.7.4, the NFA we construct has at most $2|r|$ states and at most $4|r|$ transitions. That is, in terms of the analysis in Section 3.7.3, we have $n \leq 2|r|$ and $m \leq 4|r|$. Thus, simulating this NFA on an input string $x$ takes time $O(|r| \times |x|)$. This time dominates the time taken by the NFA construction, which is $O(|r|)$, and therefore, we conclude that it is possible to take a regular expression $r$ and string $x$, and tell whether $x$ is in $L(r)$ in time $O(|r| \times |x|)$.

The time taken by the subset construction is highly dependent on the number of states the resulting DFA has. To begin, notice that in the subset construction of Fig. 3.32, the key step, the construction of a set of states $U$ from a set of states $T$ and an input symbol $a$, is very much like the construction of a new set of states from the old set of states in the NFA simulation of Algorithm 3.22. We already concluded that, properly implemented, this step takes time at most proportional to the number of states and transitions of the NFA.

Suppose we start with a regular expression $r$ and convert it to an NFA. This NFA has at most $2|r|$ states and at most $4|r|$ transitions. Moreover, there are at most $|r|$ input symbols. Thus, for every DFA state constructed, we must construct at most $|r|$ new states, and each one takes at most $O(|r|)$ time. The time to construct a DFA of $s$ states is thus $O(|r|^2 s)$.

In the common case where $s$ is about $|r|$, the subset construction takes time $O(|r|^3)$. However, in the worst case, as in Example 3.25, this time is $O(|r|^2 2^{|r|})$. Figure 3.48 summarizes the options when one is given a regular expression $r$ and wants to produce a recognizer that will tell whether one or more strings $x$ are in $L(r)$.

| AUTOMATON | INITIAL | PER STRING |
|:---:|:---:|:---:|
| NFA | $O(|r|)$ | $O(|r| \times |x|)$ |
| DFA typical case | $O(|r|^3)$ | $O(|x|)$ |
| DFA worst case | $O(|r|^2 2^{|r|})$ | $O(|x|)$ |

Figure 3.48: Initial cost and per-string-cost of various methods of recognizing the language of a regular expression

If the per-string cost dominates, as it does when we build a lexical analyzer,

we clearly prefer the DFA. However, in commands like `grep`, where we run the automaton on only one string, we generally prefer the NFA. It is not until $|x|$ approaches $|r|^3$ that we would even think about converting to a DFA.

There is, however, a mixed strategy that is about as good as the better of the NFA and the DFA strategy for each expression $r$ and string $x$. Start off simulating the NFA, but remember the sets of NFA states (i.e., the DFA states) and their transitions, as we compute them. Before processing the current set of NFA states and the current input symbol, check to see whether we have already computed this transition, and use the information if so.

### 3.7.6  Exercises for Section 3.7

**Exercise 3.7.1 :** Convert to DFA's the NFA's of:

  a) Fig. 3.26.

  b) Fig. 3.29.

  c) Fig. 3.30.

**Exercise 3.7.2 :** use Algorithm 3.22 to simulate the NFA's:

  a) Fig. 3.29.

  b) Fig. 3.30.

on input *aabb*.

**Exercise 3.7.3 :** Convert the following regular expressions to deterministic finite automata, using algorithms 3.23 and 3.20:

  a) $(\mathbf{a}|\mathbf{b})^*$.

  b) $(\mathbf{a}^*|\mathbf{b}^*)^*$.

  c) $\big((\epsilon|\mathbf{a})\mathbf{b}^*\big)^*$.

  d) $(\mathbf{a}|\mathbf{b})^*\mathbf{abb}(\mathbf{a}|\mathbf{b})^*$.

## 3.8   Design of a Lexical-Analyzer Generator

In this section we shall apply the techniques presented in Section 3.7 to see how a lexical-analyzer generator such as `Lex` is architected.  We discuss two approaches, based on NFA's and DFA's; the latter is essentially the implementation of `Lex`.

### 3.8.1 The Structure of the Generated Analyzer

Figure 3.49 overviews the architecture of a lexical analyzer generated by Lex. The program that serves as the lexical analyzer includes a fixed program that simulates an automaton; at this point we leave open whether that automaton is deterministic or nondeterministic. The rest of the lexical analyzer consists of components that are created from the Lex program by Lex itself.
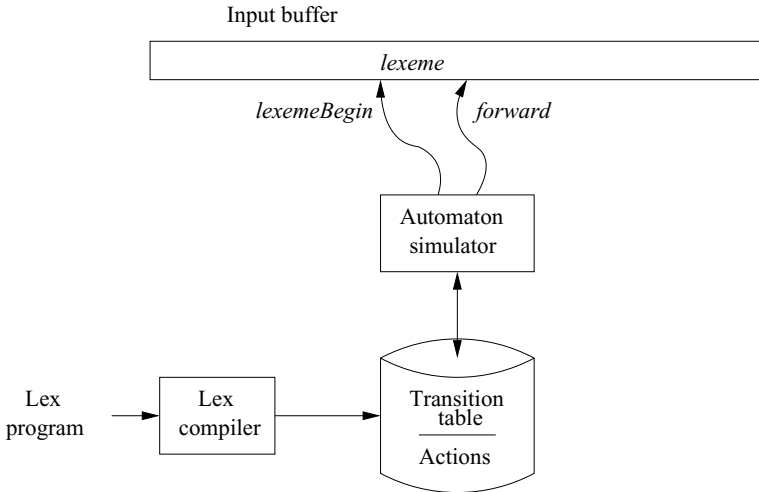


Figure 3.49: A Lex program is turned into a transition table and actions, which are used by a finite-automaton simulator

These components are:

1. A transition table for the automaton.

2. Those functions that are passed directly through Lex to the output (see Section 3.5.2).

3. The actions from the input program, which appear as fragments of code to be invoked at the appropriate time by the automaton simulator.

To construct the automaton, we begin by taking each regular-expression pattern in the Lex program and converting it, using Algorithm 3.23, to an NFA. We need a single automaton that will recognize lexemes matching any of the patterns in the program, so we combine all the NFA's into one by introducing a new start state with $\epsilon$-transitions to each of the start states of the NFA's $N_i$ for pattern $p_i$. This construction is shown in Fig. 3.50.

**Example 3.26 :** We shall illustrate the ideas of this section with the following simple, abstract example:
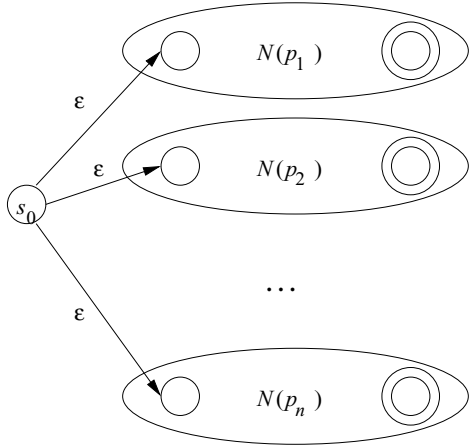
Figure 3.50: An NFA constructed from a Lex program

| | |
|---|---|
| **a** | { action $A_1$ for pattern $p_1$ } |
| **abb** | { action $A_2$ for pattern $p_2$ } |
| **a\*b$^+$** | { action $A_3$ for pattern $p_3$ } |

Note that these three patterns present some conflicts of the type discussed in Section 3.5.3. In particular, string *abb* matches both the second and third patterns, but we shall consider it a lexeme for pattern $p_2$, since that pattern is listed first in the above Lex program. Then, input strings such as *aabbb* $\cdots$ have many prefixes that match the third pattern. The Lex rule is to take the longest, so we continue reading *b*'s, until another *a* is met, whereupon we report the lexeme to be the initial *a*'s followed by as many *b*'s as there are.

Figure 3.51 shows three NFA's that recognize the three patterns. The third is a simplification of what would come out of Algorithm 3.23. Then, Fig. 3.52 shows these three NFA's combined into a single NFA by the addition of start state 0 and three $\epsilon$-transitions.  $\square$

### 3.8.2  Pattern Matching Based on NFA's

If the lexical analyzer simulates an NFA such as that of Fig. 3.52, then it must read input beginning at the point on its input which we have referred to as *lexemeBegin*. As it moves the pointer called *forward* ahead in the input, it calculates the set of states it is in at each point, following Algorithm 3.22.

Eventually, the NFA simulation reaches a point on the input where there are no next states. At that point, there is no hope that any longer prefix of the input would ever get the NFA to an accepting state; rather, the set of states will always be empty. Thus, we are ready to decide on the longest prefix that is a lexeme matching some pattern.
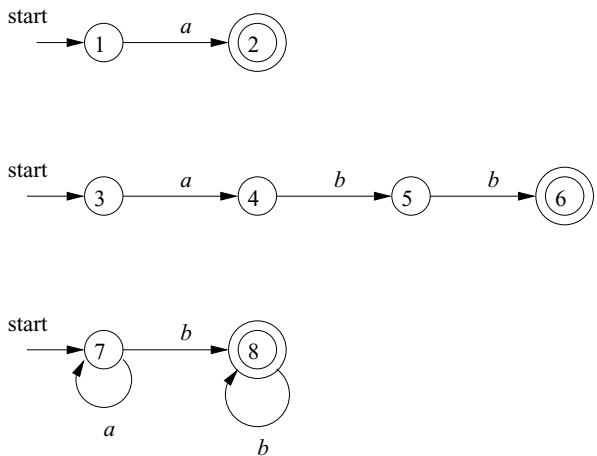
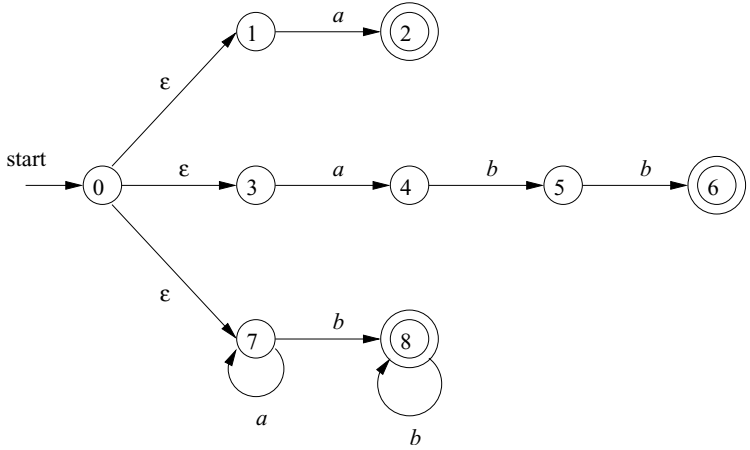Figure 3.51: NFA's for **a**, **abb**, and **a**\***b**$^+$
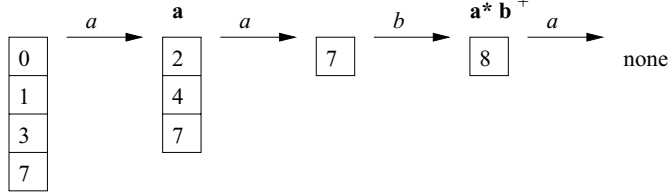


Figure 3.52: Combined NFA



Figure 3.53: Sequence of sets of states entered when processing input *aaba*

We look backwards in the sequence of sets of states, until we find a set that includes one or more accepting states. If there are several accepting states in that set, pick the one associated with the earliest pattern $p_i$ in the list from the Lex program. Move the *forward* pointer back to the end of the lexeme, and perform the action $A_i$ associated with pattern $p_i$.

**Example 3.27 :** Suppose we have the patterns of Example 3.26 and the input begins *aaba*. Figure 3.53 shows the sets of states of the NFA of Fig. 3.52 that we enter, starting with $\epsilon$-*closure* of the initial state 0, which is $\{0, 1, 3, 7\}$, and proceeding from there. After reading the fourth input symbol, we are in an empty set of states, since in Fig. 3.52, there are no transitions out of state 8 on input $a$.

Thus, we need to back up, looking for a set of states that includes an accepting state. Notice that, as indicated in Fig. 3.53, after reading $a$ we are in a set that includes state 2 and therefore indicates that the pattern **a** has been matched. However, after reading *aab*, we are in state 8, which indicates that $\mathbf{a^*b^+}$ has been matched; prefix *aab* is the longest prefix that gets us to an accepting state. We therefore select *aab* as the lexeme, and execute action $A_3$, which should include a return to the parser indicating that the token whose pattern is $p_3 = \mathbf{a^*b^+}$ has been found.   □

### 3.8.3  DFA's for Lexical Analyzers

Another architecture, resembling the output of Lex, is to convert the NFA for all the patterns into an equivalent DFA, using the subset construction of Algorithm 3.20. Within each DFA state, if there are one or more accepting NFA states, determine the first pattern whose accepting state is represented, and make that pattern the output of the DFA state.

**Example 3.28 :** Figure 3.54 shows a transition diagram based on the DFA that is constructed by the subset construction from the NFA in Fig. 3.52. The accepting states are labeled by the pattern that is identified by that state. For instance, the state $\{6, 8\}$ has two accepting states, corresponding to patterns **abb** and $\mathbf{a^*b^+}$. Since the former is listed first, that is the pattern associated with state $\{6, 8\}$.   □

We use the DFA in a lexical analyzer much as we did the NFA. We simulate the DFA until at some point there is no next state (or strictly speaking, the next state is $\emptyset$, the *dead state* corresponding to the empty set of NFA states). At that point, we back up through the sequence of states we entered and, as soon as we meet an accepting DFA state, we perform the action associated with the pattern for that state.

**Example 3.29 :** Suppose the DFA of Fig. 3.54 is given input *abba*. The sequence of states entered is $0137, 247, 58, 68$, and at the final $a$ there is no transition out of state 68. Thus, we consider the sequence from the end, and in this case, 68 itself is an accepting state that reports pattern $p_2 = \mathbf{abb}$.   □
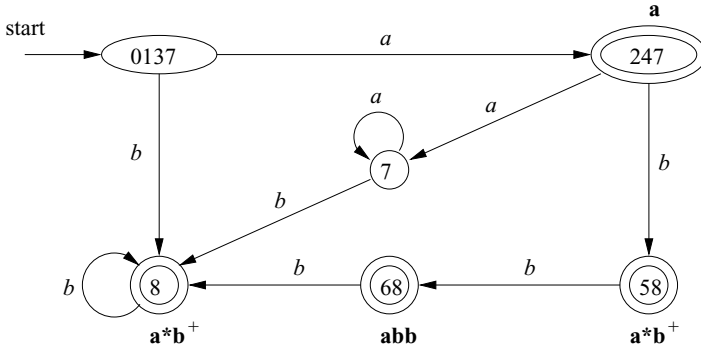
Figure 3.54: Transition graph for DFA handling the patterns $\mathbf{a}$, $\mathbf{abb}$, and $\mathbf{a^*b^+}$

### 3.8.4 Implementing the Lookahead Operator

Recall from Section 3.5.4 that the Lex lookahead operator / in a Lex pattern $r_1/r_2$ is sometimes necessary, because the pattern $r_1$ for a particular token may need to describe some trailing context $r_2$ in order to correctly identify the actual lexeme. When converting the pattern $r_1/r_2$ to an NFA, we treat the / as if it were $\epsilon$, so we do not actually look for a / on the input. However, if the NFA recognizes a prefix $xy$ of the input buffer as matching this regular expression, the end of the lexeme is not where the NFA entered its accepting state. Rather the end occurs when the NFA enters a state $s$ such that

1. $s$ has an $\epsilon$-transition on the (imaginary) /,

2. There is a path from the start state of the NFA to state $s$ that spells out $x$.

3. There is a path from state $s$ to the accepting state that spells out $y$.

4. $x$ is as long as possible for any $xy$ satisfying conditions 1-3.

If there is only one $\epsilon$-transition state on the imaginary / in the NFA, then the end of the lexeme occurs when this state is entered for the last time as the following example illustrates. If the NFA has more than one $\epsilon$-transition state on the imaginary /, then the general problem of finding the correct state $s$ is much more difficult.

**Example 3.30:** An NFA for the pattern for the Fortran IF with lookahead, from Example 3.13, is shown in Fig. 3.55. Notice that the $\epsilon$-transition from state 2 to state 3 represents the lookahead operator. State 6 indicates the presence of the keyword IF. However, we find the lexeme IF by scanning backwards to the last occurrence of state 2, whenever state 6 is entered. □

---

### Dead States in DFA's

Technically, the automaton in Fig. 3.54 is not quite a DFA. The reason
is that a DFA has a transition from every state on every input symbol in
its input alphabet. Here, we have omitted transitions to the dead state
$\emptyset$, and we have therefore omitted the transitions from the dead state to
itself on every input. Previous NFA-to-DFA examples did not have a way
to get from the start state to $\emptyset$, but the NFA of Fig. 3.52 does.

However, when we construct a DFA for use in a lexical analyzer, it
is important that we treat the dead state differently, since we must know
when there is no longer any possibility of recognizing a longer lexeme.
Thus, we suggest always omitting transitions to the dead state and elimi-
nating the dead state itself. In fact, the problem is harder than it appears,
since an NFA-to-DFA construction may yield several states that cannot
reach any accepting state, and we must know when any of these states
have been reached. Section 3.9.6 discusses how to combine all these states
into one dead state, so their identification becomes easy. It is also inter-
esting to note that if we construct a DFA from a regular expression using
Algorithms 3.20 and 3.23, then the DFA will not have any states besides
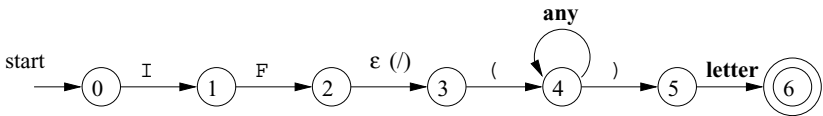$\emptyset$ that cannot lead to an accepting state.

---



Figure 3.55: NFA recognizing the keyword IF

## 3.8.5 Exercises for Section 3.8

**Exercise 3.8.1:** Suppose we have two tokens: (1) the keyword `if`, and (2) id-
entifiers, which are strings of letters other than `if`. Show:

  a) The NFA for these tokens, and

  b) The DFA for these tokens.

**Exercise 3.8.2:** Repeat Exercise 3.8.1 for tokens consisting of (1) the keyword
`while`, (2) the keyword `when`, and (3) identifiers consisting of strings of letters
and digits, beginning with a letter.

! **Exercise 3.8.3:** Suppose we were to revise the definition of a DFA to allow
zero or one transition out of each state on each input symbol (rather than
exactly one such transition, as in the standard DFA definition). Some regular

expressions would then have smaller "DFA's" than they do under the standard definition of a DFA. Give an example of one such regular expression.

!! **Exercise 3.8.4:** Design an algorithm to recognize Lex-lookahead patterns of the form $r_1/r_2$, where $r_1$ and $r_2$ are regular expressions. Show how your algorithm works on the following inputs:

a) **(abcd|abc)/d**

b) **(a|ab)/ba**

c) **aa∗/a∗**

# 3.9 Optimization of DFA-Based Pattern Matchers

In this section we present three algorithms that have been used to implement and optimize pattern matchers constructed from regular expressions.

1. The first algorithm is useful in a Lex compiler, because it constructs a DFA directly from a regular expression, without constructing an intermediate NFA. The resulting DFA also may have fewer states than the DFA constructed via an NFA.

2. The second algorithm minimizes the number of states of any DFA, by combining states that have the same future behavior. The algorithm itself is quite efficient, running in time $O(n \log n)$, where $n$ is the number of states of the DFA.

3. The third algorithm produces more compact representations of transition tables than the standard, two-dimensional table.

## 3.9.1 Important States of an NFA

To begin our discussion of how to go directly from a regular expression to a DFA, we must first dissect the NFA construction of Algorithm 3.23 and consider the roles played by various states. We call a state of an NFA *important* if it has a non-$\epsilon$ out-transition. Notice that the subset construction (Algorithm 3.20) uses only the important states in a set $T$ when it computes $\epsilon\text{-}closure\big(move(T, a)\big)$, the set of states reachable from $T$ on input $a$. That is, the set of states $move(s, a)$ is nonempty only if state $s$ is important. During the subset construction, two sets of NFA states can be identified (treated as if they were the same set) if they:

1. Have the same important states, and

2. Either both have accepting states or neither does.

When the NFA is constructed from a regular expression by Algorithm 3.23, we can say more about the important states. The only important states are those introduced as initial states in the basis part for a particular symbol position in the regular expression. That is, each important state corresponds to a particular operand in the regular expression.

The constructed NFA has only one accepting state, but this state, having no out-transitions, is not an important state. By concatenating a unique right endmarker # to a regular expression $r$, we give the accepting state for $r$ a transition on #, making it an important state of the NFA for $(r)\#$. In other words, by using the *augmented* regular expression $(r)\#$, we can forget about accepting states as the subset construction proceeds; when the construction is complete, any state with a transition on # must be an accepting state.

The important states of the NFA correspond directly to the positions in the regular expression that hold symbols of the alphabet. It is useful, as we shall see, to present the regular expression by its *syntax tree*, where the leaves correspond to operands and the interior nodes correspond to operators. An interior node is called a *cat-node*, *or-node*, or *star-node* if it is labeled by the concatenation operator (dot), union operator |, or star operator ∗, respectively. We can construct a syntax tree for a regular expression just as we did for arithmetic expressions in Section 2.5.1.

**Example 3.31 :** Figure 3.56 shows the syntax tree for the regular expression of our running example. Cat-nodes are represented by circles.   □



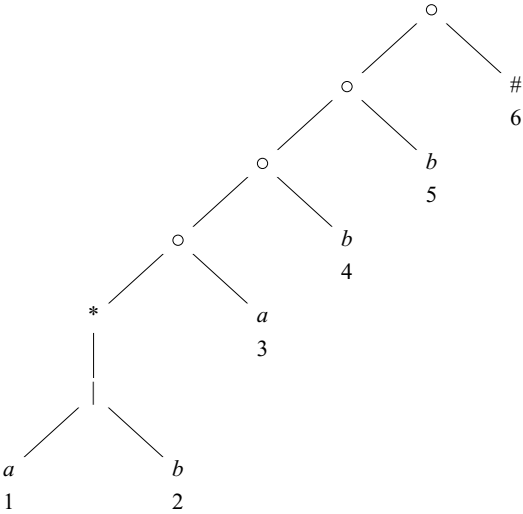Figure 3.56: Syntax tree for $(\mathbf{a}|\mathbf{b})^*\mathbf{abb}\#$

Leaves in a syntax tree are labeled by $\epsilon$ or by an alphabet symbol. To each leaf not labeled $\epsilon$, we attach a unique integer. We refer to this integer as the

*position* of the leaf and also as a position of its symbol. Note that a symbol can have several positions; for instance, $a$ has positions 1 and 3 in Fig. 3.56. The positions in the syntax tree correspond to the important states of the constructed NFA.

**Example 3.32 :** Figure 3.57 shows the NFA for the same regular expression as Fig. 3.56, with the important states numbered and other states represented by letters. The numbered states in the NFA and the positions in the syntax tree correspond in a way we shall soon see.   □
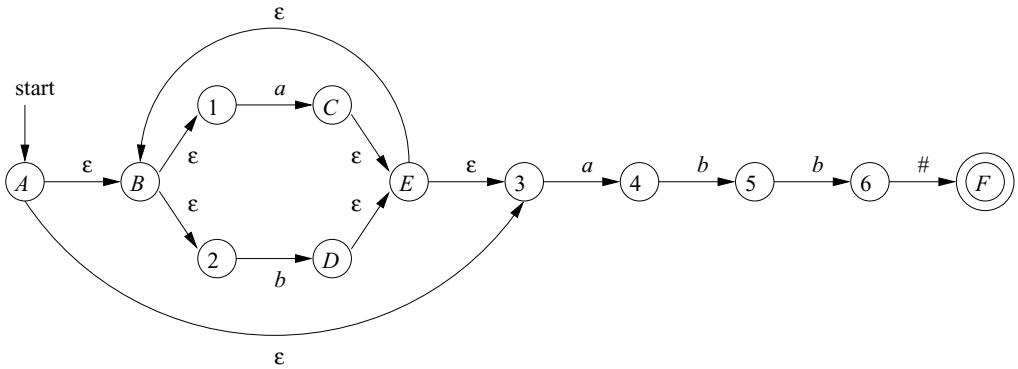


Figure 3.57: NFA constructed by Algorithm 3.23 for $(\mathbf{a}|\mathbf{b})^*\mathbf{abb}\#$

## 3.9.2   Functions Computed From the Syntax Tree

To construct a DFA directly from a regular expression, we construct its syntax tree and then compute four functions: *nullable*, *firstpos*, *lastpos*, and *followpos*, defined as follows. Each definition refers to the syntax tree for a particular augmented regular expression $(r)\#$.

1. *nullable*$(n)$ is true for a syntax-tree node $n$ if and only if the subexpression represented by $n$ has $\epsilon$ in its language. That is, the subexpression can be "made null" or the empty string, even though there may be other strings it can represent as well.

2. *firstpos*$(n)$ is the set of positions in the subtree rooted at $n$ that correspond to the first symbol of at least one string in the language of the subexpression rooted at $n$.

3. *lastpos*$(n)$ is the set of positions in the subtree rooted at $n$ that correspond to the last symbol of at least one string in the language of the subexpression rooted at $n$.

4. *followpos*(p), for a position $p$, is the set of positions $q$ in the entire syntax tree such that there is some string $x = a_1 a_2 \cdots a_n$ in $L\big((r)\#\big)$ such that for some $i$, there is a way to explain the membership of $x$ in $L\big((r)\#\big)$ by matching $a_i$ to position $p$ of the syntax tree and $a_{i+1}$ to position $q$.

**Example 3.33 :** Consider the cat-node $n$ in Fig. 3.56 that corresponds to the expression $(\mathbf{a}|\mathbf{b})^*\mathbf{a}$. We claim *nullable*(n) is false, since this node generates all strings of $a$'s and $b$'s ending in an $a$; it does not generate $\epsilon$. On the other hand, the star-node below it is nullable; it generates $\epsilon$ along with all other strings of $a$'s and $b$'s.

*firstpos*(n) = $\{1, 2, 3\}$. In a typical generated string like $aa$, the first position of the string corresponds to position 1 of the tree, and in a string like $ba$, the first position of the string comes from position 2 of the tree. However, when the string generated by the expression of node $n$ is just $a$, then this $a$ comes from position 3.

*lastpos*(n) = $\{3\}$. That is, no matter what string is generated from the expression of node $n$, the last position is the $a$ from position 3 of the tree.

*followpos* is trickier to compute, but we shall see the rules for doing so shortly. Here is an example of the reasoning: *followpos*(1) = $\{1, 2, 3\}$. Consider a string $\cdots ac \cdots$, where the $c$ is either $a$ or $b$, and the $a$ comes from position 1. That is, this $a$ is one of those generated by the $\mathbf{a}$ in expression $(\mathbf{a}|\mathbf{b})^*$. This $a$ could be followed by another $a$ or $b$ coming from the same subexpression, in which case $c$ comes from position 1 or 2. It is also possible that this $a$ is the last in the string generated by $(\mathbf{a}|\mathbf{b})^*$, in which case the symbol $c$ must be the $a$ that comes from position 3. Thus, 1, 2, and 3 are exactly the positions that can follow position 1.  $\square$

### 3.9.3   Computing *nullable*, *firstpos*, and *lastpos*

We can compute *nullable*, *firstpos*, and *lastpos* by a straightforward recursion on the height of the tree. The basis and inductive rules for *nullable* and *firstpos* are summarized in Fig. 3.58. The rules for *lastpos* are essentially the same as for *firstpos*, but the roles of children $c_1$ and $c_2$ must be swapped in the rule for a cat-node.

**Example 3.34 :** Of all the nodes in Fig. 3.56 only the star-node is nullable. We note from the table of Fig. 3.58 that none of the leaves are nullable, because they each correspond to non-$\epsilon$ operands. The or-node is not nullable, because neither of its children is. The star-node is nullable, because every star-node is nullable. Finally, each of the cat-nodes, having at least one nonnullable child, is not nullable.

The computation of *firstpos* and *lastpos* for each of the nodes is shown in Fig. 3.59, with *firstpos*(n) to the left of node $n$, and *lastpos*(n) to its right. Each of the leaves has only itself for *firstpos* and *lastpos*, as required by the rule for non-$\epsilon$ leaves in Fig. 3.58. For the or-node, we take the union of *firstpos* at the

| NODE $n$ | $nullable(n)$ | $firstpos(n)$ |
|---|---|---|
| A leaf labeled $\epsilon$ | **true** | $\emptyset$ |
| A leaf with position $i$ | **false** | $\{i\}$ |
| An or-node $n = c_1 \vert c_2$ | $nullable(c_1)$ **or** $nullable(c_2)$ | $firstpos(c_1) \cup firstpos(c_2)$ |
| A cat-node $n = c_1 c_2$ | $nullable(c_1)$ **and** $nullable(c_2)$ | **if** ( $nullable(c_1)$ ) $firstpos(c_1) \cup firstpos(c_2)$ **else** $firstpos(c_1)$ |
| A star-node $n = c_1{}^*$ | **true** | $firstpos(c_1)$ |

Figure 3.58: Rules for computing *nullable* and *firstpos*

children and do the same for *lastpos*. The rule for the star-node says that we take the value of *firstpos* or *lastpos* at the one child of that node.

Now, consider the lowest cat-node, which we shall call $n$. To compute *firstpos*($n$), we first consider whether the left operand is nullable, which it is in this case. Therefore, *firstpos* for $n$ is the union of *firstpos* for each of its children, that is $\{1, 2\} \cup \{3\} = \{1, 2, 3\}$. The rule for *lastpos* does not appear explicitly in Fig. 3.58, but as we mentioned, the rules are the same as for *firstpos*, with the children interchanged. That is, to compute *lastpos*($n$) we must ask whether its right child (the leaf with position 3) is nullable, which it is not. Therefore, *lastpos*($n$) is the same as *lastpos* of the right child, or $\{3\}$. □

### 3.9.4 Computing *followpos*

Finally, we need to see how to compute *followpos*. There are only two ways that a position of a regular expression can be made to follow another.

1. If $n$ is a cat-node with left child $c_1$ and right child $c_2$, then for every position $i$ in *lastpos*($c_1$), all positions in *firstpos*($c_2$) are in *followpos*($i$).

2. If $n$ is a star-node, and $i$ is a position in *lastpos*($n$), then all positions in *firstpos*($n$) are in *followpos*($i$).

**Example 3.35 :** Let us continue with our running example; recall that *firstpos* and *lastpos* were computed in Fig. 3.59. Rule 1 for *followpos* requires that we look at each cat-node, and put each position in *firstpos* of its right child in *followpos* for each position in *lastpos* of its left child. For the lowest cat-node in Fig. 3.59, that rule says position 3 is in *followpos*(1) and *followpos*(2). The next cat-node above says that 4 is in *followpos*(3), and the remaining two cat-nodes give us 5 in *followpos*(4) and 6 in *followpos*(5).
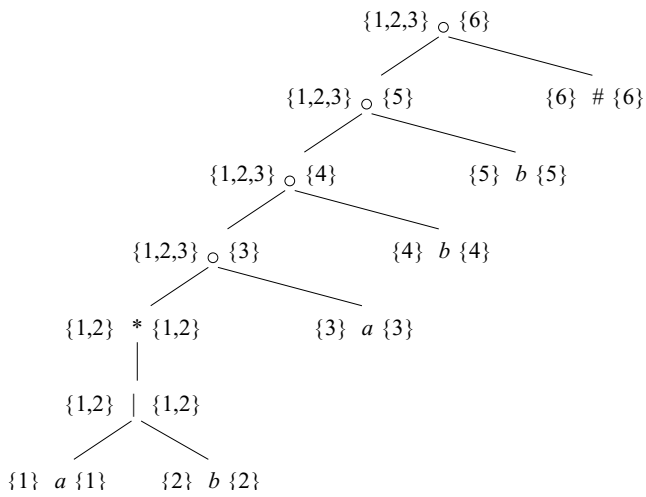
{1,2,3} o {6}

{1,2,3} o {5}                    {6}  # {6}

{1,2,3} o {4}              {5}  b {5}

{1,2,3} o {3}          {4}  b {4}

{1,2}  * {1,2}          {3}  a {3}

{1,2}  | {1,2}

{1}  a {1}      {2}  b {2}

Figure 3.59: *firstpos* and *lastpos* for nodes in the syntax tree for $(\mathbf{a}|\mathbf{b})^*\mathbf{abb}\#$

We must also apply rule 2 to the star-node. That rule tells us positions 1 and 2 are in both *followpos*(1) and *followpos*(2), since both *firstpos* and *lastpos* for this node are $\{1, 2\}$. The complete sets *followpos* are summarized in Fig. 3.60. □

| POSITION   $n$ | $followpos(n)$ |
|:---:|:---:|
| 1 | $\{1, 2, 3\}$ |
| 2 | $\{1, 2, 3\}$ |
| 3 | $\{4\}$ |
| 4 | $\{5\}$ |
| 5 | $\{6\}$ |
| 6 | $\emptyset$ |

Figure 3.60: The function *followpos*

We can represent the function *followpos* by creating a directed graph with a node for each position and an arc from position $i$ to position $j$ if and only if $j$ is in *followpos*($i$). Figure 3.61 shows this graph for the function of Fig. 3.60.

It should come as no surprise that the graph for *followpos* is almost an NFA without $\epsilon$-transitions for the underlying regular expression, and would become one if we:

1. Make all positions in *firstpos* of the root be initial states,

2. Label each arc from $i$ to $j$ by the symbol at position $i$, and

Figure 3.61: Directed graph for the function *followpos*

3. Make the position associated with endmarker # be the only accepting state.

## 3.9.5 Converting a Regular Expression Directly to a DFA

**Algorithm 3.36 :** Construction of a DFA from a regular expression $r$.

**INPUT**: A regular expression $r$.

**OUTPUT**: A DFA $D$ that recognizes $L(r)$.

**METHOD**:

1. Construct a syntax tree $T$ from the augmented regular expression $(r)\#$.

2. Compute *nullable*, *firstpos*, *lastpos*, and *followpos* for $T$, using the methods of Sections 3.9.3 and 3.9.4.

3. Construct *Dstates*, the set of states of DFA $D$, and *Dtran*, the transition function for $D$, by the procedure of Fig. 3.62. The states of $D$ are sets of positions in $T$. Initially, each state is "unmarked," and a state becomes "marked" just before we consider its out-transitions. The start state of $D$ is *firstpos*$(n_0)$, where node $n_0$ is the root of $T$. The accepting states are those containing the position for the endmarker symbol #.

□

**Example 3.37 :** We can now put together the steps of our running example to construct a DFA for the regular expression $r = (\mathbf{a}|\mathbf{b})^*\mathbf{abb}$. The syntax tree for $(r)\#$ appeared in Fig. 3.56. We observed that for this tree, *nullable* is true only for the star-node, and we exhibited *firstpos* and *lastpos* in Fig. 3.59. The values of *followpos* appear in Fig. 3.60.

The value of *firstpos* for the root of the tree is $\{1, 2, 3\}$, so this set is the start state of $D$. Call this set of states $A$. We must compute *Dtran*$[A, a]$ and *Dtran*$[A, b]$. Among the positions of $A$, 1 and 3 correspond to $a$, while 2 corresponds to $b$. Thus, *Dtran*$[A, a] = followpos(1) \cup followpos(3) = \{1, 2, 3, 4\}$,

> initialize *Dstates* to contain only the unmarked state $firstpos(n_0)$,
>        where $n_0$ is the root of syntax tree $T$ for $(r)\#$;
> **while** ( there is an unmarked state $S$ in *Dstates* ) {
>        mark $S$;
>        **for** ( each input symbol $a$ ) {
>                let $U$ be the union of *followpos*$(p)$ for all $p$
>                        in $S$ that correspond to $a$;
>                **if** ( $U$ is not in *Dstates* )
>                        add $U$ as an unmarked state to *Dstates*;
>                *Dtran*$[S, a] = U$;
>        }
> }

Figure 3.62: Construction of a DFA directly from a regular expression

and *Dtran*$[A, b] = followpos(2) = \{1, 2, 3\}$. The latter is state $A$, and so does not have to be added to *Dstates*, but the former, $B = \{1, 2, 3, 4\}$, is new, so we add it to *Dstates* and proceed to compute its transitions. The complete DFA is shown in Fig. 3.63. $\square$
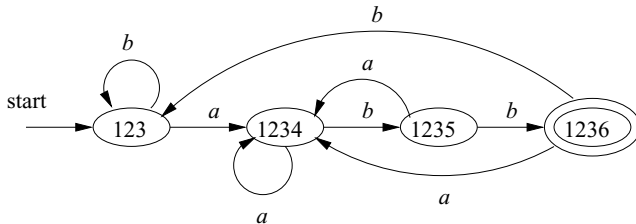


Figure 3.63: DFA constructed from Fig. 3.57

### 3.9.6   Minimizing the Number of States of a DFA

There can be many DFA's that recognize the same language. For instance, note that the DFA's of Figs. 3.36 and 3.63 both recognize language $L\big((\mathbf{a}|\mathbf{b})^*\mathbf{abb}\big)$. Not only do these automata have states with different names, but they don't even have the same number of states. If we implement a lexical analyzer as a DFA, we would generally prefer a DFA with as few states as possible, since each state requires entries in the table that describes the lexical analyzer.

The matter of the names of states is minor. We shall say that two automata are *the same up to state names* if one can be transformed into the other by doing nothing more than changing the names of states. Figures 3.36 and 3.63 are not the same up to state names. However, there is a close relationship between the

states of each. States $A$ and $C$ of Fig. 3.36 are actually equivalent, in the sense that neither is an accepting state, and on any input they transfer to the same state — to $B$ on input $a$ and to $C$ on input $b$. Moreover, both states $A$ and $C$ behave like state 123 of Fig. 3.63. Likewise, state $B$ of Fig. 3.36 behaves like state 1234 of Fig. 3.63, state $D$ behaves like state 1235, and state $E$ behaves like state 1236.

It turns out that there is always a unique (up to state names) minimum state DFA for any regular language. Moreover, this minimum-state DFA can be constructed from any DFA for the same language by grouping sets of equivalent states. In the case of $L\big((\mathbf{a}|\mathbf{b})^*\mathbf{abb}\big)$, Fig. 3.63 is the minimum-state DFA, and it can be constructed by partitioning the states of Fig. 3.36 as $\{A, C\}\{B\}\{D\}\{E\}$.

In order to understand the algorithm for creating the partition of states that converts any DFA into its minimum-state equivalent DFA, we need to see how input strings distinguish states from one another. We say that string $x$ *distinguishes* state $s$ from state $t$ if exactly one of the states reached from $s$ and $t$ by following the path with label $x$ is an accepting state. State $s$ is *distinguishable* from state $t$ if there is some string that distinguishes them.

**Example 3.38 :** The empty string distinguishes any accepting state from any nonaccepting state. In Fig. 3.36, the string $bb$ distinguishes state $A$ from state $B$, since $bb$ takes $A$ to a nonaccepting state $C$, but takes $B$ to accepting state $E$. □

The state-minimization algorithm works by partitioning the states of a DFA into groups of states that cannot be distinguished. Each group of states is then merged into a single state of the minimum-state DFA. The algorithm works by maintaining a partition, whose groups are sets of states that have not yet been distinguished, while any two states from different groups are known to be distinguishable. When the partition cannot be refined further by breaking any group into smaller groups, we have the minimum-state DFA.

Initially, the partition consists of two groups: the accepting states and the nonaccepting states. The fundamental step is to take some group of the current partition, say $A = \{s_1, s_2, \dots, s_k\}$, and some input symbol $a$, and see whether $a$ can be used to distinguish between any states in group $A$. We examine the transitions from each of $s_1, s_2, \dots, s_k$ on input $a$, and if the states reached fall into two or more groups of the current partition, we split $A$ into a collection of groups, so that $s_i$ and $s_j$ are in the same group if and only if they go to the same group on input $a$. We repeat this process of splitting groups, until for no group, and for no input symbol, can the group be split further. The idea is formalized in the next algorithm.

**Algorithm 3.39 :** Minimizing the number of states of a DFA.

**INPUT**: A DFA $D$ with set of states $S$, input alphabet $\Sigma$, start state $s_0$, and set of accepting states $F$.

**OUTPUT**: A DFA $D'$ accepting the same language as $D$ and having as few states as possible.

---

### Why the State-Minimization Algorithm Works

We need to prove two things: that states remaining in the same group in $\Pi_{\text{final}}$ are indistinguishable by any string, and that states winding up in different groups are distinguishable. The first is an induction on $i$ that if after the $i$th iteration of step (2) of Algorithm 3.39, $s$ and $t$ are in the same group, then there is no string of length $i$ or less that distinguishes them. We shall leave the details of the induction to you.

The second is an induction on $i$ that if states $s$ and $t$ are placed in different groups at the $i$th iteration of step (2), then there is a string that distinguishes them. The basis, when $s$ and $t$ are placed in different groups of the initial partition, is easy: one must be accepting and the other not, so $\epsilon$ distinguishes them. For the induction, there must be an input $a$ and states $p$ and $q$ such that $s$ and $t$ go to states $p$ and $q$, respectively, on input $a$. Moreover, $p$ and $q$ must already have been placed in different groups. Then by the inductive hypothesis, there is some string $x$ that distinguishes $p$ from $q$. Therefore, $ax$ distinguishes $s$ from $t$.

---

**METHOD**:

1. Start with an initial partition $\Pi$ with two groups, $F$ and $S - F$, the accepting and nonaccepting states of $D$.

2. Apply the procedure of Fig. 3.64 to construct a new partition $\Pi_{\text{new}}$.

   > initially, let $\Pi_{\text{new}} = \Pi$;
   > **for** ( each group $G$ of $\Pi$ ) {
   >         partition $G$ into subgroups such that two states $s$ and $t$
   >                 are in the same subgroup if and only if for all
   >                 input symbols $a$, states $s$ and $t$ have transitions on $a$
   >                 to states in the same group of $\Pi$;
   >         /* at worst, a state will be in a subgroup by itself */
   >         replace $G$ in $\Pi_{\text{new}}$ by the set of all subgroups formed;
   > }

   Figure 3.64: Construction of $\Pi_{\text{new}}$

3. If $\Pi_{\text{new}} = \Pi$, let $\Pi_{\text{final}} = \Pi$ and continue with step (4). Otherwise, repeat step (2) with $\Pi_{\text{new}}$ in place of $\Pi$.

4. Choose one state in each group of $\Pi_{\text{final}}$ as the *representative* for that group. The representatives will be the states of the minimum-state DFA $D'$. The other components of $D'$ are constructed as follows:

---

### Eliminating the Dead State

The minimization algorithm sometimes produces a DFA with one dead state — one that is not accepting and transfers to itself on each input symbol. This state is technically needed, because a DFA must have a transition from every state on every symbol. However, as discussed in Section 3.8.3, we often want to know when there is no longer any possibility of acceptance, so we can establish that the proper lexeme has already been seen. Thus, we may wish to eliminate the dead state and use an automaton that is missing some transitions. This automaton has one fewer state than the minimum-state DFA, but is strictly speaking not a DFA, because of the missing transitions to the dead state.

---

(a) The start state of $D'$ is the representative of the group containing the start state of $D$.

(b) The accepting states of $D'$ are the representatives of those groups that contain an accepting state of $D$. Note that each group contains either only accepting states, or only nonaccepting states, because we started by separating those two classes of states, and the procedure of Fig. 3.64 always forms new groups that are subgroups of previously constructed groups.

(c) Let $s$ be the representative of some group $G$ of $\Pi_{\text{final}}$, and let the transition of $D$ from $s$ on input $a$ be to state $t$. Let $r$ be the representative of $t$'s group $H$. Then in $D'$, there is a transition from $s$ to $r$ on input $a$. Note that in $D$, every state in group $G$ must go to some state of group $H$ on input $a$, or else, group $G$ would have been split according to Fig. 3.64.

☐

**Example 3.40 :** Let us reconsider the DFA of Fig. 3.36. The initial partition consists of the two groups $\{A, B, C, D\}\{E\}$, which are respectively the nonaccepting states and the accepting states. To construct $\Pi_{\text{new}}$, the procedure of Fig. 3.64 considers both groups and inputs $a$ and $b$. The group $\{E\}$ cannot be split, because it has only one state, so $\{E\}$ will remain intact in $\Pi_{\text{new}}$.

The other group $\{A, B, C, D\}$ can be split, so we must consider the effect of each input symbol. On input $a$, each of these states goes to state $B$, so there is no way to distinguish these states using strings that begin with $a$. On input $b$, states $A$, $B$, and $C$ go to members of group $\{A, B, C, D\}$, while state $D$ goes to $E$, a member of another group. Thus, in $\Pi_{\text{new}}$, group $\{A, B, C, D\}$ is split into $\{A, B, C\}\{D\}$, and $\Pi_{\text{new}}$ for this round is $\{A, B, C\}\{D\}\{E\}$.

In the next round, we can split $\{A, B, C\}$ into $\{A, C\}\{B\}$, since $A$ and $C$ each go to a member of $\{A, B, C\}$ on input $b$, while $B$ goes to a member of another group, $\{D\}$. Thus, after the second round, $\Pi_{\text{new}} = \{A, C\}\{B\}\{D\}\{E\}$. For the third round, we cannot split the one remaining group with more than one state, since $A$ and $C$ each go to the same state (and therefore to the same group) on each input. We conclude that $\Pi_{\text{final}} = \{A, C\}\{B\}\{D\}\{E\}$.

Now, we shall construct the minimum-state DFA. It has four states, corresponding to the four groups of $\Pi_{\text{final}}$, and let us pick $A$, $B$, $D$, and $E$ as the representatives of these groups. The initial state is $A$, and the only accepting state is $E$. Figure 3.65 shows the transition function for the DFA. For instance, the transition from state $E$ on input $b$ is to $A$, since in the original DFA, $E$ goes to $C$ on input $b$, and $A$ is the representative of $C$'s group. For the same reason, the transition on $b$ from state $A$ is to $A$ itself, while all other transitions are as in Fig. 3.36.   □

| STATE | $a$ | $b$ |
|:-----:|:---:|:---:|
| $A$ | $B$ | $A$ |
| $B$ | $B$ | $D$ |
| $D$ | $B$ | $E$ |
| $E$ | $B$ | $A$ |

Figure 3.65: Transition table of minimum-state DFA

### 3.9.7   State Minimization in Lexical Analyzers

To apply the state minimization procedure to the DFA's generated in Section 3.8.3, we must begin Algorithm 3.39 with the partition that groups together all states that recognize a particular token, and also places in one group all those states that do not indicate any token. An example should make the extension clear.

**Example 3.41 :** For the DFA of Fig. 3.54, the initial partition is

$$\{0137, 7\}\{247\}\{8, 58\}\{68\}\{\emptyset\}$$

That is, states 0137 and 7 belong together because neither announces any token. States 8 and 58 belong together because they both announce token $\mathbf{a}^*\mathbf{b}^+$. Note that we have added a dead state $\emptyset$, which we suppose has transitions to itself on inputs $a$ and $b$. The dead state is also the target of missing transitions on $a$ from states 8, 58, and 68.

We must split 0137 from 7, because they go to different groups on input $a$. We also split 8 from 58, because they go to different groups on $b$. Thus, all states are in groups by themselves, and Fig. 3.54 is the minimum-state DFA

recognizing its three tokens. Recall that a DFA serving as a lexical analyzer will normally drop the dead state, while we treat missing transitions as a signal to end token recognition. ☐

## 3.9.8 Trading Time for Space in DFA Simulation

The simplest and fastest way to represent the transition function of a DFA is a two-dimensional table indexed by states and characters. Given a state and next input character, we access the array to find the next state and any special action we must take, e.g., returning a token to the parser. Since a typical lexical analyzer has several hundred states in its DFA and involves the ASCII alphabet of 128 input characters, the array consumes less than a megabyte.

However, compilers are also appearing in very small devices, where even a megabyte of storage may be too much. For such situations, there are many methods that can be used to compact the transition table. For instance, we can represent each state by a list of transitions — that is, character-state pairs — ended by a default state that is to be chosen for any input character not on the list. If we choose as the default the most frequently occurring next state, we can often reduce the amount of storage needed by a large factor.

There is a more subtle data structure that allows us to combine the speed of array access with the compression of lists with defaults. We may think of this structure as four arrays, as suggested in Fig. 3.66.[5] The *base* array is used to determine the base location of the entries for state $s$, which are located in the *next* and *check* arrays. The *default* array is used to determine an alternative base location if the *check* array tells us the one given by *base*[$s$] is invalid.



Figure 3.66: Data structure for representing transition tables

To compute *nextState*($s, a$), the transition for state $s$ on input $a$, we examine the *next* and *check* entries in location $l = base[s] + a$, where character $a$ is treated as an integer, presumably in the range 0 to 127. If *check*[$l$] = $s$, then this entry

---

[5] In practice, there would be another array indexed by states to give the action associated with that state, if any.

is valid, and the next state for state $s$ on input $a$ is $next[l]$. If $check[l] \neq s$, then we determine another state $t = default[s]$ and repeat the process as if $t$ were the current state. More formally, the function $nextState$ is defined as follows:

```
int nextState(s, a) {
     if ( check[base[s] + a]  ==  s ) return next[base[s] + a];
     else return nextState(default[s], a);
}
```

The intended use of the structure of Fig. 3.66 is to make the *next-check* arrays short by taking advantage of the similarities among states. For instance, state $t$, the default for state $s$, might be the state that says "we are working on an identifier," like state 10 in Fig. 3.14. Perhaps state $s$ is entered after seeing the letters `th`, which are a prefix of keyword `then` as well as potentially being the prefix of some lexeme for an identifier. On input character `e`, we must go from state $s$ to a special state that remembers we have seen `the`, but otherwise, state $s$ behaves as $t$ does. Thus, we set $check[base[s] + e]$ to $s$ (to confirm that this entry is valid for $s$) and we set $next[base[s] + e]$ to the state that remembers `the`. Also, $default[s]$ is set to $t$.

While we may not be able to choose *base* values so that no *next-check* entries remain unused, experience has shown that the simple strategy of assigning *base* values to states in turn, and assigning each $base[s]$ value the lowest integer so that the special entries for state $s$ are not previously occupied utilizes little more space than the minimum possible.

### 3.9.9   Exercises for Section 3.9

**Exercise 3.9.1:** Extend the table of Fig. 3.58 to include the operators (a) ? and (b) $^+$.

**Exercise 3.9.2:** Use Algorithm 3.36 to convert the regular expressions of Exercise 3.7.3 directly to deterministic finite automata.

! **Exercise 3.9.3:** We can prove that two regular expressions are equivalent by showing that their minimum-state DFA's are the same up to renaming of states. Show in this way that the following regular expressions: $(\mathbf{a}|\mathbf{b})^*$, $(\mathbf{a}^*|\mathbf{b}^*)^*$, and $\big((\epsilon|\mathbf{a})\mathbf{b}^*\big)^*$ are all equivalent. *Note*: You may have constructed the DFA's for these expressions in response to Exercise 3.7.3.

! **Exercise 3.9.4:** Construct the minimum-state DFA's for the following regular expressions:

  a) $(\mathbf{a}|\mathbf{b})^*\mathbf{a}(\mathbf{a}|\mathbf{b})$.

  b) $(\mathbf{a}|\mathbf{b})^*\mathbf{a}(\mathbf{a}|\mathbf{b})(\mathbf{a}|\mathbf{b})$.

  c) $(\mathbf{a}|\mathbf{b})^*\mathbf{a}(\mathbf{a}|\mathbf{b})(\mathbf{a}|\mathbf{b})(\mathbf{a}|\mathbf{b})$.

Do you see a pattern?

**!! Exercise 3.9.5 :** To make formal the informal claim of Example 3.25, show that any deterministic finite automaton for the regular expression

$$(\mathbf{a}|\mathbf{b})^*\mathbf{a}(\mathbf{a}|\mathbf{b})(\mathbf{a}|\mathbf{b})\cdots(\mathbf{a}|\mathbf{b})$$

where $(\mathbf{a}|\mathbf{b})$ appears $n-1$ times at the end, must have at least $2^n$ states. *Hint*: Observe the pattern in Exercise 3.9.4. What condition regarding the history of inputs does each state represent?

## 3.10 Summary of Chapter 3

✦ *Tokens*. The lexical analyzer scans the source program and produces as output a sequence of tokens, which are normally passed, one at a time to the parser. Some tokens may consist only of a token name while others may also have an associated lexical value that gives information about the particular instance of the token that has been found on the input.

✦ *Lexemes*. Each time the lexical analyzer returns a token to the parser, it has an associated lexeme — the sequence of input characters that the token represents.

✦ *Buffering*. Because it is often necessary to scan ahead on the input in order to see where the next lexeme ends, it is usually necessary for the lexical analyzer to buffer its input. Using a pair of buffers cyclicly and ending each buffer's contents with a sentinel that warns of its end are two techniques that accelerate the process of scanning the input.

✦ *Patterns*. Each token has a pattern that describes which sequences of characters can form the lexemes corresponding to that token. The set of words, or strings of characters, that match a given pattern is called a language.

✦ *Regular Expressions*. These expressions are commonly used to describe patterns. Regular expressions are built from single characters, using union, concatenation, and the Kleene closure, or any-number-of, operator.

✦ *Regular Definitions*. Complex collections of languages, such as the patterns that describe the tokens of a programming language, are often defined by a regular definition, which is a sequence of statements that each define one variable to stand for some regular expression. The regular expression for one variable can use previously defined variables in its regular expression.

✦ *Extended Regular-Expression Notation.*  A number of additional opera-
tors may appear as shorthands in regular expressions, to make it easier
to express patterns.  Examples include the + operator (one-or-more-of),
? (zero-or-one-of), and character classes (the union of the strings each
consisting of one of the characters).

✦ *Transition Diagrams.*  The behavior of a lexical analyzer can often be
described by a transition diagram.  These diagrams have states, each
of which represents something about the history of the characters seen
during the current search for a lexeme that matches one of the possible
patterns.  There are arrows, or transitions, from one state to another,
each of which indicates the possible next input characters that cause the
lexical analyzer to make that change of state.

✦ *Finite Automata.*  These are a formalization of transition diagrams that
include a designation of a start state and one or more accepting states,
as well as the set of states, input characters, and transitions among
states.  Accepting states indicate that the lexeme for some token has been
found.  Unlike transition diagrams, finite automata can make transitions
on empty input as well as on input characters.

✦ *Deterministic Finite Automata.*  A DFA is a special kind of finite au-
tomaton that has exactly one transition out of each state for each input
symbol.  Also, transitions on empty input are disallowed.  The DFA is
easily simulated and makes a good implementation of a lexical analyzer,
similar to a transition diagram.

✦ *Nondeterministic Finite Automata.*  Automata that are not DFA's are
called nondeterministic.  NFA's often are easier to design than are DFA's.
Another possible architecture for a lexical analyzer is to tabulate all the
states that NFA's for each of the possible patterns can be in, as we scan
the input characters.

✦ *Conversion Among Pattern Representations.* It is possible to convert any
regular expression into an NFA of about the same size, recognizing the
same language as the regular expression defines.  Further, any NFA can
be converted to a DFA for the same pattern, although in the worst case
(never encountered in common programming languages) the size of the
automaton can grow exponentially. It is also possible to convert any non-
deterministic or deterministic finite automaton into a regular expression
that defines the same language recognized by the finite automaton.

✦ *Lex.*  There is a family of software systems, including `Lex` and `Flex`,
that are lexical-analyzer generators.  The user specifies the patterns for
tokens using an extended regular-expression notation. `Lex` converts these
expressions into a lexical analyzer that is essentially a deterministic finite
automaton that recognizes any of the patterns.

♦ *Minimization of Finite Automata.* For every DFA there is a minimum-state DFA accepting the same language. Moreover, the minimum-state DFA for a given language is unique except for the names given to the various states.

# 3.11 References for Chapter 3

Regular expressions were first developed by Kleene in the 1950's [9]. Kleene was interested in describing the events that could be represented by McCullough and Pitts' [12] finite-automaton model of neural activity. Since that time regular expressions and finite automata have become widely used in computer science.

Regular expressions in various forms were used from the outset in many popular Unix utilities such as `awk`, `ed`, `egrep`, `grep`, `lex`, `sed`, `sh`, and `vi`. The IEEE 1003 and ISO/IEC 9945 standards documents for the Portable Operating System Interface (POSIX) define the POSIX extended regular expressions which are similar to the original Unix regular expressions with a few exceptions such as mnemonic representations for character classes. Many scripting languages such as Perl, Python, and Tcl have adopted regular expressions but often with incompatible extensions.

The familiar finite-automaton model and the minimization of finite automata, as in Algorithm 3.39, come from Huffman [6] and Moore [14]. Nondeterministic finite automata were first proposed by Rabin and Scott [15]; the subset construction of Algorithm 3.20, showing the equivalence of deterministic and nondeterministic finite automata, is from there.

McNaughton and Yamada [13] first gave an algorithm to convert regular expressions directly to deterministic finite automata. Algorithm 3.36 described in Section 3.9 was first used by Aho in creating the Unix regular-expression matching tool `egrep`. This algorithm was also used in the regular-expression pattern matching routines in `awk` [3]. The approach of using nondeterministic automata as an intermediary is due Thompson [17]. The latter paper also contains the algorithm for the direct simulation of nondeterministic finite automata (Algorithm 3.22), which was used by Thompson in the text editor `QED`.

Lesk developed the first version of `Lex` and then Lesk and Schmidt created a second version using Algorithm 3.36 [10]. Many variants of `Lex` have been subsequently implemented. The GNU version, `Flex`, can be downloaded, along with documentation at [4]. Popular Java versions of `Lex` include `JFlex` [7] and `JLex` [8].

The KMP algorithm, discussed in the exercises to Section 3.4 just prior to Exercise 3.4.3, is from [11]. Its generalization to many keywords appears in [2] and was used by Aho in the first implementation of the Unix utility `fgrep`.

The theory of finite automata and regular expressions is covered in [5]. A survey of string-matching techniques is in [1].

1. Aho, A. V., "Algorithms for finding patterns in strings," in *Handbook of Theoretical Computer Science* (J. van Leeuwen, ed.), Vol. A, Ch. 5, MIT

Press, Cambridge, 1990.

2. Aho, A. V. and M. J. Corasick, "Efficient string matching: an aid to bibliographic search," *Comm. ACM* **18**:6 (1975), pp. 333–340.

3. Aho, A. V., B. W. Kernighan, and P. J. Weinberger, *The AWK Programming Language*, Addison-Wesley, Boston, MA, 1988.

4. Flex home page `http://www.gnu.org/software/flex/`, Free Software Foundation.

5. Hopcroft, J. E., R. Motwani, and J. D. Ullman, *Introduction to Automata Theory, Languages, and Computation*, Addison-Wesley, Boston MA, 2006.

6. Huffman, D. A., "The synthesis of sequential machines," *J. Franklin Inst.* **257** (1954), pp. 3–4, 161, 190, 275–303.

7. JFlex home page `http://jflex.de/` .

8. `http://www.cs.princeton.edu/~appel/modern/java/JLex` .

9. Kleene, S. C., "Representation of events in nerve nets," in [16], pp. 3–40.

10. Lesk, M. E., "Lex – a lexical analyzer generator," Computing Science Tech. Report 39, Bell Laboratories, Murray Hill, NJ, 1975. A similar document with the same title but with E. Schmidt as a coauthor, appears in Vol. 2 of the *Unix Programmer's Manual*, Bell laboratories, Murray Hill NJ, 1975; see `http://dinosaur.compilertools.net/lex/index.html` .

11. Knuth, D. E., J. H. Morris, and V. R. Pratt, "Fast pattern matching in strings," *SIAM J. Computing* **6**:2 (1977), pp. 323–350.

12. McCullough, W. S. and W. Pitts, "A logical calculus of the ideas immanent in nervous activity," *Bull. Math. Biophysics* **5** (1943), pp. 115–133.

13. McNaughton, R. and H. Yamada, "Regular expressions and state graphs for automata," *IRE Trans. on Electronic Computers* **EC-9**:1 (1960), pp. 38–47.

14. Moore, E. F., "Gedanken experiments on sequential machines," in [16], pp. 129–153.

15. Rabin, M. O. and D. Scott, "Finite automata and their decision problems," *IBM J. Res. and Devel.* **3**:2 (1959), pp. 114–125.

16. Shannon, C. and J. McCarthy (eds.), *Automata Studies*, Princeton Univ. Press, 1956.

17. Thompson, K., "Regular expression search algorithm," *Comm. ACM* **11**:6 (1968), pp. 419–422.