

Chapter 2 Lexical Analysis

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.

- When the lexical analyzer discovers a lexeme constituting an identifier, it needs to enter that lexeme into the symbol table.

- Kind of identifiers may be read from the symbol table by the lexical analyzer to assist it in determining the parser token it must pass to the parser.

- These interactions are implemented by having the parser call the lexical analyzer.

- The call, suggested by the getNextToken command, causes the lexical analyzer

The call 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.

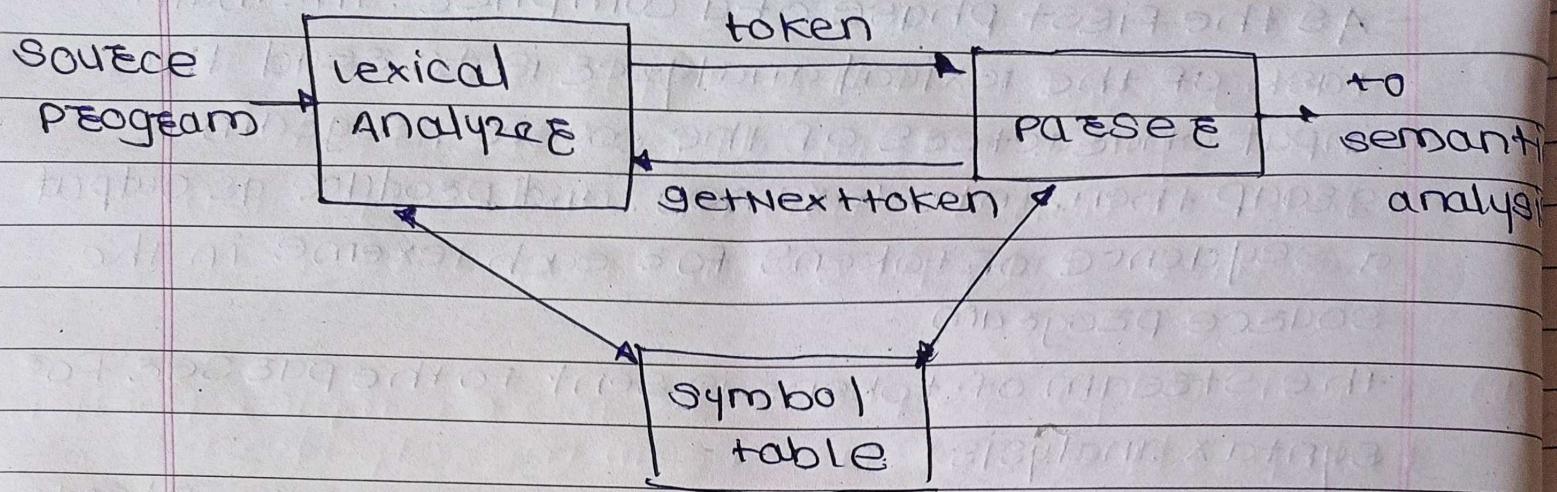


fig : interactions between the lexical analyse and the parser

lexical analyse 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.

Another task is correlating error messages generated by the compiler with the source program.

For instance, the lexical analyse may keep track of the number of newline characters seen so it can associate a line number with each error message generated by the compiler with the source program.

The lexical analyse makes a copy of the source program with the error messages inserted at the appropriate positions. if the source program uses a macro-processor, the expansion of macros may also be performed by the lexical analyse.

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.
 - c) tokens, patterns, and lexemes.
- ① 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.
- ② 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's matched by many strings.

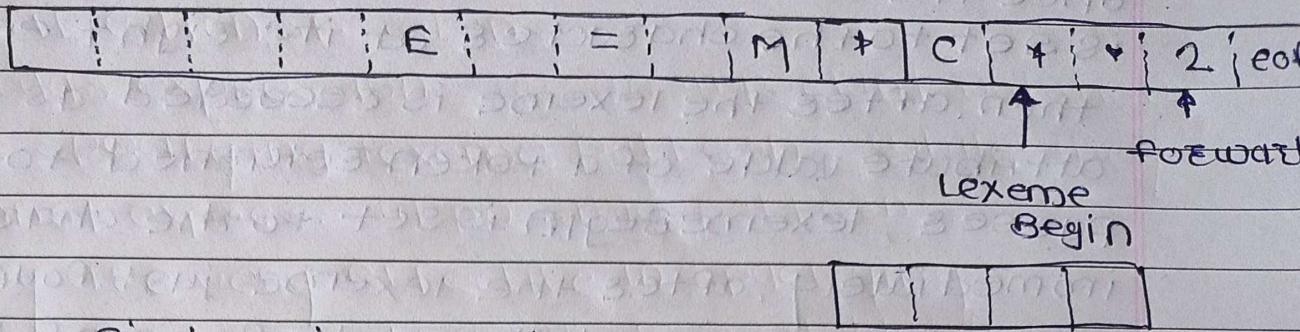
3) A lexeme is a sequence of characters in the source program that matches the pattern for a token and is identified by the lexical analyzer as an instance of that token.

* Input Buffering

① we often have to look one or more characters beyond the next lexeme.
② end of an identifier until we see character that is not a letter or digit, and therefore is not part of the lexeme for id.
In C, single-character operators like -, =, *, / could also be the beginning of a two-character operator like ->, ==, !=, <=, >=, etc.
we shall introduce a two-buffer scheme that handles large lookahead safely, involving "sentinels" that saves time checking for the ends of buffers.

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 loaded, as suggested in.



fig; using a pair of input buffers..

b) 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 in 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.

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.

Sentinels

Thus, for each character read, we make two tests: one for the end of the buffer, and one to determine what character is read. 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. Eof that appears other than at the end of a buffer means that the input is at an end.

specification and recognition of tokens :-

① 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 specifying those types of patterns that we actually need for tokens.

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.

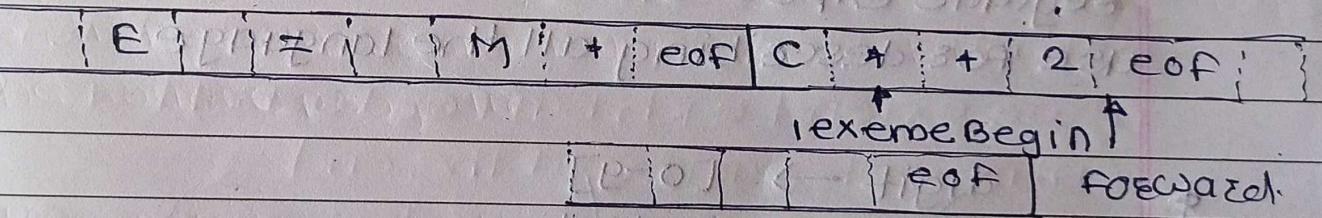


fig : sentinels at the end of each buffer.

```
switch (C*forwarded++) {
```

```
    case EOF :
```

```
        if (forwarded is at end of first buffer) {  
            Eeload second buffer;  
        }
```

```
        forwarded = beginning of second buffer;
```

```
}
```

```
    else if (forwarded is at end of second  
            buffer) {
```

```
        Eeload first buffer;
```

```
        forwarded = beginning beginning of  
                  first buffer;
```

?

else / + eof within a buffer marks the end
input + /
terminate lexical analysis;

break ;

Cases for the other characters

}

fig lookahead code with sentinels

* Regular Expressions

Suppose we wanted to describe the set of valid identifiers. It is almost exactly the language described in item

the only difference is that the underscore is included among the letters.

This process is so useful that a notation called regular expressions has come into common use for describing all the language

*

digit \rightarrow [0-9]

digits \rightarrow digit⁺

number \rightarrow digits (digits)? (E [+-] digits)?

letter \rightarrow [A-zA-Z]

id \rightarrow letter (letter | digit)⁺

if \rightarrow if

then \rightarrow then

else \rightarrow else

brace \rightarrow { } < > | <= > = | < >

patterns for tokens of example

else if eof within a buffer marks the end

input + /

terminate lexical analysis

break;

Cases for the other characters

{}

fig lookahead code with sentinels

+ Regular Expressions

Suppose we wanted to describe the set of valid identifiers. It is almost exactly the language described in item 1. The only difference is that the underscore is included among the letters.

This process is so useful that a notation called regular expressions has come into common use for describing all the language.

digit → [0-9]

digits → digit⁺

number → digits (digits)? (E [+-] digits)?

letter → [A-Za-z]

id → letter (letter | digit)⁺

if → if

then → then

else → else

token → < | > | <= | > = | = | < >

patterns for tokens of example

Recognition of TOKENS:-

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.

stmt → if expE then stmt

| if expE then stmt else stmt

expE → term eelop term

term → id

| number

Fig: A grammar for branching

of nesting selection statements.

for eelop, 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, eelop, id, and number, are the names of tokens as far as the lexical analysis is concerned.

The patterns for these tokens are described using regular definitions.

The patterns for id and number are similar to what we saw in.

digit → [0-9]

digits → digit⁺

number → digits (, digits)? (E (+-)?)? digits)

letter → [A-Za-z]

id → letter (letter | digit)⁺

if → if

then → then

else → else

loop → <|> | <= | > = | <>

fig Patterns for tokens

The lexical analyzer will recognize the keywords if, then, and else, as well as lexemes that

match the patterns for loop, id, and number.

To simplify matters, we make the common assumption that keywords are also reserved words:

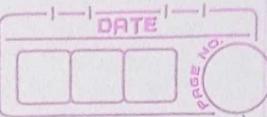
that is, they are not identifiers, even though their lexemes match the pattern for

for identifiers.

In addition, we assign the lexical analyzer the job of stripping out whitespace, by recognizing the "token" we defined by:

ws → (blank | tab | newline) +

* Finite automata implications



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 ϵ , 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.

Nondeterministic Finite Automata

A nondeterministic finite automaton (NFA) consists of:

- 1) A finite set of states S .
- 2) A set of input symbols Σ , the input alphabet. We assume that ϵ , which stands for the empty string, is never a member of Σ .

- DATE _____
- PAGE NO. _____
- 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:

- a) The same symbol can label edges from one state to several different states, and
- b) An edge may be labeled by ϵ , the empty string, instead of, or in addition to, symbols from the input alphabet.

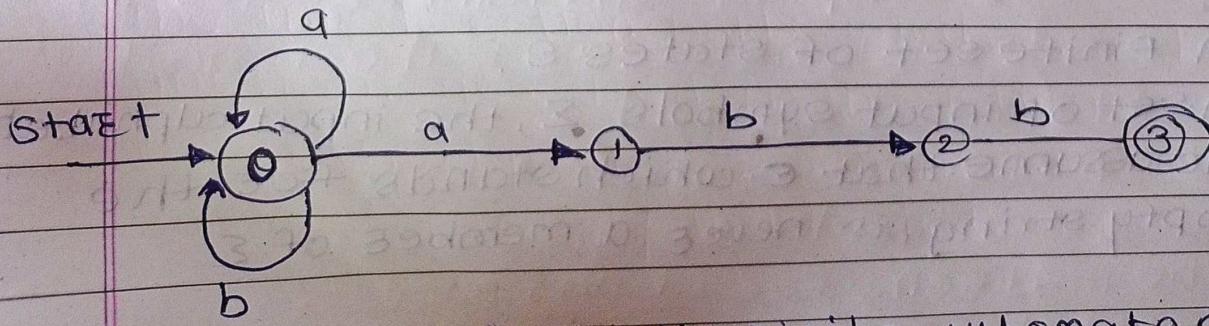


fig:- A nondeterministic finite automaton.

* The lexical-Analyse Generator Lex

Lexical analyse by specifying regular expressions to describe patterns for tokens.

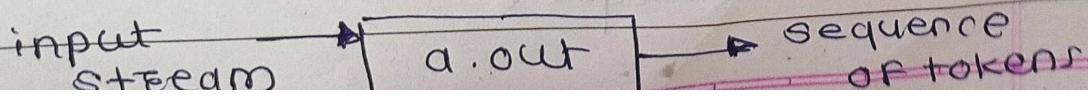
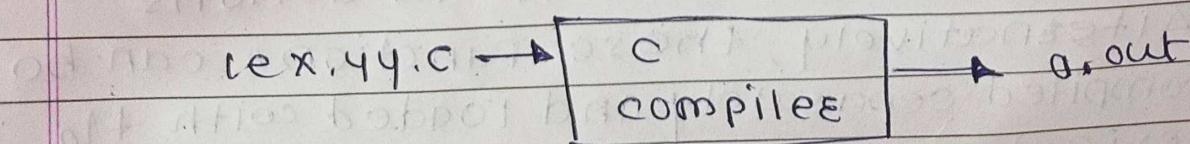
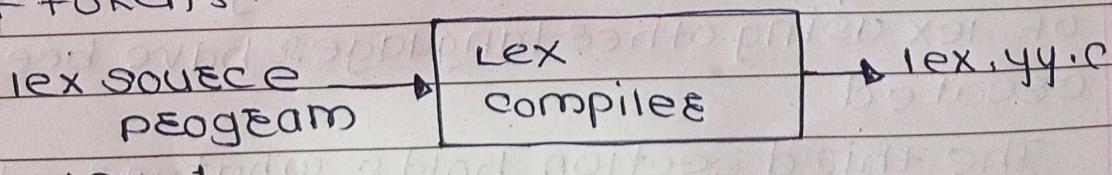
The input notation for the Lex tool is referred to as the Lex language and the tool itself is the Lex compiler.

Behind 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.

use of Lex

Lex is used. An input file, which we call `lex.l`, is written in the Lex language, and describes the lexical analyse 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 analyse that can take a stream of input characters and produce a stream of tokens.



STRUCTURE OF LEX PROGRAMS.

A Lex Program has the following form.

declarations

/* */

translation rules

/* */

auxiliary functions.

The declaration section includes declaration 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

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.