

Study material for Compiler Design:

Lexical Analysis

Reference: Compilers: Principles, Techniques, and Tools by Aho, Sethi & Ullman

3.1 OVER VIEW OF LEXICAL ANALYSIS

- To identify the tokens we need some method of describing the possible tokens that can appear in the input stream. For this purpose we introduce regular expression, a notation that can be used to describe essentially all the tokens of programming language.
- Secondly, having decided what the tokens are, we need some mechanism to recognize these in the input stream. This is done by the token recognizers, which are designed using transition diagrams and finite automata.

3.2 ROLE OF LEXICAL ANALYZER (LA)

The LA is the first phase of a compiler. Its main task is to read the input character and produce as output a sequence of tokens that the parser uses for syntax analysis.

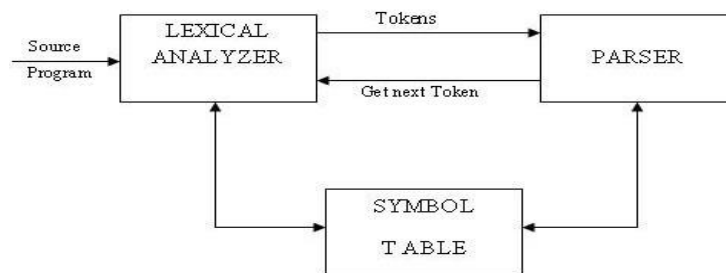


Fig. 3.1: Role of Lexical analyzer

Upon receiving a 'get next token' command from the parser, the lexical analyzer reads the input character until it can identify the next token. The LA returns to the parser representation for the token it has found. The representation will be an integer code, if the token is a simple construct such as parenthesis, comma or colon.

LA may also perform certain secondary tasks as the user interface. One such task is stripping out from the source program the comments and white spaces in the form of blank, tab and new line characters. Another is correlating error message from the compiler with the source program.

3.3 TOKEN, LEXEME, PATTERN:

Token: Token is a sequence of characters that can be treated as a single logical entity. Typical tokens are,

1) Identifiers 2) keywords 3) operators 4) special symbols 5) constants

Pattern: A set of strings in the input for which the same token is produced as output. This set of strings is described by a rule called a pattern associated with the token.

Lexeme: A lexeme is a sequence of characters in the source program that is matched by the pattern for a token.

Token	lexeme	pattern
const	const	const
if	if	If
relation	<<=, <, >, >=	< or <= or = or <> or >= or letter followed by letters & digit
i	pi	any numeric constant
nun	3.14	any character b/w "and "except"
literal	"core"	pattern

Fig. 3.2: Example of Token, Lexeme and Pattern

3.4. LEXICAL ERRORS:

Lexical errors are the errors thrown by lexer when unable to continue. Which means that there's no way to recognise a *lexeme* as a valid *token* for you lexer. Syntax errors, on the other side, will be thrown by the scanner when a given set of **already** recognised valid tokens don't match any of the right sides of your grammar rules. simple panic-mode error handling system requires that we return to a high-level parsing function when a parsing or lexical error is detected.

Error-recovery actions are:

- i. Delete one character from the remaining input.
- ii. Insert a missing character in to the remaining input.
- iii. Replace a character by another character.
- iv. Transpose two adjacent characters.

3.5. REGULAR EXPRESSIONS

Regular expression is a formula that describes a possible set of string. Component of regular expression..

- X** the character x
- .** any character, usually accept a new line
- [x y z]** any of the characters x, y, z,
- R?** a R or nothing (=optionally as R)
- R*** zero or more occurrences.....
- R+** one or more occurrences
- R1R2** an R1 followed by an R2
- R1|R1** either an R1 or an R2.

A token is either a single string or one of a collection of strings of a certain type. If we view the set of strings in each token class as an language, we can use the regular-expression notation to describe tokens.

Consider an identifier, which is defined to be a letter followed by zero or more letters or digits. In regular expression notation we would write.

Identifier = letter (letter | digit)*

Here are the rules that define the regular expression over alphabet .

- ϵ is a regular expression denoting $\{ \epsilon \}$, that is, the language containing only the empty string.
- For each 'a' in Σ , a is a regular expression denoting $\{ a \}$, the language with only one string consisting of the single symbol 'a' .
- If R and S are regular expressions, then

$(R) | (S)$ means $L(r) \cup L(s)$

$R.S$ means $L(r).L(s)$

R^* denotes $L(r^*)$

3.6. REGULAR DEFINITIONS

For notational convenience, we may wish to give names to regular expressions and to define regular expressions using these names as if they were symbols.

Identifiers are the set or string of letters and digits beginning with a letter. The following regular definition provides a precise specification for this class of string.

Example-1,

$Ab^*|cd?$ Is equivalent to $(a(b^*)) | (c(d?))$

Pascal identifier

Letter - $A | B | \dots | Z | a | b | \dots | z$

Digits - $0 | 1 | 2 | \dots | 9$

Id - $\text{letter} (\text{letter} / \text{digit})^*$

Recognition of tokens:

We learn how to express pattern 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.

Stmt \rightarrow if expr then stmt
 | If expr then else stmt
 | ϵ

Expr \rightarrow term relop term
 | term

Term \rightarrow id
 | number

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

The terminal of grammar, which are if, then , else, relop ,id and numbers are the names of tokens as far as the lexical analyzer is concerned, the patterns for the tokens are described using regular definitions.

digit $\rightarrow [0,9]$

digits $\rightarrow \text{digit}^+$

number $\rightarrow \text{digit}(\text{.digit})?(e.[+-]?digits)?$

letter \rightarrow [A-Z,a-z]
 id \rightarrow letter(letter/digit)*
 if \rightarrow if
 then \rightarrow then
 else \rightarrow else

 relop \rightarrow < | > | <= |

 >= | = | < >

In addition, we assign the lexical analyzer the job stripping out white space, by recognizing the “token” we defined by:

WS \rightarrow (blank/tab/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 parser ,but rather restart the lexical analysis from the character that follows the white space . It is the following token that gets returned to the parser.

Lexeme	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

3.7. TRANSITION DIAGRAM:

Transition Diagram has 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 .

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 one state s, and the next input symbol is a, we look for an edge out of state s labeled by a. if we find such an edge ,we advance the forward pointer and enter the state of the transition diagram to which that edge leads.

Some important conventions about transition diagrams are

1. Certain states are said to be accepting or final .These states indicates that a lexeme has been found, although the actual lexeme may not consist of all positions b/w the lexeme Begin and forward pointers we always indicate an accepting state by a double circle.
2. In addition, if it is necessary to return the forward pointer one position, then we shall additionally place a * near that accepting state.

- One state is designed the state ,or initial state ., it is indicated by an edge labeled “start” entering from nowhere .the transition diagram always begins in the state before any input symbols have been used.

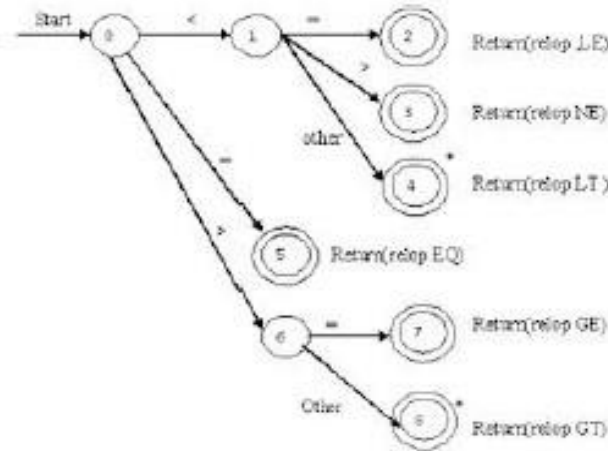


Fig. 3.3: Transition diagram of Relational operators

As an intermediate step in the construction of a LA, we first produce a stylized flowchart, called a transition diagram. Position in a transition diagram, are drawn as circles and are called as states.

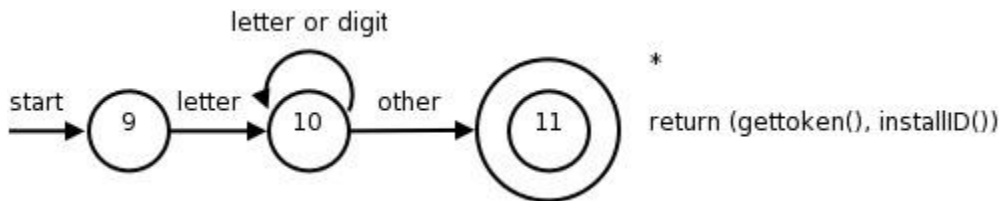


Fig. 3.4: Transition diagram of Identifier

The above TD for an identifier, defined to be a letter followed by any no of letters or digits.A sequence of transition diagram can be converted into program to look for the tokens specified by the diagrams. Each state gets a segment of code.

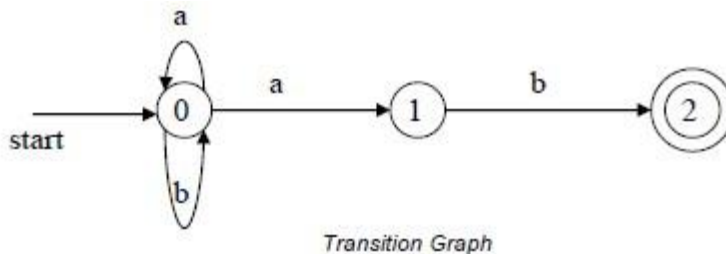
3.8. FINITE AUTOMATON

- A *recognizer* for a language is a program that takes a string x, and answers “yes” if x is a sentence of that language, and “no” otherwise.
 - We call the recognizer of the tokens as a *finite automaton*.
 - A finite automaton can be: *deterministic (DFA)* or *non-deterministic (NFA)*
 - This means that we may use a deterministic or non-deterministic automaton as a lexical analyzer.
 - Both deterministic and non-deterministic finite automaton recognize regular sets.
 - Which one?
 - deterministic – faster recognizer, but it may take more space
 - non-deterministic – slower, but it may take less space
- Deterministic automata are widely used lexical analyzers.

- First, we define regular expressions for tokens; Then we convert them into a DFA to get a lexical analyzer for our tokens.

3.9. Non-Deterministic Finite Automaton (NFA)

- A non-deterministic finite automaton (NFA) is a mathematical model that consists of:
 - S - a set of states
 - Σ - a set of input symbols (alphabet)
 - move - a transition function move to map state-symbol pairs to sets of states.
 - s_0 - a start (initial) state
 - F - a set of accepting states (final states)
 - ϵ - transitions are allowed in NFAs. In other words, we can move from one state to another one without consuming any symbol.
 - A NFA accepts a string x , if and only if there is a path from the starting state to one of accepting states such that edge labels along this path spell out x .
- Example:



0 is the start state s_0
 {2} is the set of final states F
 $\Sigma = \{a,b\}$
 $S = \{0,1,2\}$

Transition Function:

	a	b
0	{0,1}	{0}
1	\emptyset	{2}
2	\emptyset	\emptyset

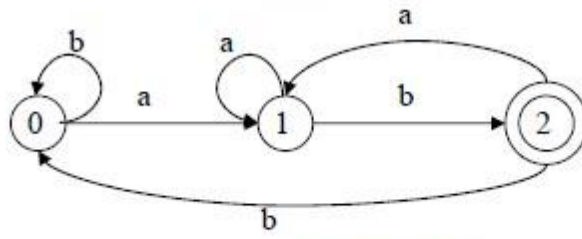
The language recognized by this NFA is $(a|b)^*ab$

3.10. Deterministic Finite Automaton (DFA)

- A Deterministic Finite Automaton (DFA) is a special form of a NFA.
- No state has ϵ - transition
- For each symbol a and state s , there is at most one labeled edge a leaving s .
 i.e. transition function is from pair of state-symbol to state (not set of states)

Example:

The DFA to recognize the language $(a|b)^* ab$ is as follows.



Transition Graph

0 is the start state s_0
 {2} is the set of final states F
 $\Sigma = \{a,b\}$
 $S = \{0,1,2\}$

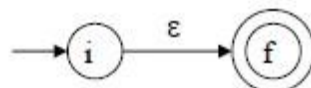
Transition Function:

	a	b
0	1	0
1	1	2
2	1	0

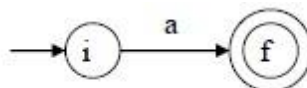
Note that the entries in this function are single value and not set of values (unlike NFA).

3.11. Converting RE to NFA

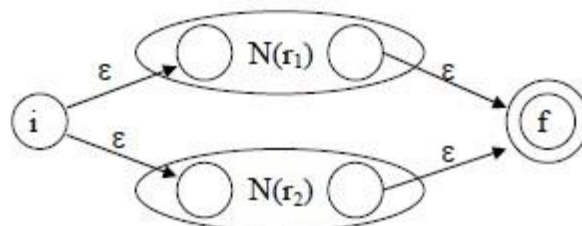
- This is one way to convert a regular expression into a NFA.
- There can be other ways (much efficient) for the conversion.
- Thomson's Construction is simple and systematic method.
- It guarantees that the resulting NFA will have exactly one final state, and one start state.
- Construction starts from simplest parts (alphabet symbols).
- To create a NFA for a complex regular expression, NFAs of its sub-expressions are combined to create its NFA.
- To recognize an empty string ϵ :



- To recognize a symbol a in the alphabet Σ :

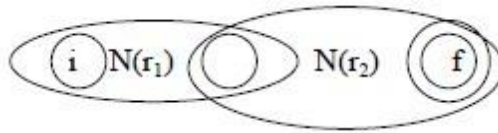


- For regular expression $r_1 | r_2$:

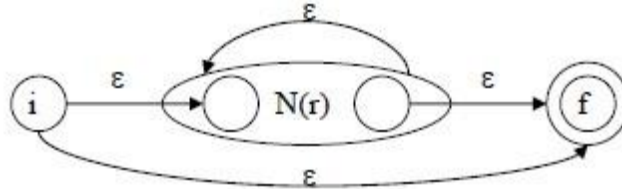


$N(r_1)$ and $N(r_2)$ are NFAs for regular expressions r_1 and r_2 .

- For regular expression $r_1 r_2$

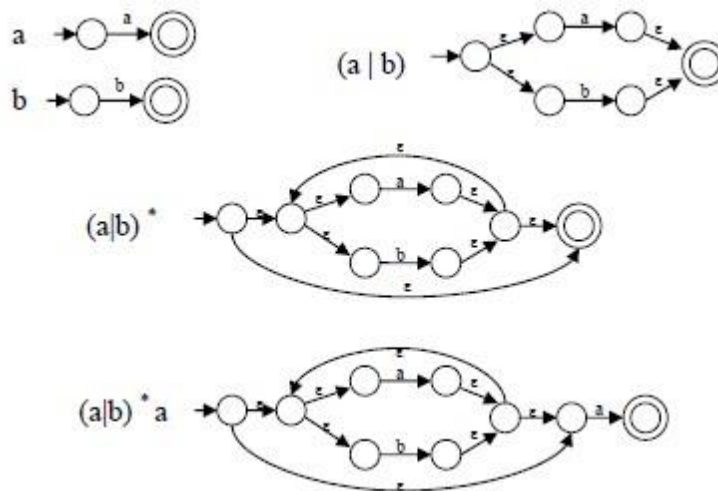


Here, final state of $N(r_1)$ becomes the final state of $N(r_1r_2)$. • For regular expression r^*



Example:

For a RE $(a|b)^* a$, the NFA construction is shown below.



3.12. Converting NFA to DFA (Subset Construction)

We merge together NFA states by looking at them from the point of view of the input characters:

- From the point of view of the input, any two states that are connected by an ϵ -transition may as well be the same, since we can move from one to the other without consuming any character. Thus states which are connected by an ϵ -transition will be represented by the same states in the DFA.
- If it is possible to have multiple transitions based on the same symbol, then we can regard a transition on a symbol as moving from a state to a set of states (ie. the union of all those states reachable by a transition on the current symbol). Thus these states will be combined into a single DFA state.

To perform this operation, let us define two functions:

- The **-closure** function takes a state and returns the set of states reachable from it based on (one or more) ϵ -transitions. Note that this will always include the state itself. We should be able to get from a state to any state in its ϵ -closure without consuming any input.
- The function **move** takes a state and a character, and returns the set of states reachable by one transition on this character.

We can generalise both these functions to apply to sets of states by taking the union of the application to individual states.

For Example, if A, B and C are states, $\text{move}(\{A,B,C\},`a') = \text{move}(A,`a') \cup \text{move}(B,`a') \cup \text{move}(C,`a')$.

The Subset Construction Algorithm is as follows:

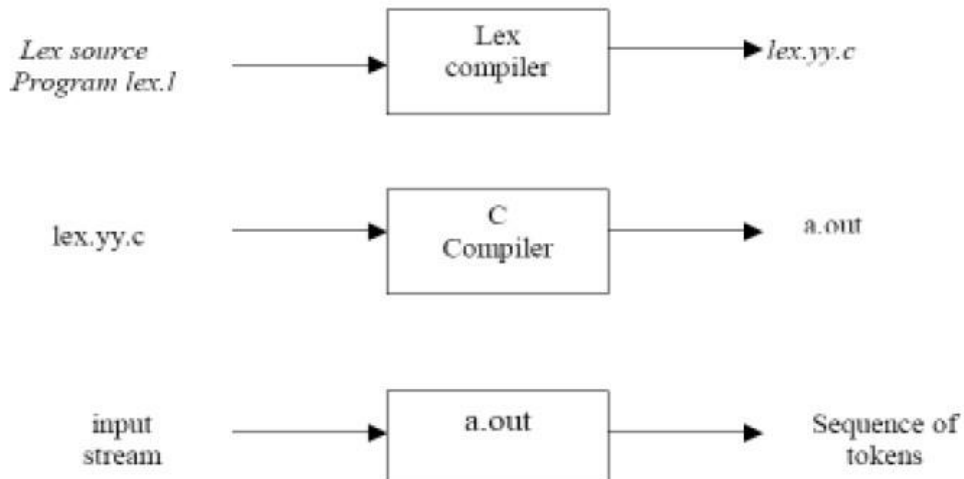
```

put  $\epsilon$ -closure( $\{s_0\}$ ) as an unmarked state into the set
of DFA (DS) while (there is one unmarked S1 in DS)
do begin mark S1
    for each input
        symbol a do
            begin
                 $S_2 \leftarrow \epsilon$ -
                closure( $\text{move}(S_1,a)$ ) if
                ( $S_2$  is not in DS) then
                add  $S_2$  into DS as an
                unmarked state
                 $\text{transfunc}[S_1,a] \leftarrow S_2$ 
            end
        end
    end
end

```

- a state S in DS is an accepting state of DFA if a state in S is an accepting state of NFA
- the start state of DFA is ϵ -closure($\{s_0\}$)

3.13. Lexical Analyzer Generator



3.18. Lex specifications:

A Lex program (the .l file) consists of three parts:

declarations

%%

translation rules

%%

auxiliary procedures

1. The *declarations* section includes declarations of variables, manifest constants (A manifest constant is an identifier that is declared to represent a constant e.g. # *define PIE 3.14*), and regular definitions.
2. The *translation rules* of a Lex program are statements of the form :

```
p1
{action1}
p2
{action2}
p3
{action 3}
```

```
... ..
... ..
```

Where, each *p* is a regular expression and each *action* is a program fragment describing what action the lexical analyzer should take when a pattern *p* matches a lexeme. In Lex the actions are written in C.

3. The third section holds whatever *auxiliary procedures* are needed by the *actions*. Alternatively these procedures can be compiled separately and loaded with the lexical analyzer.