

1) Role of lexical analysis and its issues.

- ✓ The lexical analyzer is the first phase of compiler.
- ✓ Its main task is to read the input characters and produce as output a sequence of tokens that the parser uses for syntax analysis.
- ✓ It may also perform secondary task at user interface.
- ✓ One such task stripping out from the source program comments and white space in the form of blanks, tab, and newline characters.
- ✓ Some lexical analyzer are divided into cascade of two phases, the first called scanning and second is "lexical analysis".
- ✓ The scanner is responsible for doing simple task while lexical analysis does the more complex task.

Issues in Lexical Analysis:

There are several reason for separating the analysis phase of compiling into lexical analysis and parsing:

- ✓ Simpler design is perhaps the most important consideration. The separation of lexical analysis often allows us to simplify one or other of these phases.
- ✓ Compiler efficiency is improved.
- ✓ Compiler portability is enhanced.

2) Explain token, pattern and lexemes.

Token: Sequence of character having a collective meaning is known as *token*.

Typical tokens are,

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

Lexeme: The character sequence forming a token is called *lexeme* for token.

Pattern: The set of rules by which set of string associate with single token is known as *pattern*

Token	Lexeme	Pattern
id	x y n0	letter followed by letters and digits
number	3.14159, 0, 6.02e23	any numeric constant
If	If	if
relation	<, <=, =, >, >=, >	< or <= or = or > or >= or letter followed by letters & digit
Literal	"abc xyz"	anything but ", surrounded by " 's

1. **if(x<=5)** :Token – if (keyword),
X (id),
<= (relation),
5 (number)

Lexeme - if , x , <= , 5

2. total = sum + 12.5

Token – total (id),

= (relation),

Sum (id),

+ (operator)

12.5 (num)

Lexeme – total, =, sum, +, 12.5

3) What is input buffering? Explain technique of buffer pair.

The speed of lexical analysis is a concern in compiler design.

1. Buffer pair:

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.

Two pointers to the input are maintained:

1. Pointer *Lexeme Begin*, 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.

code to advance forward pointer is given below

if forward at end of first half then begin

reload second half;

forward := forward + 1

end

else if forward at end of second half then begin

reload first half;

move forward to beginning of first half

end

else forward := forward + 1;

Once the next lexeme is determined, *forward* is set to character at its right end. Then, after the lexeme is recorded as an attribute value of a token returned to the parser, *Lexeme Begin* is set to the character immediately after the lexeme just found.

2. Sentinels:

If we use the scheme of Buffer pairs 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**. Look ahead code with sentinels is given below:

```
forward := forward + 1;
if forward = eof then begin
    if forward at end of first half then begin
        reload second half;
        forward := forward + 1
    end
else if forward at the second half then begin
    reload first half;
    move forward to beginning of first half
end
else terminate lexical analysis
end;
```

4) Regular expression & regular definition

Regular Expression

- 0 or 1
 $0+1$
- 0 or 11 or 111
 $0+11+111$
- Regular expression over $\Sigma=\{a,b,c\}$ that represent all string of length 3.
 $(a+b+c)(a+b+c)(a+b+c)$
- String having zero or more a.
 a^*
- String having one or more a.
 a^+
- All binary string.
 $(0+1)^*$
- 0 or more occurrence of either a or b or both
 $(a+b)^*$
- 1 or more occurrence of either a or b or both
 $(a+b)^+$
- Binary no. end with 0
 $(0+1)^*0$
- Binary no. end with 1
 $(0+1)^*1$
- Binary no. starts and end with 1.

- 1(0+1)*1
12. String starts and ends with same character.
 $0(0+1)^*0$ or $a(a+b)^*a$
 $1(0+1)^*1$ $b(a+b)^*b$
13. All string of a and b starting with a
 $a(a/b)^*$
14. String of 0 and 1 end with 00.
 $(0+1)^*00$
15. String end with abb.
 $(a+b)^*abb$
16. String start with 1 and end with 0.
 $1(0+1)^*0$
17. All binary string with at least 3 characters and 3rd character should be zero.
 $(0+1)(0+1)0(0+1)^*$
18. Language which consist of exactly Two b's over the set $\Sigma = \{a, b, \}$
 $a^*ba^*ba^*$
19. $\Sigma = \{a, b\}$ such that 3rd character from right end of the string is always a.
 $(a/b)^*a(a/b)(a/b)$
20. Any no. of a followed by any no. of b followed by any no of c.
 $a^*b^*c^*$
21. It should contain at least 3 one.
 $(0+1)^*1(0+1)^*1(0+1)^*1(0+1)^*$
22. String should contain exactly Two 1's
 $0^*10^*10^*$
23. Length should be at least be 1 and at most 3.
 $(0+1)^+ (0+1)^+ (0+1)^+ (0+1)^+ (0+1)^+$
24. No.of zero should be multiple of 3
 $(1^*01^*01^*01^*)^*+1^*$
25. $\Sigma = \{a, b, c\}$ where a are multiple of 3.
 $((b+c)^*a (b+c)^*a (b+c)^*a (b+c)^*)^*$
26. Even no. of 0.
 $(1^*01^*01^*)^*$
27. Odd no. of 1.
 $0^*(10^*10^*)^*10^*$
28. String should have odd length.
 $(0+1)((0+1)(0+1))^*$
29. String should have even length.
 $((0+1)(0+1))^*$
30. String start with 0 and has odd length.
 $0((0+1)(0+1))^*$
31. String start with 1 and has even length.

- 1(0+1)((0+1)(0+1))*
32. Even no of 1
(0*10*10*)*
33. String of length 6 or less
(0+1+[^])⁶
34. String ending with 1 and not contain 00.
(1+01)⁺
35. All string begins or ends with 00 or 11.
(00+11)(0+1)*+(0+1)*(00+11)
36. All string not contains the substring 00.
(1+01)* (^+0)
37. Every 0 is followed immediately by 1.
1*(011*)*
38. Language of all string containing both 11 and 00 as substring.
((0+1)*00(0+1)*11(0+1)*)+ ((0+1)*11(0+1)*00(0+1)*)
39. Language of C identifier.
(_+L)(_+L+D)*
40. The set of all string not containing 101 as a substring.
0*(1+000*)*0*
41. The language of all 0's and 1's having 1 at odd position.
(1(0+1)*)(1+[^])

Regular Definition

A regular definition gives names to certain regular expressions and uses those names in other regular expressions.

Here is a regular definition for the set of Pascal identifiers that is defined as the set of strings of letters and digits beginning with a letter.

letter $\rightarrow A \mid B \mid \dots \mid Z \mid a \mid b \mid \dots \mid z$

digit $\rightarrow 0 \mid 1 \mid 2 \mid \dots \mid 9$

id $\rightarrow \text{letter} (\text{letter} \mid \text{digit})^*$

The regular expression id is the pattern for the Pascal identifier token and defines letter and digit.

Where letter is a regular expression for the set of all upper-case and lower case letters in the alphabet and digit is the regular for the set of all decimal digits.

5) Explain Finite automata (NFA & DFA)

- It is a mathematical model- state transition diagram
- It is a Recognizer for a given language
- 5-tuple $\{Q, \Sigma, \delta, q_0, F\}$
 - Q is a finite set of states
 - Σ is a finite set of input
 - δ transition function $Q \times \Sigma$

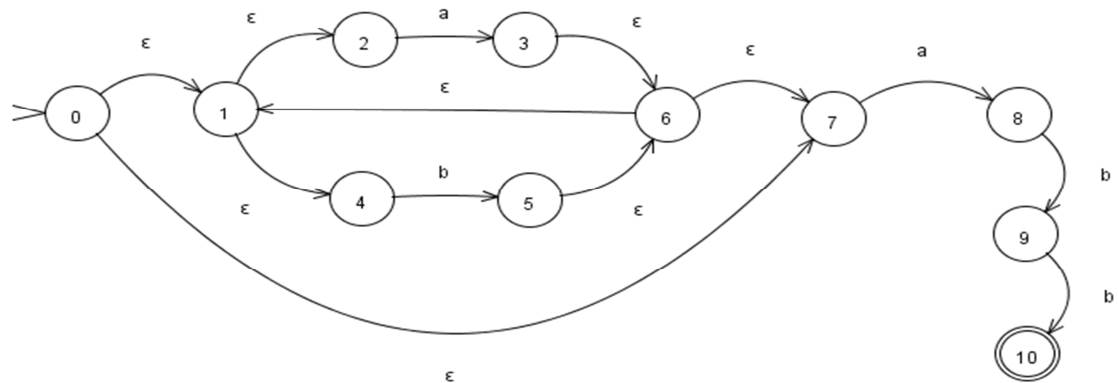
- q_0 , F initial and final state respectively

There are mainly two types of finite automata, Non deterministic automata and Deterministic automata. Difference between NFA and DFA.

Parameter	NFA	FA
Transition	Non deterministic	Deterministic
No. of states	NFA has a fewer number of states.	More, if NFA contains Q states then the corresponding FA will have $\leq 2^Q$ states.
Power	NFA is as powerful as DFA	FA is powerful as an NFA
Design	Easy to design due to non-determinism	More difficult to design
Acceptance	It is difficult to find whether $w \in L$ as there are several paths. Backtracking is required to explore several parallel paths.	It is easy to find whether $w \in L$ as transition are deterministic.

6) Conversion from NFA to DFA using Thompson's rule

Example: $(a/b)^*abb$



➤ ϵ - closure (0) = {0,1,2,4,7} ---- Let A

➤ Move(A,a) = {3,8}

ϵ - closure (Move(A,a)) = {1,2,3,4,6,7,8} ---- Let B

Move(A,b) = {5}

ϵ - closure (Move(A,b)) = {1,2,4,5,6,7} ---- Let C

➤ $\text{Move}(B,a) = \{3,8\}$

ϵ - closure ($\text{Move}(B,a)$) = $\{1,2,3,4,6,7,8\}$ ----- B

$\text{Move}(B,b) = \{5,9\}$

ϵ - closure ($\text{Move}(B,b)$) = $\{1,2,4,5,6,7,9\}$ ----- Let D

➤ $\text{Move}(C,a) = \{3,8\}$

ϵ - closure ($\text{Move}(C,a)$) = $\{1,2,3,4,6,7,8\}$ ----- B

$\text{Move}(C,b) = \{5\}$

ϵ - closure ($\text{Move}(C,b)$) = $\{1,2,4,5,6,7\}$ ----- C

➤ $\text{Move}(D,a) = \{3,8\}$

ϵ - closure ($\text{Move}(D,a)$) = $\{1,2,3,4,6,7,8\}$ ----- B

$\text{Move}(D,b) = \{5,10\}$

ϵ - closure ($\text{Move}(D,b)$) = $\{1,2,4,5,6,7,10\}$ ----- Let E

➤ $\text{Move}(E,a) = \{3,8\}$

ϵ - closure ($\text{Move}(E,a)$) = $\{1,2,3,4,6,7,8\}$ ----- B

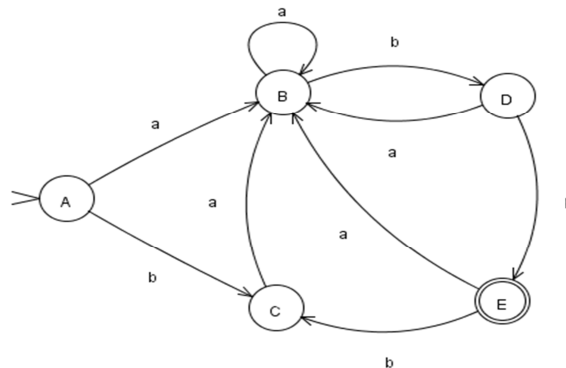
$\text{Move}(E,b) = \{5\}$

ϵ - closure ($\text{Move}(E,b)$) = $\{1,2,4,5,6,7\}$ ----- C

Transition Table:

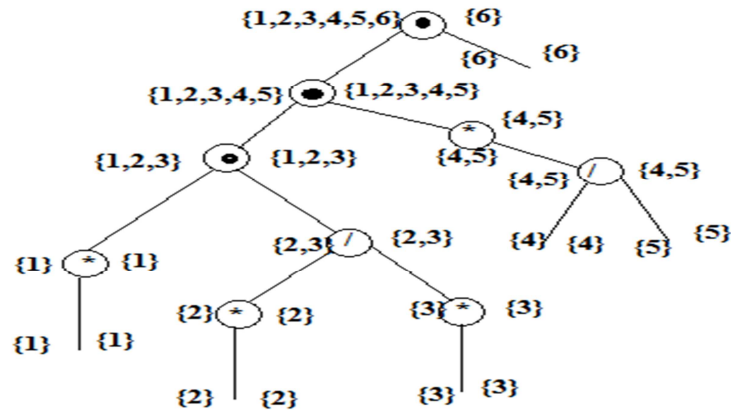
	a	b
A	B	C
B	B	D
C	B	C
D	B	E
E	B	C

DFA:



7) Conversion from Regular Expression to DFA without constructing NFA

Example: $a^*(b^*/c^*)(a/c)^*$



➤ For root node

$i = \text{lastpos}(C_1) = \{1,2,3,4,5\}$

$\text{followpos}(i) = \text{firstpos}(C_2) = \{6\}$

$\text{followpos}(1) = \{6\}$

$\text{followpos}(2) = \{6\}$

$\text{followpos}(3) = \{6\}$

$\text{followpos}(4) = \{6\}$

$\text{followpos}(5) = \{6\}$

$i = \text{lastpos}(C_1) = \{1,2,3\}$

$\text{followpos}(i) = \text{firstpos}(C_2) = \{4,5\}$

$\text{followpos}(1) = \{4,5\}$

$\text{followpos}(2) = \{4,5\}$

$\text{followpos}(3) = \{4,5\}$

$i = \text{lastpos}(n) = \{4,5\}$

$\text{followpos}(i) = \text{firstpos}(n) = \{4,5\}$

$\text{followpos}(4) = \{4,5\}$

$\text{followpos}(5) = \{4,5\}$

$i = \text{lastpos}(C_1) = \{1\}$

$\text{followpos}(i) = \text{firstpos}(C_2) = \{2,3\}$

$\text{followpos}(1) = \{2,3\}$

$i = \text{lastpos}(n) = \{2\}$
 $\text{followpos}(i) = \text{firstpos}(n) = \{2\}$
 $\text{followpos}(2) = \{2\}$

$i = \text{lastpos}(n) = \{3\}$
 $\text{followpos}(i) = \text{firstpos}(n) = \{3\}$
 $\text{followpos}(3) = \{3\}$

i	Followpos(i)
1	$\{1,2,3,4,5,6\}$
2	$\{2,4,5,6\}$
3	$\{3,4,5,6\}$
4	$\{4,5,6\}$
5	$\{4,5,6\}$

Construct DFA:

Initial node = firstpos (root node)
 $= \{1,2,3,4,5,6\}$

$\delta((1,2,3,4,5,6), a) = \text{followpos}(1) \cup \text{followpos}(4)$
 $= \{1,2,3,4,5,6\}$

$\delta((1,2,3,4,5,6), b) = \text{followpos}(2)$
 $= \{2,4,5,6\}$

$\delta((1,2,3,4,5,6), c) = \text{followpos}(3) \cup \text{followpos}(5)$
 $= \{3,4,5,6\}$

$\delta((2,4,5,6), a) = \text{followpos}(4)$
 $= \{4,5,6\}$

$\delta((2,4,5,6), b) = \text{followpos}(2)$
 $= \{2,4,5,6\}$

$\delta((2,4,5,6), c) = \text{followpos}(5)$
 $= \{4,5,6\}$

$\delta((3,4,5,6), a) = \text{followpos}(4)$
 $= \{4,5,6\}$

$\delta((3,4,5,6), b) = \Phi$

$\delta((3,4,5,6), c) = \text{followpos}(3) \cup \text{followpos}(5)$
 $= \{3,4,5,6\}$

$\delta((4,5,6), a) = \text{followpos}(4)$
 $= \{4,5,6\}$

$\delta((4,5,6), b) = \Phi$

$\delta((4,5,6), c) = \text{followpos}(5)$
 $= \{4,5,6\}$

