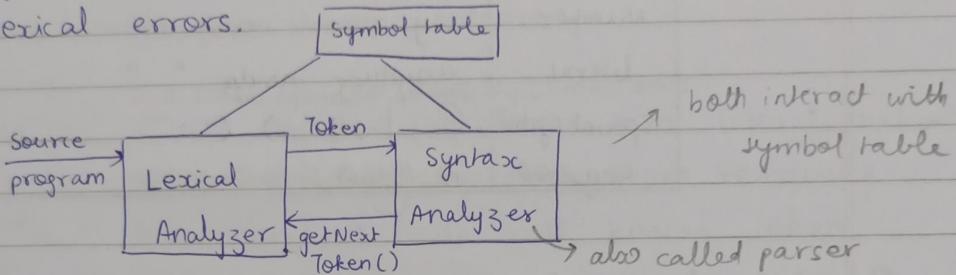


CH. 3 LEXICAL ANALYSER (SCANNER)

- 1st phase of compiler
- reads input (source program)
- identifies lexemes
- generates tokens
- creates symbol table for entries for identifiers (variables)
- detects lexical errors.



i.e. reads syntax analyzer invokes lexical analyzer, only then lexical analyzer reads a set of characters and identifies the next lexeme, converts it to token.

Once syntax analyzer processes a token it asks for next token, only then lexical analyzer reads next set of chars.

- Discards comments & whitespaces (tab, blank, \n)
- Correlates error message with source program
(i.e. identifies which position error takes place)
- Lexical analyser is divided into 2 parts
 - Scanning
 - Tokenization

→ Token is a pair consisting of a token name and optional attribute value.

→ Pattern is a description of the form that the lexeme of a token may take (i.e. possible values for token)
They are represented using regular expressions.

• Expansion macros
 - (#define etc.)
 • preprocessing directives
 • discarding whitespace

eg:

key words ↗ ↗

Tokens	Patterns
if	‘;’ followed by ‘f’
else	‘e’ followed by ‘l’ followed by ‘s’ followed by ‘e’
comparison	<, >, <=, >=, ==, !=
id	letter / underscore followed by letters / digits
Number	any numeric constant
literal	anything inside “ ”
special symbol, @, [, ;,), (common also special
reserved keyword	int, float, const etc.

- Q) write lexeme & token for the following C instruction
 printf ("Total = %d\n", score);

Lexeme	Token
printf	id
(LCB or LP or (→ left common bracket
"Total = %d\n"	literal
,	comma or ,
score	id
)	RCB or RP or) → left parenthesis
;	SC → semicolon

LEXICAL ERROR

- Lexical Errors are detected when none of the patterns listed for the given language ^{are} matched.

eg: int 1abc ;

here, 1abc is wrong but it won't be considered a lexical error.

lexeme	token
int	< int >
1	< num, 1 >
abc	< id >
;	< special symbol >

1abc is a syntax error.

- Any language other than english is a lexical error.
eg: int ~~gfg~~;

LEXICAL ERROR RECOVERY TECHNIQUES

- It recovers using panic mode recovery.
- Panic Mode - if it finds any not matching (listed) chars
↳ basically deletes successive characters from remaining input unit it finds well known token

→ Other recovery strategies are:

- ① Delete 1 char from remaining input
 - ② insert missing char to remaining input
 - ③ Replacing char by another char
 - ④ Transferring this char to the next char, abc becomes ab1c.
- how is this lexical error?

INPUT BUFFERS

- Instead of reading each char from hardware again & again we use input buffers to store input (source) code.
- This helps with faster retrieval and processing.
- After reading a set of characters we still need to look ahead to understand if it is a keyword ~~or~~, variable or something else.

eg: int a;
float interested;

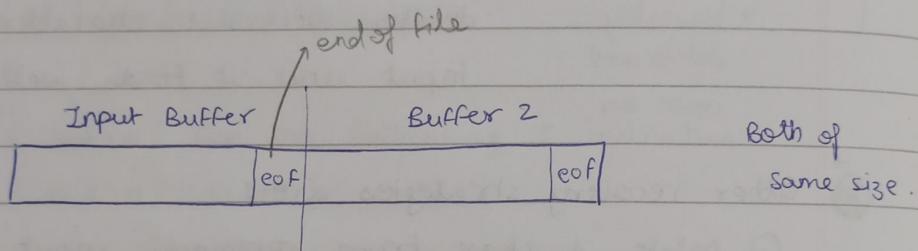
here after reading "int" of interested we still need to look ahead since it is a variable else it will be misunderstood

eg: if ($a < b$)
if ($a \leq b$)

in second statement we need to look ahead after $<$
as it is \leq

- Input buffers help with making lookahead easier.

Buffer Pairs



- In single buffer scheme once the buffer is full, tokens are generated for it. Once that is done and eof is encountered, the whole buffer is reloaded with the remaining input data that could not be stored earlier. Now tokens are generated for the remaining data.
- Due to this we face issues with lookahead and wrong tokens may be generated. (i.e wrong lexeme identified)

eg: F|l|o|a|t|i|n|t|eof float interest;

$\langle \text{float} \rangle \langle \text{int} \rangle$

after eof it is flushed and rest of data produced

e|r|e|s|t|eof,

$\langle \text{id}, 1 \rangle$

it thinks rest is a variable

∴ Wrong tokens $\langle \text{float} \rangle \langle \text{int} \rangle \langle \text{id}, 1 \rangle$

but correct was $\langle \text{float} \rangle \langle \text{id}, 1 \rangle$ where id1 = interest.

- ∵ whenever EOF encountered buffer is flushed and ~~new~~ new remaining data is loaded.
- To solve this problem we have 2 buffers.
- In buffer pairs alternately two buffers are loaded. When EOF of 1st buffer encountered, 2nd buffer is loaded.
- EOF is a sentinel used to denote end of buffer. To differentiate b/w ~~end~~ end of input and end of buffer (since both denoted by EOF) we tally no. of bytes read with no. of ~~bytes~~ bytes scanned so far. If scanned = read then end of input scanned < read then end of buffer
- Two pointers are used to maintain buffer. i.e.
lexeme begin & forward pointer.
- Lexeme begin marks the beginning of the current lexeme.
- Forward pointer scans ahead until a pattern match is found.

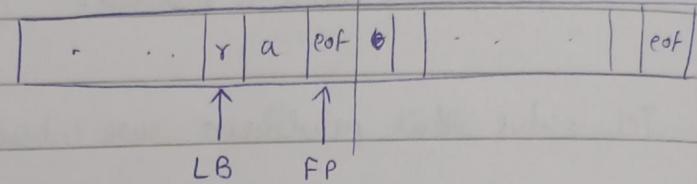
eg) crate EL pos * * *

- ^{lexeme} patterns are denoted by regular expressions. Each regex has a finite automata (FA)
when a char is read it is also put into the FA for each of these regex by lexical analyser.
- As each char is read it might get rejected by some FA while it ~~is~~ is still being read by other FAs.
- This goes on until all except one ~~reject~~ FA reject the input. And that final FA also accepts the lexeme.
(if no FA accepts it then it is a lexical error)

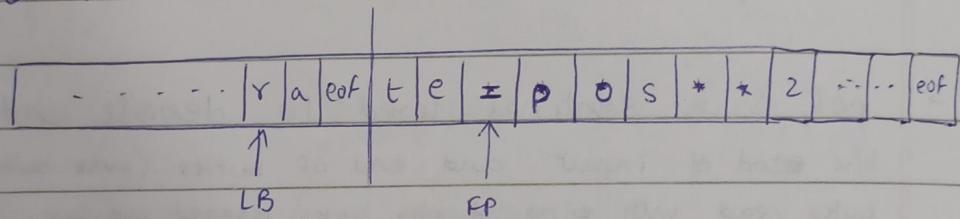


→ eg: rate = pos ** 2

initially lets say,



FP keeps reading each character. On reaching eof
the second buffer is loaded.



When FP reads = it retracts to previous index containing
'e' & matches 'rate' to identifier pattern.

Note: Retraction is not required for single characters like (,), comma, ;, etc. because it is obvious that they don't need lookahead.

REGULAR DEFINITION

regular definition

letter → A | B | ... | z | a | b | ... | z | - underline

digit → 0 | 1 | ... | 9 λ or ε → null

id → letter (letter | digit) * written as superscript { * → one or more occurrence}

RegEx : [A-Z a-z_][A-Z a-z_0-9]* written as normal sized { + → 1 or more occurrence

for identifier ? → 0 or 1 occurrence

Q Write regular definition that matches all of these, minus 4
 5280, 0.01234, 6.336E4, 5E5, 1.89E-4

sol

digit $\rightarrow [0-9]$

or digit digit* both are same

digits \rightarrow digit⁺

~~fraction~~ \rightarrow digits . digits

optional fraction \rightarrow . digits | E

optional Exponent \rightarrow E (+|-|E) digits | E

number \rightarrow digits optional fraction optional exponent

null is used
in optional frac

& optional exp

as you may

or may

not have

decimal part or
exponent

\therefore number accepts all above eg's.

The above can be written ^{regular} using extensions in a shorter way

digits $\rightarrow [0-9]^+$

0 or 1 occurrence

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

RECOGNITION OF TOKENS

Grammar

i) Grammar for if statement

only if

if else

stmt \rightarrow if else if statement

stmt \rightarrow if expr then stmt | if expr then stmt else stmt

| e

equal

relational operator like <, >, <=, >=, =, <

expr \rightarrow term relop term | term

in FORTRAN

\neq is written as \neq

not equal

term \rightarrow id | num

\neq " " " \neq "

atomic i.e. don't
have production rule

Terminals: id, num, if, then, else, relop

Non-Terminals: stmt, expr, term

E is not considered
as terminal here
as it is not
considered for tokens

- every terminal ^{we write} ~~is~~ regular definition
- Every compiler is given grammar of the language.
The terminals in this grammar are used as token names.

→ every one

Regular Definition

→ using ~~same~~ eg: we wrote

- Writing regular definition for ^{all} terminals

$$\text{eg: digit} \rightarrow [0-9]$$

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

$$\text{number} \rightarrow \text{digits} \cdot (\text{digits})? (E (+-)?) \text{ digits}?)$$

$$\text{letter} \rightarrow [A-Z] [a-z]$$

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

$$\text{relOp} \rightarrow < | > | <= | >= | = | <> \text{ } \xrightarrow{\text{not equal}}$$

$$\text{if} \rightarrow \text{i followed by f}$$

$$\text{else} \rightarrow \text{e followed by l followed by s followed by e}$$

$$\text{white space} \rightarrow (\text{blank} | \text{tab} | \text{newline})^+$$

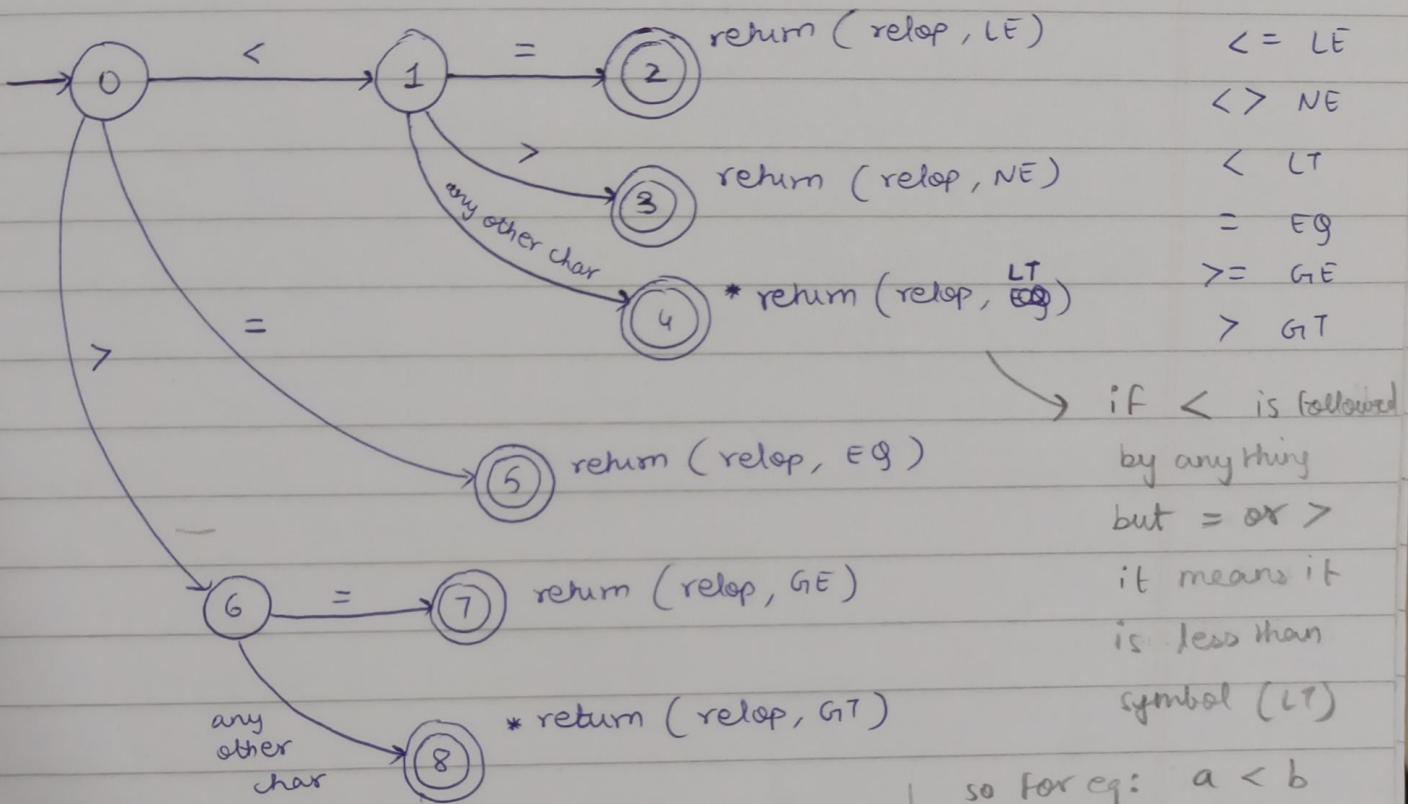
When white space is encountered, token is generated by lexical analyzer but not sent to parser.

TRANSITION DIAGRAM

- Similar to finite Automata (FA)
- ~~so~~ Every accepting state returns a token. (FA does not return anything)
- Whenever retract operation occurs it must be represented using *
- accepting state denoted by  return (<tokenname, attribute>)

TRANSITION DIAGRAM FOR RELOP

- Relop transition diagram is only invoked if LB (lexeme begin) is at <, > or =



TRANSITION DIAGRAM FOR NUMBER

check slides

& some

more

transition

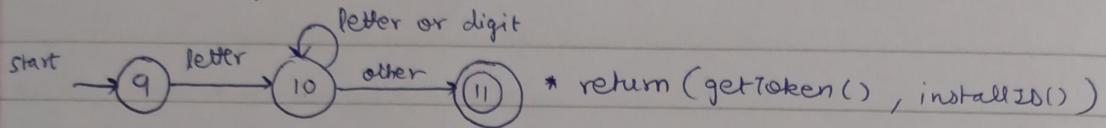
diagrams

in slides

RECOGNITION OF RESERVED WORDS & IDENTIFIERS

(transition diagram for id's & keywords)

case 1:



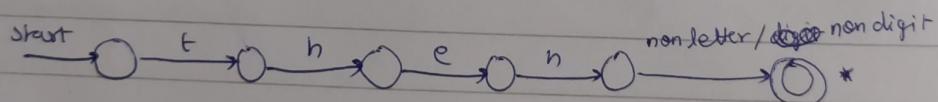
→ Only when an identifier is encountered for the first time it is added to symbol table.

If it is encountered again, it is found in symbol table and its type is checked. (if keyword then keyword is returned else, ^{if identifier then} serial no. of identifier returned)

case 2:

→ Another method is to create individual transition diagrams for each keyword.

eg: then is a keyword



here after ~~then~~ ^{now} after "then" if ~~anything else~~ ^{any} non-letter or nondigit is encountered it means keyword has ended. and we backtrack that one symbol and accept "then"

eg: then 1

upon reaching 1 it realises keyword has ended and then is accepted

eg: thence

here after "then" c is not non-letter / non-digit and DFA will not accept this string as then was just a substring in thence not whole word.

ARCHITECTURE OF A TRANSITION DIAGRAM BASED LEXICAL ANALYZER

Implementation of relop transition diagram

→ drawn earlier refer to it for states

TOKEN getRelop()

{

TOKEN retoken = new (RELOP);

while(1) { /* repeat char processing until a return or failure occurs */

switch (state) {

case 0: c = nextChar();

start state

if (c == '<') state = 1;

elseif (c == '=') state = 5;

else if (c == '>') state = 6;

else fail(); /* retracts the chars read during this getRelop() fn as the lexeme was not a relop */

break;

→ e.g. in C++

cout <<

here, < means getRelop() is

called but actually it is not a relop so fail().

case 2, 3, 4, 5, 7, 8
are accepting

case 1: ..

..

case 8: retract();

retoken.attribute = GT;

return (retoken);

States
acc. to transition
diagram
accepting
state

}

}

}

→ Code not complete we must write on our own.

case 1: if (c == '<')



case 1 : $c = \text{nextChar}();$
if ($c == '='$) state = 2;
else if ($c == '>'$) state = 3;
else {
~~if (c <= 'z')~~
~~retract();~~
~~ret~~
else state = 4;
break;

case 2 : ~~@@@~~ retToken.attribute = LE;
return (retToken);

case 3 : retToken.attribute = NE;
return (retToken);

case 4 : retract();
retToken.attribute = LT;
return (retToken);

case 5 : retToken.attribute = EQ;
return (retToken);

case 6 : $c = \text{nextChar}();$
if ($c == '='$) state = 7;
else state = 8;
break;

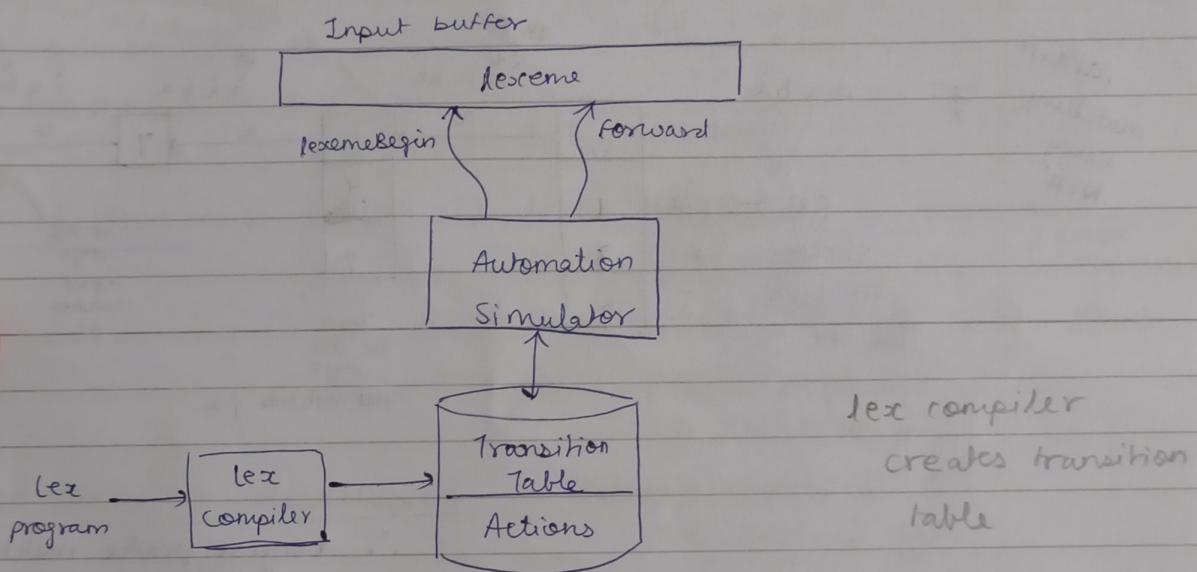
case 7 : retToken.attribute = GF;
return (retToken);

FINITE AUTOMATA

- Doesn't return anything
- 2 types: NFA & DFA

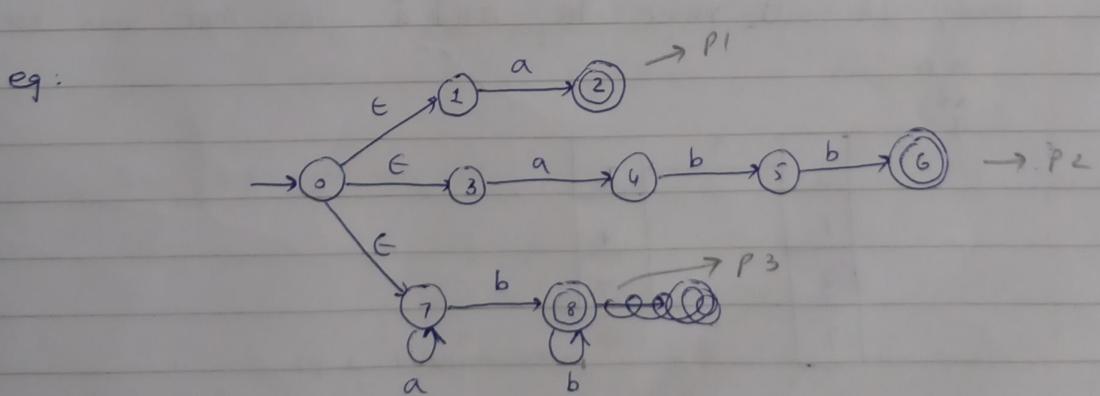
lexical analyzer

DESIGN OF LA generator using FA



- We will only use NFA

↳ \in Transitions (check FAL)



If here, abb is accepted by both p_2 & p_3

This is a conflict.

If also out of abb, when a is read, that matches p_1 even before reading bb.

To solve problem ii)

- To prevent conflict we match the longest prefix during pattern matching.

To solve problem i)

- Also if multiple patterns match then we choose the first pattern that matches.

∴ b/w p_2 & p_3 , p_2 will be chosen.

pattern
matching
using
NFA

eg: aaba

precess

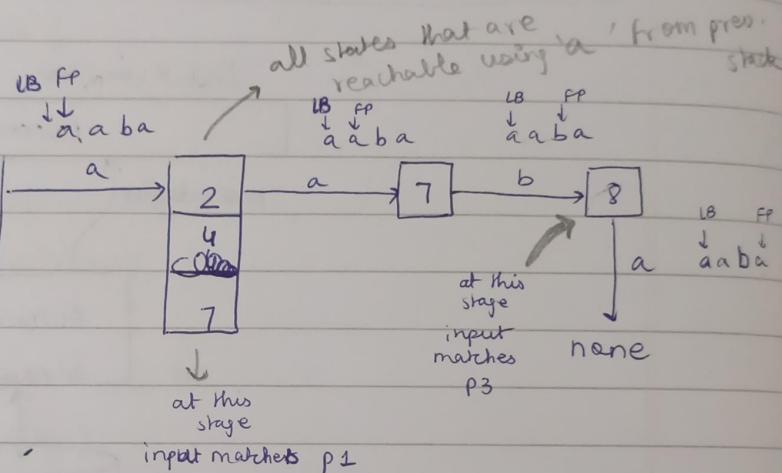
stack →
with states

This one is
epsilon closure
of start state 0

i.e all states reachable

using ϵ from state 0

$$\epsilon\text{-closure}(0) = \{0, 1, 3, 7\}$$

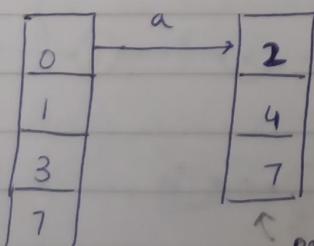


∴ Longest prefix matched is ~~aab~~ aab
pattern for that is p_3 .

∴ lexeme = aab

now, LB will reset to $aab\downarrow a$ as that was not matched.

LB
FP
↓
a



∴ 'a' is another lexeme.