

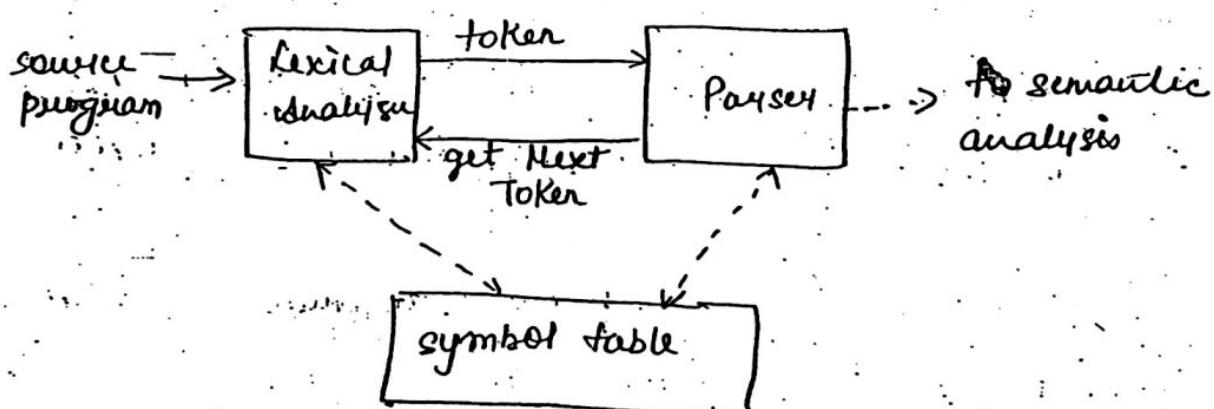
Lexical Analyser (Scanner) -

It is the first phase of a compiler. Its main task is to read the input characters and produce a sequence of tokens that the parser (next phase) uses for syntax analysis.

Interaction b/w Lexical Analyser and Parser

Whenever the parser need a token it calls the lexical analyser by "get next token" command, the lexical analyser then reads characters until it can identify the next token.

Instead of analyzing the entire input string, the lexical analyser sees enough of the IP string to return a single token.



Interaction b/w Grammer and Parser

Task of a scanner

- 1) Recognizing reserved words (Keywords)
- 2) Finding identifiers (variables)
- 3) Finding string and character constants
- 4) Ignoring comments and white spaces.
- 5) Counting the number of lines.
- 6) Reporting error messages.
- 7) Correlating error messages generated by the compiler with the source program.
eg. associating line number with error messages.

Some Useful Terms

- 1) Pattern - Pattern is a rule which describes a set of lexemes that can represent a particular token in the source program.
for eg. identifier can be described as a letter followed by letters or digits.
- 2) Lexeme - Lexemes are the smallest logical units of a program. It is a sequence of characters in the source program from which a token is produced.
eg. - 10.0, +, etc.

(3)

Tokens - classes of lexeme
identified by the basic tokens.

e.g. identifier, keyword etc.

following constructs are treated as tokens.

- keywords
- operators
- identifiers
- constants
- punctuation symbols etc.

each token is pair consisting of two values:
{ token name, attribute value }

Attribute for token - more than one lexeme can match a same pattern. So the lexical analyzer must provide additional ~~value~~ information about that lexeme. This attribute is a pointer to the lexeme in the symbol table entry.

for operators, punctuation, keywords there is no need for an attribute value. The first component itself is sufficient to identify the lexeme.

Sample lexemes	Tokens	Description of patterns
pi, count	id	letters followed by letters or digits
3.14, 286	num	any numeric constant,
"made easy"	literal	characters b/w " "
if	if	if

ex of pattern, lexeme & token

ex: sum_two (float num1, float num2)

```
1 float sum ;
sum = num1 + num2 ;
return (sum);
```

Tokens	Lexeme
Keywords	float, return
Identifiers	sum_two, num1, num2, sum
Delimiters	, ;, (,), {, }
Assign op	=
Addop	+

add +
m. ref.
<asspp>
<const>

Scanning the Input

lexical analysis needs the semantic profile and
character by character. The reading of characters
from records storage is very costly.

So, block of data is first read into buffer and then scanned by the lexical analyzer.

Two types of buffering are

1) The buffer scheme - uses extra link buffer scheme, but it has some problems, let the message is longer than the length of a buffer.

2) Two buffer scheme - uses two buffers of same size that are scanned alternately.

lexical analysis uses two pointers

lb - scene beginning painter

(look-ahead) fp - forward pointer - to keep track of the pen

When the token is identified the token and the attribute of the token are returned. Both pointers are then set to the beginning of the next token.

e.g. $\text{sum} = \text{num1} + \text{num2}$
return (sum);

$\boxed{\text{sum} = \text{num1} + \text{num2}} \quad \text{N}$

$\boxed{2; \text{return sum}; \quad \text{off}} \quad \text{N}$

if there are less than n characters in the IIP while refilling a buffer then off is read into the buffer after the IIP characters so each time fp have to make three tests

- 1) checking the 1 buffer is full
- 2) checking the second before is full
- 3) check the end of Ifp.

if ($fp == \text{eof}$)

{ unload buf 2

$fp = fp + 1;$

}

elseif ($fp == \text{eof} 2$)

{ unload buf 1

$fp = fp + 1;$

}

elseif ($fp == \text{eof}$)

{ terminate scanning

}

else

$fp = fp + 1;$

"

Sentinel - These source tests can be reduced to
use if we add a special character sentinel at
the end of each buffer. Sentinel is not a
part of the source program.

Let choose eof as sentinel

↓
fp

sum = num / + num eof

2 ; return sum() ; eof

↑
fp

Language of file for input UP
Input stream

```
fp = fp + 1;
if (fp == eof)
{
    if (fp == eob1)
        { unload buf1;
         fp = fp + 1;
        }
    elseif (fp == eob2)
        { unload buf1;
         fp = fp + 1;
        }
    else
        terminate scanning;
}
else
    fp = fp + 1;
```

SPECIFICATION OF TOKENS (Patterns)

Regular expressions are an important notation for specifying patterns. A regular exp. is defined by a set of strings which it matches.

Language (Λ) - Set of strings generated by the regular expression Λ .

Alphabet - is a finite set of symbols.

{0, 1} is a binary alphabet

String - finite sequence of symbols drawn from that alphabet.

Empty string is denoted by ϵ .

ϵ - string of length 0.

Ex., six is a string over an alphabet containing letters.

Terms for parts of strings

Prefix of s - obtained by removing some trailing symbols (0 or more)

eg. ~~take, table, e~~
tab, ~~for~~ table

Suffix of s - formed by deleting some leading symbols
eg. ~~able, but, e~~
~~able, for~~ table

Substring of s - deleting a suffix and a prefix from string s.

eg. abl, e, table etc for table

Proper prefix, suffix and substring of s - do not include s and ϵ

can be obtained by taking
portions of s .

e.g. baan is a subsequence of banana.

Operations on languages

Union

$$L_1 \cup L_2 = \{s \mid s \in L_1 \text{ or } s \in L_2\}$$

Intersection

$$L_1 \cap L_2 = \{s \mid s \in L_1 \text{ and } s \in L_2\}$$

Concatenation

$$L_1 L_2 = \{sp \mid s \in L_1 \text{ and } p \in L_2\}$$

Exponentiation

$$L^i = L^{i-1} \text{ for } i > 1$$

$$L^0 = \epsilon$$

Kleene closure

$$L^* = \bigcup_{i=0}^{\infty} L^i$$

(zero or more concatenations of L)

Positive closure

$$\text{(one or more concatenations of } L\text{)} \quad L^+ = \bigcup_{i=1}^{\infty} L^i$$

$$\text{Let } L = \{A, \dots, Z, a, \dots, z\}$$

$$D = \{0, 1, \dots, 9\}$$

then $(LUD)^*$ - is the set of all strings of letters and digits

and $(LUD)^+$ - is set of all strings of letters and digits beginning with a letter

L_1 - denotes set of strings of length two where one letter followed by one digit.

L^4 - set of strings having 4 length.

L^{OD} - set of letters and digits with 62 strings of length one.

D^+ - set of all strings of one or more digits.

Regular Expressions

A regular expression is defined by the set of strings it matches and is built up of basic regular expressions using a set of rules.

A regular expression is defined by the following rules:

- e is a regular expression.

- if symbol ' a ' is in alphabet then

- a is a regular expression.

- if r_1 and r_2 are regular expressions over the alphabet then

- a) r_1r_2 is a regular expression

- b) r_1^* is a regular expression.

- if r_1 is a regular expression

- a) r_1^* is a regular expression.

- b) (r_1) is a regular expression.

$a(ba)^*$

operator precedence *

concatenation

e.g. $(a) | ((b)^*(c)) \rightarrow a | b^*c$

Algebraic laws for Regular Expressions

some laws

description

1) $r|s = s|r$

is commutative

2) $r|(s|t) = (r|s)|t$

is associative

3) $r(st) = (rs)t$

concatenation is associative

4) $r(s|t) = rs|r t$

concatenation distributes over |

5) $\epsilon r = r\epsilon = r$

ϵ is the identity for concatenation

6) $r^* = (r|\epsilon)^*$

ϵ is guaranteed in closure

7) $r^{*\#} = r^*$

* is idempotent

the languages of all strings of length two over the alphabet $\{a, b\}$.

2) a^* denotes $\{\epsilon, a, aa, aaa, \dots\}$

all strings of two or more a's

3) $(a|b)^*$ denotes $\{\epsilon, a, b, aa, ab, ba, bb, aaa, \dots\}$

all strings of a and b.

4) $a|a^*b$ denotes the language -

$\{a, b, ab, aab, aab^*, \dots\}$

string a and all strings consisting of one or

more a's and ending in b.

5) $(a|b)^*c(a|b)^*$ - set of all strings over this alphabet that contains exactly one c.

some strings matched by above regular exp are
 $c, acaba, caaab, bbabb, aabbc$

Regular Definitions - It is good to give name to a regular expression so that we can refer it when we want to use it.

A regular definition is a ~~sequence~~ of definiteness of the form -

$$d_1 \rightarrow s_1$$

$$d_2 \rightarrow s_2$$

$$d_3 \rightarrow s_3$$

$$\vdots$$

$$d_m \rightarrow s_m$$

where each s_i is a regular exp defined over the alphabet $\Sigma \cup \{d_1, d_2, \dots, d_m\}$, d_i is name given to that regular exp.

eg. digit = 0|1|2|...|9

regular definition for defining the digit

$$\text{natural} \Rightarrow \text{digit} \cdot \text{digit}^*$$

$$\text{natural} \Rightarrow \text{digit}^+$$

eg. letters - $\rightarrow A|B|...|z|a|b|...|x| -$

$$\text{digit} \rightarrow 0|1|2|...|9$$

$$\text{id} \Rightarrow \text{letter} - (\text{letter} - 1 \text{ digit})^*$$

Design of lexical analyzer

As we all know that flow chart plays a main role in the writing of programs.

Transition diagram is a kind of flow chart for the lexical analysis.

In the construction of lexical analyzer first the specification of pattern is done. Then the transition diagram for each token is drawn.

Transition Diagram - A transition diagram is a directed graph with states drawn as circles (nodes) and edges represent transitions on input symbols. Each transition diagram has a start state and a final state.

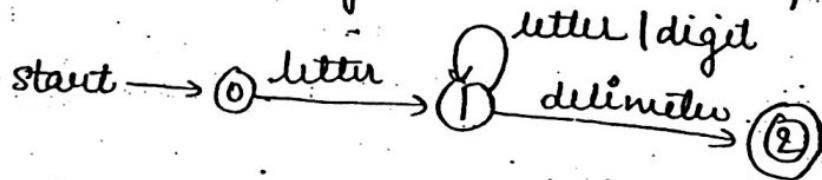
A transition diagram is used to keep track of information about characters that are seen as the forward pointer fp scans the input.

Ex- for identifier

the pattern in form of regular expression

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

the transition diagram for above pattern is



transition diagram

code for state 0

c = getchar();

if letter(c) then goto state 1

else fail();

code for state 1

c = getchar();

if letter(c) or digit(c) then goto state 1

else if delimiter(c) then goto state 2

else fail();

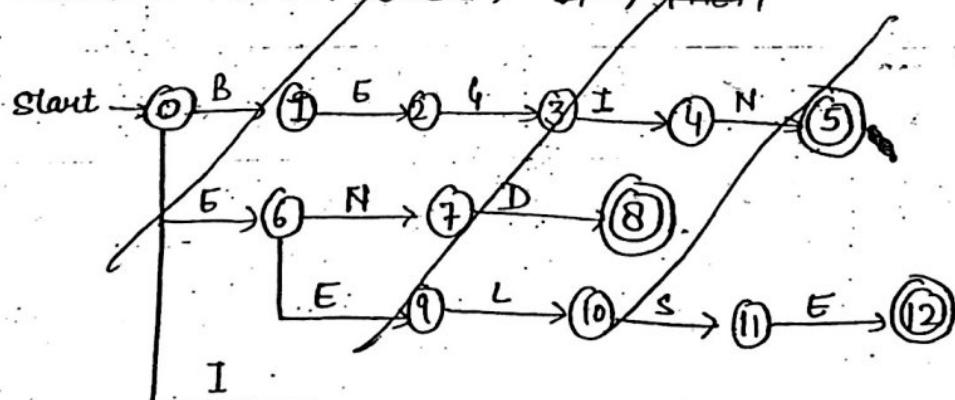
code for state 2

Retract();

return (id, install);

now let us have to draw the transition diagram
for keywords BEGIN, END, ELSE, IF, THEN

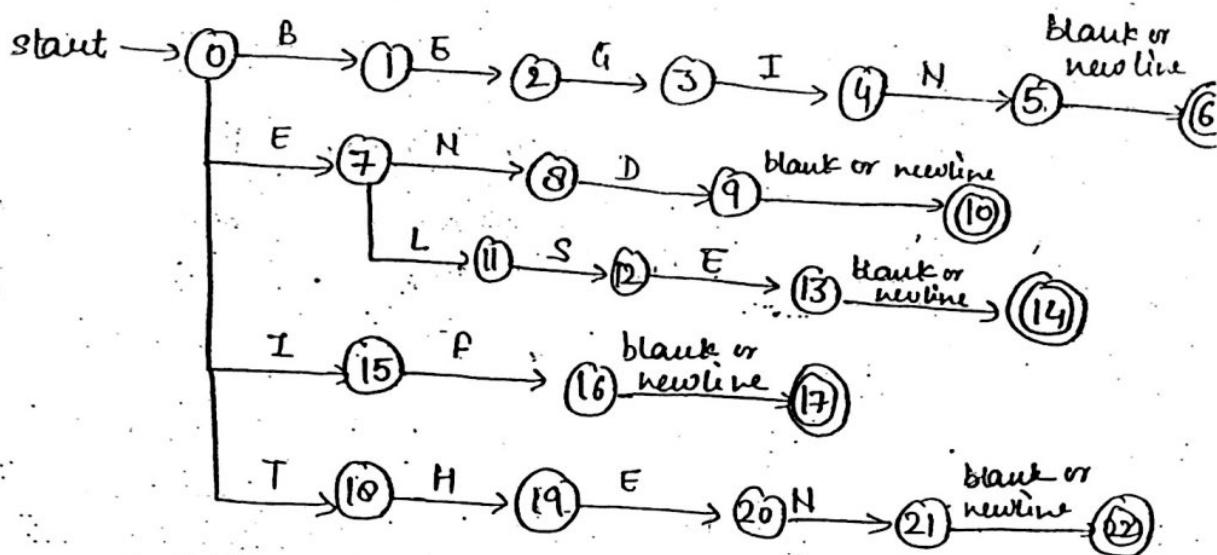
Transition
Diagram



B P.T.O

Now let we have to draw the transition diagram for keywords BEGIN, END, ELSE, IF, THEN

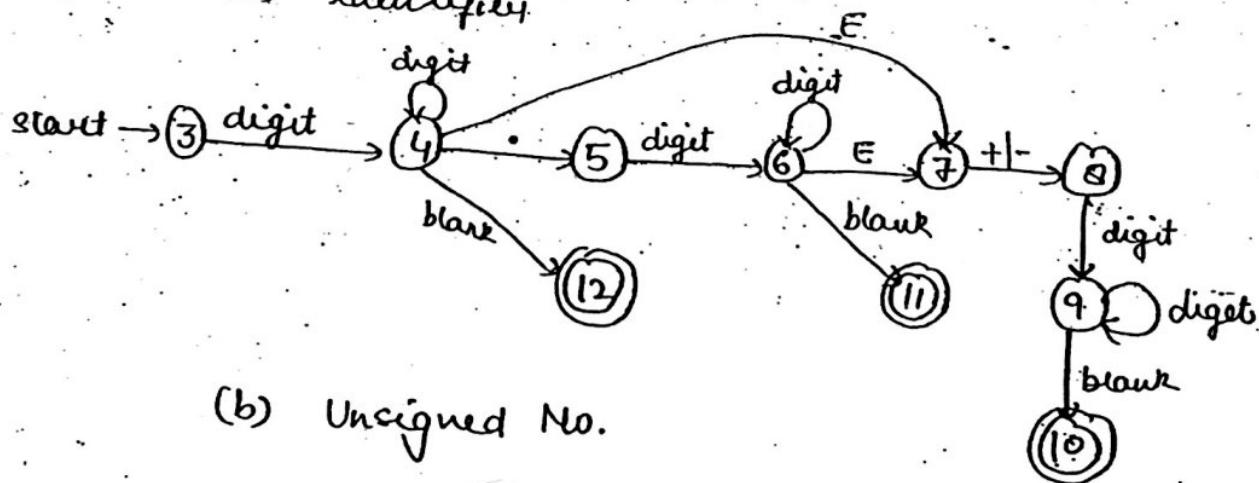
Transition diagram →



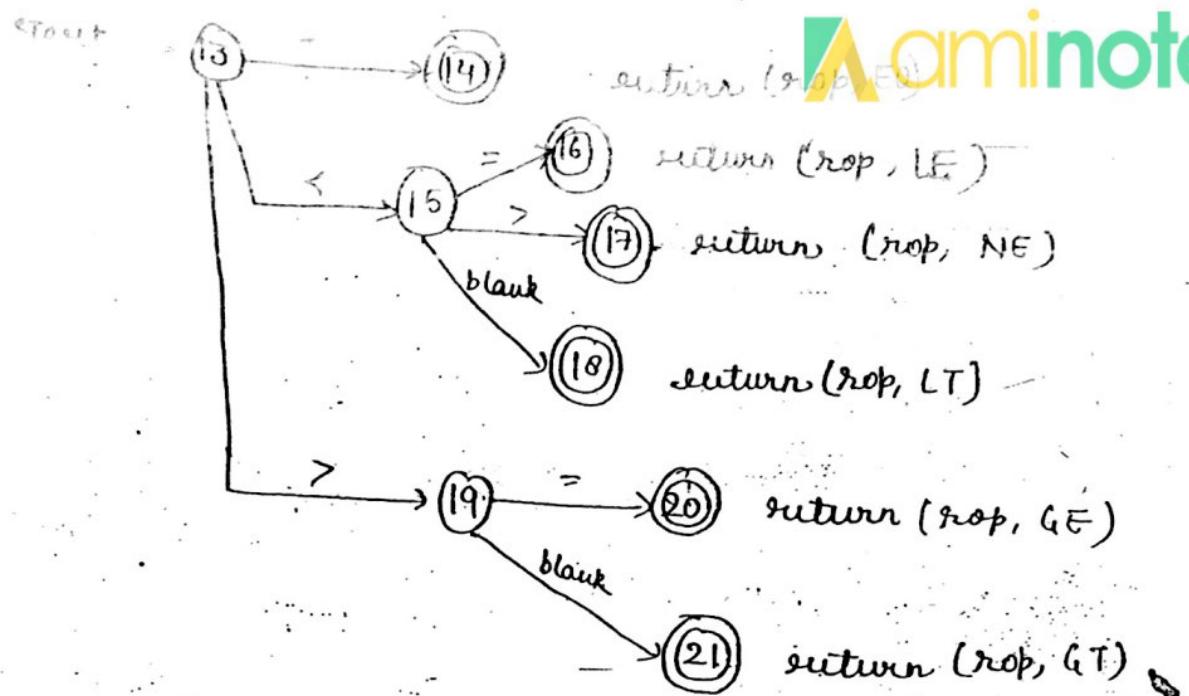
Now we have a transition diagram for identifier, unsigned number and relational operator as follows-



(a) identifier



(b) Unsigned No.



(c) transition diagram for relational operators

Converting above transition diagrams into code

function fail()

{ #p = 1b;

case start of

{ 0: start = 3;

3: start = 13;

13: recover();

default : error;

}

}

function nextoken()

{ state=0; start=0;

while(1)

{ case state of

0: {

lb = fp;

c = nextchar();

if letter(c) state = 1;

else ~~state~~ fail();

break;

}

1: {

c = nextchar();

if (letter(c) or digit(c)) state = 1;

else state = 2;

break;

}

2: {

fp = fp - 1;

return (id, input[lb, fp])

break;

}

/* ... write cases 3 to 12 state similarly ... */

13: {

c = nextchar();

if ~~else~~ (c == '=') state = 14;

elseif (c == '<') state = 15;

else if (c == '>') state = 19;

}

break;

15: { c = next char();
 if (c == '=') state = 16;
 else if (c == '>') state = 17;
 else state = 18;
 break;
}

16: { return (rop, LE);
 break;
}

17: { return (rop, NE);
 break;
}

18: { fp = fp - 1;
 return (rop, LT);
 break;
}

19: write cases for states 19 to 21 similarly.

Finite automation - lexical analyzers are
one simple representation which is called finite automata.

Finite automaton

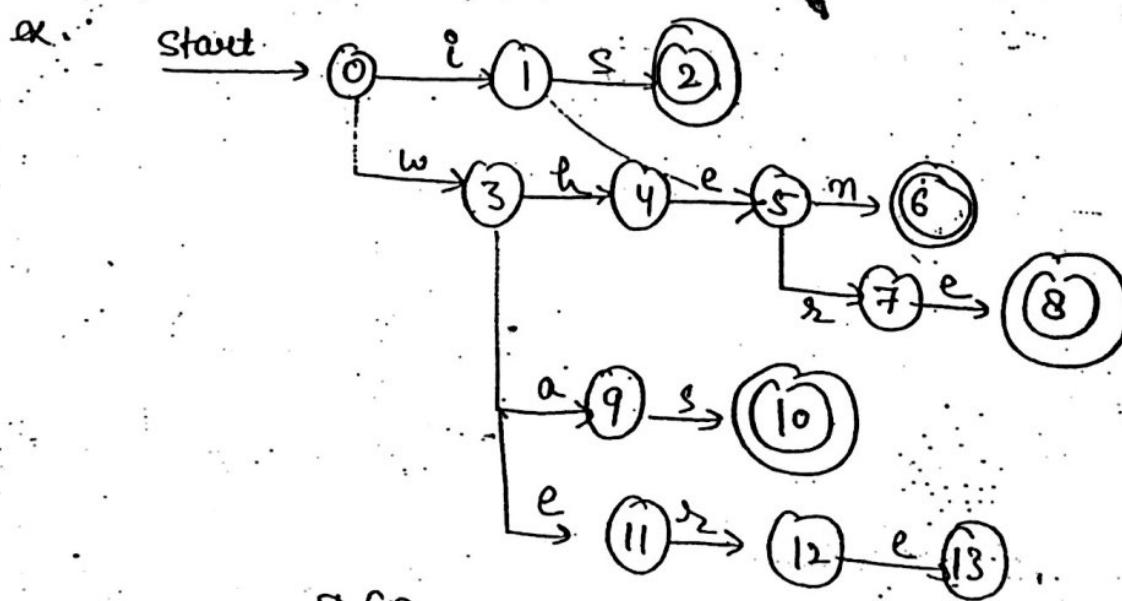
which is the part of the lexical analysis that
finite automaton are recognizers, they simply
say "yes" or "no" about each possible IIP string.

The automaton takes a sequence of inputs
and performs a sequence of actions for the
IIP string.

Finite automata come in two flavours:

a) NFA - non deterministic finite automaton.
any IIP alphabet can have several edges
out of the same state.

b) DFA deterministic finite automaton -
for each state and for each input alphabet
there can be only one edge leaving that state.



encountered the rule among
would and were.

IIP alphabet is { i, s, w, h, e, n, a, g }

set of states = { 0, 1, 2, 3, 4, ..., 13 }

Accepting state = { 2, 6, 8, 10, 13 }

start

DFA is defined by 5 tuples

s_0 start state

f final state

Σ IIP alphabet

S_0 set of states

Transition function $\rightarrow S \times \Sigma \rightarrow S$

In DFA the next state can be determined by knowing the current state and the current IIP symbol.

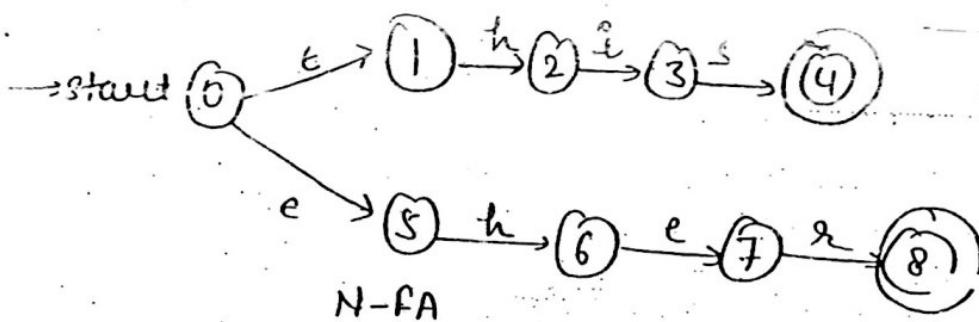
In a DFA there are no two edges leaving a given state have the same label.

Transition table

current state	IIP alphabet -							
	OIP state							
0	-	-	-	-	-	-	-	-
1	-	-	-	-	-	-	-	-
2	-	-	-	-	-	-	-	-
3	-	-	-	-	-	-	-	-
4	-	-	-	-	-	-	-	-
5	-	-	-	-	-	-	-	-
6	-	-	-	-	-	-	-	-
7	-	-	-	-	-	-	-	-
8	-	-	-	-	-	-	-	-
9	-	-	-	-	-	-	-	-
10	-	-	-	-	-	-	-	-
11	-	-	-	-	-	-	-	-
12	-	-	-	-	-	-	-	-
13	-	-	-	-	-	-	-	-

NFA

e.g. automaton that recognizes strings like abc



In this NFA two outgoing edges have the same label with ϵ .

NFA can have more states than the ~~DFA~~ equivalent DFA and it is difficult to decide which path to follow in the NFA if the transition table shows the transition to more than one state.

and also it is difficult to construct a DFA directly from a regular expression.

So Lex uses the two step process

- 1) Construct a NFA from the regular exp
- 2) Convert the NFA to a DFA.

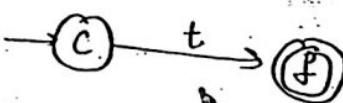
1) Connecting a regular expression

This construction process is guided by a algorithm known as Thompson's construction.

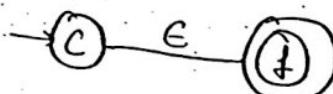
Four basic NFA structures are as follows -
 c - current state, f - final state

1) Let $t \in \Sigma$ symbol in Σ alphabet

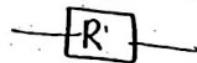
For regular expression t , NFA is



For ϵ



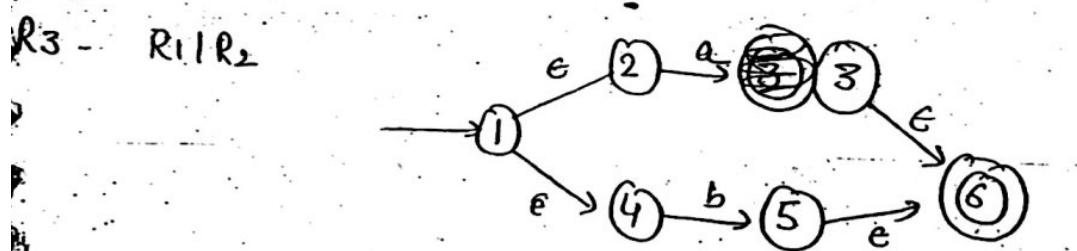
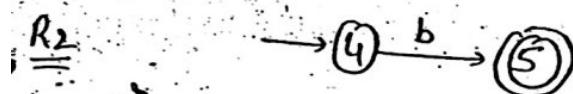
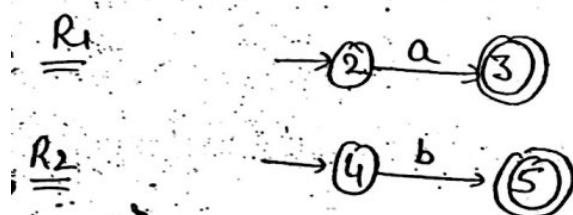
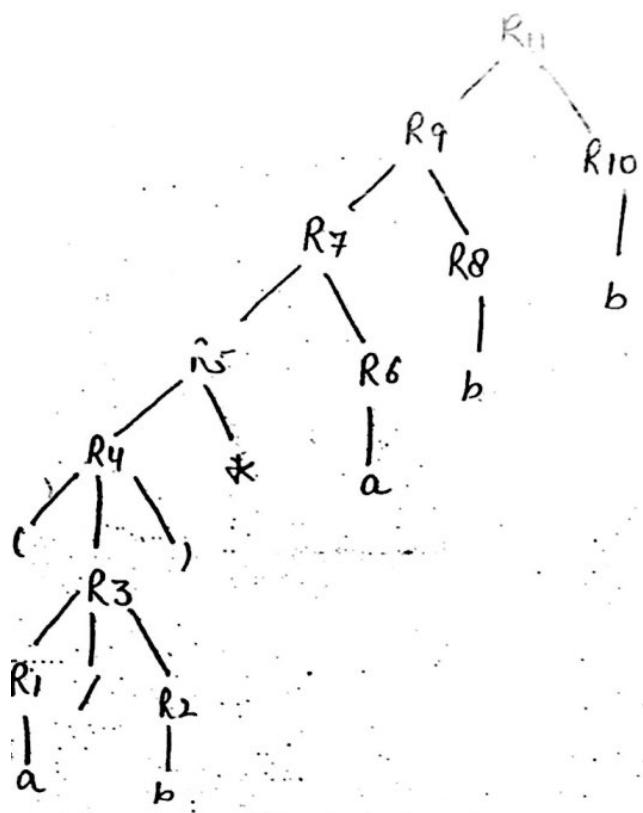
2) Let NFA for regular exp R is



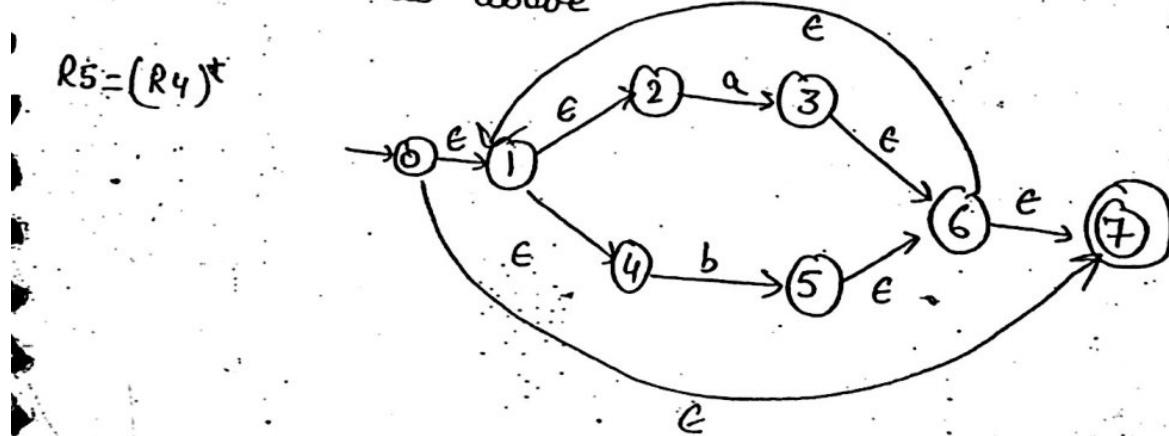
then concatenation $(R).(S)$ is represented by

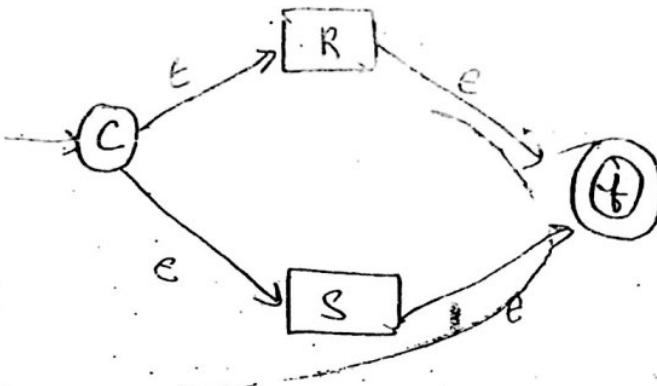


here the current state of R becomes the current state of S NFA and final state of R becomes the final state of NFA.



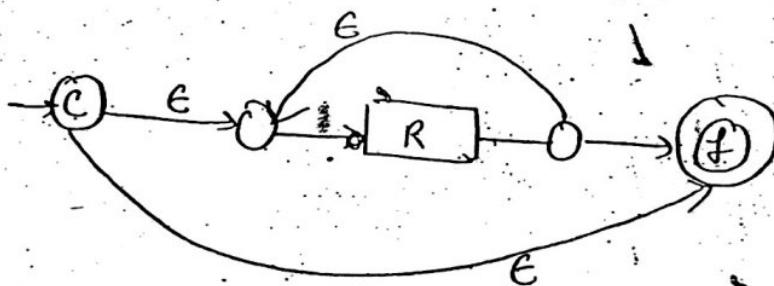
R4 = $(R_3)^*$ same as above





here the current and final states of R and S
are not the current and final states of NFA.
These are new current and final states.

4) NFA for regular expression RT



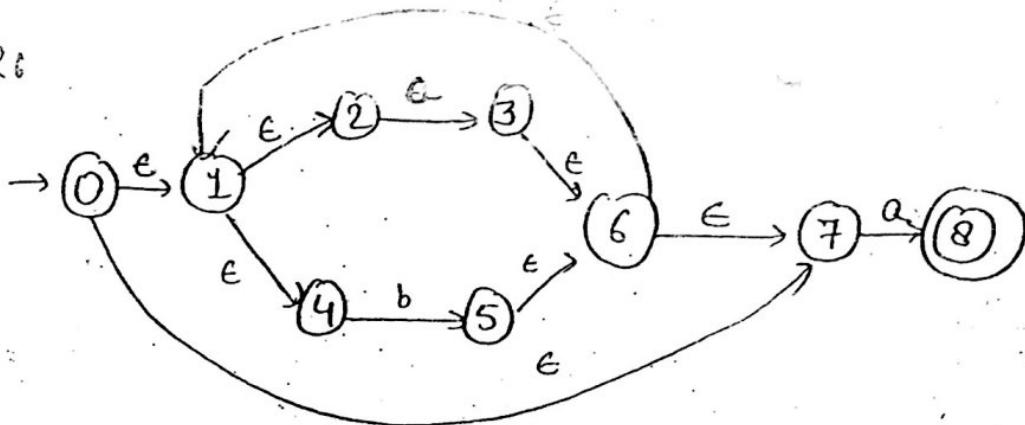
e.g. construction of NFA for regular exp.
 $(a/b)^*$, abb

let decompose this R into its primitive components.

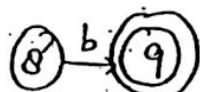
R6



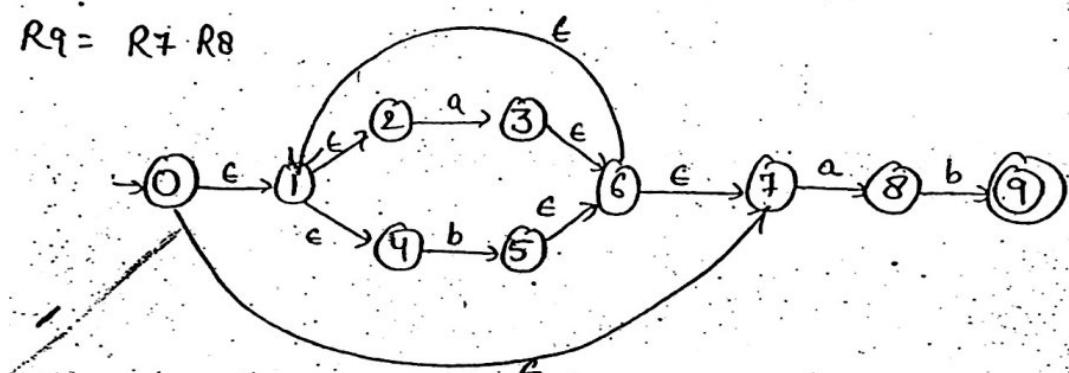
$$R7 = R5 \cdot R6$$



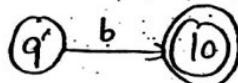
$$R8 =$$



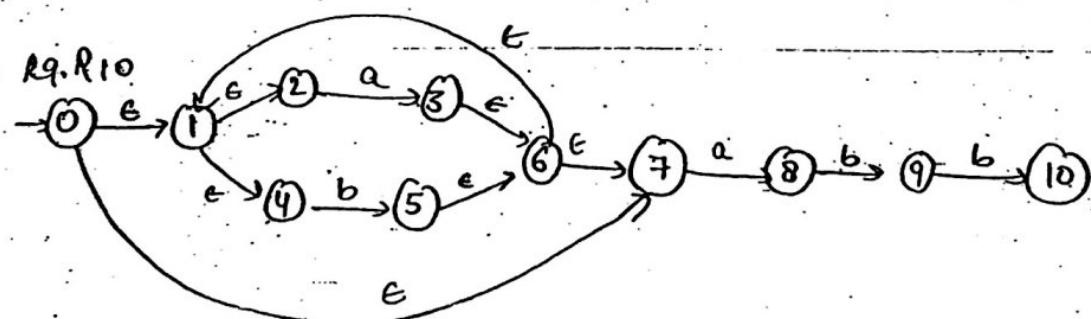
$$R9 = R7 \cdot R8$$



$$R10 =$$



$$R10 = R9 \cdot R10$$



representing $\overline{\text{exp}} (a/b)^*abb$

To obtain a DFA from an NFA we have to find out the states of an NFA that can be reached for a particular prefix and put together these states to form one state of a DFA.

steps for construction of a DFA from NFA called subset construction.

Let S be a set of states equiv. DFA state equivalent to set of states $\{s_1, \dots, s_i\}$ of NFA and T another DFA state eq. to set of states $\{t_1, t_2, \dots, t_j\}$ of NFA.

If $\{s_i, a, t_j\}$ exists in the NFA, then $\{S, a, T\}$ will exist in the DFA.

means for the IIP symbol a a state transition from S to T will occur if there exists a state transition from s_i to t_j labelled a .

now define ϵ -closure to be the set of NFA states by following rule,

- 1) s is added to ϵ -closure(s)
- 2). If t is in ϵ -closure(s), and there is an edge labeled ϵ from t to u then u is added

to ϵ -closure(2) repeat 2 until no new state can be added.

ex. see NFA for $(a|b)^*$ abb

The initial state for equivalent DFA is

initial state = ϵ -closure(0)

$$A = \{0, 1, 2, 4, 7\}$$

→ now find out those states from A who have transitions on a
these are 2 and 7 and have such transitions to 3 and 8.

$$\text{so } \epsilon\text{-closure}(\{3, 8\}) = \{3, 6, 7, 1, 2, 4, 8\}$$

$$B = \{1, 2, 3, 4, 6, 7, 8\}$$

→ now find out those states from A who have transitions on b. These are 4 and transition to 5.

$$\epsilon\text{-closure}(\{5\}) = \{5, 6, 7, 1, 2, 4\}$$

$$C = \{1, 2, 4, 5, 6, 7\}$$

Now find the states from B having transitions on a $2 \rightarrow 3$ and $7 \rightarrow 8$

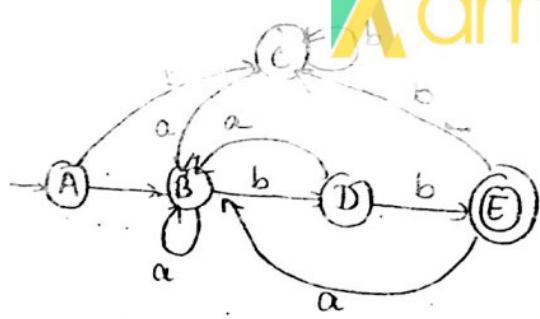
$$\epsilon\text{-closure}(\{3, 8\}) = \{1, 2, 3, 4, 6, 7, 8\}$$

for a b $4 \rightarrow 5, 8 \rightarrow 9$

$$\epsilon\text{-closure}(\{5, 9\}) = \{5, 6, 7, 1, 2, 4, 9\}$$

$$D = \{1, 2, 4, 5, 6, 7, 9\}$$

STATE	A	B	C
A	B	C	
B	B	D	
C	B	C	
D	B	E	
E	B	C	



Transition table for $(a/b)^*$ abb

Mindnising the no. of states of DFA

1) Construct a partition π of the set of states

two groups F - final states

S-F - non final states

eg:

two groups (ABCD), (E)

2) Now again split them

(E) can not be partitioned again so

we place (E) in π_{new}

now partitioned - (ABCD) on the basis of IIP symbols

let on input a each of these states go

to B but on b A B C go to members of group

ABCD while D goes to a member of another

group E.

so new value of π is (ABC) (D) (E)

(ABC) on IIP a → no split

on IIP b (AC) (B) (D) (E)

(AC) on IIP a no split

on IIP b no split

so final partition is. (AC) (B) (D) (E)

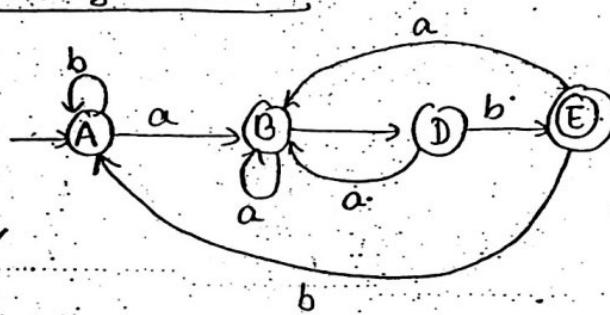
— now superset (AC) by A

then we get

States	IIP	
	a	b
A	B	A
B	B	D
D	B	E
E	B	A

Reduced automaton

minimized DFA is



minimized DFA accepting $(a/b)^* abb$

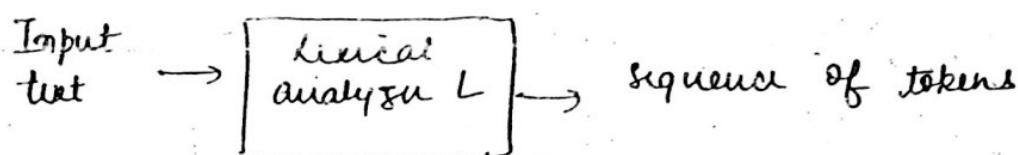
~~Automatic generation of lexical analyzer~~ **aminotes**
lexical analyzer generator is used to generate
lexical analyzer automatically.

LEX - is a user tool which can automatically
generate a lexical analyzer.

This LEX is a lexical analyzer generator
which takes as Input the precise specification
of tokens of the language, and a specification
of the action to be performed on identification
of each token. As an output it generates
a lexical analyzer.

IIP of LEX - Input of lex compiler is a
source program which is a specification of a
lexical analyzer, consisting of a set of regular
expressions together with an action for each
regular expression. This action is a piece of
code which is to be executed whenever a token
specified by the corresponding regular expression
is recognized. This action will pass an
indication of the token found to the parser,
with the side effects such as making an entry
in symbol table. ~~and many definitions~~
~~in function table (actions)~~

O/P of LEX - is lexical analyzer program
constructed from the LEX source specification.



LEX source programme consists of two parts

(a) Auxiliary Definitions - these are the statements

of the form -

$$D_1 \rightarrow R_1$$

$$D_2 \rightarrow R_2$$

:

$$D_n \rightarrow R_n$$

where each D_i is a distinct name and each R_i is a regular expression whose symbols are chosen from $\Sigma \cup \{D_1, D_2, \dots, D_{i-1}\}$.

Let four identified sequence of auxiliary definitions are -

$$\text{letter} \rightarrow A|B| \dots |Z$$

$$\text{digit} \rightarrow 0|1| \dots |9$$

(b) Translation Rules, these are the statements of the form

$$P_1 \quad \{A_1\}$$

$$P_2 \quad \{A_2\}$$

:

$$P_m \quad \{A_m\}$$

each P_i is a regular expression called Pattern over the alphabet Σ and the auxiliary definition names

Actions describes the form of tokens each AI is to recognize the Action true the next analysis should take when PC is found.

This action can be a returning of position of token to point or can be a making of entry into the symbol-table.

Ex-

Let we have following tokens

begin

end

if

then

else

identifier

~~constants~~

<

<=

=

<>

>

>=

We have to generate the Lex source program for these tokens.

Syntax Definitions

letter → A/B.../z

digit → 0/1/.../9

Translation Rules

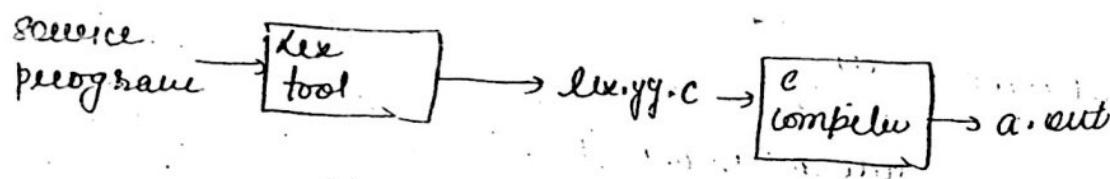
BEGIN	{ return (BEGIN); }
END	{ return (END); }
IF	{ return (IF); }
THEN	{ return (THEN); }
ELSE	{ return (ELSE); }
letter (letter digit)*	{ LEXVAL = install(); return (ID); }
digit ⁺	{ LEXVAL = install(); return (NUM); }
<	{ LEXVAL = LT, return (ROP); }
<=	{ LEXVAL = LE, return (ROP); }
=	{ LEXVAL = EO, return (ROP); }
<>	{ LEXVAL = NE, return (ROP); }
>	{ LEXVAL = GT, return (ROP); }
>=	{ LEXVAL = GE, return (ROP); }

LEX ProgramConflict Resolution in LEX

If several prefixes of the input match one or more patterns then

- 1) Always prefers a longer prefix to a shorter prefix. e.g. if < is encountered then token can be < or <=
- 2) If the longest possible prefix matches two or more patterns, prefers the pattern listed first in the LEX program. e.g. BEGIN will be treated as keyword not identifier

Automatic creation of lexical analysis
using lex



→ The binary executable file contains the assembly language code.

→ The assembly language code is converted into machine language code.

→ Machine language code is converted into binary code.

→ Binary code is converted into machine language code.

→ Machine language code is converted into assembly language code.

→ Assembly language code is converted into binary code.

→ Binary code is converted into machine language code.

→ Machine language code is converted into assembly language code.

→ Assembly language code is converted into binary code.

→ Binary code is converted into machine language code.

→ Machine language code is converted into assembly language code.

→ Assembly language code is converted into binary code.

→ Binary code is converted into machine language code.

→ Machine language code is converted into assembly language code.

→ Assembly language code is converted into binary code.

→ Binary code is converted into machine language code.

→ Machine language code is converted into assembly language code.

→ Assembly language code is converted into binary code.

→ Binary code is converted into machine language code.

→ Machine language code is converted into assembly language code.

Compiler Construction

C86

SYNTAX ANALYSIS

In previous step phase lexical analysis we had learned about the specification of patterns using regular expressions and the construction of lexical analyzer to recognize tokens using automaton.

Now the syntax analysis phase is the next phase after lexical analysis in compiler design.

Syntax Analyzer (Parser) takes a token stream produced by lexical analyzer as input and create a tree like intermediate representation (Parse tree) as O/P that depicts the grammatical structure of tokens stream.

We all know that every programming language has some precise rules to perceive the syntactic structure of well formed program.

Eg. In C, a program is made up of functions, a function out of declarations and statements, a statement out of expressions and so on.

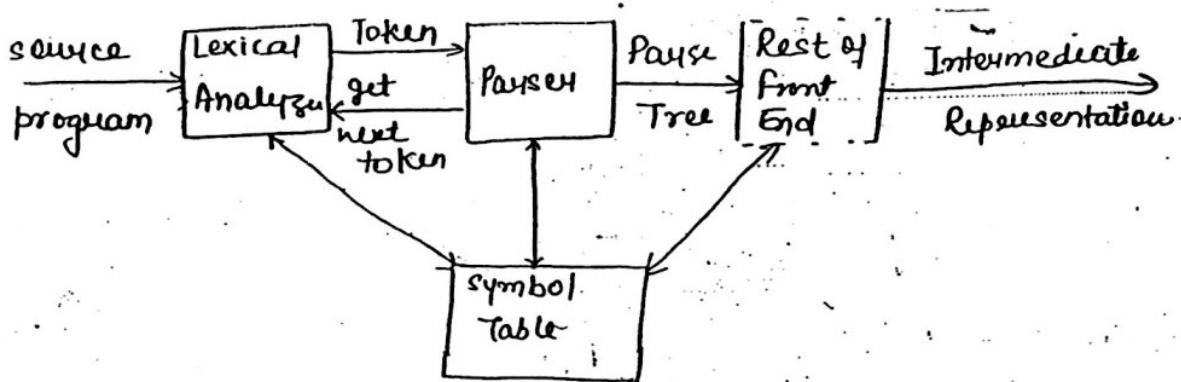
This syntax of programming language constructs can be specified by CFG (context free grammar) which is also called BNF (Backus-Naur Form) description.

grammar is used to specify the language. This grammar has following advantages -

- 1) Grammar gives a precise, easy to understand syntactic specification of programming language.
- 2) An efficient parser can be constructed automatically from a properly designed grammar.
- 3) A grammar imparts a structure to a program that is useful for its translation into object code and for the detection of errors.
- 4) A grammar allows a language to be evolved or developed iteratively by adding new constructs to perform new tasks.

ROLE OF PARSER

- It obtains a string of tokens from the lexical analyser.
- It groups the tokens to identify larger structures in the program. This process is done to verify that the string can be generated by the grammar for the source language.
- It should report any syntax error in the program.
- It should also recover from the errors so that it can continue to process the rest of the I/P.



CONTEXT FREE GRAMMAR

To start for the

The syntax of a programming language is specified

by a notation called context free grammar (CFG)
which is also called as BNF (Backus Naur Form).

1 A CFG involves four quantities

→ Terminals - these are the basic symbols of
which strings in the language are composed.

Tokens can be considered as terminals symbols.

→ Non terminals - these are the variables
that denote a set of strings. They do not exist
in the source program, they only help in
defining the language generated by the
grammar.

→ start symbol - It is a ^{by} non terminal and the set of strings denoted by this start symbol is the language defined by the grammar. Means it denotes that language in which we are truly interested.

→ Production Rules → It defines the way in which the terminals and non terminals can be combined to form strings.

Each production consists of

- a ~~non~~ terminal on left hand side of production
- A symbol →
- Zero or more terminals and non terminals on the right hand side

we may define a CFG as follows

$$G = (T, N, S, P)$$

where $T(G) \rightarrow$ non empty set of terminals

$N(G) \rightarrow \dots \text{ " " } \text{ non terminal}$

$S(G) \rightarrow \text{ single non terminal } \in N(G)$

$P(G) \rightarrow \text{ set of production rules}$

like $A \rightarrow \alpha$

where $A \in N(G)$

& $\alpha \in (N(G) \cup T(G))^*$

St Derivation produces a new string from given string. So derivation can be used repeatedly to obtain a new string from a given string.

Let $\alpha_1 \rightarrow \alpha_2 \rightarrow \alpha_3 \rightarrow \alpha_4 \dots \rightarrow \alpha_n$

then we can say

$\alpha_1 \xrightarrow{*} \alpha_n$ means α_1 derives α_n and derivation comprise of zero or more steps.

If $\alpha \xrightarrow{*} \beta$ and $\beta \xrightarrow{*} \gamma$

then we can say that

$\alpha \xrightarrow{*} \gamma$

$\xrightarrow{*}$ is used to say that the derivation comprises of one or more steps.

Let $s \xrightarrow{*} \alpha$ where $\alpha \in (N(G) \cup T(G))^*$

means α contains non terminals, then α is called sentential form of grammar.

as follows

$$S \rightarrow E \text{ op } E \mid (E) \mid -E \mid \text{id}$$

$$\text{op} \rightarrow + \mid - \mid * \mid / \mid \cdot$$

Now we have to derive a string ($\text{id} + \text{id}$)
from the above grammar.

So

$$E \rightarrow (E)$$

$$\rightarrow (\underline{E} \text{ op } E)$$

$$\rightarrow (\text{id} \text{ op } E)$$

$$\rightarrow (\text{id} + \underline{E})$$

$$\rightarrow (\text{id} + \text{id})$$

LEFT most and RIGHT most derivation

Each derivation step needs to choose a non-terminal to replace and a production to apply.

If at each step we replace the left most non terminal, derivation is known as left most derivation.

2 If at each step we replace the right most non terminal, derivation is known as Right most derivation.

Ex-1 The grammar for ⁵⁸ that describe the syntax of an expression consisting of +, - and digits -

$$P(G) \rightarrow \left\{ \begin{array}{l} \text{exp} \rightarrow \text{exp} + \text{digit} \\ \text{exp} \rightarrow \text{exp} - \text{digit} \\ \text{exp} \rightarrow \text{digit} \\ \text{digit} \rightarrow 0/1/2/\dots/9 \end{array} \right\}$$

here

$$T(G) = \{0, 1, \dots, 9, +, -\}$$

$$N(G) = \{\text{exp, digit}\}$$

$$S(G_1) = \{\text{exp}\}$$

Ex-2 Now expression has following terminal symbols

id, +, -, *, /, (,), ;

then the grammar is

$$\text{exp} \rightarrow \text{exp} + \text{term}$$

$$\text{exp} \rightarrow \text{exp} - \text{term}$$

$$\text{exp} \rightarrow \text{term}$$

$$\text{term} \rightarrow \text{term} * \text{factor}$$

$$\text{term} \rightarrow \text{term} / \text{factor}$$

$$\text{term} \rightarrow \text{factor}$$

$$\text{factor} \rightarrow (\text{exp})$$

$$\text{factor} \rightarrow \text{id}$$

grammar for $+,-,*,/,\frac{ab}{c}$, in can also
be modified to

$$Exp \rightarrow Exp \ op \ Exp$$

$$Exp \rightarrow (Exp)$$

$$Exp \rightarrow -Exp$$

$$Exp \rightarrow id$$

$$Op \rightarrow +|-|*|/|\frac{ab}{c}|$$

by using shorthands it can be written as,

$$Exp \rightarrow Exp \ op \ Exp \mid (Exp) \mid -Exp \mid id$$

$$Op \rightarrow +|-|*|/|\frac{ab}{c}|$$

DERIVATIONS AND PARSE TREES

DERIVATION - derivation provides a mean for generating the sentences of a language.

Derivation refers to replacing an instance of a non-terminal in a given string's non-terminal by R.H.S of the Production rule whose L.H.S contains the non-terminal to be replaced.

e.g. $A \rightarrow \gamma$ is a production rule

$$\text{then } \alpha A \beta \Rightarrow \alpha \gamma \beta.$$

means $\alpha A \beta$ derives $\alpha \gamma \beta$ in a derivation of one step.

eg: derivation of string "id + (id * id)" by the same grammar - 59

Left most
derivation

$$\begin{aligned}
 E &\Rightarrow E \text{ op } E \\
 &\rightarrow id \text{ op } E \\
 &\rightarrow id + E \\
 &\rightarrow id + (E) \\
 &\rightarrow id + (E \text{ op } E) \\
 &\rightarrow id + (id \text{ op } E) \\
 &\rightarrow id + (id * E) \\
 &\rightarrow id + (id * id)
 \end{aligned}$$

Right most
derivation

$$\begin{aligned}
 E &\Rightarrow E \text{ op } E \\
 &\rightarrow E \text{ op } (E) \\
 &\rightarrow E \text{ op } (E \text{ op } E) \\
 &\rightarrow E \text{ op } (E \text{ op } id) \\
 &\rightarrow E \text{ op } (E * id) \\
 &\rightarrow E \text{ op } (id * id) \\
 &\rightarrow id + (id * id)
 \end{aligned}$$

Parse Trees A parse tree is a graphical representation of a derivation that filters out the order in which productions are applied to replace non terminals.

Properties of parse tree

- 'Root contains' the start symbol.
- Each interior node contains a non terminal
- If there is a production rule $A \rightarrow Y_1 Y_2 Y_3$ and there is a interior node containing a non terminal A then children of that node contains $Y_1 Y_2 Y_3$ from left to right.

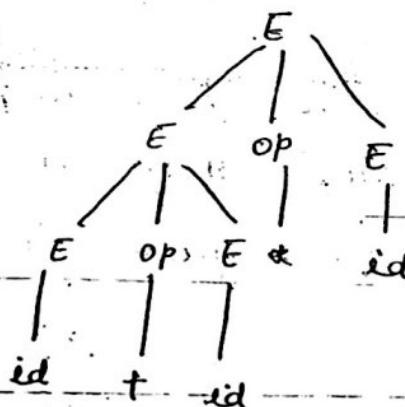
fact leaf node contains a symbol or an E (Epsilon)

60
=

Ex- we have to generate a string id + id * id from the previous grammar -

$$\begin{aligned} E &\rightarrow E \text{ op } E \\ &\Rightarrow E \text{ op } E \text{ op } E \\ &\Rightarrow id \text{ op } E \text{ op } E \\ &\Rightarrow id + E \text{ op } E \\ &\Rightarrow id + id \text{ op } E \\ &\Rightarrow id + id * E \\ &\Rightarrow id + id * id \end{aligned}$$

Parse tree for the above derivation is as



Parse tree for string id + id * id

Ambiguity - A grammar is ambiguous if it produces more than one parse tree for some sentence. Mostly a grammar is ambiguous when the same non-terminal appears twice on a R.H.S of the production.

For eg. string $id - id / id$ has two different left most derivations.

$$E \rightarrow E \text{ op } E$$

$$\rightarrow id \text{ op } E$$

$$\rightarrow id - E$$

$$\rightarrow id - E \text{ op } E$$

$$\rightarrow id - id \text{ op } E$$

$$\rightarrow id - id / E$$

$$\rightarrow id - id / id$$

$$E \rightarrow E \text{ op } E$$

$$\rightarrow E \text{ op } E \text{ op } E$$

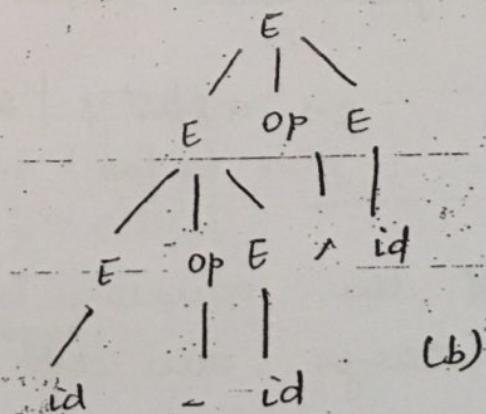
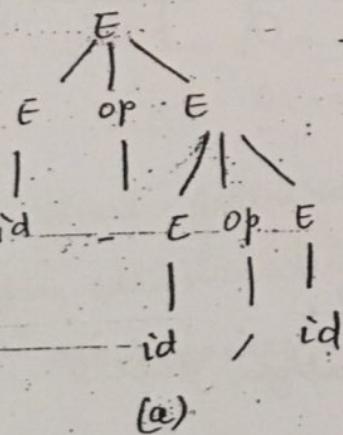
$$\rightarrow id \text{ op } E \text{ op } E$$

$$\rightarrow id - E \text{ op } E$$

$$\rightarrow id - id \text{ op } E$$

$$\rightarrow id - id / E$$

$$\rightarrow id - id / id$$



Ambiguous grammars are not good for parsing because we can not determine which production should be preferred for the expansion.

So, grammar should be rewritten to be unambiguous.

Rules for disambiguating a grammar

Generally productions are ambiguous when they have more than one occurrence of a given non-terminal on their right hand side.

Rules :-

- check the precedence of the operators involved in the productions.
- Different precedence operators treated differently for removing ambiguity
- first remove the ambiguity for minimum precedence & then go on.
- ~~if the operator is left~~
if we have a production rule like

$$S \rightarrow \alpha S \beta S \gamma \mid \alpha_1 \mid \alpha_2 \dots \alpha_n$$

then ambiguity can be removed by reworking the production rule as

$$\begin{aligned} S &\rightarrow \alpha S \beta S' \gamma \mid S' \\ S' &\rightarrow \alpha_1 \mid \alpha_2 \dots \alpha_n \end{aligned}$$

If the operator is left associative then
change the ~~right~~ most symbol

eg. $E \rightarrow E * E \mid id$
 $\quad \quad \quad \left\{ \begin{array}{l} E \rightarrow E' * E \mid E' \\ E' \rightarrow id \end{array} \right.$

and if the operator is right associative then
change the left most symbol

eg. $E \rightarrow E * E \mid id \rightarrow \left\{ \begin{array}{l} E \rightarrow G * E' \mid E' \\ E' \rightarrow id \end{array} \right.$

eg:

$$E \rightarrow E+E \mid E-E \mid E/E \mid E * E \mid E \uparrow E \mid (E) \mid id$$

here the precedence of operator is as follows:

$+, -$ } lowest

$*, /$ } and $+, -, *, /$ are left associative
 \uparrow ↓ highest and \uparrow is right associative.

so first remove the ambiguity of $+ & -$ grammar

$$\begin{array}{l} \left\{ \begin{array}{l} E \rightarrow E+T \mid E-T \mid T \\ T \rightarrow E/E \mid E * E \mid E \uparrow E \mid (E) \mid id \end{array} \right. \\ \quad \left[\begin{array}{l} T \rightarrow T/F \mid T \leftarrow T \mid T \uparrow T \mid (E) \mid id \end{array} \right] \\ \quad \left\{ \begin{array}{l} T \rightarrow T/F \mid T \leftarrow F \mid F \\ F \rightarrow T \uparrow T \mid (E) \mid id \end{array} \right. \\ \quad \quad \left[\begin{array}{l} F \rightarrow F \uparrow F \mid (E) \mid id \\ F \rightarrow G \uparrow F \mid G \end{array} \right] \\ \quad \quad G \rightarrow (E) \mid id \end{array}$$

so final grammar is

after eliminating

ambiguity \rightarrow

$E \rightarrow E+T \mid E-T \mid T$
$T \rightarrow T/F \mid T \leftarrow F \mid F$
$F \rightarrow G \uparrow F \mid G$
$G \rightarrow (E) \mid id$

$E \rightarrow E+T \mid E-T \mid T$
 $T \rightarrow T/F \mid T \leftarrow F \mid F$
 $F \rightarrow G \uparrow F \mid G$
 $G \rightarrow (E) \mid id$

Eliminating left recursion
 Syntactic structure of a programming languages are expressed using recursive rules. These recursions can have left excursion & right excursion.

Left recursive grammars are not suitable for top down parsers.

Left recursive grammars. - A grammar is left recursive if the first symbol in the right hand side of a rule is the same non terminal as that in the left hand side.

e.g. $S \rightarrow S\alpha$

Top down parsing methods can not handle left recursive grammars, so the elimination of left excursion is needed.

Rules: ex if we have a production like

$S \rightarrow S\alpha | \beta$

then

$S \rightarrow B_1 S'$

$S' \rightarrow \alpha_1 S' | \epsilon$

In General

$S \rightarrow S\alpha_1 | S\alpha_2 | \dots | S\alpha_n | B_1 | B_2 | \dots | B_m$

replaced by

$S \rightarrow B_1 S' | B_2 S' | \dots | B_m S'$

$S' \rightarrow \alpha_1 S' | \alpha_2 S' | \dots | \alpha_n S'$

e.g. $E \rightarrow E + T | T$
 $T \rightarrow T * F | F$
 $F \rightarrow (E) | id$

then

$E \rightarrow TE'$

$E' \rightarrow +TE' | \epsilon$

$T \rightarrow FT'$

$T' \rightarrow \alpha F T' | \epsilon$

$F \rightarrow (E) | id$

65

But if the elimination of left recursion involve
indirectly recursive grammar - eg. -

$$S \rightarrow A \rightarrow Ad \mid b$$

$$A \rightarrow Sd \mid b$$

written as

$$S \rightarrow Sd \mid b$$

$$S \rightarrow bS'$$

$$S \rightarrow dS' \quad (\text{E})$$

$$A \rightarrow A\alpha \mid \beta$$

$$A \rightarrow BA'$$

$$A' \rightarrow \alpha A' \mid \epsilon$$

eg-2

$$A \rightarrow Ba \mid b$$

$$B \rightarrow bd \mid B' \mid B'$$

$$A \rightarrow Ba \mid b$$

$$B \rightarrow Bc \mid Ad \mid E$$



$$A \rightarrow Ba \mid b$$

$$B \rightarrow Bc \mid Bad \mid bd \mid E$$



$$A \rightarrow Ba \mid b$$

$$B \rightarrow bd \mid B' \mid B'$$

$$B' \rightarrow cB' \mid ad \mid b' \mid E$$

$$B \rightarrow Bc \mid Bad \mid bd \mid E$$

$$B \rightarrow bdA$$

$$B \rightarrow bdA' \mid A'$$

$$A' \rightarrow CA' \mid adA' \mid E$$

I/P - grammar G with no cycles or ϵ production
O/P - equivalent grammar with no left recursion 66

1. Arrange non terminals in some order A_1, A_2, A_n
2. for $i=1$ to n do begin
 - for $j=1$ to $i-1$ do begin
 - replace each production of the form $A_i \rightarrow A_j \gamma$ by the production $A_i \rightarrow \gamma_1 | \gamma_2 | \dots | \gamma_k$ where $A_j \rightarrow \gamma_1 | \gamma_2 | \dots | \gamma_k$ are all current A_j productions.
 - end
 - eliminate the immediate left recursion among the A_i productions.
 - end.

Elimination of left factoring

67

Left factoring is a process which isolates the common parts of two productions into a single production. After elimination of left factoring the produced grammar is suitable for predictive parsing.

Rule

$$\text{eg } A \rightarrow \alpha\beta_1 | \alpha\beta_2 | \dots | \alpha\beta_n$$

can be replaced by

$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta_1 | \beta_2 | \dots | \beta_n | \epsilon$$

$$\text{eg } S \rightarrow iETs | iETse | a$$

$$E \rightarrow b$$

left factored grammar is

$$S \rightarrow iETs' | a$$

$$S' \rightarrow es | e$$

$$E \rightarrow b$$

Precedence of Operators

will be calculated as

9 + 5 = 14

68

$$9 + (5 \times 2)$$

because $*$, / has higher precedence than +, -

Associativity of Operators

9-5-2

$+, -, *, /$ are left associative

so - calculated as $(9-5)^{-2}$

$$a = b = c$$

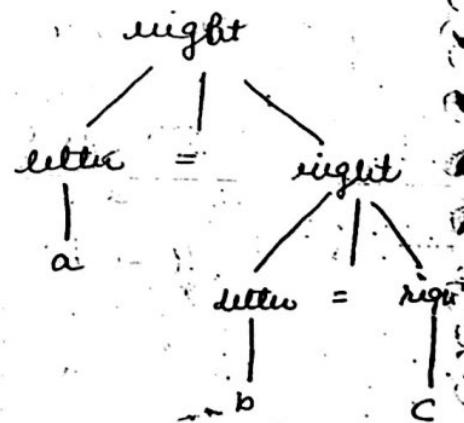
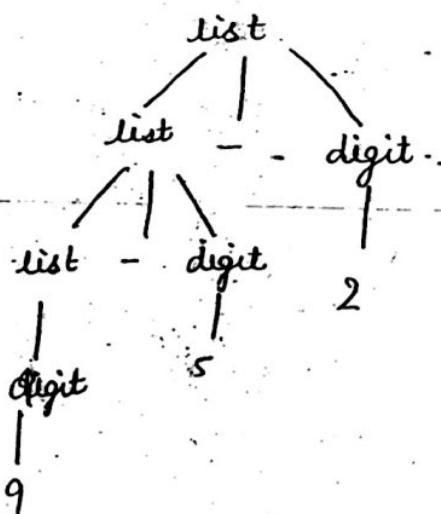
= is eight associative

Calculated

$$a = (b = c)$$

$\text{list} \rightarrow \text{list-digit} \mid$
 $\quad \quad \quad \text{digit}$
 $\text{digit} \rightarrow 0 \mid 1 \cdots \mid 9$

right \rightarrow letter = right /
letter
letter \rightarrow a/b/c



Parser - A parser scans an input token stream from left to right and groups those tokens into larger structures in order to check the syntax of an IIP string by identifying a derivation using production rules of the grammar.

Parser takes a input string w and a context free grammar G and produce output -

- a) a parse tree for w if w is a sentence of grammar G i.e. $w \in L(G)$.
- b) error message if w is not a sentence of grammar G and also provide some information like type of error and location of error.

So, output of a Parser is a representation of a parse tree for the input, if the input is syntactically well formed.

⇒ Parsing Strategies - There are two basic strategies for parsing on the basis of the way in which a parse tree can be generated.

By this strategy parse trees can be generated from root to leaves.

2) Bottom up Parsing

By this strategy parse trees can be generated from leaves to root.

⇒ TOP DOWN PARSING

Top down parsing involves the construction of parse tree from root to leave using left most derivation.

Here the main task is the finding of a suitable production rule for each left most non terminal.

Top down parsing can be done in two ways

- without Backtracking →
- with Backtracking

Here we can see that ⁷¹ at each step we have only one unique production rule for each non terminal.

Top down Parsing with Backtracking

If a left most non terminal has more than one production rule

like $A \rightarrow \alpha_1 | \alpha_2 | \dots | \alpha_n$

then there is a problem to decide which

A -production should be used for derivation here

Then parser will select one of the A -production for derivation and if the derivation is done it enters the successful completion of parsing otherwise parser backtracks & goes to the previous step and apply other production

Ex → examine $S \rightarrow cAd$

$A \rightarrow ab$

$A \rightarrow a$

IIP string

cad

If we top down parsing for each left-most non-terminal has only one unique production rule, then backtracking is not needed for derivation.

Eg.

$D \rightarrow \text{type list};$

$\text{type} \rightarrow \text{int}$

$\text{type} \rightarrow \text{float}$

$\text{list} \rightarrow \text{id}, \text{list}$

$\text{list} \rightarrow \text{id}$

now $\text{id} \rightarrow a/b$

we have to derive a sentence

int a, b;

so

$D \rightarrow \text{type list};$

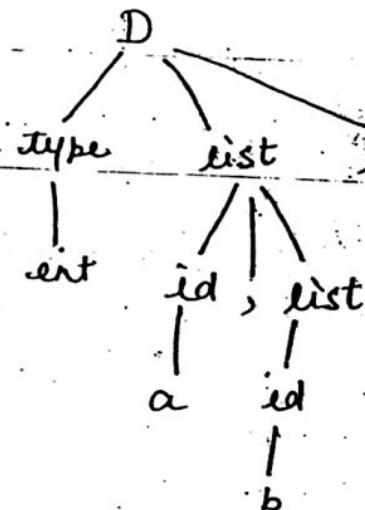
$\rightarrow \text{int list};$

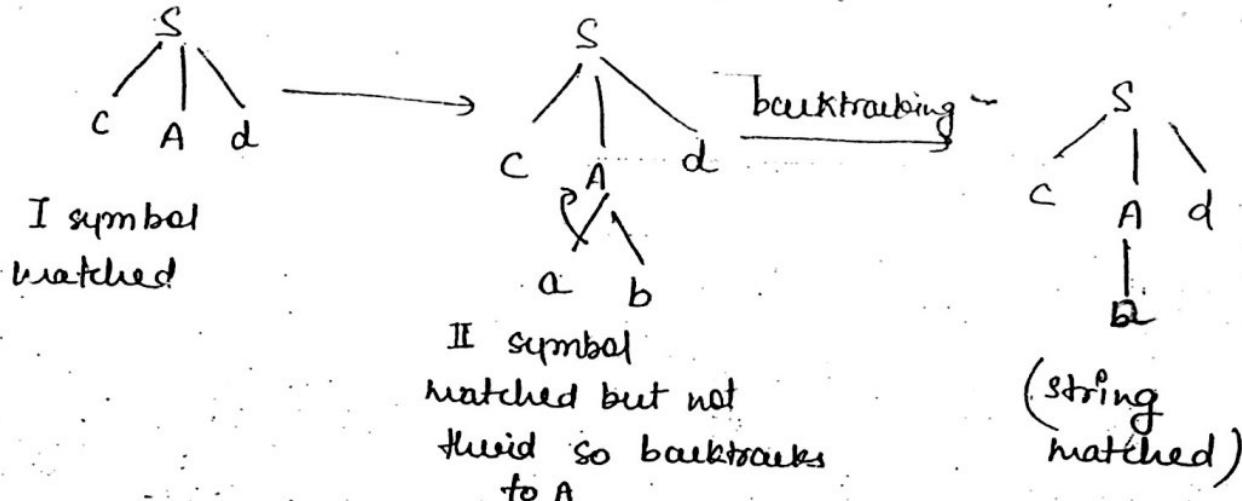
$\rightarrow \text{int id, list};$

$\rightarrow \text{int } a, \text{ list};$

$\rightarrow \text{int } a, \text{ id};$

$\rightarrow \text{int } a, b;$





Disadvantages of backtracking

Backtracking is not a good way of parsing because when backtracking occurs we have to undo some changes like entries made in the symbol table have to be removed which creates a overhead. So a top down parser having no backtracking is more reasonable.

Difficulties with Top down Parsing

- 1) Left Recursion - a left recursive grammar can cause a top down parser to go into an infinite loop. So, elimination of left recursion is necessary.
- 2) Backtracking → As already said above backtracking is not reasonable for parsing.

Recursive Descent Parsing

A recursive descent parser is a top down method of syntax analysis which uses a collection of recursive procedures to process the input.

so

A parser that uses a set of recursive procedures to recognize its input with no backtracking is called a 'Recursive Descent Parser'.

Recursive descent parser accepts

a left recursion free grammar

and a left factoring free grammar.

Elimination of left recursion

$A \rightarrow A\alpha | B$ (left recursive grammar)

$A \rightarrow BA'$

$A' \rightarrow \alpha A' | \epsilon$

] left recursion free grammar

Elimination of left factoring

$A \rightarrow \alpha\beta | \alpha\gamma$

$A \rightarrow \alpha A'$

$A' \rightarrow \beta | \gamma$] left factored free grammar

Predictive Parser → Predictive parser is a special kind of recursive descent parsing with no backtracking.

A predictive parser can perform parsing by implementing a recursive descent parsing by handling the stack of activation records. It parses a left recursion free & a left factored tree grammar only.

Ex - we have grammar

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * F \mid F$$

$$F \rightarrow (E) \mid \text{id}$$

After eliminating left recursion from above grammar we get

$$\underline{E \rightarrow TE'}$$

$$E' \rightarrow +TE' \mid \epsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT' \mid \epsilon$$

$$F \rightarrow (E) \mid \text{id}$$

pausing because it implements the exclusive descent pausing by maintaining a stack of activation records.

Predictive Parser has four main parts -

1) Input Buffer

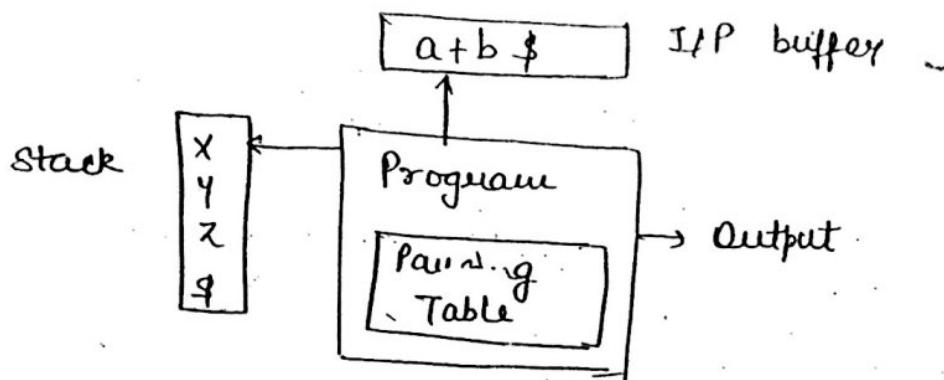
An input buffer contains the input string to be parsed followed by symbol \$.

2) Stack - Stack contains a sequence of grammar symbols with \$ on the bottom. Initially stack contains start symbol.

3) Pausing Table - It is a 2 dimensional array $M[A, a]$

where A is a non terminal
and a is a terminal

4) Output - the actions performed



(Model of Predictive Parser)

Parser is controlled by a program. This program determines (X) top element of stack and (a) current input symbol and performs following actions -

1) If $x = a = \$$ then parser halts & enters successful completion of parsing.

2) If $x = a \neq \$$ then parser pops x off the stack and advance the input pointer to next input symbol.

3) If x is a non terminal then it checks in parsing table. If $M[x, a] = \{x \rightarrow vuv\}$ then parser replace x by wvu with v on top of stack.

	id	$+$	$*$	$($	$)$	\sim
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'		$E' \rightarrow +TE'$			$E' \rightarrow e$	$E' \rightarrow e$
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'		$T' \rightarrow e$	$T' \rightarrow *FT'$		$T' \rightarrow e$	$T' \rightarrow e$
F	$F \rightarrow \text{id}$			$F \rightarrow (E)$		

and the input string is $\text{id} + \text{id} * \text{id}$
 that we can derive in following manner.

moves by predictive parser \rightarrow	Stack	Input	Action
	\$ E	$\text{id} + \text{id} * \text{id} \$$	
	\$ E' T	$\text{id} + \text{id} * \text{id} \$$	$E \rightarrow TE'$
	\$ E' T F	$\text{id} + \text{id} * \text{id} \$$	$T \rightarrow FT'$
	\$ E' T' id	$\text{id} + \text{id} + \text{id} \$$	$F \rightarrow \text{id}$
	\$ E' T'	$+ \text{id} * \text{id} \$$	POP
	\$ E'	$+ \text{id} * \text{id} \$$	$T' \rightarrow e$
	\$ E' T +	$+ \text{id} * \text{id} \$$	$E' \rightarrow +TE'$
	\$ E' T	$\text{id} * \text{id} \$$	POP
	\$ E' T' F	$\text{id} * \text{id} \$$	$T \rightarrow FT'$
	\$ E' T' id	$\text{id} * \text{id} \$$	$F \rightarrow \text{id}$
	\$ E' T'	$\text{id} * \text{id} \$$	POP
	\$ E' T' F *	$* \text{id} \$$	$T' \rightarrow *FT'$
	\$ E' T' F	$\text{id} \$$	
	\$ E' T' id	$\text{id} \$$	$F \rightarrow \text{id}$
	\$ E' T'	$\$$	$T' \rightarrow e$
	\$ E'	$\$$	$E' \rightarrow e$
	④	④	successful parser

Construction of Predictive Parsing Table



Finding FIRST and FOLLOW

Rules to calculate FIRST(x)

① If x is a terminal then

$$\text{FIRST}(x) = \{x\}$$

② If $x \rightarrow e$ is a production
then $\text{FIRST}(x) = \{e\}$

③ If x is a non terminal and $x \rightarrow a\alpha$
then $\text{FIRST}(x) = \{a\}$

④ If $x \rightarrow y_1 y_2 \dots y_k$ is a production for
some y_i and if

$$\text{FIRST}(y_i) = \{a\} \text{ and}$$

$$\text{FIRST}(y_1), \text{FIRST}(y_2), \dots, \text{FIRST}(y_{i-1})$$

contains e

$$\text{then } \text{FIRST}(x) = \{a\}$$

If $\text{FIRST}(y_j)$ contains e
where $j = 1, 2, \dots, i-1$

then add every non e symbol in $\text{FIRST}(y_i)$ to
 $\text{FIRST}(x)$.

If e is in $\text{FIRST}(y_j)$ for all $j = 1, 2, \dots, k$
then add G to $\text{FIRST}(x)$

1) If S is a start symbol

$$\text{FOLLOW}(S) = \{\$\}$$

2)

If there is a production $A \rightarrow \alpha B \beta$

$$\text{then } \text{FOLLOW}(B) = \text{FIRST}(\beta) - \{\$ \}$$

3)

If there is a production

$$A \rightarrow \alpha B$$

$$\text{or } A \rightarrow \alpha B \beta \quad \text{if } \text{FIRST}(B) \text{ contains } \$ \quad \text{else } \text{FOLLOW}(B) = \{\$ \}$$

$$\text{then } \text{FOLLOW}(B) = \text{FOLLOW}(A)$$

Ex

Let we have a following grammar.

$$E \rightarrow TE'$$

$$E' \rightarrow +TE' \mid \epsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT' \mid \epsilon$$

$$F \rightarrow CE \mid id$$

and we have to construct a parsing table for this grammar.

First calculate FIRST & FOLLOW for each non terminal in grammar.

1) $\text{FIRST}(E) = \text{FIRST}(T) = \text{FIRST}(R) = \{\$, \text{id}\}$
by rule 4 and 3 of FIRST

2) $\text{FIRST}(E') = \{+, \epsilon\}$ by rule 1 & 2

3) $\text{FIRST}(T') = \{*, \epsilon\}$ by rule 1 & 2

4) $\text{Follow}(E) = \{\$\}$ by rule 1 of follow

$\text{Follow}(E) = \text{Follow}(E')$ by rule 3

$\text{Follow}(E) = \{\)\}$ by rule 2

$A \rightarrow \alpha B,$
 $E \rightarrow T E'$
 $A \rightarrow \alpha B \beta$
 $F \rightarrow (E)$

so $\text{Follow}(E) = \text{Follow}(E') = \{\$,)\}$

5) $\text{Follow}(T) = \text{Follow}(T')$ by rule 3

$T \rightarrow FT'$

$\text{Follow}(T) = \text{FIRST}(E') - \epsilon$ (by rule 2)

$= \{+\}$

$A \rightarrow \alpha B \beta$
 $E' \rightarrow +TE'$

and

$A \rightarrow \epsilon B \beta$

$E' \rightarrow +TE'$ by rule 3

because $\text{FIRST}(E')$ contains ϵ

so $\text{Follow}(T) = \text{Follow}(E') = \{\$,)\}$

so $\text{Follow}(T) = \text{Follow}(T') = \{+,), \$\}$

$$A \rightarrow \alpha BB$$

$T' \rightarrow *FT'$ by rule 3

because FIRST(T') contains ϵ

$$\text{so } \text{FOLLOW}(F) = \text{FOLLOW}(T') = \{+,), \$\}$$

and

by rule 2

$$A \rightarrow \alpha BB$$

$$T' \rightarrow *FT'$$

~~because first~~

$$\text{FOLLOW}(F) = \text{FIRST}(T') - \epsilon$$

$$= \{ *\}$$

$$\text{so } \text{FOLLOW}(F) = \{ +, *,), \$\}$$

Algorithm for construction of a predictive parsing table ① remove left recursion

② ~~remove~~ left factor the grammar

i) for each production $A \rightarrow \alpha$ do following steps

a)

calculate FIRST(A)

if $\text{FIRST}(A) = \{\epsilon\}$ then add

$A \rightarrow \alpha$ to $M[A, \epsilon]$

b) if $\text{FIRST}(A)$ contains ϵ then for each terminal b in $\text{FOLLOW}(A)$ add $A \rightarrow \epsilon$ to $M[A, b]$

c) if ϵ is $\text{FIRST}(A)$ and $\$$ in $\text{FOLLOW}(A)$ then add $A \rightarrow \epsilon$ to $M[A, \$]$ also.

Now construct the parsing table

I production $E \rightarrow TE'$

$$\text{FIRST}(E) = \{\text{id}, \text{id}\}$$

so add $E \rightarrow TE'$ in table for terminal id, id

II production $E' \rightarrow +TE' | \epsilon$

$$\text{FIRST}(E') = \{+, \epsilon\}$$

because contains ϵ add $E' \rightarrow \epsilon$ to $\text{Follow}(E') = \{\), \$\}$

III production $T \rightarrow FT'$

$$\text{FIRST}(T) = \{\text{id}, \text{id}\} \text{ add } T \rightarrow FT' \text{ for } \{\text{id}, \text{id}\}$$

	id	$+$	$*$	$($	$)$	$\$$
E	$E \rightarrow TE'$				$E \rightarrow TE'$	
E'		$E' \rightarrow +TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow \text{id}$			$F \rightarrow (E)$		

IV Production $T' \rightarrow *FT' | \epsilon$

$$\text{FIRST}(T') = \{* \epsilon\} \text{ add } T' \rightarrow *FT' \text{ to } *$$

because it contains *

add $T' \rightarrow *FT'$ to $\text{Follow}(T') = \{+, \), \$\}$

V Production $F \rightarrow (E) | \text{id}$

$$\text{FIRST}(F) = \{\text{id}, \text{id}\}$$

add $F \rightarrow (E) | \text{id}$ for $($ and id

respectively

grammar if its LL(1) parsing table has no multiply defined entries. Ambiguous and left recursive grammar can not be a LL(1).

So LL(1) parsers are those parsers which uses only one input symbol to determine which production should be apply, and the grammar for which such parsers can be generated are called LL(1) grammar.

In LL(1) first L stands for left to right scanning of IPP.

Second L stands for producing left most derivation.

1 indicates that the next IPP symbol is used to decide next parsing process (i.e. length of lookahead is 1).

Properties of LL(1) Grammar

- 1) ~~No ambiguous or left recursive grammar can be LL(1)~~
- 2) ~~All the entries in the parsing-table are unique for LL(1) grammar~~

For a grammar to be LL(1) following conditions must be satisfied -

For production $A \rightarrow \alpha | \beta$:

{ FIRST(α) \cap FIRST(β) = \emptyset

and

if FIRST(β) contains ϵ and FIRST(α) does not contain ϵ
then

FIRST(α) \cap FOLLOW(A) = \emptyset

e.g.

$$\begin{array}{l} S \xrightarrow{\alpha} B\beta \\ S \xrightarrow{\beta} i c t s' \mid a \\ S' \xrightarrow{\epsilon} e \mid \epsilon \\ C \xrightarrow{\beta} b \end{array}$$

Parsing table

$$\begin{array}{l} \text{FIRST}(S) = \{ \text{fa} \} \\ \text{FIRST}(S') = \{ e, \epsilon \} \\ C \Rightarrow b \end{array}$$

$$\begin{array}{l} \text{Follow}(S) = \{ \text{eb} \} \\ \text{Follow}(S') = \{ \text{e} \} \end{array}$$

	a	b	e	i	t	\$
S	$S \rightarrow a$			$S \rightarrow i c t s'$		
S'			$S' \rightarrow e$			$S' \rightarrow e$
C		$C \Rightarrow b$				

above grammar is not a LL(1) grammar because its parsing table has multiple defined entries.

$$\text{FIRST}(S) = \{i, a\}$$

$$\text{FIRST}(S') = \{e, \epsilon\}$$

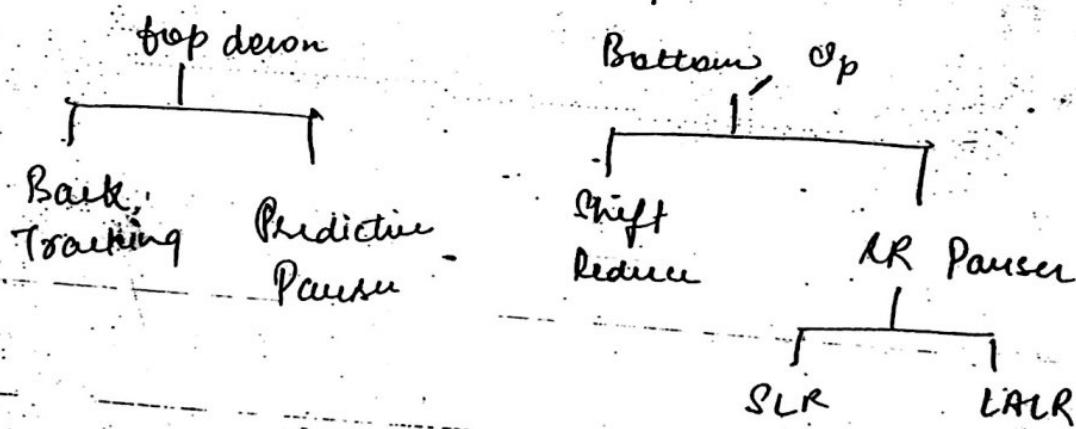
$$\text{FIRST}(C) = \{b\}$$

$$\text{Follow}(S) = \text{Follow}(S') = \{e, \$\}$$

$$\text{Follow}(S') = \{e, \$\}$$

$$\text{Follow}(C) = \{t\}$$

Types of Parser



LR parser

can be used for LR(0) grammar.

context free grammars.

This technique is called LR(k) parsing where

L stands for left to Right scanning of IIP.

R stands for constructing right most derivation.

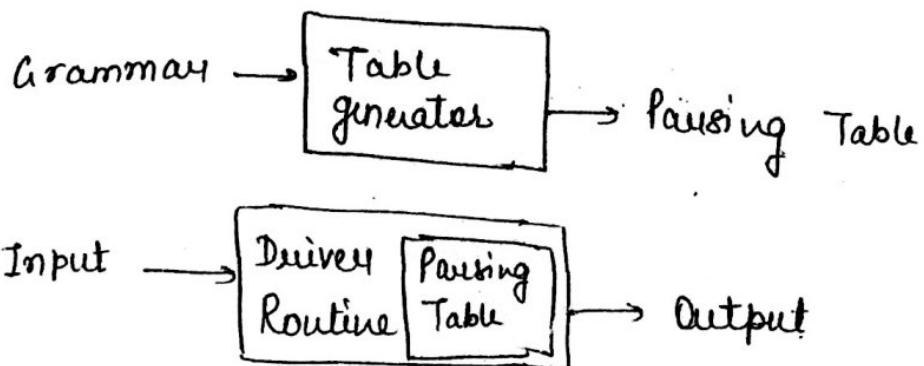
In reverse and k stands for no. of input symbols of look ahead.

So, LR parsers are bottom up parsers which scan the input from left to right and construct a right most derivation in reverse.

LR parser consists of two parts

- a driver routine
- a pausing table

driver routine is same for all parsers but pausing table changes from one parser to another.



There are three different techniques for generating LR parsing tables.

88

aminoNotes

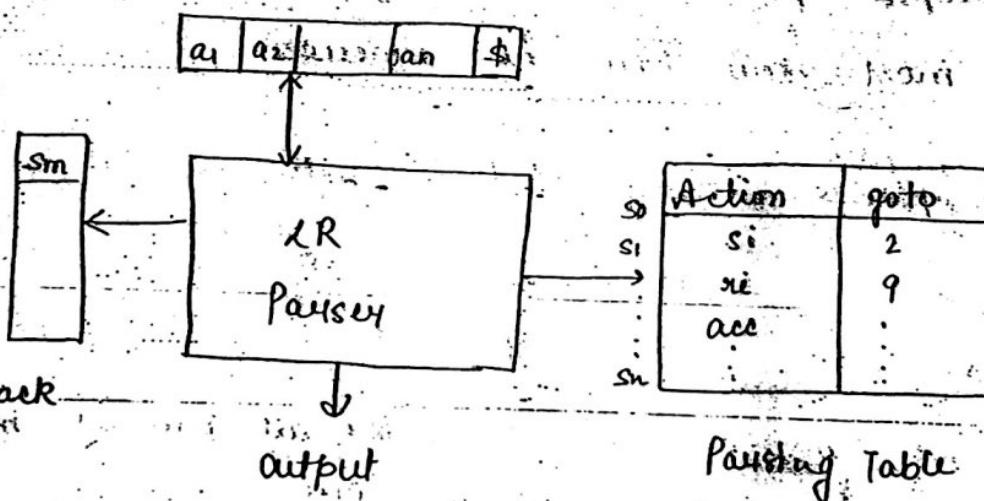
1) SLR - Simple LR

easiest to implement but fail for certain grammars.

2) Canonical LR - can work on a very large class of grammars, but expensive to implement.

3) LALR - Lookahead LR is intermediate in power b/w SLR and canonical LR.

Structure of LR Parser -



Configuration of LR Parser

- 1) Input
- 2) Stack
- 3) Parsing Table

Input is read from left to right one symbol at a time.

Stack can be in the following configuration:

$$s_0 X_1 s_1 X_2 s_2 X_3 \dots X_m s_m$$

where each X_i is a grammar symbol

and each s_i is a state

and s_m is on top of stack - top of stack

always contain a ~~stack~~ state.

Parsing table has two parts:

- 1) Action
- 2) Goto

Parsing action table entry: action $[s_m, a_i]$

for state s_m and input a_i can have four possible values

- a) s_i which means shift to state i .
- b) r_j reduce by using the j th rule.
- c) acc means accept
- d) error

GOTO takes a state and generates ~~some~~ arguments and produces a state.

Configuration of LR Parsers.

I component is stack contents &

II component is input

$(S_0 X_1 X_2 \dots X_m s_m, a_i a_{i+1} \dots a_n \$)$

next move of parser depends on current input symbol a_i and top of stack s_m .

ACTION $[s_m, a_i]$

Four types of moves can be done

1) If $\text{ACTION } [s_m, a_i] = \text{shift } s$

means parser executes a shift move and generate the following configuration

$(S_0 X_1 X_2 \dots X_m a_i s, a_{i+1} \dots a_n \$)$

here current IIP symbol a_i and state s is shifted to stack where $s = \text{GOTO } [s_m, a_i]$

then configuration will be

91 ~

21

$(S_0 X_1 S_1 X_2 \dots X_{m-2} S_{m-2} A \in, a_1 a_2 \dots a_n \$)$

here parser popped off symbols of
the stack and push A, left side of production
on to stack and state s.

where $s = \text{ACTION}[s_{m-2}, A]$

③ If $\text{ACTION}[s_m, a_i] = \text{accept} \Rightarrow$ parsing completed

④ If $\text{ACTION}[s_m, a_i] = \text{error}$

calls an error recovery routine.

$E \rightarrow T$
 $T \rightarrow T * F$
 $F \rightarrow F$
 $F \rightarrow (E)$
 $F \rightarrow id$

Let the given parsing table is

state	Action						Goto		
	id	+	*	()	\$	E	T	F
0	s5			s4			1	2	3
1				s6					
2					s2	s2			
3				s4	s4	s4			
4	s5			s4			8	2	3
5		s6	s6			s6			
6	s5			s4					13
7	s5			s9					10
8	s6			s11					
9		s7		s10	s10				
10		s3	s3		s3	s3			
11		s15	s15		s15	s15			

Now you have to parse a string

id * id * id using above ~~pseudo~~ parsing table by LR parser.

stack

9

o id 5

0 f 3

072

O. T. 2 - * 7

OT 2 * 1 id 5

OT 2 * 1 F 10

OT2

O. E. J.

DE 1 + 6

0 E 1 + 6 id 5

0'E 1 + 6 F 3

0 E 1 + 6 T 9

OE 1

accept

Input

$\text{id} \cdot * \text{id} \neq \text{id}$ \$

* id + id *

* id + id \$

* id + id \$

$\text{id} + \text{id} \$$

+ id.

+ 24

卷之三

卷之三

14

卷之三

卷之三

卷之三

卷之二

- LR parsing tables can be generated by three methods
- 1) SLR
 - 2) canonical LR
 - 3) LALR

① SLR Parser and construction of SLR parsing tables

For the construction of SLR parser, first we

→ make the ~~augmented grammar~~

means if a grammar has a start symbol S then we make a augmented grammar with new production

$$S' \rightarrow S$$

This new production indicates the parser when it should stop parsing

→ generate canonical collections of LR(0) items

Ex. $A \rightarrow xyz$ generates four items

$$A \rightarrow \cdot xyz$$

$$A \rightarrow x \cdot yz$$

$$A \rightarrow xy \cdot z$$

$$A \rightarrow xy \cdot z$$

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * F \mid F$$

$$F \rightarrow (E) \mid id$$

Follow (E) = { ,), }, {

(T) = *, +,), }

Augmented grammar for above grammar is

- 0 $\rightarrow E' \rightarrow E$
- 1 $\rightarrow E \rightarrow E + T$
- 2 $\rightarrow E \rightarrow T$
- 3 $\rightarrow T \rightarrow T * F$
- 4 $\rightarrow T \rightarrow F$
- 5 $\rightarrow F \rightarrow (E)$
- 6 $\rightarrow F \rightarrow id$

Now find I_0

$$E' \rightarrow E$$

$$E \rightarrow E + T$$

$$E \rightarrow T$$

$$T \rightarrow T * F$$

$$T \rightarrow F$$

$$F \rightarrow (E)$$

$$F \rightarrow id$$



$I_1 (I_0, E)$

$$E' \rightarrow E$$

$$E \rightarrow E + T$$

$I_2 (I_0, T) (I_4, T)$

$$E \rightarrow T$$

$$T \rightarrow T * F$$

I₃ (I₀, F) (I₄, F)
 (I₅, F)

(T → F.)

(I₄, F)

(F → id.)

—
—
—

I₄ (I₀, C) (I₉, C) (I₆, C)

F → (C, E.)

E → .E + T

- E → .T

T → .T * F

T → .F

F → .(E)

F → .id

I₆ (I₁, +) (I₈, +)

E → E + T

T → .T * F

T → .F

F → .(E)

F → .id

—
—
—

I₇ (I₂, *) (I₉, *)

T → T * F

F → .(E)

F → .id

—
—
—

I₈ (I₄, E)

F → (E, .)

E → E + T

(T, F)

I₁₁ (I₈,)

(F → (E), .)

~~I₁₀~~

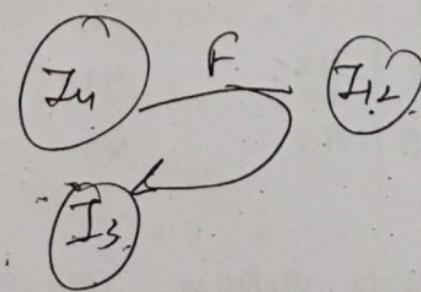
I₉ (I₆, T)

(E → E + T)

T → T * F

I₁₀ (I₇, F)

(T → T * F.)



Ex. Consider the following grammar and generate the SIR(1) table

$$S \rightarrow \lambda = R$$

$$S \rightarrow R$$

$$\lambda \rightarrow *R$$

$$\lambda \rightarrow \text{id}$$

$$R \rightarrow \lambda$$

Augmented grammar is

<u>$S' \rightarrow S$</u>	<u>I_0</u>	<u>$I_1(I_0, S)$</u>	<u>$I_3(I_0, R)$</u>
$S \rightarrow \lambda = R$	$S' \rightarrow .S$	$S' \rightarrow S.$ ✓	$S = .$ ✓
$S \rightarrow R$	<u>$S \rightarrow .\lambda = R$</u>	<u>$I_2(I_0, L)$</u>	<u>$I_4(I_3, *) (I_6, *)$</u> $(+, +)$
$\lambda \rightarrow *R$	$S \rightarrow .R$	$S \rightarrow \lambda = R$	$\lambda \rightarrow * .R$
$\lambda \rightarrow \text{id}$	$\lambda \rightarrow .\lambda$	$R \rightarrow \lambda .$ ✓	$R \rightarrow .L$
$R \rightarrow \lambda$	$R \rightarrow .\lambda$		$L \rightarrow * .R$
			$L \rightarrow .\text{id}$
		<u>$I_5(I_0, \text{id}) (I_6, \text{id})$</u>	
		<u>$I_7(I_4, R)$</u>	
	$\lambda \rightarrow \text{id} .$ ✓	$\lambda \rightarrow *R.$	
	<u>$I_6(I_2, =)$</u>	<u>$I_8(I_4, \lambda) (I_6, \lambda)$</u>	
	$S \rightarrow \lambda = .R$	$R \rightarrow \lambda .$	
	$R \rightarrow .\lambda$		
	$\lambda \rightarrow .*R$	<u>$I_9(I_6, R)$</u>	
	$\lambda \rightarrow .\text{id}$	$S \rightarrow \lambda = R .$	

Bottom Up Parsing

Bottom up Parsing - In this technique parsers build parse trees from the bottom (leaves) to the top (root).

IIP to the parser is being scanned from left to right ~~in one symbol at a time~~.

The Bottom up parsing method is called Shift-Reduce Parsing because it consists of shifting IIP symbols onto ~~my~~ stack until the right side of the production appears on top of the stack. The right side then can be replaced (reduced) by the symbol on the left side of the production.

One kind of shift Reduce Parser (SR-parser) is an operator precedence Parser.

Shift Reduce Parsing

SR parsing follows a process of right most derivation of a string w . This is a process of reducing a string w to the start symbol of a grammar.

e.g. let we have a grammar

$$S \rightarrow aAcBc$$

$$A \rightarrow Ab/b$$

$$B \rightarrow d$$

and we have to derive the string abcded

so in bottom up parsing the string abcded will be scanned from L to R. For some substrings that can match the right side of some production.

$\rightarrow a \underline{\bar{A}bcde}$
 $\rightarrow a \bar{A} cde$
 $\rightarrow a \bar{A} cBe$
 $\rightarrow \underline{s}$

Each replacement of the right side of production by left side is called Reduction.

Handle :- A handle of a right sentential form s is a production $A \rightarrow \beta$

because $s \xrightarrow{\text{RHS}} \alpha Aw \xrightarrow{\text{RHS}} \alpha \beta w$

e.g. $abbcde$ is a right sentential form whose handle is $A \rightarrow b$ at position 2.

If a grammar is unambiguous then every right sentential form of the grammar has exactly one handle.

Handle Pruning

Right most derivation in reverse is often called Canonical Reduction sequence is obtained by handle pruning.

means Handle pruning is the process of obtaining R.M.D in reverse by shift and reduction sequences.

Stack implementation of shift Reduce Parsing

In automating the SR parsing process by hand, there are two problems:

- how to locate a handle in right sentential form
- what production to choose if there are more than one production with the same right side

To implement SR parser use - stack

& IIP buffer

stack IIP
\$ w\$

IIP buffer holds the IIP string to be parsed with \$ at the end.

during L to R scanning of IIP parser shifts zero or more IIP symbols onto the stack, until it is ready to reduce a string B , on top of the stack. Then reduces string B by the left side of production and repeat this cycle until the stack contains start symbol & IIP is empty or an error has occurred.

in the grammar is

$$E \rightarrow E+E$$

$$E \rightarrow G \cdot E$$

$$E \rightarrow id_1 | id_2 | id_3$$

and we have to derive a string
 $id_1 = id_2 * id_3$ by LR parser

stack

TOP

action

\$

$id_1 + id_2 * id_3 \$$

shift

\$ id₁

$* id_2 * id_3 \$$

reduce, by $E \rightarrow id_1$

\$ E

$* id_2 * id_3 \$$

shift

\$ E +

$id_2 * id_3 \$$

shift

\$ E + id₂

$* id_3 \$$

reduce by $E \rightarrow id_2$

\$ E + E

$* id_3 \$$

shift

\$ E + E ~~id₃~~

$id_3 \$$

shift

\$ E + E ~~id₃~~

\$

reduce by $E \rightarrow id_3$

\$ E + E ~~G~~

\$

reduce by $E \rightarrow id_3$

- \$ E + E

\$

reduce by $E \rightarrow id_3$

\$ E

\$

reduce by $E \rightarrow E +$

Accept

Operations of LR Parsers

- ① Shift - In this operation next IIP symbol is shifted to the top of the stack.
- ② Reduce - If the parser knows the right end of the hand at top of the stack and if it locate the left end of the hand then it decides that with what non terminal the stack should be replaced.
- ③ Accept - This operation shows the successful parsing.
- ④ Error - Parser discovers that syntax error has occurred and then calls an error recovery routine.

Example :-

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * F \mid F$$

$$F \rightarrow id$$

for IIP string: id + id * id

the R.H.D is: $E \rightarrow E + T$

$$\rightarrow E + T * F$$

$$\rightarrow E + T * id$$

$$\rightarrow E + F * id$$

$$\rightarrow E + id * id$$

$$\rightarrow T + id * id$$

$$\rightarrow F + id * id$$

$$\rightarrow id + id * id$$

stack

\$
\$ id
\$ F
\$ T
\$ E
\$ E+
\$ E+id
\$ EFF
\$ E+T
\$ E+T*
\$ E+T* id
\$ E+T* F
\$ E+T
\$ E

SIP buffer

id + id * id \$
id * id \$
* id \$
* id \$
id \$
\$
\$
\$

Action

shift
Reduce by F → id
Reduce by T → F
Reduce by E → T
shift
shift
Reduce F → id
Reduce T → F
shift
shift
Reduce F → id
Reduce by T → T*F
Reduce by E → E+T
Accept

Conflicts During Shift Reduce Parsing

There are grammars for which shift reduce parser can not be used. For such grammar SR parser reach in a situation of conflict.

There are two types of conflict

(1) Shift Reduce Conflict - This conflict arises when the parser reach in a configuration in which it can not decide whether to shift or to reduce.

e.g. $\text{stmt} \rightarrow \text{if expr then stmt}$
 $\qquad\qquad\qquad \rightarrow \text{if expr then stmt else stmt}$
 $\qquad\qquad\qquad \rightarrow \text{other}$

Let the parser is in following configuration

Stack

if expr then stmt

if

else \$

It is a situation of shift reduce conflict because we can not say whether we should reduce if expr then stmt or shift else.

(2) Reduce-Reduce Conflict - This conflict arises when SR parser cannot decide that which of the several reduction should be used for.

eg. 1. Stmt \rightarrow id (parameter-list)

2. Stmt \rightarrow expr := expr

3. parameter-list \rightarrow parameter-list , parameter

4. parameter-list \rightarrow parameter-

5. parameter \rightarrow id

6. expr \rightarrow id (expr-list)

7. expr \rightarrow id

8. expr-list \rightarrow expr-list , expr

9. expr-list \rightarrow expr

eg. A(I,J)

token stream id(id, id)

stack

id(id

IIP

, id) \$

In this situation we can't say that whether we should reduce id by production 5 or 7.

Correct choice is production 5 if A is a procedure

7 if A is an array

Operator Precedence Parsing

This is a form of shift reduce parsing which is used to parse operator grammar. This parsing is especially used for parsing expressions.

Operator precedence grammar

These grammars should have property

- that no production right side is ϵ
- no production has two adjacent non terminals.

e.g. $E \rightarrow EAE \mid (E) \mid -E \mid id$
 $A \rightarrow + \mid - \mid * \mid / \mid \%$

Is not an operator grammar, because the right side EAE has two consecutive non terminals.

Operator grammar can be obtained by substituting the value of A in L production

we get

$$E \rightarrow E+E \mid E-E \mid E*E \mid E/E \mid E\%E \mid (E) \mid -E \mid id$$

Operator Precedence parsing uses the information about the precedence & associativity of operators to guide the parse tree

$$\begin{array}{l} q = v \\ qtv = v \end{array}$$



In operator precedence parsing we use three disjoint precedence relations between pair of terminals.

\leftarrow
 \equiv
 \rightarrow

These relations guide the selection of handles and have following meaning.

- $a \leftarrow b$ a yields precedence to b
 $a \equiv b$ a has same precedence as b
 $a \rightarrow b$ a takes precedence over b.

There are two methods of generating precedence relations

- ① By determining the associativity & precedence of operators
- ② By finding unambiguous grammar

(will be discussed later)

Using Precedence Relations for Parsing

Let we have operator precedence relations as

	id	+	*	\$
id		>	>	>
+	<	>	<	>
*	<	>	>	>
\$	<	<	<	

These precedence relations are used to delimit the handle of a right sentential form.

e.g. Let the string is $\$ \leftarrow id \rightarrow + \leftarrow id \rightarrow * \leftarrow id \rightarrow \$$

then the string with precedence relations inserted is

$\$ \leftarrow id \rightarrow + \leftarrow id \rightarrow * \leftarrow id \rightarrow \$$

Process of finding handle

- 1) Scan the string from L to R until \rightarrow is encountered.
- 2) Then scan backwards (to the left) over any \equiv until a $<$ is encountered. Then scan backwards to $\$$.
- 3) Handle contains everything between $<$ and \rightarrow .

In above string on scanning from L to R until \rightarrow and then by scanning backwards to Left until $<$ we get id.

So the handle is the first id

we then reduce id to E

and we get right sentential form $E + id \rightarrow id$

Similarly after reducing the remaining id's to E

we get $E + E \# E$

now consider the string $\$ + * \$$ by deleting non-terminal from above string

then

$\$ \leftarrow + \leftarrow * \rightarrow \$$

Rules to calculate Operator Precedence Relations

I Method

Relations from Associativity and Precedence

Rule 1:- If operator θ_1 has higher precedence than operator θ_2

then generate $\theta_1 > \theta_2$

e.g. b/w + and *

$\theta_2 < \theta_1$

relations will be $+ < *$ and $* > +$

Rule 2:- If θ_1 and θ_2 are operators of equal precedence,

then generate $\theta_1 > \theta_2$

and $\theta_2 > \theta_1$ {if operators are left associative}

or make

$\theta_1 < \theta_2$

$\theta_2 < \theta_1$ {if they are right associative}

e.g. relation b/w + and + and + and - can be

$+ > +$

$+ > -$

$- > +$

$- > -$

because they are left associative

$\uparrow < \uparrow$
is right associative so

it shows that in $E-E+E$, $(E-E)$ will be selected first as handle.

and in $E \uparrow E+E$ will have the last ~~E~~
 $E \uparrow E$ selected first.

Rule 3: Make $\theta < id$
 $id > \theta$
 $\theta < ($
 $(< \theta$
 $) > \theta$
 $\theta >)$
 $\theta > \$$
 $\$ < \theta$

and also

(\div)	$\$ < ($	$\$ < id$
$(< ($	$id > \$$	$) > \$$
$(< id$	$id >)$	$) >)$

By these rules operator precedence relations by following grammar can be created

Grammar $E \rightarrow E+E \mid E-E \mid E * E \mid E/E \mid E \uparrow E \mid (\underline{E}) \mid -E \mid id$

	$+$	$-$	$*$	$/$	\uparrow	$($	$)$	id	$\$$
$+$	$>$	$>$	$<$	$<$	$<$	$>$	$<$	$<$	$>$
$-$	$>$	$>$	$<$	$<$	$<$	$>$	$<$	$<$	$>$
$*$	$>$	$>$	$>$	$<$	$<$	$>$	$<$	$<$	$>$
$/$	$>$	$>$	$>$	$<$	$<$	$>$	$<$	$<-$	$>$
\uparrow	$>$	$>$	$>$	$<$	$<$	$>$	$<$	$<->$	$>$
$($	$<$	$<$	$<$	$<$	$<$	$<$	$=$	$<$	$>$
$)$	$>$	$>$	$>$	$>$	$>$	$>$	$>$	$>$	$>$
id	$>$	$>$	$>$	$>$	$-$	$>$	$-$	$-$	$>$
$\$$	$<$	$<$	$<$	$<$	$<$	$<$	$=$	$<$	$-$

<u>Stack</u>	<u>Input</u>
\$	id * (id + id) - id / id \$
\$ < id	* (id + id) - id / id \$
\$ E:	* (id + id) - id / id \$
\$ < E *	(id + id) - id / id \$
\$ < E * < ((d + id) - id / id \$
\$ < E * < (< cd	+ id) - id / id \$
\$ < E * < (E	+ (cd) - id / id \$
\$ < E * < (< E *	(d) - id / id \$
\$ < E * < (< E + < id) - id / id \$
\$ < E * < (< E + < E) - id / id \$
\$ < E * < (E) - cd / id \$
\$ < E * < (= E	- id / id \$
\$ < E * E	- id / id \$
\$ E	- id / id \$
\$ < E -	id / id \$
\$ < E - < id	/ id \$
\$ < E - E	/ id \$
\$ < E - < E /	cd \$
\$ < E - < E / < id	\$
\$ < E - < E / E	\$
\$ < E - E	\$
\$ (E)	\$

Accept

Rules

126

- ① $a \doteq b$ if there is a right side of a production of the form $\alpha a \beta b \gamma$ where β is either ϵ or a single non-terminal.

means $a \doteq b$ if a appears immediately to the left of b or if they appear separated by one non-terminal.
eg. $S \rightarrow c t S t \Rightarrow \varnothing \doteq t$ & $t \doteq \epsilon$

- ② $a \lessdot b$ if for some non-terminal A there is a right side of the form $\alpha a A B$ and $A \rightarrow \gamma b \delta$ where γ is either ϵ or single non-terminal.

means $a \lessdot b$ if a non-terminal A appears immediately to the right of a and derives a string in which b is the first terminal symbol.

eg. $S \rightarrow c t S t \& C \rightarrow b$

so $\varnothing \lessdot b$

- ③ $a \succ b$ if for some non-terminal A there is a right side of the form $\alpha A b \beta$ and $A \rightarrow \gamma a \delta$, where γ is either ϵ or single non-terminal.

means $a \succ b$ if a non-terminal appears immediately to the left of b derives a string whose last terminal is a .

eg. $S \rightarrow c t S t$ and $C \rightarrow b$.

$b \succ t$

relation among all the terminals involved in the grammar.

An operator precedence grammar is an ϵ free grammar in which the precedence relations $<$, \doteq and $>$ are disjoint.

means for any pair of terminals a and b , never more than one of the relations $a < b$, $a \doteq b$ and $a > b$ is true.

e.g. $E \rightarrow E+E \mid E \cdot E \mid (E) \mid id$

This grammar is not an operator precedence grammar as by the rule (ii) and (iii) two precedence relations hold b/w + and \cdot

by rule (ii) in $E \rightarrow E+E$

$\alpha A B$

$a = +$
 $A = E$

and $E \rightarrow E+E$

so E derives string. so
first terminal is $+$

so $+ < +$

by rule (iii) in $E \rightarrow E+E$

$\alpha A B \beta$

so $+ > +$

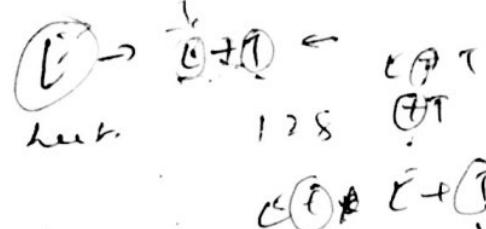
$A = E$ $E \rightarrow E+E$

$b = +$ last of E is $+$

Therefore above grammar is not an operator precedence grammar.

transform this grammar into operator precedence
unambiguous grammar

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow id \mid (E) \end{aligned}$$



Now find the precedence relations by these three rules. These rules can be modified as follows:-
Calculate FIRST or LEADING & LAST or TRAILING all the non terminals.

LEADING(A) = {a | A $\Rightarrow^* \alpha \beta$, where β is either ϵ or a single non terminal}

LNG or LAST(A) = {a | A $\Rightarrow^* \alpha \beta$, where β is either ϵ or a single non terminal}

now find precedence relations

1) = in above grammar and by rule (i) this relation holds b/w (ans) only
so (=)

2) < - we find a terminal a immediately to the left of a non terminal to play the role of a & A then make a <. LEADING(A) also make & <. LEAD(^{start symbol})

3) > - we find a terminal b immediately to the right of a non terminal A then make

$$b \rightarrow \text{LAST}(A)$$

$$\text{LAST}(A) \rightarrow b \quad \text{also make } \text{LAST}(A \text{ out symbol}) \rightarrow b$$

Further, if a is in LEADING(B), and there is a production of the form $A \rightarrow B\alpha$

- 29

then a is in LEADING(A)

Consider the grammar

$$E \rightarrow E + T \mid T$$

$$\begin{aligned} T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid id \end{aligned}$$

and find preed relations

so by rule 1 \Rightarrow (\Leftarrow)

by rule 2 \Rightarrow $<$

$+ < \cdot \text{leading}(T)$

$* < \cdot \text{leading}(F)$

$(< \cdot \text{leading}(E)$

$$ < \cdot \text{leading}(F)$

so we get

$+ < \cdot *$

$+ < \cdot ($

$+ < \cdot id$

$* < \cdot *$

$* < \cdot *$

$* < \cdot +$

$* < \cdot ($

$* < \cdot id$

$(< \cdot *$

$(< \cdot +$

$(< \cdot ($

$(< \cdot id$

	LEADING	TRAILING
E	$*, +, (, id$	$+, +,)$
T	$*, (, id$	$+,), id$
F	$(, id$	$id,$

by rule 3 \Rightarrow

$\text{last}(E) \Rightarrow +$

$\text{last}(T) \Rightarrow *$

$\text{last}(F) \Rightarrow)$

we get

$* \Rightarrow +$	$* \Rightarrow *$	$* \Rightarrow)$	$* \Rightarrow *$
$+ \Rightarrow +$	$+ \Rightarrow *$	$+ \Rightarrow)$	$+ \Rightarrow *$
$) \Rightarrow +$	$) \Rightarrow *$	$) \Rightarrow)$	$) \Rightarrow *$
$id \Rightarrow +$	$id \Rightarrow *$	$id \Rightarrow)$	$id \Rightarrow *$

now we get the table as

130

row\col	+	*	()	id	\$
+	>	<	<	>	<	>
*	>	>	<	>	<	>
(<	<	<	=	<	
)	>	>		>		
id	>	>			>	
\$	<	<	<	>	<	

Operator Precedence Parsing Algorithm

repeat forever

if only \$ is on the stack and \$ is on the input then accept and break

else

begin

let a be the top most terminal symbol on the stack, and let b be the current input symbol

if $a < b$ or $a = b$ then shift b on to the stack

else if $a > b$ then /* reduce */

repeat pop the stack

until the top stack terminal is related by $<$ to the terminal most recently popped

else call the error correcting routine

end

Operator precedence parsers have no need to store one large table of precedence relations because this table can be encoded by only two functions, f and g , which map terminal symbols to the integers.

We select f and g for symbol a and b like

$$f(a) < g(b) \text{ whenever } a \leq b$$

$$f(a) = g(b) \quad " \quad a \stackrel{?}{=} b$$

$$f(a) > g(b), \quad " \quad a \geq b$$

Method for finding precedence functions for a tail

① Create symbols fa and ga for each a that is a terminal or $\$$.

② Partition the created symbols as into as many groups as possible.

means if $a \stackrel{?}{=} b$ then fa and gb are in same group.

if $a \stackrel{?}{=} b$ and $c \stackrel{?}{=} b$

then fa and fc are in same group as they are same groups as gb .

③ Create a directed graph whose nodes are the groups found in step 2.

if $a \leq b$ then create an edge from gb to fa

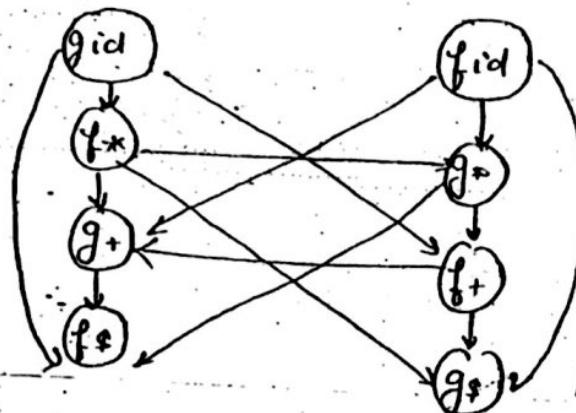
if $a \geq b$ then create an edge from fa to gb

Q) If the graph constructed in 3 has a cycle then no precedence function exists. If there are no cycles then let $f(a)$ be the length of the longest path beginning at the group of fa ; let $g(a)$ be the length of the longest path from the group of ga .

e.g. Generate the precedence function table for

id	+	*	\$
id	>	>	>
+	<	>	<
*	<	>	<
\$	<	<	<

because there are no relationships, so each symbol is in group by itself



then generate directed graph.

as there are no cycles, so precedence function exists now. find out the longest path from each group and create the precedence function table.

	id	+	*	\$
f	4	2	4	0
g	5	1	3	0

(6)

As we know that after syntax analysis, the semantic analysis of input program is needed.

To perform semantic analysis in a syntax directed fashion, we attach some semantic rules with each production in context free grammar.

So, we can say

"Syntax directed translation scheme is a process of allowing semantic actions to be attached to the productions of context free grammar."

Two notations are mainly used for associating semantic rules with productions

- Syntax directed Definition
- Translation Scheme

Syntax Directed Definition (SDD)

Syntax directed definition is a generalisation of a context free grammar in which each grammar symbol either terminal or non terminal has associated set of attributes.

The attribute can be anything like a string, a number, a type or a memory location etc.

Definition - A syntax directed definition is a generalization of a context free grammar where each production $A \rightarrow \alpha$ is associated with a set of semantic rules of the form

$$x := f(a_1, a_2, \dots, a_k)$$

where x is an attribute

These attributes are of two types:

Attributes are associated with grammar symbols and rules are associated with production.

1) Synthesized attributes

If $A \rightarrow \alpha$ is a production and

$$x := f(a_1, a_2, \dots, a_k) \text{ is a semantic rule}$$

then x is a synthesized attribute of A and a_1, a_2, \dots, a_k are attributes belonging to the grammar symbols of α .

We can say that synthesized attribute x can be evaluated during a single bottom up traversal of the parse tree.

Production

$$\text{eg. } E \rightarrow E^{(1)} + E^{(2)}$$

$$T \rightarrow F$$

Semantic Rule

$$\{ E.\text{VAL} := E^{(1)}.VAL + E^{(2)}.VAL \}$$

$$\{ T.VAL := F.VAL \}$$

A synthesized attribute is defined in terms of attribute values at the children of node and node itself.

135

The evaluation of the semantic rules proceed along with the construction of the parse tree. It defines the values of the attributes at the nodes in the parse tree for the input string.

Annotated Parse tree - A parse tree containing the values of attributes at each node of the tree is called an annotated parse tree.

Productions

eg.

$$E \rightarrow E_1 + T$$

$$E \rightarrow E_1 - T$$

$$E \rightarrow T$$

$$T \rightarrow T_1 * F$$

$$T \rightarrow T_1 / F$$

$$T \rightarrow F$$

$$F \rightarrow (E)$$

$$F \rightarrow \text{num}$$

Semantic Rules

$$\{ E.\text{VAL} := E_1.\text{VAL} + T.\text{VAL} \}$$

$$\{ E.\text{VAL} := E_1.\text{VAL} - T.\text{VAL} \}$$

$$\{ E.\text{VAL} := T.\text{VAL} \}$$

$$\{ T.\text{VAL} := T_1.\text{VAL} * F.\text{VAL} \}$$

$$\{ T.\text{VAL} := T_1.\text{VAL} / F.\text{VAL} \}$$

$$\{ T.\text{VAL} := F.\text{VAL} \}$$

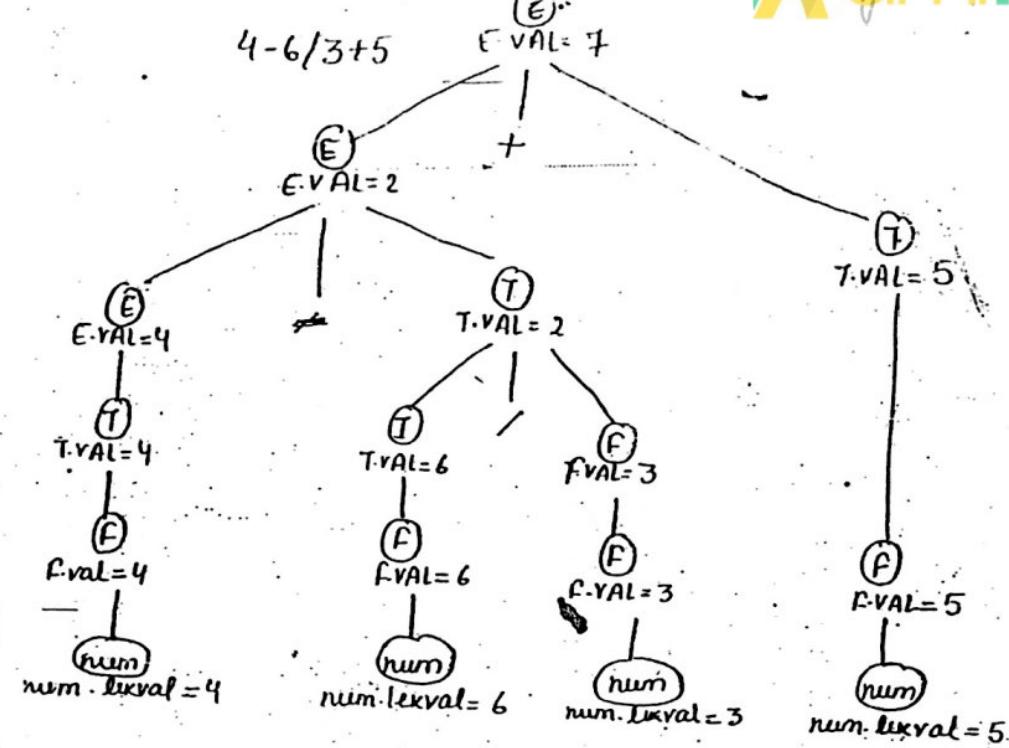
$$\{ F.\text{VAL} := E.\text{VAL} \}$$

$$\{ F.\text{VAL} := \text{num, textual} \}$$

S-attributed definition - In above syntax directed definition, every grammar symbol has synthesized attributes only so it is called S-attributed definition.

tree G-vAt is called the translation

now annotate the parse tree for string



Annotated Parse tree for string $4-6/3+5$

2) Inherited Attributes -

If $A \rightarrow \alpha$ is a production and

$\alpha = f(a_1, a_2, \dots, a_r)$ is a semantic rule

then α is an inherited attribute of one of the grammar symbols in α and a_1, a_2, \dots, a_r are attributes belonging to A or α .

e.g. $A \rightarrow X Y Z \quad \{ Y \cdot \text{VAL} = 2 * A \cdot \text{VAL} \}$

we can say that inherited attributes are those whose initial value at a node in the parse tree is defined in terms of the attributes of the parent and/or siblings.

Inherited attribute at node N is defined only
in terms of attribute values at N's parent, N itself
and N's siblings

- 137

eg. $T \rightarrow FT'$

Semantic Rules

$$T.\text{inh} = F.\text{val}$$

$$T.\text{val} = T'.\text{syn}$$

$$T' \rightarrow *FT'$$

$$T'.\text{inh} = T.\text{inh} \times F.\text{val}$$

$$T'.\text{syn} = T'.\text{syn}$$

$$T' \rightarrow \epsilon$$

$$T'.\text{syn} = T'.\text{inh}$$

$$F \rightarrow \text{digit}$$

$$F.\text{val} = \text{digit}.$$

Implementation of Syntax directed translators

S-attribute definitions may be implemented
with a bottom up parsing strategy.

Syntax directed translators can be implemented
by using extra fields in the parser stack
entries corresponding to the grammar symbols.

So stack is implemented by two pair of arrays -

STATE and VAL

If ith STATE symbol is E then VAL[i] will
hold the value of the translation E.VAL

TOP is a pointer to the current TOP of the
stack.

eg. $A \rightarrow xyz$
 before the reduction of xyz to A - the value
 of the translation of x is in $VAL[TOP]$
 and that of y is in $VAL[TOP-1]$,
 x is in $VAL[TOP-2]$
 after the reduction the value of A .VAL,
 will appear in $VAL[TOP]$.

state	val
X	X.VAL
Y	Y.VAL
Z	Z.VAL
:	:

stack before reduction

state	val
A	A.VAL
:	:

stack after reduction top=top

Ex- let we have a syntax directed translation scheme of "desk calculator".

$$\begin{array}{l}
 S \rightarrow E \\
 E \rightarrow E+E \\
 E \rightarrow E * E \\
 E \rightarrow (E) \\
 E \rightarrow I \\
 I \rightarrow I \text{ digit} \\
 I \rightarrow \text{digit}
 \end{array}$$

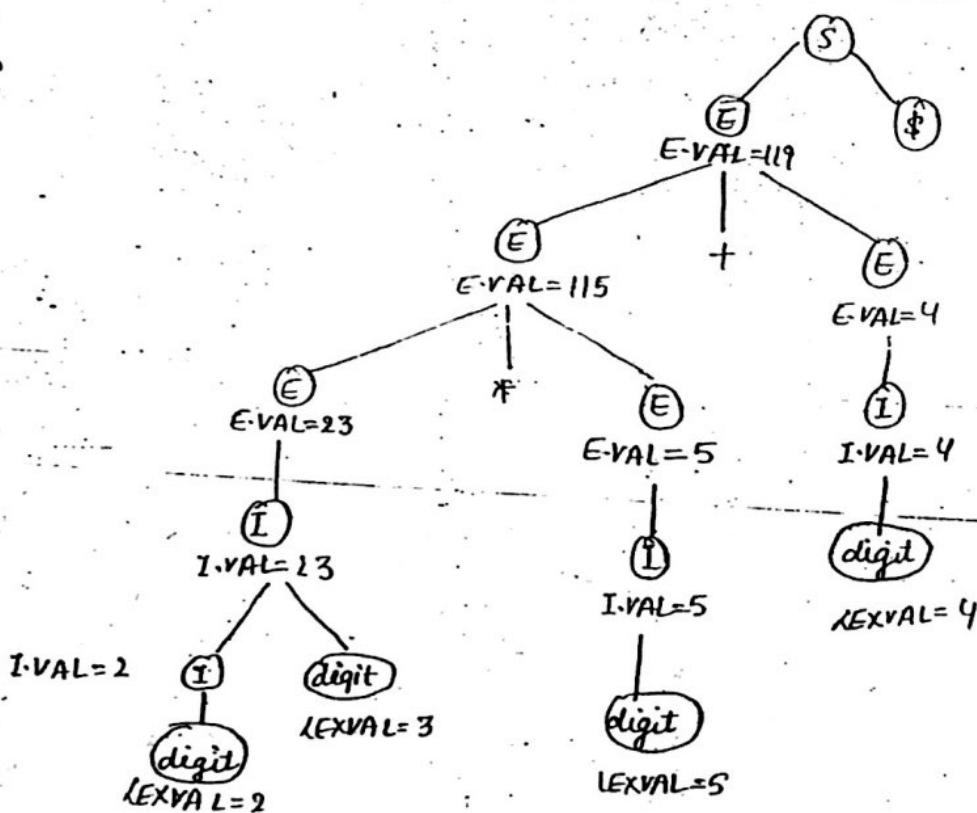
calculator

Production

- | | | |
|-----|-----------------------------------|--|
| (1) | $S \rightarrow E \$$ | Semantic action
{ print E.VAL } |
| (2) | $E \rightarrow E^{(1)} + E^{(2)}$ | { E.VAL := $E^{(1)}.VAL + E^{(2)}.VAL$ } |
| (3) | $E \rightarrow E^{(1)} * E^{(2)}$ | { E.VAL := $E^{(1)}.VAL * E^{(2)}.VAL$ } |
| (4) | $E \rightarrow (E^{(1)})$ | { E.VAL := $E^{(1)}.VAL$ } |
| (5) | $E \rightarrow I$ | { E.VAL := I.VAL } |
| (6) | $I \rightarrow I^{(1)} digit$ | { I.VAL := $10 * I^{(1)}.VAL + LEXVAL$ } |
| (7) | $I \rightarrow digit$ | { I.VAL := LEXVAL } |

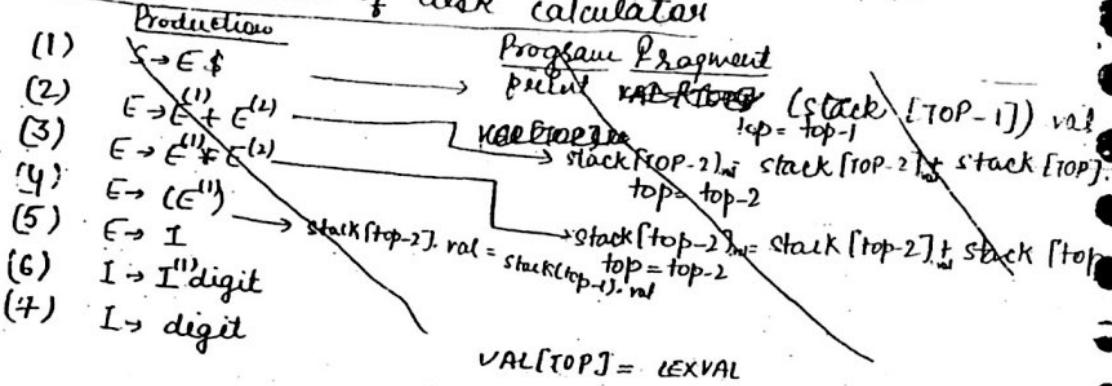
23 * 5 + 4

Parse tree with translation



Implementation of desk calculator

140

Sequence of moves

Input	state	Val	Productions used
(1) 23 * 5 + 4 \$	-	-	
(2) 3 * 5 + 4 \$	2	-	
(3) 3 * 5 + 4 \$	I	2	
(4) * 5 + 4 \$	I3	2	$I \rightarrow \text{digit}$
(5) * 5 + 4 \$	I	23	
(6) * 5 + 4 \$	E*	23	$I \rightarrow I \text{ digit}$
(7) 5 + 4 \$	E*	23	$E \rightarrow I$
(8) + 4 \$	E*5	23	
(9) + 4 \$	E*I	23.5	$I \rightarrow \text{digit}$
(10) + 4 \$	E*E	23.5	$E \rightarrow I$
(11) + 4 \$	E	115	$E \rightarrow E * E$
(12) 4 \$	E+	115	
(13) \$	E+4	115	
(14) \$	E+I	115.4	$I \rightarrow \text{digit}$
(15) \$	E+E	115.4	$E \rightarrow I$
(16)	E	119	$E \rightarrow E + E$
(17)	E\$	119	
(18)	-	119	

Implementation of desk calculatorProductions

- (1) $S \rightarrow E\$$
 - (2) $E \rightarrow E^{(1)} + E^{(2)}$
 - (3) $E \rightarrow E^{(1)} * E^{(2)}$
 - (4) $E \rightarrow (E^{(1)})$
 - (5) $E \rightarrow I$
 - (6) $I \rightarrow I^{(1)} \text{ digit}$
 - (7) $I \rightarrow \text{ digit}$
-

Program fragment

```

print (stack[top-1].val)
stack[top-2].val = stack[top-2].val + stack[1]
stack[top-2].val = stack[top-2].val + stack[top-1]
stack[top-2].val = stack[top-1].val
stack[top-1].val = 10 * stack[top-1] + stack[top]

```



Q Give the parse tree and translation for the expression $(4*7+19)*2$ by following grammar

$$\begin{aligned} S &\rightarrow E\$ \\ E &\rightarrow E+E \\ E &\rightarrow E*E \\ E &\rightarrow (E) \\ E &\rightarrow I \\ I &\rightarrow I \text{ digit} \\ I &\rightarrow \text{ digit} \end{aligned}$$

Syntax directed translation scheme -

$$\begin{aligned} S &\rightarrow E\$ \\ E &\rightarrow E^{(1)} + E^{(2)} \\ E &\rightarrow E^{(1)} * E^{(2)} \\ E &\rightarrow (E) \\ E &\rightarrow I \\ I &\rightarrow I^{(0)} \text{ digit} \\ I &\rightarrow \text{ digit} \end{aligned}$$

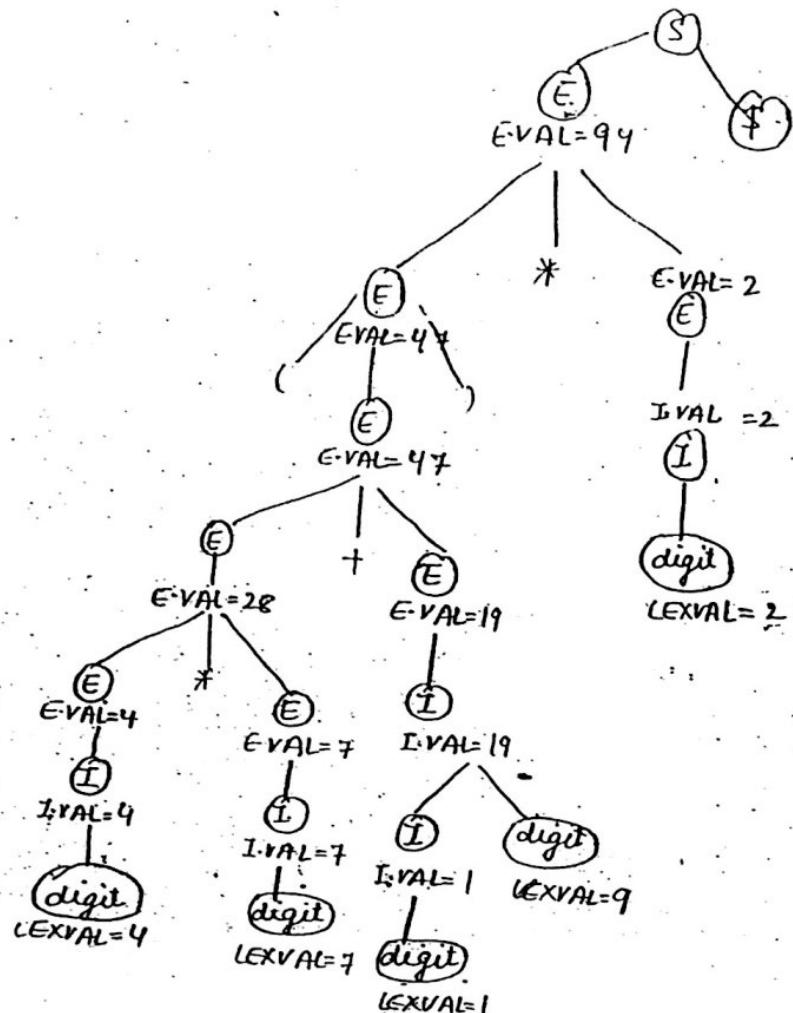
$$\begin{aligned} &\{ \text{put } E.\text{VAL} \} \\ \{ E.\text{VAL} := E^{(1)}. \text{VAL} + E^{(2)}. \text{VAL} \} \\ \{ E.\text{VAL} := E^{(1)}. \text{VAL} * E^{(2)}. \text{VAL} \} \\ \{ E.\text{VAL} := E^{(0)}. \text{VAL} \} \\ \{ E.\text{VAL} := I.\text{VAL} \} \\ \{ I.\text{VAL} := \text{op} I^{(0)}. \text{VAL} + \text{LEXVAL} \} \\ \{ I.\text{VAL} = \text{LEXVAL} \} \end{aligned}$$

Implementation

$$\begin{aligned} S &\rightarrow E\$ \\ E &\rightarrow E+E \\ E &\rightarrow E*E \\ E &\rightarrow (E) \\ E &\rightarrow I \\ I &\rightarrow I \text{ digit} \\ I &\rightarrow \text{ digit} \end{aligned}$$

Program fragment.

$$\begin{aligned} &\text{point stack [TOP]} \\ \text{stack [TOP]} &= \text{stack [top-2].t}, \text{stack [top].val} \\ \text{stack [top-2].val} &= \text{stack [top-2].val} + \text{stack [top-1].val} \\ \text{stack [top-2].val} &= \text{stack [top-1].val} \\ \text{stack [top-1].val} &= \text{stack [top-1].val} * \text{stack [top].val} \end{aligned}$$



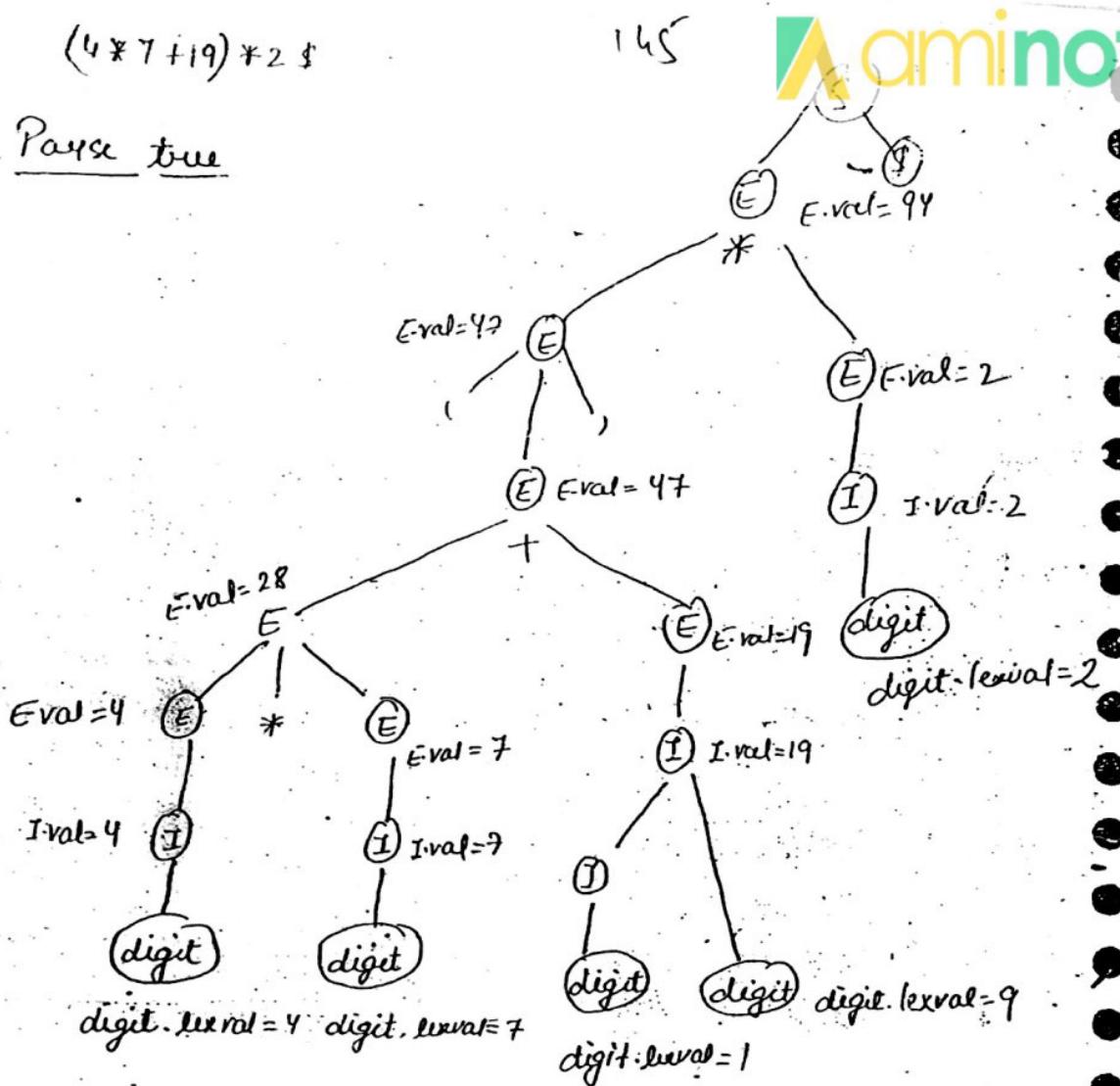
Sequence of movesEvaluate $(4 * 7 + 19) * 2\$$

Input	state	value	Productions used
$(4 * 7 + 19) * 2\$$	-	-	-
1) $(4 * 7 + 19) * 2\$$	(-	-
2) $* 7 + 19) * 2\$$	(4	-	-
3) $* 7 + 19) * 2\$$	(I	-4	$I \rightarrow \text{digit}$
4) $* 7 + 19) * 2\$$	(E	-4	$E \rightarrow I$
5) $7 + 19) * 2\$$	(E*	-4	-
6) $+ 19) * 2\$$	(E+E)	-4	-
7) $+ 19) * 2\$$	(E+E*)	-4	-
8) $+ 19) * 2\$$	(E+E*	4-7	$I \rightarrow \text{digit}$
9) $+ 19) * 2\$$	(E	4-7	$E \rightarrow I$
10) $19) * 2\$$	(E+	28	$E \rightarrow E * E$
11) $9) * 2\$$	(E+I	28	-
12) $9) * 2\$$	(E+I	28-1	-
13) $) * 2\$$	(E+I9	28-1	$E \rightarrow \text{digit}$
14) $) * 2\$$	(E+I	28-1	-
15) $) * 2\$$	(E+E	28-19	$I \rightarrow \text{I digit}$
16) $) * 2\$$	(E	28-19	$E \rightarrow I$
17) $* 2\$$	(E)	47	$E \rightarrow E+E$
18) $* 2\$$	E	47	$E \rightarrow (E)$
19) $2\$$	E*	47-	-
20) $2\$$	E*2	47-	-
21) $2\$$	E*I	47-2	$I \rightarrow \text{digit}$
22) $2\$$	E*E	47-2	$E \rightarrow I$
23) $2\$$	E	94	$E \rightarrow E * E$
24) $2\$$	E\$	94	-
25) $2\$$	E	-	-

$$(4 * 7 + 19) * 2 \$$$

145

Parse tree



Intermediate code - In most of the compilers source code is translated into a language which is intermediate in complexity between a high level programming language and machine code. Such code is called intermediate code or intermediate text.

Most common representations of source code are -

- postfix notation
- syntax trees
- quadruples
- Triplet

Q Postfix Notation \Rightarrow postfix expression contains operators at their right end.

e.g. for expression infix expression

postfix expression is $a b c * d e$

Q. $(a+b)*c$

$a b + c *$ is postfix expression for above infix expression

Q. If x then y else z

by using ? ternary postfix operator

Postfix form for above is $x y ? z$

e.g. If a then if c-d then a+c else a*c else a/b
 postfix for above ab expression is

acd-act+act? abt?

Syntax directed translation to postfix code

Production

$$E \rightarrow E^{(1)} \text{ op } E^{(2)}$$

$$E \rightarrow (E^{(1)})$$

$$E \rightarrow \text{id}$$

Semantic Action

$$E\text{-code} = E^{(1)}\text{-code} || E^{(2)}\text{-code} || \text{op}$$

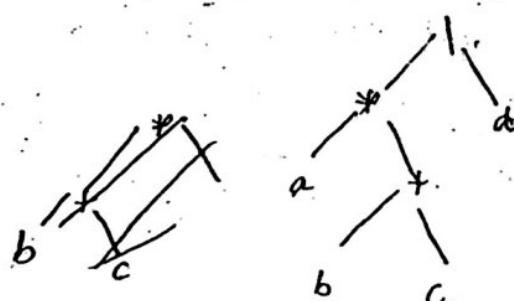
$$E\text{-code} = E^{(1)}\text{-code}$$

$$E\text{-code} = \text{id}$$

22 SYNTAX TREES

Syntax tree is basically a variant of parse tree in which each leaf represents an operand and each interior node represents an operator.

e.g. Syntax tree for expression $a * (b + c) / d$ is



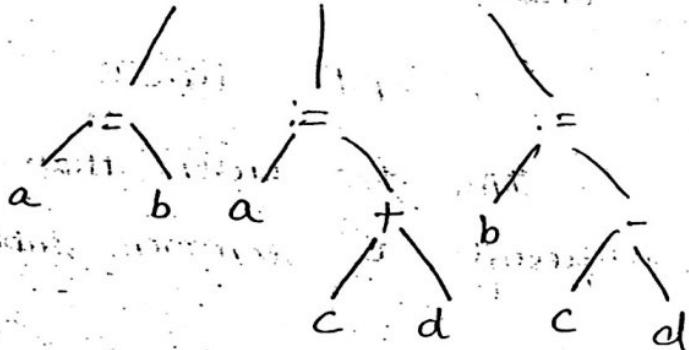
$a * (b + c) / d$

J. if $a=b$ then $a:=c+d$ else $b:=c-d$

syntax tree

(48)

if - then - else



Syntax directed translation scheme to construct
syntax tree

Semantic actions

$E \rightarrow E^{(1)} + T$

$E \rightarrow E^{(1)} - T$

$E \rightarrow T$

$T \rightarrow (E)$

$T \rightarrow id$

$\rightarrow num$

$\rightarrow - E^{(1)}$

$E\text{-Node} := \text{newNode}(+, E^{(1)}.node, T.node)$

$E\text{-Node} := \text{newNode}(-, E^{(1)}.node, T.node)$

$E\text{-Node} := T\text{-Node}$

$T\text{-Node} := E\text{-Node}$

$T\text{-Node} := \text{newLeaf}(id, id\text{-entry})$

$T\text{-Node} := \text{newLeaf}(num, num\text{-val})$

$E\text{-Node} := \text{UNARY}(-, E^{(1)}.node)$

Directed Acyclic graph for expressions

DAG is similar to syntax trees having leaves corresponding to operands and interior nodes corresponding to operators.

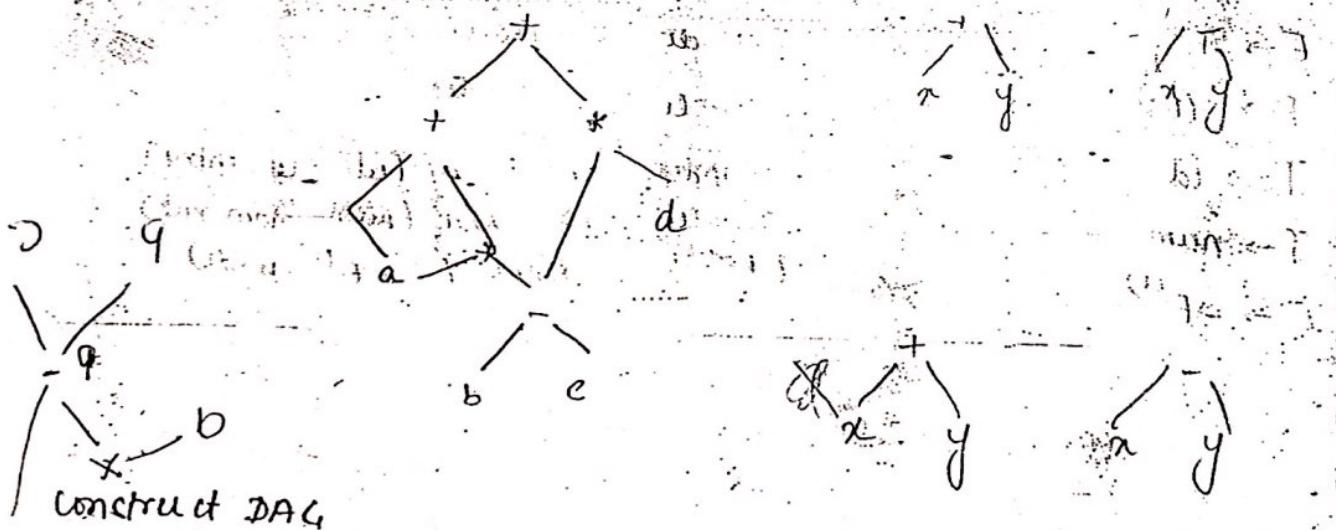
Difference b/w DAG and syntax tree is

that a node N in a DAG has more than one parent if N represents a common subexpression in a syntax tree.

means DAG represents the generation of efficient code to evaluate the expressions.

$$(a + a + (b - c)) * (b - c) * d$$

its DAG



$$((x+y) - (x+y) * (x-y))) + ((x+y) * (x-y))$$

Three address code is a linearized representation of a syntax tree or a DAG in which explicit names corresponding to the interior nodes of the graph are used.

Three address code is a sequence of statements of the form $A = B \text{ op } C$. Each statement generally contains 3 addresses two for operands and one for the result. Each statement does not involve, not more than three addresses. Therefore it is called three address statement, and a sequence of such statements is referred as three address code.

Eg. Three address code for expression

$$a + b * c + d \text{ is}$$

$$T_1 := b * c$$

$$T_2 := a + T_1$$

$$T_3 = T_2 + d$$

Here T_1, T_2 and T_3 are temporary names which are generated by the compiler.

1) Arithmetic or logical statement 15/2

$$A = B \text{ op } C$$

2) Assignment instructions of the form

$$A := \text{op } B$$

where op is a unary operator like unary minus, logical negation, shift operators and conversion operators etc.

e.g. ~~Truncate float T1~~

$$T1 := \text{int to float } T2$$

3) Copy statements of the form

$$A := B$$

where B is assigned to the value of A.

4) Unconditional jump statement

goto L

means the three address statement with label L is next to be executed.

5) Conditional jump statement

if A $\text{rel op } B$ goto L

means if A $\text{rel op } B$ is true statement with label L is executed else the next statement in the normal flow sequence is executed.

2) conditional jump statement

— if a goto L

means if a is non zero then statement with label L is executed and if a is false means go to the next statement in normal flow sequence executed.

Procedure calls and return statement

form parameters param x

to procedure and function call p,n

call p,u

return y

in returning

eg. p(x₁, x₂, ..., x_n)

these address → param x₁
param x₂

param x_n

call p,u

Indexed assignment statement

x = y[i]

y[i] = y

Address and relative assignment

x = *y

*x = y

D Implementation of three address statements

A three address code is an abstract form of intermediate code. These statements can be implemented as record with fields for the operator and operands. Two such representations are -

Quadruples

Triples

Indirect Triples

Quadruples - Quadruple is a record structure using four fields -

OP ARG1 ARG2 RESULT

OP field contains an internal code for the operator and ARG1, ARG2, RESULT are pointers to the symbol table entries.

$A = B \text{ OP } C$ is represented by
B in ARG1
C in ARG2
A in RESULT

$A = OPB$

B in ARG1

nothing in ARG2

A in RESULT

Param. use neither ARG2 nor RESULT

conditional and unconditional jumps put the target label in result field.

$$\underline{\text{Ex}} \quad A := -B * (C + D)$$

Op	ARG1	ARG2	RESULT
(0) unminus	B	-	T ₁
(1) +	C	D	T ₂
(2) *	T ₁	T ₂	T ₃
(3) :=	T ₃		A

three address code is

$$T_1 := -B$$

$$T_1 \rightarrow -B$$

$$T_2 := C + D$$

$$T_2 = C + D$$

$$T_3 := T_1 * T_2$$

$$T_3 = T_1 * T_2$$

$$A := T_3$$

$$A = T_3$$

To avoid entering temporary names into the symbol table, temporary value can be referred by the position of the statement that computes it.

So, in this representation, three address statements can be represented by record structures with only three fields

OP ARG1 ARG2

where ARG1, ARG2 are either pointers to the symbol table or pointers into the structure itself

e.g. $A := -B * (C+D)$

Three address code

$$\begin{aligned}T_1 &:= -B \\ T_2 &:= C+D \\ T_3 &:= T_1 * T_2 \\ A &:= T_3\end{aligned}$$

	OP	ARG1	ARG2
(0)	minus	T ₁	-
(1)	+	C	D
(2)	*	(0)	(1)
(3)	:=	A	(2)

Indirect Triple - consists of a listing of pointers to triples rather than a listing of triples themselves in desired order.

156

e.g. $A = -B * (C + D)$

	statement	OP	ARG1	ARG2
(0)	(9)	(9)	uminus	B
(1)	(10)	(10)	+	C
(2)	(11)	(11)	*	(9) (10)
(3)	(12)	(12)	:=	A (11)

Indirect Triple representation

Comparison of quadtuples, Triples, and Indirect Triples

Quadtuples

- ↳ direct access of the location for temporaries
- ↳ easier for optimization

Triples

Space efficiency

Indirect Triples

- ↳ easier for optimization
- ↳ space efficiency

→ Assignment statements of integer types.

e.g. $a = b + c + d$

then address code is

$$T_1 := b + c$$

$$T_2 := T_1 + d$$

$$a := T_2$$

→ Assignment statements with mixed types

e.g. $x = y + I * J$

where x and y are real and I, J are of integer types.

$$T_1 := I * J$$

$$T_2 := \text{int to float } T_1$$

$$T_3 := y + T_2$$

$$x := T_3$$

e.g. $x := y * z + y * u$

where y, z are integers and x, u are of type float.

$$T_1 := \text{int to float } y$$

$$T_2 := T_1 * z$$

$$T_3 := y * u$$

$$T_4 := \text{int to float } T_3$$

$$T_5 := T_2 + T_4$$

$$x := T_5$$

Boolean expressions

For generating the three address code for Boolean expressions, some branching statements of the following forms are used.

goto L

if A goto L

if A \neq B goto L

eg. 1 Three address code for expression

A or B and C

$T_1 := B \text{ and } C$

$T_2 := A \text{ or } T_1$

eg. 2 $A < B$

similar to if $A < B$ then 1 else 0

three address code is

if $A < B$ goto L_1 ,

$T = 0$

goto L_2

$L_1: T = 1$

$L_2:$

or

(1) if $A < B$ goto (4)

(2) $T = 0$

(3) goto (5)

(4) $T = 1$

(5)

eg. 3

if $A < B$ or C

(1) if $A < B$ goto (4)

(2) $T_1 := 0$

(3) goto (5)

(4) $T_1 \neq 1$

(5) $T_2 := T_1 \text{ or } C$

4 if ($A < B \text{ or } C > D$) $x = y + z$

if $A < B$ goto

4 if $a > b$ then $x = a + b$

if $a > b$ goto L₁

goto L₂

L₁:

$T_1 := a + b$

$X := T_1$

L₂:

5 if $a > b$ then $x = a + b$
else $x = a - b$

if $a > b$ goto L₁

goto L₂

4:

$T_1 := a + b$

$X := T_1$

L₂: goto L₃

$T_2 := a - b$

$X := T_2$

L₃:

Three address code for Boolean Expressions

eg.1 A or B and C

TAC) Three address code is

$$T_1 := B \text{ and } C$$

$$T_2 := A \text{ or } T_1$$

eg.2 A < B

it is similar to if $A < B$ then 1
else 0

TAC is

if $A < B$ goto L1

$T = 0$

goto L2

L1: $T = 1$

L2:

eg.3 if $A < B$ or C

TAC is

if $A < B$ goto L1

$T_1 = 0$

goto L2

L1: $T_1 = 1$

L2: $T_2 = T_1 \text{ or } C$

L3:

Q. Three address code for ~~array~~ ~~looping~~ ~~array~~ ~~reference~~ ~~constant~~ ~~reference~~

Relative location of $a[i]$ is

$$\text{base} + i \times w$$

Relative location of $a[i][j]$ is

$$\text{base} + [i \times n_2 + j] \times w$$

eg. 1

three address code for $a[i][j]$ where a is a 2×3 array, where size of each element is 4 bytes
location of $a[i][j]$ can be represented by
 $\text{base} + [i \times n_2 + j] \times w$

$$= 0 + [2 \times 3 + j] \times 4 = 12i + 4j$$

$$T_1 = 12 * i$$

$$T_2 = 4 * j$$

$$T_3 = T_1 + T_2$$

$$x := a[T_3]$$

- ~~Answer~~

Third
will address code for looping constructs.
while statements

while $a \geq b$ do $a = a + b$

L1: If $a \geq b$ goto L2

goto L3

L2: $T_1 := a + b$

$a := T_1$

goto L1

L3

for statements

$$\begin{array}{l} a=3 \\ b=4 \end{array}$$

for ($i=0$; $i < n$; $i++$)

$a = b + 1$

$a = a * a$

$c = a$

TAC is $a=3$

$$b=4$$

go

L1: if $i < n$ goto L2

goto L3

L2: $T_1 := b + 1$

$a = T_1$

$T_2 = a * a$

$a = T_2$

goto L3

L3: $T_3 = a$ $T_3 = i + 1$
 $c = T_3$ $i = T_3$

goto L1

L4: $T_4 := a$ $c = T_4$

e.g. 2

$$a = 0$$

$$b = 2$$

for ($i=1$; $i < n$; $i++$)

{ $a = a + b \neq x[i] + y[i]$

}

TAC is

$$a = 0$$

$$b = 2$$

$$i = 1$$

L₁: $\text{if } i < n \text{ goto L}_2$

 goto L₄

L₂: $T_1 := a + b$

$T_2 := i * 4$

$T_3 := x[T_2]$

$T_4 := y[T_2]$

$T_5 := T_3 + T_4$

$T_6 := T_1 + T_5$

$a := T_6$

goto L₃

L₃:

$T_7 := i + 1$

$i = T_7$

goto L₁

L₄:

Q) Generate Quaduple Representation of

$$a = b * -c + b * -c$$

Three address code by syntax tree representation,

$$T_1 := -c$$

$$T_2 := b * T_1$$

$$T_3 := -c$$

$$T_4 := b * T_3$$

$$T_5 := T_2 + T_4$$

$$a = T_5$$

	OP	ARG1	ARG2	RESULT
0	uminus	c	-	T ₁
1	*	b	T ₁	T ₂
2	uminus	c	-	T ₃
3	*	b	T ₃	T ₄
4	+	T ₂	T ₄	T ₅
5	:=	T ₅	-	a

Three address code by DAG representation,

$$T_1 := -c$$

$$T_2 := b * T_1$$

$$T_3 := T_2 + T_2$$

$$a := T_3$$

	OP	ARG1	ARG2	RESULT
(1)	uminus	c	-	T ₁
(2)	*	1.b	T ₁	T ₂
(3)	+	T ₂	T ₂	T ₃
(4)	:=	T ₃	-	a

$$x[i] = y$$

OP	ARG1	ARG2
[]	x	i
:=	(0)	y

$$x = y[i]$$

	OP	ARG1	ARG2
(0)	[]=	y	i
(1)	:=	x	(0)

Quadruple Representation of

$a = b * c + d * e$

aminotes

	OP	ARG1	ARG2	
0	uminus	c	-	
1	*	b	(0)	
2	uminus	c	-	
3	*	b	(2)	
4	+	(1)	(3)	
5	:=	a	(3)	

	OP	ARG1	ARG2	
(0)	uminus	c	-	
(1)	*	b	(0)	
(2)	+	(1)	(1)	
(3)	:=	a	(2)	

Using syntax

Tree

Using DAT

Generate quadruple representation of

$x = a * b * c + c * d * e$

use address code is

$$T_1 := b * c$$

$$T_2 := a + T_1$$

$$T_3 := c * d$$

$$T_4 := T_3 * e$$

$$T_5 := T_2 + T_4$$

$$x = T_5$$

	OP	ARG1	ARG2	RESULT
(0)	*	b	c	T ₁
(1)	+	a	T ₁	T ₂
(2)	*	c	d	T ₃
(3)	*	T ₃	e	T ₄
(4)	+	T ₂	T ₄	T ₅
(5)	:=	T ₅	-	x

Q. Generate indirect triple representation

$$x := a + b * c + c * d * e \quad \text{and} \quad y := y + a + b * c$$

→ Triple representation

	OP	ARG1	ARG2
(10)	*	b	c
(11)	+	a	(10)
(12)	*	c	d
(13)	*	(12)	e
(14)	+	(11)	(13)
(15)	:=	x	(14)
16	+	y	(11)

Indirect Triple Representation:

Statement table for expression

$$x := a + b * c + c * d * e$$

Statement table for expression

$$y + a + b * c$$

(10)	0	(10)
(11)	1	(11)
(12)	2	(16)
(13)		
(14)		
(15)		

generate three address code for

$$x := y + I * J$$

where

x and y are real and
 I & J are integers

$$T_1 := I * J$$

$$T_2 := \text{int to float } T_1$$

$$T_3 := y + T_2$$

$$x := T_3$$

Q Generate three address code for

$$x := y * z + y * w$$

where

y, z are integers

x, z, w are float

$$T_1 := \text{int to float } y$$

$$T_2 := T_1 * z$$

$$T_3 := y * w$$

$$T_4 := \text{int to float } T_3$$

$$T_5 := T_2 + T_4$$

$$x := T_5$$

Q

TACif $A < B$ goto L₁

T=0

goto L₂L₁: T=1L₂: T or C.L₃:Q. If $A < B$ then $x = a+b$ if $A < B$ goto L₁goto L₂L₁: T₁ := a+bx := T₁L₂:(3) if $a < b$ then $x = a+b$ else $x = a-b$ if $a < b$ goto L₁goto L₂L₁: T₁ := a+bx := T₁goto L₃L₂: T₂ := a-bx := T₂L₃:

Q.4 $a \leq b \& c < d$

if $c < d$ goto L₁

T₁ = 0

goto L₂

L₁: T₁ = 1

L₂: T₂ := b & T₁

T₃ := a || T₂

L₃:

Q.5 $a < b \& c < d$

if $a < b$ goto L₁

goto L₂

L₁: if $c < d$ goto L₃

L₂: T = 0 goto L₄

L₃: T = 1

L₄:

Q.6 if ($x < 100 \text{ || } x > 200 \& x \neq y$) x = 0

if $x < 100$ goto L₁

~~goto L₂~~

L₂: if $x > 200$ goto L₃

goto L₂

L₃: if $x \neq y$ goto L₁

goto L₂

L₁: x = 0

L₂:

Q7. If $a < b$ then $x = a+b$
 else if $a = b$ then $x = a$
 else $x = b$

TAC. i.e.

If $a < b$ goto L1

goto L2

L1: $T_1 := a+b$

$x := T_1$

goto L3

L2: If $a = b$ goto L4

goto L5

L4: $t_1 := a$

$x := t_1$

goto L3

$t_1 := b$

$x := t_1$

L3:

Q8. If $(a < b \text{ || } c < d)$ $x = y+z$

TAC. i.e.

If $a < b$ goto L1

If $c < d$ goto L1
 goto L2

L1: $T_1 := y+z$

$x := T_1$

L2:

Generate three address code for

$$x = c + a[i][j]$$

where a is a 2×3 array and size of each element is 4 bytes.

e.g. $a[i][j]$ can be represented by

$$\begin{aligned} & \text{base} + (i \times n_2 + j) \times w \\ &= 0 + (i \times 3 + j) \times 4 \\ &= 12i + 4j \end{aligned}$$

$$\text{so TAE is } T_1 := 12i$$

$$T_2 := 4j$$

$$T_3 := T_1 + T_2$$

$$T_4 := a[T_3]$$

$$T_5 := c + T_4$$

$$x := T_5$$

Generate three address code for

$$x[i, j] = y$$

where x is 5×10 array and size of each element is 2 bytes.

$x[i, j]$ can be represented by

$$\text{base} + (i \times n_2 + j) \times w$$

$$\text{e.g. } 0 + (i \times 10 + j) \times 2 = 20i + 2j$$

$$T_1 := 20i$$

$$T_2 := 2j$$

$$T_3 := T_1 + T_2$$

Q1. while $a \geq b$ do $a = a + b$

L1: if $a \geq b$ goto L2
 goto L3

L2: $T_1 = a + b$
 $a = T_1$
 goto L1

L3:

~~initial~~
 $i = 1;$ $a = 3$
 $b = 4$

for ($i = 0; i < n; i++$)
{ $a = b + i;$
 $a = a * a;$
}
 $c = a$

TAC U

$a = 3$
 $b = 4$
 $i = 0$

L4: if $i < n$ goto L2
 goto L4

L2: $T_1 := b + i$
 $a = T_1$
 $T_2 = a * a$
 $a = T_2$
 goto L3

L3: $T_3 := i + 1$
 $i = T_3$
 goto L1

L4: $T_4 = a$
 $c = T_4$

Q $a=0$
 $b=2$

$\text{for } (i=1; i < n; i++)$
 $\{ a = a + b + x[i] + y[i]$
 $\}$

TAC is $a=0$
 $b=2$
 $i=1$

L1: $\text{if } i < n \text{ goto L2}$
 goto L4

L2: $T_1 := a + b$
 $T_2 := i + 1$
 $T_3 := x[T_2]$
 $T_4 := y[T_2]$
 $T_5 := T_3 + T_4$
 $T_6 := T_1 + T_5$
 $a = T_6$
 goto L3

L3: $T_7 := i + 1$
 $i = T_7$
 goto L1

L4:

do $i+1$
 $\text{while } (i < n);$

TAC is L1: $T_1 := i + 1$
 $i = T_1$
 $\text{if } i < n \text{ goto L1}$

L1:

do L+1
while ($a[i^*] < v$)

L1: $T_1 := i^*$
 $i^* = T_1$

if T_2

L2: $i^* = i^* + 4$

$T_3 := a[T_2]$

if $T_3 < v$ goto L1

L2:

$n = f(a[i^*])$

$T_1 := i^* * 4$

$T_2 := a[T_1]$

param T2

$t_3 = \text{call } f, 1$

$n = t_3$

int i;

i=1;

while a<10 do

if. x>y then a= x+y

else a=x-y

TAC
t₁=1i=t₁L₁: if {a<10 goto L₂
goto L₄L₂: if x>y goto L₃T₂ := x-y
a = T₂; goto L₁L₃: T₁ := x+y
a = T₁; goto L₁L₄:

int c; int p; int i=1;

c=i

p=i

for (i=2; i=n; i++)

c=c+p;

p = c-p;

y

TACT₁=1i=T₁T₂=iC=T₂p=T₂T₃=2i=T₃T₄=ni=T₄T₅=iL₁: if T₅ >= n goto L₂
goto L₃T₆=C+pC=T₆T₇=C-pD=T₇goto L₃L₃: T₈ = L₁T₁ = T₈goto L₁

Q. ~~Q~~ while ($a < c \& b < d$) do

 if $a=1$ then $c=c+1$

 else

 while $a \leq d$ do

$a=a+3$

 L1: if $a < c$ goto L2
 goto L5

 L2: if $b < d$ goto L3
 goto L5

 L3: if $a \leq d$ goto L4

 L4: $T_1 = a+3$
 $a = T_1$ goto L3

L5:

Q. P(A+B, C)

P is a procedure

do while statement

do i++

while ($i < n$);

L1: $T_1 := i + 1$

$i = T_1$

if $i < n$ goto L1
goto L2

L2:

β_2

do i++

while ($a[i] < v$)

L3: $T_1 := i + 1$

$i = T_1$

$T_2 := i * 4$

$T_3 := a[T_2]$

if $T_3 < v$ goto L1

goto L2

L2:

TAC for Programs

$n = f(a[i]);$

TAC is

$t_1 := L * 4$

$t_2 := a[t_1]$

param t_2

$t_3 := \text{call } f, 1$

$n = t_3$

TOP

$f(x_1, x_2, \dots, x_n)$

TAC is _____

param x_1

param x_2

⋮

param x_n

$t_3 := \text{call } f, n$

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮

⋮</p

sum = 0;

for ($i=1$; $i \leq 20$; $i++$)

 sum = sum + a[i] + b[i];

 true

$T_1 = 0$

 sum = T_1

$T_2 = 1$

$i = T_2$

L1: if $i \leq 20$ goto L2
 goto L4

L2: $T_3 = 4 * i$

$T_4 = a[T_3]$

$T_5 = b[T_3]$

$T_6 = T_4 + T_5$

$T_7 = T_6 + 8000$

$T_8 = T_6 + T_7$

 sum = T_8

 goto L3

L3: $T_9 = i + 1$

$i = T_9$

 goto L1

L4:

2

```

int b;
int a;
b=3
a=12;
a- (b+2) - (a*3)/6

```

TAC

$$\begin{aligned}
 t_1 &= 3 \\
 b &= t_1 \\
 t_2 &= 12 \\
 a &= t_2 \\
 t_3 &= b + 2 \\
 t_4 &= a * 3 \\
 t_5 &= t_4 / 6 \\
 t_6 &= t_4 / t_5 \\
 t_7 &= t_3 - t_6 \\
 a &= t_7
 \end{aligned}$$

void main

```

{
    int a;
    a = 23;
    if (a == 23)
        a = 10;
    else
        a = 19;
}

```

TAC →

Begin func

$$\begin{aligned}
 t_0 &= 23; \\
 a &= t_0; \\
 \text{if } (a == 23) \text{ goto L1} \\
 t_1 &= 19; \\
 a &= t_1; \\
 \text{L1: } t_1 &= 10; \\
 a &= t_1;
 \end{aligned}$$

L2:

switch A+B

begin

case 2: $x := 4$

case 5:

switch X

begin

case 0: $A := B+1$

case 1: $A := B+3$

default: $A := 2$

end

case 9: $x = 4-1$

default: $x = 4+1$

end

switch A+B

begin

case v_1 : $x = a$

case v_2 : $x = b$

default:

$x = c$

end

$T_1 := A+B$

if $T_1 = v_1$, goto L₁

if $T_1 = v_2$, goto L₂

if

$T_1 = A+B$

if $T_1 = 2$, goto L₁

goto L₅

if $T_1 = 5$, goto L₂

goto L₅

if $T_1 = 9$, goto L₉

L₁: $T_2 = 4$

$X = T_2$

goto NEXT

L₂: $T_3 = X$ gotoNEXT

if $T_3 = 0$, goto L₁

goto L₆

if $T_3 = 1$, goto L₇

goto L₆

L₄: $A = B+1$ gotoNEXT

L₅: $A = B+3$ gotoNEXT

L₆: $A = 2$ gotoNEXT

L₅: $x = 4+1$
goto NEXT

L₉: $x = 4-1$

NEXT

$T_1 = A+B$

if $T_1 = v_1$, goto L₁

goto L₅

if $T_1 = v_2$, goto L₂

goto L₅

L₁: $x = a$

L₂: $x = b$

L₅: $x = c$

Code optimization

Code optimization refers to the techniques used by the compiler to improve the execution efficiency of the generated object code.

Code optimization includes various transformations on intermediate code, but these transformations must preserve the semantics of the program.

Criteria to apply optimizing transformations -

- 1) The optimization should capture most of the potential improvements without an unreasonable amount of effort.
- 2) Optimization should be such that the meaning of the source program is preserved.
- 3) The optimization should on average, reduce the time and space expended by the object code.

Code optimization can be done in three ways -

1) Local optimization

which performed in a straight line (basic block of code means a block with ~~except~~ no jumps except at the beginning and no jumps except at the end).

- 2) loop optimization
- 3) data flow analysis

↳ transmission of useful information from all parts of the program to the places where the information can be of use.

BASIC BLOCKS AND FLOW GRAPHS

A graph representation of three address statements called a flow graph. Nodes in flow graph represent computations and the edges represent the flow of control. Flow graphs can be used to collect information about the intermediate program and to find inner loops where a program is expected to spend most of its time.

BASIC BLOCK

Basic block is a sequence of consecutive statements in which

a) flow of control can only enter at the beginning. No jumps are allowed into the middle of the block.

b) flow of control can leave the block at the end without halt or possibility of branching except at the end

Algorithm for the construction of Basic Blocks

Step 1 - First determine the set of leaders, the first statements of basic blocks.

Rules for finding leaders are:-

- First three address statement in intermediate code is a leader.
- Any statements that immediately follows a conditional or unconditional jump (goto) is a leader.
- Any statement that is the target of a conditional or unconditional jump is a leader.

Step 2 - For each leader, its basic block consists of all instructions upto the next start of a block.

Example ①

$$t_1 = a * a$$

$$t_2 = a * b$$

$$t_3 = 2 * t_2$$

$$t_4 = t_1 + t_3$$

$$t_5 = b * b$$

$$t_6 = t_4 + t_5$$

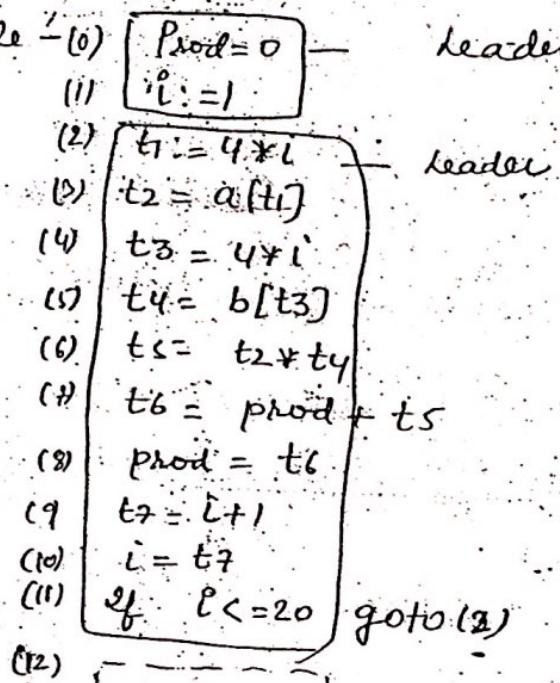
These sequence of statements forms a basic block.

```

begin
prod=0
l:=1
do begin
prod = prod + a[i]*b[i]
i := i+1
end
while i <= 20
end

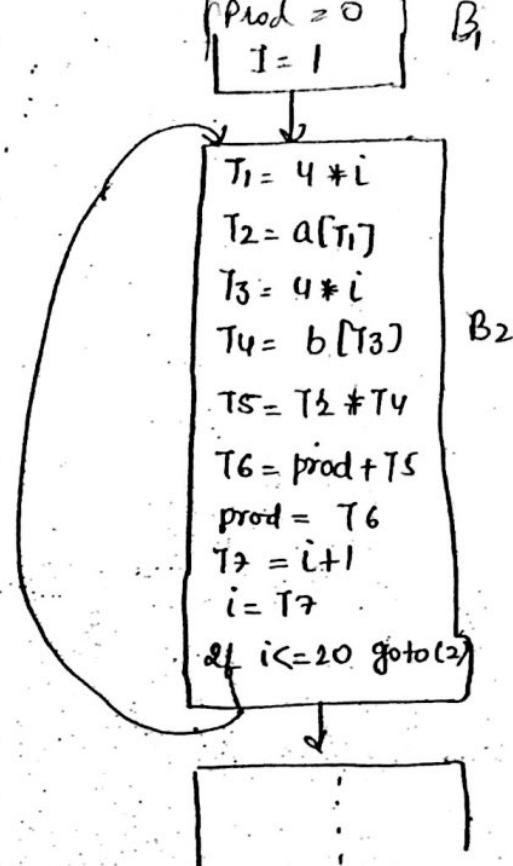
```

Three address code - (0)



Sequence of Basic Block.

Flow graph :- The relationship among basic block is represented by a directed graph called a flow graph. There is a directed edge from block B_1 to B_2 if B_2 immediately follows B_1 during execution. we can say B_1 is a predecessor of B_2 & B_2 is a successor of B_1 .



B2

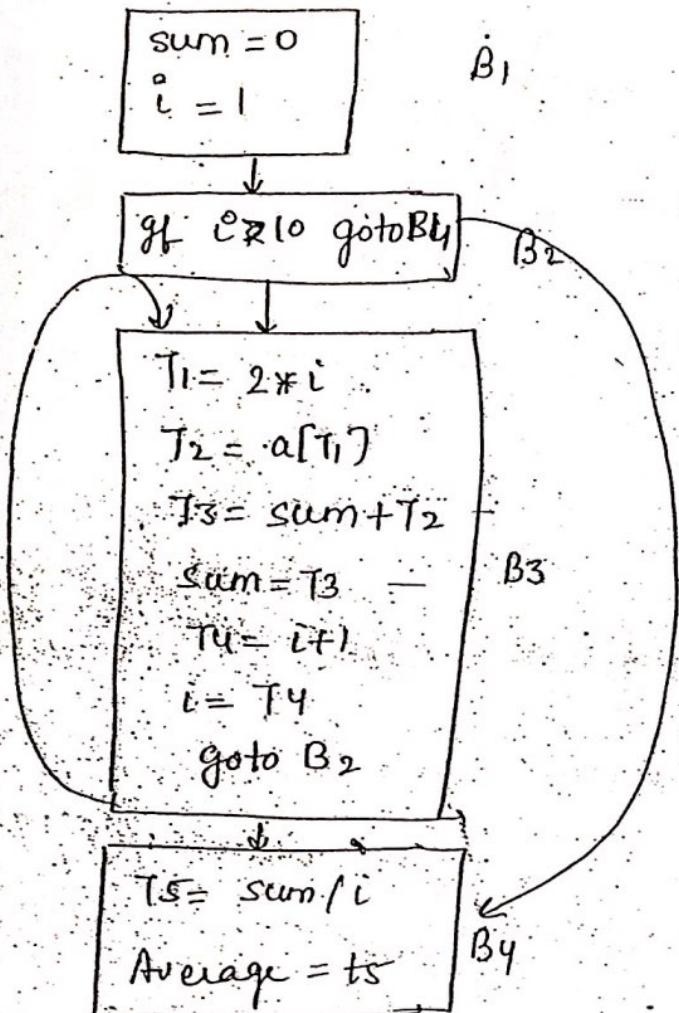
Ex-3

```

sum=0
i=1
while i<=10 do
{
    sum= sum+ a[2*i]
    i= i+1
}
average = sum/i
  
```

Three address statement

- (1) sum=0
- (2) i=k
- (3) if i>=10 goto 11
- (4) T1 = 2*i
- (5) T2 = a[T1]
- (6) T3 = sum+T2
- (7) sum = T3
- (8) t4 = i+1
- (9) i = T4
- (10) goto(3)
- (11) T5 = sum/i
- (12) average = T5

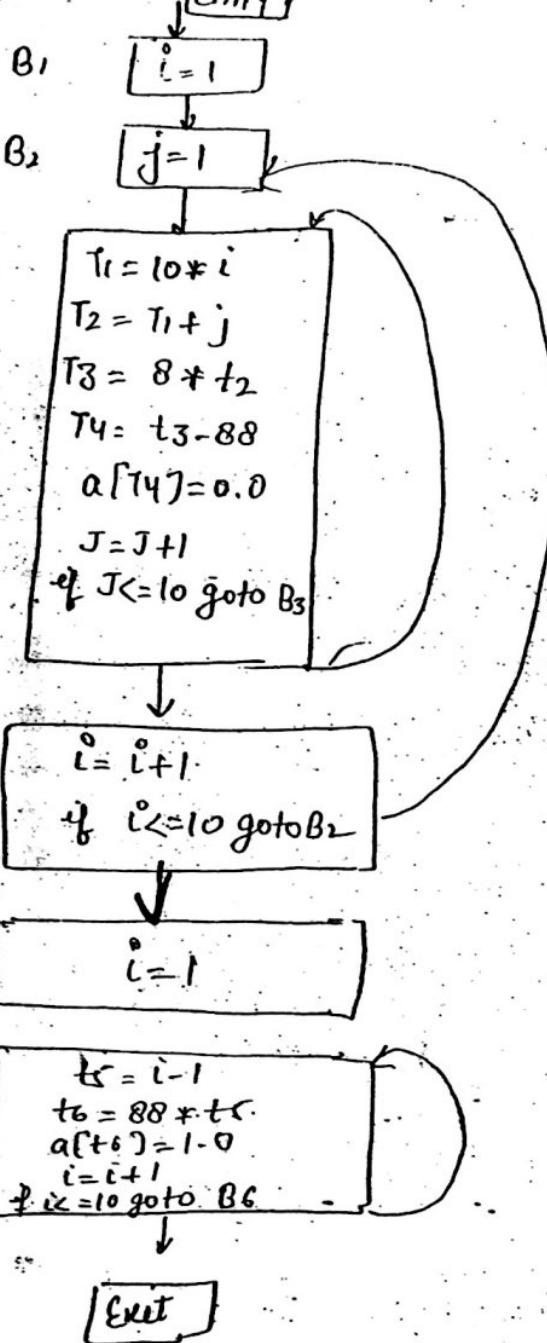


Ex-4 Consider the following three address code and find basic blocks.

- (1) $i = 1$
- (2) $j = 1$
- (3) $t1 = 10 * i$
- (4) $t2 = t1 * j$
- (5) $t3 = 8 * t2$
- (6) $t4 = t3 - 88$
- (7) $a[t4] = 0.0$
- (8) $j = j + 1$
- (9) $\text{if } j \leq 10 \text{ goto (3)}$
- (10) $i = i + 1$
- (11) $\text{if } i \leq 10 \text{ goto (2)}$
- (12) $i = 1$
- (13) $t5 = i - 1$

- (14) $t6 = 88 * t5$
- (15) $a[t6] = 1.0$
- (16) $i = i + 1$
- (17) $\text{if } i \leq 10 \text{ goto (3)}$

blocks &
flow chart



Local Optimization

A basic block computes a set of expressions. Two basic blocks are said to be equivalent if they compute the same set of expressions.

Optimization can be achieved by applying a no. of transformations within each basic block.

Two classes of local transformations can be applied to basic blocks -

1) Structure Preserving transformations

Common subexpression elimination

Dead code elimination

~~Renaming, temporary variables~~ Copy Propagation

~~Renaming of statements~~ Constant folding

2) Algebraic transformations

Data Representation of Basic Blocks

- DAG (directed acyclic graph) is a useful data structure for analysing basic blocks. ~~for~~ Transforming basic blocks into DAG helps in various techniques of local optimization.

Application of DAG's representation

- Using DAG we can perform following actions:
 - 1) we can eliminate local common subexpressions, i.e. instructions that compute a value that has already been computed.
 - 2) we can eliminate dead code, statements that compute a value that is never used.
 - 3) we can reorder statements that do not depend on one another. this reordering may reduce the time a temporary value needs to be preserved in a register.
 - 4) we can apply algebraic laws to simplify the computation.

we have three types of statements

- (i) $A = B \text{ op } C$
- (ii) $A = \text{op } B$
- (iii) $A \leftarrow B$

If $C \leftarrow B$ will be treated as case (i) with A undefined.

1) Construct node in the DAG for each initial values of the variables appearing in the basic block.

2) For statement (i) check if we have a node

labeled op with whose left child is node B and right child is node C. If not then create such node. Let n is the node.

3) For statement (ii) check if we have a node labeled op with lone child node B. If not then create such node. Let n is the node.

4) In case (iii) make n as node(B).

5) Attach A to node N.

6) An assignment from an array like

$x = a[i]$ is represented by creating

a node with operator $=[]$ with two children a & variable x becomes the label of the new node.

7) For statement like $a[j] = y$ replaced by

two node with operator $[] =$ and three children

representing a, j & y. It has no variable

labeling this node. Creation of this node kills

currently constructed nodes whose value

depends on a. This node will not receive any label.

► local optimization

► structure preserving optimization

► A) common sub expression elimination

► e.g. consider the following basic block

1) $a = b + c$

2) $b = a - d$

3) $c = b + c$

4) $d = a - d$

► here 2nd and 4th statement computes the same expression. So transformed basic block is

$a = b + c$

$b = a - d$

$c = b + c$

$d = b$

► Also 1st and 3rd statement have same expression but the value of b in statement 3rd is different from statement 1st. so 1st and 3rd statement do not compute the same expression and hence cannot be eliminated.

Ex-2

$$T_6 = 4 * i,$$

$$x = a[T_6]$$

$$T_7 = 4 * i$$

$$T_8 = 4 * j \quad \text{after elimination}$$

$$T_9 = a[T_8] \quad \text{we get}$$

$$a[T_7] = 19$$

$$T_{10} = 4 * j$$

$$a[T_{10}] = ?$$

$$T_6 = 4 * i$$

$$x = a[T_6]$$

$$T_8 = 4 * j$$

$$T_9 = a[T_8]$$

$$a[T_6] = T_9$$

$$a[T_8] = x$$

Common subexpression elimination using DAG

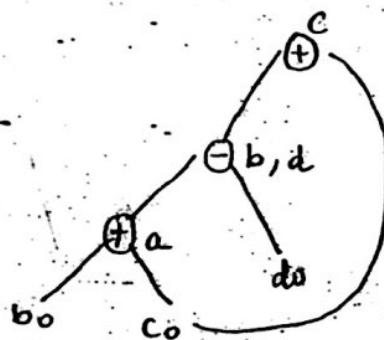
ex-

$$a = b + c$$

$$b = a - d$$

$$e = b + c$$

$$d = a - d$$

DAG

because the DAG has only three non leaf nodes, the transformed basic block will contain only statements

If i is not alive on exit then no need to compute b . we get

$$a = b + c$$

$$d = a - d$$

$$c = d + c$$

$$a = b + c$$

$$b = a - d$$

$$c = b + c$$

$$d = b$$

013

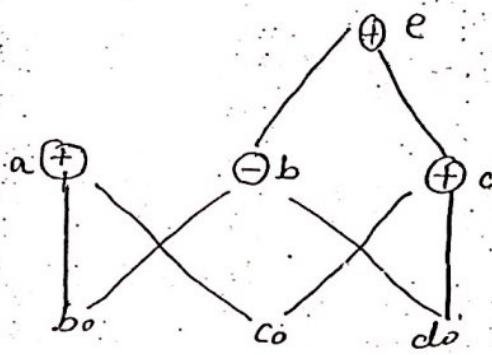
Ex-3 Construct DAG for following basic block

$$a = b + c$$

$$b = b - d$$

$$c = c + d$$

$$e = -b + c$$

DAG

In statement 1 and 4 the value of both variables b and c changes so they does not exhibit common subexpressions.

Ex-4 construct DAG for following basic block and minimize the code.

~~$a = x + y$~~

~~$x = x$~~

~~$b = x + y$~~

~~$x = b$~~

~~$c = x + y$~~

A variable is said to be dead if its value is not used anywhere in the program. Subsequently in assignment statement made to a dead variable will be safely removed.

Ex $x = t_3$

$a[t_2] = t_5$

$a[t_4] = t_2$

Print $a[t_4]$

here $x=t_3$ is a dead code.

(c) Copy Propagation

Let we have following expression

$$\begin{array}{l} a = d + e \\ b = d + e \\ \downarrow \quad \downarrow \\ c = d + e \end{array}$$

here common sub expression in $c = d + e$ is eliminated by using new variable t to hold the value of $d + e$.

$$\begin{array}{l} t = d + e \\ a = t \\ \downarrow \quad \downarrow \\ b = t \\ c = t \end{array}$$

► constant folding

- Deducing at compile time that the value of an expression is a constant and using the constant instead of known as constant folding

eg. $a = 36 * 2$

point a

Here $a = 72$ will be computed at compile-time.

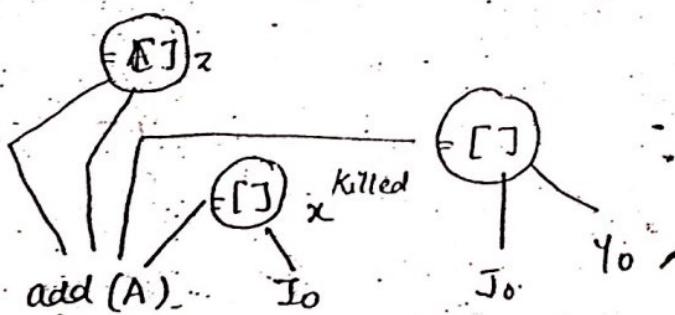
DAG construction

Q construct DAG for

$$x = A[I]$$

$$A[J] = 4$$

$$z = A[I]$$



acc to the last rule in DAG construction algo

$A[J] = 4$ will have 3 child and has no label

this node will kill the previous node having add A.

so for statement

$$z = A[I]$$

new node will be constructed.

$$S_1 = 4 * I$$

$$S_2 = \text{add}(A) - 4$$

$$S_3 = S_2 [S_1]$$

$$S_4 = 4 * I$$

$$S_5 = \text{add}(B) - 4$$

$$S_6 = S_5 [S_4]$$

$$S_7 = S_3 * S_6$$

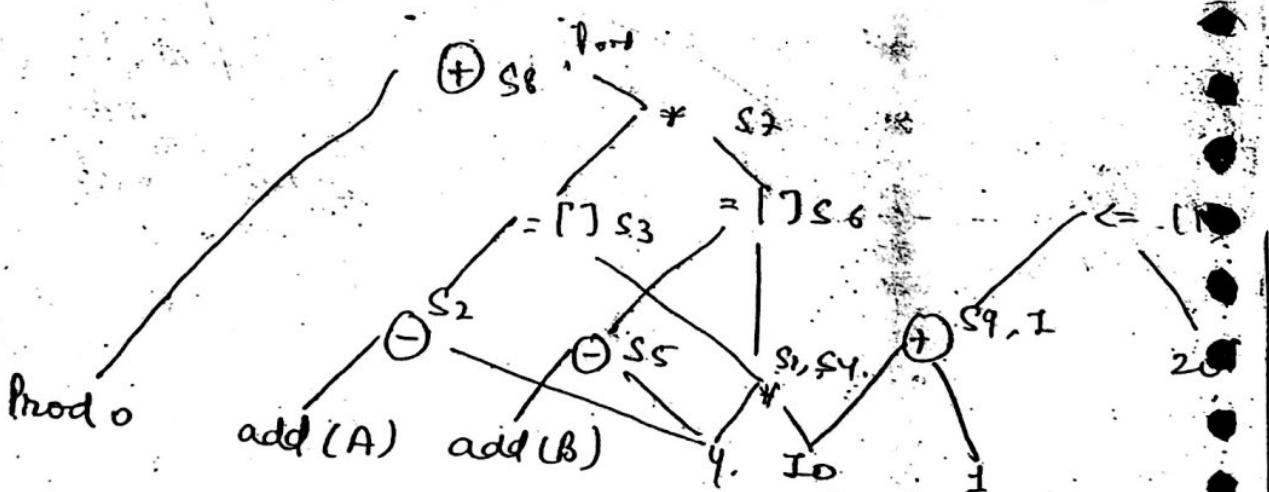
$$S_8 = \text{Prod} + S_7$$

$$\text{Prod} = S_8$$

$$S_9 = I + 1$$

$$I = S_9$$

~~if~~ $I \leq 20$ goto(1)



after optimization we get

$$a = b * 5$$

$\text{for } (i=1; i \leq 10; i++)$

$$\left\{ \begin{array}{l} x = a + b + i \\ \end{array} \right.$$

}

eg. 2 $\text{while } (i \leq \text{limit}-2)$

after code motion we get

$$t = \text{limit} - 2$$

$\text{while } (i \leq t)$

Reduction in strength

Replacement of an expensive operation by a cheaper one is called as reduction in strength.

eg. $x = 1$

$\text{for } (i=1; i \leq 10; i++)$

$$\left\{ \begin{array}{l} x = i * 2 \\ \end{array} \right.$$

}

can be optimized as

$$x = 0$$

$\text{for } (i=1; i \leq 10; i++)$

$$\left\{ \begin{array}{l} x = x + 2 \\ \end{array} \right.$$

}

Induction variable elimination

The induction variable in the loop are those variables whose assignments with in the loop are of the form $I = I + C$ where C is a constant.

Elimination of induction variable sometimes has performance strength in reduction.

e.g. we have following

code in loop

$$L: T_1 = 4 * i$$

$$T_3 = T_2[T_1]$$

$$I = I + 1$$

if $I \leq 20$ goto L

here we can see that I & T_1 are in lock step.

so when value of I increases from 1...20

then T_1 will increase from 4..8..16..20

so I & T_1 are induction variables

by induction variable elimination

we get

$$T_1 = 0$$

$$L: T_1 = T_1 + 4$$

$$T_3 = T_2[T_1]$$

if $T_1 \leq 76$ goto L

② Algebraic transformations

many algebraic transformations can be used
to change the set of expressions computed by
basic stock into simplified expressions or
to replace expensive operations into cheaper operations

eg. $x = x + 0$

and

$$x = x * 1$$

can be eliminated

eg. $x = y * * 2$

can be replaced by $x = y * y$

eg. $2 * x$

$$\hookrightarrow x + x$$

eg. $x / 2$

$$\hookrightarrow x * 0.5$$

eg. x^2

$$\hookrightarrow x * x$$

② Loop Optimization

Generally programs spend most of their time in the execution of inner loops. The running time of program may be improved by decreasing the no. of instructions in an inner loop, or by increasing the amount of code outside the loop.

Three techniques for loop optimization

- ↳ Code motion / loop invariant
- ↳ Induction variable elimination
- ↳ induction in strength.

Code motion / loop invariant

Code motion is a technique of moving loop invariant computations outside the loop.

A loop invariant computation is one that computes the same value every time a loop is executed.

e.g. $\text{for } (i=1; i \leq 10; i++)$

$$\left\{ \begin{array}{l} a = b * 5 \\ x = a + b + i; \end{array} \right.$$

```

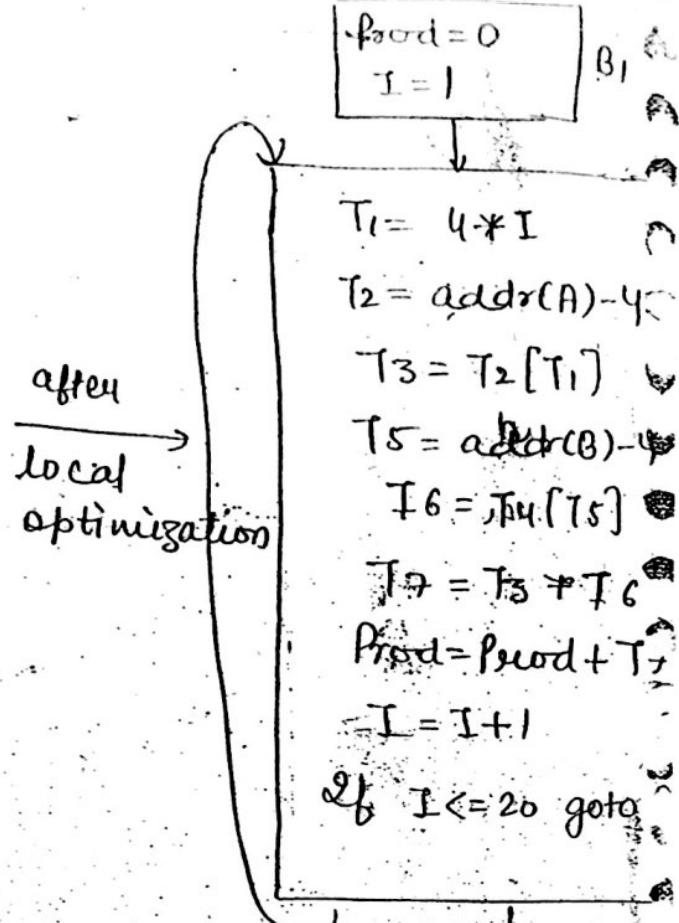
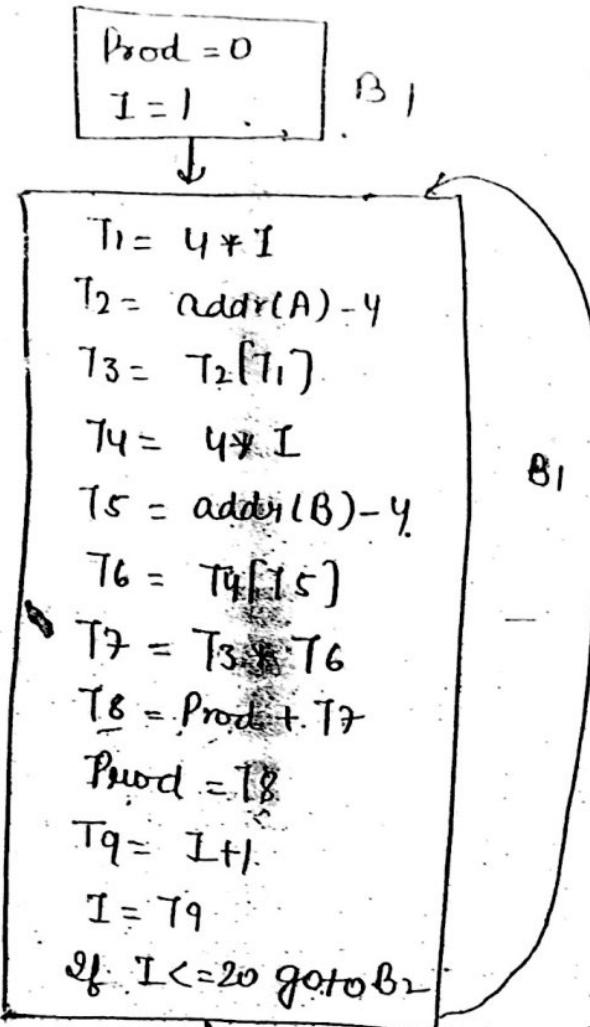
begin
    Prod = 0
    I = 1
do
    begin
        Prod = Prod + A[I] * B[I]
        I = I + 1
    end
    while I <= 20
end

```

- 1) Generate three address code
- 2) Create basic blocks & flow graph
- 3) Perform local & loop optimization

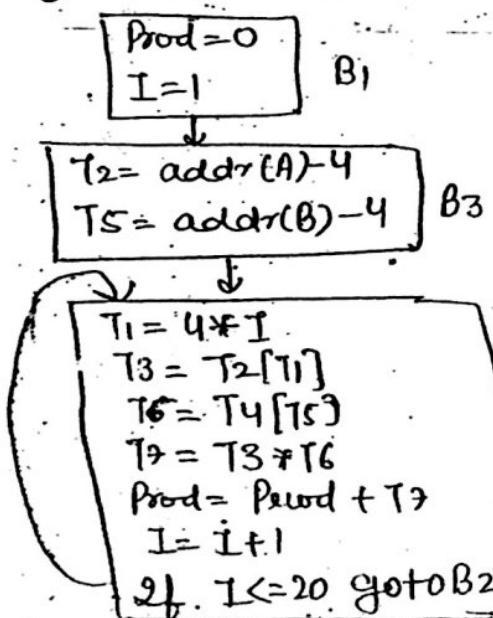
3AC Three address code

- 1) Prod = 0
- 2) I = 1
- 3) T₁ = 4 * I
- 4) T₂ = Adder(A) - 4
- 5) T₃ = T₂[T₁]
- 6) T₄ = 4 * I
- 7) T₅ = Adder(B) - 4
- 8) T₆ = T₅[T₄]
- 9) T₇ = T₃ * T₆
- 10) T₈ = Prod + T₇
- 11) Prod = T₈
- 12) T₉ = I + 1
- 13) I = T₉
- 14) If I <= 20 goto (3)

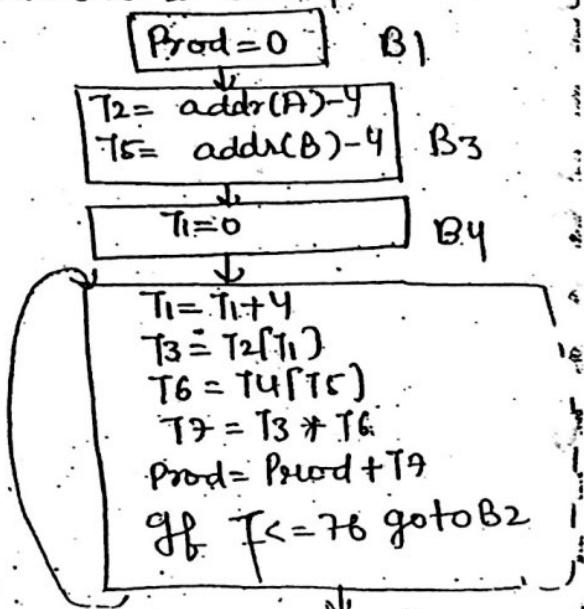


Loop Optimization

After Code motion:



After elimination, induction variable & strength reduction:



```

        j=n; v=a[n]
        while(i)
            do i = i+1 ; while(a[i] < v)
            do j = j-1 ; while(a[j] > v)
            if (i >= j) break;
            x = a[i];
            a[i] = a[j]
            a[j] = x
    
```

y.

```

x = a[i];
a[i] = a[n];
a[n] = x
    
```

Three address code

1)	$T_1 = m - 1$
2)	$i = T_1$
3)	$T_2 = n$
4)	$j = T_2$
5)	$T_3 = 4 * n$
6)	$v = a[T_3]$

B₁

7)	$T_4 = i + 1$
8)	$i = T_4$
9)	$T_5 = 4 * i$
10)	$T_6 = a[T_5]$
11)	$\text{if } T_6 < v \text{ goto(7)}$

B₂

12)	$T_7 = j - 1$
13)	$j = T_7$
14)	$T_8 = 4 * j$
15)	$T_9 = a[T_8]$
16)	$\text{if } T_9 > v \text{ goto(12)}$
17)	$\text{if } i >= j \text{ goto(28)}$

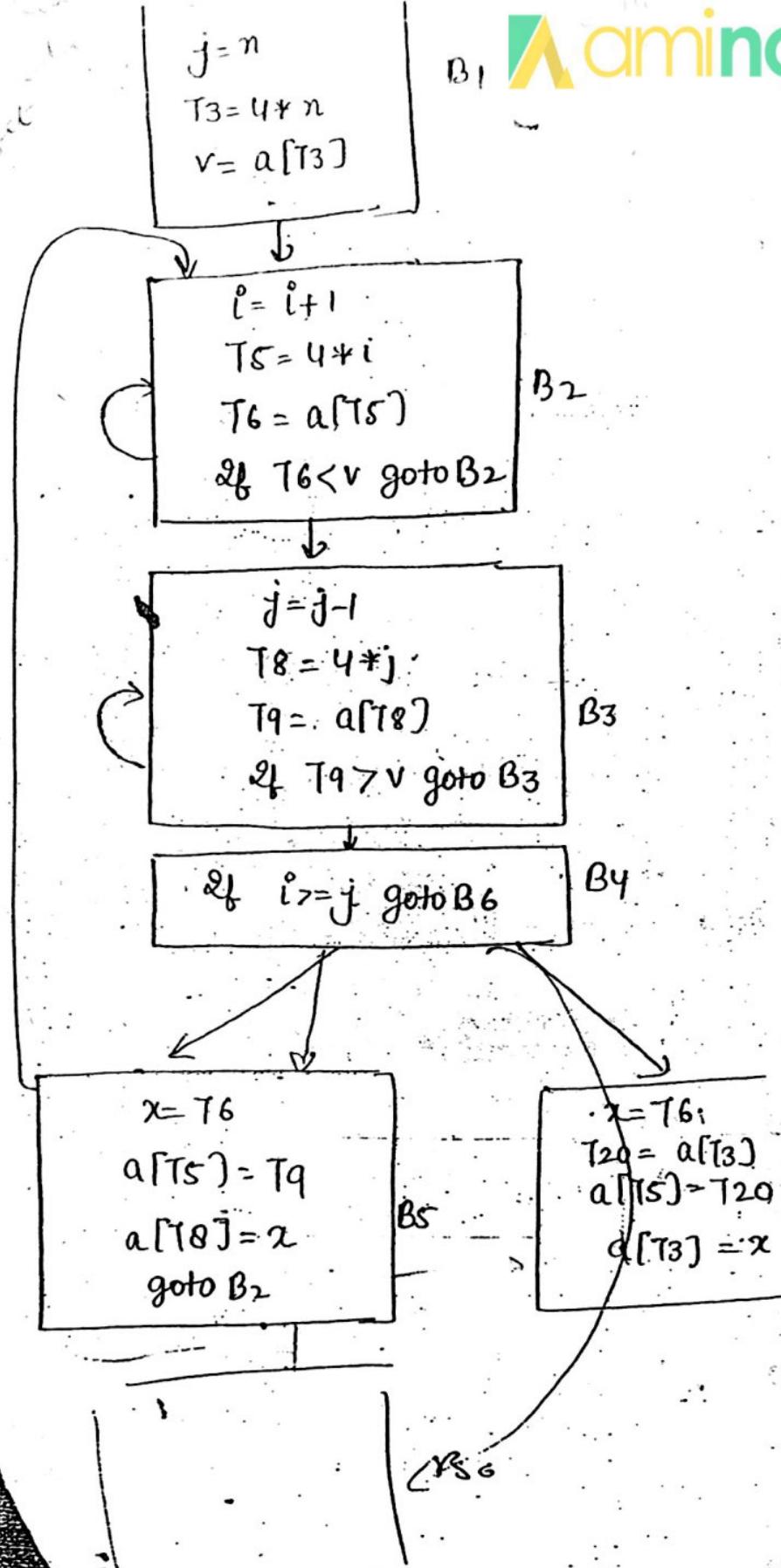
B₄

18)	$T_{10} = 4 * i$
19)	$T_{11} = a[T_{10}]$
20)	$x = T_{11}$
21)	$T_{12} = 4 * i$
22)	$T_{13} = 4 * j$
23)	$T_{14} = a[T_{13}]$
24)	$a[T_{12}] = T_{14}$
25)	$T_{15} = 4 * j$
26)	$a[T_{15}] = x$
27)	goto(7)

B₅

28)	$T_{16} = 4 * i$
29)	$T_{17} = a[T_{16}]$
30)	$x = T_{17}$
31)	$T_{18} = 4 * i$
32)	$T_{19} = 4 * n$
33)	$T_{20} = a[T_{19}]$
34)	$a[T_{18}] = T_{20}$
35)	$T_{21} = 4 * n$
36)	$a[T_{21}] = x$

B₆



(B)

loop optimization

$i = m - 1$
 $j = n$
 $T_3 = 4 + n$
 $v = a[T_3]$
 $T_5 = 4 + i$
 $T_8 = 4 + j$

B₁

$T_8 = T_5 + 4$
 $T_6 = a[T_5]$
 $\text{if } T_6 < v \text{ goto } B_2$

B₂

$T_8 = T_8 - 4$
 $T_9 = a[T_8]$
 $\text{if } T_9 > v \text{ goto } B_3$

B₃

$\text{if } T_5 > T_8 \text{ goto } B_6$

B₄

$a[T_5] = T_9$

$a[T_6] = T_8$

goto B₂

B₅

$T_{20} = a[T_3]$

$a[T_5] = T_{20}$

$a[T_3] = T_6$

B₆

Peephole Optimization

Peephole optimization is a method for buying improving the performance of the target program examining a short sequence of target instructions (called the peephole) and replacing these instructions by a shorter or faster sequence. Peephole is a small moving window on the target program.

Some peephole optimization characteristics are -

- 1) Redundant instructions elimination
- 2) Flow of control optimization
- 3) Strength reduction
- 4) Use of machine codes

Unreachable code

DEBUG = 0

If DEBUG = 1 then
print debugging information

TAC

DEBUG = 0

If DEBUG = 1 goto L1

goto L2

L1: print debugging information

L2: -