



RAJASTHAN TECHNICAL UNIVERSITY, KOTA

B.Tech Computer Science and Engineering (Artificial Intelligence)

RAJASTHAN TECHNICAL UNIVERSITY, KOTA

Syllabus

III Year-V Semester: B.Tech. Computer Science and Engineering

5CAI4-02: Compiler Design

Credit: 3

Max. Marks: 100(IA:30, ETE:70)

3L+0T+0P

End Term Exam: 3 Hours

SN	Contents	Hours
1	Introduction: Objective, scope and outcome of the course.	01
2	Introduction: Objective, scope and outcome of the course. Compiler, Translator, Interpreter definition, Phase of compiler, Bootstrapping, Review of Finite automata lexical analyzer, Input, Recognition of tokens, Idea about LEX: A lexical analyzer generator, Error handling.	06
3	Review of CFG Ambiguity of grammars: Introduction to parsing. Top down parsing, LL grammars & passers error handling of LL parser, Recursive descent parsing predictive parsers, Bottom up parsing, Shift reduce parsing, LR parsers, Construction of SLR, Conical LR & LALR parsing tables, parsing with ambiguous grammar. Operator precedence parsing, Introduction of automatic parser generator: YACC error handling in LR parsers.	10
4	Syntax directed definitions; Construction of syntax trees, S-Attributed Definition, L-attributed definitions, Top down translation. Intermediate code forms using postfix notation, DAG, Three address code, TAC for various control structures, Representing TAC using triples and quadruples, Boolean expression and control structures.	10
5	Storage organization; Storage allocation, Strategies, Activation records, Accessing local and non-local names in a block structured language, Parameters passing, Symbol table organization, Data structures used in symbol tables.	08
6	Definition of basic block control flow graphs; DAG representation of basic block, Advantages of DAG, Sources of optimization, Loop optimization, Idea about global data flow analysis, Loop invariant computation, Peephole optimization, Issues in design of code generator, A simple code generator, Code generation from DAG.	07
	Total	42

**Programme: B.Tech. CS(AI)****Semester: V****CourseName(CourseCode):COMPILER DESIGN (5CAI4-02)****CourseOutcomes****After completion of this course, students will be able to -**

5CAI4-02.1	Discuss the major phases of compilers and to use the knowledge of Lex tool
5CAI4-02.2	Develop the parsers and experiment the knowledge of different parsers design without automated tools.
5CAI4-02.3	Describe intermediate code representations using syntax trees and DAG's as well as use this knowledge to generate intermediate code in the form of three address code representation.
5CAI4-02.4	Classify various storage allocation strategies and explain various data structures used in symbol tables
5CAI4-02.5	Summarize various optimization techniques used for data flow analysis and generate machine code from the source code of an ovel language.

Name of Faculty:**(Signature)****Verified by Course Coordinator****Signature****(Name:)****Verified by Verification and Validation Committee, DPAQIC****Signature****(Name:)**



RAJASTHAN TECHNICAL UNIVERSITY, KOTA

B.Tech Computer Science and Engineering (Artificial Intelligence)



Swami Keshvanand Institute of
Technology, Management & Gramothan, Jaipur

C O	Outcomes	Bloom's Level	PO Indicators	PSO Indicators
Upon successful completion of this course, students should be able to:				
5CAI0 4-02.1	Discuss the major phases of compilers and to use the knowledge of Lex tool.	2,3	PO(1.3.1, 1.4.5, 2.1.1, 2.1.2, 2.1.3, 2.2.1, 2.2.2, 2.2.3, 2.2.4, 5.1.1, 5.1.2, 5.2.1, 12.1.1, 12.1.2, 12.2.1, 12.2.2)	PSO 1.1.1, 1. 1.3
5CAI04 -02.2	Develop the parsers and experiment the knowledge of different parsers design without automated tools.	3,6	PO(1.3.1, 1.4.1, 2.1.1, 2.1.2, 2.1.3, 2.2.1, 2.2.2, 2.2.3, 2.2.4, 2.3.2, 2.4.3, 2.4.4, 3.1.1, 3.1.6, 3.2.3, 3.4.2, 4.1.1, 4.1.2, 4.2.1, 4.3.1, 4.3.2, 4.3.3, 4.3.4, 12.1.1, 12.1.2, 12.2., 12.2.2)	
5CAI0 4-02.3	Describe intermediate code representations using syntax trees and DAG's as well as use this knowledge to generate intermediate code in the form of three address code representation.	1,3	PO(1.3.1, 1.4.1, 2.1.1, 2.1.2, 2.1.3, 2.2.1, 2.2.2, 2.2.3, 2.2.4, 2.3.2, 2.4.3, 2.4.4, 3.1.1, 3.1.6, 3.2.1, 3.2.3, 3.4.2, 4.1.1, 4.1.2, 4.2.1, 4.3.1, 4.3.2, 4.3.3, 4.3.4, 12.1.1, 12.1.2, 12.2.1, 12.2.2)	PSO1.1.1

UNIT-I

COURSE DESCRIPTION

The goal of the course is to provide an introduction to the system software like assemblers, compilers, and macros. It provides the complete description about inner working of a compiler. This course focuses mainly on the design of compilers and optimization techniques. It also includes the design of Compiler writing tools. This course also aims to convey the language specifications, use of regular expressions and context free grammars behind the design of compiler.

COURSE OBJECTIVES

The objective of this course is to:

1. Provide an understanding of the fundamental principles in compiler design
2. Provide the skills needed for building compilers for various situations that one may encounter in a career in Computer Science.
3. Learn the process of translating a modern high-level language to executable code required for compiler construction.

COURSE OUTCOMES

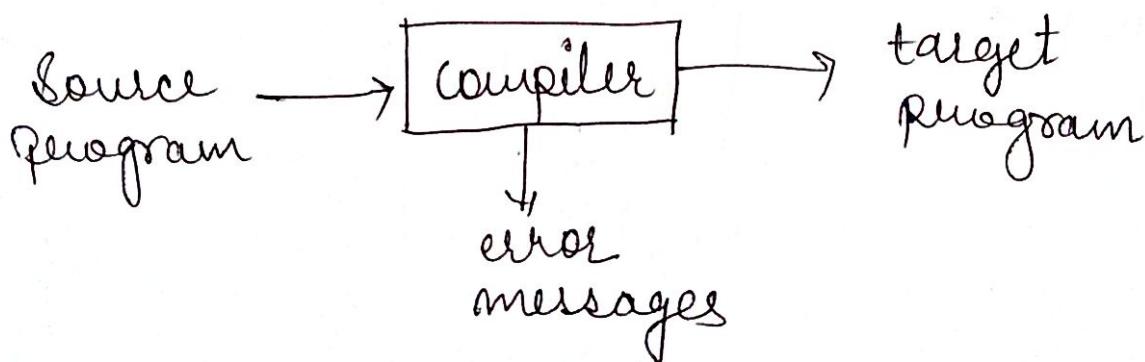
At the end of course students will be able to:

1. Discuss the major phases of compilers and to use the knowledge of Lex tool
2. Develop the parsers and experiment the knowledge of different parsers design without automated tools.
3. Describe intermediate code representations using syntax trees and DAG's as well as use this knowledge to generate intermediate code in the form of three address code representation.
4. Classify various storage allocation strategies and explain various data structures used in symbol tables
5. Summarize various optimization techniques used for dataflow analysis and generate machine code from the source code of a novel language.

Introduction to compilers:-

Translators :- a translator is a generic term that could refer to a compiler, assembler or interpreter; anything that converts higher level code into another high-level code or lower-level language.

Compiler :- is a program that reads a program written in one language - the source language - and translates it into an equivalent program in another language - the target language. In this process, compiler reports to its user the presence of errors in source code.

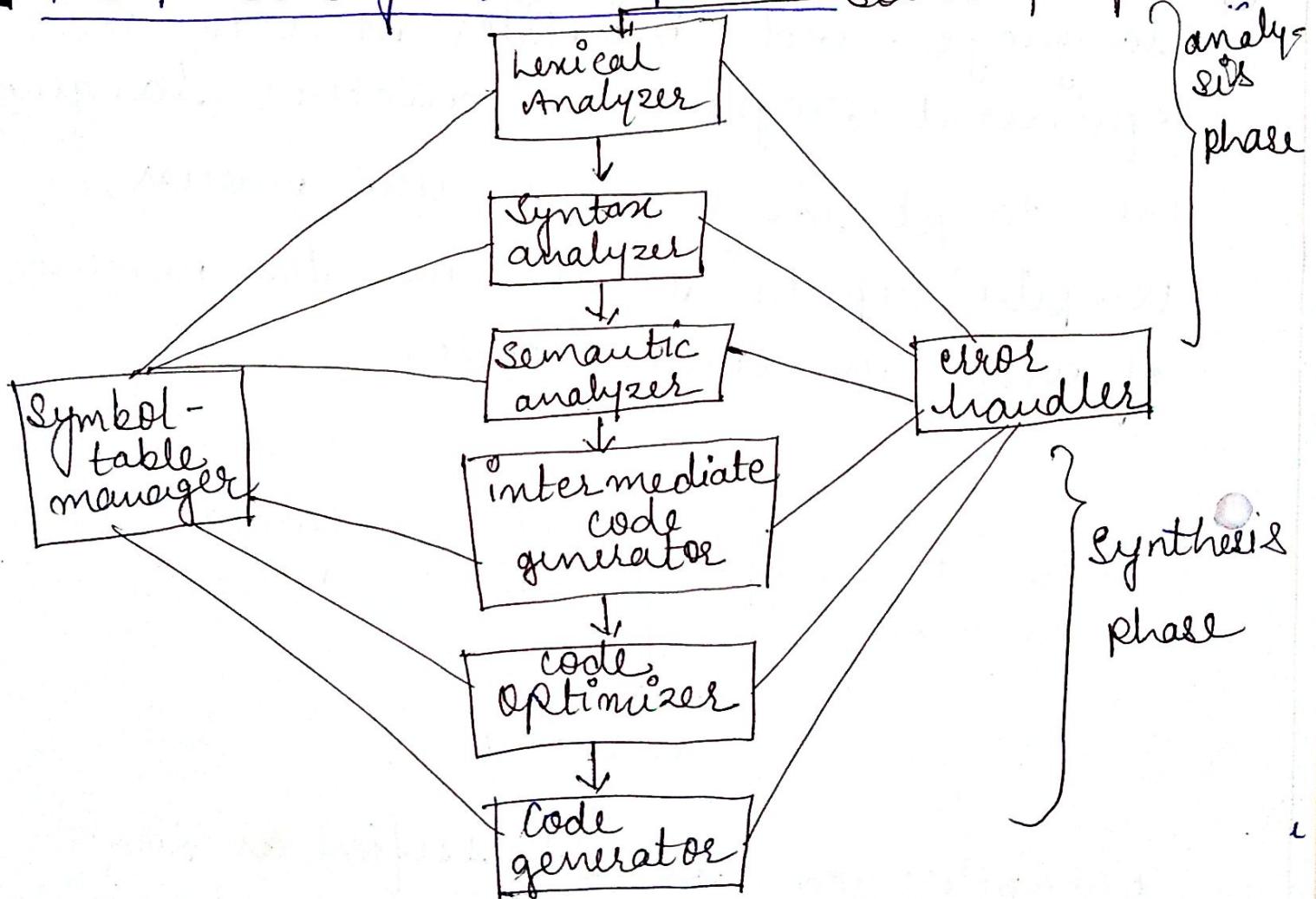


* compilers can be classified as single-pass, multi-pass, load-and-go, debugging etc.

Interpreter:-

An interpreter translates code like a compiler but reads the code and immediately executes on that code and therefore is initially faster than a compiler. Interpreters execute a single line of code at a time.

The phases of a compiler:- source program

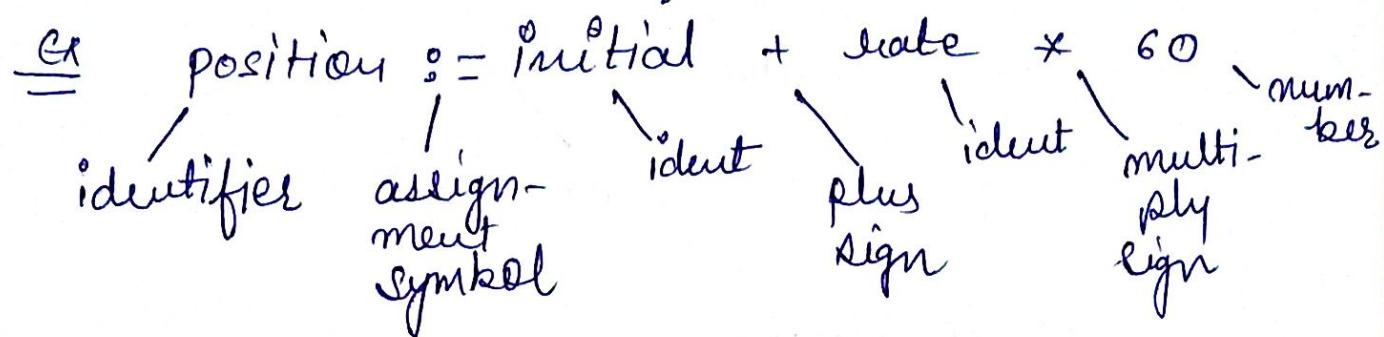


The Analysis Phase:-

(3)

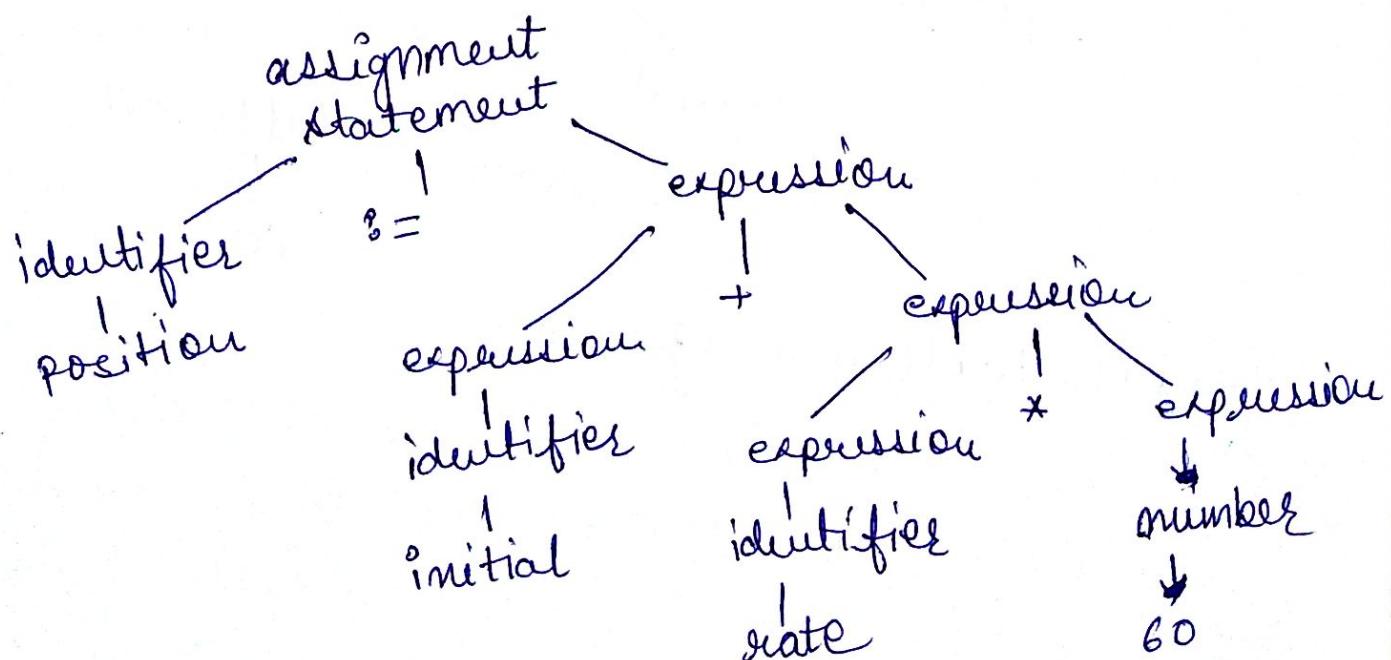
Syntactical Analysis:-

- Also called linear analysis.



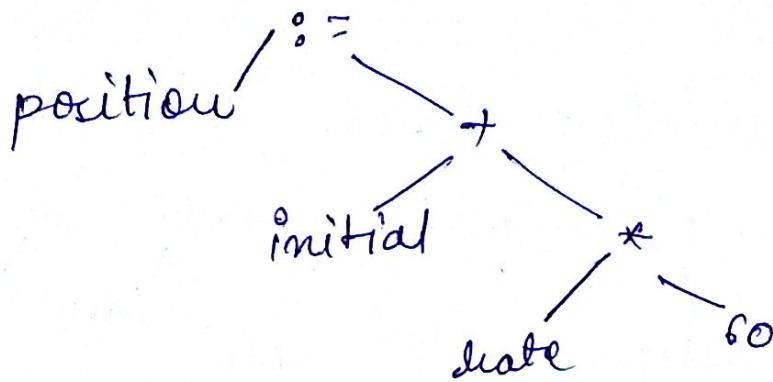
Syntactic Analysis:-

- Also called hierarchical analysis / parsing.
- It involves grouping the tokens of the source program into grammatical phrases that are used by the compiler to synthesize output.
- Represented by parse tree.



- A more common internal representation of above tree,

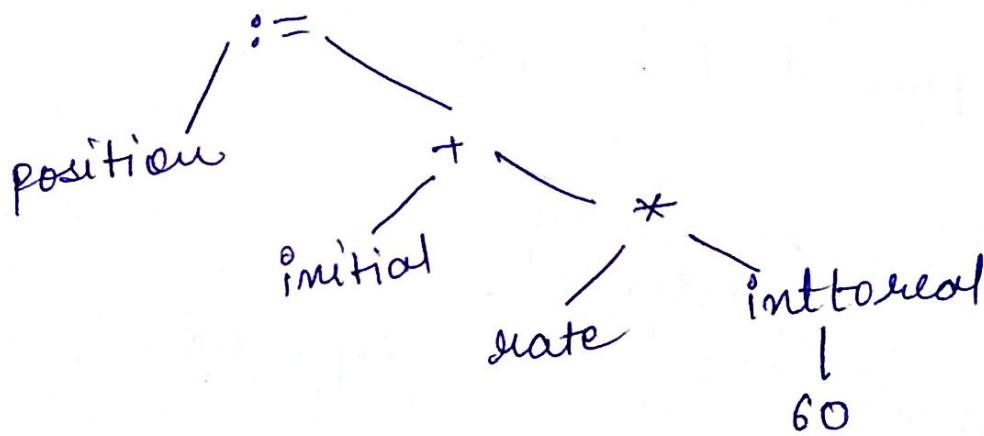
(4)



3) Semantic Analysis:-

- Checks the source program for semantic errors and type checking.

Ex : Assume in above example all operands/identifiers are in real and 60 is an integer.



Symbol-table management:-

- A symbol table is a data structure containing a record for each identifier, with fields for the attributes of the identifier.

- fields can store attributes like storage allocated for identifier, its type, its scope, and in case of procedures it can store number/type of arguments etc. ⑤
- Symbol table allows us to find the record for each identifier quickly and to store/retrieve data from record quickly.
- When an identifier is detected in source program by lexical analyzer, the identifier is entered into the symbol table. The attributes of identifier is generally not known at this time.

Ex:- (Pascal) Var position, initial,
rate : real;

When lexical analyzer identifies the variable 'position', its data type will not known and entered into the symbol table by remaining phases.

Error Detection and Reporting:-

(6)

- Each phase can encounter errors.
- Syntax & semantic analysis phases usually handle a large fraction of errors.
- lexical phase can detect errors where the characters remaining in the I/P do not form any token.
- Errors where token stream violates the structure rules of the language are determined by the syntax analysis phase.
- Semantic phase detect constructs that have right syntactic structure but no meaning to the operation involved.
Ex:- Add two identifiers. One is array and one is procedure.

Intermediate Code Generation:-

- The intermediate representation should have two important properties:-
 - 1) it should be easy to produce
 - 2) it should be easy to translate into the target program.
- It can has multiple forms/methods.
Figure on page 7 has intermediate form called "Three-address code".

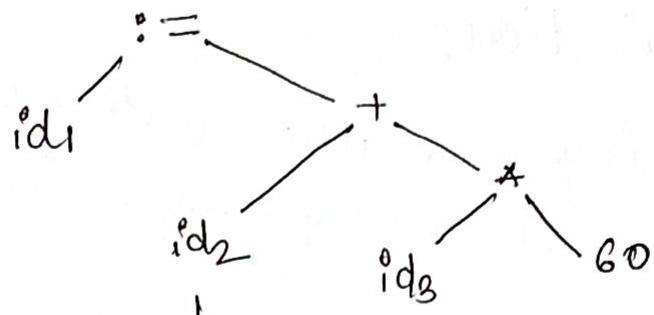
Position := initial + rate * 60

⑦

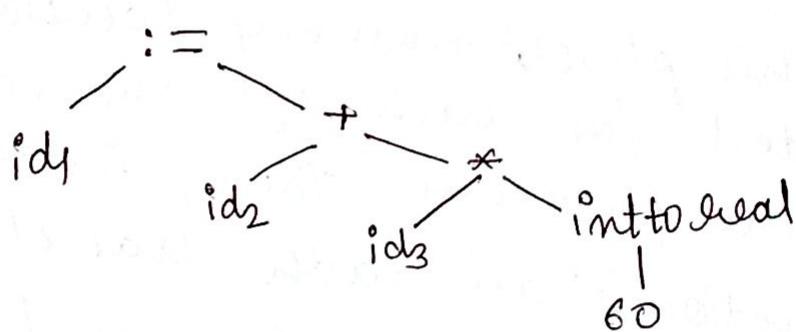
↓
lexical analyzer

↓
 $id_1 := id_2 + id_3 * 60$

↓
syntax analyzer



↓
semantic analyzer



↓
intermediate code generator

$temp1 := \text{inttoreal}(60)$

$temp2 := id_3 * temp1$

$temp3 := id_2 + temp2$

$id_1 = temp3$

↓
code optimizer

$temp1 := id_3 * 60.0$

$id_3 := id_2 + temp1$

↓

Code generator

```

MOVF id3, R2
MULF # 60.0, R2
MOVF id2, R1
ADDF R2, R1
MOVF R1, id1

```

Code Optimization:-

- Attempts to improve the intermediate code, so that faster-running machine code will result.

Code generation:-

- In this phase, memory locations are selected for each of the variables used by the program. Then, intermediate instructions are each translated into a sequence of machine instructions that perform the same task.

Bootstrapping:-

- Bootstrapping is the technique for producing a self-compiling compiler, that is compiler written in the source programming language that it intends to compile.

- Using the facilities offered by a language to compile itself is the essence of bootstrapping.
- Use of bootstrapping to create compilers and to move them from one machine to another by modifying the back end.
- For bootstrapping purposes,

$S \rightarrow$ source language

$T \rightarrow$ Target language

$I \rightarrow$ in which language compiler is written

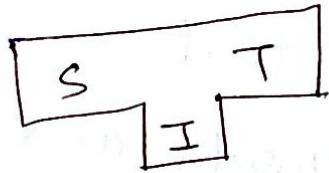
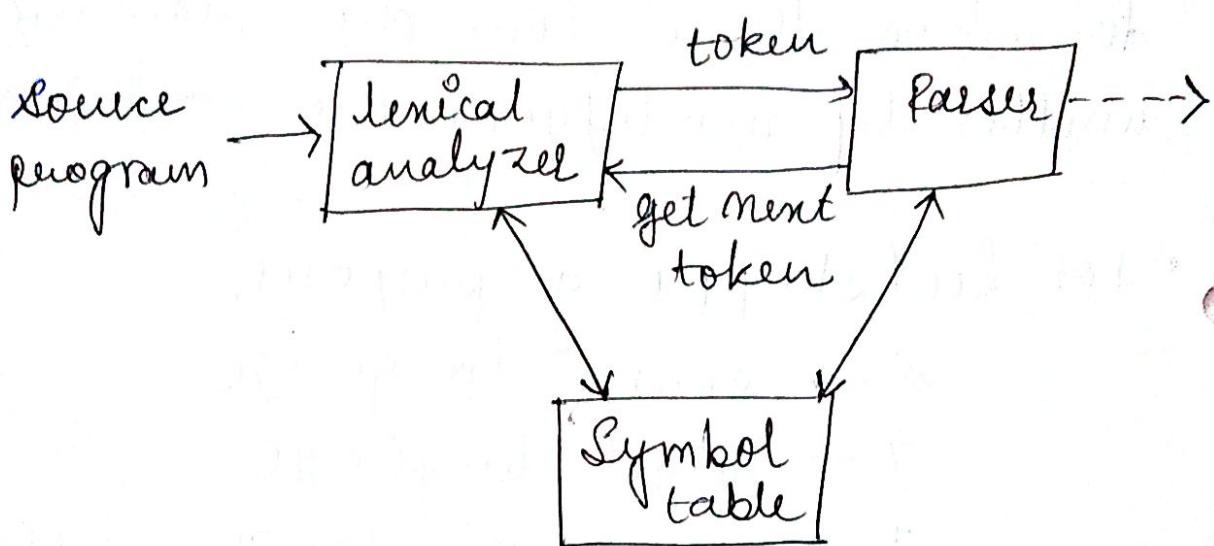


Fig: T-Diagram

- $S, T \& I$ can be different.
- A compiler may run on one machine and produce target code for another machine. Such a compiler is often called a cross-compiler.

- Lexical Analysis:-** (Role of lexical analyzer)
- Its main task is to read the I/P characters and produce as O/P a sequence of tokens that the parser uses for syntax analysis.



Lexical analyzer's tasks:-

- 1) Reads Source text
- 2) Strip out comments and white spaces from source program
- 3) correlating error messages from the compiler with the source program.
(Keep track of newline character for giving line number in error message)
- 4) Makes a copy of source program with the error messages marked in it.
- 5) Preprocessor functions may be implemented

as lexical analysis takes place.

(11)

Issues in lexical analysis:-

- Reasons for separating analysis phase i.e. lexical analysis and parsing,
 - It simplifies the design of lexical analyzer already eliminate the comments/white spaces before passing)
 - Compiler efficiency is improved (specialized buffering techniques for reading I/P characters and processing tokens can significantly speed up the performance of a compiler)
 - compiler portability is enhanced (I/P alphabet peculiarities and other device-specific anomalies can be restricted to the lexical analyzer.)

Tokens, Patterns, hexemes:-

- Following constructs are treated as tokens:
keywords, operators, identifiers, constants, literal strings and punctuation symbols.

- (12)
- Set of strings is described by a rule called a pattern associated with the token.
 - A lexeme is a sequence of characters in the source program that is matched by the pattern for a token.
Ex:- const pi = 3.1416;
pi is a lexeme for the token "Identifier".

<u>Token</u>	<u>sample lexemes</u>	<u>Informal description of pattern</u>
1) const	const	const
2) if	if	if
3) relation	<, <=, =, <=, >, >=	any kind of relational operators
4) id	pi, count, R2	letter followed by letters & digits.
5) num	3.14, 0, 6.02	any numeric constant

~~Attribute~~ Attributes for Tokens:-

- When more than one pattern matches a lexeme, the lexical analyzer must provide additional information about it, so that additional attributes are provided to token. → so that code generator can know what string was actually matched.
- Ex:- the pattern num matches both the string 0 and 1.

int a = 2 free additional
 float a = 2 attribute (int, float)
 tells which kind of no a is.

Ex:- Tokens and associated attributes - values,

$e = M * C ^ { \star \star 2 }$

(id, pointer to symbol-table entry for E)

< assign-op, >

< id, pointer _____ M >

< mult-op, >

< id, pointer _____ C >

< exp-op, >

< num, integer value 27 >

Lexical errors:-

(K4)

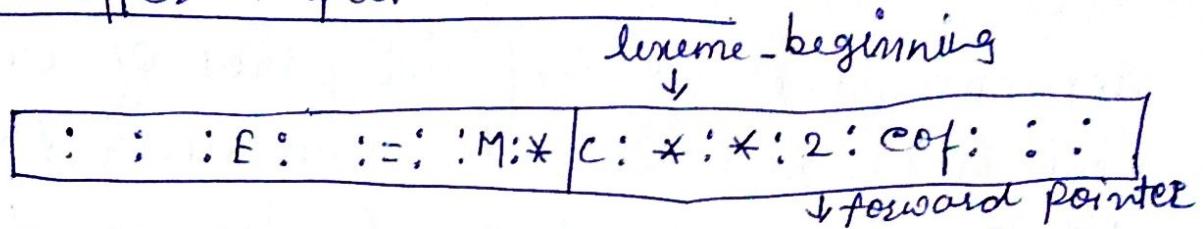
- Few errors are seen at the lexical level alone.

Ex. $f^o(a = - f(x)) \dots$

- Lexical analyzer will return 'f^o' as identifier.
- Suppose a situation in which lexical analyzer is unable to proceed bcoz none of the patterns for tokens matches a prefix of the remaining input.
- possible error-recovery actions,
 - ~~deleting~~
 - "panic mode" recovery method,
we delete successive characters from the remaining i/p until the lexical analyzer can find a well-formed token.
 - delete an extra char
 - insert a missing char
 - replace an incorrect char by correct char
 - transposing two adjacent char.

Input Buffering:-

Two-buffer input scheme:-



- we divide the buffer in two halves, each having capacity of N . N is the no of chars on one disk block, eg, 1024 or 4096.
- Read N chars ~~as~~ into each half.
- Two pointers extreme - beginning and forward. The string of chars b/w two pointers is the current sentence.
- If the forward pointer is at the end of first half then right buffer is filled with N chars and if forward pointer is at the end of second half then left buffer is filled with N chars.
- Problem is that each time; ^(while moving) we have to check that we have not moved off one half of the buffer, so that we must reload other half. This will increase the overload. So we use a special sentinel char to ~~check whether~~

at end of ~~the~~ each half ~~are~~ ~~be~~ ~~not~~.

(1)

:	:	E:	:	=	:	M:	*	:	eob		c:*	*	:	2: eof:	!	:	eob
---	---	----	---	---	---	----	---	---	-----	--	-----	---	---	---------	---	---	-----

- we can take as eof in place of eob, but we have put special mechanism to check whether it is end of file or end of buffer.

specification of tokens:-

- regular expressions are used to specifying patterns.

String and Languages:-

- Alphabet or Char class can denotes any finite set of symbols.
Ex: - {0,1} binary class
- A string over some alphabet is a finite sequence of symbols drawn from that alphabet.

Ex:- Symbol $\rightarrow s$

Empty string $\rightarrow \epsilon$

length of string $\rightarrow |s|$

- The term language denotes any set of strings over some fixed alphabet.

= Abstract languages $\rightarrow \emptyset$ (17)
empty set $\rightarrow \{ \}$

$$S \subseteq E^* = S$$

= Ex: - $x = \text{dog}$, $y = \text{house}$
 $xy = \text{doghouse}$

Operations on languages:-

1) union of L and M (LUM)

$$LUM = \{ s \mid s \text{ is in } L \text{ or } s \text{ is in } M \}$$

2) concatenation of L and M (LM)

$$LM = \{ st \mid s \text{ is in } L \text{ and } t \text{ is in } M \}$$

3) Kleene closure of L (L^*)

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

4) Positive closure of L (L^+)

$$L^+ = \bigcup_{i=1}^{\infty} L^i$$

Regular Expressions:-

• Ex: - identifier: - a letter followed by zero or more letters or digits.
 $\text{letter} (\text{letter/digit})^*$

Regular Definitions:-

We may wish to give names to regular expressions and to define regular expressions using these names as if they were symbols.

If Σ is an alphabet of basic symbols, then a regular definition is a sequence of definitions of the form,

$$d_1 \rightarrow e_1$$

$$d_2 \rightarrow e_2$$

:

$$d_n \rightarrow e_n$$

Ex:- letter $\rightarrow A | B | \dots | z | a | b | \dots | z |$

digit $\rightarrow 0 | 1 | \dots | 9 |$

id \rightarrow letter (letter | digit)*

Notational Shorthands:-

1) One or more instances: - + is used for this.

2) Zero or one instance: - ? is used for this.

3) Character classes:-

Ex:- $[A - z a - z] [A - z a - z 0 - 9]^*$

Non regular Sets :-

18

- Some languages cannot be described by any regular expression.

Ex:- Repeating strings cannot be described by regular expressions.

Push Down Automata \rightarrow If $w \in \{a, b\}^*$ is a string of a, b then w is a string of a, b ?

- Regular expressions can be used to denote only a fixed no of repetitions.

■ Recognition of Tokens :-

Ex:- $\text{stmt} \rightarrow \text{if expr then stmt}$
 $\quad \quad \quad \quad | \quad \text{if expr then stmt else stmt}$
 $\quad \quad \quad \quad | \quad \epsilon$

$\text{expr} \rightarrow \text{term eelop term}$
 $\quad \quad \quad \quad | \quad \text{term}$

$\text{term} \rightarrow \text{id} / \text{num}$

Regular definitions

$\text{if} \rightarrow \text{if}$

$\text{then} \rightarrow \text{then}$

else or else

$\text{eelop} \rightarrow < | \leq | = | \geq | > | \geq$

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

$\text{num} \rightarrow \text{digit}^+ (.\text{digit}^+)? (\text{E} (+/-) \text{digit}^+)?$

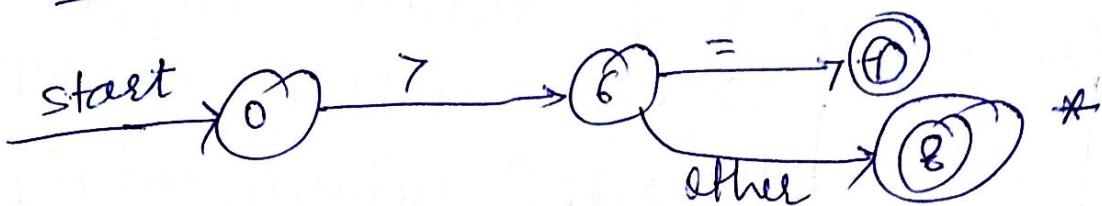
- In addition, we assume tokens are ~~are~~
separated by white space,
 $\text{delim} \rightarrow \text{blank} \mid \text{tab} \mid \text{newline}$
 $\text{ws} \rightarrow \text{delim}^+$
- Our goal is to construct a lexical analyzer that will indicate the lexeme for the next token in the i/p buffer and produce as o/p a pair consisting of the appropriate token and attribute value, using the translation table,

Regular Expression	Token	Attribute Value
ws	-	-
if	if	-
then	then	-
else	else	pointer to tables
id	id	"
num	num	"
	develop	LT
<	"	LB
<=	"	EQ
=	"	NC
<>	"	"
>	"	GT
>=	"	GE

Transition Diagram:-

- Transition diagrams are intermediate step in the construction of a lexical analyzer.
- These diagrams depict the actions that take place when a lexical analyzer is called by the parser to get the next token.
 - positions in diagram - states (circles)
 - connected by - edges (arrow)
- We assume these diagrams are deterministic.

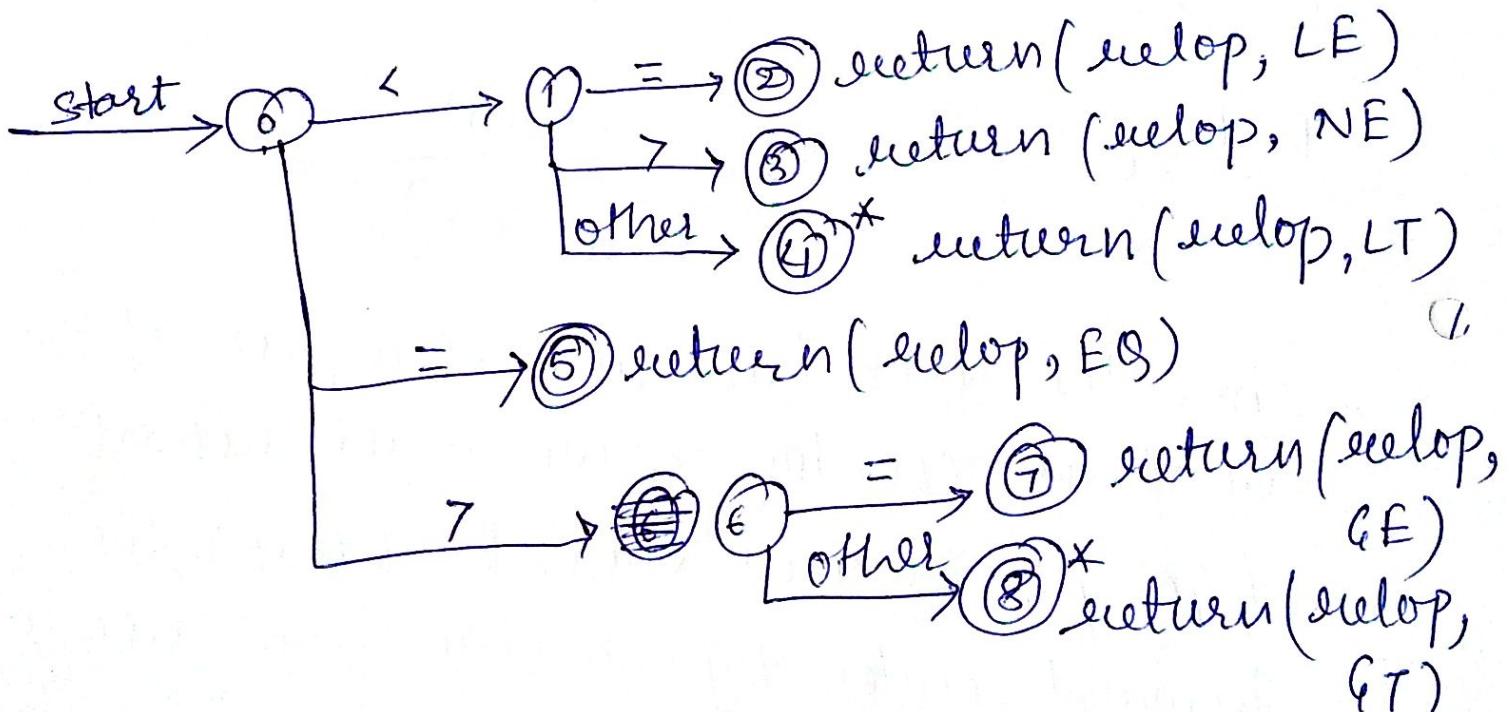
Ex $\geq =$



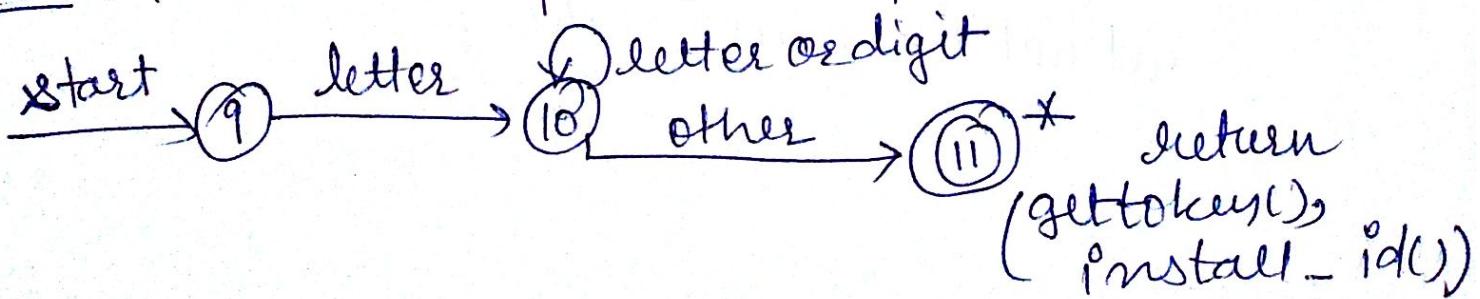
- Since the extra char is not a part of the relational operator \geq , we must extract (going to previous state, backtract) the forward pointer by one char. We use a * to indicate states on which this i/p extraction must take place.

- In general, there may be several transition diagrams, each specifying a group of tokens. If failure occurs while we are following one transition diagram, then we extract the forward pointer to where it was in the start state of this diagram, and activate the next transition diagram. (Forward pointer is extracted to the position marked by the leneme - beginning pointer). (22)

Ex :- Transition diagram for $\ell\ell\ell\phi$,



ex : transition diagram for keywords & identifiers



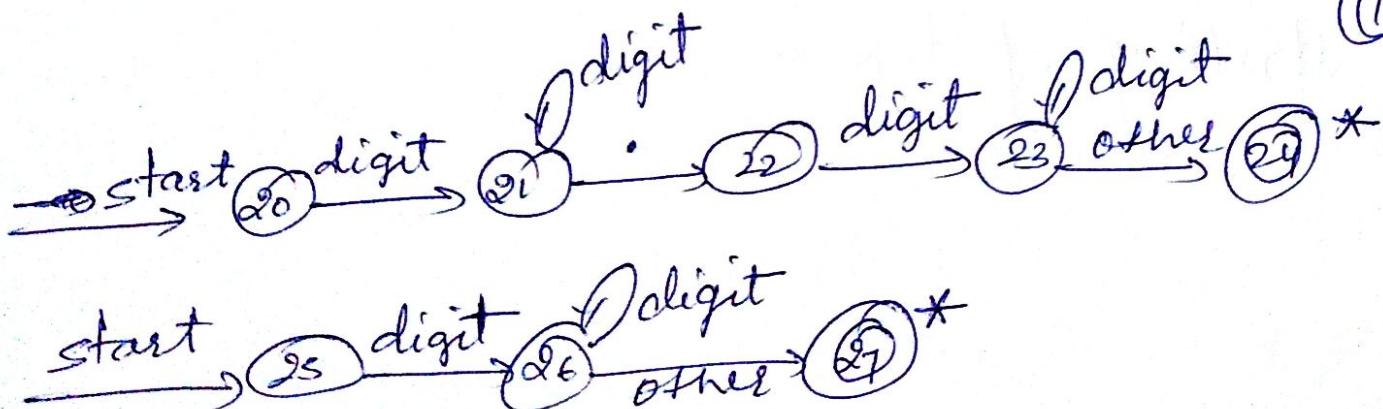
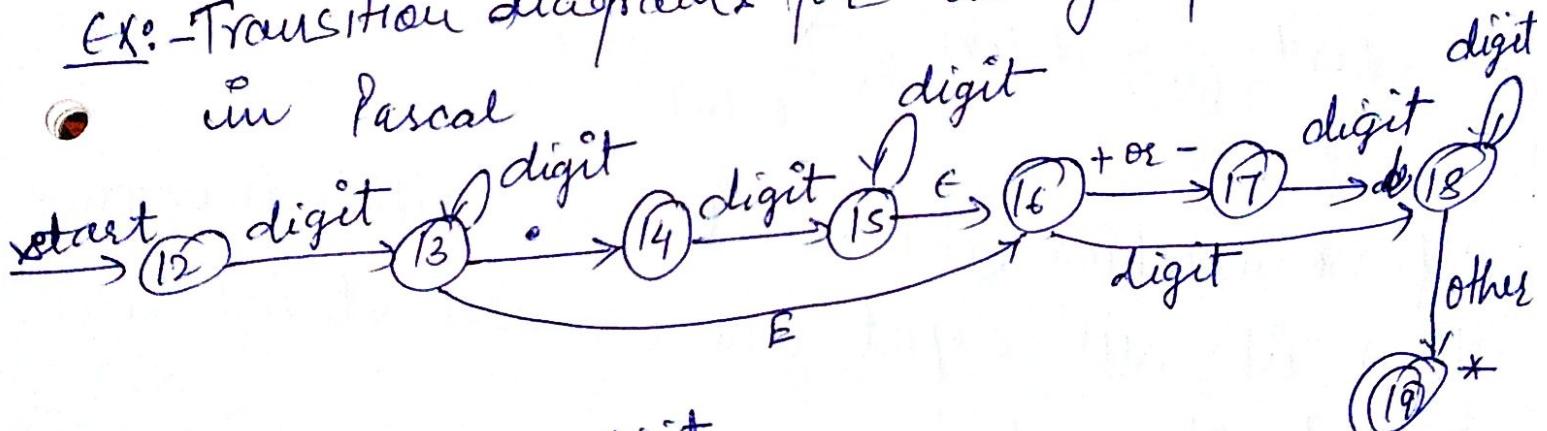
(23)

In starting all keywords / identifiers are detected by above transition diagrams. Then matched with symbol table, if exist then it will be a keyword otherwise identifier. All keywords are already placed in symbol table before i/p any char.

- install-id() returns 0 if i/p is keyword according to symbol table. Otherwise it will returns a pointer to the symbol table entry.
- getToken() returns the token if i/p is token, otherwise token id is returned.

Ex:- Transition diagrams for unsigned numbers

in Pascal



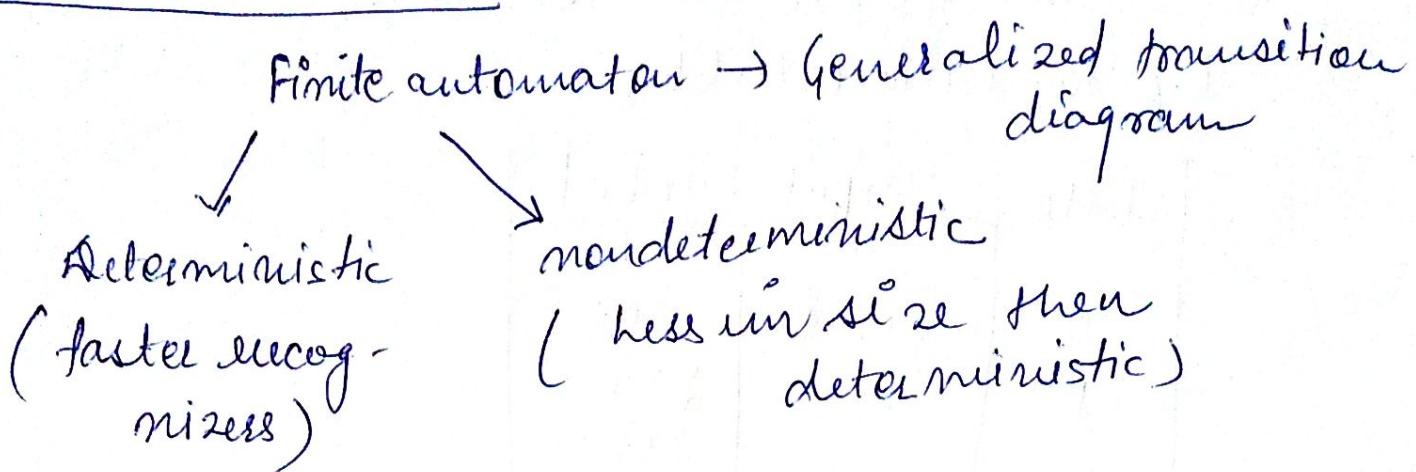
- Error-handling: - Info^r about the language
that is not in the regular definition of the
tokens can be used to pinpoint errors in
i/p. Ex:- i/p $\rightarrow 1.0 \langle x,$
we fail in states 14 and 22
 - ~~with next i/p char~~
 - Rather than returning the no 1, we may
wish to report an error and continue as if
the i/p were $1.0 \langle x$. Such knowledge can
be used to simplify the transition dia-
gram, bcoz error-handling may be
used to recover from the situations
that would otherwise lead to failure.

Ex:- white Space



If we combine above all transition discourses, then it will depict the example stated in the starting of topic.

Finite Automata:-



Ex:- $(a|b)^*abb$

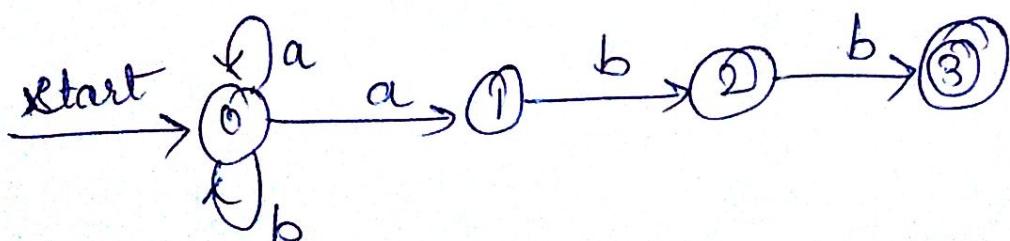
nondeterministic finite Automata:-

NFA is a mathematical model that consists of,

- 1) $S \rightarrow$ states
- 2) $P \Sigma \rightarrow$ set of i/p symbols
- 3) Transition function
- 4) $s_0 \rightarrow$ start state
- 5) $F \rightarrow$ set of accepting states

Represented by transition graph, in this same graph can have two or more transitions out of one state and edges can be labeled by ϵ as well.

Ex:- $(a|b)^*abb$



- Transition graph can be implemented by transition table,

26

State	Input Symbol	
	a	b
0	{0, 1}	{0}
1	-	{2}
2	-	{3}

Advantage: - fast access to transitions

Disadvantage: - take lot of space

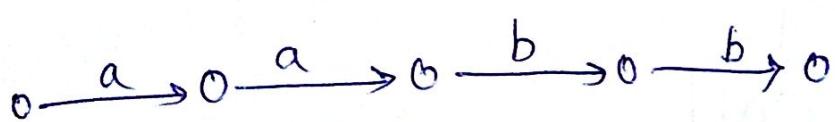
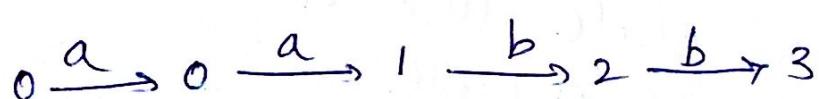
- Adjacency list representation can also be used,

Advantage: - compact

Disadvantage: - access is slower

- In a NFA, same i/p string can lead to accepting state by different sequences,

Ex: - aabb



Deterministic Finite Automata:-

- DFA is a special case of NFA,
 - no state has a transition
 - for each state s and i/p symbol a , there is at most one edge labeled a leaving s

Design of a lexical analyzer generator:-

- Aim is to design a s/w tool that automatically constructs a lexical analyzer from a program in the lex language.

- Specifications $P_1 \{ \text{action}_1 \}$

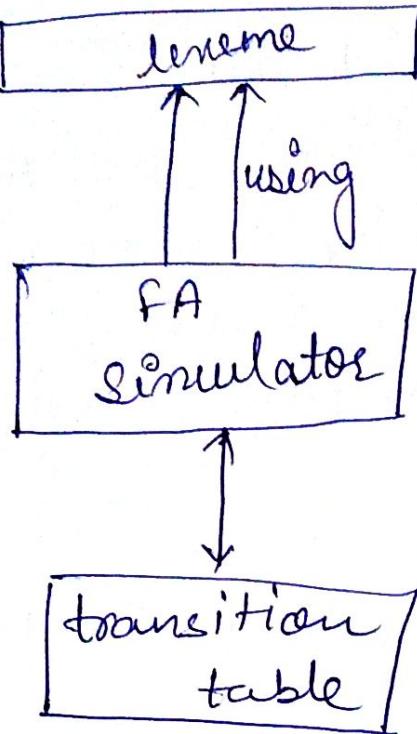
$P_2 \{ \text{action}_2 \}$

:

$P_m \{ \text{action}_m \}$

P_i = pattern i is a regular expression & each action_i is a program fragment which is to be executed whenever a lexeme matched by P_i .





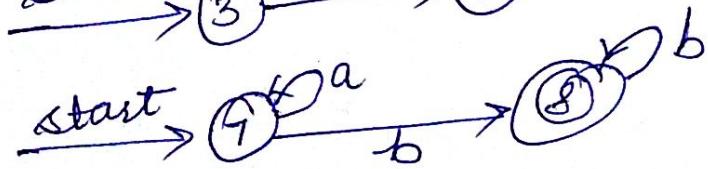
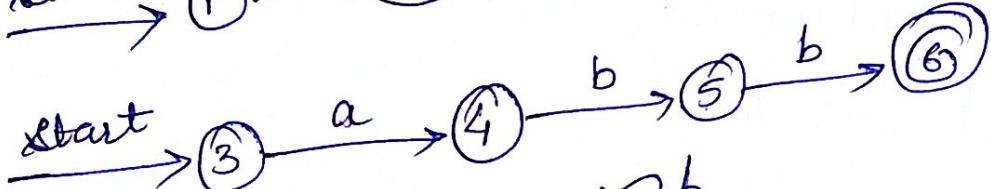
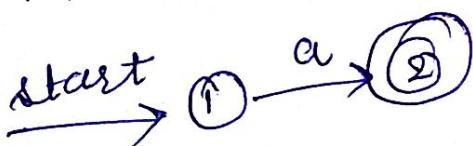
- Finite automata simulator uses transition table to look for the regular expression patterns in the i/p buffer.

Pattern matching based on NFA's :-

Ex a^* { }
 = abb { }
 a* b* { }

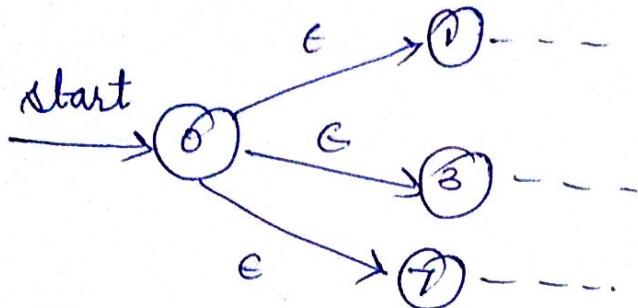
* actions are omitted here*

a) NFA for each regular exp

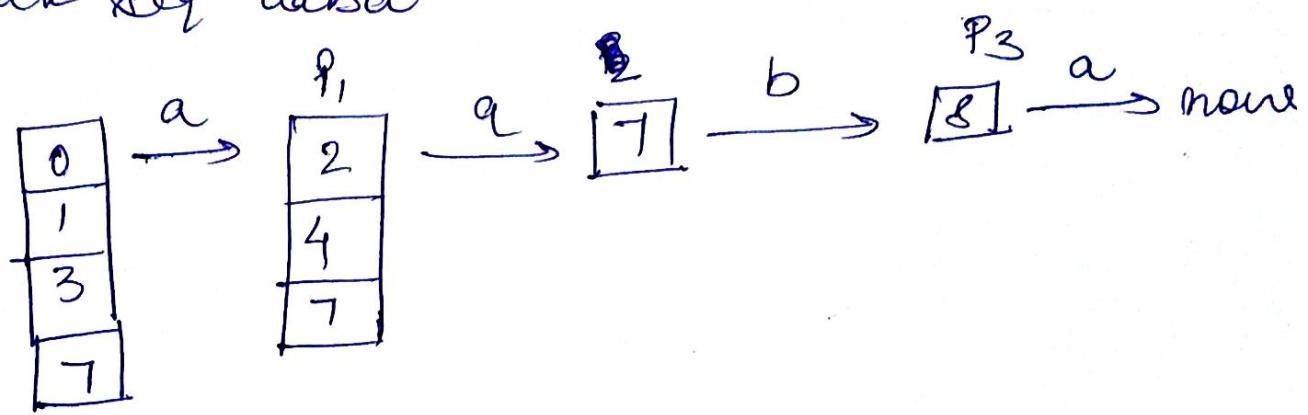


b) Combined NFA

(29)



c) Check seq aaba



NFA for lexical analysis:-

- Convert NFA into DFA (transition table).
- Use transition table to check the seq.

state	I/P symbol		pattern
	a	ab	
0137	241	8	none
247	7	58	a
8	-	8	a^*b^+
7	7	8	none
58	-	BB	a^*b^+
68	-	8	abb

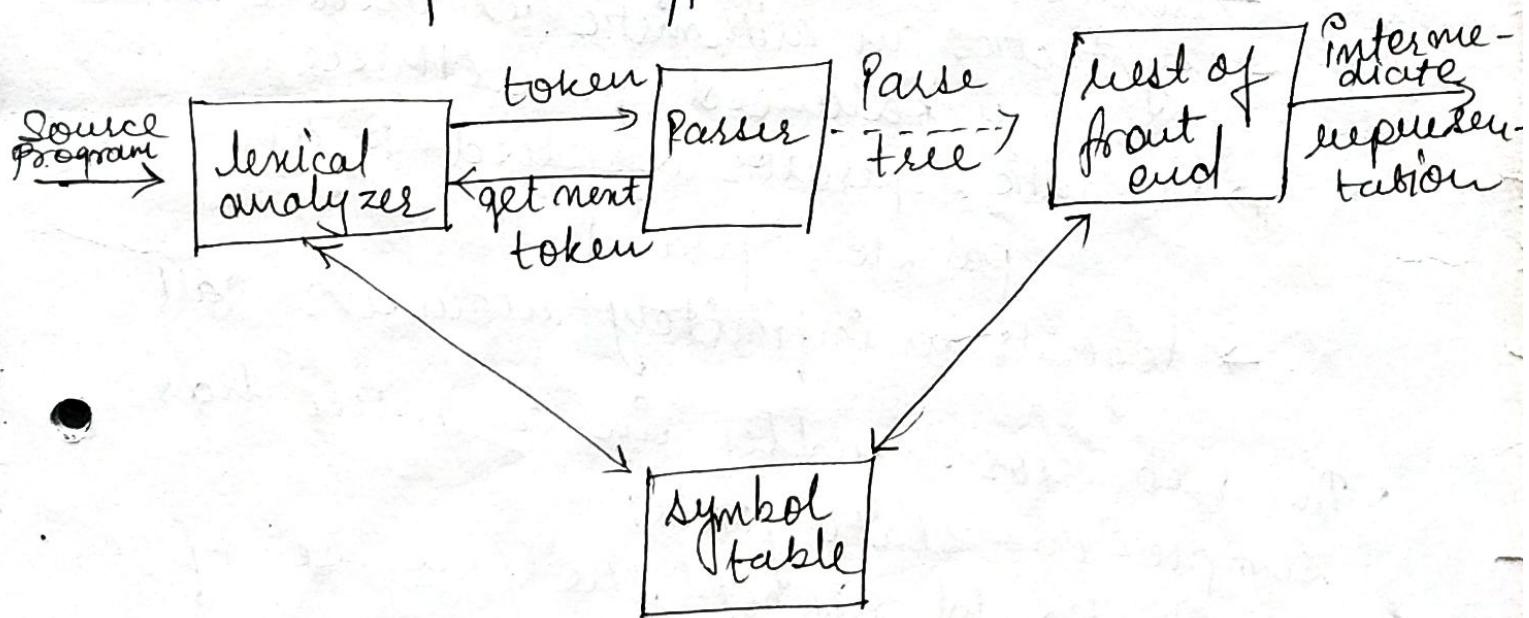
UNIT-III

30

Syntax Analysis

■ Introduction (Role of parser):-

- parser obtains a string of tokens from the lexical analyzer and verifies that the string can be generated by the grammar for the source language.
- It should also recover from commonly occurring errors so that it can continue processing the remainder of its I/P.



• Grammars

- top-down (build parse tree from top to bottom)
- bottom-up

Syntax Error Handling :-

- A good compiler should assist the programmes in identifying and locating errors.
- Few languages have been designed with error handling.
- Programming language do not describe how a compiler should respond to errors; the response is left to the compiler design.
- Errors can be,
 - lexical: misspelling an identifier, keyword
 - syntactic: an arithmetic expression with unbalanced parentheses
 - semantic: operator applied to an incompatible operand
 - logical: an infinitely recursive call
- The ~~error~~ error handler in a parser has simple-to-state goals:
 - it should report the presence of errors clearly and accurately.
 - should recover from each error quickly
 - should not significantly slow down the processing of correct programs.

- Common errors are simple ones.

- In LL and LR methods, they have the viable-prefix property, meaning they detect that an error has occurred as soon as they see a prefix of I/p that is not a prefix of any string in the language.

Ex: common errors can be,

- use a comma at place of semicolon.
- Forgot semicolon.

- Rules can vary from language to language. A compiler should report the place in source program where an error is detected.
- In some parsers, when want to recover from error, some data can be unread. So that some variables can be unknown.

Error-Recovery strategies:-

→ Panic Mode:-

- simplest method & can be used by most parsing methods.
- On discovering an error, the parser discards I/P symbols one at a time until one of a designated set of synchronizing tokens (semicolon, end) is found.
- This method often skips a considerable amount of I/p without checking it for additional errors.

2) Phrase-level Recovery:-

- (33)
- On discovering an error, this method performs local corrections on the remaining I/P; that is, it may replace a prefix of the remaining I/P by some string that allows the parser to continue.

Ex:- replace a comma by semicolon,
delete an extra semicolon, insert a
missing semicolon.

- We must be careful to choose replacement that do not lead to infinite loops.

3) Error Productions:-

- By using common errors we augment grammar and make parser for erroneous constructs. So that we can recognize the erroneous constructs in the I/P.

4) Global correction:-

- Given an incorrect I/P string x and grammar G , these algorithms will find a parse tree for a related string y , such that no of insertions, deletions and changes of tokens required to transform x into y is as small as possible.

- They try to obtain a globally least-cost correction.

- These methods are in general too costly in terms of time and space, so these techniques are currently only of theoretical interest.

Content-Free Grammars:-

- Some type of statements → In Regular exp → right side of production will be start from terminal only.
- cannot be specified using regular exp, so that CFG is used. → In CFG → right side of production can be terminal/non-terminal

Ex:- start → if expr then
stmt else stmt

→ CSG → It will have

- Terminals, non-terminals, same count in a start symbol & productions left and right side
- are diff. parts of CFG. → unrestricted - no rules

→ Terminals :- symbols from which strings are formed. Also called token.

- lower-case letters early in alphabet a, b, c
- operators +, -, *
- Punctuation like parentheses, comma
- Digits 0, 1, ..., 9
- Boldface string id or if.

→ Non terminals:- syntactic variables that denote set of strings.

- Upper-case letters early in alphabet

A, B, C

- letter s (start state)
- Lower-case italic names like expr / stmt .

→ Start state:- one nonterminal, and the set of strings it denotes is the language defined by grammar.

→ Productions specify the manner in which the terminals & nonterminals can be combined to form strings.

- Upper-case letters late in alphabet x, y, z represent grammar symbols, that is either non-terminals or terminals.
- lower-case letters late in alphabet u, v, w represent strings of terminals.
- lower-case Greek letters, α, β, γ represent strings of grammar symbols.
- if $A \rightarrow \alpha_1, A \rightarrow \alpha_2 \dots A \rightarrow \alpha_k$
then $A \rightarrow \alpha_1 | \alpha_2 \dots | \alpha_k$ called alternatives of A .
- The ~~start~~ left side of first production is start symbol.

Ex:- $E \rightarrow E \cdot O \cdot E$

$E \rightarrow (E)$

$E \rightarrow -E$

$E \rightarrow id$

$O \rightarrow +$

$O \rightarrow -$

$O \rightarrow *$

$O \rightarrow /$

$O \rightarrow 1$

According to above rules it can be stated
as,

$E \rightarrow EAB \mid (E) \mid -E \mid id \quad \text{--- (1)}$

$A \rightarrow + \mid - \mid * \mid / \mid 1$

Derivations :-

- A production is treated as a rewriting rule in which the nonterminal on the left is replaced by the string on the right side of the production.

Ex:- $E \rightarrow E+E \mid E \cdot E \mid (E) \mid -E \mid id \quad \text{--- (2)}$

$E \Rightarrow -E$ (means E derives $-E$)

$E \cdot E \Rightarrow (E) \cdot E$ or $E \cdot (E)$ (Both are correct)

$E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(id)$ (derivation of id from E)

$-(id)$ from E)

- \Rightarrow derives
- $\xrightarrow{*}$ derives in zero or more steps
- $\xrightarrow{+}$ derives in one or more steps
- $L(G) \rightarrow$ language generated by G . Strings in $L(G)$ may contain only terminal symbols of G .
- if $s \xrightarrow{*} \alpha$, where α may contain nonterminals, then we say that α is a sentential form of G .
- if $s \xrightarrow{+} w$, here w is a string of terminals, then w is called a sentence of G .

Ex:- String $-(id + id)$ is a sentence of grammar 2.

$$E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow -(id+id)$$

Means $E \xrightarrow{*} -(id+id)$

One more choice

$$-(E+E) \Rightarrow -(E+id) \Rightarrow -(id+id)$$

- In general leftmost nonterminal in any sentential form is replaced at each step.

such derivations are termed leftmost.

38

$$\begin{aligned} E &\xrightarrow{\text{lm}} -E \xrightarrow{\text{lm}} -(E) \xrightarrow{\text{lm}} -(E+E) \xrightarrow{\text{lm}} -(Id+E) \\ &\xrightarrow{\text{lm}} -(Id+Id) \end{aligned}$$

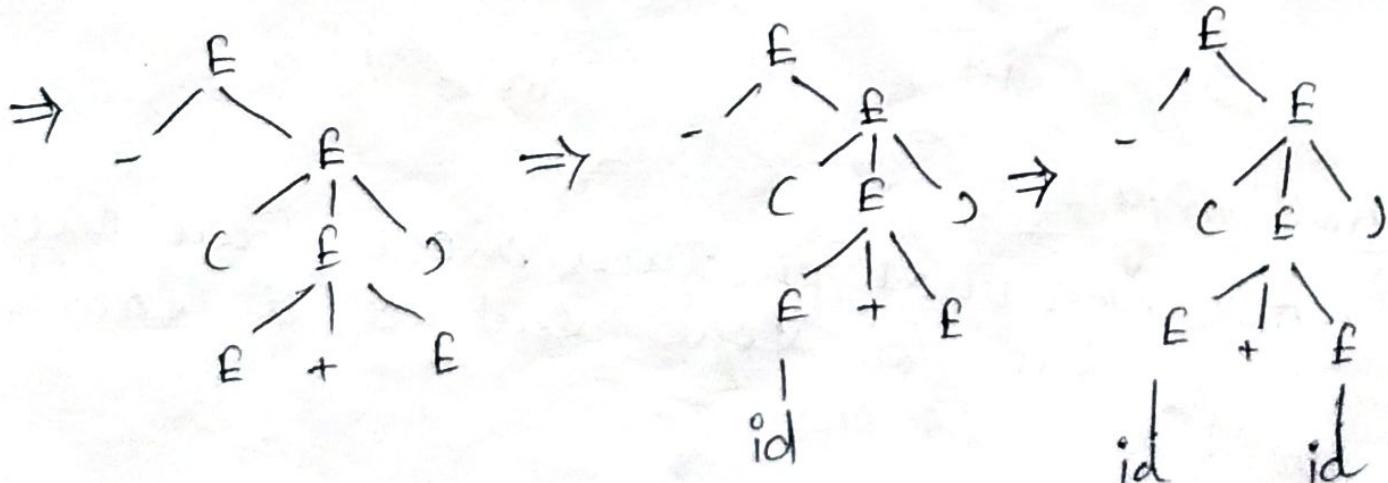
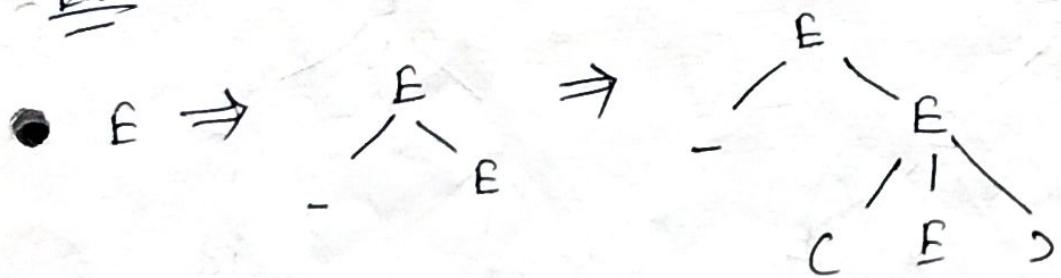
- By which highest nonterminal is replaced
is called canonical derivations.

Parse Trees and Derivation:-

- Parse Trees and Derivation:-

 - A parse tree may be viewed as a graphical representation for a derivation that filters out the choice regarding replacement order. Each interior node has nonterminals while leaves has ~~nonterminals~~ called yield of tree.

Ex:- derive - (id + id)

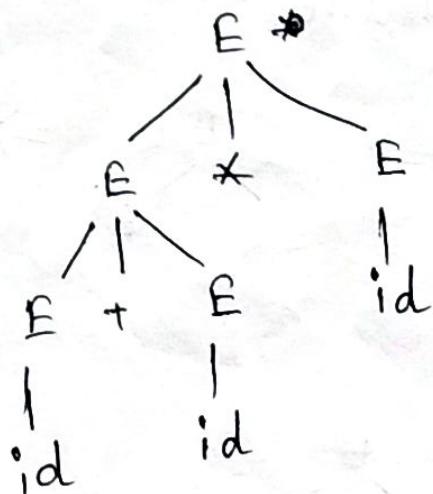
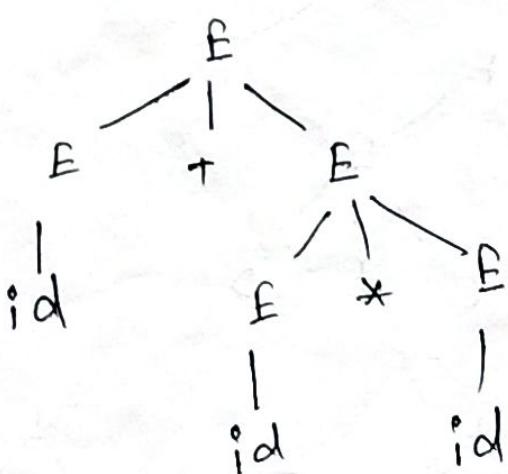


- The variations in the order in which productions are applied can also be eliminated by considering only leftmost derivations.
- Generally every parse tree has associated with a unique leftmost and unique rightmost derivation. sometimes it may be wrong.

Ex $id + id * id$

$$\begin{aligned} ① \quad E &\Rightarrow E+E \\ &\Rightarrow id+E \\ &\Rightarrow id+E*E \\ &\Rightarrow id+id*E \\ &\Rightarrow id+id*id \end{aligned}$$

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



Ambiguity:-

- A grammar that produces more than one parse tree for some sentence is said to be ambiguous.

Writing a grammar:-

- Regular expressions vs. Context free grammar-
 - constructs that can be described by a regular exp can also be described by a grammar.

Ex Regular exp $(a|b)^*abb$
 grammar $A_0 \rightarrow aA_0 | bA_0 | aA_1$
 $A_1 \rightarrow bA_2$
 $A_2 \rightarrow bA_3$
 $A_3 \rightarrow \epsilon$

→ why use regular exp:-

- * lexical rules are simple and do not require a notation as powerful as grammars.
- * Regular exp provide a more concise & easier to understand notation for tokens.
- * More efficient lexical analyzers can be constructed automatically from regular exp

Eliminating Ambiguity:-

Ex

~~stmt → if expr then stmt
 | if expr than stmt
 | other~~

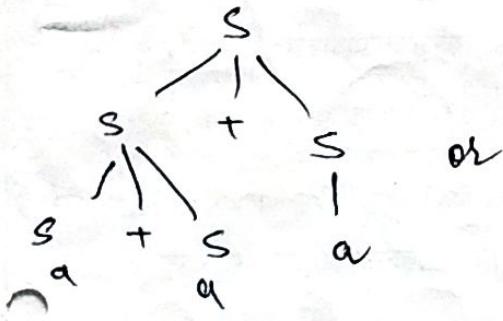
Eliminating Ambiguity:-

Ex:- $E \rightarrow E + E$
 $E \rightarrow E * E$
 $E \rightarrow (E) | id$
 $E \rightarrow \text{#}$

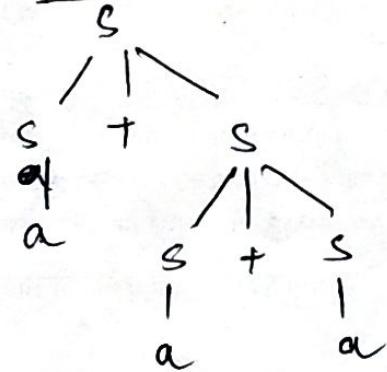
Sol:- $E \rightarrow E + T / T$
 $T \rightarrow T * F / F$
 $F \rightarrow (E) | id$

Ex:- $S \rightarrow S + S / a$

$$w = a + a + a$$



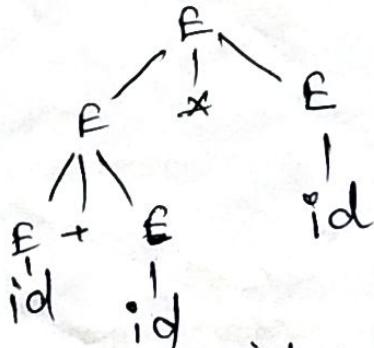
or



In this case we keep highest precedence operator far from start symbol.

]
precedence ambiguity,
solve this using level

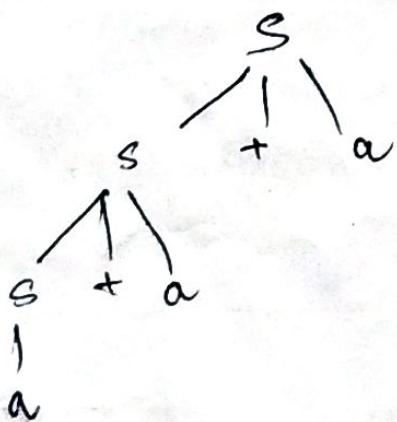
$$w = id + id * id$$



]
Associativity ambiguity,
solve by using left/right recursion

Sol:- $S \rightarrow S + a / a$

$$w = a + a + a$$



(left recursion)

• Elimination of left Recursion:-

→ A grammar is left recursive if it has a nonterminal A such that there is a derivation

$$A \xrightarrow{+} A\alpha.$$

→ Rule to eliminate left recursion

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

then $A \rightarrow \beta A'$

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

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

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

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



$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

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

consider

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

here $A = E$

$$\alpha = +T$$

$$\beta = T$$

• Left factoring :-

→ when it is not clear which of two alternative productions to use to expand a non-terminal A, we may be able to rewrite the A-productions.

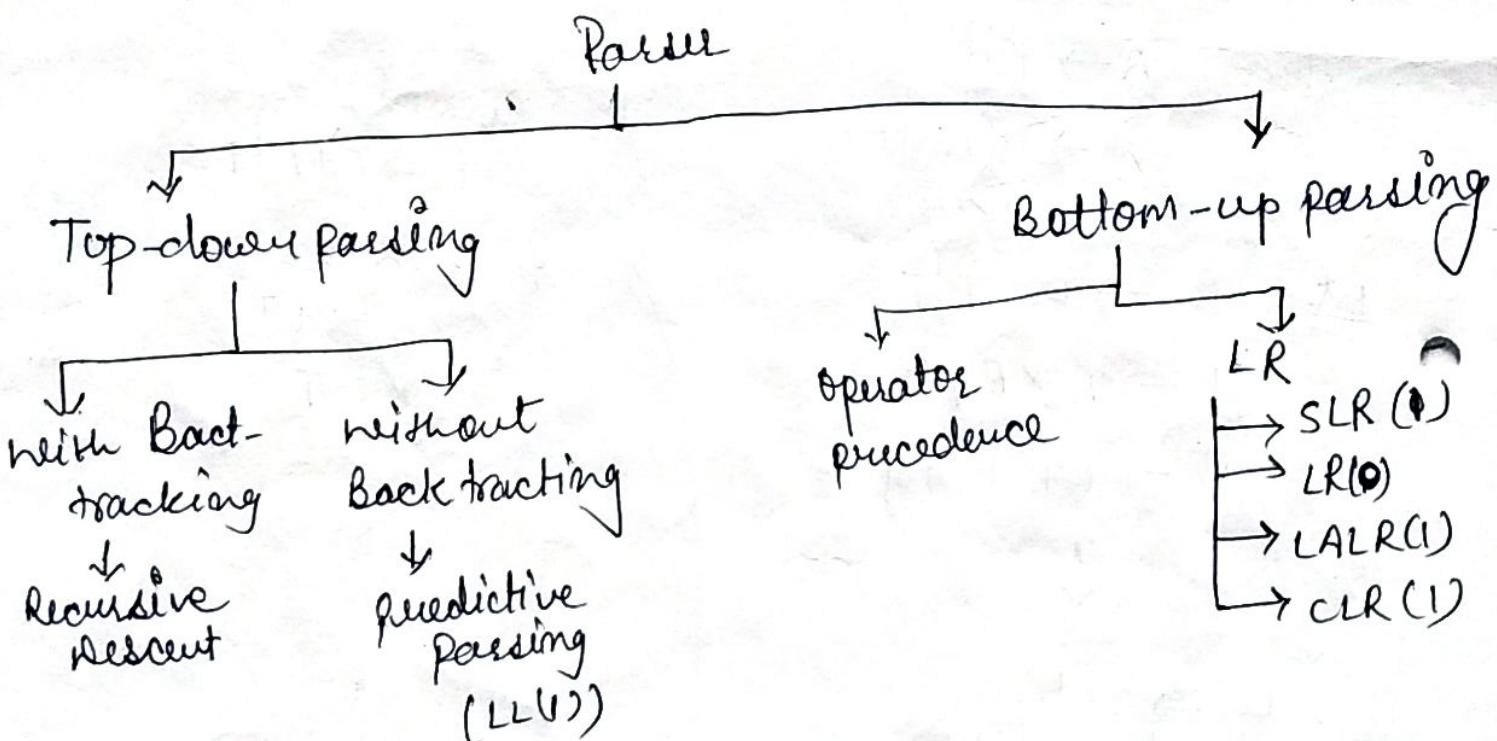
Ex :-

$$\text{if } A \rightarrow \alpha\beta_1 \mid \alpha\beta_2$$

$$\text{then } A \rightarrow \alpha A'$$

$$A' \rightarrow \beta_1 \mid \beta_2$$

~~Top - Down Parsing~~ :-



Recursive-Descent Parsing! -

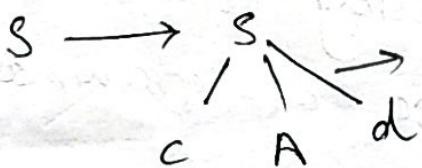
(44)

- Top-down parsing can be viewed as an attempt to find a leftmost derivation for an I/P string.
- This may involve backtracking.

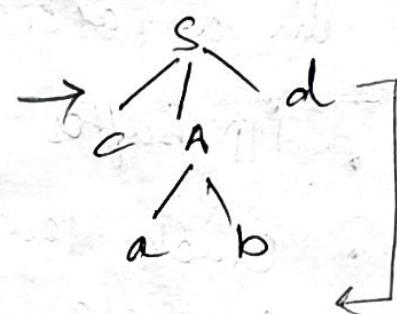
$$\text{Ex: } S \rightarrow cAd \\ A \rightarrow ab \mid a$$

Sol $w = cad$

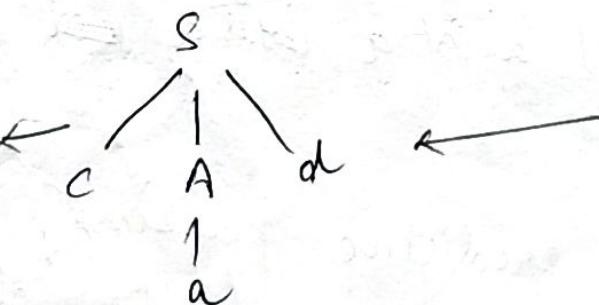
i/p pointer
points to c



c is matched, so
i/p pointer will
points to A,
replace with
first option



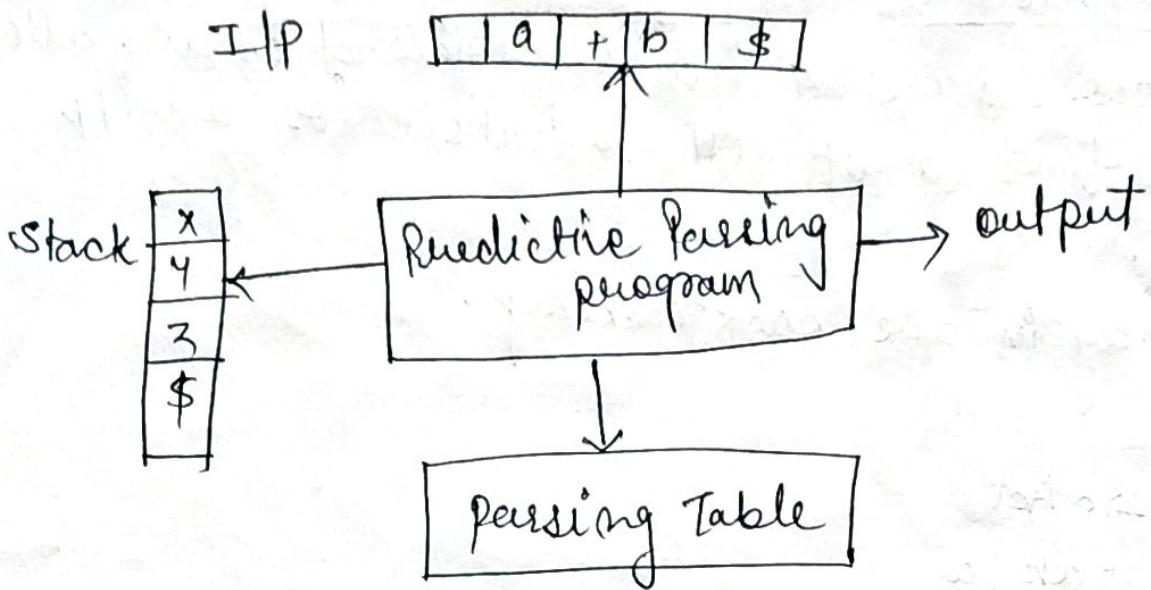
matched,
so halt



a is matched
with second i/p,
so i/p pointer will
points to b, but
b does not matches
with d, so
backtracking

Predictive Parser: -

- For this type of parsing we have to remove the left recursion and left factoring from the grammar.



- parser looks up the production to be applied in a parsing table.
- I/P buffer contains the string to be parsed.
- stack contains a sequence of grammar symbols.
- Parsing table is a 2-D array $M[A, a]$, where A is a terminal and a is a ~~is a~~ nonterminal or \neq symbol.
- Brief steps for predictive parsing.
 - 1) Remove left recursion / left factoring
 - 2) Find first and follow
 - 3) Parsing table
 - 4) Stack implementation
 - 5) Parse tree

~~S~~ $\rightarrow AGB \mid Cbb \mid Ba$

A $\rightarrow da \mid BC$

B $\rightarrow g \mid ? \mid \epsilon$

C $\rightarrow h \mid ? \mid \epsilon$

first(S) = d, g, ?, h, ε, b, a
first(A) = d, g, ?, ε, h
first(B) = g, ?, ε
first(C) = h, ?, ε

Follow(S) = { \$ }
follow(A) = first(C) = h, ?, ε = follow(B)
= d, ?, g, ε = follow(S)
= h, ?, g, \$

follow(B) = first(C) = h, ?, ε = follow(A)
= h, ?, g, \$

follow(C) = b, follow(A)
= b, h, ?, g, \$

Ex :- $S' \rightarrow S$
 $S \rightarrow CC$
 $C \rightarrow cCc \mid d$

$$\begin{aligned} \text{first}(S') &= c, d \\ (\text{S}) &= c, d \\ (\text{C}) &= c, d \end{aligned}$$

$$\begin{aligned} \text{follow}(S') &= \{ \$ \} \\ (\text{S}) &= \{ \$ \} \\ (\text{C}) &= \text{first}(C) = c, d \end{aligned}$$

Ex $S \rightarrow Aa$
 $A \rightarrow BD$
 $B \rightarrow b \mid \epsilon$
 $D \rightarrow d \mid \epsilon$

$$\begin{aligned} \text{first}(S) &= b, \epsilon, d, a \\ (\text{A}) &= b, \epsilon, d \\ (\text{B}) &= b, \epsilon \\ (\text{D}) &= d, \epsilon \end{aligned}$$

$$\begin{aligned} \text{follow}(S) &= \{ \$ \} \\ (\text{A}) &= \{ a \} \\ (\text{B}) &= \text{first}(D) = d, \epsilon = \text{follow}(A) \\ &\quad = d, a \\ (\text{D}) &= \text{Follow}(A) = \{ a \} \end{aligned}$$

Ex consider Grammer:-

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

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

After removing left recursion ambiguity,

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

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

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

After removing left excursion,

$$E \rightarrow T \cdot E'$$

$$E' \rightarrow + T E' \mid \epsilon$$

$$T \rightarrow F \cdot T'$$

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

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

Calculating first and follow,

- first and follow is calculated only for nonterminals.
- for calculating first we see the left non-terminals and follow the path till it generate a terminal production. Now first char of terminal production will be first.

$$\text{first}(E) = \text{first}(\tau) = \text{first}(f) = \{c, id\} \quad (47)$$

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

$$\text{first}(\tau') = \{\ast, \epsilon\}$$

→ For calculating follow we see nonterminals in right side and see which nonterminal follows it. Now ~~follow~~^{first} of following nonterminal will also be follow of followed nonterminal.

$$\text{follow}(E) = \{ \}, \$ \}$$

[By the production
 $f \rightarrow (E)$
and $\$$ bcoz its start state]

$$\text{follow}(E') = \text{follow}(E)$$

[By $E \rightarrow T E'$]

$$= \{ \}, \$ \}$$

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

But ϵ can't come in follow
so follow of ϵ

$$\text{follow}(\tau) = \{ +, \}, \$ \}$$

$$\text{follow}(\tau') = \text{follow}(\tau) = \{ +, \}, \$ \}$$

$$\text{follow}(f) = \text{first}(\tau') = \{ \ast, \epsilon \}$$

$$= \{ \ast, +, \}, \$ \}$$

- Construction of Parsing Table:- Doesn't contain ϵ . 48

Non-terminals	Input Symbols					
	id	+	*	()	\$
E	$E \rightarrow Ee'$			$E \rightarrow Te'$		
E'		$E' \rightarrow Te'$			$E' \rightarrow G$	$E' \rightarrow e$
T	$T \rightarrow PT'$			$T \rightarrow FT'$		
T'		$T' \rightarrow E$	$T' \rightarrow *FT'$		$T' \rightarrow E$	$T' \rightarrow e$
F	$f \rightarrow id$			$f \rightarrow (E)$		

- See first of E, that are {e, id}, entries will be here. but entries that are producing these first symbols
- If a ϵ occurs in first of some non-terminal, then see the follow and put entry in that place with ϵ production.

• Doing entries in stack:- $w = id + id * id$

Stack

~~(\$E) \$ E \$~~
~~\$ E' T \$ T e' \$~~
~~\$ E' T' F \$ F T' E' \$~~
~~\$ E' T' id id \$ T' e' \$~~
~~\$ E' T' T' e' \$~~
~~* P T' E' \$~~
~~\$ e' \$~~
~~\$ E' T \$ T E' \$~~

Input

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

Output

$E \rightarrow Te'$

$T \rightarrow FT'$

$f \rightarrow id$

~~$T' \rightarrow *FT'$~~

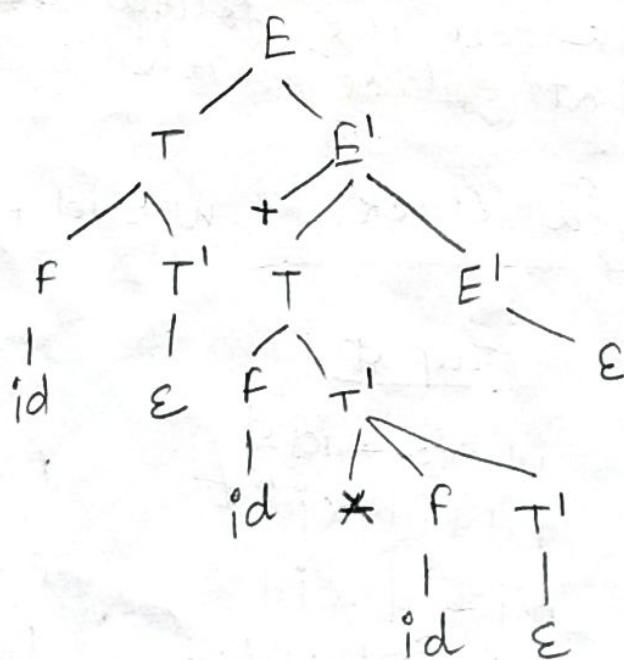
$T' \rightarrow E$

$E' \rightarrow TE'$

(49)

<u>Stack</u>	<u>Input</u>	<u>Output</u>
\$ E' T	id * id \$	
\$ E' T' F	id * id \$	T → f T'
\$ E' T' id	id * id \$	f → id
\$ E' T'	* id \$	
\$ E' T' F	* id \$	T' → * f T'
\$ E' T' F	id \$	
\$ E' T' id	id \$	f → id
\$ E' T'	\$	
\$ E'	\$	T' → ε
\$	\$	E' → ε

Making Parse Tree



Error Recovery in Predictive Parsing:- (LL Parser)

* Panic-mode error recovery:-

- We insert synchronizing tokens in "first" and "follow" of non-terminals.
- if parser looks up entry $M[A, a]$ and finds that it is blank, then i/p symbol 'a' is skipped.
- if the entry is synch, then the non-terminal on top of the stack is popped in an attempt to resume parsing.

NonTerminated	I/P Symbol					
	Id	+	*	'c')	\$
E	$E \rightarrow Ee'$			$E \rightarrow Fe'$	synch	synch
E'		$E' \rightarrow +Te'$			$E \rightarrow E$	$E \rightarrow E$
T	$T \rightarrow PT'$	Synch		$T \rightarrow FT'$	Synch	Synch
T'		$T' \rightarrow E$	$T' \rightarrow *FT'$		$T' \rightarrow E$	$T' \rightarrow E$
F	$F \rightarrow id$	Synch	Synch	$F \rightarrow (E)$	Synch	Synch

Input \rightarrow id * + id \$

Stack	Input	Remark
\$ E) id * + id \$	error, skip (As E is first symbol and we can't remove this)
\$ E	id * + id \$	
\$ E' T	id * + id \$	
\$ E' T' F	id * + id \$	
\$ E' T' id	id * + id \$	
\$ E' T'	* + id \$	
\$ E' T' *	* + id \$	error, M[F, +] = sync
\$ E' T' F	+ id \$	F has been popped
\$ E' T'	+ id \$	
\$ E'	+ id \$	
\$ E' T +	+ id \$	
\$ E' T	id \$	
\$ E' T' F	id \$	
\$ E' T' id	id \$	
\$ E' T'	\$	
\$ E'	\$	
\$	\$	

* Please-level recovery:- it fills table entries of table with pointers to error routines. These routines may change, insert or delete symbols on I/P and issue appropriate error msg.

LL(1) Grammars :-

- A grammar whose parsing table has no multiple-defined entries is said to be LL(1).
- First 'L' stands for scanning the input from left to right, second 'L' for producing a leftmost derivation and the '1' for using one I/P symbol of lookahead at each step to make parsing action decisions.

Bottom-up Parsing :-

- Also known as shift-reduce parsing.
- Attempts to construct a parse tree for an I/P string beginning at the leaves and working up towards the root. We can think of this process as one of "reducing" a string w to the start symbol of a grammar. At each reduction step a particular substring matching the right side of a production is replaced by the symbol on the left of that production.

$$\text{Ex:- } \begin{array}{l} S \rightarrow aABe \\ A \rightarrow Abc/b \\ B \rightarrow d \end{array}$$

$$\begin{aligned} w &= ab\cancel{bcde} \\ &= a\cancel{A}\underline{bcde} \\ &= \underline{aAde} \\ &= S \end{aligned}$$

$$S \xrightarrow{\text{erm}} aABe \xrightarrow{\text{erm}} aA\cancel{bcde} \xrightarrow{\text{erm}} aAbcde \xrightarrow{\text{erm}} abcdde$$

Handles:-

- A handle of a string is a substring that matches the right side of a production, and whose reduction to the nonterminal on the left side of the production represents one step along the reversal of a rightmost derivation.

Ex

$$\begin{aligned} E &\rightarrow E+E \\ E &\rightarrow E \times E \\ E &\rightarrow (E) \\ E &\rightarrow \underline{id} \end{aligned}$$

Stack implementation of shift-reduce

passing:-

- On start, the stack is empty and the string w is on the I/P, as follows,

<u>Stack</u>	<u>Input</u>
\$	$w\ $$

- The parser operates by shifting zero or more I/P symbols onto the stack until a handle β is on top of the stack. The parser then reduces β .

(52)

to the left side of the appropriate production. The parser expects this cycle until it has detected an error or until the stack contains the start symbol and i/p is empty.

Ex: - $w = id_1 + id_2 * id_3$

<u>Stack</u>	<u>Input</u>	<u>Action</u>
1) \$	$id_1 + id_2 * id_3 \$$	shift
2) \$ id_1	$+ id_2 * id_3 \$$	reduce $E \rightarrow id$
3) \$ E	$+ id_2 * id_3 \$$	shift
4) \$ $E +$	$id_2 * id_3 \$$	shift
5) \$ $E + id_2$	$* id_3 \$$	reduce $E \rightarrow id$
6) \$ $E + E$	$id_3 \$$	shift
7) \$ $E + E *$	\$	reduce $E \rightarrow id$
8) \$ $E + E * id_3$	\$	reduce $E \rightarrow E * F$
9) \$ $E + E * E$	\$	reduce $E \rightarrow E + E$
10) \$ $E + E$	\$	accept
11) \$ E		

• There are four possible actions,

(53)

1) Shift: - next I/P symbol is shifted onto the top of stack

2) reduce: - parser replace handle with nonterminal

3) accept: - successful

4) error: - syntax error

Operator - Precedence Parsing:-

operator - precedence parser constructs parse tree

for operator grammar only,

→ no production in right side is ϵ

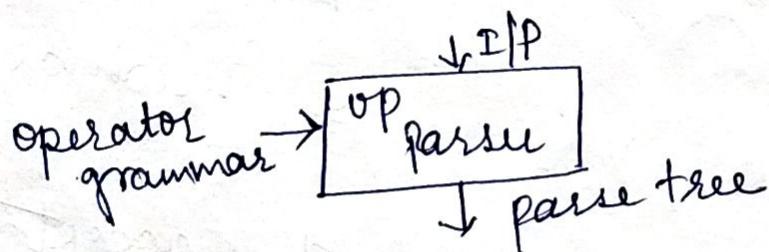
→ no production in right side has two adjacent nonterminals.

Ex:- $E \rightarrow EAE | (E) | -E | id$

$A \rightarrow + | - | * | / | ^$

change in operator grammar,

$E \rightarrow E+E | E-E | E * E | E/E | E^E | (E)-E | id$



Ex:- $T \rightarrow T+T / T*T / id$

$$w = id + id * id$$

- steps -
- 1) Check given grammar is operator precedence grammar or not
 - 2) Operator precedence relation table.
 - 3) Parse the given string
 - 4) Generate parse tree

operator precedence relation table:-

	+	*	id	\$
+	>	<	<	>
*	>	>	<	>
id	>	>	-	>
\$	<	<	<	A

Basics

id, a, b, c → High

\$ → low

+ > +

* > *

id ≠ id

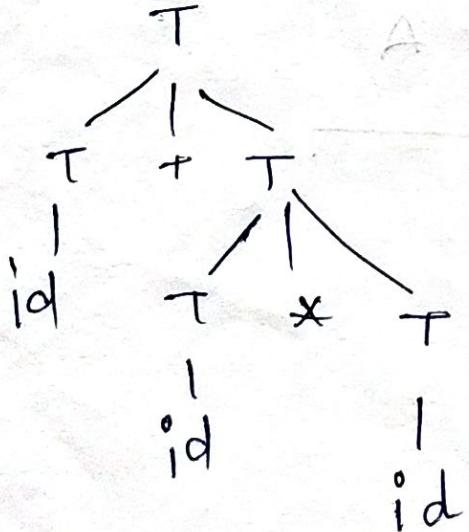
\$ A \$

Passing string

55

<u>Stack</u>	<u>Relation</u>	<u>Input</u>	<u>Comment</u>
\$	<	id + id * id \$	shift id
\$ id	>	+ id * id \$	Reduce id ($T \rightarrow id$)
\$ T	<	+ id * id \$	shift +
\$ T +	<	id * id \$	shift id
\$ T + id	>	* id \$	Reduce $T \rightarrow id$
\$ T + T	<	* id \$	shift *
\$ T + T *	>	id \$	shift id
\$ T + T * id	>	\$	Reduce $T \rightarrow T * T$
\$ T + T * T	>	\$	Reduce $T \rightarrow T + T$
\$ T + T	>	\$	
\$ T	A	\$	

Parse Tree



LR Passes:-

LR(k)

LR(0)

(56)

- 'L' stands for left to right scanning of input.
- 'R' stands for constructing a rightmost derivation of ~~to~~ in reverse
- 'k' stands for ~~constructing~~ no of ~~1/p~~ symbols of lookahead.

Input

$a_1 \dots a_i \dots a_n \$$

Stack

S_m
 x_m
 S_{m-1}
 x_{m-1}
...
 S_0

LR
Passing
Program

Output

action goto

x_i^0 = grammar symbol

s_i^0 = symbol called a state

Each state symbol summarizes the information contained in the stack below it.

SLR Passes :- SLR(1)

E1 :- $E \rightarrow E + T \mid T$

$T \rightarrow T * F \mid F$

$F \rightarrow (e) \mid id$

Step1 : Augmented expression grammar

$E \rightarrow E$

$E \rightarrow E + T \mid T$

$T \rightarrow T * F \mid F$

$F \rightarrow (E) \mid id$

I₀: E' → E

E → E' + T — ①

E → · T — ②

T → · T * F — ③

T → · F — ④

F → · (E) — ⑤

F → · id — ⑥

I₆: E → E + T

T → · T * F

T → · F

F → · (E)

F → · id

I₇: T → T * F

F → · (E)

F → · id

I₁: E' → E ·

E → E · + T

I₂: E → T ·

T → T · * F

I₃: T → F ·

I₈: F → (E ·)

E → E · + T

I₉: E → E + T ·

T → T · * F

I₁₀: T → T * F ·

I₄: F → (· E)

E → · E + T

E → · T

T → · T * F

T → · F

F → · (E)

F → · id

I₁₁: F → (E) ·

first(E) = { (, id}

first(T) = { (, id}

first(F) = { (, id}

follow(E) = { +,), \$ }

follow(T) = { *, +,), \$ }

follow(F) = { *, +,), \$ }

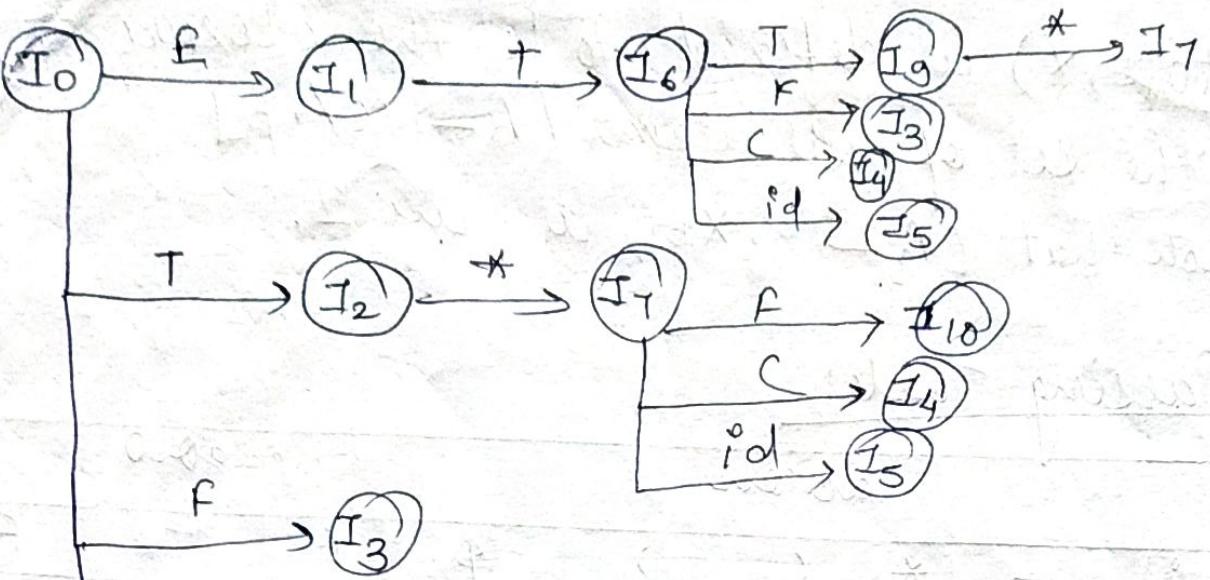
I₅: F → id

After removing left recursion

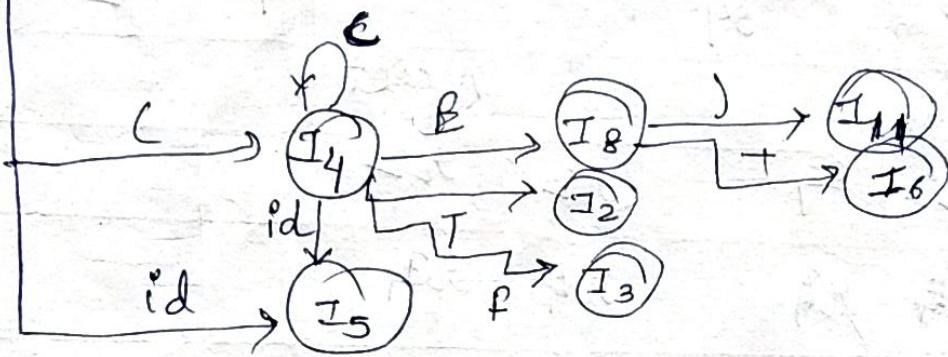
~~E → TE' E' → TE' | ε~~

~~T → FT' T' → *FT' | ε~~

~~F → (E) | id~~



(58)



- The closure operation
 → if I is a set of items for a grammar G ,
 then $\text{closure}(I)$ is the set of items
 constructed from I by two rules:
 a) Initially, every item in I is added
 to $\text{closure}(I)$.
 b) if $A \rightarrow \alpha \cdot B\beta$ is in $\text{closure}(I)$ and
 $B \rightarrow \gamma$ is a production, then add the
 item $B \rightarrow \cdot\gamma$ to I , if it is not already
 there.
- The goto operation: - $\text{goto}(I, x)$ where I is a
 set of items and x is a grammar symbol.

69) $\text{goto}(I, x)$ is defined to be the closure of the set of all items $[A \rightarrow \alpha X \cdot \beta]$ such that $[A \rightarrow \alpha \cdot X \beta]$ is in I .

Make Parsing Table:-

State	action							goto		
	id	+	*	()	\$	E	T	F	
0	S5			S4			1	2	3	
1		S6				Accept				
2		r2	S7		r2	r2				
3		r4	r4		r4	r4				
4	S5			S4			8	2	3	
5		r6	r6		r6	r6				
6	S5			S4			9	3		
7	S5			S4						10
8		S6			S11					
9		r1	S7		r1	r1				
10		r3	r3		r3	r3				
11		r5	r5		r5	r5				

- Check which symbol is starting from a state and going to which state.
- For finding reduce find such state which are final (having a . in the end) states. Then

(60)

check the follow of symbol just before '•'.
 Now place reduce action in such places
 in the parsing table with production
 numbers from 10.

Stack Implementation:-

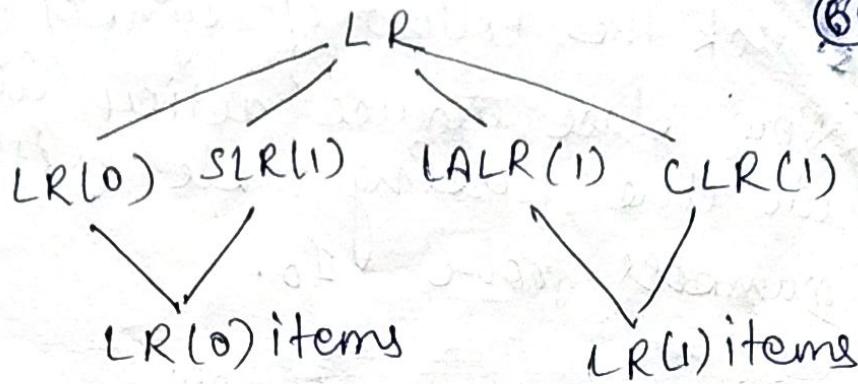
$$w = \text{id} * \text{id} + \text{id}.$$

Stack	Input	Action
1) 0	$\text{id} * \text{id} + \text{id} \$$	shift
2) 0 id 5	$* \text{id} + \text{id} \$$	reduce $F \rightarrow \text{id}$
3) 0 F 3	$* \text{id} + \text{id} \$$	reduce $T \rightarrow F$
4) 0 T 2	$* \text{id} + \text{id} \$$	shift
5) 0 T 2 * 7	$\text{id} + \text{id} \$$	shift
6) 0 T 2 * 7 id 5	$+ \text{id} \$$	reduce $F \rightarrow \text{id}$
7) 0 T 2 * 7 F 10	$+ \text{id} \$$	reduce $T \rightarrow T * F$
8) 0 T 2	$+ \text{id} \$$	reduce $E \rightarrow T$
9) 0 E 1	$+ \text{id} \$$	shift
10) 0 B 1 + 6	$\text{id} \$$	shift
11) 0 E 1 + 6 id 5	$\$$	reduce $F \rightarrow \text{id}$
12) 0 B 1 + 6 F 3	$\$$	reduce $T \rightarrow F$
13) 0 E 1 + 6 T 9	$\$$	reduce $E \rightarrow E + T$
14) 0 E 1	$\$$	Accept

- In line 1, 0 and id will be processed and we look at the parsing table, it is shift to state 5. So shift 'id' to stack with state no 5 at top.
- In line 2, B to * is reduce to production 5 which is $F \rightarrow \text{id}$, now put F into state and then see 0 to F, that is 3. So stack is 0F3.

LR(0) :-

Example :-



Steps :-

- 1) Augment the given grammar
- 2) Draw canonical collection of LR(0) items / DFD
- 3) Number the production
- 4) Create the Parsing table
- 5) Stack implementation
- 6) Draw parse tree

Example :-

$$E \rightarrow BB$$

$$B \rightarrow cB/d$$

Step 1 :- augment grammar

$$E' \rightarrow E$$

$$E \rightarrow BB$$

$$B \rightarrow cB/d$$

Step 2: Draw conical collection

(62)

$$I_0 : E^1 \rightarrow E$$

$$E \rightarrow \cdot BB$$

$$B \rightarrow \cdot cB$$

$$B \rightarrow \cdot d$$

$$I_5 : E \rightarrow BB.$$

$$I_6 : B \rightarrow cB.$$

$$I_1 : E^1 \rightarrow E.$$

$$I_2 : E \rightarrow B.B$$

$$B \rightarrow \cdot cB$$

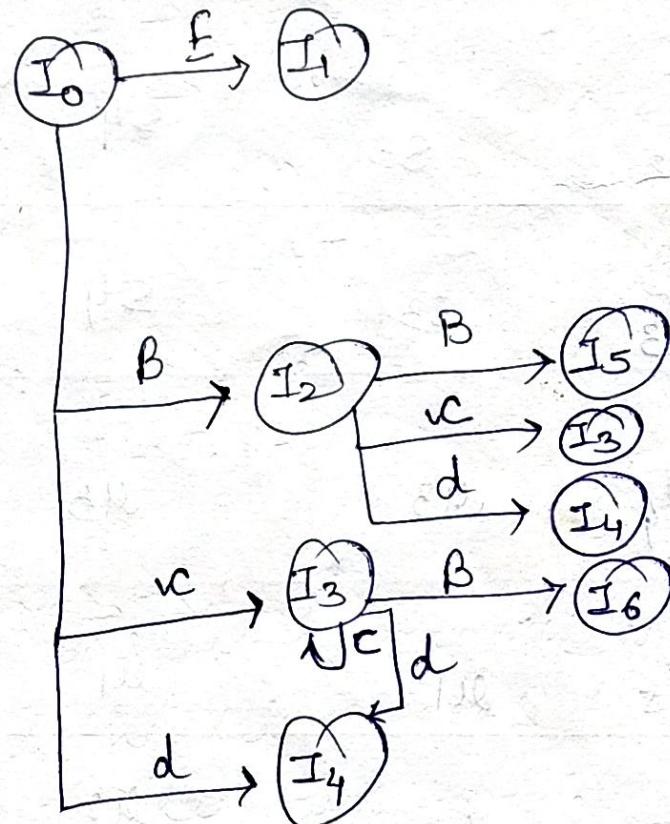
$$B \rightarrow \cdot d$$

$$I_3 : B \rightarrow c.B$$

$$B \rightarrow \cdot cB$$

$$B \rightarrow \cdot d$$

$$I_4 : B \rightarrow d.$$



Step 3: No other productions

$$E^1 \rightarrow E$$

$$E \rightarrow BB \rightarrow \textcircled{1}$$

$$B \rightarrow cB \rightarrow \textcircled{2}$$

$$B \rightarrow d \rightarrow \textcircled{3}$$

State	Action				goto
	c	d	\$	E	
I ₀	s ₃	s ₄		1	2
I ₁			Accept		
I ₂	s ₃	s ₄			5
I ₃	s ₃	s ₄			6
I ₄	m ₃	m ₃	m ₃		
I ₅	m ₁	m ₁	m ₁		
I ₆	m ₂	m ₂	m ₂		

As we have (o) lookahead symbols, so we can't find follow, as it requires next symbol. So conflict can occur. ~~To remove conflicts~~ To remove conflicts we use SLR(1) parse with one lookahead symbol.

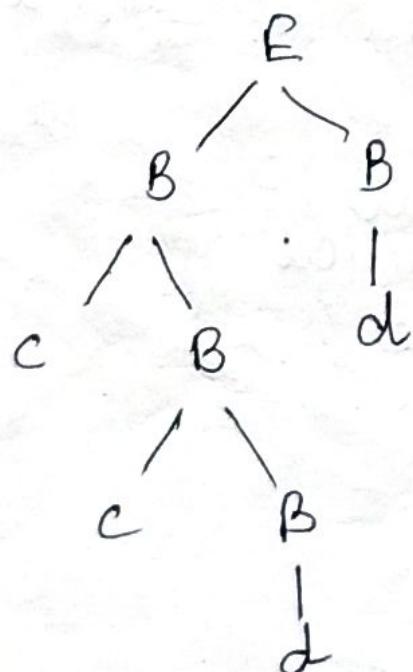
Step 5 : Start implementation

(64)

string = ccdd

0	cc dd \$	shift
0 C3	cc dd \$	shift
0 C3C3	dd \$	shift
0 C3C3d4	d \$	reduce B → d
0 C3C3B6	d \$	reduce B → cB
0 C3B6	d \$	reduce B → cB
0 B2	d \$	shift
0 B2d4	\$	reduce B → cB d
0 B2B5	\$	reduce E → BB
0 E1	\$	Accept

Step 6: Draw parse tree



CLR(1) :- Canonical LR(1)

(65)

- Based on LR(1) items

Final lookahead:-

$$\underline{\text{Ex:}} \quad E \rightarrow BB \\ B \rightarrow cB/d$$

$E \rightarrow \cdot B, \$$ (as start symbol do not has any lookahead)

$$E \rightarrow \cdot BB, \$$$

~~for B~~

$$B \rightarrow \cdot cB/d, c/d$$

$$\underline{\text{Ex:}} \quad A \rightarrow \alpha \cdot B \beta ; a'$$

Suppose we are now opening B, so ($B \rightarrow L$) we take lookahead as first ($B\alpha$)

• suppose first of B is c

$$\text{then } B \rightarrow \cdot L, c$$

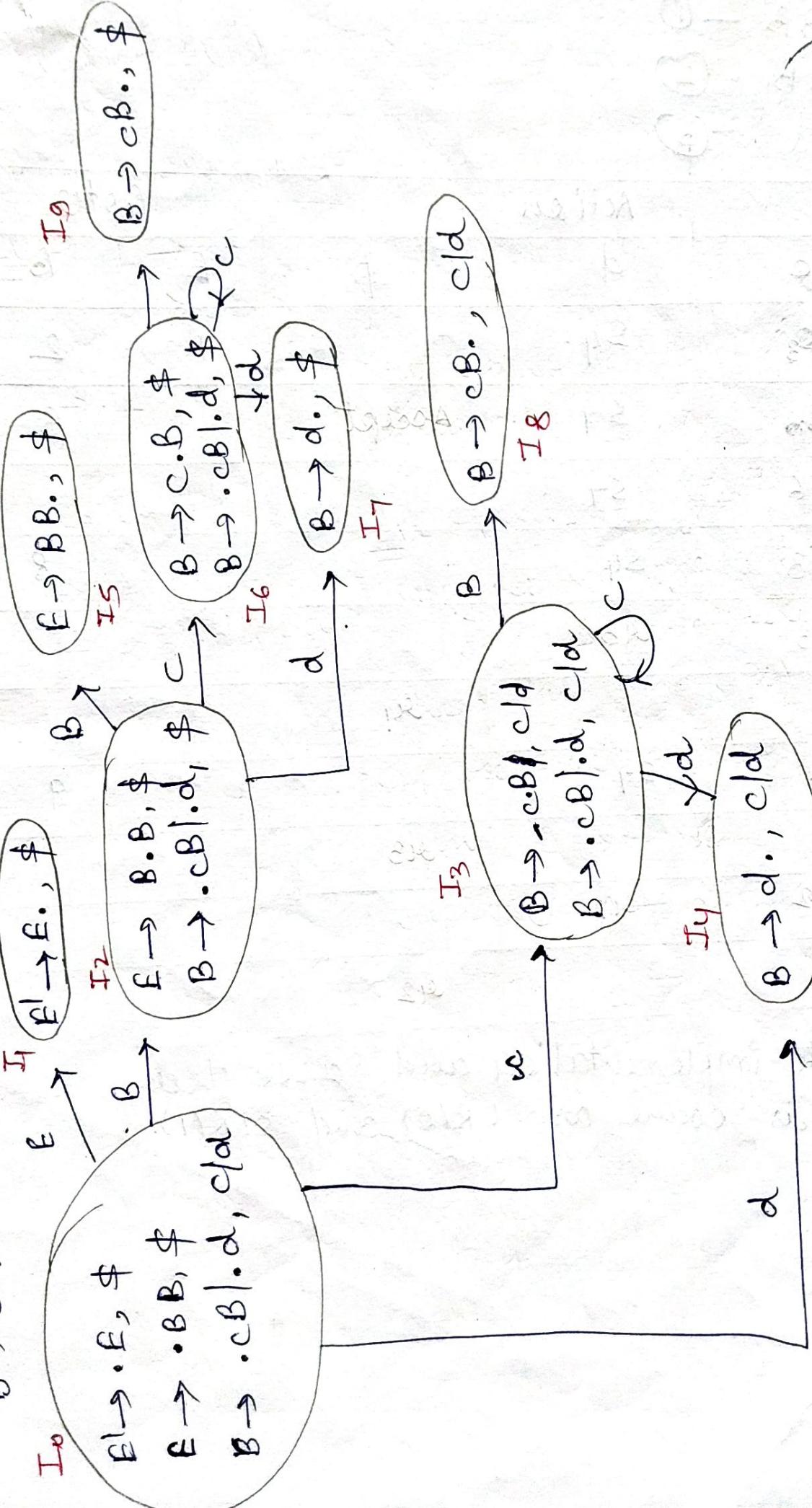
• suppose $A \rightarrow \alpha \cdot B, a$

$$\text{then } B \rightarrow \cdot L, a$$

Steps: Augment grammar
 Draw canonical collection
 Number the production
 Create passing table
 Stack implementation
 Draw fail table

where we do get these look-ahead changes when we open some symbols.

$$\text{Ex: } E \rightarrow BB \\ B \rightarrow cB|d$$



(67)

$$E' \rightarrow E$$

$$E \rightarrow BB \rightarrow ①$$

$$B \rightarrow cB \rightarrow ②$$

$$B \rightarrow d \rightarrow ③$$

State	Action			Goto	
	c	d	\$	E	B
I ₀	s ₃	s ₄		t	2
I ₁			Accept		5
I ₂	s ₆	s ₇			
I ₃	s ₃	s ₄			8
I ₄	a ₃	a ₃			
I ₅			a ₄		
I ₆	s ₆	s ₇			9
I ₇			a ₃		
I ₈	a ₂	a ₂			
I ₉			a ₂		

Now stack implementation and parse tree creation is same as LR(0) and SLR(1).

LALR(1) :- Lookahead LR(1)

(68)

- using LR(1) items
- combine those states which are having same production but different lookahead.

$$I_8 + I_6 = I_{36}$$

$$I_4 + I_7 = I_{47}$$

$$I_8 + I_9 = I_{89}$$

state	Action			goto	
	C	d	\$	E	B
I ₀	S ₃₆	S ₄₇		1	2
I ₁			Accept		
I ₂	S ₃₆	S ₄₇			5
I ₃₆	S ₃₆	S ₄₇			89
I ₄₇	M ₃	M ₃	M ₃		
I ₅			M ₄		
I ₈₉	M ₂	M ₂	M ₂		

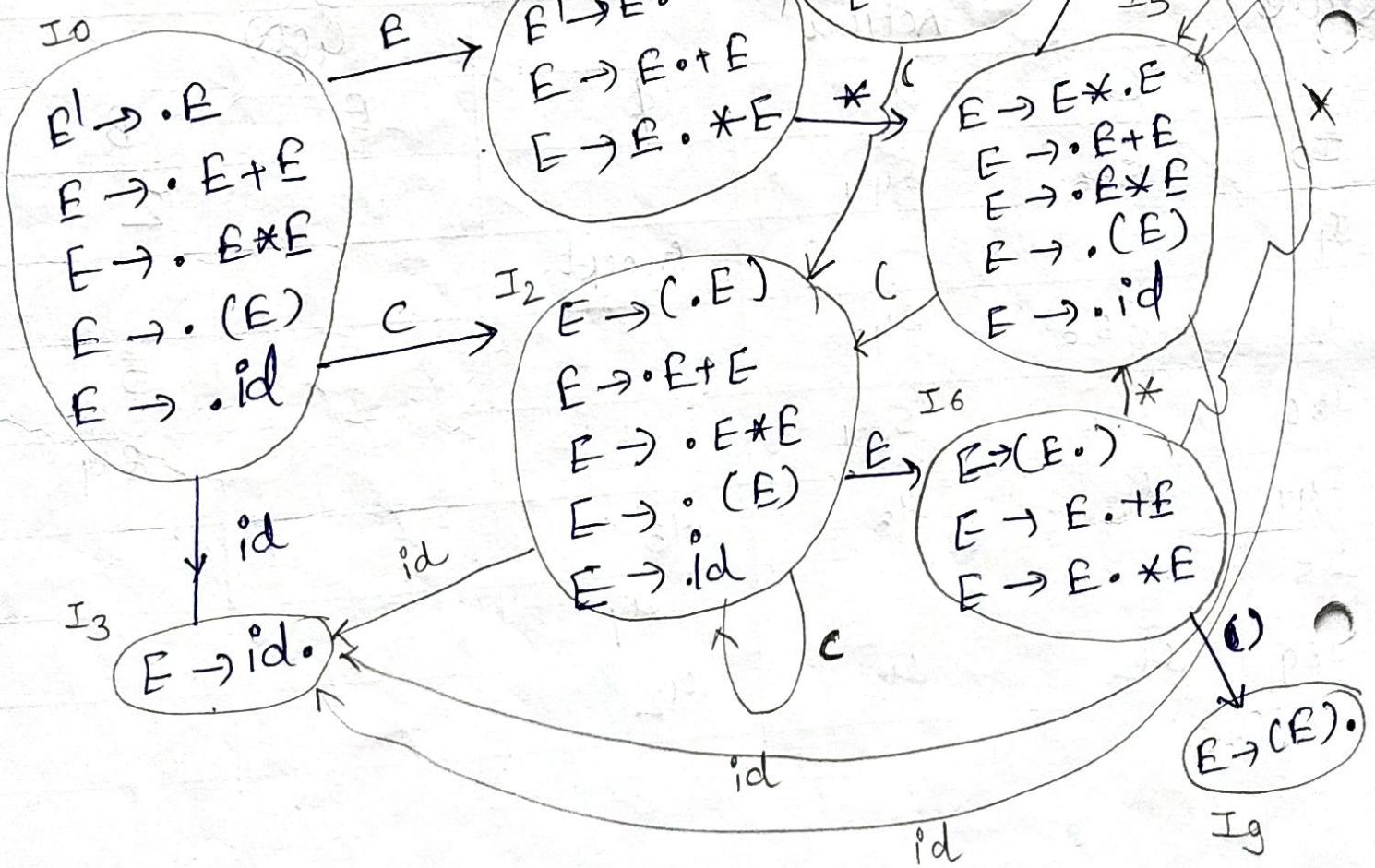
Using Ambiguous Grammars:-

(6.g)

- This type of grammar uses precedence and associativity to resolve parsing action conflicts.

Ex:- $E \rightarrow E + E \quad ①$
 $E \rightarrow E * E \quad ②$
 $E \rightarrow (E) \quad id \quad ④$
 $E \rightarrow .id \quad ③$

LR(0) items



State	Action							Goto
	id	+	*	()	\$	E	
I ₀	S ₃			S ₂				1
I ₁		S ₄	S ₅				Accept	
I ₂	S ₃			S ₂				6
I ₃		a ₄	a ₄			a ₄	a ₄	
I ₄	S ₃			S ₂				7
I ₅	S ₃			S ₂				8
I ₆		S ₄	S ₅	a₄	S ₉			
I ₇		a ₄ /S ₄	S ₅ /a ₁			a ₄	a ₄	
I ₈		S ₄ /a ₂	S ₅ /a ₂			a ₂	a ₂	
I ₉		a ₃	a ₃			a ₃	a ₃	

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

$$\text{follow}(E) = \{ +, *,), \$ \}$$

state	Input	Action
0	id + id * id \$	shift
0 id 3	+ id * id \$	reduce $\epsilon \rightarrow id$
0 E 1	+ id * id \$	shift
0 E 1 + 4	+ id * id \$	shift
0 E 1 + 4 id 3	* id \$	reduce $\epsilon \rightarrow id$
0 E 1 + 4 E 7	* id \$	
Now as precedence of * is higher on +, so shift * onto stack, otherwise reduce the stack.		

Syntax Directed Definitions:-

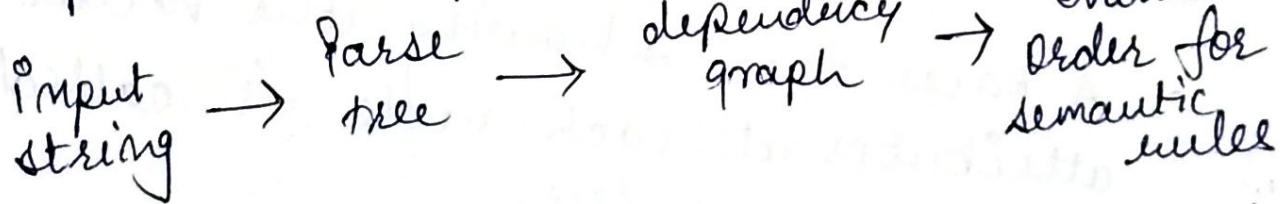
- In this chapter, we are going to attach attributes to the grammar symbols. Values for attributes are computed by "semantic rules" associated with the grammar productions.
- Notations for semantic rules

syntax-directed
definitions → translation
scheme.

(High level specifications
for translation schemes
and hide many implementa-
tion details)

(indicate the order
in which semantic
rules are to be
evaluated)

- conceptual view of syntax-directed translation,



Syntax - directed definitions:-

- A syntax-directed definition is a generalization of a context-free grammar in which each grammar symbol has an associated set of attributes, called synthesized and inherited attributes.

- (72)
- An attribute can represent anything a string, number, type or memory location.
 - Synthesized attribute :- computed from the values of attributes at children
 - Inherited attribute :- computed from the values of attributes at the siblings and parent of that node.
 - A dependency graph is computed from attributes, from this graph an evaluation order for semantic rules is derived. Evaluation of these rules defines the values of attributes at the node in parse tree for I/P string.
 - A parse tree showing the values of attributes at each node is called an annotated parse tree.

form of a syntax-Directed Definition:-

- In a syntax-directed definition, each grammar production $A \rightarrow \alpha$ has associated with it a set of semantic rules of the form $b := f(g_1, g_2, \dots, g_k)$ where f is a function, and either,

→ b is a synthesized attribute of A (73) and $a_1, a_2 \dots a_k$ are attributes belonging to the grammar symbols of the production, or

→ b is an inherited attribute of one of the grammar symbols on the right side of the production and $a_1, a_2 \dots a_k$ are attributes belonging to the grammar symbols of the production.

Ex:-

Production

dummy attri $L \rightarrow E_n$ newline char
for printing $E \rightarrow E_1 + T$
 $E \rightarrow T$
 $T \rightarrow T_1 * F$
 $T \rightarrow F$
 $F \rightarrow (E)$
 $F \rightarrow \text{digit}$

Semantic rules

print(E.val)

$E.val := E_1.val + T.val$

$E.val := T.val$

$T.val := T_1.val * F.val$

$T.val := F.val$

$f.val := E.val$

$f.val := \text{digit}.\text{lexical}$

Value is assumed
to be supplied
by lexical analyzer

- Terminals are assumed to have synthesized attributes only.

Synthesized attributes:-

- called S-attributed definition. A parse tree for an S-attributed definition can always

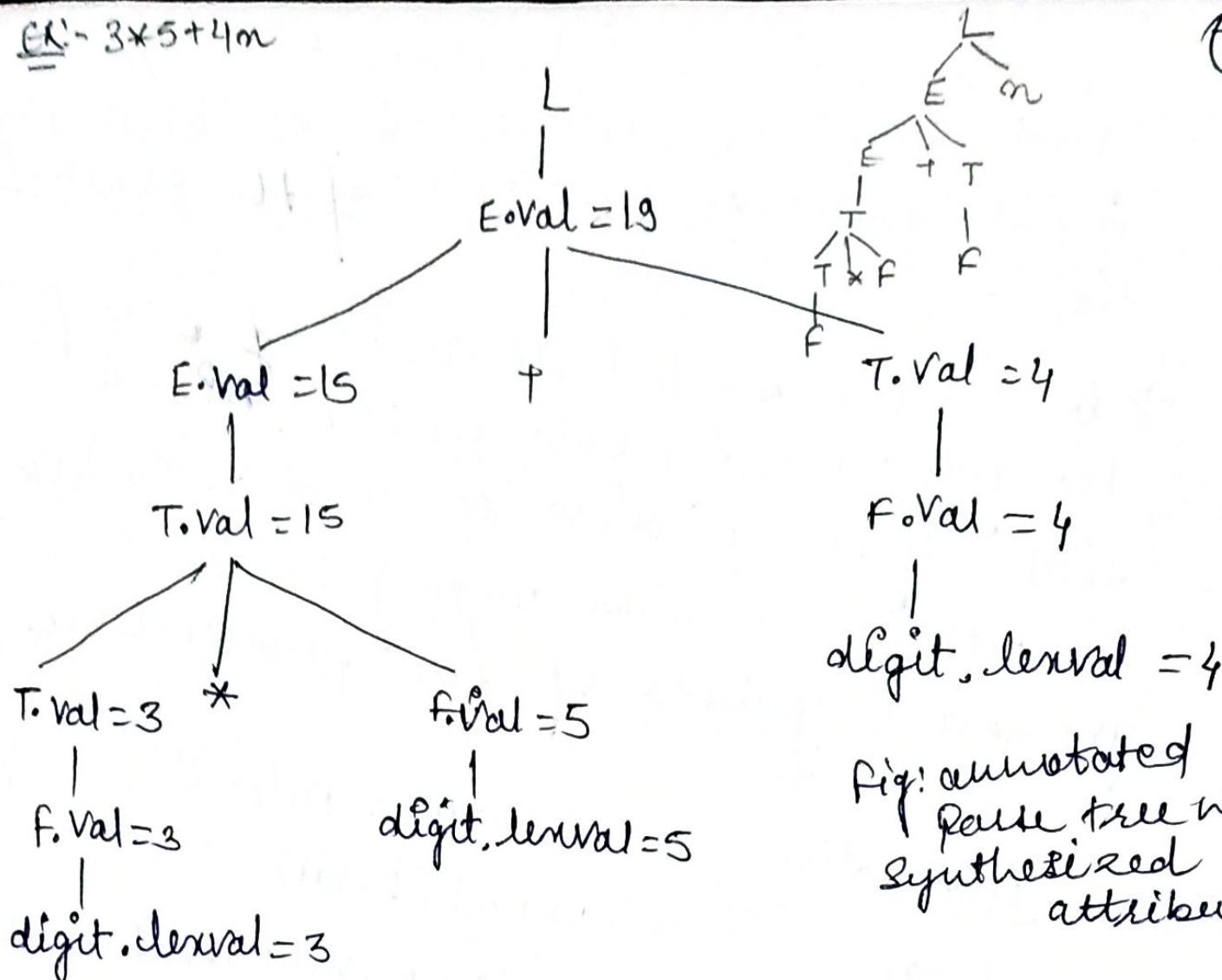


Fig: annotated
parse tree with
synthesized
attributes

be annotated by evaluating the semantic rules for the attributes at each node bottom up, from the leaves to the root.

Inherited attributes:-

- Inherited attributes are convenient for expressing the dependence of a programming language construct on the context in which it appears. For ex, we can use it to keep track whether an identifier appears on the left or right side of an assignment in order to decide whether the address or the value of identifier is needed.

Ex:-

Production

$$D \rightarrow TL$$

Synthetic attribute
T → int

$$L \rightarrow L_1, id$$

$$L \rightarrow id$$

Semantic Rules

(15)

L.iM := T.type

T.type := integer

T.type := real

L.iM := L.iM

↓ addtype(id.entry, L.iM)

addtype(id.entry, L.iM)

add the type of each identifier
to its entry in the symbol
table.

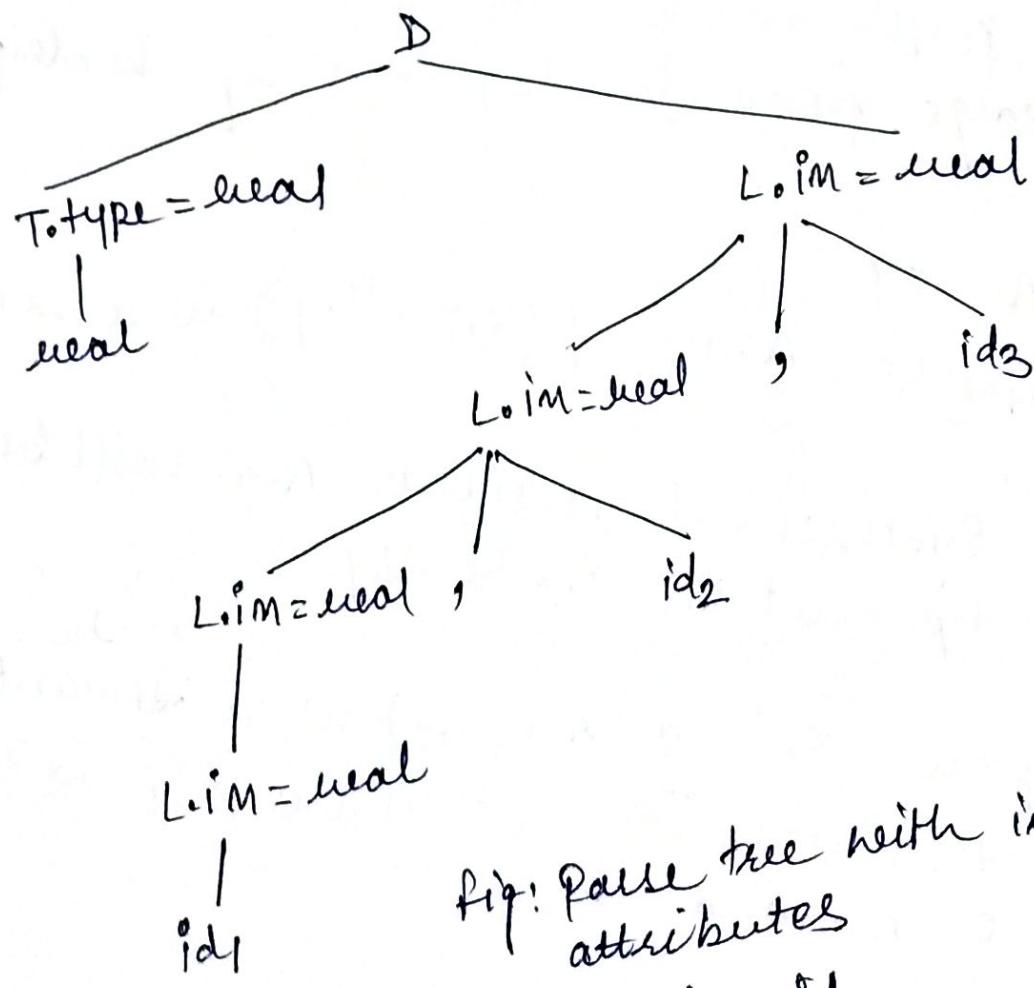


fig: Parse tree with inherited
attributes

Ex:- id1, id2, id3

Dependency Graphs:-

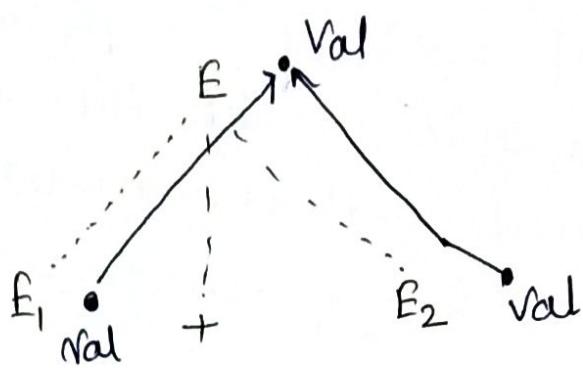
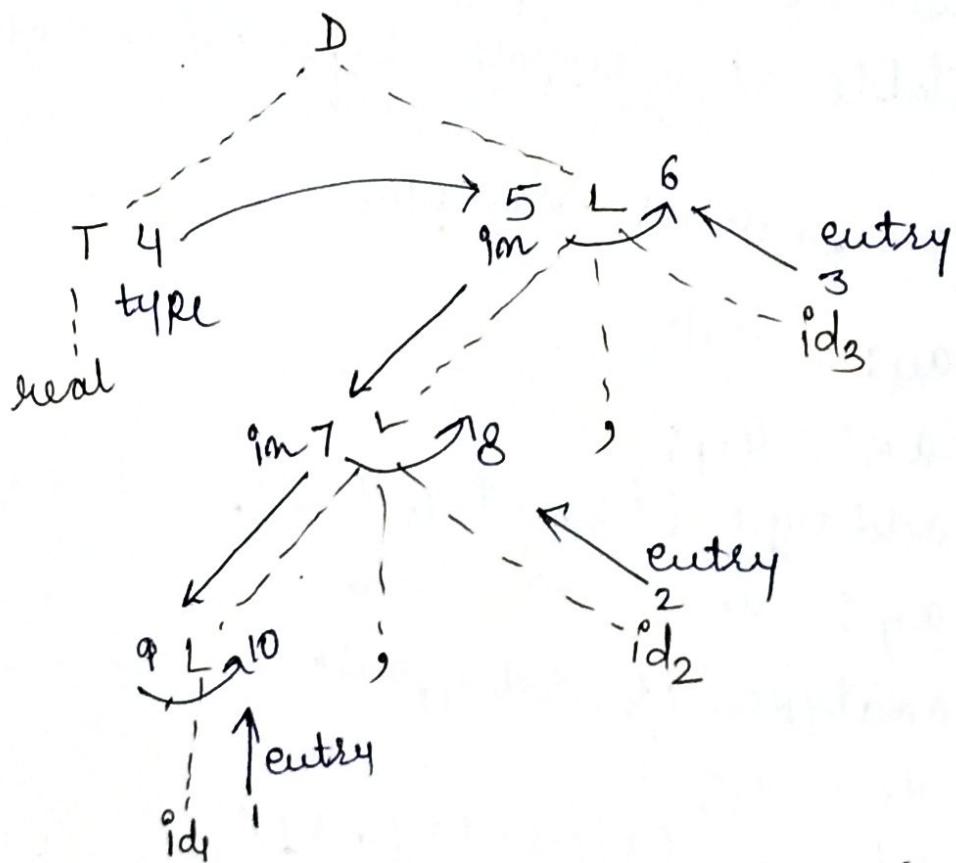
- the interdependencies among the inherited and synthesized attributes at the nodes in a parse tree can be depicted by a directed graph called a dependency graph.
- Before constructing a dependency graph for a parse tree, we put each semantic rule into the form $b := f(c_1, c_2, \dots, c_k)$, by introducing a dummy synthesized attribute b for each semantic rule that consists of a procedure call.
- The graph has a node for each attribute and an edge ~~from~~ for b from c if b depends on c .

Ex:- $A \rightarrow x.y$

Suppose $A.a := f(x.x, Y.y)$ is a semantic rule.
Synthesized attribute $A.a$ will be dependent on $x.x$ & $Y.y$.

Suppose $x.i := g(A.a, Y.y)$ is a semantic rule
then inherited attribute $x.i$ will depend on $A.a$ and $Y.y$

(77)

Ex:- $E \rightarrow E_1 + E_2$ semantic rule $f.\text{Val} := E_1.\text{Val} + E_2.\text{Val}$ Ex:- $\text{id}_1, \text{id}_2, \text{id}_3$ 

- Each of the semantic rules { add type (Id.entry, L.qm) associated with the L-production leads to the creation of a dummy attribute. Nodes 6, 8, 10 are constructed for these dummy attributes.

Evaluation Order:-

- Topological sort is used for this.
- Any topological sort of a dependency graph gives a valid order in which the semantic rules associated with the nodes in a parse tree can be evaluated. That is the dependent attributes $c_1, c_2 \dots c_k$ in a semantic rule $b := f(c_1, c_2 \dots c_k)$ are available at a node before f is evaluated.

Ex:- from previous example,

$a_4 := \text{real};$

$a_5 := a_4;$

addtype ($\text{id}_3.\text{entry}$, a_5);

$a_7 := a_5;$

addtype ($\text{id}_2.\text{entry}$, a_7);

$a_9 = a_7;$

addtype ($\text{id}_1.\text{entry}$, a_9);

Methods for evaluating semantic rules:-

- 1) Parse-tree methods:- These methods obtain an evaluation order from a topological sort of the dependency graph constructed from the parse tree for each input. This method will fail, if dependency graph contain cycle.

2) Rule-based methods:- At compilation time, the semantic rules associated with productions are analyzed, either by hand, or by a special tool. (79)

3) oblivious methods:- an evaluation order is chosen without considering the semantic rules.

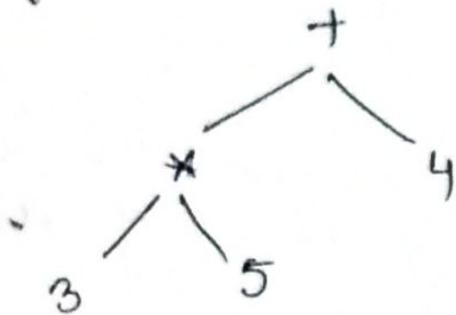
■ construction of syntax trees:-

- The use of syntax trees as an intermediate representation allows translation to be decoupled from parsing.
- Translation routines that are invoked during parsing have two restrictions:
 - 1) grammar that is suitable for parsing may not reflect natural hierarchical structure of the constructs in the language.
 - 2) Parsing method contains the order in which nodes in a parse tree are considered. This order may not match the order in which info about a construct becomes available.

Syntax Tree:-

- A syntax tree is a condensed form of parse tree useful for representing language constructs.

Ex:- Syntax tree for * the tree given (20)
iii) Synthesized attribute is,

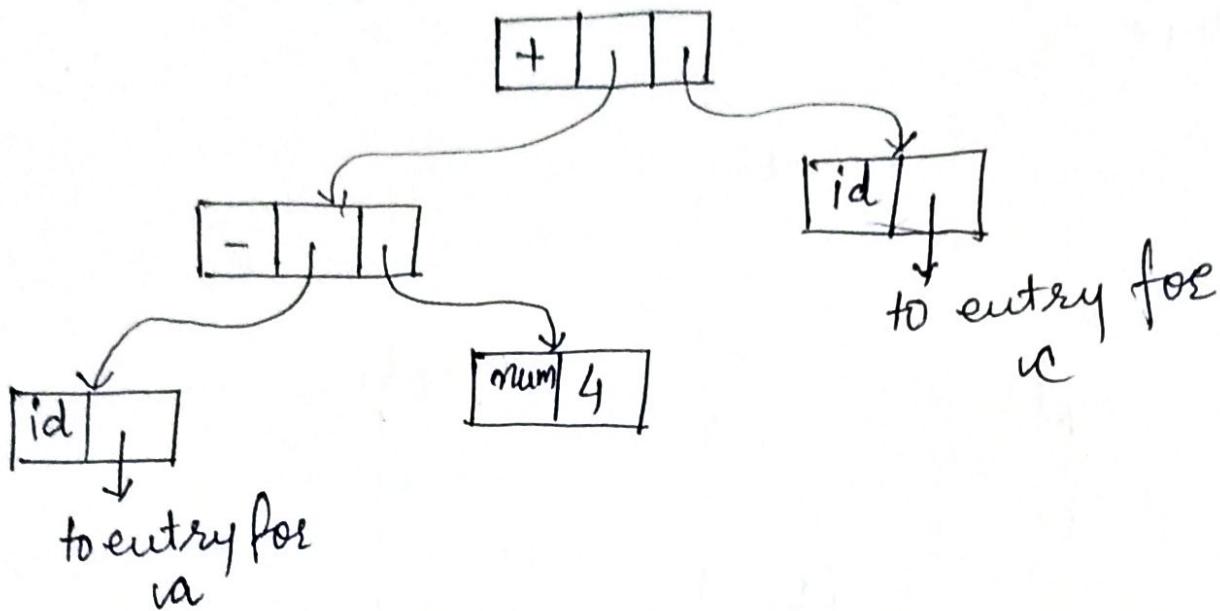


Constructing syntax trees for expressions.
• Similar to translation of expression into postfix form.

- a) `mknode(op, left, right)`: creates an operator node with label op and two fields containing pointers to left and right.
- b) `mkleaf(id, entry)`: identifies node with label id, and a field containing a pointer to symbol-table entry for identifier.
- c) `mkleaf(num, val)`: number node with label num and a field containing value.

Ex:- exp: - $a - 4 + c$

- 1) $P_1 := \text{mkleaf}(\text{id}, \text{entry}_a);$
- 2) $P_2 := \text{mkleaf}(\text{num}, 4);$
- 3) $P_3 := \text{mknode}('-', P_1, P_2);$
- 4) $P_4 := \text{mkleaf}(\text{id}, \text{entry}_c);$
- 5) $P_5 := \text{mknode}('+', P_3, P_4);$



A Syntax-directed definition for constructing
Syntax trees:- (mptr = node pointer)

Production

$$E \rightarrow E_1 + T$$

$$E \rightarrow E_1 - T$$

$$E \rightarrow T$$

$$T \rightarrow (B)$$

$$T \rightarrow id$$

$$T \rightarrow num$$

Semantic Rules

$E.mptr := \text{mknode}('+' , E_1.mptr, T.mptr)$

$E.mptr := \text{mknode}('-', E_1.mptr, T.mptr)$

$E.mptr = T.mptr$

$T.mptr = E.mptr$

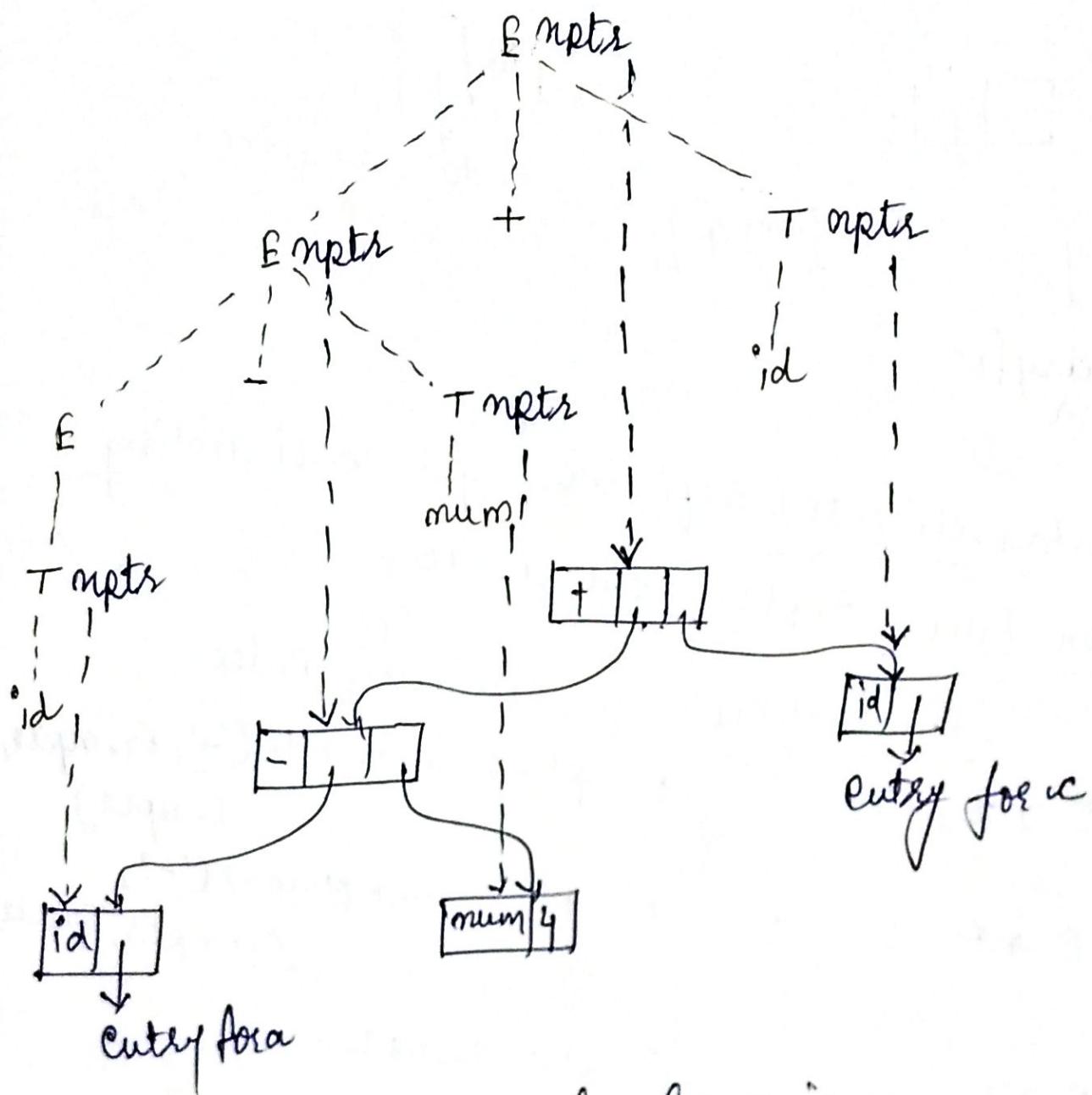
$T.mptr := \text{mkleaf}(id, id.entry)$

$T.mptr := \text{mkleaf}(num, num.val)$

Fig:- S-attributed definition

Ex:- a-4+c

(82)

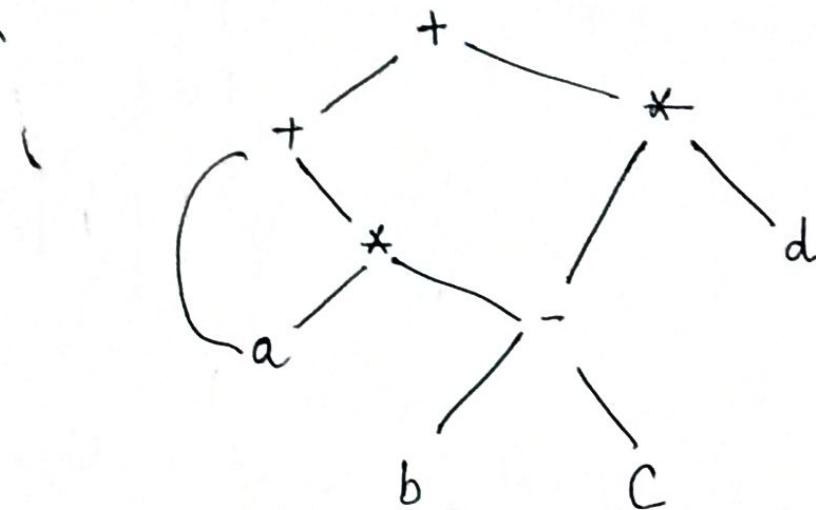


Directed Acyclic Graphs for Expressions:-

- A dag for an expression identifies the common subexpressions in the expression. A dag has a node for every subexpression; an interior node represents an operator and its children represent its operands.
- A node in a dag representing a common subexpression has more than one parent.

Ex:-

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



- 1) $P_1 := \text{mkleaf}(\text{id}, a);$
- 2) $P_2 := \text{mkleaf}(\text{id}, a);$
- 3) $P_3 := \text{mkleaf}(\text{id}, b);$
- 4) $P_4 := \text{mkleaf}(\text{id}, c);$
- 5) $P_5 := \text{mknode}(' - ', P_3, P_4);$
- 6) $P_6 := \text{mknode}(' * ', P_2, P_5);$
- 7) $P_7 := \text{mknode}(' + ', P_1, P_6);$

- 8) $P_8 := \text{mkleaf}(\text{id}, b);$
- 9) $P_9 := \text{mkleaf}(\text{id}, c);$
- 10) $P_{10} := \text{mknode}(' - ', P_8, P_9);$
- 11) $P_{11} := \text{mkleaf}(\text{id}, d);$
- 12) $P_{12} := \text{mknode}(' * ', P_{10}, P_{11});$
- 13) $P_{13} := \text{mknode}(' + ', P_7, P_{12});$

- When the call $\text{mkleaf}(\text{id}, a)$ is repeated in line 2, the node constructed by the previous call $\text{mkleaf}(\text{id}, a)$ is returned, so $P_1 = P_2$
- lines 8 & 9 are same as 3 & 4
- line 10 is equal to line 5.
- In many applications dag nodes are implemented as records also.

Ex:-

Assignment

 $i := i + 10$ 

Representation

1	id		→ entry for i
2	num	10	
3	+	1 2	
4	:=	1 3	
5			

value - number method for constructing a node in a dag:-

Input :- label op, node l, node r

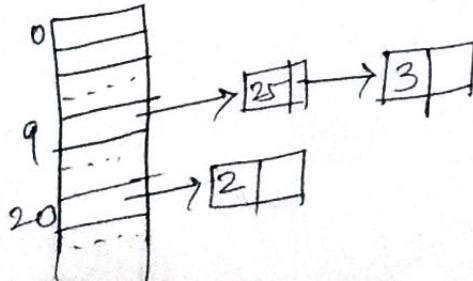
Output:- A node with signature $\langle op, l, r \rangle$

Method:- Search the array for same node m and return it. If same node does not exist then create a new one.

- The search for m can be made more efficient by using k lists, called buckets and using a hashing function h to determine which bucket to search.

- Hash function h computes the number of a bucket from the value of op, l and r.

If m is not in the bucket ~~l, r~~, then a new node m is $\langle h(op, l, r), l, r \rangle$, created and added to this bucket.



Bottom-up Evaluation of S-attributed definitions :-

(85)

→ we have already seen how to use SDD to specify translations, now we will see how to implement translators for them.

→ In this section, we will study only a class of SDD: the S-attributed definitions, means SDD with only synthesized attributes.

● Synthesized attributes can be evaluated by a bottom-up parser as I/P is being passed. The parser can keep the values of synthesized attributes associated with the grammar symbols on its stack. whenever a reduction is made, the values of the new synthesized attributes are computed from the attributes appearing on the stack for the grammar symbols on the right side of the reducing production.

Synthesized attributes on the parser stack:-

• Implemented using LR-parser generator

• A bottom-up parser uses a stack to hold info about subtrees that have been passed. we can use extra fields in the parser stack to hold the values of

Synthesized attributes.

(86)

state	val
-	-
X	X.Z
Y	Y.Y
Z	Z.Z
-	-

top →

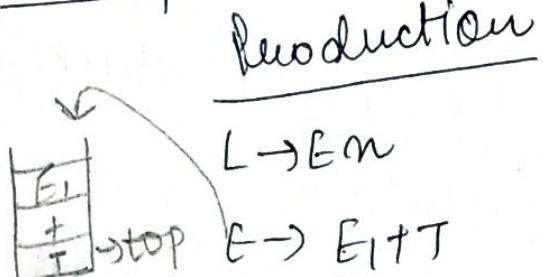
- if the i th state symbol is A, then $\text{val}[i]$ will hold the value of attribute associated with the parse tree node corresponding to this A.

Suppose $A \rightarrow X Y Z$

then semantic rule $A.a = f(X.x, Y.y, Z.z)$

Before $X Y Z$ is reduced to A, the value of the attribute $Z.z$ is in $\text{val}[\text{top}]$, $Y.y$ in $\text{val}[\text{top}-1]$ & $X.x$ in $\text{val}[\text{top}-2]$. After reduction, top is decremented by 2 and A is put in state [top].

Example:-



$E \rightarrow T$

$T \rightarrow T, * F$

$T \rightarrow F$

$F \rightarrow (E)$

$F \rightarrow \text{digit}$

Code fragment

print ($\text{val}[\text{top}]$)

$\text{val}[\text{ntop}] := \text{val}[\text{top}-2] + \text{val}[\text{top}]$

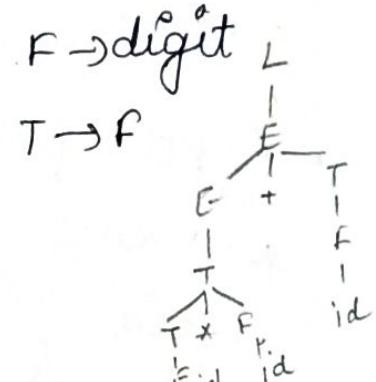
$\text{val}[\text{ntop}] := \text{val}[\text{top}-2] * \text{val}[\text{top}]$

$\text{val}[\text{ntop}] := \text{val}[\text{top}-1]$

- When a production with n symbols on the right side is reduced, the value of mtop is set to top - n + 1. After each code fragment is executed, top is set to ctop.

Ex:- $3 * 5 + 4m$

<u>Input</u>	<u>state</u>	<u>val</u>	<u>Production Used</u>
$3 * 5 + 4m$	-	-	
$* 5 + 4m$	3	3	$F \rightarrow \text{digit}$
$* 5 + 4m$	F	3	$T \rightarrow F$
$* 5 + 4m$	T	3	
$5 + 4m$	$T *$	3 -	
$+ 4m$	$T * 5$	$3 - 5$	$F \rightarrow \text{digit}$
$+ 4m$	$T * F$	$3 - 5$	
$+ 4m$	$T *$	15	$T \rightarrow T * F$
$+ 4m$	E	15	$E \rightarrow T$
$4m$	$E +$	15 -	
m	$E + 4$	$15 - 4$	$F \rightarrow \text{digit}$
m	$E + F$	$15 - 4$	
m	$E + T$	$15 - 4$	$T \rightarrow F$
m	E	19	$E \rightarrow E + T$
$E m$		19 -	
L		19	$L \rightarrow Em$



■ L-Attributed definitions:-

(88)

- when translation takes place during parsing, the order of evaluation of attributes is linked to order in which nodes of a parse tree are created by the parsing method. If visit is a method for many top to down and bottom to up translation.

Depth-first evaluation order:-

procedure dvisit (n : node);

begin for each child m of n , from left to right do begin
evaluate inherited attributes of m ;
dvisit (m)

end; evaluate synthesized attributes of n

end

L-attributed definitions:

- A syntax-directed definition is L-attributed if each inherited attribute of x_j , $1 \leq j \leq n$, on the right side of $A \rightarrow x_1 x_2 \dots x_n$, depends only on
 - the attributes of the symbols x_1, x_2, \dots ,

x_{j-1} , to the left ~~to~~ of x_j in the production and

(89)

e. the Inherited attributes of A.

Ex:- $A \rightarrow LM$

Here L's inherited attributes will depend on A and

A & L.

M's

Translation Schemes:-

→ A translation scheme is a context-free grammar in which attributes are associated with grammatical symbols and semantic actions enclosed between braces {} are inserted within the right side of production.

→ For every production $A \rightarrow x_1 \dots x_n$, the semantic rule formed the string for A by concatenating the strings for $x_1 \dots x_n$. Now we could perform translation by simply printing the literal strings in the order they appeared in the semantic rules.

Ex:- A simple translation scheme that maps infix exp into postfix exp.

$E \rightarrow TR$

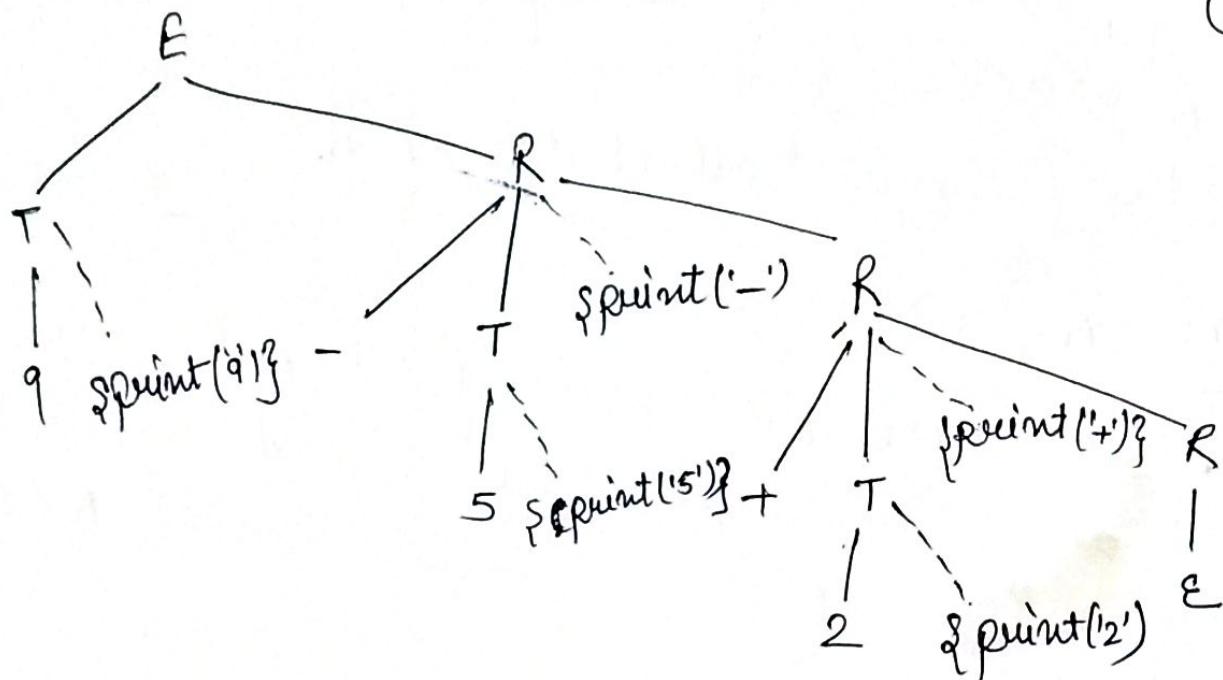
$R \rightarrow \text{addop } T \quad \{ \text{print (addop. lexeme) } \} R | E$

$T \rightarrow \text{num } \quad \{ \text{print (num. Val) } \}$

$E \rightarrow TR$

$R \rightarrow +TR / -TR / E$

$T \rightarrow \text{num}$



Input $\rightarrow 9 - 5 + 2$ (infix)
 output $\rightarrow 95-2+$ (postfix)

- When designing a translation scheme, we must ensure that an attribute value must available when an action refers to it.
- The easiest case occur when only synthesized attributes are needed. We place each semantic rule at the end of the right side of the associated production.

<u>Rule</u>	<u>Production</u>	<u>Semantic Rule</u>
	$T \rightarrow T_1 * F$	$T.\text{Val} := T_1.\text{Val} \times F.\text{Val}$

Production & semantic rule

$$T \rightarrow T_1 * F \quad \{ \quad T.\text{Val} := T_1.\text{Val} \times F.\text{Val} \quad \}$$

- If we have both inherited & synthesized attributes, we must be careful:

- 1.> An inherited attribute for a symbol on the ~~(1)~~ right side of a production must be computed in an action before that symbol
- 2.> An action must not refer to a synthesized attribute of a symbol to the right * of the action.
- 3.> A synthesized attribute for the nonterminal on the left can only be computed after all attributes it references have been computed. The action computing such attributes can usually be placed at the end of the right side of the production.

Ex1 L-attributed SDD for which following alone three rules,

Production

$$S \rightarrow B$$

$$B \rightarrow B_1, B_2$$

$$B \rightarrow \text{tent}$$

Semantic Rules

$$B.\text{ps} := 10 \quad \text{inherited att. (point size)}$$

$$S.\text{ht} := B.\text{ht} \quad \text{synthesized, height}$$

$$B_1.\text{ps} := B.\text{ps}$$

$$B_2.\text{ps} := B.\text{ps}$$

$$B.\text{ht} := \max(B_1.\text{ht}, B_2.\text{ht})$$

$$B.\text{ht} := \text{tent.h} \times B.\text{ps}$$

Fig: - SDD for size and height of boxes.

$$S \rightarrow \{B.ps := 10\}$$

$$B \{ s.ht := B.ht \}$$

$$B \rightarrow \{B_1.ps = B.ps\}$$

$$B_1 \{ B_2.ps = B.ps \}$$

$$B_2 \{ B.ht := \max(B_1.ht, B_2.ht) \}$$

$$B \rightarrow \text{text} \{ B.ht := \text{text}.h \times B.ps \}$$

Top-Down Translation :-

→ Here we implement 2-attributed definitions during predictive parsing. Here we work with translation schemes rather than SDD.

Eliminating left recursion from a translation scheme:-

Here we apply transformations to translation schemes with synthesized attributes.

$$\begin{array}{ll} \text{Ex:-} & E \rightarrow E_1 + T \\ & E \rightarrow E_1 - T \\ & E \rightarrow T \\ & T \rightarrow (E) \\ & T \rightarrow \text{num} \end{array}$$

$$\{E.Val := E_1.Val + T.Val\}$$

$$\{E.Val := E_1.Val - T.Val\}$$

$$\{E.Val := T.Val\}$$

$$\{T.Val := E.Val\}$$

$$\{T.Val := \text{num}.Val\}$$

Fig:- Translation scheme with left recursive grammar (synthesized attributes)

After Removing left Recursion,

(98)

$$E \rightarrow T R$$
$$R \rightarrow + TR | - TR | \epsilon$$

~~E, T, R~~

$$T \rightarrow (E)$$
$$T \rightarrow \text{num}$$

Fig:- Transformed translation scheme with right-recursive grammar.

$$E \rightarrow T \quad \left\{ \begin{array}{l} R.i := T.\text{val} \\ E.\text{val} := R.S \end{array} \right\}$$

$$R \rightarrow + \quad \left\{ \begin{array}{l} R_1.i := R.i + T.\text{val} \\ R_1.S := R.i S \end{array} \right\}$$

$$R \rightarrow - \quad \left\{ \begin{array}{l} R_1.i := R.i - T.\text{val} \\ R_1.S := R.i S \end{array} \right\}$$

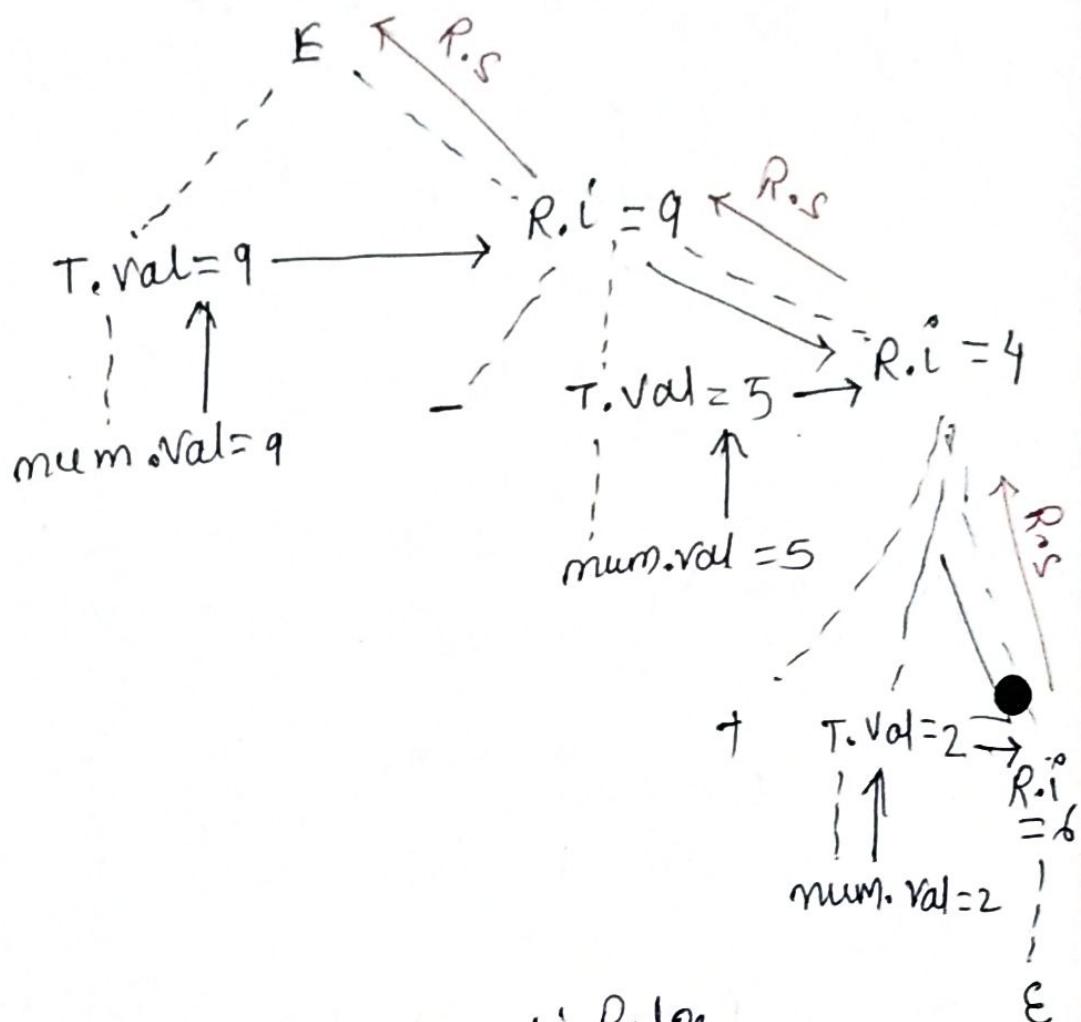
$$R \rightarrow \epsilon \quad \left\{ \begin{array}{l} R.S := R.i^2 \end{array} \right\}$$

$$T \rightarrow (\quad \left\{ \begin{array}{l} E \\ T \end{array} \right\} \quad \left\{ \begin{array}{l} T.\text{val} := E.\text{val} \end{array} \right\}$$

$$T \rightarrow \text{num} \quad \left\{ \begin{array}{l} T.\text{val} := \text{num}. \text{val} \end{array} \right\}$$

Fig: Evaluation of exp $9 - 5 + 2$

(94)



Semantic Rules

$E.\text{mpte} = \text{mknode}('+', E_1.\text{mpte}, T_1.\text{mpte})$

$E.\text{mpte} = \text{mknode}(' - ', E_1.\text{mpte}, T_1.\text{mpte})$

$E.\text{mpte} = T.\text{mpte}$

$T.\text{mpte} = E.\text{mpte}$

$T.\text{mpte} = \text{mkleaf}(\text{id}, \text{id}. \text{entry})$

$T.\text{mpte} = \text{mkleaf}(\text{num}, \text{num}. \text{val})$

Ex:- Production

$$E \rightarrow E_1 + T$$

$$E \rightarrow E_1 - T$$

$$E \rightarrow T$$

$$T \rightarrow (E)$$

$$T \rightarrow \text{id}$$

$$T \rightarrow \text{num}$$

After removing left recursion,

$$E \rightarrow E \cdot TR$$

$$R \rightarrow +TR \quad | -TR/E$$

$$T \rightarrow (E)$$

$$T \rightarrow id$$

$$T \rightarrow num$$

Fig:- Transformed translation scheme for
constructing syntax trees.

$$\begin{array}{l} E \rightarrow T \\ R \end{array} \quad \left\{ \begin{array}{l} R.i = T.mptr \\ E.mptr = R.S \end{array} \right\}$$

$$\begin{array}{l} R \rightarrow +T \\ R_1 \end{array} \quad \left\{ \begin{array}{l} R_1.i = mknode (+, R.i, T.mptr) \\ R_1.S = R.i \cdot S \end{array} \right\}$$

$$\begin{array}{l} R \rightarrow -T \\ R_1 \end{array} \quad \left\{ \begin{array}{l} R_1.i = mknode (-, R.i, T.mptr) \\ R_1.S = R.i \cdot S \end{array} \right\}$$

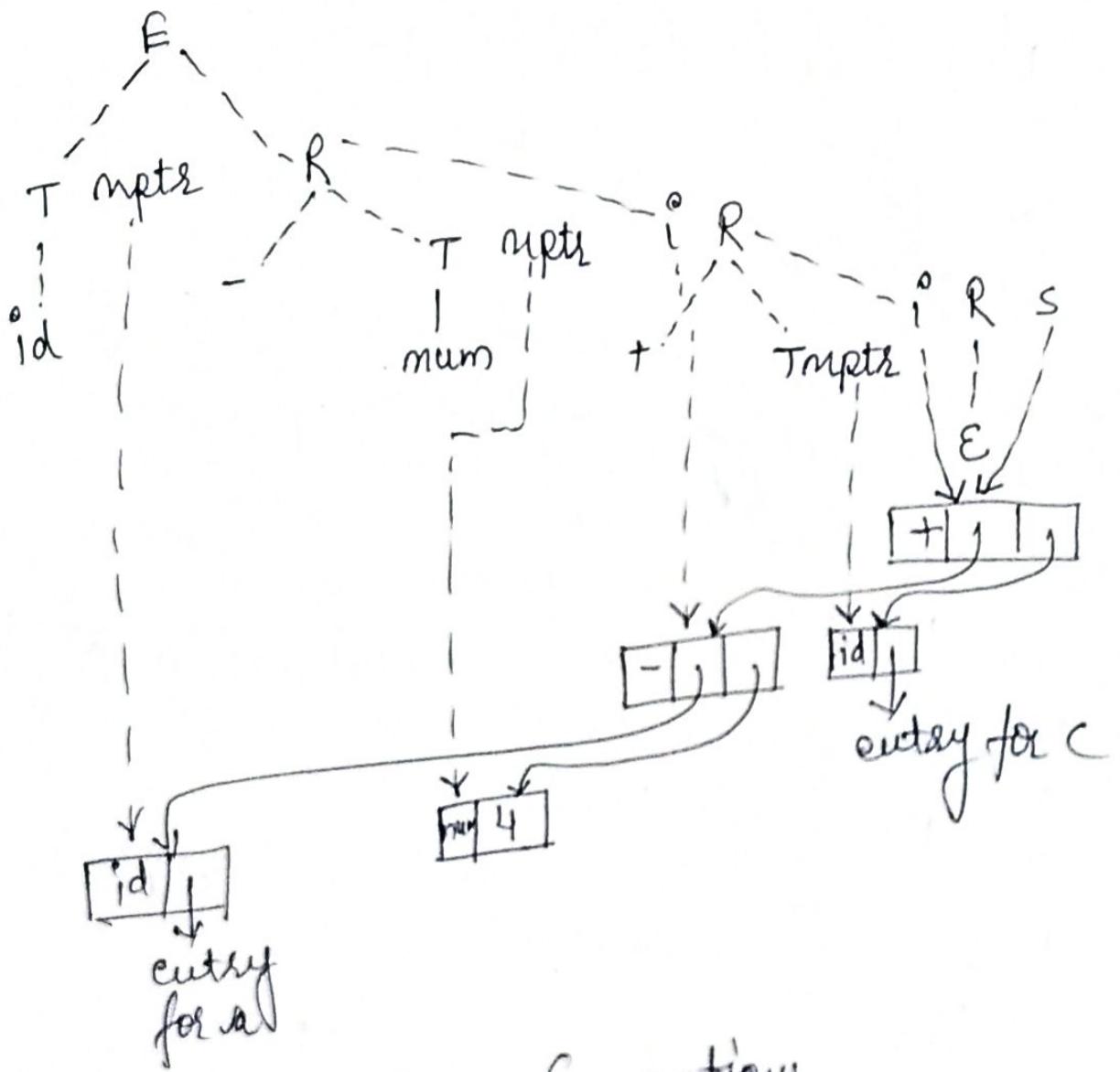
$$R \rightarrow \epsilon \quad \left\{ \begin{array}{l} R.S = R.i \end{array} \right\}$$

$$\begin{array}{l} T \rightarrow (E) \\ T \end{array} \quad \left\{ \begin{array}{l} T.mptr = E.mptr \end{array} \right\}$$

$$T \rightarrow id \quad \left\{ \begin{array}{l} T.mptr = mkleaf (id, id.entry) \end{array} \right\}$$

$$T \rightarrow num \quad \left\{ \begin{array}{l} T.mptr = mkleaf (num, num.val) \end{array} \right\}$$

Fig:- Use of inherited attributes to construct syntax trees. (9-4+1) (96)



Intermediate Code Generation:-

Intermediate Code Generation:-

Although a source program can be translated directly into the target language, but there are some benefits of using a machine-independent intermediate form all. For a different machine can back end.

- benefits of using:

 - a) A compiler for a different machine can be created by attaching a back end.
 - b) A machine-independent code optimizer can be applied to the intermediate representation.

Three-Address Code:-

(21)

- General form

$$x := y \text{ op } z$$

where x, y and z are names, constants or compiler-generated temporaries. Op is an operator like arithmetic, logical or boolean operator.

- No built-up arithmetic expressions are permitted in right side, as there is only one operator on the right side of a statement.

Ex:- $x + y * z$

$$t_1 := y * z$$

$$t_2 := x + t_1$$

- complicated arithmetic expressions and of nested flow-of-control statements makes three-address code desirable for target code generation and optimization.
- Three-address code is a linearized representation of a syntax tree or a dag.

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

$$t_1 := -c$$

$$t_2 := b * t_1$$

$$t_3 := -c$$

$$t_4 := b * t_3$$

$$t_5 := t_2 + t_4 \quad (\text{code for syntax tree})$$

$$a := t_5$$

$$t_1 := -c$$

$$t_2 := b * t_1$$

$$t_3 := t_2 + t_2$$

$$a := t_5$$

code for the
dag.

Types of Three-Address statements:-

- 1) Assignment statement $\rightarrow X := Y \text{ op } Z$
- 2) Assignment statement $\rightarrow X := \text{op } Y$ (op can be unary operator like unary minus, logical negation)
- 3) Copy stmts $\rightarrow X := Y$
- 4) Unconditional jump goto L.
- 5) conditional jumps like if $X < y$ goto L.
($<, =, >=$ etc.)
- 6) param X and call p, n
 param x_1
 param x_2

 param x_n
 call p, n

parameters

call to procedure p with n parameters
- 7) Indexed assignments $\rightarrow X[i] := Y$ or $X[i^j] := Y$
- 8) Address and pointers assignments $\rightarrow X := \&Y$, $X := *Y$, $*X := Y$

Implementation of Three-address statements:-

- Abstract form of intermediate code and implemented as records with fields for operator and operands

Q) Quadruples:-

(9)

- A record structure with four fields, op, arg1, arg2 & result.

$$\text{Ex:- } x_1 = y \text{ op } z$$

put y in arg1, z in arg2 and x in result

Ex:- Quadruples for the previous ex in TAC

	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

- Statements like many operators like $x_1 = -y$ with or $x_1 = y$ do not use arg2
- Operators like param use neither arg2 nor result.
- Conditional/unconditional jumps put the target label in result.

Triples:-

(160)

- In this, three-addressees are represented by records with only three fields : op, arg1 and arg2. The fields arg1 & arg2, for an op, all either pointers to the symbol table or pointers into the triple structure.

	<u>op</u>	<u>arg1</u>	<u>arg2</u>	
(0)	unminus	c		
(1)	*	b		(0) — pointers to the triple structure
(2)	unminus	c		
(3)	*	b		(2)
(4)	+		(1)	
(5)	assign	a		(3)
				(4)

symbol-table pointers

Ex:- $x[i] := y$

	<u>op</u>	<u>arg1</u>	<u>arg2</u>
(0)	[] =	x	i
(1)	assign	(0)	y

	<u>op</u>	<u>arg1</u>	<u>arg2</u>
(0)	= []	y	i
(1)	assign	a	(0)

Indirect Triples:-

- In this approach we list pointers to triples, rather than listing the triples themselves.

statement		op	arg1	arg2	(10)
(0)	(14)	(14)	uminus	c	
(1)	(15)	(15)	*	b	(14)
(2)	(16)	(16)	uminus	c	
(3)	(17)	(17)	*	b	(16)
(4)	(18)	(18)	+		(15) (17)
(5)	(19)	(19)	assign	a	(18)

Comparison of Representations: The use of imm-direction:-

- * Quaduple → using temporary can immediately access the location for that temp via the symbol table.
- In an optimizing compiler, statements are often moved around for doing optimization. If we move a stmt computing a, the stmt does not require changes. But in triples notation, moving a stmt that defines a temp value requires us to change all its references.

* Tuples:- save space

* Direct triples: - a stmt can be moved by reordering the statement list.

Boolean Expressions:-

- Boolean expressions have two primary purposes,
 - They used to compute logical values.
 - Also used in conditional expressions in statements that alter the flow of control such as if-then, if-then-else, while-do.

$E_1 \text{ eulop } E_2$
 ↴
 arithmetic
 exp

- Boolean exp are composed of boolean operators (and, or, not) applied to boolean variables or relational exp.

Ex:- $E \rightarrow E \text{ or } E | E \text{ and } E | \text{not } E | (E) | \text{id eulop id}$
 true | false

- eulop - $<, \leq, =, \neq, \gamma, \gamma$

Methods of translating Boolean expressions:-

- Encode true and false numerically and to evaluate a boolean exp analogously to an arithmetic exp. 1 represent true / 0 represent false.

d) The 2nd method of implementing
boolean exp is by flow-of-control,
that is, representing the value of a boolean
exp by a position reached in a program.
For ex, given exp $E_1 \text{ or } E_2$, if we determine
 E_1 is true, then we can conclude that
the entire exp is true without having to
evaluate E_2 .

(103)

Numerical Representation:-

- 1 denote true and 0 denote false.
- exp will be evaluated completely, from left to right.

Ex:- a or b and not c

Three address sequence:-

$t_1 := \text{not } c$

$t_2 := b \text{ and } t_1$

$t_3 := a \text{ or } t_2$

Ex:- if $a < b$ then 1 else 0

Three address code sequence:-

100: if $a < b$ goto 103

101: $t := 0$

102: goto 104

103: $t := 1$

104:

Short-Circuit Code:-

(104)

- In this, it is not necessary to evaluate the entire expression, this is called "short circuit" or "jumping" code.

Ex:- 100: if a < b goto 103

101: $t_1 := 0$

102: goto 104

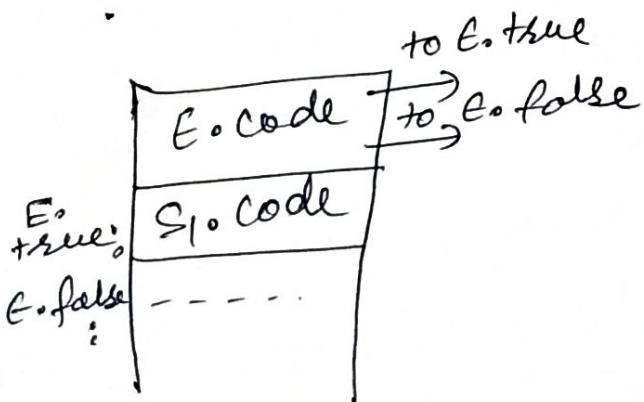
103: $t_1 := 1$

104: ---

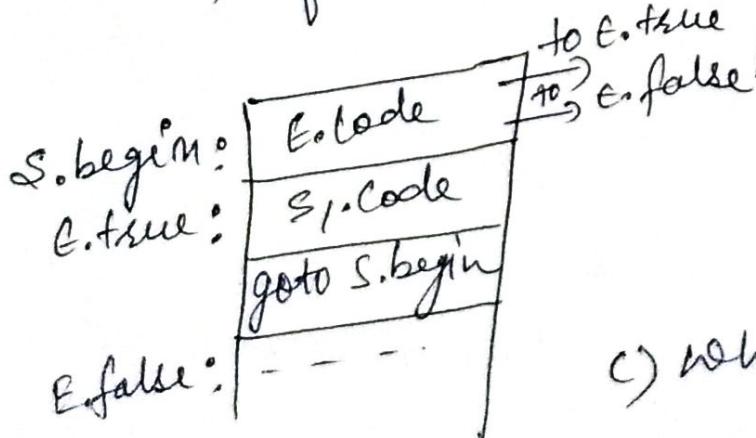
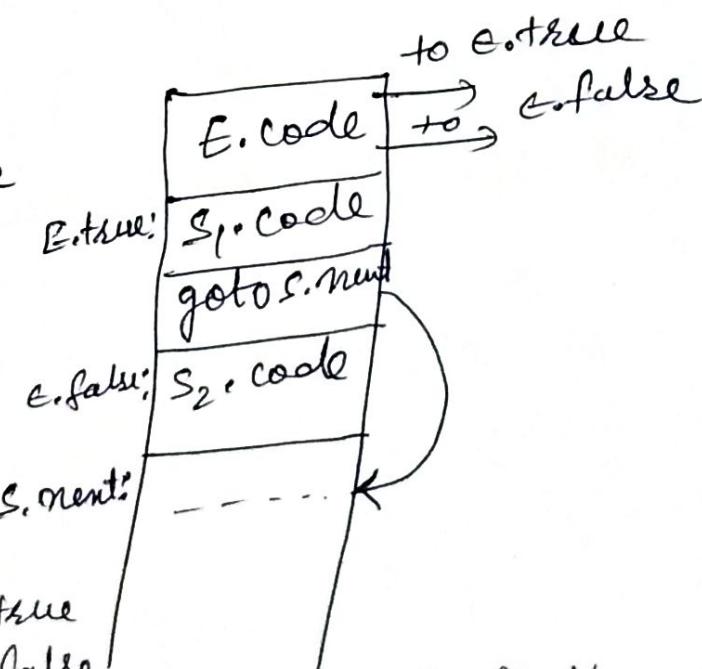
Here we can tell what value t_1 will have by whether we reach stmt 101 or 103.

control statements:-

$S \rightarrow \text{if } E \text{ then } S_1 \mid \text{if } E \text{ then } S_1 \text{ else } S_2$
 $\mid \text{while } E \text{ do } S_1$



a) if - then



c) while - do

b) if - then - else

Examples of three addrees code:-

(105)

1) $x = 1;$
 $\text{if}(x > 1)$
 $y = y + 3;$

1. $x = 1$
2. if $x > 1$ goto(4)
3. goto(6)
4. $t_1 = y + 3$
5. $y = t_1$
6. stop

2.) $x = 5;$
 $\text{if}(x < 5)$
 $z = z + 2;$
 else
 $y = y * 2;$

1. $x = 5$
2. if $x < 5$ goto(4)
3. goto(7)
4. $t_1 = z + 2$
5. $z = t_1$
6. stop goto(9)
7. $t_2 = y + 2$
8. $y = t_2$
9. stop

3.) int i
 $i = 1$
 while $i < 10$ do
 $\text{if } x > y \text{ then}$
 $i = x + y$
 else
 $i = x - y$

1. $i = 1$
2. if $i < 10$ goto(4)
3. goto(12)
4. if $x > y$ goto(6)
5. goto(9)
6. $t_1 = x + y$
7. $i = t_1$
8. goto(2)
9. $t_2 = x - y$
10. $i = t_2$
11. goto(2)
12. stop

Backpatching :- leaving the labels as empty and filling them later is called backpatching.

4.) $c = 0$
 do
 {
 if ($a < b$) then
 $x++$;
 else
 $x--$;
 $c++$;
} while ($c < 5$);

1. $c = 0$
2. if ($a < b$) goto (4)
3. goto (9)
4. $t_1 = x + 1$
5. $x = t_1$
6. $t_2 = c + 1$
7. $c = t_2$
8. goto (14)
9. $t_3 = x - 1$
10. $x = t_3$
11. ~~$t_4 = c + 1$~~
12. $c = t_4$
13. goto (14)
14. if $c < 5$ goto (2)
15. goto (15)
16. stop

5.) if ($a < b$ & $c < d$ or $e < f$)

101. if ($a < b$) goto (104)
102. $t_1 = 0$
103. goto (105)
104. $t_1 = 1$
105. if ($c < d$) goto (108)
106. $t_2 = 0$
107. goto (109)
108. $t_2 = 1$
109. if ($e < f$) goto (112)
110. $t_3 = 0$
111. goto (113)
112. $t_3 = 1$
113. $t_4 = t_1$ and t_2 .
114. $t_5 = t_4$ or t_3

6) $\text{for } (i=0; i<10; i++)$
 $a = b+c;$

$i=0$
L: if ($i < 10$) goto L1
goto last

(107)

L1: $t_1 = b+c$
 $a = t_1$
 $t_2 = i+1$
 $i = t_2$
goto L

Last: stop

7) Switch ($i+j$)

{

case 1:

$a = b+c;$

break;

case 2:

$P = Q+R;$

break;

default:

$x = y+z;$

break;

$t = i+j$
goto test

test: if ($t == 1$) goto L1
if ($t == 2$) goto L2
goto L3

L1: $t_1 = b+c$
 $a = t_1$
goto last

L2: $t_2 = Q+R$
 $P = t_2$
goto last

L3: $t_3 = y+z$
 $x = t_3$
goto last

Last: stop

8) $\text{for}(\overset{\circ}{i=1}; \overset{\circ}{i<10}; i++)$

(108)

{ $a[\overset{\circ}{i}] = x * 5;$
3

$i=1$

L: $t_1 = x * 5$

$t_2 = &a$

$t_3 = \text{sizeof}(\overset{\circ}{int})$

$t_4 = t_3 * 1$

$t_5 = t_2 + t_4$

$*t_5 = t_1$

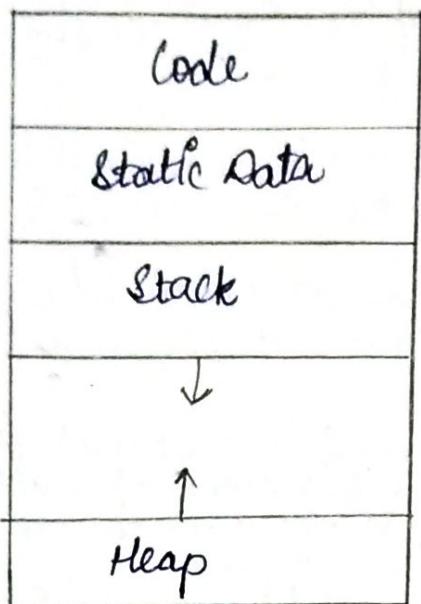
$i^{\circ} = i^{\circ} + 1$

if ($i^{\circ} < 10$) goto L

goto last

last: stop

Storage Organization:-



- Generated target code: compiler can place it in a statically determined area.
- Static Data: which data whose size is known at compile time.
- Stack: when a call occurs, execution of an activation is interrupted and info about the status of the machine, such as the value of the program counter and machine registers, is saved on the stack. When control returns from the call, this activation can be restarted after restoring the values of relevant registers and setting the program counter to the point immediately after the call.
- Heap: used for dynamic memory allocation.
- Stack and heap grows in opposite side.

Activation Records: info needed by a single execution of a procedure is managed using a contiguous block of storage called an activation record or frame, consisting of the collection of fields.

returned value	used by called procedure to return a value to the calling procedure.
actual parameters	used by calling procedure to supply parameters to the called procedure.
optional control link	points to the activation record of the callee.
optional access link	non-local data held in other activation records.
saved machine status	info about the status (program counter & registers) of machine just before the procedure call.
local data	local data to a procedure.
temporaries	temporary values during exp evaluation.

Compile-Time layout of local data:-

- Run-time storage comes in blocks of contiguous bytes, (smallest unit of addressable memory).
- The amount of storage needed for a name is determined from its type. An elementary data type, such as a char, int or real can usually be stored in an integral number of bytes.
- For aggregate, such as array & record uses large space to hold its all components. For easy access, aggregates also use contiguous block of bytes.

- For local data memory locations are fixed at compile time. Variable-length data is kept outside this field. We keep a count of the memory locations that have been allocated for previous declarations. From the count we determine a relative address of the storage for a local. The relative address is offset w.r.t the difference b/w the addresses of the position and the data object.

- The storage layout for data objects is strongly influenced by the addressing constraints of the target machine. Like some machines have 16 bits for int and some have 32 bits for int.

Storage Allocation Strategies:-

Static Allocation:-

- In static allocation, names are bound to storage as the program is compiled, so there is no need for a run-time support package. Since the bindings do not change at run time, every time a procedure is activated, its name are bound to the same storage locations. This property allows the values of local names to be retained.

across activations of a procedure. That is, when control returns to a procedure, the values of the locals are the same as they were when control left the last time.

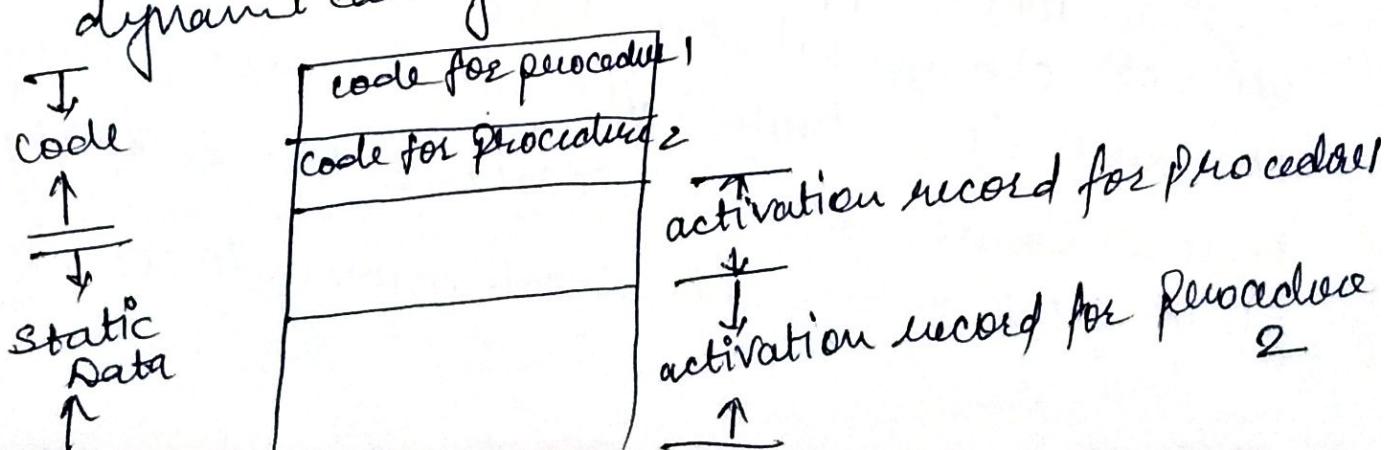
(112)

- From the type of Name, compiler determine the amount of storage to set aside for that name.
- In this scheme, at compile time we can therefore fill in addresses at which target code can find the data it operates on. Similarly, compiler has to know on which address, which info is saved.

Limitations:-

- 1) Size of data object and constraints on its position in memory must known at compile time.
- 2) Recursive procedures are restricted, bcoz all activations of a procedure use the same bindings for local names.
- 3) Data structures cannot be created dynamically.

Ex:- Fortran 77



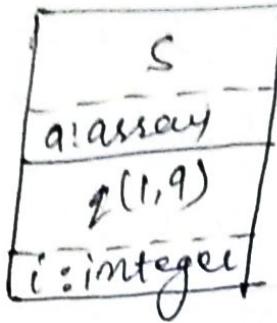
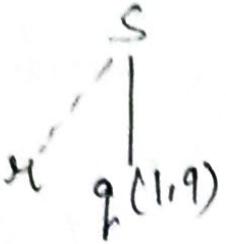
Stack Allocation:-

(113)

- Based on the idea of a control stack; storage is organised as a stack, and activation records are pushed and popped as activations begin and end respectively.
- storage for locals are freshly allocated at each time when a call made, bcoz the storage for locals disappears when the activation record is popped.
- Register top marks top of the stack.
- Suppose procedure q has an activation record of size a is pushed, top will be ↑ by a and after control returns from q, top is decremented by a.

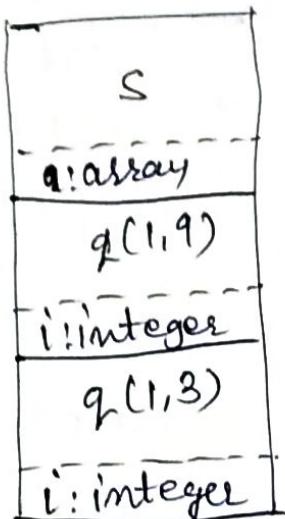
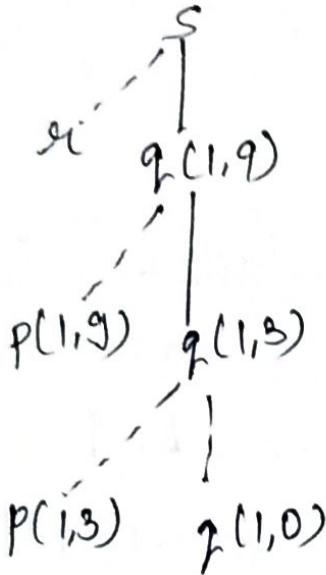
Ex:-

position in activation tree	Activation Records on the stack	Remarks
s	<div style="border: 1px solid black; padding: 5px; display: inline-block;"> s a:array </div>	frame for s
r	<div style="border: 1px solid black; padding: 5px; display: inline-block;"> s a:array r i:integer </div>	e is activated



114

frame for `r`
has been popped
and `q(1,9)`
pushed

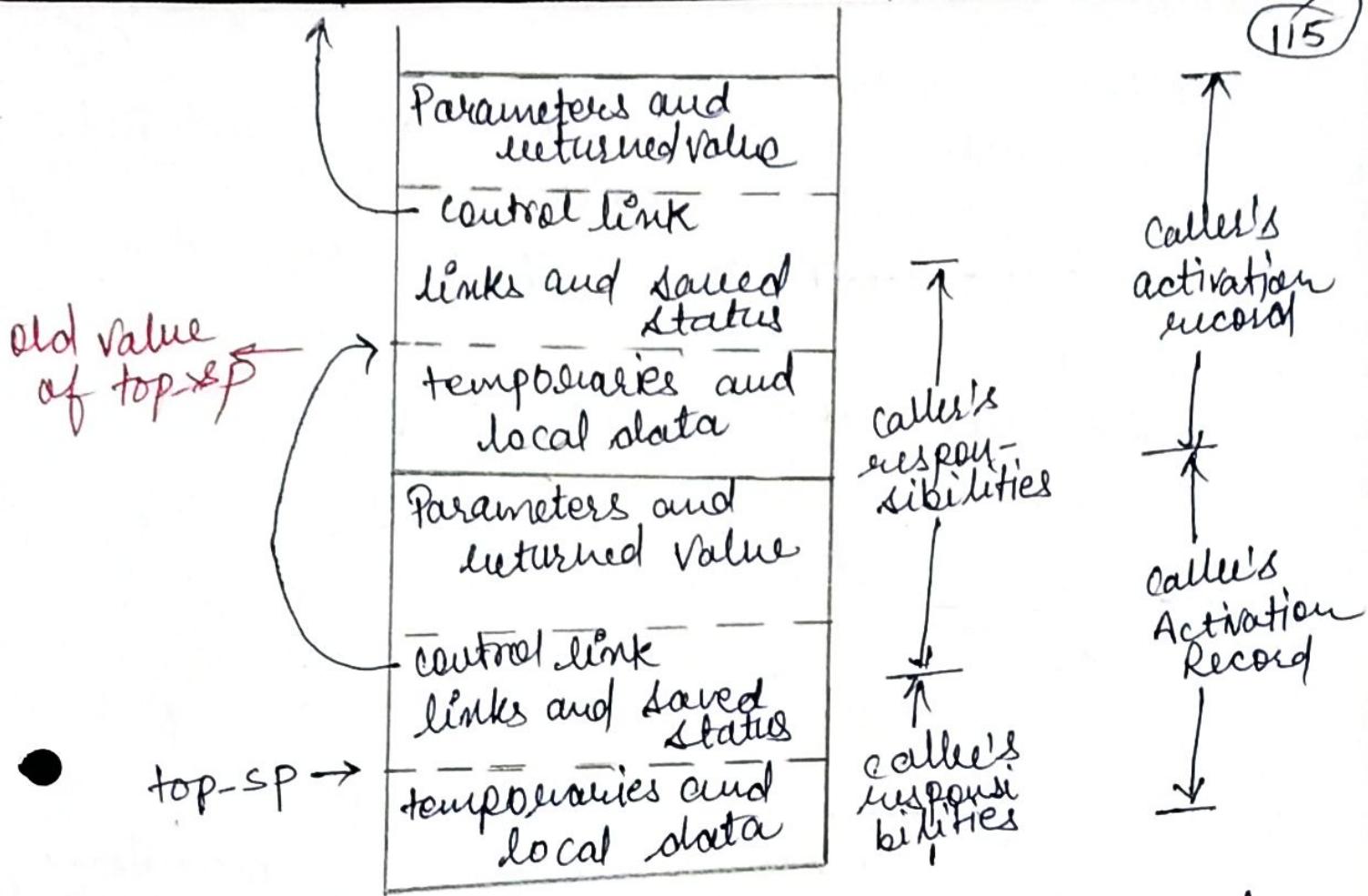


control has
just returned
to `q(1,3)`

Calling Sequences :-

- A call sequence allocates an activation record and enters info into its fields.
- A return sequence restores the state of the machine so the calling procedure can continue execution.

fixed-length data:-



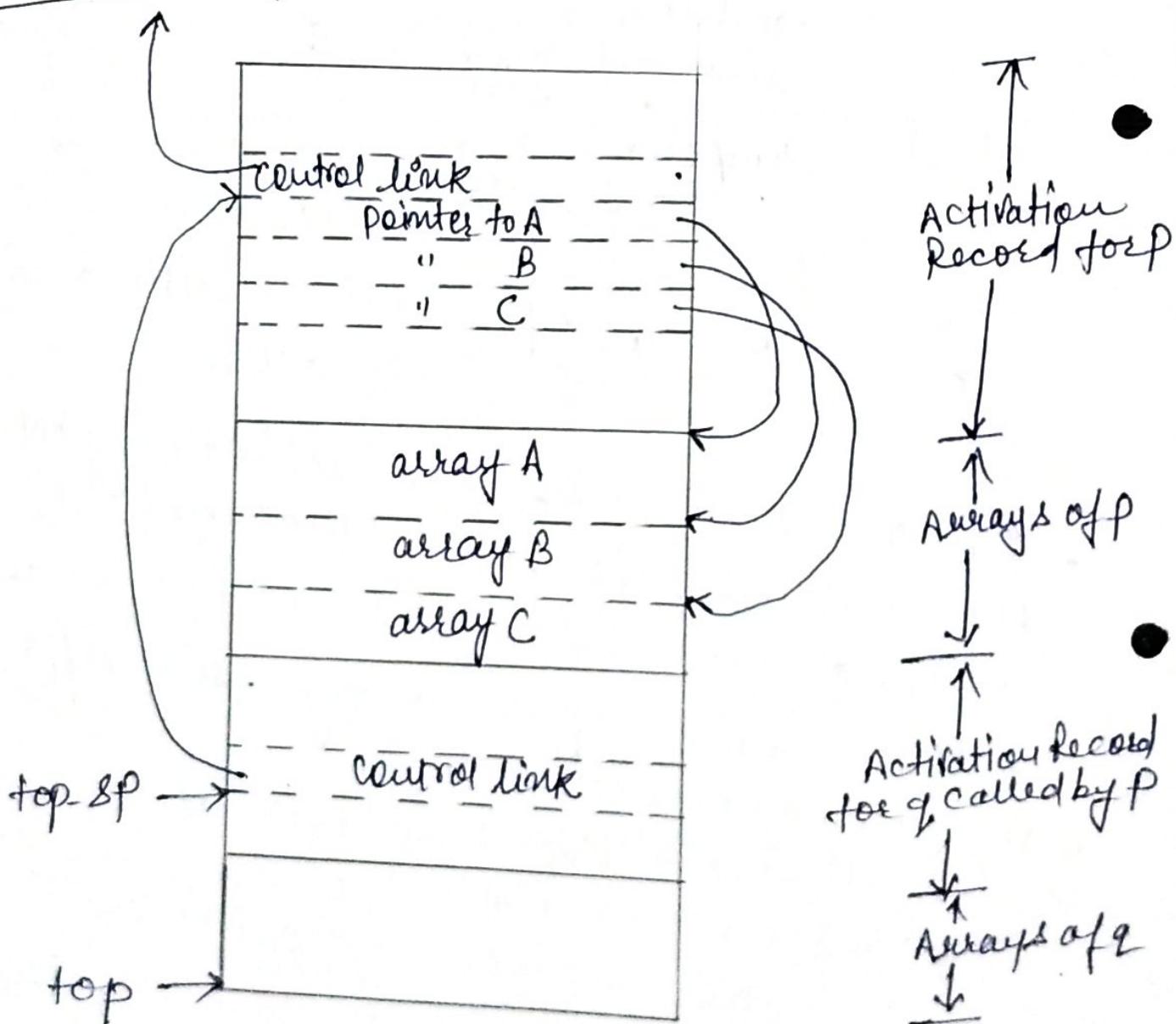
Division of tasks b/w caller and callee.

- In run-time stack, the activation record of the callee is just below that for the caller. So that callee can access parameters & returned values of caller by using offset.
- Sequence of tasks:-
 - * Caller evaluates actual parameters
 - * Caller stores old value of top-SP into the activation record of callee.
 - * top-SP is moved to new position.
 - * Caller saves register values and other status info.
 - * Callee initializes its local data and begins execution.

• Return Sequence:-

- * callee places a return value next to the activation record of caller.
- * Caller restores top-sp and other registers
- * Caller copy the value of top-sp

Variable-length Data:-



- P has three local arrays. storage for 117
these arrays are not part of activation record, only a pointer to the beginning of each array appears in the AR.
- q is called by P. AR for q begins after array for P, and variable length arrays of q begin beyond that.
- Access data using $\begin{cases} \text{top} & \text{actual top of stack} \\ \text{top-sp} & \text{used to find local data} \end{cases}$
- When q returns.
 - * $\text{top} = \text{top-sp}$ - length of machine status and parameter fields in q's activation record
 - * top-sp - copied from the control link of q.

Dangling references:-

- A dangling reference occurs when there is a reference to storage that has been deallocated.

Ex: `main()`

```

  {
    int *p;
    p = dangle();
  }
  
```

```

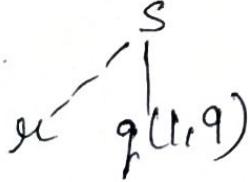
int *dangle()
{
  int i=23;
  return &i;
}
  
```

• Here when control returns to main from dangle, storage for local is freed. Since P in main refers to this storage, the use of P is a dangling reference. (118)

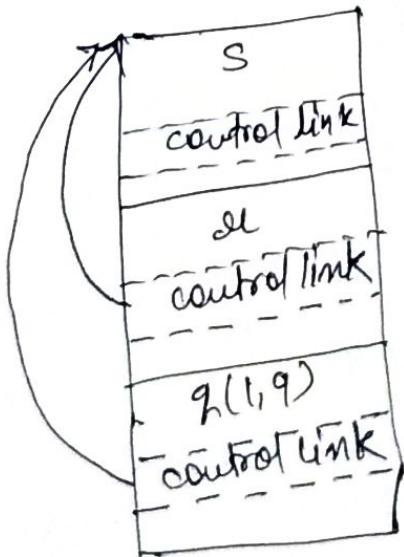
Heap Allocation:-

- Stack allocation will not work, if
 - * Values of local names must be retained when an activation ends.
 - * A called activation outlives the caller.
- Heap allocation can be used for these. It parcels out pieces of contiguous storage, as needed for activation records or other objects.

Ex:- Activation Tree



AR in the Heap



Remarks

retained activation record for r.

- for each size of interest, keep a linked list of free blocks of that size.
- fill a request for size s with block of size s' , where s' is the smallest size greater than or equal to s .
- For large blocks of storage use the heap manager.

Access to nonlocal names :-

Blocks :-

- A block is a statement containing its own local data declarations.
- { declarations statements }

Ex:-

main()

```
{ int a=0;
int b=0;
```

```
{ int b=1;
```

```
{ int a=2;
printf("%d %d", a, b); }
```

B₀

B₁

```
{ B2
```

```
int b=3;
printf("%d %d \n", a, b); }
```

printf ("%d %d \n", a, b);

(20)

printf ("%d %d \n", a, b);

Fig: Blocks in a C program

Declaration

int a=0;

int b=0;

int b=1;

int a=2;

int b=3;

Scope

B0-B2

B0-B1

B1-B3

B2

B3

Output a, b

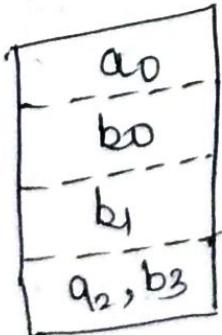
2 1

0 3

0 1

0 0

- Can be implemented using stack. Space for the declared name can be allocated when the block is entered and deallocated when control leaves the block.
- An alternative implementation is to allocate storage for a complete procedure body at one time.



a₂ and b₃ may be assigned the same storage becoz they are in blocks that are not alive at the same time.

Lexical scope without nested procedures (121)

Ex:- `int a[10];
readarray() { ... a ... }
int partition(y,z){int y, z; ... q...}
quicksort(m,n){int m,n; ... e...}
main() { ... a ... }`

- Here non-local occurrences of `a` in `readarray`, `partition` and `main` refer to the array declared on line 1.
- In absence of nested procedures, the stack allocation strategy for local names can be used directly for a lexically scoped language like C.
- Any non-local names can be allocated statically. Position for this storage is known at compile time, so we can use it directly.
- Local names to an activation at the top of stack, easily accessible through the top pointer.
- Benefit of static allocation for nonlocals is that declared procedures can freely be passed as parameters and returned as results.

Lexical scope with Nested Procedures:-

(122)

- A nonlocal occurrence of a name 'a' in a Pascal procedure is in the scope of the most closely nested declaration of 'a' in the static program text.

Ex:- Program Sort (input, output);
Var a: array [0..10] of integer;
x:

procedure readarray;

Var i --

procedure exchange (- -);

procedure quicksort(- -);

Var k, l --

function partition(- -)

Var i, j --

begin --- a ---

--- exchange (- -)

end { - - }

begin --- end { quicksort }

begin --- end { sort }

This 'a' is first search
in partition, then in
quicksort, after that
in sort.

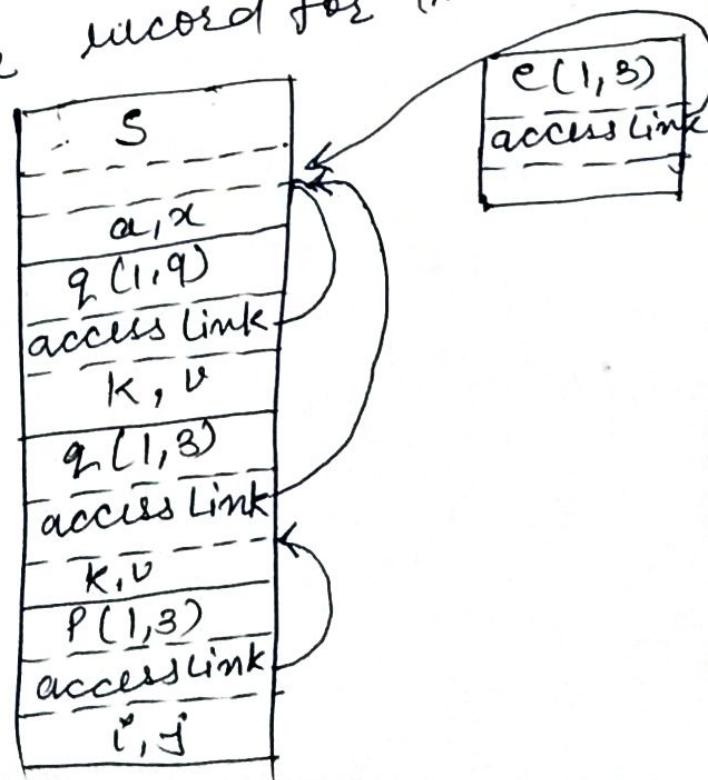
This exchange is
first searched in
quicksort, then in
sort.

1) Nesting Depth:-

- Nesting depth of a procedure is used to implement lexical scope.
- Let the name of the main program be at nesting depth 1; we add 1 to the nesting depth as we go from an enclosing to an enclosed procedure.
- In the above example 'quicksort' is at nesting depth 2, while partition ~~at~~ is at nesting depth 3.

2) Access Links:-

- A direct implementation of lexical scope for nested procedures is obtained by adding a pointer called an access link to each activation record. If procedure P is nested immediately within Q in the source text, then the access link in an activation record for P points to the access link in the record for most recent activation of Q.



3) Procedure Parameters:-

(124)

Ex-1. Program param (input, output)

2. procedure b (function h (n: integer): integer;
er);

3. begin ---.

4. procedure c;

5. var m: integer;

6. function f (n: integer): ---

7. begin f: ---.

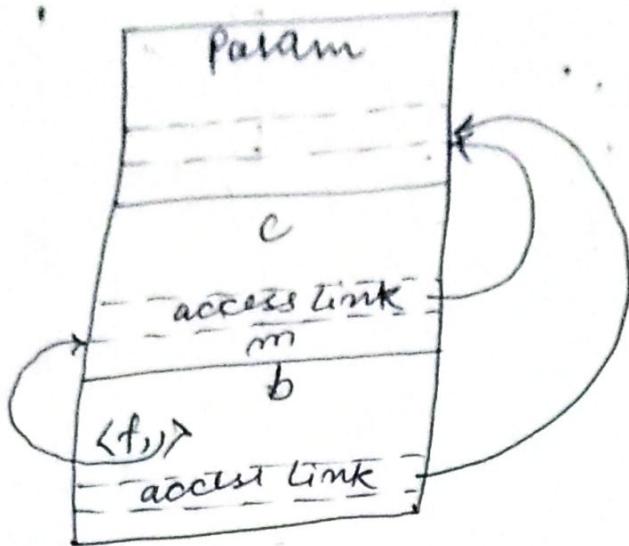
8. begin m:=0; b(f) end {c};

9. begin

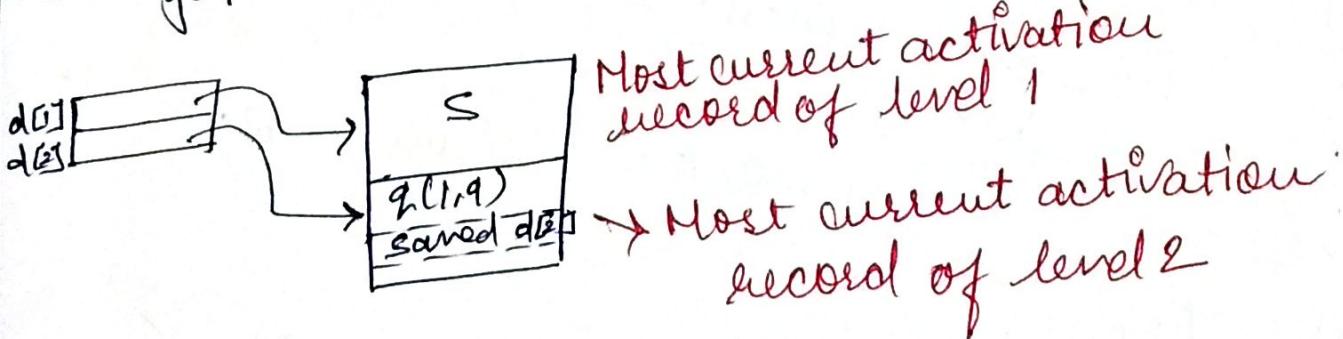
10. c

11. end

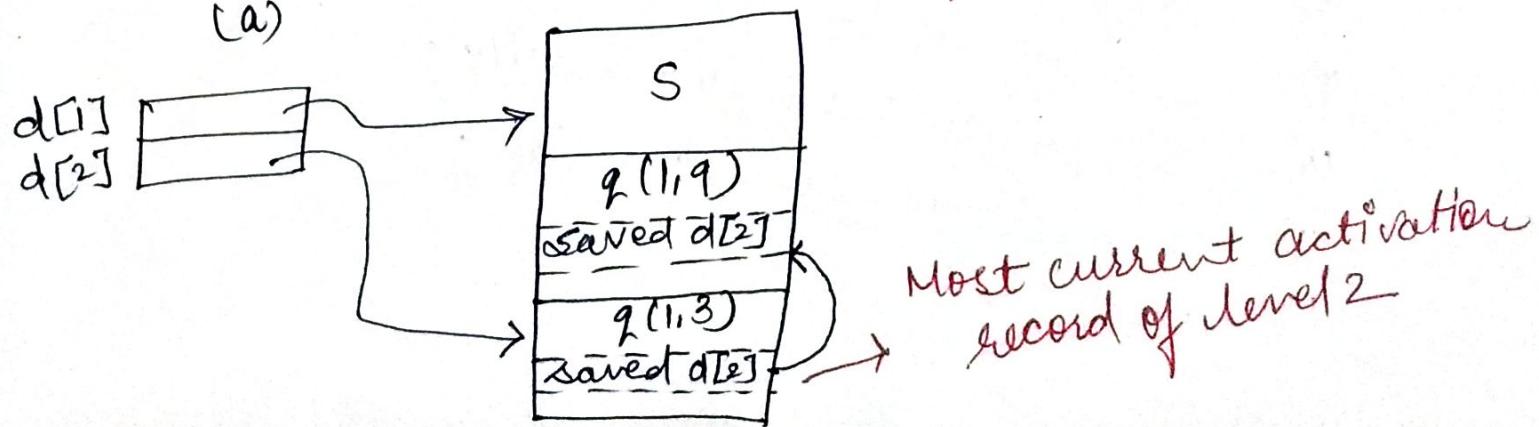
- Enclosed scope rules apply even when a nested procedure is passed as a parameter.
- If procedure c at line 8, f (function) is passed as parameter to procedure b. So when a nested procedure is passed as a parameter must take its access link along with it.
- When procedure c passes f, it will also pass access link of f to b.



- 4) Displays:
- faster access to nonlocals than with access links can be obtained using an array d of pointers to activation records, called a display.
 - Suppose control is in an activation of a procedure p at nesting depth j . Then, the first $j-1$ elements of the display point to the most recent activations of the procedures at depth $j-1$.



(a)



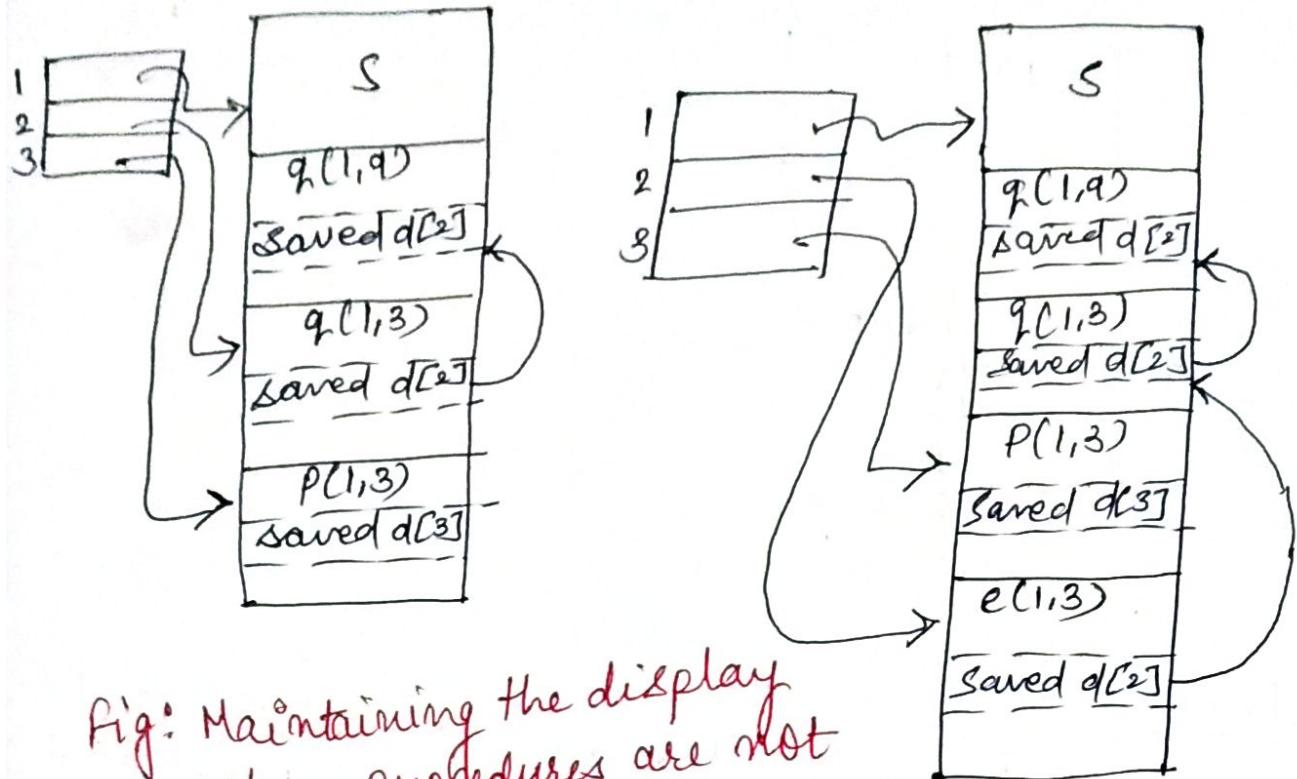


Fig: Maintaining the display
when procedures are not
passed as parameters

Dynamic Scope:-

- Ex:- 1) program dynamic (—);
 2) Var a —
 3) Procedure Show;
 — write (a) —
 4) Procedure Small;
 5) Var a —
 6) — a = 0.125 ; Show end;
 7)
 8) begin
 9) a = 0.25
 10) Show; Small; write;
 11) Show; Small; write;
 12) end

- In above example, when show is called on lines 10-11 in main program, 0.250 is written bcoz x is local to main function. But when show is called $x = 0.125$ is used as this x is local to ~~show~~. small.
- This is called dynamic scope.
- Two approaches for this,
 - Deep access:- in this control links are used rather than access links. In this search for a nonlocal can go deep into the stack, maybe for the first activation record.
 - Shallow access:- hold the current value of each name in statically allocated storage. When a new activation of a procedure p occurs, a local name m in p takes over the storage statically allocated for M.

Parameter Passing:-

- Both nonlocals and parameters are used by the procedure below,

ex:- procedure exchange ($i, j; integer$);

Var x: integer;

begin

$x := a[i]; a[i] := a[j]; a[j] := x$

end

• Here i^o and j^o are parameters and 128
a is nonlocal name.

- Different parameter passing methods exists like call by value, call by reference, copy restore etc.
- Diff b/w parameter-passing methods are based primarily on whether an actual parameter represents an r-value, an l-value or the text of the actual parameter itself.

Call-by-Value :-

- Simplest method of passing parameters.
- Formal parameters are said from called procedure. And the storage for formals is in the activation record of the called procedure.
- Called uses these parameters, modify them but these modifications are not seen at called procedure.

Ex:- Swap (x, y)
{
 int x, y ;
 int temp;
 temp = x ;
 $x = y$;
 $y = temp$;
}

main()

{
 int a=1, b=2;
 swap (a, b);
 cout ("a=%d, b=%d",
 a, b);
}

a and b will not be modified at main function.

call-by- reference:-

(129)

- when parameters are passed by reference
(call-by-address or call-by-location)

Ex:- swap(x, y)

{ int *x, *y }

int temp

temp = *x

*x = *y

*y = temp

main()

{ int a=1, b=2;

swap(&a, &b);

printf("a is %d, b is %d", a, b);

}

(Footman uses this)

- Copy-Retore:- (Footman uses this)
 - Hybrid between call-by-value and call-by-reference.
 - Also known as copy-in copy-out or value-result.
 - Before control flows to the called procedure, the actual parameters are evaluated. The l-values of the actuals are passed to the called procedure as in call by value. In addition, the l-values of those actual parameters having l-values are determined before the call.
 - When control returns, the current l-values of the formal parameters are copied back into the l-values of the actuals, using the l-values computed before the call.

Ex:- program copyout (input, output);

(130)

```

var a: integer;
procedure unsafe (var x: integer);
begin x := 2; a := 0 end;
begin
a := 1; unsafe(a); writeln(a)
end.

```

Here a's value is passed to x at 1. But in procedure unsafe x will be assigned 2 and this value 2 is assigned to a by restoring it - this value 2 is returned from function.

call-by-Name :-

- It is traditionally defined by the copy - rule of Algol, which is:
 - The procedure is treated as if it were a macro that is, its body is substituted for the call in the caller, with the actual parameters literally substituted for the formals. Such a literal substitution is called macro-expansion or in-line expansion.

Ex:- #define f(var)
do {
 val = 3;
 i + = 1;
 var = 4;
 while (0)

```

void g(void) {
    int a[2];
    int i;
    i = 0;
    f(a[i]);
}

```

• Here as we define F as a function then it will work like call-by-value. But macro will replace its calling in a function and works like call-by-name and give correct output.

(131)

■ Symbol Tables:-

- A compiler uses a symbol table to keep track of scope and binding information about names. The symbol table is searched every time a name is encountered in the source text. Changes to the table occur if a new name or new info about an existing name is discovered.
- It is useful for a compiler to be able to grow the symbol table ~~dynamically~~ dynamically if necessary at compile time, rather than fixing the size of symbol table at compile time.

Symbol-Table Entries:-

- Each entry in the symbol table is for the declaration of a name. The format of entries does not have to be uniform, bcoz the info saved about a name depends on the usage of the name.
- Each entry can be implemented as a record.

- For uniformity, it can be possible to keep max info outside the table entry, with only a pointer to this info stored in the record.
- Info about a name can be entered into the symbol table at various times. like keywords are entered into the symbol table initially. Further entries are done by lexical analyzer. Attributes values are filled as info being available.

Ex:- `int x;`
`struct x {float y,z;};`

Here two entries will be done in symbol table for x, one as int and one as struct.

- Attributes of a name are entered in response to declarations.

characters in a Name-

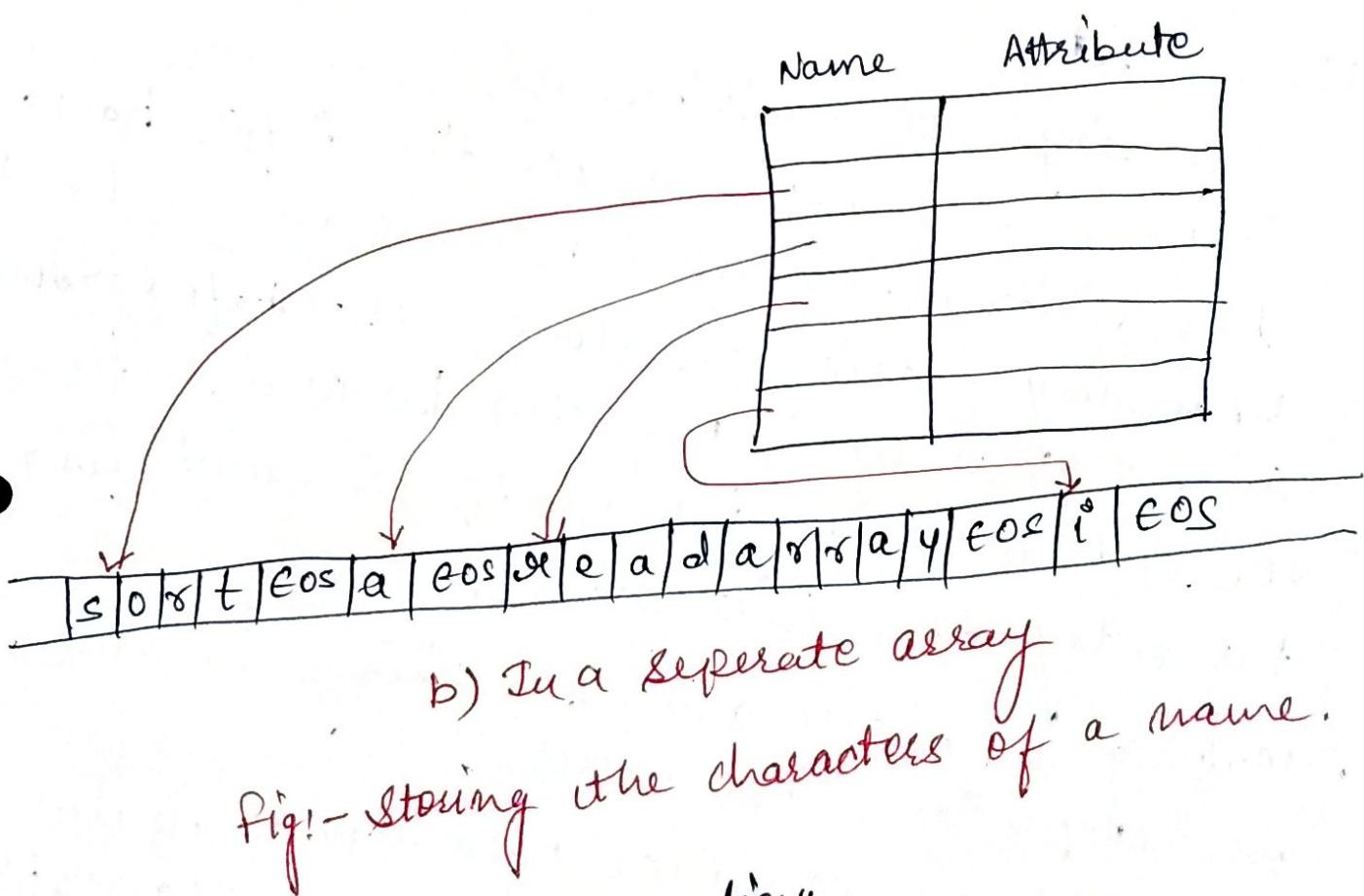
When we are doing entry for first time it is considered as lexeme, but when a lexeme is input found in symbol table it is considered as a name. A common representation of a name is a pointer to a symbol-table entry for it.

- If there is a modest upper bound on the length of a name, then the characters in the name can be stored in the symbol-table entry. If there is no limit on the length of a name, the indirect scheme can be used.

(133)

Name	Attributes
s o u t a g e a d a g a r a y i	

a) Be fined -
size & place
within a
record



Storage Allocation Information:-

- Storage Allocation Information:-

 - Info about the storage locations that will be bound to names at run time it kept in

the symbol table.

- Static storage:- if target code is in assembly language then assembler will take care of storage. If target code is in machine language, then the position of each data object is relative to fixed origin.
- Names whose storage is allocated on a stack or heap, the compiler does not allocate storage at all.

■ The List of Data structure for Symbol Tables.

1) Linear List :-

- Use single or multiple arrays to store names and their associated info. Array's last position is marked as available and upcoming entries are done here. For searching an entry, array is processed backwards from the end to the beginning. If we don't find the entry then fault occurs.
- Making an entry and searching a name are independent operations.
- In general, an occurrence of a name is in the scope of most closely nested declaration of the name. It can be implemented using list by making a fresh entry for a name every time it is declared.

- By searching from available towards the beginning of the array, we are sure to find the most recently created entry.

(35)

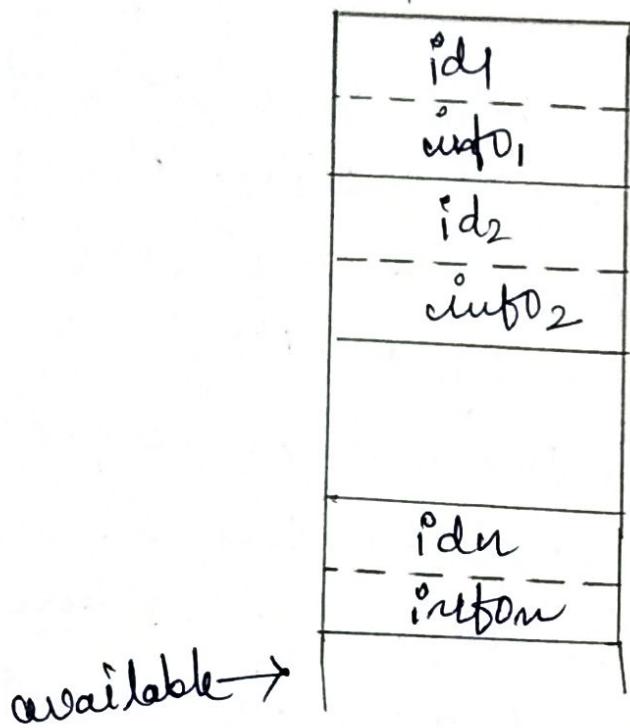


Fig: linear list of records.

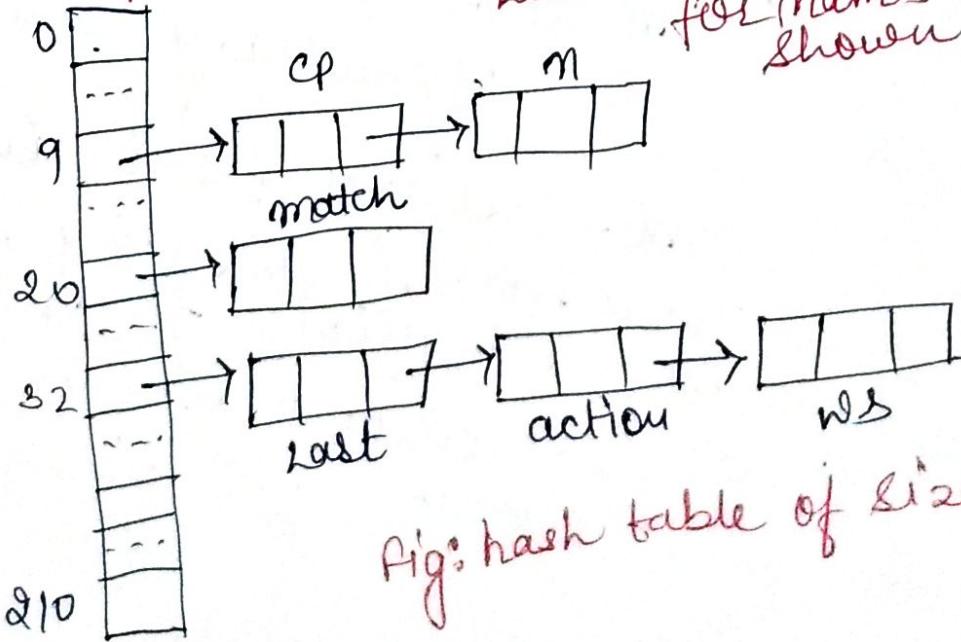
- Symbol table contains $\rightarrow m$ names
 - * if we do not check for multiple entries in symbol table and make a new entry, it will require constant work
 - * if multiple entries are not allowed, m - require work to check whether entry already available or not.
- m - also cost for an inquiry
- So Total work for inserting m names and making e inquiries = $C(m + e)$
where C = constant time requires for machine operations.

- As m and e increases too much time is spent for performing operations.

Q) Hash Tables:-

- ~~Hash~~ Hashing is another data structure used for symbol table creation. Here we are considering "open hashing", where "open" refers to the property that there need be no limit on the number of entries that can be made in the table.
- To perform inquiries on names requires $= \frac{m(\alpha e)}{m}$ time, for constant α .
- As m increases, the space requirement grows, so a time-space tradeoff is involved.

- Array of list headers, indexed by hash value



- Two parts of hash table
 - * A hash table consisting of a fixed array of m pointers to table entries.
 - * Table entries organized into m separate linked lists, called buckets.
- To determine whether there is an entry for string s in the symbol table, we apply a hash function h to s , such that $h(s)$ returns an integer between 0 and $m-1$. If s is in the symbol table, then it is on the list numbered $h(s)$. If s is not yet in the symbol table, it is entered by creating a record for s that is linked at the front of the list numbered $h(s)$.
- The choice of m depends on the intended application for a symbol table. Choosing m to be large enough, make table lookup to a negligible time.
- Approach for designing hash function,
 - * Determine a set of integers from the characters c_1, c_2, \dots, c_k in string s by adding ASCII values of its characters.
 - * Convert m into b/m and $m-1$. (Divide by m and taking remainder. Generally m is taken as prime no.)

- Variations for calculating h_i
 - * add its character values.
 - * $h_i = dh_{i-1} + c_i$ (d is a constant)
 - *

138

Representing Scope Information:-

- When an occurrence of a name in the source text is looked up in the symbol table, the entry for the appropriate declaration of that name must be returned. (Means name from correct scope must be returned.)
 - Simple approach - Maintain separate symbol table for each scope.
 - 2nd Approach - every procedure must provide a unique number. And this no will be appear with declaration of that name.
 - 3rd Approach:- a pointer "front" points to the most recently created entry in the list. Entries will be made from front and will be searched from front to last. So that we can access most recent name (local name)
 - Deletion can also be performed in same way.

front

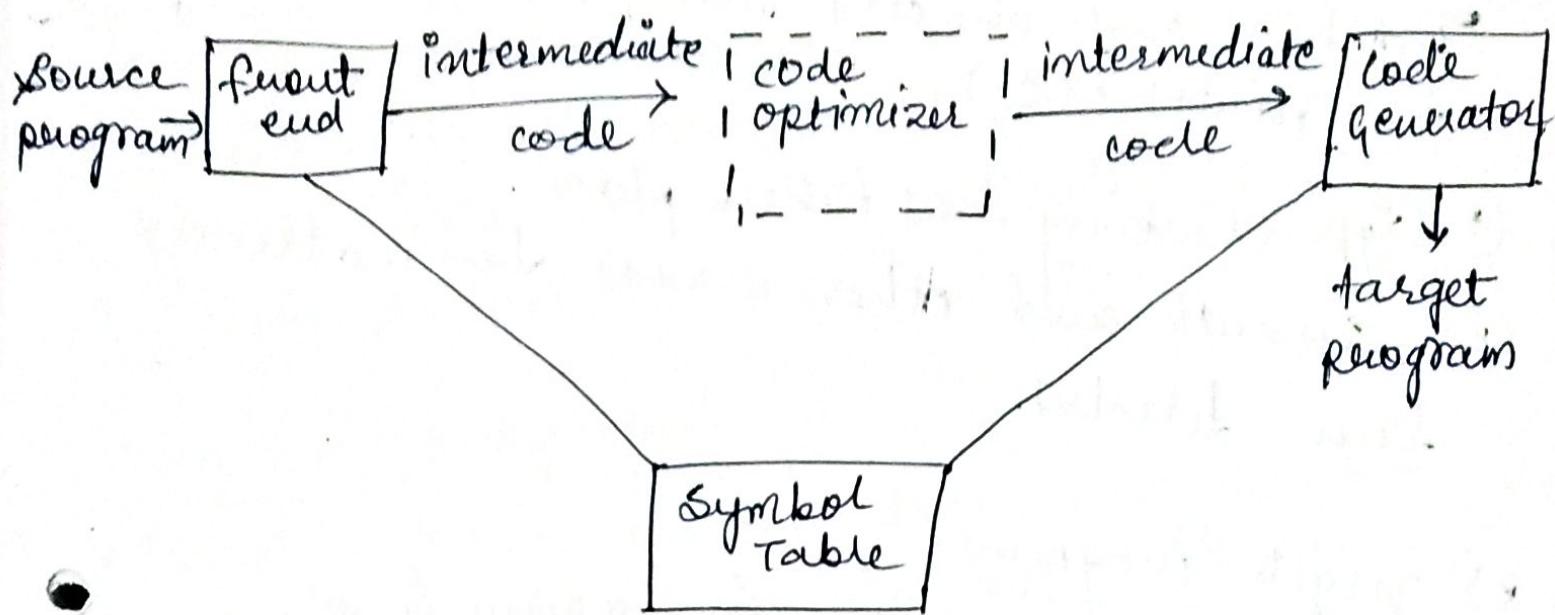
• 4th Approach:- A hash table consists of 139 m lists accessed through an array. Since a name always hashes to the same list, individual lists are maintained as previous example.

• Deletion of entries from the ~~list~~ hash table must be done with care. If we delete i^{th} entry, then $i-1^{\text{st}}$ entry must point to $i+1^{\text{st}}$ entry. For this hash lists must be circular linked list, so that $i-1^{\text{st}}$ entry can be accessed.

UNIT-VI

(140)

Code Generation



- Code optimizer only exists in "optimizing compilers".
- Requirements imposed on a code generator,
 - * output code must be correct and of high quality
 - * it should make effective use of resources
 - * code generator itself should run efficiently.

Issues in the design of a code generator:-

1) Input to the code generator:-

- consists of intermediate representation of the source program and info in symbol table that is used to determine run time addresses.
- code generation assumes flat, ① front end has scanned, parsed and translated the source program into a reasonably detailed intermediate

representation. So the values of names can be represented by quantities that the target machine can directly manipulate (bit, integer, real, pointer etc.)

- ① Type checking has taken place.
- ② Semantic and other errors have already been detected.

2) Target Programs:-

Target code can be in many forms,

- ① Absolute machine language programs - it can be placed in a fixed location in memory and immediately executed.

Ex:- PLC

② Relocatable machine language programs - it allows subprograms to be compiled separately. A set of relocatable object modules can be linked together and loaded for execution by a linking loader. For this we will have added expense. It provides great flexibility in being able to compile subroutines separately and to call other previously compiled programs from an object module.

⑧ Assembly language program:- makes code generation easier. We can generate symbolic instructions and use the macro facilities of assembler to help generate code. Price paid is the assemble step after code generation.

37 Memory Management:-

- Mapping names in the source program to addresses of data objects in run-time memory is done co-operatively by the front end and the code generator.
- A name in a three-address refers to a symbol-table entry for the name. Symbol-table entries were created as the declaration in a procedure were examined. The type in a declaration determines the width, i.e., the amount of storage needed for the declared name. From the symbol-table info, a relative address can be determined for the name in a data area for the procedure.
- If machine code is being generated, labels in three-address statements have to be converted to address of instruction.

4) Instruction Selection:-

- The nature of the instruction set of the target machine determines the difficulty of instruction selection. The uniformity and completeness of the instruction set are imp factors.
- Instruction speeds and machine idioms are other imp factors.
- Instructions must be chosen like this that the target program is efficient.

Ex:- $a := b + c$
 $d := a + e$

```

Mov b, R0
add e, R0
Mov R0, a
Mov a, R0
add e, R0
Mov R0, d

```

Here $start$ is redundant, and so is the third if a is not subsequently used.

- Different instruction sets leads to different implementation of given operation and may have significant cost difference.

- A naive (lack of experience) translation may lead to correct but unacceptable target code.

Ex:- $a := a + 1$

MOV a, R0 .

Adl #1, R0

MOV R0, a

This must be like,
Inc a

5) Register Allocation:-

- Instructions involving register operands are usually shorter and faster than those involving operands in memory. Efficient utilization of registers is particularly imp in generating good code. The use of registers is often sub-divided into two subproblems:

① During "register allocation", we select the set of variables that will reside in registers at a point in the program

② During a subsequent "register assignment" phase, we pick the specific register that a variable will reside in.

- finding an optimal assignment of registers to variables is difficult.

- certain machines require register-pairs (even-odd) for some operands and results.

Ex:- IBM System/370 multiplication operation
 $M \times Y$ (odd) \times will be in a register and
 $(even) Y$ will be in another register.

⑥ Choice of Evaluation Order:-

(145)

- The order in which computations are performed can affect the efficiency of the target code. Some computation orders require fewer registers.

7) Approaches to Code Generation:-

- The most imp criterion for a code generator is that it produce correct code.

Basic Blocks and Flow Graphs:-

- Graphical representation of three-address statements is called flow graph. Nodes in the flow graph represent computations and the edges represent the flow of control.

Basic Blocks:- a basic block is a sequence of consecutive statements in which flow of control enters at the beginning and leaves at the end without halt or possibility of branching except at the end.

$$\text{Ex} \quad a^2 + 2ab + b^2$$

$$t_1! = a * a$$

$$t_2! = a * b$$

$$t_3! = 2 * t_2$$

$$t_4! = t_1! + t_3$$

$$t_5! = b * b$$

$$t_6! = t_4! + t_5!$$

- Algo: Partition into basic blocks 146
 I/p: Sequence of 3-address stmts.
 O/p: A list of basic blocks
 Method:

- 1) A stmt will be a leader:
 - i) A first stmt is a leader
 - ii) Any stmt that is the target of a conditional/unconditional goto is a leader
 - iii) Any stmt that immediately follows a goto or conditional goto stmt is a leader.
- 2) for each leader, its basic block consists of the leader and all stmts up to but not including the next leader or the end of program.

Ex:- begin
 period := 0;
 $i^{\circ} := 1^{\circ}$
 do begin
 $period := period + a[i] * b[i]$
 $i^{\circ} := i^{\circ} + 1$
 end
 while $i < 20$
 end

- first block
1. $\text{perod} = 0 \rightarrow \text{leader by rule 1}$
 2. $i = 1$
 3. $t_1 = 4 * i \rightarrow \text{leader by rule 2}$
 4. $t_2 = a[t_1]$
 5. $t_3 = 4 * i$
 6. $t_4 = b[t_3]$
 7. $t_5 = t_2 * t_4$
 8. $t_6 = \text{perod} + t_5$
 9. $\text{perod} = t_6$
 10. $t_7 = i + 1$
 11. $i = t_7$
 12. $\text{if } i <= 20 \text{ goto 3}$
 13. $\rightarrow \text{leader by rule 3}$

Transformations on Basic Blocks:-

- A basic block computes a set of expressions.
- Two basic blocks are said to be equivalent if they compute the same set of expressions.
- A number of transformations can be applied to a basic block without changing the set of expressions computed by the block. These transformations can improve the quality of code.

1) Structure - Preserving Transformations:-

(i) common subexpression elimination:-

These can not be eliminated as end stmt redefine b.

$$\begin{array}{l} a := b + c \\ b := a - d \\ c := b + c \\ d := a - d \end{array} \rightarrow \begin{array}{l} a := b + c \\ b := a - d \\ c := b + c \\ d := -b \end{array}$$

(ii) Dead-code elimination:- suppose x is dead, that is, never subsequently used, at the point where stmt $x = y + z$ appears. So we can eliminate this stmt.

(iii) Renaming temporary variables:-

$t = b + c$ we can replace all occurrence of a & t by u , this will not create a difference.

(iv) Interchange of statements:-

Ex- $t_1 = b + c$ } can be interchanged as
 $t_2 = x + y$ } independent of each other.

2) Algebraic Transformations:-

Countless algebraic transformations can be used to change the set of expressions computed by a basic block into an algebraically equivalent set.

Ex:- $x = x_0 + 0$ } can be eliminated
 $x = x - x_1$

(149)

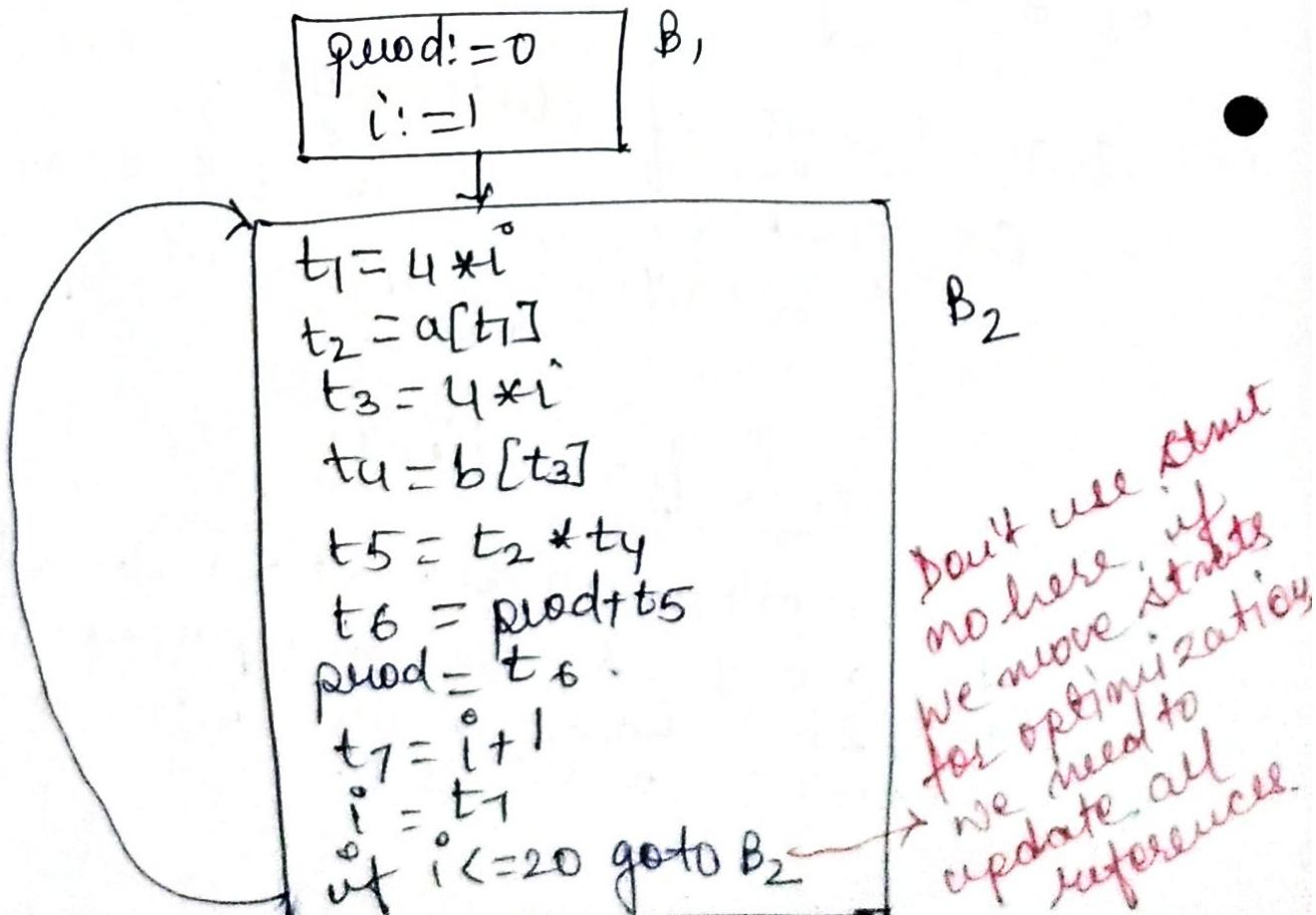
Ex:- $x = y * x^2$
 can be replaced with,

$$x = y * y$$

Flow Graphs:-

- we can add the flow-of-control info to the set of basic blocks making up a program by constructing a directed graph called a "flow graph".
- The nodes of the flow graph are the basic blocks.

Ex:- Flow graph for the prog given in basic block, (vector product)



Ex:- TAC

```

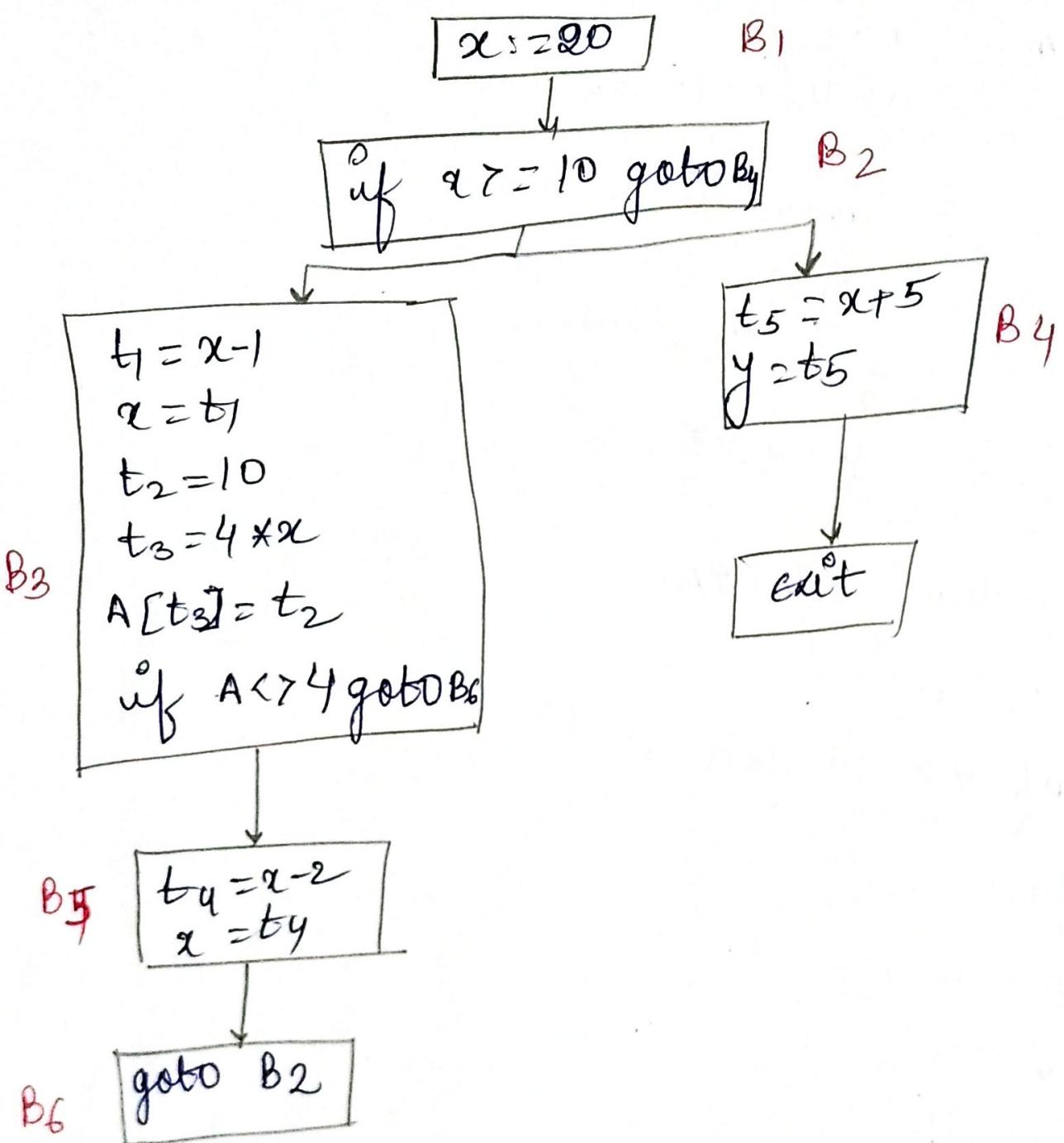
 $x := 20$ 
while( $x < 10$ ) {
     $x := x - 1$ ;
    A[x] = 10;
    if( $x == 4$ )
         $x := x - 2$ 
}
 $y = x + 5$ 

```

Intermediate Code (TAC):-

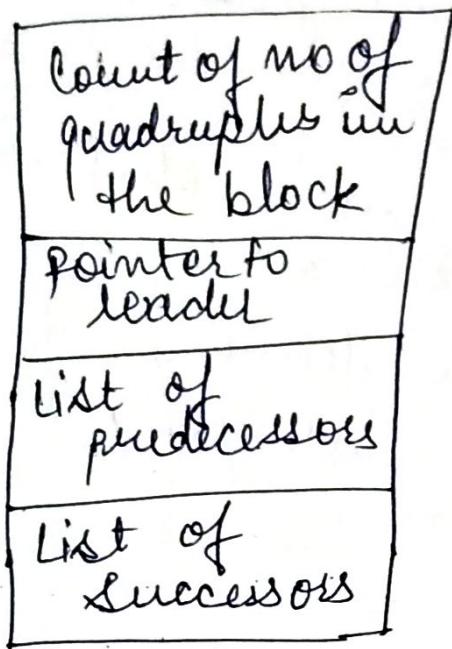
- 1) $x := 20$
- 2) $\text{if } x >= 10 \text{ goto } \underline{12}$
- 3) ~~goto —~~
- 4) $t_1 = x - 1$
- 5) $t_2 = 10$
- 6) $t_3 = H * x$
- 7) $A[t_3] = t_2$
- 8) $\text{if } x < H \text{ goto } \underline{11}$
- 9) $t_4 = x - 2$
- 10) $x = t_4$
- 11) $\text{goto } \underline{2}$
- 12) $t_5 = x + 5$
- 13) $y = t_5$

(151)



Representation of Basic Block:-

- Can be represented by variety of data structures.
- By Record:- (In this each basic block is represented by following record)



- By linked list :- make linked list of the quadruples in each block.

Loops:- ex- One previous ex there is a loop consisting of b2.

- Loop is a collection of nodes such that,
 - (i) All nodes in the collection are strongly connected. (There is a path from each node to other)
 - (ii) The collection of nodes must have a unique entry.

A Simple Code Generation -

- Code generation strategy consider each statement in turn, remembering if any of the operands of the statmt are currently in registers, and taking its advantage.
- Computed results are generally kept in registers as long as possible, store them only,
 - if register is needed for another computation or
 - just before leaving a basic block
(procedure call, jump or labeled stmt)

Ex: $a = b + c$

Add R_j, R_i
Result in R_i , cost=1 (Both operands exists in registers)

Add c, R_i
Result in R_i , cost=2 (one operand exist in memory)

MOV c, R_j

Add R_j, R_i

Result in R_i , cost=3 (c can be used again & again, as it exist in R_j)

Substituted
statement
live

- Many more cases to consider,
 - * where b, c are located.
 - * whether b 's value will be subsequently used.
 - * either b or c is a constant (both)
 - * whether $+$ is commutative

(154)

We have to consider a large no of cases before code generation.

Register and Address Descriptors:-

- Uses descriptors to keep track of register contents and addresses for names, and address for names,
- (i) A register descriptor keeps track of what is currently in each register
- (ii) An address descriptor keeps track of the location where the current value of the name can be found at run time. The location might be a register, a stack location, a memory address or some set of these.

A Code-generation algorithm:-

- Input: - sequence of 3-address stmts from a basic block.
- For each stmt $x := y \ op z$ do following,
 - (i) Invoke a func getreg to determine location L where result will be stored. (Register/Memory)
 - (ii) Get address^(*) of y from address descriptor. If y is at both register/memory. Then choose register.

- (i) if value of y is not already in L , then $\text{Mov } y, L$.
 - (ii) generate instruction $\text{Op } z, L$ (z is location of z). Update address descriptor of x to indicate that x is in location L . Remove x from all other register descriptors.
 - (iv) If after the current uses, y and z are not live then alter registers, that they no longer will contain y or z .
- If stmt has a unary operator then same steps will be followed.
- if stmt $x := y$,
 - (i) if y is in a register, then simply change the register and address descriptor, and x will now contain y only in the register containing y .
 - (ii) if y is in a memory location, then use getreg to find a register in which to load y and make that register location of x .
- After processing of all TAC in a basic block, we store, by Mov instructions, those names that are live on exit and not in their locations.

The function getreg:-

(36)

- This function returns the location L to hold the value of x in $x := op2$. Scheme for selecting the location,
 - (i) if y is in register that holds the value of no other names and y is not live after this stmt then return register of y for L.
 - (ii) failing (i), return an empty register for L.
 - (iii) failing (ii), if x has next use in the block, then find an occupied register R, store the value of R in memory and return R.
 - (iv) if x is not used in the block, or no suitable occupied register then select memory location for x as L.

Ex:- $d := (a-b) + (a-c) + (a-c)$

TAC

$$\begin{aligned}t &:= a - b \\u &:= a - c \\v &:= t + u \\d &:= v + u\end{aligned}$$

<u>Statements</u>	<u>code generated</u>	<u>Register Descriptor</u>	<u>Address Descriptor</u>
$t := a - b$	MOV a, R0 SUB b, R0	R0 contains t	t in R0
$u := a - c$	MOV a, R1 SUB c, R1	R0 contains u R1 contains u	t in R0 u in R1
$v := t + u$	ADD R1, R0	R0 contains v R1 contains v	u in R1 v in R0
$d := v + u$	ADD R1, R0 MOV R0, d	R0 contains d	d in R0 d in memory

Fig:- Code Sequence

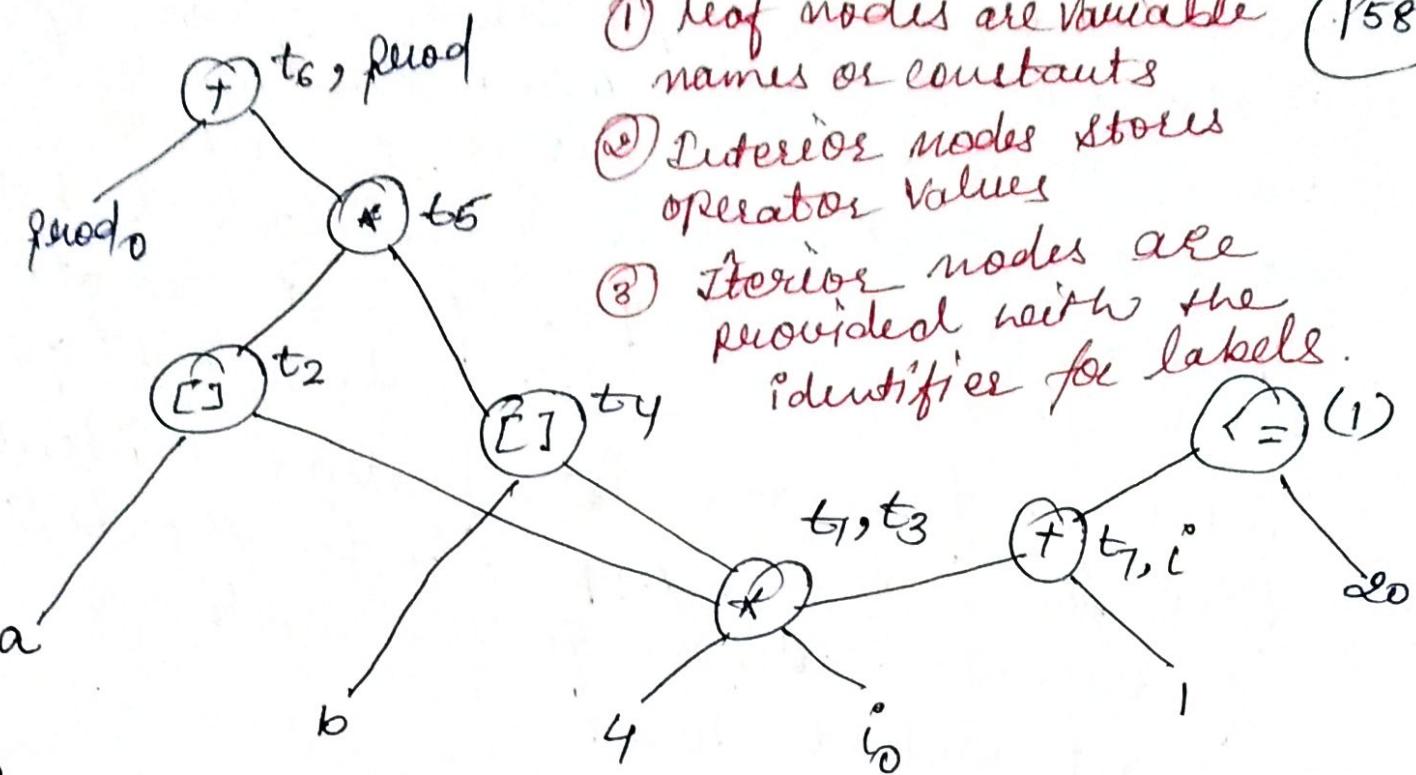
The Dag Representation of Basic Blocks:-

- Directed acyclic graphs (dag) are useful data structures for implementing transformations on basic blocks. A dag gives a picture of how the value computed by each statement in a basic block is used in subsequent statements of the block.

Ex:-

- (1) $t_1 := 4 * i$
- (2) $t_2 := a[t_1]$
- (3) $t_3 := 4 * i$
- (4) $t_4 := b[t_3]$
- (5) $t_5 := t_2 * t_4$

- (6) $t_6 := \text{prod} + t_5$
- (7) $\text{prod} := t_6$
- (8) $t_7 := i + 1$
- (9) $i := t_7$
- (10) if $i <= 20$ goto(1)



Dag construction:-

Algo: constructing a dag

Ip: A basic block

Opp: A dag

Method:

case 1: $x := y \text{ op } z$

case 2: $x := \text{op } y$

case 3: $x := f$

i. If $\text{node}(y)$ is undefined, create a ~~leaf~~ leaf labeled y , then set this to $\text{node}(y)$. In case 1, if $\text{node}(z)$ is also undefined, then create $\text{node}(z)$.

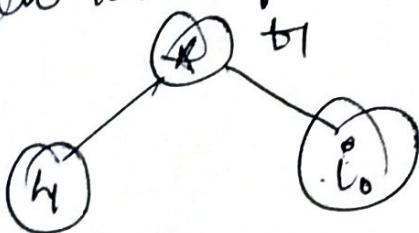
ii. In case 1, if there is a node labeled op with left child (y) and right child (z) .

(to check common subexpressions.) If not, create such a node. Assume m is the node found or created. In case 2, determine whether there is a node labeled op , and alone child (y). If not, create node m . In case 3, let n be $\text{node}(y)$.

3. Delete x from the list of attached identifiers for node (x). Append x to the list of attached identifiers for the node m and set $\text{node}(x)$ to n .

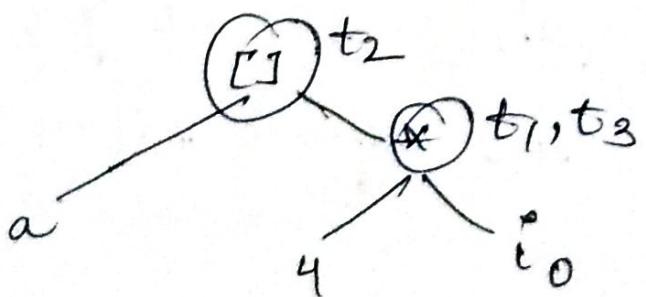
Ex:- $t_1 := 4 * i$

- create leaves for 4 and i .
- create a node labeled $*$.
- attach identifier t_1 to it.



(a)

$$t_2 := a[t_1] \text{ and } t_3 = 4 * i$$



Applications of Day:-

(160)

- 1) We automatically detect common subexpressions.
- 2) We can determine which identifiers have their values used in the block. (those for which a leaf is created in step 1)
- 3) We can determine which statements compute values that could be used outside the block. (those statements whose node is constructed or found in step 2)
- 4) List of quadruples are simplified by eliminating the common subexpressions. Also the assignment of the form $a := b$ cannot be performed good code until and unless it is necessary to do this.

Peephole Optimizations:-

- The quality of target code can be improved by applying "optimizing" transformations to the target program.
- A effective technique is peephole optimization.
- It is a method for trying to improve the performance of the target program by examining a short sequence of target

(161)

instructions (called peephole) and replacing these instructions by a shorter or ~~faster~~ faster sequence, whenever possible.

- Can be applied directly after intermediate representation.
- Peephole is a small, moving window on the target program. The code in peephole need not be contiguous.

Examples of program transformations using peephole optimization:-

1) Redundant loads and stores:- (medium-daut Instructions)

Ex:- ~~Mov R0, A~~
~~Mov A, R0~~ → can be deleted.

But if second stmt has a label then we can't remove it.

2) Unreachable code:- instructions.

• Remove unreachable

Ex:- #define debug 0

if(debug) { point debugging into }

Intermediate code for this,
 if debug = 1 goto L1

goto L2

L1: print debugging info^r

L2:

After eliminating jumps over jumps,

if debug ≠ 1 goto L2

print debugging info^r

L2:

Replacing 'debug' with its defined value,

if 0 ≠ 1 goto L2

print debugging info^r

L2:

As in above example, 1st statement
 always evaluated to true, it can be
 replaced by goto L2. Now all statements
 "printing debugging info^r" are unreach-
 able and can be eliminated one at a
 time.

3) flow-of-control optimization:-

~~Unnecessary~~ unnecessary jumps are eliminated.

Ex:- goto L1 goto L2
 - - - - - - - - - -
L1: goto L2 → L1: goto L2

Ex:- if $a < b$ goto L1 if $a < b$ goto L2
 -- -- -- -- -- --
 L1: goto L2 → L1: goto L2

4) Algebraic Simplifications:-

- Algebraic Simplification

 - There is no end to the amount of algebraic simplification that can be attempted through peephole optimization.

Ex:- $x := x + 0$ } can be eliminated
 or $x := x * 1$ } directly

5) Reduction in Strength:-

- This replaces expensive operations by equivalent cheaper ones.

Ex- x^2 is invariably cheaper to implement as $x * x$ than as a call to an exponentiation routine.

6) Use of machine Idioms:-

- Use of machine fallows:-

 - The target machine may have nine instructions to implement certain specific operations efficiently. Detecting situations that permit

the use of these instructions can reduce execution time significantly. 164

Ex:- Some machines have auto-increment and auto-decrement addressing modes.

Generating Code from Dags :-

- This section shows generating code for a basic block from its dag representation.

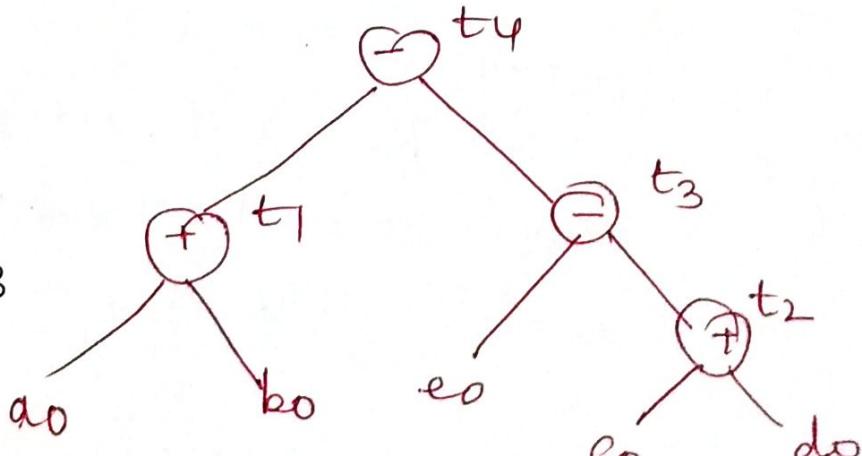
Rearranging the order:-

- In the following example, it is shown, how the order in which computations are done can affect the cost of resulting object code.

$$\text{stmt} := (a+b)-(e-(c+d))$$

Ex:-

$$\begin{aligned}t_1 &:= a+b \\t_2 &:= c+d \\t_3 &:= e-t_2 \\t_4 &:= t_1 - t_3\end{aligned}$$



Code Sequence:-

$R_0 - a$	$\text{MOV } a, R_0$
$R_0 - a+b$	$\text{ADD } b, R_0$
$R_1 - c$	$\text{MOV } c, R_1$
$R_1 - c+d$	$\text{ADD } d, R_1$
$t_1 - a+b$	$\text{MOV } R_0, t_1$

$\text{MOV } e, R_0$	$R_0 - e$
$\text{SUB } R_1, R_0$	$R_0 - e - (c+d)$
$\text{MOV } t_1, R_1$	$R_1 - t_1$
$\text{SUB } R_0, R_1$	$R_1 - t_1 - (e - (c+d))$
$\text{MOV } R_1, t_4$	$t_4 \rightarrow$

After rearranging the order of stmts,

$$\begin{aligned}t_2 &:= c + d \\t_3 &:= e - t_2 \\t_1 &:= a + b \\t_4 &:= t_1 - t_3\end{aligned}$$

Code sequence :-

MOV C, R0	MOV a, R0
ADD D, R0	ADD b, R0
MOV E, R1	SUB R1, R0
SUB R0, R1	MOV R0, t4

A Heuristic ordering for Page :-

A node listing algorithm is presented in this section, it gives the order of statements of three-address code.

(1) while unlisted interior nodes remain do

begin

(2) Select an unlist node n , all of whose parents have been listed

(3) list n ;

(4) while the leftmost child m of n has no unlisted parents and is not a leaf do,

begin

 list m ;

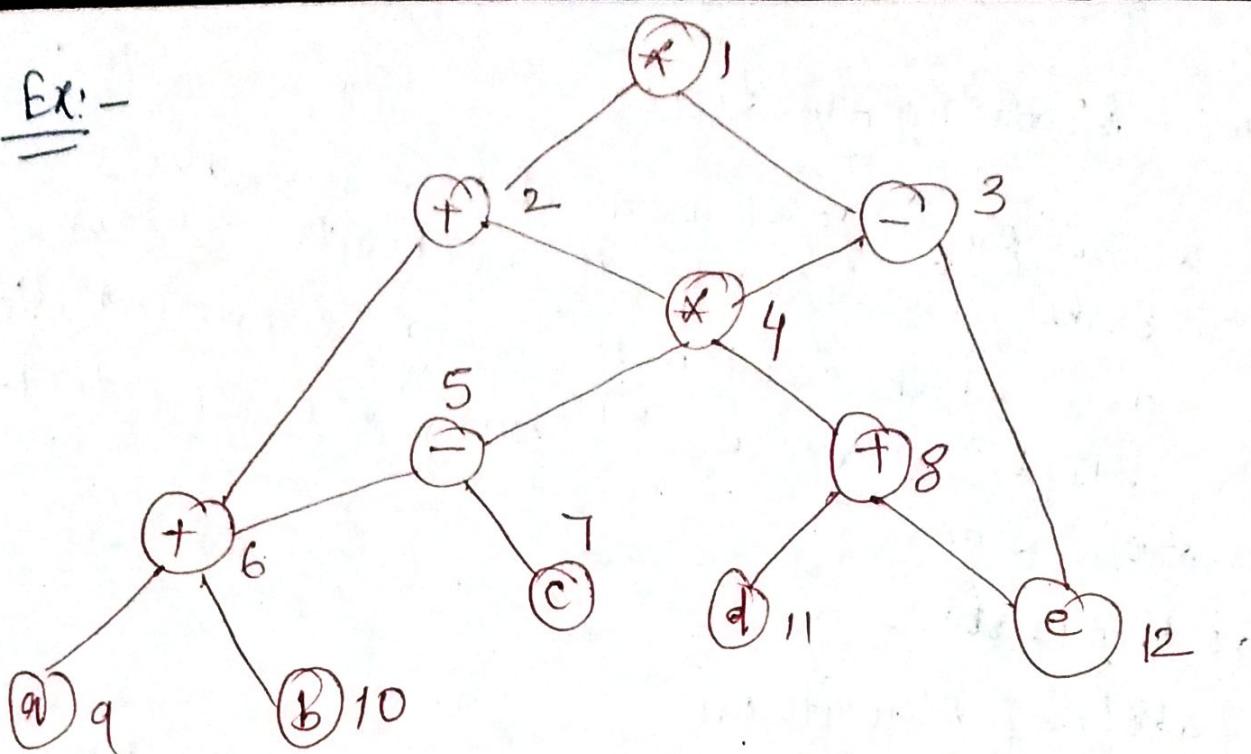
$m := m'$

(5)

end

end

Fig: Node listing Algorithm

Ex:-

- Initially only node 1 is there which has no unlisted parents, so $n=1$
- left child of 1 is 2, has its parent listed, so $n=2$
- now left child of 2 is 6 which has an unlisted parent 5, so we select new mat line ② of algo as 3. List 3
- Now proceed down its left chain 4, 5, 6.
- This leaves only 8. So list it.
- Resulting list is 1234568. So the suggested order of evaluation is 8654321.

TAC :-

$$\begin{aligned}
 t_8 &= d + e \\
 t_6 &= a + b \\
 t_5 &= t_6 - c \\
 t_4 &= t_5 * t_8
 \end{aligned}$$

$$\begin{aligned}
 t_3 &= t_4 - e \\
 t_2 &= t_6 + t_4 \\
 t_1 &= t_2 * t_3
 \end{aligned}$$

Optimal Ordering for trees:-

- It is a simple algorithm to determine the optimal order in which to evaluate statements in a basic block when the dag representation of the block is atree. It gives the shortest instruction sequence.
- It has two parts:-

1) The Labeling Algorithm:-

- Labeling is done in bottom-up order.
- left leaf:- a node that is a leaf and leftmost descendant of its parent.
- Right leaves:- other leaves

Algo:- Label computation

```

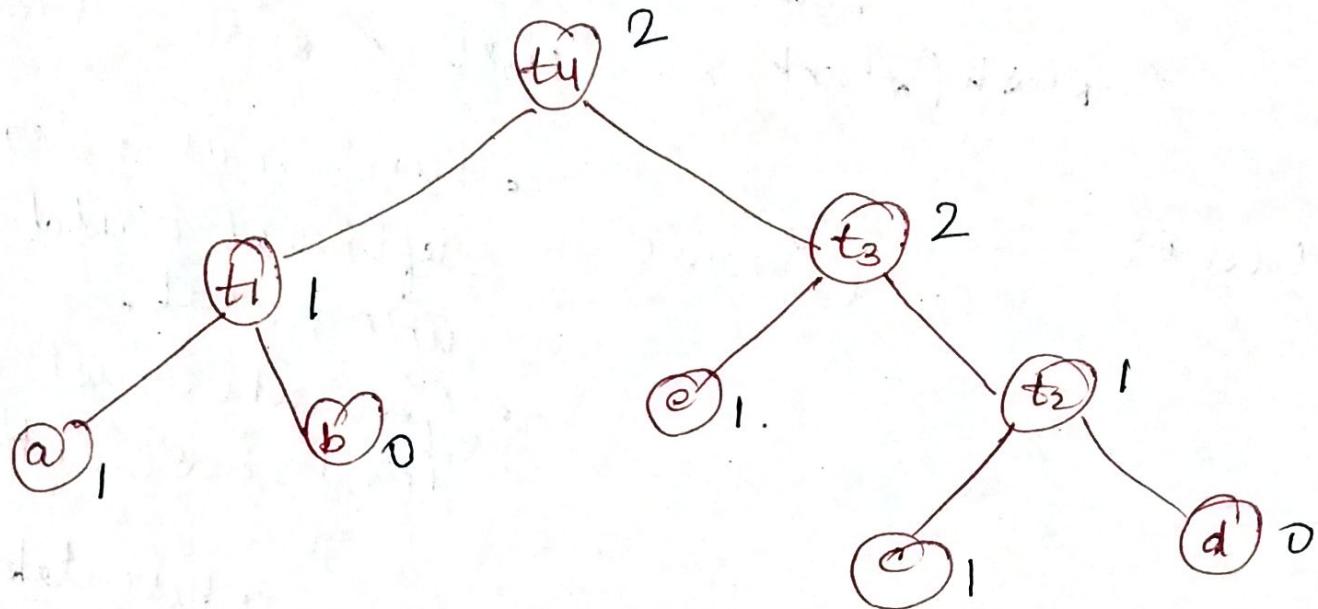
(1) if n is a leaf then
    label(n)=1
(2) if n is the leftmost child of its
    parent then
        label(n)=1
(3) else label(n)=0
(4) else begin /* n is an interior node */
        let  $n_1, n_2, \dots, n_k$  be the children of
        n ordered by label
        so  $\text{label}(n_1) \geq \text{label}(n_2) \geq \dots \geq \text{label}(n_k)$ 
(5)  $\text{label}(n) = \max_i (\text{label}(n_i) + i - 1)$ 
(6) end

```

- If n is a binary node then its children have labels l_1 and l_2 , the formula of line (6) reduces to,

$$\text{label}(n) = \begin{cases} \max(l_1, l_2) & \text{if } l_1 \neq l_2 \\ l_1 + 1 & \text{if } l_1 = l_2 \end{cases}$$

Ex:-



- post order visit = $a \ b \ t_1 \ e \ c \ d \ t_2 \ t_3 \ t_4$
(order for doing labeling)

& Code Generation from a labeled Tree:-

- Takes as input a labeled tree T and produces as output a machine code sequence that evaluates T in R_0 .
- Procedures used by the algorithm,
 - * $\text{genode}(n)$ - to produce machine code evaluating the subtree of T with root n into a register.

- * stack - $\text{stack}(\text{To allocate registers})$
available register initially $R_0, R_1 \dots R_{(3t-1)}$ (169)
- * swap(stack) - Interchanges the top two registers on stack.
- * tstack - temporary memory locations.
- * $x = \text{pop}(\text{stack})$ - pop stack and assign value to x .
- * $\text{push}(x\text{stack}, x)$ - push x onto stack.

case 0:-



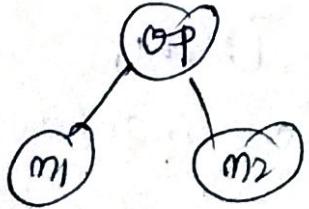
- Here m is a leaf & leftmost child of its parent.
- Generate just a load instruction

case 1:-



into register $R = \text{top}(\text{tstack})$, followed by
instruction $\text{OP name, } R$.

case 2:-



- here label(m₂) > label(m₁)
Means m_2 is harder to evaluate.
- Swap top two registers on stack, evaluate m_2 into $R = \text{top}(\text{tstack})$. Now remove R from stack & evaluate m_1 into $s = \text{top}(\text{tstack})$.

R from stack

& evaluate m_1 into $s = \text{top}(\text{tstack})$

- Then generate the instruction
OP, R, S.

case 3:- similar to case 2, except that here left tree is harder to evaluate. So it is evaluated first. ($\text{label}(m_2) \leq \text{label}(m_1)$)

case 4:- when both subtrees require n (total no of registers) or more registers to evaluate without stores.

- Use a temp location to evaluate right subtree, then left subtree and finally the root.

Ex:- Generate code for labeled tree in previous example,

genocode(t ₄)	[R ₄ R ₀]
genocode(t ₃)	[R ₀ R ₁]
genocode(e)	[R ₀ R ₁]
	print MOV e, R ₁
genocode(t ₂)	[R ₀]
	print MOV c, R ₀
	print ADD d, R ₀
genocode(c)	[R ₀]
	print MOV c, R ₀
	print ADD d, R ₀
	print SUB R ₀ , R ₁
genocode(t ₁)	[R ₀]

// case 2	$t_3 - 2$
// case 3	$e + t_2 - 1$
// case 0	
// case 1	$c - 1$
// case 0	$d - 0$
// case 0	
// case 1	$a - 1$
// case 0	$b - 0$

gen code(a) [R0] // Case 0
171

print NOV a, R₀

print ADD b, R₀

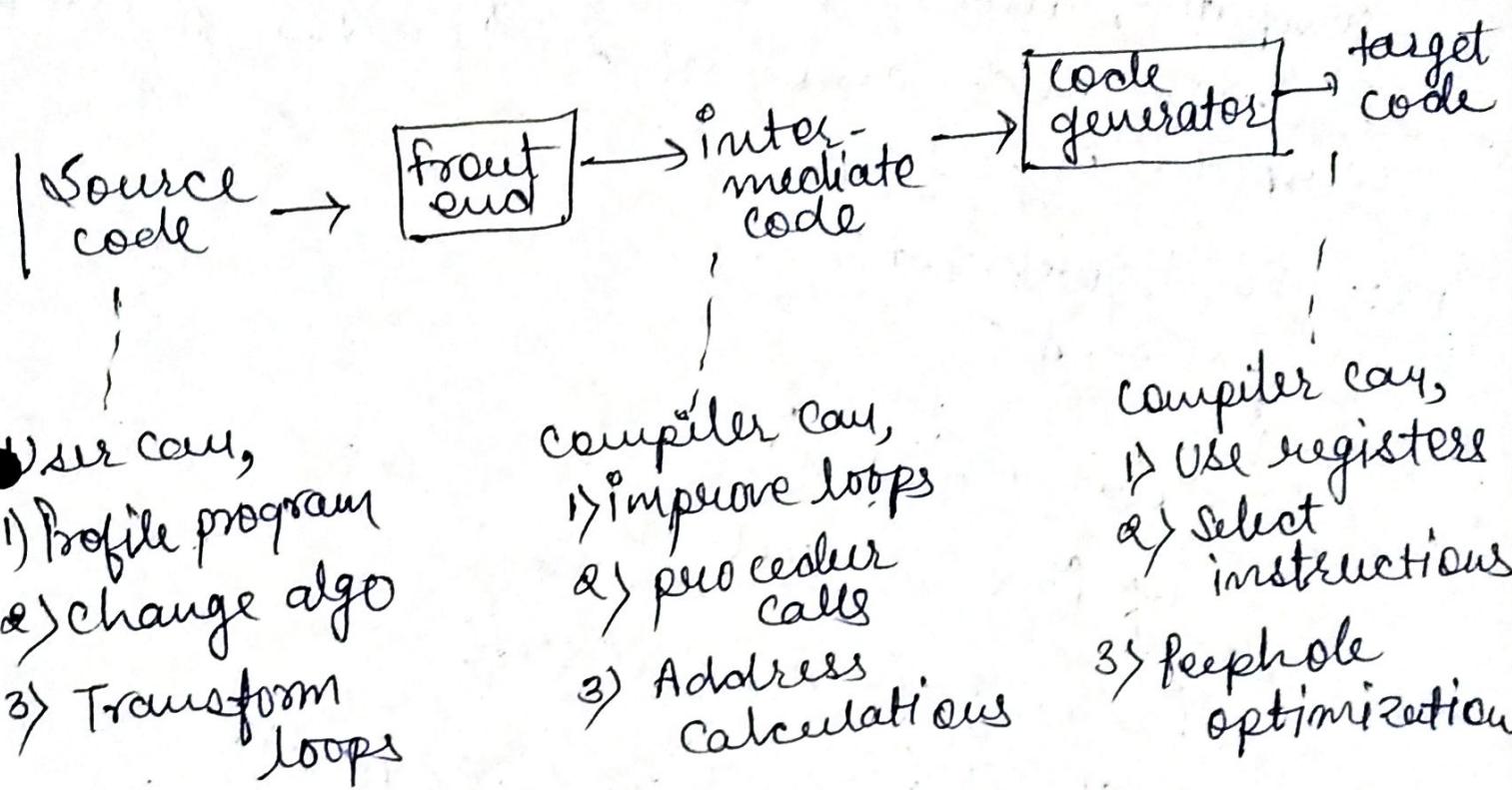
print SUB R₁, R₀

common Subexpressions:

- When there are common subexpressions in a basic block, the corresponding dag will no longer be a tree. The common subexpressions will correspond to nodes with more than one parent, called shared nodes. We can no longer apply the labelling algorithm or gencode directly.
- In practice, we can obtain a reasonable solution if we partition the dag into a set of trees. ~~if~~

Code Optimization:-

- Compilers that apply code-improving transformations are called optimizing compilers.
 - ↳ Machine-independent optimization (program transformations that improve the target code without taking into consideration any properties of target machine)
 - ↳ Machine-dependent optimization (improves register allocation & utilization of special machine-instruction sequence.)



The principle sources of optimization:-

A transformation of a program is called local if it can be performed by looking only at the stmts in a basic block; otherwise it is called global.

Function-Preserving Transformations:-

- By using various ways the compiler can improve the program without changing the properties of the function it computes.

1) Common Subexpressions:

(173)

- An occurrence of an expression E is called a common subexpression if E was previously computed, and the values of variables in E have not changed since the previous computation.

Ex:- void quicksort(m,n)

int m, n;

{

int i, j, v, x;

if ($m <= n$) return;

i = m - 1; 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;

$x = a[i]; a[i] = a[j]; a[j] = x;$

}

$x = a[i]; a[i] = a[n]; a[n] = x;$

quicksort(m, j); quicksort(i + 1, n);

}

TAC for above code:

(174)

$$1) i = m - 1$$

$$2) j = n$$

$$3) t_1 = 4 * n$$

$$4) v = a[t_1]$$

$$5) i = i + 1$$

$$6) t_2 = 4 * i$$

$$7) t_3 = a[t_2]$$

$$8) \text{if } t_3 < v \text{ goto } 5$$

$$9) j = j - 1$$

$$10) t_4 = 4 * j$$

$$11) t_5 = a[t_4]$$

$$12) \text{if } t_5 > v \text{ goto } 9$$

$$13) \text{if } i >= j \text{ goto } 23$$

$$14) t_6 = 4 * i$$

$$15) x = a[t_6]$$

$$16) t_7 = 4 * i$$

$$17) t_8 = 4 * j$$

$$18) t_9 = a[t_8]$$

$$19) a[t_7] = t_9$$

$$20) t_{10} = 4 * j$$

$$21) a[t_{10}] = x$$

$$22) \text{goto } 5$$

$$23) t_{11} = 4 * i$$

$$24) x = a[t_{11}]$$

$$25) t_{12} = 4 * i$$

$$26) t_{13} = 4 * n$$

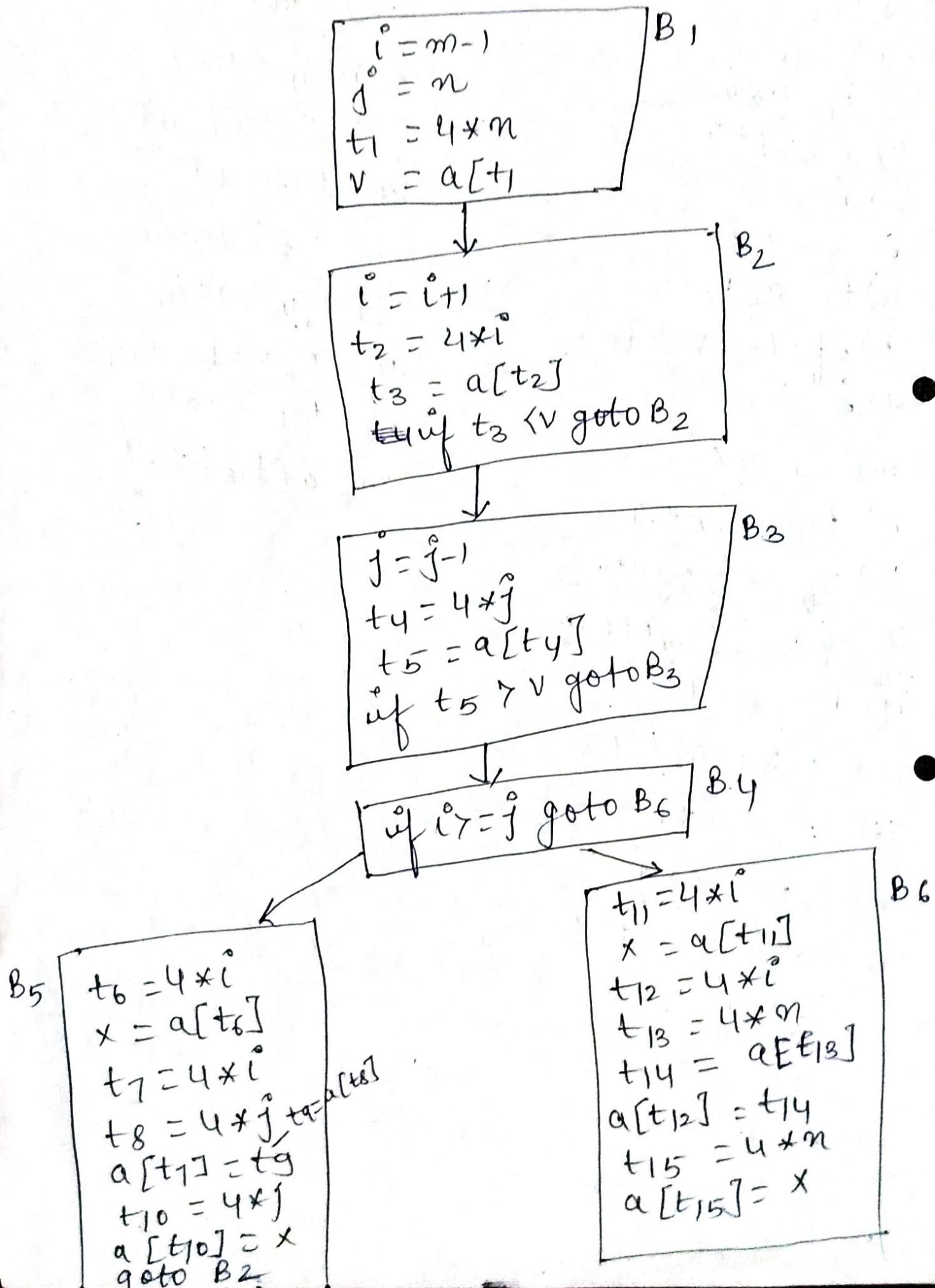
$$27) t_{14} = a[t_{13}]$$

$$28) a[t_{12}] = t_{14}$$

$$29) t_{15} = 4 * n$$

$$30) a[t_{15}] = x$$

Flow graph:-



H6

B5 after local common subexpression elimination,

B5

$t_6 = 4 * i$
$x = a[t_6]$
$t_8 = 4 * j$
$t_9 = a[t_8]$
$a[t_6] = t_9$
$a[t_8] = x$
goto B ₂

Here t_6 is used in place of t_7 and t_8 in place of t_{10} .

B₅ and B₆ after global and local CSE,

B₅

$x = t_3$
$a[t_2] = t_5$
$a[t_0] = x$
goto B ₂

B₆

$x = t_3$
$t_{14} = a[t_2]$
$a[t_2] = t_{14}$
$a[t_3] = x$

2) Dead-code Elimination:-

- A variable is live at a point in a program if its value can be used subsequently; otherwise, it is dead at that point.

$k = 0$
 if ($k = 1$) } dead code
 { $a = a + 5;$
 }
 ?

3) Copy Propagation:-

- Assignment statement $f = g$ called copy statements or copies for short.
- The idea behind the copy-propagation transformation is to use g for f , whenever possible after the copy stmt $f = g$.

Ex:- B_5 block in previous ex

read code now

$x = t_3$
 $a[t_2] = t_5$
 $a[t_4] = x$
 goto B_2



$x = t_3$
 $a[t_2] = t_5$
 $a[t_4] = t_3$
 goto B_2

Now we can eliminate the assignment to x .

4) Loop Optimization:-

loop is a very imp place for optimization.

- loop is a very imp place for optimization.
- a) Code Motion: it takes an exp that yields the same result independent of the number of times a loop is executed and places the exp before the loop.

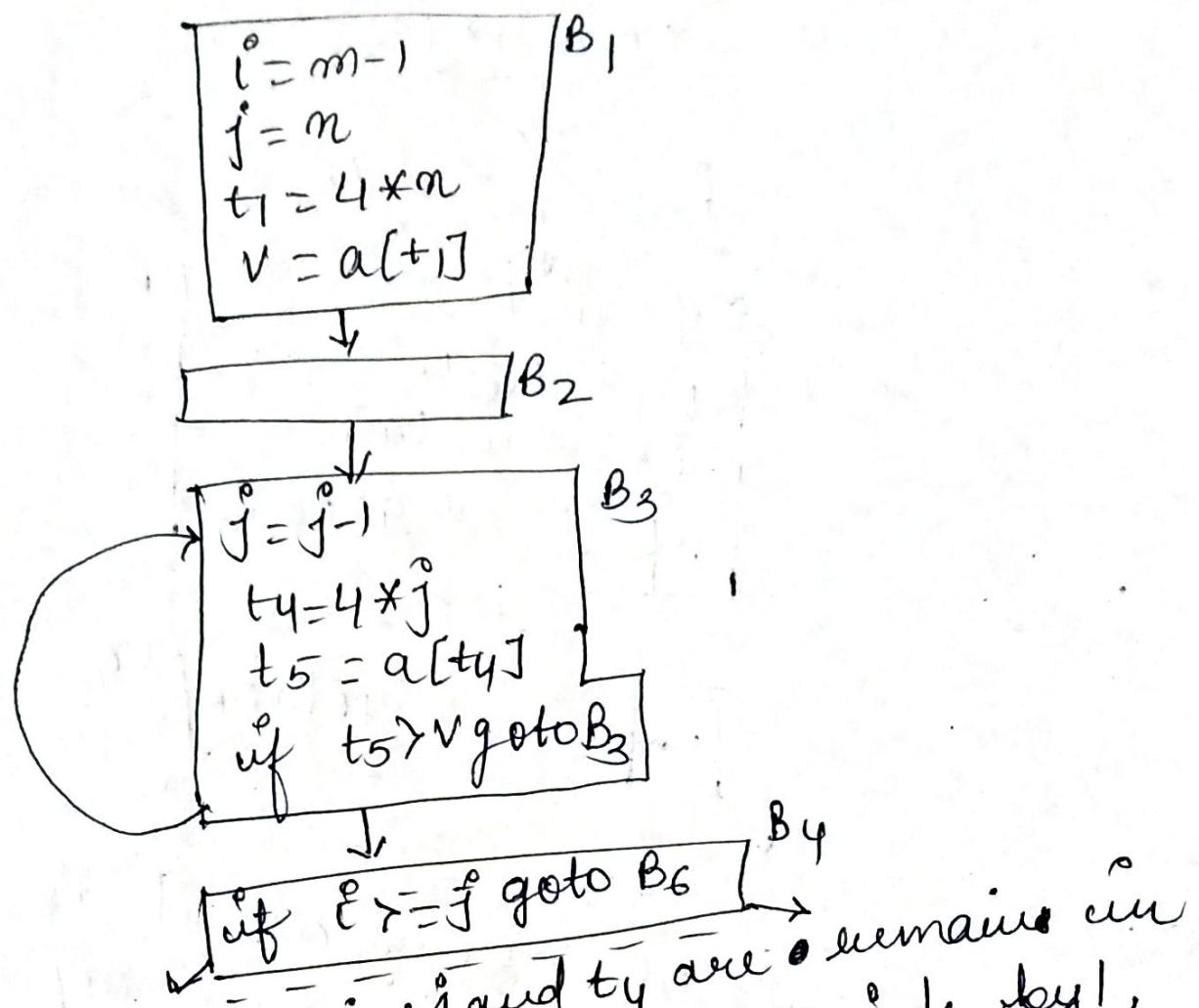
Ex:- while ($i \leq \text{limit} - 2$)

\downarrow
 $t = \text{limit} - 2;$
 while ($i \leq t$)

b) Induction Variables and Reduction
in strength:

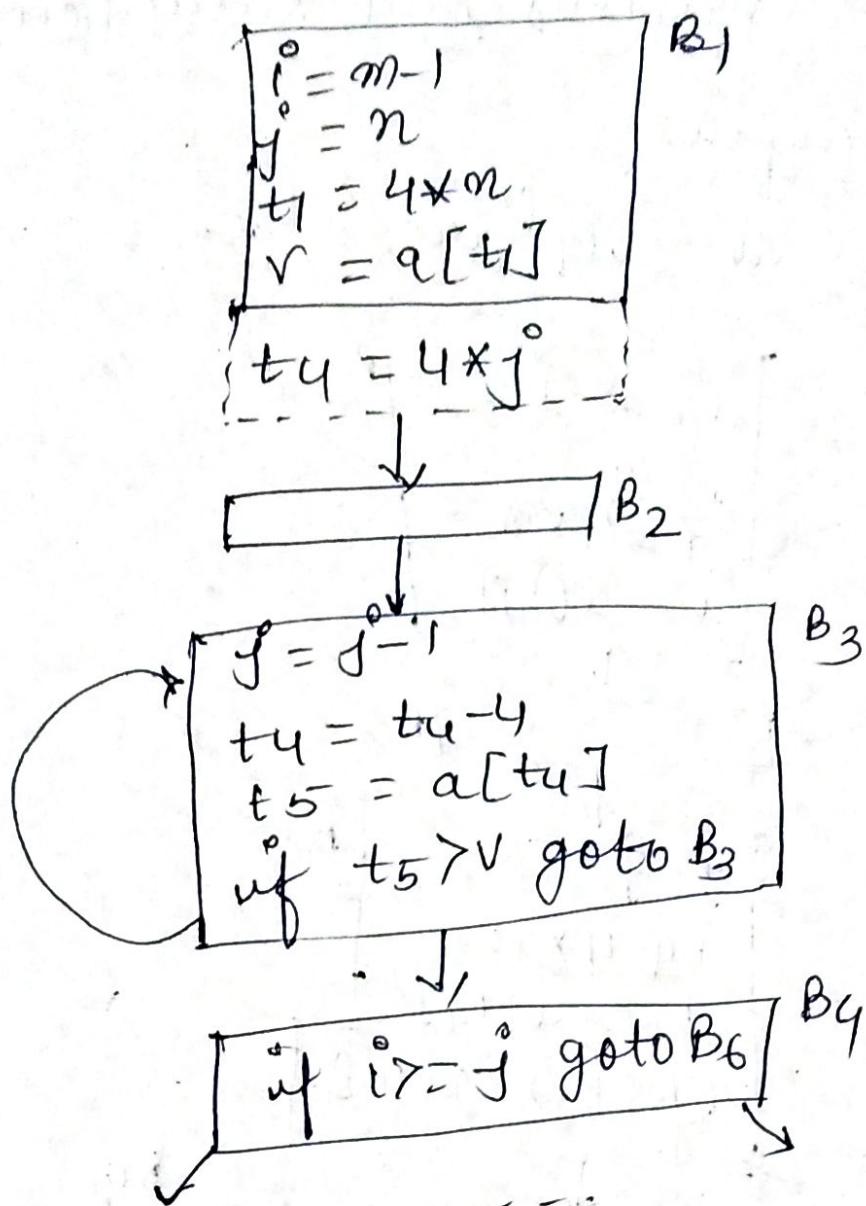
(128)

Ex:- Consider the loop around B_3 ,



- In this, block B_3 contains j and t_4 are remain in lock-step; every time the value of $j \downarrow$ by 1, t_4 decreases by 4. Such identifiers are called induction variables.
- Now we can replace $t_4 = t_4 - 4$. But here we have to initialize t_4 too, otherwise it is not having a initial value, we place an initialization of t_4 at the end of block where j itself is initialized. (In Block B_1)

(79)



Here we can't remove j , as it is using in block B4. But in t_4 we have replaced a multiplication with subtraction. (Mul takes more time than sub)

Loops in Flow Graphs -

Now i and j are only in use after B4. We know that the values of i and t_2 satisfy the relationship $t_2 = 4 \times i$. Same as with j and t_4 ($t_4 = 4 \times j$).

So the test $t_2 > t_4$ is equivalent
to $i > j$ in B₄. After eliminating
these induction-variable,

(180)

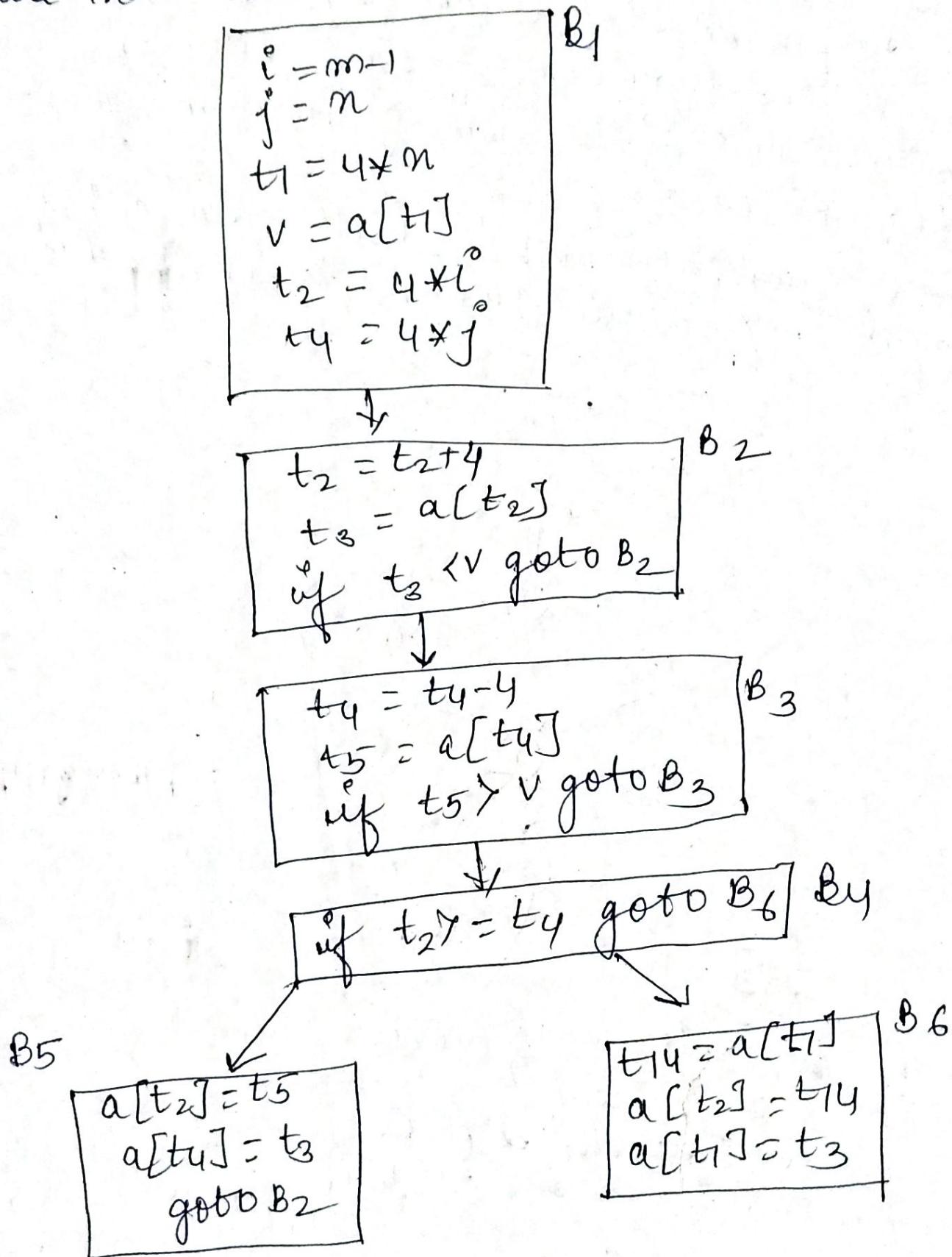


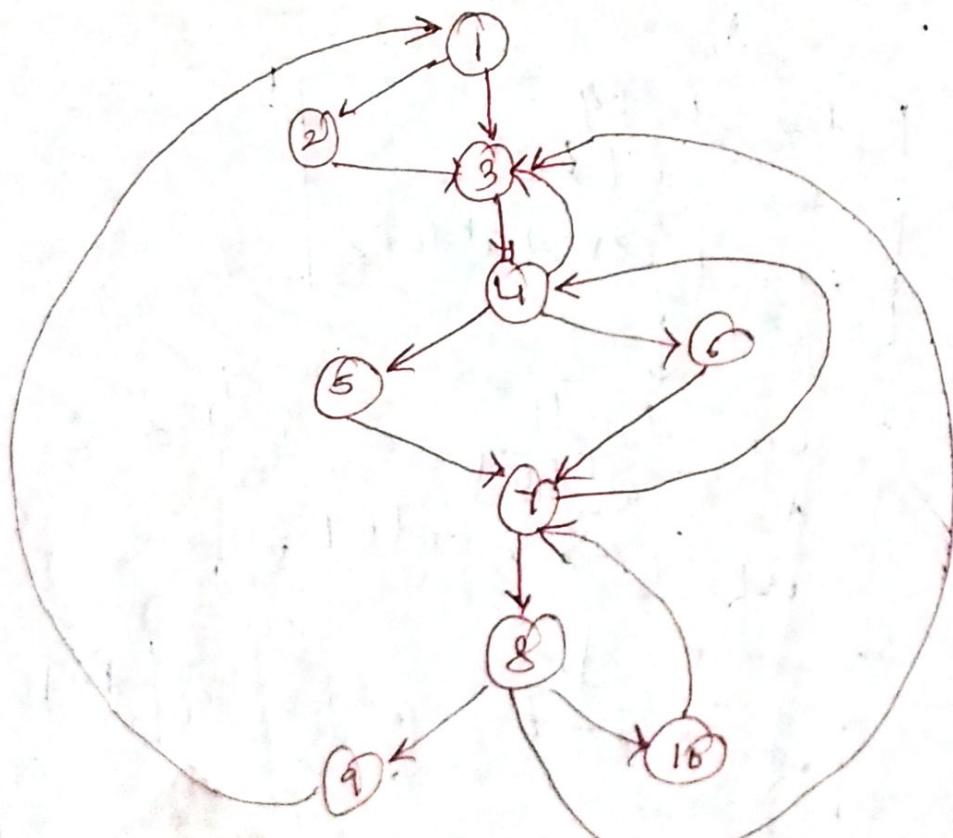
Fig: final optimized TAC

Loops in Flow Graphs:-

(181)

Dominators:-

- we say node d of a flow graph dominates node n, written $d \text{ dom } n$, if every path from the initial node of the flow graph to n goes through d.
- every node dominates itself and entry of a loop dominates all nodes in the loop.



1 - dominates every node

2 - " itself

3 - dominates all, except 1 & 2
" 1, 2, 3

4 - "

5, 6, 9, 10 - itself

7 - 7, 8, 9, 10

8 - 8, 9, 10

Natural Loops:-

- Two essential properties,
 - * A loop must have a single entry point called the header.
 - * There must be at least one way to iterate the loop, i.e. at least one path back to the header.

Ex:- In previous example there is an edge

$7 \rightarrow 4$ and $4 \text{ dom } 7$.

• Similarly $10 \rightarrow 7$ & $7 \text{ dom } 10$

$4 \rightarrow 3$ & $3 \text{ dom } 4$

$8 \rightarrow 3$ & $3 \text{ dom } 8$

$9 \rightarrow 1$ & $1 \text{ dom } 9$

} natural loops

- Given a back edge $m \rightarrow d$, we define the natural loop of the edge to be d plus the set of nodes that can reach m without going through d .

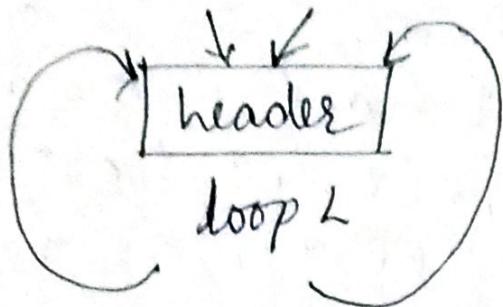
Inner loops:-

- A loop which contains no loop within it is called inner loop. In previous ex ~~7, 8, 10~~ is a inner loop

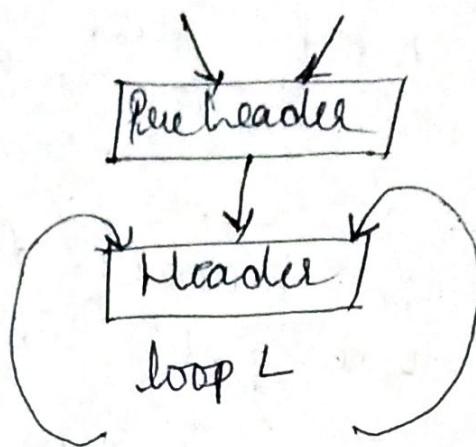
Pre-Headers:-

- When we move many starts outside of the loop (before the header), we create a new

block, called preheader and more stnts
inside this. (183)



a) before



b) after

Reducible Flow Graphs :-

- Structured flow-of-control statements such as if-then-else, while-do, continue and break stnts produces programs whose flow graphs are always reducible.
- A flow graph G is reducible if and only if we can partition the edges into two disjoint groups, often called the forward edges and back edges with following two properties:
1) The forward edges form an acyclic graph in which every node can be reached from the initial node of G .

Q) The back edges consist of only of edges whose heads dominate their tails. (184)

- Ex:- • The graph in previous example is reducible. In general, we can find and remove all back edges. If all remaining forward edges make an acyclic graph, then it is reducible.
• Like if we remove five back edges $4 \rightarrow 3$, $1 \rightarrow 4$, $8 \rightarrow 3$, $9 \rightarrow 1$ and $10 \rightarrow 7$, whose heads dominate their tails, the remaining graph is acyclic. So it is reducible.

loop Invariant Computations:-

- It can be obtained by moving some amount of code outside the loop and placing it just before entering the loop. Code which is moved outside the loop is the expression which is computed independent of the no of times.

Ex:- $\text{while } (k <= l-1)$

{ $\text{sub} = \text{sub} - a[k];$ }

$\text{temp} = l-1;$

$\text{while } (k <= \text{temp})$

{ $\text{sub} = \text{sub} - a[k];$ }

}

Global Data-flow analysis:-

185

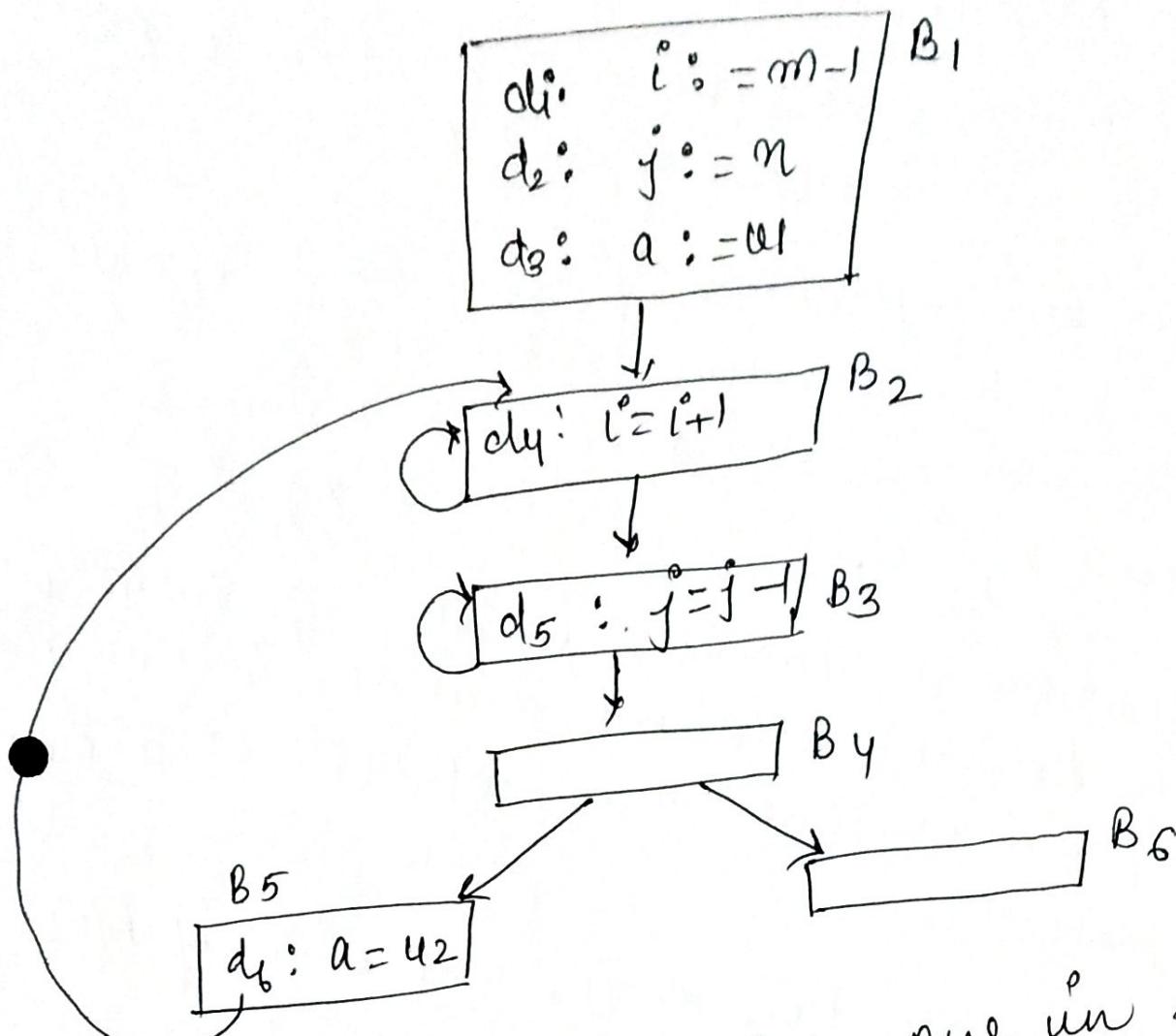
- The process of collecting data-flow info by an optimizing compiler is known as data-flow analysis. This will help in code optimization.
- Ex:-
 - ① To know what variables are live on exit from each block.
 - ② knowledge of global common sub-expressions for elimination.
 - ③ knowledge of unreachable code for dead-code elimination.
- Data-flow information can be collected by setting up and solving systems of equations that relate info at various points in a program.

$$\text{Ex:- } \text{out}[s] = \text{gen}[s] \cup (\text{in}[s] - \text{kill}[s])$$

the info at the end of stmt is either generated within the stmt or enters at the beginning and is not killed as the control flows through the statement. Such equations are called data-flow equations.

Points and Paths:-

- within a basic block, we talk of the point between two adjacent stmts, as well as the point before the first stmt and after the last.



In block B_1 we have 4 points. one in starting and one after each of the three assignments.

- A path from p_1 to p_m is a sequence of points p_1, p_2, \dots, p_m such that for each i b/w 1 & $m-1$ either,
 - p_i is the point immediately preceding a start & p_{i+1} is following the start in same block or
 - p_i is end of some block and p_{i+1} is beginning of a successive block.
- ex:- path $\rightarrow B_5 - B_2 - B_3 - B_4 - B_6$
(from beginning of B_5 to beginning of B_6)

Reaching Definitions:-

(187)

- A definition of a variable x is a stmt that assigns, or may assign, a value to x .
- Unambiguous definition:- if assignment operator or I/O device assign value to x .
- Ambiguous definition:- assignment through pointer ($*q = y$, where q points to x)
- A definition d reaches a point p if there is a path from the point immediately following d to p , such that d is not "killed" along that path.
- we kill a definition of a variable a' if b/w two points along the path there is a definition of a .

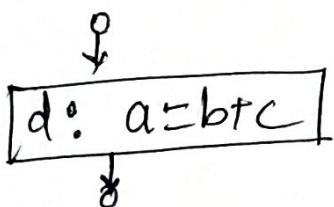
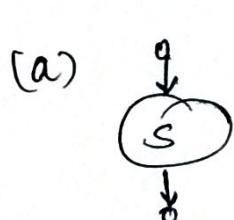
Ex:- In last example, both definitions $i = m - 1$ and $j = m$ in B_1 reach to the beginning of B_2 , but after B_2 $j = j + 1$ will reach to B_5 , B_2 kills the definition $j = m$, so the latter does not reach B_4, B_5 or B_6 .

Data-Flow Analysis of Structured Programs 188

- flow graphs for control-flow constructs have a useful property: there is a single beginning point at which control enters and a single end point that control leaves from when execution of the start is over.

Ex:- $s \rightarrow id := E \mid S; S \mid \text{if } E \text{ then } S \text{ else } S \mid$
 $\text{do } S \text{ while } E$

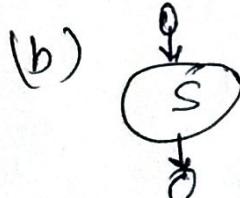
$E \rightarrow id + id \mid id$



$$\begin{aligned} \text{gen}[s] &= \{d\} \\ \text{kill}[s] &= \{a - d\} \\ \text{out}[s] &= \text{gen}[s] \cup (\text{in}[s] - \text{kill}[s]) \end{aligned}$$

d-definition

Da - all definition of a



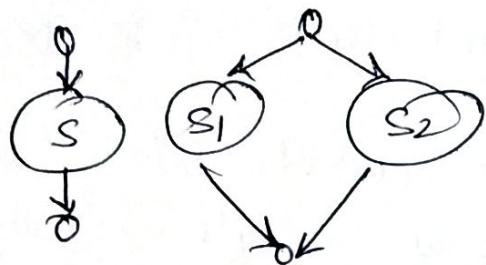
$$\begin{aligned} \text{gen}[s] &= \text{gen}[s_2] \cup (\text{gen}[s_1] - \text{kill}[s_2]) \\ \text{kill}[s] &= \text{kill}[s_2] \cup (\text{kill}[s_1] - \text{gen}[s_2]) \end{aligned}$$

$$\text{in}[s_1] = \text{in}[s]$$

$$\text{in}[s_2] = \text{out}[s_1]$$

$$\text{out}[s] = \text{out}[s_2]$$

(c)



$$\text{gen}[s] = \text{gen}[s_1] \cup \text{gen}[s_2]$$

(189)

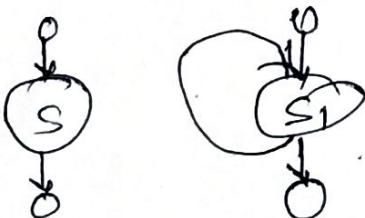
$$\text{kill}[s] = \text{kill}[s_1] \cap \text{kill}[s_2]$$

$$\text{in}[s_1] = \text{in}[s]$$

$$\text{in}[s_2] = \text{in}[s]$$

$$\text{out}[s] = \text{out}[s_1] \cup \text{out}[s_2]$$

(d)



$$\text{gen}[s] = \text{gen}[s_1]$$

$$\text{kill}[s] = \text{kill}[s_1]$$

$$\text{in}[s_1] = \text{in}[s] \cup \text{gen}[s_1]$$

$$\text{out}[s] = \text{out}[s_1]$$

Fig: Data-flow equations for reaching definitions.



आसतो मा राद्गमय

Tutorial Sheet -1

1. Remove left-recursion from following grammars,

(Medium)

$$A \rightarrow ABd / Aa / a$$

$$B \rightarrow Be / b$$

2. Remove left factoring from following grammars,

(Medium)

$$S \rightarrow iEtS / iEtSeS / a$$

$$E \rightarrow b$$

3. Remove left-recursion from following grammars,

(Easy)

$$E \rightarrow E+E$$

$$E \rightarrow E^*E$$

$$E \rightarrow a$$

4. Show that following grammar is LL(1) or not

(Difficult)

$$S \rightarrow EF$$

$$E \rightarrow a \mid \epsilon$$

$$F \rightarrow abF \mid ac$$

Note: Upper case letters are non-terminals and lower case letters are terminals.

Name :- Raghav Bhandari
Subject :- Computer Design
Mam: Mam. Deepa Modi

Section : D G I L
Branch : C.S.E

Tutorial Sheet +
Answers :-

Ans 1

i) $A \rightarrow ABd / Aa / a$

$$A \rightarrow ABd / a$$

$$A \rightarrow aA'$$

$$A' \rightarrow BdA' / \epsilon$$

$$A \rightarrow Aa / a$$

$$A \rightarrow aA'$$

$$A' \rightarrow aA' / \epsilon$$

So After removing left recursion we get.

$$A \rightarrow aA'$$

$$A' \rightarrow BdA' / aA' / \epsilon$$

ii) $B \rightarrow Be / b$

$$B \rightarrow bB'$$

$$B' \rightarrow eB' / \epsilon$$

So After removing left recursion we get

$$B \rightarrow bB'$$

$$B' \rightarrow eB' / \epsilon$$

Aus 2

$$S \rightarrow iEtS \mid iEtSeS/a$$

$$E \rightarrow b$$

After removing left factoring we get

$$S \rightarrow iEtSS'/a$$

$$S' \rightarrow eS/\epsilon$$

$$E \rightarrow b.$$

Aus 3.

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow a$$

$$E \stackrel{\text{or}}{\rightarrow} E + E \mid E * E \mid a$$

After removing left recursion we get.

$$E \rightarrow aE'$$

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

$$E' \rightarrow *EE' \mid \epsilon$$

or

$$E \rightarrow aE'$$

$$E' \rightarrow +EE' \mid *EE' \mid \epsilon$$

Aus 4 LL(1)

First (S) = {a, ε}

First (E) = {a, ε}

First (F) = {a}

First (F') = {b, c}

Follow (S) = {\\$}

Follow (E) = {a}

Follow (F) = {\\$}

Follow (F') = {\\$}

Grammer:-

$$\begin{array}{l} S \rightarrow a F \\ E \rightarrow a \quad | \quad \epsilon \\ F \rightarrow a F' \\ F' \rightarrow b F \quad | \quad c \end{array}$$

-: Parse Table :-

	a	b	c	\$
S	(1)			☒
E	(2)			
F	(3)			
F'	(4)			
		(5)	(6)	

As we get production (2), (3) in one cell of E row and a column so the given grammer does not support LL(1) parsing.



आसतो मा रात्गमय

Swami Keshvanand Institute of Technology, Management & Gramothan,

Ramnagar, Jagatpura, Jaipur-302017, INDIA

Approved by AICTE, Ministry of HRD, Government of India

Recognized by UGC under Section 2(f) of the UGC Act, 1956

Tel. : +91-0141- 5160400 Fax: +91-0141-2759555

E-mail: info@skit.ac.in Web: www.skit.ac.in

Tutorial Sheet -2

1. Remove left-recursion from following grammars, (Easy)

$S \rightarrow (L) / a$

$L \rightarrow L, S / S$

$S \rightarrow A$

2. Remove left factoring from following grammars, (Medium)

$A \rightarrow aAB / aBc / aAc$

3. Remove left-recursion from following grammars, (Medium)

$S \rightarrow A$

$A \rightarrow Ad / Ae / aB / ac$

$B \rightarrow bBC / f$

4. Find the first and follow for the following grammar: (Difficult)

$S \rightarrow ACB | CbB | Ba$

$A \rightarrow da | BC$

$B \rightarrow g | \epsilon$

$C \rightarrow h | \epsilon$

Note: Upper case letters are non-terminals and lower case letters are terminals.

Name - Sakshi Khandelwal

Roll No - 196SKCS812

Tutorial Sheet - 2

① $S \rightarrow (L) / a$

$L \rightarrow L, S / \epsilon$

$S \rightarrow A$

Removing left Recursion:

$S \rightarrow (L) / a$

$L \rightarrow SL'$

$L' \rightarrow ; SL' / \epsilon$

$S \rightarrow A$

② left factoring:

$A \rightarrow aAB / aBc / aAc$

$A \rightarrow aA'$

$A' \rightarrow AB / Bc / Ac$

③

$S \rightarrow A$

$A \rightarrow Ad / Ae / aB / ac$

$B \rightarrow bBC / f$

$S \rightarrow A$

$A \rightarrow Ad / aB$

$A \rightarrow Ae / ac$

$B \rightarrow bBc / f$

Now Removing left Recursion

$$S \rightarrow A$$

$$A \rightarrow aBA'$$

$$A' \rightarrow dA'/\epsilon$$

$$A \rightarrow acA'$$

$$A' \rightarrow eA'/\epsilon$$

$$B \rightarrow bBc/f$$

4) $S \rightarrow ACB/cbB/Ba$

$$A \rightarrow da/Bc$$

$$B \rightarrow g/\epsilon$$

$$c \rightarrow h/\epsilon$$

$$\text{first}(S) = \{d, g, h, b, a, \epsilon\}$$

$$\text{first}(A) = \{d, g, h, \epsilon\}$$

$$\text{first}(B) = \{g, \epsilon\} \quad \text{first}(C) = \{h, \epsilon\}$$

$$\text{follow}(A) = \{h, g, \$\} \quad \text{follow}(S) = \{\$\}$$

$$\text{follow}(B) = \{\$, a, h, g\}$$

$$\text{follow}(C) = \{g, \$, b, h\}$$



Tutorial Sheet -3

1. Write three address code for following programming segments, (Medium)
 - a) $-(a * b) + (c + d) - (a + b + c + d)$
 - b) If $A < B$ and $C < D$ then $t = 1$ else $t = 0$
 - c) $c = 0$
do
{
 if ($a < b$) then
 $x++;$
 else
 $x-;$
 $c++;$
} while ($c < 5$)
2. Translate the arithmetic expression $a + b * c / e \uparrow f + b * c$ into, (Difficult)
 - a) Three address code
 - b) Quadruple
 - c) Triples
 - d) Indirect triples

Name: Shiv Maheshwari
Roll no: 19ESKCS819

Tutorial Sheet

1) Three Address Code:

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

1. $t_1 = a * b$

2. $t_2 = -t_1$

3. $t_3 = c + d$

4. $t_4 = a + b$

5. $t_5 = c + d$

6. $t_6 = t_4 + t_5$

7. $t_7 = -t_6$

8. $t_8 = t_5 + t_7$

9. $t_9 = t_2 + t_8$

2) if $A < B$ and $C < D$ then $t = 1$ else

1. If $A < B$ goto (3)

2. goto (7)

3. If $C < D$ goto (5)

4. goto (7)

5. $t = 1$

6. goto (8)

7. $t = 0$

8. Stop

c) $C = 0$

do

i

if ($a < b$) then
 $x + x$

else

$n - i$

$C + i$

\downarrow while ($C < 5$)

1. $C = 0$

2. if $a < b$ goto (4)

3. goto (7)

4. $t_1 = n + 1$

5. $x = t_1$

6. goto (9)

7. $t_2 = n - 1$

8. $x = t_2$

9. $t_2 = C + 1$

10. $C = t_2$

11. if $C < 5$ goto (2)

12. STOP

$$a + b \times c / e \uparrow f + b \times c$$

a) Three Address Code :

$$t_1 = c \uparrow f$$

$$t_2 = b \times c$$

$$t_3 = t_2 / t_1$$

$$t_4 = t_3 + t_2$$

$$t_5 = a + t_4$$

b) quadruple :-

S.No.	operator	arg1	arg2	Result
0	\uparrow	c	f	t ₁
1	*	b	c	t ₂
2	/	t ₂	t ₁	t ₃
3	+	t ₃	t ₂	t ₄
4	\downarrow	a	t ₄	t ₅

c) triples :-

S.No	operator	arg1	arg2
(0)	\uparrow	c	f
(1)	*	b	c
(2)	/	(1)	(0)
(3)	+	(2)	(1)
(4)	\downarrow	a	(3)

c) Induced Triples:-

Statement		operator	arg1	arg2
(0)	31	(31)	\uparrow	e
(1)	32	(32)	*	b
(2)	33	(33)	/	(32) (31)
(3)	34	(34)	+	(33) (32)
(4)	35	(35)	+	a (34)



Question Bank:

Unit-I

Q1. What is a compiler?

- A. A compiler does a conversion line by line as the program is run
- B. A compiler converts the whole of a higher level program code into machine code in one step
- C. A compiler is a general purpose language providing very efficient execution
- D. All of the Above

Ans.: B

Q2. Users write the programs in which language?

- A. Low-level Language
- B. High-Level Language
- C. Decimal-Format
- D. Middle-Level Language

Ans: B

Q3. Which computer program accepts the high-level language and converts it into assembly language?

- A. Interpreter
- B. Linker
- C. Assembler
- D. Compiler

Ans: D

Q4. Does the compiler program translate the whole source code in one step?

- A. No
- B. Depends on the Compiler
- C. Don't Know
- D. Yes

Ans:D

Q5. Which tool is used for grouping of characters in tokens in the compiler?

- A. Parser
- B. Code optimizer
- C. Code generator
- D. Scanner

Ans: D

Q6. What is the linker?

- A. It is always used before the program execution.
- B. It is required to create the load module.
- C. It is the same as the loader
- D. None of the above

Ans: B



Unit-II

Q1. What are the stages in the compilation process?

- A. Feasibility study, system design, and testing
- B. Implementation and documentation
- C. Analysis Phase, Synthesis Phase
- D. None of These

Ans.: C

Q2. What is the definition of an interpreter?

- A. An interpreter does the conversion line by line as the program is run
- B. An interpreter is a representation of the system being designed
- C. An interpreter is a general purpose language providing very efficient execution
- D. All of the Above

Ans.: A

Q3. Symbol table can be used for

- A. Checking type compatibility
- B. Storage allocation
- C. Suppressing duplication of error messages
- D. All of the Above

Ans.: D

Q4. Assembly language

- A. is usually the primary user interface
- B. requires fixed format commands
- C. is a mnemonic form of machine language
- D. is quite different from the SCL interpreter

Ans.: C

Q5. Every symbolic reference to a memory operand has to be assembled as

- A. (offset, index base)
- B. (segment base, offset)
- C. (index base, offset)
- D. offset

Ans.: B

Q6. Which translator program converts assembly language program to object program

- A. Assembler
- B. Compiler
- C. Microprocessor
- D. Linker

Ans.: A

Q7. A compiler for a high level language that runs on one machine and produces code for a different machine is called



- A. Optimizing compiler
- B. One pass compiler
- C. Cross compiler
- D. Multipass Compiler

Ans.: C

Q8. The action of passing the source program into the proper syntactic classes is known as

- A. syntax analysis
- B. lexical analysis
- C. interpretation analysis
- D. general syntax analysis

Ans.: B

Q9. Grouping of characters into tokens is done by the

- A. scanner
- B. parser
- C. code generator
- D. code optimizer

Ans.: A

Q10. Lexical analysis phase uses

- A. regular grammar
- B. context free grammar
- C. context sensitive grammar
- D. none of the above

Ans.: A



Unit-III

Q1. Which one of the following is not a syntax error

- A. Semantic
- B. Lexical
- C. Arithmetic
- D. Logical

Ans.: A

Q2. CFG can be recognized by a

- A. Push-down automata
- B. 2-way linear bounded automata
- C. Both (a) and (b)
- D. None of these

Ans. C

Q3. The ambiguous grammar can have

- A. Only one parse tree
- B. More than one parse tree
- C. Parse trees with l-values
- D. None of the above

Ans. B

Q4. Which of the following suffices to convert an arbitrary CFG to an LL(1) grammar

- A. Removing left recursion alone
- B. Factoring the grammar alone
- C. Removing left recursion and factoring the grammar
- D. None of the above

Ans. C

Q5. The 'k' in LR(k) cannot be

- A. 0
- B. 1
- C. 2
- D. None of these

Ans. D

Q6. Non backtracking form of the top-down parser are called

- A. Recursive-descent parsers
- B. Predictive parsers
- C. Shift reduce parsers
- D. None of these

Ans. B

Q7. Parsing table in the LR parsing contains

- A. action, goto



- B. state, action
- C. input, action
- D. state, goto

Ans. A

Q8. The condensed form of the parse tree is called

- A. L-attributed
- B. Inherited attributed
- C. DAG
- D. Syntax tree

Ans. D

Q9. Shift-reduce parsers are

- A. top-down parsers
- B. bottom-up parsers
- C. both (a) and (b)
- D. none of these

Ans. B

Q10. In some programming languages, an identifier is permitted to be a letter followed by any number of letters or digits. If L and D denote the sets of letter and digits respectively, which expression defines an identifier?

- A. (LUD)+
- B. L(LUD)*
- C. (LD)*
- D. L(LD)*

Note: Here U stands for Union

Ans. B

Q11. Backtracking is possible in

- A. LR parsing
- B. Predictive parsing
- C. Recursive descent parsing
- D. None of the above

Ans. C

Q12. Consider the following grammar

expr op expr

expr->id

op-> + | *

Which of the following is true?

- A. op and expr are start symbols
- B. op and id are terminals
- C. expr is start symbol and op is terminal
- D. none of these

Ans. C



Q13. To compute FOLLOW(A) for any grammar symbol A

- A. We must compute FIRST of some grammar symbols
- B. No need of computing FIRST of some symbol
- C. May compute FIRST of some symbol
- D. None of these

Ans. A

Q14. Which language is generated by the given grammar $S \rightarrow 0S1 \mid 01$

- A. 0011
- B. 00*11
- C. 001*1
- D. 00*1*1

Ans. D

Q15. The space consuming but easy parsing is

- A. LALR
- B. SLR
- C. LR
- D. Predictive parser

Ans. A

Q16. A top down parser generates

- A. left-most derivation
- B. right-most derivation
- C. right-most derivation in reverse
- D. left-most derivation in reverse

Ans. A

Q17. Choose the false statement

- A. LL(k) grammar has to be a CFG
- B. LL(k) grammar has to be unambiguous
- C. There are LL(k) grammars that are not Context Free
- D. LL(k) grammars cannot have left recursive non-terminals

Ans. C

Q18. If a grammar is unambiguous then it is surely be

- A. regular
- B. LL(1)
- C. Both (a) and (b)
- D. Cannot say

Ans. D

Q19. Predictive parsing is a special case of

- A. top down parsing



- B. bottom up parsing
- C. recursive descent parsing
- D. none of the above

Ans. C



Unit-IV

Q1. Type checking is normally done during

- A. lexical analysis
- B. syntax analysis
- C. syntax directed translation
- D. code optimization

Ans.: C

Q2. An intermediate code form is

- A. Postfix notation
- B. Syntax trees
- C. Three address code
- D. All of these

Ans.: D

Q3. Semantic errors can be detected

- A. at compile time only
- B. at run time only
- C. both at compile and run time
- D. none of these

Ans.: C

Q4. Undeclared name is error

- A. syntax
- B. lexical
- C. semantic
- D. not an error

Ans. C

Q5. Intermediate code generator is used between

- A. symbol table and code generator
- B. symbol table and code optimizer
- C. code optimizer and code generator
- D. semantic analyzer and code optimizer

Ans.: D

Q6. ‘Divide by 0’ is a

- A. lexical error
- B. syntactic error
- C. semantic error
- D. internal error

Ans. C

Q7. In which of the following has attribute values at each node

- A. Associated parse trees



- B. Postfix parse tree
- C. Annotated parse tree
- D. Prefix parse tree

Ans.: C

Q8. The post fix form of $A-B/(C*D$E)$ is

- A. ABCDE\$*/-
- B. AB/C*D\$
- C. ABCDE\$-/*
- D. ABCDE/-*\$

Ans. A

Q9. Synthesized attribute can easily be simulated by an

- A. LL grammar
- B. Ambiguous grammar
- C. LR grammar
- D. None of the above

Ans. C

Q10. The prefix form of $(A+B)^*(C-D)$ is

- A. +-AB*C-D
- B. *+-ABCD
- C. *+AB-CD
- D. *AB+CD

Ans. C

Q11. In a syntax directed translation scheme, if the value of an attribute of a node is a function of the values of the attributes of its children, then it is called a

- A. Synthesized attribute
- B. Inherited attribute
- C. Canonical attribute
- D. None of the above

Ans. A

Q12. Which of the following is not an intermediate code form?

- A. Postfix notation
- B. Syntax trees
- C. Three address code
- D. Quadruples

Ans. D

Q13. Three address codes can be implemented by

- A. indirect triples
- B. direct triples
- C. both (a) and (b)



D. none of the above

Ans. A

Q14. In a bottom up evaluation of a syntax directed definition, inherited attributes can be

- A. always be evaluated
- B. be evaluated only if the definition is L-attributed
- C. be evaluated only if the definition has synthesized attributes
- D. none of the above

Ans. C

Q15. Three address code involves

- A. at the most 3 address
- B. exactly 3 address
- C. no unary operators
- D. none of the above

Ans. A

Q16. Inherited attribute is a natural choice in

- A. keeping track of variable declaration
- B. checking of the correct use of L-values and R-values
- C. both (a) and (b)
- D. none of the above

Ans. C



Unit-V

Q1. The idea of cache memory is based _____

- A. On the property of locality of reference
- B. On the heuristic 90-10 rule
- C. On the fact that references generally tend to cluster
- D. All of the mentioned

Ans. A

Q2. Which of the following known as the text part of a program that does not change at runtime. Its memory requirements are known at the compile time?

- A. Code
- B. Procedures
- C. Variables
- D. All of the above

Ans: A

Q3. A procedure has a start and an end delimiter and everything inside it is called the body of the procedure.

- A. TRUE
- B. FALSE
- C. A procedure has a start with delimiter but not end with delimiter.
- D. A procedure has a not start with delimiter but end with delimiter.

Ans: A

Q4. In activation record, Which of the following Stores the address of activation record of the caller procedure?

- A. Access Link
- B. Actual Parameters
- C. Control Link
- D. Temporaries

Ans: C

Q5. Whenever a procedure is executed, its activation record is stored on the stack, also known as?

- A. Access Stack
- B. Control stack
- C. Formal Stack
- D. Return Stack

Ans : B

Q6. _____ are known at the runtime only, unless they are global or constant.

- A. Values
- B. Object
- C. Variables
- D. All of the above

Ans : C



Q7.The location of memory (address) where an expression is stored is known?

- A. r-value
- B. k-value
- C. l-value
- D. t-value

Ans : C

Q8.What is true about Formal Parameters?

- A. These variables are declared in the definition of the called function.
- B. These variables are specified in the function call as arguments.
- C. Variables whose values or addresses are being passed to the called procedure are called Formal Parameter.
- D. All of the above

Ans: A

Q9.In which mechanism, the calling procedure passes the r-value of actual parameters and the compiler puts that into the called procedure's activation record?

- A. Pass by Reference
- B. Pass by Name
- C. Pass by Copy-restore
- D. Pass by Value

Ans: D

Q10.In which mechanism, the name of the procedure being called is replaced by its actual body?

- A. Pass by Reference
- B. Pass by Name
- C. Pass by Copy-restore
- D. Pass by Object

Ans: B

Q11.What will be error?

$7 = x + y;$

- A. l-value error
- B. r-value error
- C. Infinite loop
- D. Both A and B

Ans: A



Unit-VI

Q1. A basic block can be analyzed by

- A. A DAG
- B. A graph which may involve the cycles
- C. Flow Graph
- D. None of These

Ans.: A

Q2. Which of the following is true for the flow of control among procedures during execution of program?

- A. Control flows randomly
- B. Control flows line by line without jumping
- C. Control flows sequentially
- D. None of these

Ans. C

Q3. The identification of common sub-expression and replacement of run time computations by compile-time computations is-----

- A. Local optimization
- B. Constant folding
- C. Loop optimization
- D. Data flow analysis

Ans: B

Q4. Which of the following class of statement usually produces no executable code when compiled?

- A. Declaration
- B. Assignment statements
- C. Input and output statements
- D. Structural statements

Ans: A

Q5. Code optimization is responsibility of-----

- A. Application programmer
- B. System programmer
- C. Operating system
- D. All of the above.

Ans: B

Q6. Dead-code elimination in machine code optimization refers to -----

- A. Removal of all labels.
- B. Removal of values that never get used.
- C. Removal of function which are not involved.
- D. Removal of a module after its use.

Ans: B



Q.7. Which of the following statement is false?

- A. Flow graph is used to represent DAG.
- B. Three address code is the input to the code generator.
- C. The first statement of three address code is always leader of the first basic block.
Transformation of block is needed for code optimization.

Ans: A

Q.8. ----- is the final phase of compiler.

- A. Semantic analysis
- B. Code generation
- C. Target code generation
- D. Syntax analysis

Ans: B

Q.9. Some code optimizations are carried out on the intermediate code because-----

- A. They enhance the portability of the compiler to other target processors
- B. Program analysis is more accurate on intermediate code than on machine code
- C. The information from dataflow analysis cannot otherwise be used for optimization
- D. The information from the front end cannot otherwise be used for optimization

Ans: A

Q.10 Peephole optimization is form of-----

- A. Loop optimization
- B. Local optimization
- C. Constant folding
- D. Data flow analysis

Ans: B



List of Text and Reference Books

Text Books:

1. Aho, Alfred V., Ravi Sethi, and Jeffrey D. Ullman. "Compilers, principles, techniques." Addison wesley 7.8 (1986): 9.
2. Ullman, Jeffrey D., and Alfred V. Aho. "Principles of compiler design." Reading: Addison Wesley (1977).

Reference Books:

1. Srikant, Y. N., and Priti Shankar, eds. The compiler design handbook: optimizations and machine code generation. CRC Press, 2018.
2. Mogensen, TorbenÆgidius. Introduction to compiler design. Springer, 2017.
3. Grune, Dick, et al. Modern compiler design. Springer Science & Business Media, 2012.
4. Aho, Alfred V. Compilers: principles, techniques and tools (for Anna University), 2/e. Pearson Education India, 2003.
5. Louden, Kenneth C. "Compiler construction." Cengage Learning (1997).
6. Wilhelm, Reinhard, Dieter Maurer, and Stephen S. Wilson. Compiler design. Reading: Addison-Wesley Publishing Company, 1995.