



# PIE Tech

**POLLACHI INSTITUTE OF ENGINEERING AND TECHNOLOGY**

(Approved by **AICTE** and Affiliated to **Anna University**)

*sky is the limit*

**Department of Computer Science and Engineering**

**Regulation 2021**

**III Year – V Semester**

**CS3501 – Compiler Design**

## **CS3501 - COMPILER DESIGN**

**2021-Regulation Anna University**

**CS3501**

**COMPILER DESIGN**

**LPTC**

**3 0 24**

### **OBJECTIVES:**

- To learn the various phases of compiler.
- To learn the various parsing techniques.
- To understand intermediate code generation and run-time environment.
- To learn to implement front-end of the compiler.
- To learn to implement code generator.
- To learn to implement code optimization.

### **UNIT I INTRODUCTION TO COMPILERS & LEXICAL ANALYSIS 8**

Introduction- Translators- Compilation and Interpretation- Language processors -The Phases of Compiler – Lexical Analysis – Role of Lexical Analyzer – Input Buffering – Specification of Tokens – Recognition of Tokens – Finite Automata – Regular Expressions to Automata NFA, DFA – Minimizing DFA - Language for Specifying Lexical Analyzers – Lex tool.

### **UNIT II SYNTAX ANALYSIS 11**

Role of Parser – Grammars – Context-free grammars – Writing a grammar Top Down Parsing - General Strategies - Recursive Descent Parser Predictive Parser-LL(1) - Parser-Shift Reduce Parser-LR Parser- LR (0)Item Construction of SLR Parsing Table - Introduction to LALR Parser - Error Handling and Recovery in Syntax Analyzer-YACC tool - Design of a syntax Analyzer for a Sample Language

### **UNIT III SYNTAX DIRECTED TRANSLATION & INTERMEDIATE CODE GENERATION 9**

Syntax directed Definitions-Construction of Syntax Tree-Bottom-up Evaluation of S-Attribute Definitions- Design of predictive translator - Type Systems-Specification of a simple type Checker Equivalence of Type Expressions-Type Conversions. Intermediate Languages: Syntax Tree, Three Address Code, Types and Declarations, Translation of Expressions, Type Checking, Back patching.

### **UNIT IV RUN-TIME ENVIRONMENT AND CODE GENERATION 9**

Runtime Environments – source language issues – Storage organization – Storage Allocation Strategies: Static, Stack and Heap allocation - Parameter Passing-Symbol Tables - Dynamic Storage Allocation - Issues in the Design of a code generator – Basic Blocks and Flow graphs - Design of a simple Code Generator - Optimal Code Generation for Expressions– Dynamic Programming Code Generation.

### **UNIT V CODE OPTIMIZATION 8**

Principal Sources of Optimization – Peep-hole optimization - DAG- Optimization of Basic Blocks - Global Data Flow Analysis - Efficient Data Flow Algorithm – Recent trends in Compiler Design.

**45 PERIODS**



### **LIST OF EXPERIMENTS:**

1. Using the LEX tool, Develop a lexical analyzer to recognize a few patterns in C. (Ex. identifiers, constants, comments, operators etc.). Create a symbol table, while recognizing identifiers.
2. Implement a Lexical Analyzer using LEX Tool
3. Generate YACC specification for a few syntactic categories.
  - a. Program to recognize a valid arithmetic expression that uses operator +, -, \* and /.
  - b. Program to recognize a valid variable which starts with a letter followed by any number of letters or digits.
  - c. Program to recognize a valid control structures syntax of C language (For loop, while loop, if-else, if-else-if, switch-case, etc.).
  - d. Implementation of calculator using LEX and YACC
4. Generate three address code for a simple program using LEX and YACC.
5. Implement type checking using Lex and Yacc.
6. Implement simple code optimization techniques (Constant folding, Strength reduction and Algebraic transformation)
7. Implement back-end of the compiler for which the three address code is given as input and the 8086 assembly language code is produced as output.

### **OUTCOMES: On Completion of the course, the students should be able to:**

- Understand the different phases of compiler.
- Design a lexical analyzer for a sample language.
- Apply different parsing algorithms to develop the parsers for a given grammar.
- Understand syntax-directed translation and run-time environment.
- Learn to implement code optimization techniques and a simple code generator.
- Design and implement a scanner and a parser using LEX and YACC tools.

### **TEXT BOOK:**

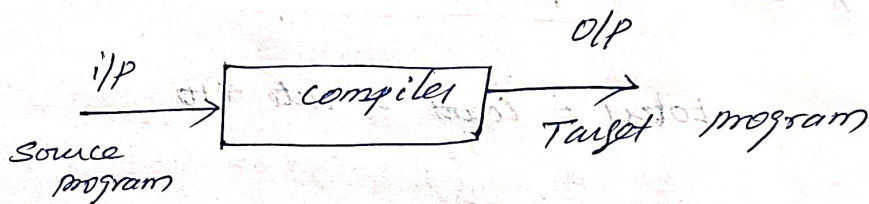
1. Alfred V. Aho, Monica S. Lam, Ravi Sethi, Jeffrey D. Ullman, Compilers: Principles, Techniques and Tools, Second Edition, Pearson Education, 2009.

### **REFERENCES:**

1. Randy Allen, Ken Kennedy, Optimizing Compilers for Modern Architectures: A Dependence based Approach, Morgan Kaufmann Publishers, 2002.
2. Steven S. Muchnick, Advanced Compiler Design and Implementation, Morgan Kaufmann Publishers - Elsevier Science, India, Indian Reprint 2003.
3. Keith D Cooper and Linda Torczon, Engineering a Compiler, Morgan Kaufmann Publishers Elsevier Science, 2004.
4. V. Raghavan, Principles of Compiler Design, Tata McGraw Hill Education Publishers, 2010.
5. Allen I. Holub, Compiler Design in C, Prentice-Hall Software Series, 1993.

## Structure of a Compiler

Compiler is a program which takes one language (source program) as input and translates it into an equivalent another language (target program).



→ during this process of translation if some errors are encountered then, compiler displays them as error message.

→ Compiler takes a source program as higher level language such as C, PASCAL, converts it into low level language or machine level language.

## Analysis - synthesis model

- Analysis part the source program is read and broken down into constituent pieces.

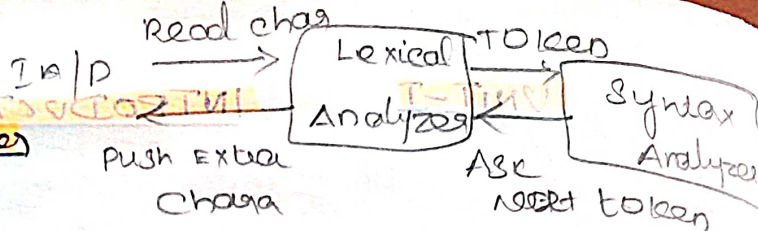
The syntax and the meaning of the source string is determined

→ synthesis part this intermediate form of the source language is taken and converted into an equivalent target program.



2

## Phases of Compiler



### 1. Lexical Analysis

- The Lexical Analysis is also called scanning. It converts the high level program into a seq. of Tokens.
- This phase of compilation is where the complete source code is scanned and your source program is broken up into group of strings called tokens.

$$\text{Total} = \text{Count} + \text{rate} * 10$$

id op id op id op const

1. The identifier total
2. assignment symbol
3. identifier count
4. The plus sign
5. The identifier rate
6. The multiplication sign
7. The constant number 10.

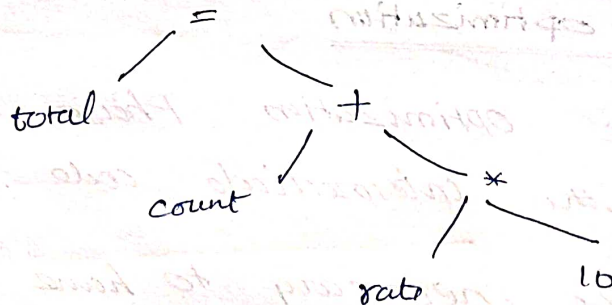
token → I/P

### 2. Syntax Analysis (parser) (syntax tree)

- The syntax analysis is also called parsing.
- The syntax analysis determines the structure of the source string by grouping the tokens together.

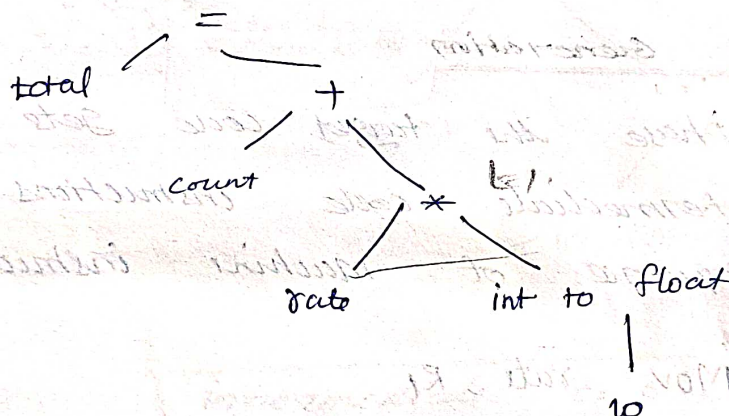
- expression  $\text{Total} = \text{Count} + \text{rate} * 10$

$\langle id_1 \rangle = \langle id_2 \rangle \langle + \rangle \langle id_3 \rangle \langle * \rangle \langle id_4 \rangle$



### 3. Semantic Analysis (Type checking)

- The Semantic Analysis determines the meaning of the source string.



### 4. Intermediate code Generation

- It is a kind of code which is easy to generate and this code can be easily converted to target code.

- This code is in variety of forms such as three address code, quadruple, triplet.

$t_1 := \text{int to float}(10)$

$t_2 := \text{rate} \times t_1$

$t_3 := \text{count} + t_2$

$\text{total} := t_3$



## 5. code optimization

- The code optimization phase attempts to improve the intermediate code.
- This is necessary to have a faster execution code or less consumption of memory.
- Thus by optimizing the code the overall running time of the target program can be improved.

$$t_i = rate \times 60$$

## 6. code Generation

This phase the target code gets generated. The intermediate code instructions are translated into sequence of machine instructions.

MOV rate, R<sub>1</sub>

MUL #100, R<sub>1</sub>

MOV Count, R<sub>2</sub>

ADD R<sub>2</sub>, R<sub>1</sub>

MOV R<sub>1</sub>, total

Write down the output of each phases  
for the expression

$$C = a + b * 5$$



Lexical analyzer



$\langle id, 1X = Xid, 2X + Xid, 3X * X5 \rangle$



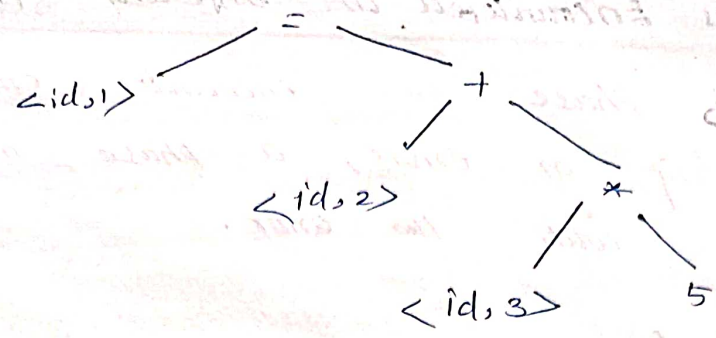
Syntax analyzer



a	...
b	...
c	...
...	...

Symbol  
Table

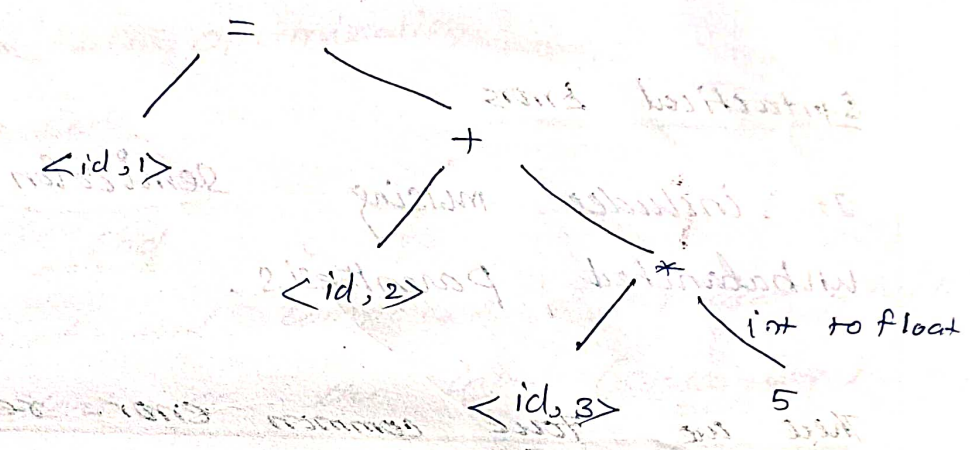
Tokenizer



Syntax Tree

Semantic analyzer

Semantic Tree



intermediate code generator

$t_1 = \text{int of float}(5)$   
 $t_2 = id_3 * t_1$   
 $t_3 = id_2 + t_2$   
 $id_1 = t_3$

Intermed code

Code Optimizer

$t_1 = id_3 * (5)$   
 $id_1 = id_2 * t_1$

optimized code

code generator

STF id<sub>1</sub>, R<sub>1</sub>  
LDF R<sub>2</sub>, id<sub>3</sub>  
MULF R<sub>2</sub>, # 5.0  
LDF R<sub>1</sub>, id<sub>2</sub>  
ADDF R<sub>1</sub>, R<sub>2</sub>

machine code



(6)

## Errors Encountered in Different Phases

Each Phase can encounter errors. After detecting an error, a phase must somehow deal with the error.

### Lexical Errors

It includes incorrect or misspelled name of some identifier.

### Syntactical Errors

It includes missing semi colon or unbalanced parentheses.

### There are four common error-recovery strategies

- Panic mode
- Statement level
- Error productions
- Global correction

### Semantical Errors

- These errors are a result of incompatible value assignment.

- Type mismatch
- undeclared variable
- Reserved identifier misuse
- multiple declaration of variable in a scope

## GROUPING OF PHASES

(7)

### Front end

- Lexical analysis
- syntax analysis
- semantic analysis
- intermediate code generation

### Back end

- code optimization
- code generation

### Front End

→ Front end comprises of phases which are dependent on the input (source lang) and independent on the target machine.

Lexical an

### Passes

→ The phases of compiler can be implemented in a single pass by making the primary action.

- \* reading of input file
- \* writing to the output file

→ Several phases of compiler are grouped into one pass in such a way that the operations in each and every phases are incorporated during the pass.

### Reducing the Number of Passes

→ Minimizing the number of passes improves the time efficiency as reading from and writing to intermediate files can be reduced.



## Compiler Construction Tools

1. Parser generators
2. Scanner generators
3. Syntax-directed translation engines
4. Automatic code generators
5. Data-flow analysis engines
6. Compiler construction toolkits

### Parser generators

inp : Grammatical description of a programming language

oup : Syntax analyzers

Parser generator takes the grammatical description of a programming language and produces a syntax analyzer.

### Scanner Generators

inp : Regular expression description of the tokens of a language.

oup : Lexical analyzers

Scanner generator generates lexical analyzers from a regular expression description of its tokens of a language.

### Syntax-directed Translation Engines

inp : Parse tree

oup : Intermediate code

Syntax-directed translation engines produces collections of routines that walk a parse tree and generates intermediate code.



## Automatic code Generators

i/p : Intermediate language

o/p : Machine language

Code generator takes a collection of rules that define the translation of each operation of the intermediate language into the machine language for a target machine.

## Data Flow Analysis Engines

- Data flow analysis engines gathers the information.
- Data flow analysis is a key part of code optimization.

## Compiler construction Toolkits

The toolkits provide integrated set of routines for various phases of compiler.

## LEXICAL ANALYSIS

- Lexical analysis is the process of converting a sequence of characters from source program into a sequence of tokens.

It consists of two stages.

- Scanning
- Tokenization

## Token, Pattern and Lexeme

### Token

Token is a valid sequence of character which are given by lexeme.

- keywords (if, for, while, ...)
- constant (1, 2, 3, ...)
- identifiers (var, sum, ...)
- Numbers (1, 2, 3, ...)
- operators (+, -, \*, /, ...)
- punctuation symbols (., ,, {, }, ...)



(10)

## Pattern

Pattern describes a rule that must be matched by sequence of characters, to form a token.

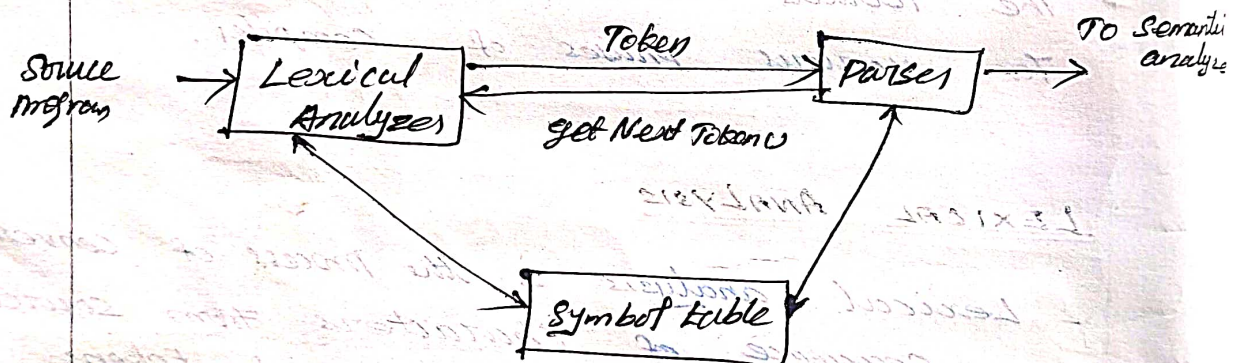
It can be defined by regular expression or grammar rules.

## Lexeme

Lexeme is a sequence of characters that matches the pattern for a token.

identifier  $\uparrow$  id  $\uparrow$   $\uparrow$  const  
 $C = a + b * 5$

## Role of Lexical Analyzer



## Interaction b/w lexical analyzer and parser

- Reads the source program, scans the input characters, group them into lexemes and produce the tokens as output.
- Enters the identified tokens into the symbol table.
- Strips out white spaces and comments from source program.
- Correlates error messages with the source program.
- Expands the macros if it is found in the source program.



## Need of Lexical Analyzer

(11)

- Simplicity of design of compiler
- Compiler efficiency is improved
- Compiler portability is enhanced.

## Issues in Lexical Analysis

### - Lookahead

It is required to decide where one token will end and the next token will begin.

### - Ambiguities (REP)

- \* The longest match is preferred
- \* Among rules which matched the same number of characters, the rule given first is preferred.

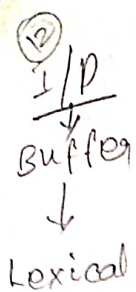
## Lexical Errors (REP)

- A character sequence that cannot be scanned into any valid token is a lexical error.
- Lexical errors are uncommon, but they still must be handled by a scanner.
- Misspelling of identifiers, keywords or operators are considered as lexical errors.

## Error Recovery Schemes

- panic mode recovery
- Local correction
- Global correction.



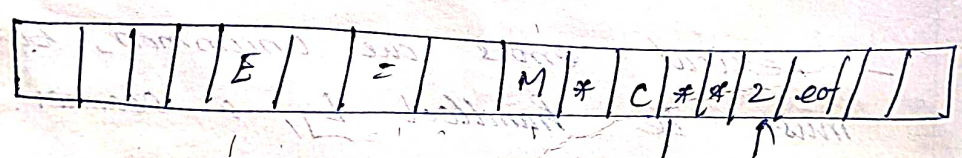


- Input Buffering → is a location that holds all input before it continues to process.
- To ensure that a right lexeme is found one or more characters have to be looked up beyond the next lexeme.
  - A two-buffer scheme is introduced to handle large lookaheads safely.
  - Techniques for speeding up the process of lexical analyzer such as the use of sentries to mark the buffer end have been adopted.

### Buffer Pairs

Because of large amount of time consumption in moving characters.

Specialized buffering techniques have been developed to reduce the amount of overhead required to process input characters.



Lexeme begins. Forward

Buffer Pairs

### Scheme

- consists of two buffers, each consists of N-characters size which are reloaded alternately.
- N-Number of character on one disk block
- N characters are read from the input file to the buffer using one system read command
- eof is inserted at the end of the number of character is less than N.



## Pointers

13

Two pointers

- lexemeBegin
- forward

### \* lexemeBegin

Points to the beginning of the current lexeme which is yet to be found

### \* Forward

Scans ahead until a match for a pattern is found.

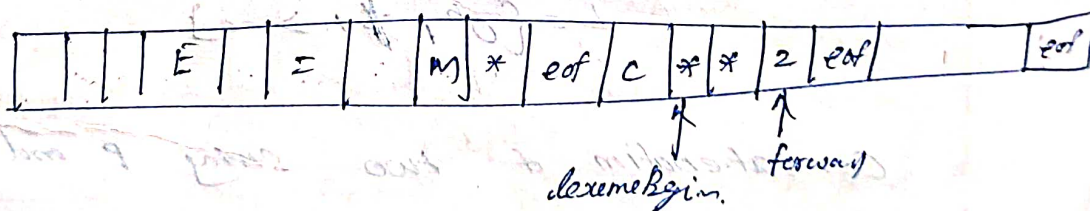
Sentinels - special character - eof.

Test 1 : For end of buffer

Test 2 : To determine what character is read.

→ The usage of sentinel reduces the two tests to one by extending each buffer half to hold a sentinel character at the end

→ The sentinel is a special character that cannot be part of the source program.



## Advantages

- Most of the times, it performs only one test to see whether Forward pointer points to an eof
- More tests are performed only when it reaches the end of the buffer half or eof.



## Specification of Token - Reg. EXP.

- Regular expressions are a notation to represent lexeme patterns for a token
- They are used to represent the language for lexical analyzers
- They assist in finding the type of token that accounts for a particular lexeme

## Strings and languages

Alphabets are finite set of symbols  $\Sigma$  (signal)

$\Sigma = \{0, 1\}$  - binary alphabet

String represents  $\{a, b, c, \dots\}$  lower case letters

$\Sigma = \{a, b\} \rightarrow$  generate n number  
 $w = \{01, 00, 01, 10, 11, 001, 010, \dots\}$

length  
empty  
prefix  
suffix  
sub string  
proper sub string

$a, b, ab, aab, ba$   
 indicates the set of possible strings for the binary alphabet  $\Sigma$

which are generated from all  
 Language (L) is the collection of strings which are accepted by finite automata

$L = \{0^n \mid n \geq 0\}$ ,  $\Sigma = \{a, b\} \rightarrow$   
 $L = \{a, b, ab, ba\}$   
 Concatenation of two string P and Q

$P = 010$

$Q = 001$

$PQ = 010001$

$QP = 001010$

$PQ \neq QP$



### Prefix

A prefix of any string  $s$ , is obtained by removing zero or more symbols from the end of  $s$ .

( $s = \text{balloon}$ .  $s = abc \Rightarrow \epsilon, abc$  )  
 $\Rightarrow \text{ball, balloon}$  )  
 prefix =  $\epsilon, abc$  (does not contain any char)

Suffix: A suffix of any string  $s$ , is obtained by removing zero or more symbols from the beginning of  $s$ .

( $s = \text{balloon}$ .  
 possible prefixes are: loon, balloon.)

### Operations on Languages

- Union
- Concatenation
- Closure

#### Union

Union of two language  $L$  and  $M$  produces the set of strings which may be either in language  $L$  or in language  $M$  or in both.

$L \cup M = \{p \mid p \text{ is in } L \text{ or } p \text{ is in } M\}$  → pipe symbol

#### Concatenation

$L = \{0, 1\}$   $M = \{00, 11\}$   
 $L \cup M = \{0, 1, 00, 11\}$

Combining The language  $L$  and  $M$ , produces a set of strings which are formed by merging the strings in  $L$  with strings in  $M$ .

$L \cup M = \{p_1 p_2 \mid p_1 \text{ is in } L \text{ and } p_2 \text{ is in } M\}$

$x_1 = abc$   
 $x_2 = hlo$   
 $x_4 = abede$



## Closure

### (i) Kleene closure ( $L^*$ )

Kleene closure refers to zero or more occurrences of input symbols in a string, and includes empty string  $\epsilon$ .

$$L^* = \bigcup_{i=0}^{\infty} L^i = L^0 \cup L^1 \cup L^2 \dots$$

### (ii) Positive closure ( $L^+$ )

It indicates one or more occurrence of input symbols in a string.

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

$$L^* = \{ \epsilon, a, b, aa, ab, ba, aab, aba, aaba, \dots \}$$

$$L^+ = \{ a, b, aa, ab, ba, bb, aab, aaba \}$$

$$L^3 = \{ aaa, aba, abb, bba, bab, bbb, \dots \}$$

### Precedence of operators

1. unary operator ( $*$ )
2. Concatenation operator ( $\cdot$ )
3. Union operator ( $|$  or  $\cup$ )

### Regular Expressions

It is a combination of input symbols and language operators such as union, concatenation, and closure.

It can be used to describe the identifiers for a language.

- The identifier is a collection of letters, digits or underscore which must begin with a letter.

letter - (letter | digit)\*

S.No	Regular expression	Languages
1	$r$	$L(r)$
2	$a$	$L(a)$
3	$r   s$	$L(r)   L(s)$
4	$rs$	$L(r) L(s)$
5	$r^*$	$(L(r))^*$

Language for regular expression

### Regular Definition

Regular definition  $d$ , gives aliases to regular expression  $r$  and uses it for convenience

$$d_1 \rightarrow r_1$$

$$d_2 \rightarrow r_2$$

$$d_3 \rightarrow r_3$$

$$d_n \rightarrow r_n$$

Regular expressions for identifiers and numbers given as follows

$$\text{letter} \rightarrow A | B | \dots | z | a | b | \dots | z$$

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

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

$$\text{num} \rightarrow \text{digit} (\text{digit})^*$$



## Recognition of Tokens

It also generates code for examining the input string and to find the prefixes (lexeme) that matches with any one of the patterns.

### Rules for conditional statement

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

expr  $\rightarrow$  term relop term  
 | term

term  $\rightarrow$  id | number

### Conditions for branching statements

The terminals of the grammar which are if, then, else, relop, id and number are the names of tokens for lexical analyzer.

Lexical analyzer also performs stripping out of white spaces.

ws  $\rightarrow$  (blank | tab | newline)<sup>+</sup>

### Token names with their attribute value

Lexemes	Token names	Attribute value
Any ws	-	-
if	if	-
Then	then	-
Else	else	-
Any i'd	id	Pointer to symbol table entry
Any number	number	Pointer to symbol table entry
<	relop	LT
<=	relop	LE
=	relop	EQ
>	relop	NE
>=	relop	GT
		GE



## Transition Diagrams

19

→ Transition diagrams are pictorial representation of transition from one state to another one.

→ Start State

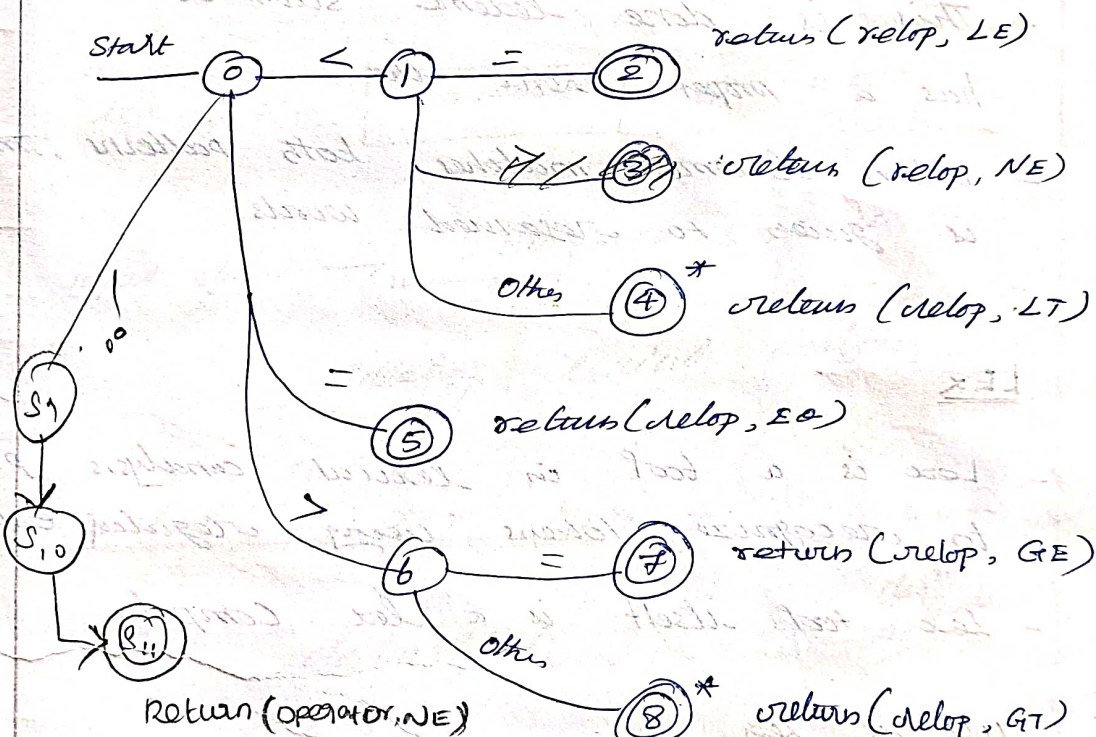


→ Final State



→ To indicate the retraction of forward pointer

(\*)

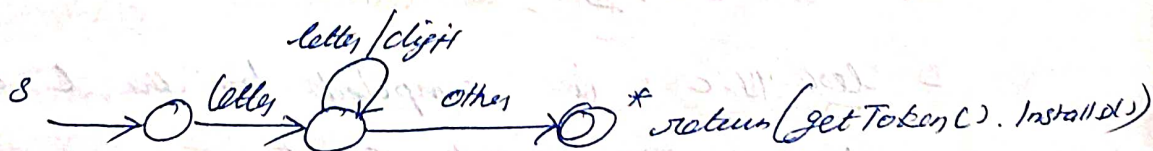


Transition diagram for relational operations.

## Recognition of Reserved words and Keywords

- Keyword patterns match that of identifiers.

- But they should be recognized differently



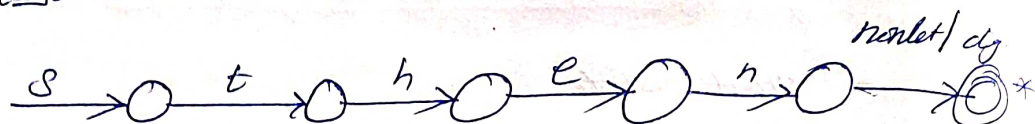
Transition diagram for identifier & keyword.



⑧

To handle reserved words different from identifiers

- Create separate transition diagram for each keyword.



Transition diagram for keyword.

- A test for 'non-letter or digit' to check the character for identifiers.
- If it reaches the accepting state, it is recognized as a keyword; else, it is an identifier.
- This is done for lexemes such as the next value.
- When a lexeme matches both patterns, priority is given to reserved words.

## LEX

- Lex is a tool in the lexical analysis phase to recognize tokens using regular expressions.
- Lex tool itself is a lex compiler.

## Use of Lex

- lex.l is an ASCII file written in a language which describes the generation of lexical analyzers.

The lex compiler transforms lex.l to a C program known as lex.yy.c.

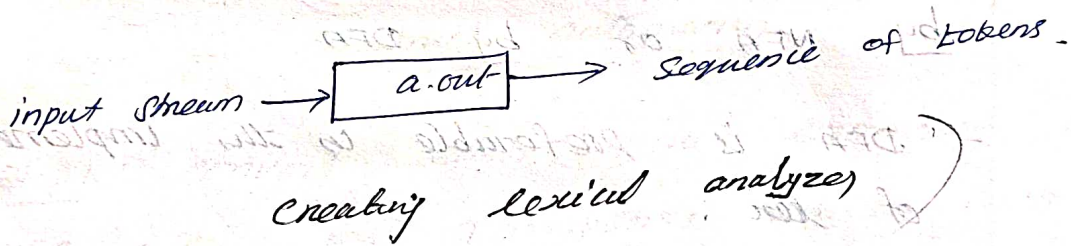
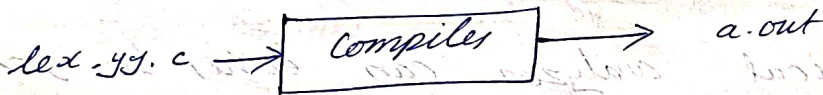
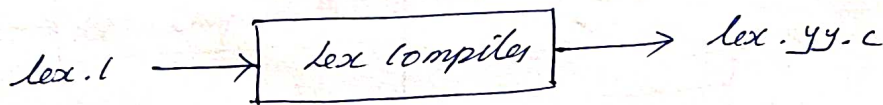
- lex.yy.c is compiled by the C compiler to a file called a.out.



- (21)
- The output of C compiler is the working lexical analyzer which takes stream of input character and produces a stream of tokens.

yyval is a global variable which is shared by lexical analyzer and parser to return the name and an attribute value of token.

- The attribute value can be numeric code, pointer to symbol table or nothing.



### Structure of Lex Programs

declarations

%. %.

translation rules

%. %.

auxiliary functions.

- declare constant
- declare variable
- regular expression and code segment.
- hold the additional functions
- These functions are compiled separately and loaded with lexical analyzer.

### Conflict Resolution in Lex

Conflict arises when several prefixes of input matches one or more patterns.



- Always prefer a longer prefix than a shorter prefix.
- If two or more patterns are matched for the longest prefix, then the first pattern listed in lex program is preferred.

### Lookahead operator

- Lookahead operator is the additional operator that is read by lex in order to distinguish additional patterns for a token.

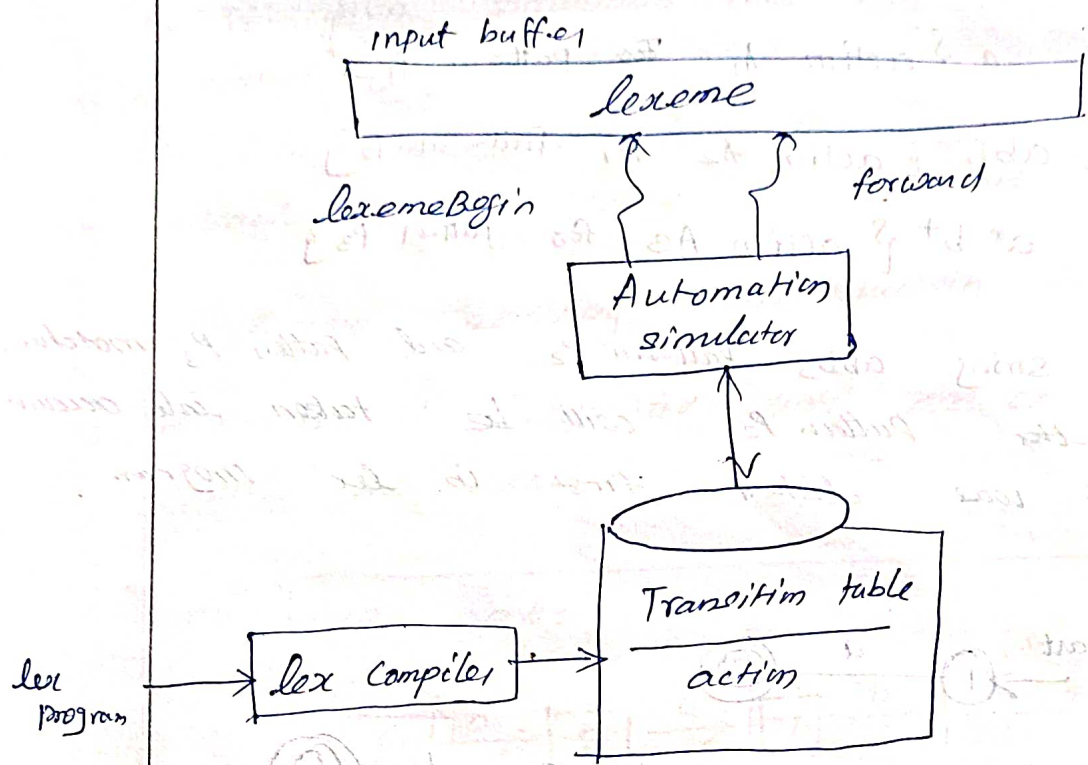
### DESIGN OF LEXICAL ANALYZER

- Lexical analyzer can either be generated by NFA or by DFA.
- DFA is preferable in the implementation of lex.

### Structure of Generated Analyzer

- Program to simulate automaton.
- Components created from lex program by lex itself which are listed as follows:
  - \* A transition table for automaton.
  - \* Functions that are passed directly through lex to the output.
  - \* Actions from input programs (fragments of code) which are invoked by automaton simulator when needed.



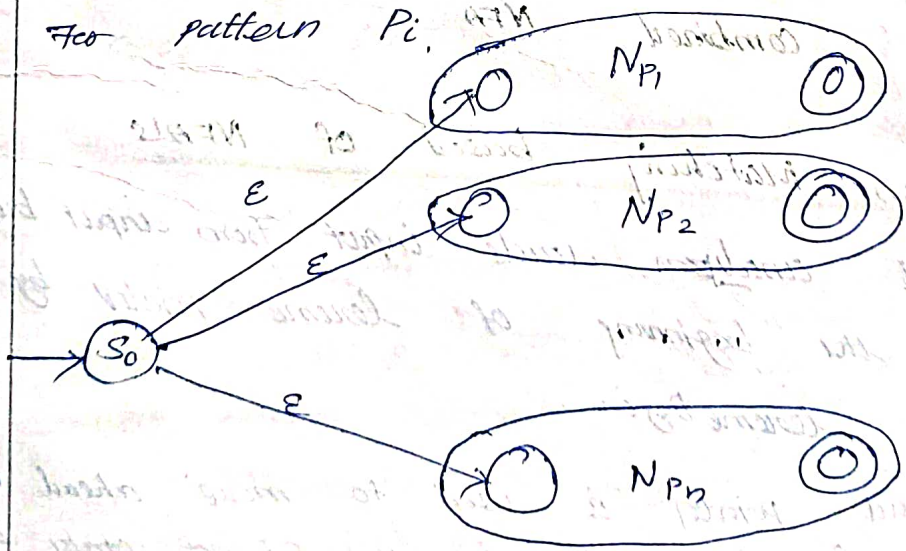


Lex program used by finite automaton simulator.

Steps to construct automation

Step 1: Convert each regular expression into NFA either by Thompson's sub-set construction or direct method.

Step 2: Combine all NFA into one by introducing new start state with  $\epsilon$ -transitions to each of start states of NFA's  $N_i$  for pattern  $P_i$ .



construction of NFA from Lex program

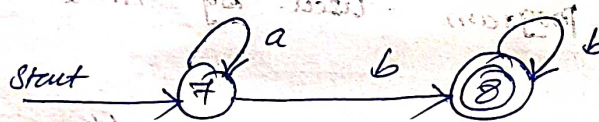
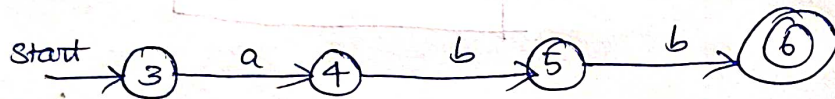
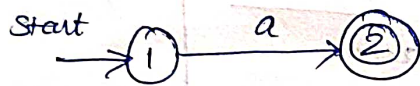


ex: a {action A<sub>1</sub> for pattern P<sub>1</sub>}

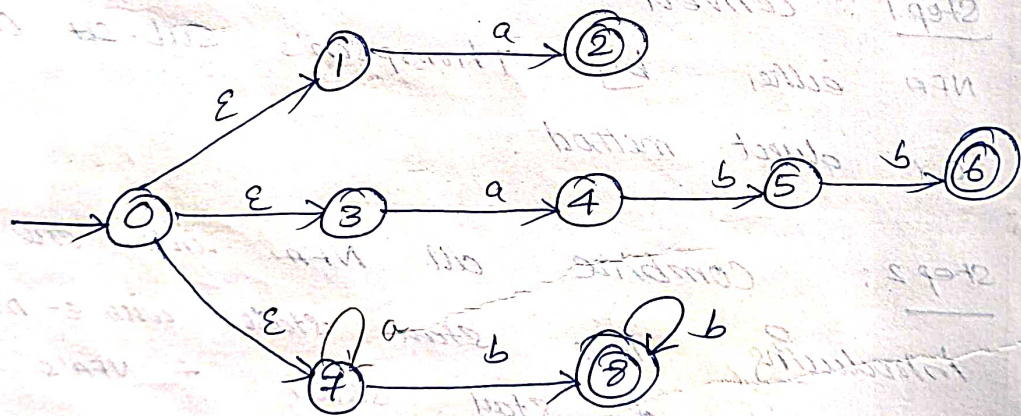
abb {action A<sub>2</sub> for pattern P<sub>2</sub>}

a<sup>+</sup>b<sup>+</sup> {action A<sub>3</sub> for pattern P<sub>3</sub>}

For string abb, pattern P<sub>2</sub> and pattern P<sub>3</sub> match.  
 But the pattern P<sub>2</sub> will be taken into account  
 as it was listed first in lex program.



NFA's for a, abb, a<sup>+</sup>b<sup>+</sup>



Combined NFA

### Pattern Matching based on NFAs

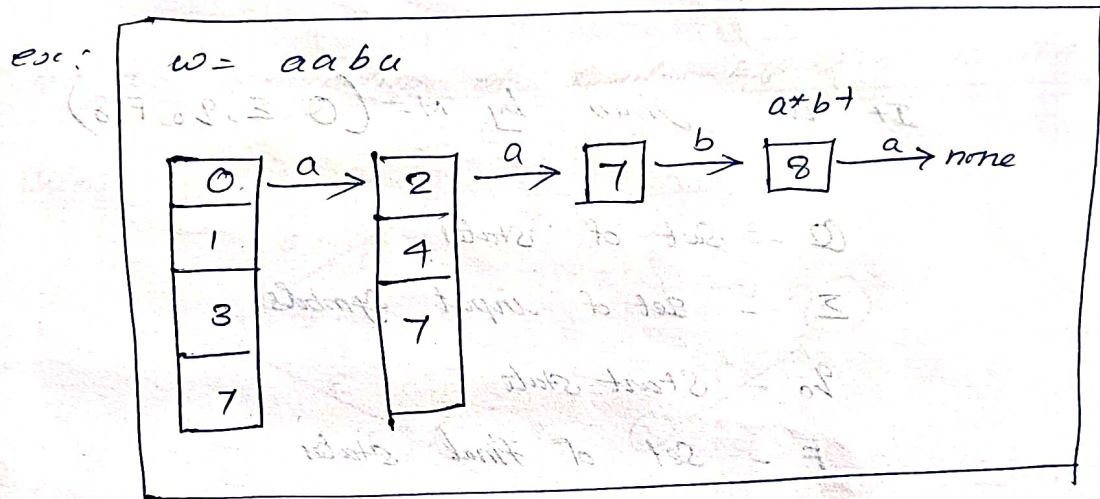
- Lexical analyzer reads input from input buffer from the beginning of lexeme pointed by the pointer lexeme begin.
- Forward pointer is used to move ahead of input symbols, calculates the set of states it is in at each point.



- If NFA simulation has no next state for some input symbol, then there will be no longer prefix which reaches the accepting state exists.

- Lexeme matching some pattern.

- process is repeated until one or more accepting states are reached.



processing input aaba

- Process starts with & closure of initial state 0.

- After processing all the input symbols, no state is found as there is no transition out of state 8 on input a.

- accepting is looked by retracting to previous state.

- From state 2 which is an accepting state is reached after reading i/p symbol a and therefore the pattern a has been matched.

- At state 8, string aab has been matched with pattern  $a+b^+$ .

- By lex rule the longest matching prefix should be considered.

- Action A3 corresponding to pattern P3 will be executed for the string aab.



(20)

## Finite Automata

A recognizer for a language is a program that takes as input a string  $x$  and answers yes if  $x$  is a sentence of the language or no otherwise.

⇒ A regular expression is compiled into a recognizer by constructing a generalized transition diagram called a Finite Automaton (FA).

It is given by  $M = (Q, \Sigma, q_0, F, \delta)$

$Q$  - Set of states

$\Sigma$  - set of input symbols

$q_0$  - Start state

$F$  - Set of Final states

$\delta$  - Transition Function

### Deterministic Finite Automata DFA

- For each state and for each input symbol exactly one transition occurs from the state.

### Non-deterministic Finite Automata (NFA)

- More than one transition occurs for any input symbol from a state.

- Transition can occur even on empty string ( $\epsilon$ )

Regular expression can be converted into DFA

#### (1) Thompson's Sub-set Construction

- Given regular expression is converted into NFA

- Resultant NFA is converted into DFA



## (ii) Direct method

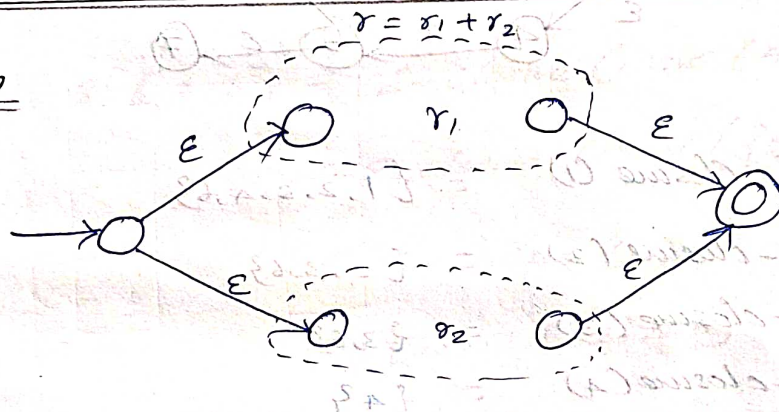
- In direct method, given regular expression is converted directly into DFA.

## Regular Expressions to DFA

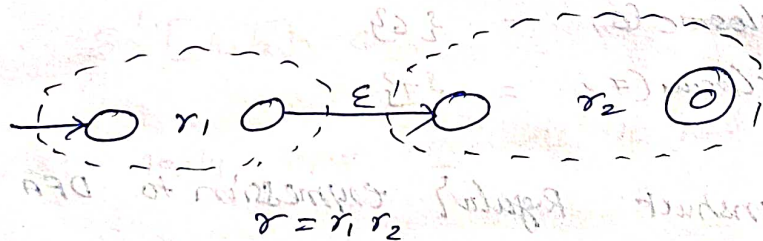
Regular expression is used to represent the language (lexeme) of finite automata (lexical analyzer)

## Rules for conversion of Regular Expression to NFA

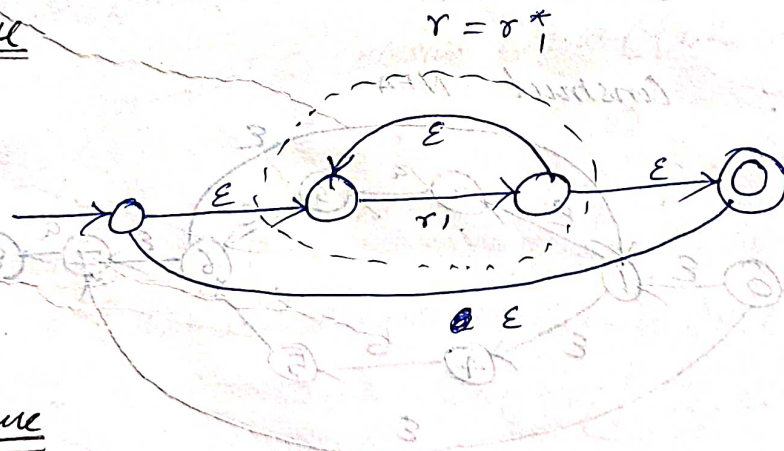
### \* Union



### \* Concatenation



### \* Closure



## E-closure

E-closure is the set of states that are reachable from the state concerned on taking empty string as input.

It describes the path that consumes empty string (E) to reach some state of NFA.



ex:

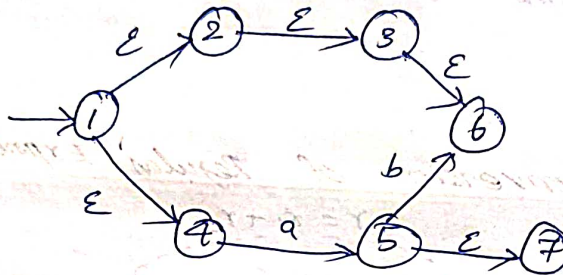


$$\epsilon\text{-closure}(q_0) = \{q_0, q_1, q_2\}$$

$$\epsilon\text{-closure}(q_1) = \{q_1, q_2\}$$

$$\epsilon\text{-closure}(q_2) = \{q_2\}$$

ex:



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

$$\epsilon\text{-closure}(2) = \{2, 3, 6\}$$

$$\epsilon\text{-closure}(3) = \{3, 6\}$$

$$\epsilon\text{-closure}(4) = \{4\}$$

$$\epsilon\text{-closure}(5) = \{5, 7\}$$

$$\epsilon\text{-closure}(6) = \{6\}$$

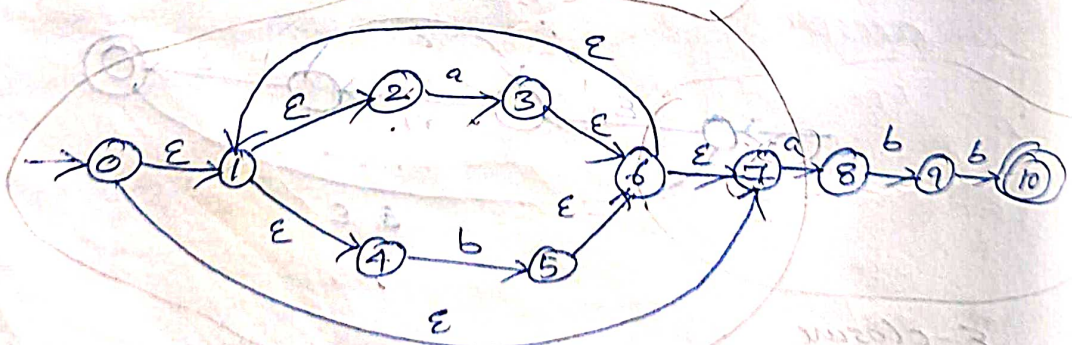
$$\epsilon\text{-closure}(7) = \{7\}$$

Ⓐ

ex: Construct Regular expression to DFA

$$RE = (a+ b)^* abb$$

Step 1: Construct NFA



Step 2: Start with finding  $\epsilon$ -closure of state 0

$$\epsilon\text{-closure}(0) = \{0, 1, 2, 4, 7\} = A$$



Step 3: Apply input symbol  $a, b$  to  $A$

(21)

$$\begin{aligned} \delta' [A, a] &= \varepsilon\text{-closure}(\text{move}(A, a)) \\ &= \varepsilon\text{-closure}(\text{move}(\{0, 1, 2, 4, 7\}, a)) \\ &= \varepsilon\text{-closure}(3, 8) \quad (0, a) \quad (1, a) \quad (2, a) \quad (4, a) \quad (7, a) \\ &= \{3, 6, 7, 1, 2, 4, 8\} \\ &= \{1, 2, 3, 4, 6, 7, 8\} = B \end{aligned}$$

$$\delta' [A, a] = B$$

$$\begin{aligned} \delta' [A, b] &= \varepsilon\text{-closure}(\text{move}(A, b)) \\ &= \varepsilon\text{-closure}(\text{move}(\{0, 1, 2, 4, 7\}, b)) \\ &= \varepsilon\text{-closure}(5) \\ &= \{5, 6, 7, 1, 2, 4, 8\} \\ &= \{1, 2, 4, 5, 6, 7\} = C \end{aligned}$$

$$\delta' [A, b] = C$$

Step 4: Apply input symbols to new state  $B$

$$\begin{aligned} \delta' [B, a] &= \varepsilon\text{-closure}(\text{move}(B, a)) \\ &= \varepsilon\text{-closure}(\text{move}(\{1, 2, 3, 4, 6, 7, 8\}, a)) \\ &= \varepsilon\text{-closure}(3, 8) \\ &= \{1, 2, 3, 4, 6, 7, 8\} = B \end{aligned}$$

$$\delta' [B, b] = \varepsilon\text{-closure}(\text{move}(B, b))$$

$$\begin{aligned} &= \varepsilon\text{-closure}(\text{move}(\{1, 2, 3, 4, 6, 7, 8\}, b)) \\ &= \varepsilon\text{-closure}(5, 9) \rightarrow \{6, 1, 2, 4, 7, 8\} \\ &= \{1, 2, 4, 5, 6, 7, 9\} = D \end{aligned}$$

$$\delta' [B, b] = D$$

(30)

Step 5: Apply input symbols to new state c

$$\begin{aligned}
 \delta'[c, a] &= \varepsilon\text{-closure}(\text{move}(c, a)) \\
 &= \varepsilon\text{-closure}(\text{move}(\{1, 2, 4, 5, 6, 7\}, a)) \\
 &= \varepsilon\text{-closure}(3, 8) \\
 &= \{1, 2, 3, 4, 6, 7, 8\} = B
 \end{aligned}$$

$$\delta'[c, a] = B$$

$$\begin{aligned}
 \delta'[c, b] &= \varepsilon\text{-closure}(\text{move}(c, b)) \\
 &= \varepsilon\text{-closure}(\text{move}(\{1, 2, 4, 5, 6, 7\}, b)) \\
 &= \varepsilon\text{-closure}(5) \\
 &= \{1, 2, 4, 5, 6, 7\} = c
 \end{aligned}$$

$$\delta'[c, b] = c$$

Step 6: Apply input symbols to new state d

$$\begin{aligned}
 \delta'[d, a] &= \varepsilon\text{-closure}(\text{move}(d, a)) \\
 &= \varepsilon\text{-closure}(\text{move}(\{1, 2, 4, 5, 6, 7, 9\}, a)) \\
 &= \varepsilon\text{-closure}(3, 8) \\
 &= \{1, 2, 3, 4, 6, 7, 8\} = B
 \end{aligned}$$

$$\delta'[d, a] = B$$

$$\begin{aligned}
 \delta'[d, b] &= \varepsilon\text{-closure}(\text{move}(d, b)) \\
 &= \varepsilon\text{-closure}(\text{move}(\{1, 2, 4, 5, 6, 7, 9\}, b)) \\
 &= \varepsilon\text{-closure}(5, 10) \\
 &= \{1, 2, 4, 5, 6, 7, 10\} = E
 \end{aligned}$$

$$\delta'[d, b] = E$$



Step 7: Apply input symbols to new State E (31)

$$\begin{aligned}
 \delta'[E, a] &= \epsilon\text{-closure}(\text{move}(E, a)) \\
 &= \epsilon\text{-closure}(\text{move}(\{1, 2, 4, 5, 6, 7, 10\}, a)) \\
 &= \epsilon\text{-closure}(3, 8) \\
 &= \{1, 2, 3, 4, 6, 7, 8\} = B
 \end{aligned}$$

$$\delta'[E, a] = B$$

$$\begin{aligned}
 \delta'[E, b] &= \epsilon\text{-closure}(\text{move}(E, b)) \\
 &= \epsilon\text{-closure}(\text{move}(\{1, 2, 4, 5, 6, 7, 10\}, b)) \\
 &= \epsilon\text{-closure}(5) \\
 &= \{1, 2, 4, 5, 6, 7\} = C
 \end{aligned}$$

$$\delta'[E, b] = C$$

Step 8: Construct transition table

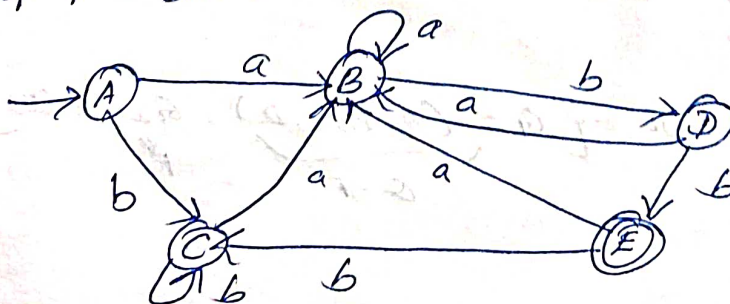
	a	b
→ A	B	C
B	B	D
C	B	C
D	B	E
* E	B	C

Note:

→ Start state is the  $\epsilon$ -closure (0), i.e. A.

→ Final state is the state that contains Final state of drawn NFA.

Step 9: Construct transition diagrams

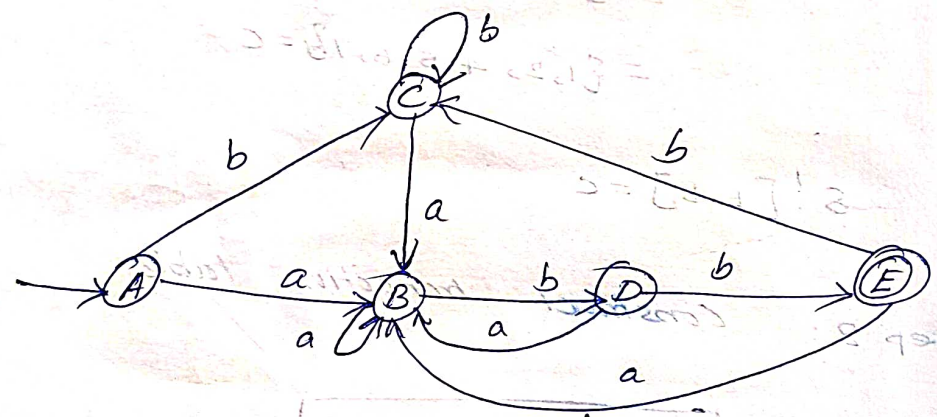




# Minimization of DFA

- Lexical analyzer can be implemented by
- The number of states of constructed DFA must be minimized as for each state is made in a table which describes the lexical analyzer.
- Minimization of DFA focuses on reducing the number of states in the given finite automata.

Ex: 1



Step 1:

Construct transition table

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

Step 2: partition the states into two groups F and Q-F, final and non-final states respectively

$$P = \{ G_1 = \underbrace{(A, B, C, D)}_{Q-F}, G_2 = \underbrace{(E)}_F \}$$

Step 3: Construct  $\pi$ , by partitioning groups in  $\pi$

Compare (A, B)

	a	b
A	B	C
B	B	D

$G_1$

$G_1$

Both are in same group

$\Rightarrow$  1 - equivalent

Compare (A, C)

	a	b
A	B	C
C	B	C

$G_1$

$G_1$

Both are in same group

$\Rightarrow$  1 - equivalent

Compare (A, D)

	a	b
A	B	C
D	B	E

$\downarrow$

Not in same group

$\Rightarrow$  Not 1 - equivalent

$$\pi_1 = \{ (A, B, C), (D), (E) \}$$

Step 4: Construct  $\pi_2$  from  $\pi_1$

For (A, B)

	a	b
A	B	C
B	B	D

Not in same group

$\Rightarrow$  Not 2 - equivalent

For (A, C)

	a	b
A	B	E
C	B	C

$G_1$

$G_2$

$\Rightarrow$  2 - equivalent

$$\pi_2 = \{ (A, C), (B), (D), (E) \}$$



24

Step 5: Construct  $\pi_3$  from  $\pi_2$

For  $(A, C)$

	a	b
A	B	C
C	B	C

$\downarrow$   $G_2$        $\downarrow$   $G_1$

Both are in same group

$\Rightarrow$  3-equivalent

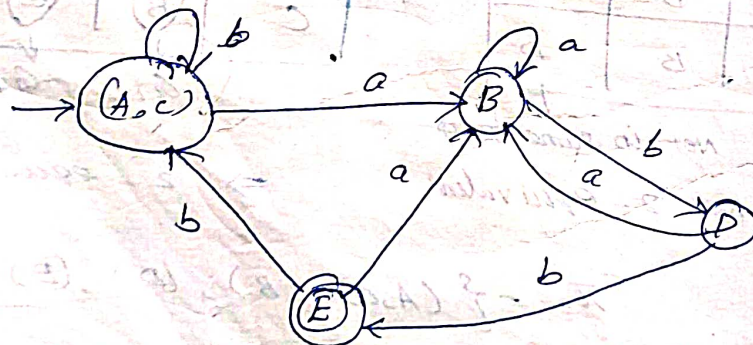
$$\pi_3 = \{(A, C), (B), (D), (E)\}$$

Step 6: Stop partitioning since  $\pi_3 = \pi_2$

Step 7: Construct minimized transition table

State \ input	a	b
$\rightarrow (A, C)$	B	$(A, C)$
B	B	D
D	B	E
* E	B	$(A, C)$

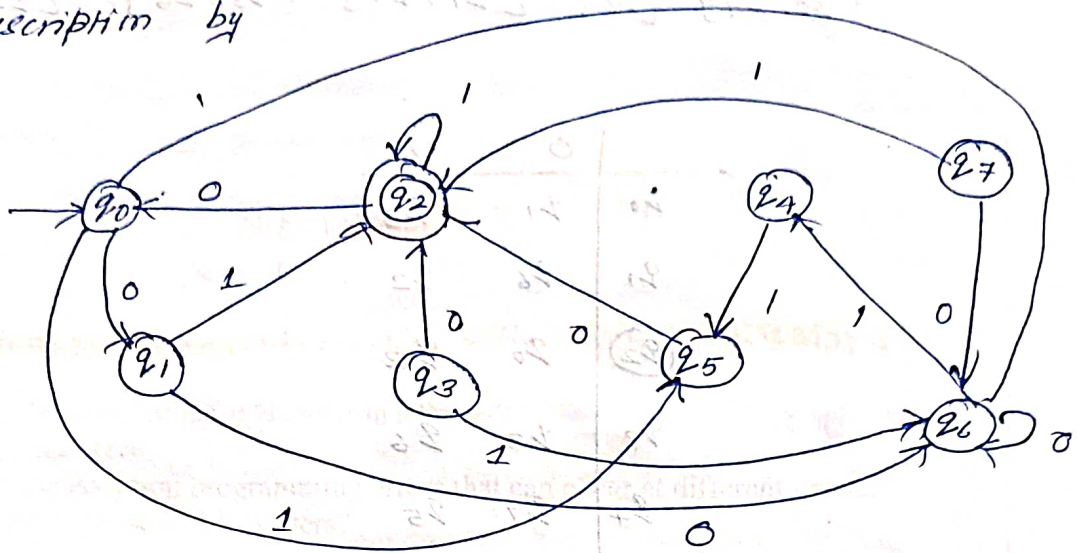
Step 8: Draw minimized finite automata





Ex: 2

Construct a minimize DFA equivalent to the DFA description by



	0	1
→ q0	q1	q5
q1	q6	q2
q2	q0	q2
q3	q2	q6
q4	q7	q5
q5	q2	q6
q6	q6	q4
q7	q6	q2

0 - Equivalence

$\{q0, q1, q3, q4, q5, q6, q7\}$   $\{q2\}$

1 - Equivalence

$\{q0, q4, q6\}$

$\{q1, q7\}$

$\{q3, q5\}$   $\{q2\}$

2 - Equivalence

$\{q0, q4\}$   $\{q6\}$   $\{q1, q7\}$   $\{q3, q5\}$   $\{q2\}$



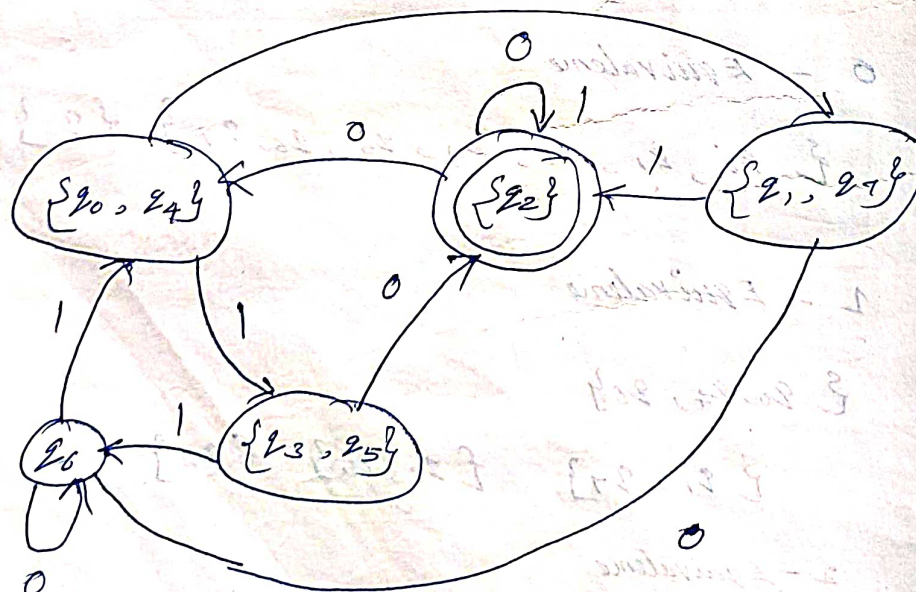
(36)

## 3- Equivalence

 $\{20, 24\}$   $\{26\}$   $\{21, 27\}$   $\{23, 25\}$   $\{22\}$ 

	0	1
20	21	25
21	26	22
22	20	22
23	22	26
24	27	25
25	22	26
26	26	24
27	26	22

	0	1
$\{20, 24\}$	$\{21, 27\}$	$\{23, 25\}$
$\{26\}$	$\{26\}$	$\{20, 24\}$
$\{21, 27\}$	$\{26\}$	$\{22\}$
$\{23, 25\}$	$\{22\}$	$\{26\}$
$\{22\}$	$\{20, 24\}$	$\{22\}$





## Unit-I

# Introduction to Compiler and Lexical Analysis

### Introduction:

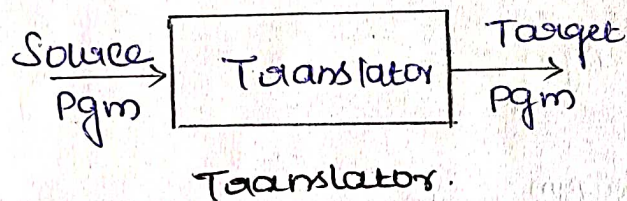
Compilers are basically translators. Designing a compiler for some language is a complex and time consuming process.

### Translators:

A translator is one kind of pgm that takes source code as i/p and converts it into another form.

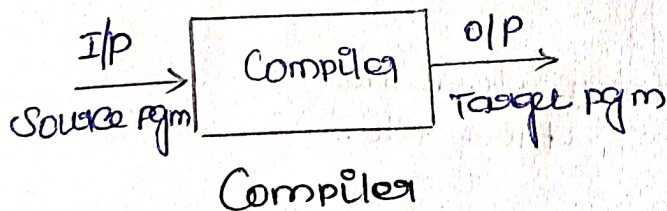
The i/p pgm is called source languages and the o/p pgm is called target language.

eg: C, C++, Fortran



### Process of compilation and Interpretation:

Compiler pgm which takes one language (source pgm) as i/p and translates it into an equivalent another language (target pgm)



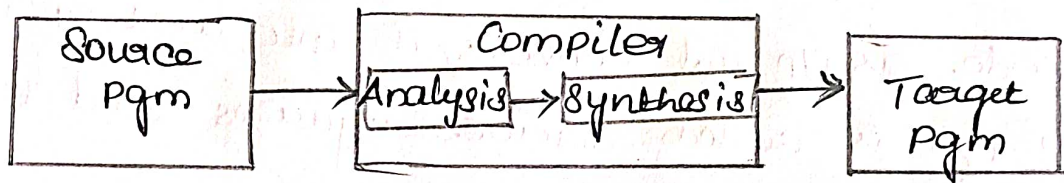
⇒ Source Prgms as higher level language such as C, PASCAL, FORTRAN and converts into low level language or a machine level language such as assembly language.



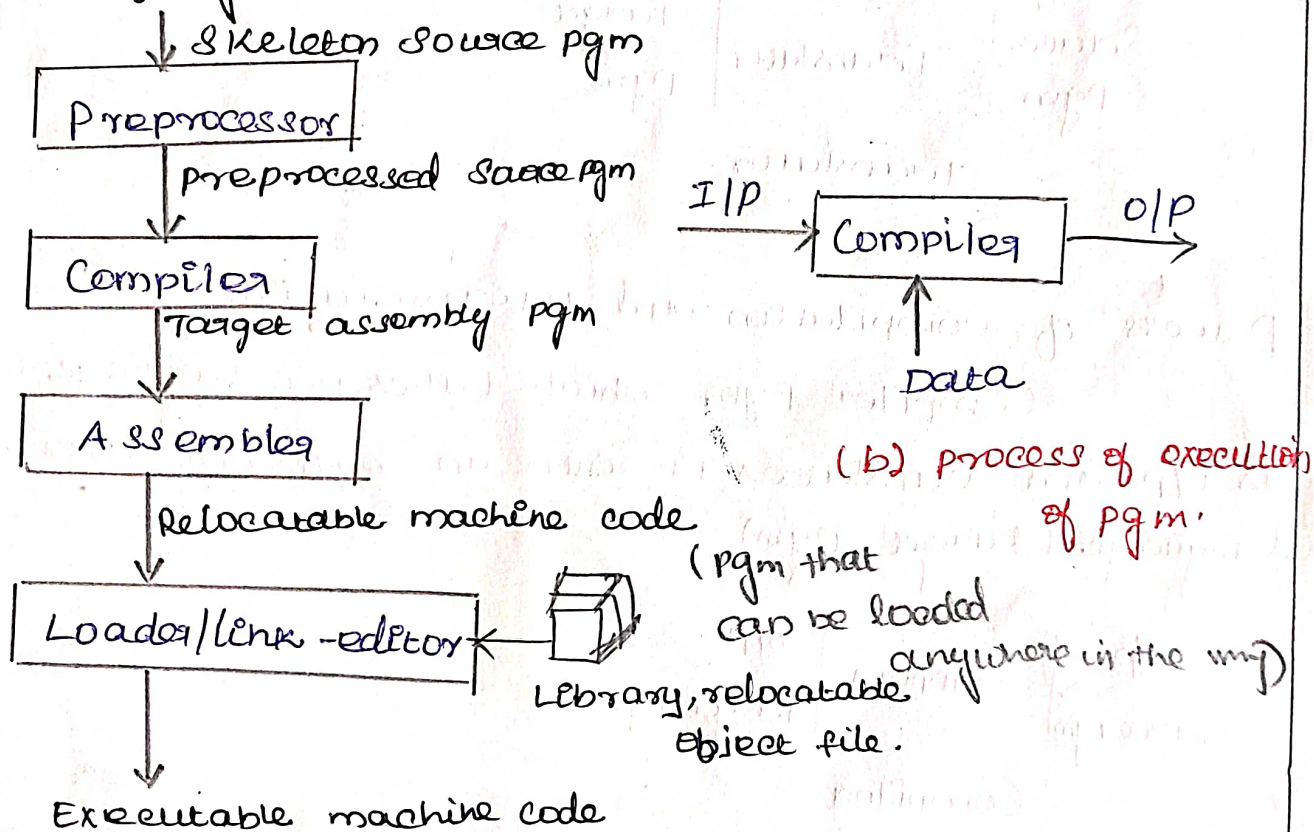
# Analysis - Synthesis Model:

Compilation can be done in two parts:

- ① Analysis part: Source pgm is read and broken down into constituent pieces
- ② Synthesis part: The intermediate form of the source language is taken and converted into an equivalent target pgm.



## Execution of pgm:



## (a) process of execution of pgm.

→ Generate executable file.  
Linker - combines obj file  
loader - loads <sup>exec. file</sup> to main mem

## Properties of Compiler:

(2)

- 1) The Compiler itself must be bug-free
- 2) It must generate correct machine code
- 3) The generated machine code must run fast
- 4) The compiler itself must run fast (Compilation time must be proportional to pgm size)
- 5) The Compiler must be portable.
- 6) It must give good diagnostics and error msg.
- 7) The generated code must work well with existing debugger.
- 8) It must have consistent optimization.

Why should we study compilers?

1. Performance can be measured.
2. The abilities of programming languages can be understood
3. S/m building Tools such as LEX and can be developed by analytical study of compilers.



SYNTAX ANALYSISPARSER

Parser is a program that obtains tokens from lexical analyzer and constructs the parse tree which is passed to the next phase of compiler for further processing.

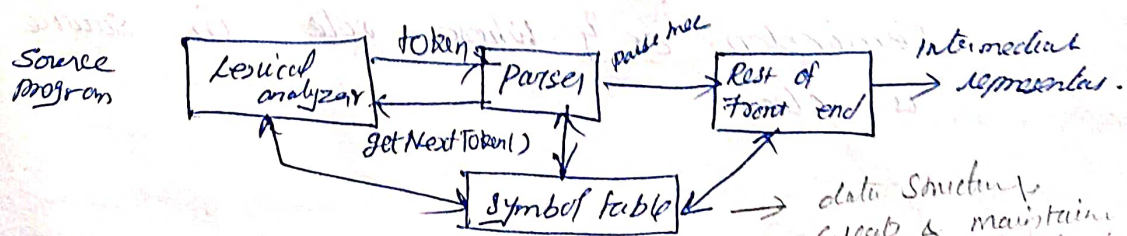
Parser implements CFG for performing error checks.

Types of Parser

- Top down parsers - Construct parse tree from root to leaves
- Bottom up parsers - Construct parse tree from leaves to root.

Role of Parser

- \* Once a token is generated by the lexical analyzer, it is passed to the parser.
- \* On receiving a token, the parser verifies the string of token names that can be generated by the Grammar of source language.
- \* It calls the function getNextToken(), to notify the lexical analyzer to yield another token.
- \* It scans the token one at a time from left to right to construct the parse tree.
- \* It also checks the syntactic construct of the grammar.



Position of parser in compiler

data structure  
create & maintain  
by compiler in order to  
store info

## Need For parser

- Parser is needed to detect syntactic errors efficiently.
- Error is detected as soon as a prefix of the input cannot be completed to form a string in the language.
- This process of analyzing the prefix of input is called Viable - prefix property.

Syntax Analysis Error: ① Errors in structure

Error Recovery Strategies ② Missing operators

③ unbalanced parentheses

\* It is used by the parser to recover from errors once it is detected.

\* The simplest recovery strategy is to quit parsing with an error message for the the itself.

### 1. Panic Mode Recovery

1. It prevents the parser from developing infinite loops.

2. When parser finds an error in the statement it ignores the rest of the statement by not processing the input.

3. The parser intends to find designated set of synchronizing tokens by discarding input symbols one at a time.

4. Synchronizing tokens may be delimiters, Semicolon or { whose role in source program is clear.



## 2. Phrase Level Recovery

- When a parser finds an error, it tries to take corrective measures so that the rest of inputs of statements allow the parser to parse ahead.
- one wrong correction will lead to an infinite loop.

## 3. Error production

- Productions which generate erroneous constructs are augmented to the grammar by considering common errors that occur.
- These productions detect the anticipated errors during parsing.
- Error diagnostic about the erroneous constructs are generated by the parser.

## - Global correction

- \* There are algorithms which make changes to modify an incorrect string into a correct string.

## GRAMMARS

It is used to describe the syntax of a programming language. It specifies the structure of expression and statement.

Start  $\rightarrow$  if (expr) then start

## Types of Grammar

- Type 0 Grammar - Turing machine
- Type 1 Grammar - Context sensitive language
- Type 2 Grammar - Context free Grammar
- Type 3 Grammar - Finite Automata

## Context Free Grammar

A Context Free Grammar  $G$  is defined by four tuples as

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

Where

- $G$  - Grammar
- $V$  - Set of variables
- $T$  - Set of terminals
- $P$  - Set of productions
- $S$  - Start symbol

It produces Context Free Language (CFL) which is defined as

$$L(G) = \{w \mid w \text{ is in } T^* \mid S \xrightarrow{*}_G w\}$$

$L$  - Language

$G$  - Grammar

$w$  - input string

$S$  - Start symbol

$T$  - Terminal

### Conventions

#### Terminals

- Lowercase letters i.e.  $a, b, c$
- Operators  $+, -, *$
- Punctuation symbols, comma, parentheses
- Digits  $0, 1, 2, \dots, 9$
- Bold face letters i.e.  $\mathbf{a}, \mathbf{b}, \mathbf{c}$



Non-terminals -  $A, B, C$

Start symbols -  $\rightarrow S$

Production -  $LHS \rightarrow RHS$ , (or) head  $\rightarrow$  body

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

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

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

Solution.

$$V = \{E, T, F\}$$

$$T = \{+, -, *, /, (, ), id\}$$

$$S = \{E\}$$

P:

$$E \rightarrow E + T$$

$$E \rightarrow E - T$$

$$E \rightarrow T$$

$$T \rightarrow T * F$$

Writing a Grammar

- Grammars are more powerful than regular expressions.

$$RE = (a|b)^* abb$$

$$A_0 \rightarrow aA_0 \mid bA_0 \mid aA_1$$

$$A_1 \rightarrow bA_2$$

$$A_2 \rightarrow bA_3$$

$$A_3 \rightarrow \epsilon \quad (\text{Set of string ending with } abb)$$

## Derivations

Derivation is used to find whether the string belongs to a given grammar.

Types:

- Leftmost derivation
- Rightmost derivation

### \* Leftmost derivation

At each and every step the leftmost non-terminal is expanded by substituting its corresponding production to derive a string.

ex:

$$E \rightarrow E + E \mid E * E \mid id$$

let:

$$w = id + id * id$$

$$E \rightarrow E + E$$

lm

$$E \rightarrow id + E$$

lm

$$E \rightarrow id + E * E$$

lm

$$E \rightarrow id + id * E$$

lm

$$E \rightarrow id + id * id$$

lm

ex:

$$S \rightarrow SS + \mid SS * \mid a$$

$$w = aa + a *$$

$$S \rightarrow SS *$$

$$S \rightarrow SS +$$

$$S \rightarrow aS +$$

$$S \rightarrow aa +$$

$$S \rightarrow aa + a *$$

$$[S \rightarrow SS +]$$

$$[S \rightarrow a]$$

$$[S \rightarrow a]$$

$$[S \rightarrow a]$$

H.W



## Rightmost derivation

(43)

ex:

$$E \rightarrow E + E \mid E * E \mid id$$

Let

$$w = id + id * id$$

$$E \rightarrow E + E$$

$$E \rightarrow E + E * E$$

$$E \rightarrow E + E * id$$

$$E \rightarrow E + id * id$$

$$E \rightarrow id + id * id$$

$$(E \rightarrow E * E)$$

$$[E \rightarrow id]$$

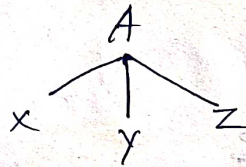
$$[E \rightarrow id]$$

$$[E \rightarrow id]$$

## Parse Tree

- parse tree is a hierarchical structure which represents the derivation of the grammar to yield input strings.

- Root node of parse tree has the start symbol of the given grammar from where the derivation proceeds.

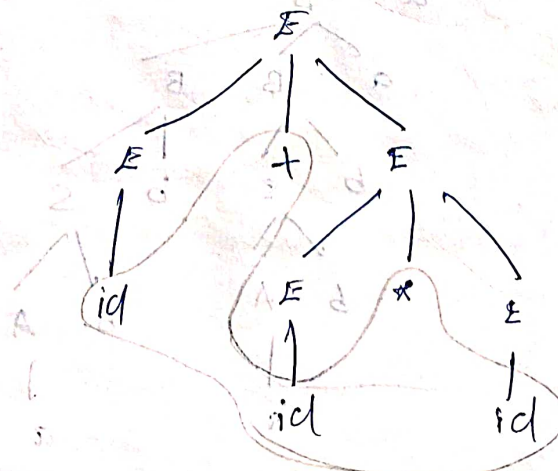


ex: Construct parse tree

$$E \rightarrow E + E \mid E * E \mid id$$

Let

$$w = id + id * id$$



42

ex! Let  $G$  be the Grammar

$$S \rightarrow aB / bA$$

$$A \rightarrow a / aS / bAA$$

$$B \rightarrow b / bS / aBB$$

$$S \rightarrow aABe$$

$$A \rightarrow Abc / b$$

$$B \rightarrow d$$

$$w = abbcd$$

H.W!

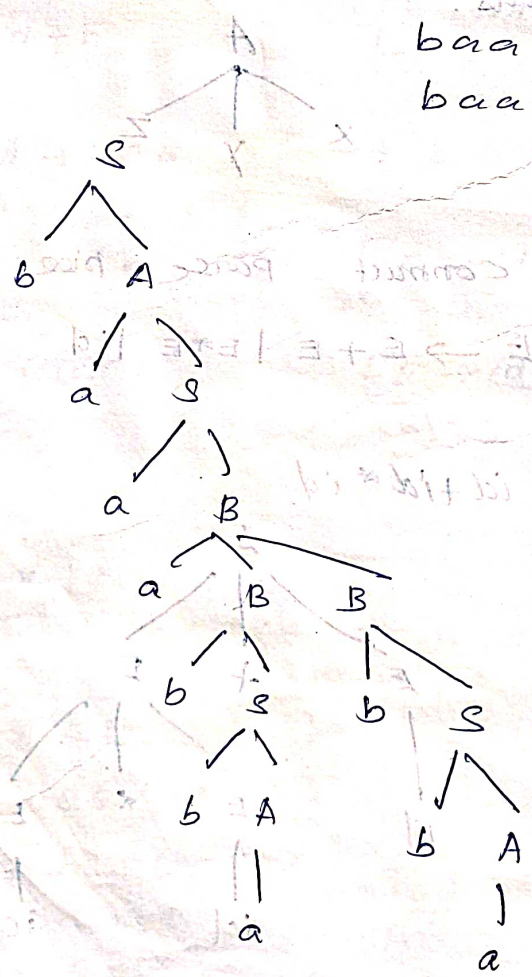
For the string baaabbbabba Find leftmost derivation, rightmost derivation and parse tree.

Leftmost

Rightmost

$S$   
 $bA$   
 $baS$   
 $baaB$   
 $baaaBB$   
 $b a a a b S B$   
 $b a a a b b A B$   
 $b a a a b b a B$   
 $b a a a b b a b S$   
 $b a a a b b a b b A$   
 $b a a a b b a b b a$

$S$   
 $bA$   
 $baS$   
 $baaB$   
 $b a a a B B$   
 $b a a a B b S$   
 $b a a a B b b A$   
 $b a a a B b b a$   
 $b a a a b b A b b a$   
 $b a a a b b a b b a$





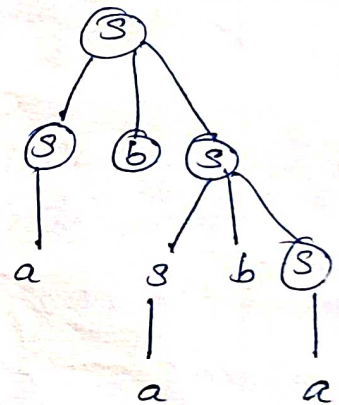
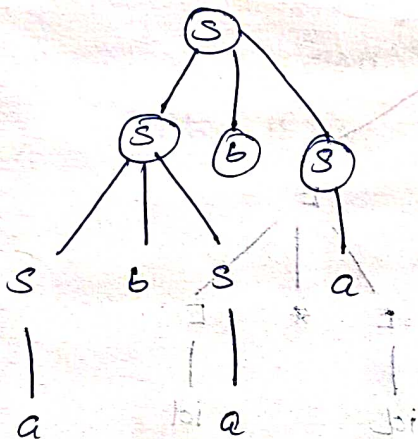
## AMBIGUITY

Grammar  $G$  is said to be ambiguous if it produces different parse trees for the same sentence or yield  $w$  which has been derived from the same start symbol.

ex: 1. if  $G$  is the grammar  $S \rightarrow Sbs \mid a$ . show that  $G$  is ambiguous

$w = ababa$  be the string

$$\begin{array}{l} S \rightarrow Sbs \\ S \rightarrow a \end{array}$$



LM  $\Rightarrow$  ababa  $\Leftarrow$  RM

ex 2:

The CFG given by  $G = (V, T, P, S)$  where

$$V = \{E\}$$

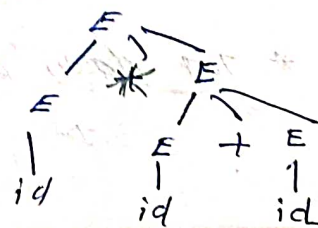
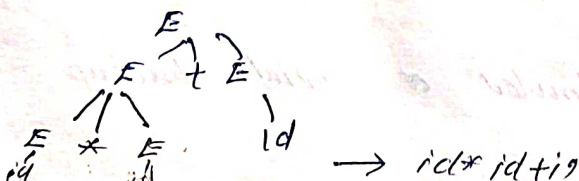
$$T = \{id\}$$

$$P = \{E \rightarrow E + E, E \rightarrow E * E, E \rightarrow id\}$$

$$S = \{E\}$$

Is the grammar ambiguous?

$$G = id * id + id$$



(46)

Top-down Parsing

Top-down parsing constructs parse tree for the input string.

Starting from root node and creating the nodes of parse tree in pre-order.

ex:

$$E \rightarrow E + E \mid E * E \mid id$$

$$w = id + id * id$$

$$E \rightarrow E + E$$

lm

$$E \rightarrow id + E$$

$$E \rightarrow id + E * E$$

$$E \rightarrow id + id * E$$

$$E \rightarrow id + id * id$$

$$E \rightarrow id$$

$$E \rightarrow E * E$$

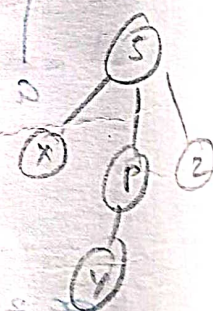
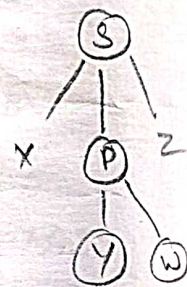
$$E \rightarrow id$$

$$E \rightarrow id$$

$$S \rightarrow x p z$$

$$p \rightarrow y w \mid y$$

$$w = \boxed{x \mid y \mid z}$$

General strategies\* Brute-force method

All possible combinations are attempted before the failure to parse is recognized.

\* Recursive descent

It is parsing techniques which does not allow backup.

\* Involves backtracking and left recursion

\* Top down parsing

with limited or partial backup.



## Recursive Descent Parser

- \* Recursive descent parser is a top-down parser.
- \* It requires backtracking to find the correct production to be applied.
- \* The parsing program consists of a set of procedures, one for each non-terminal.
- \* process begins with the procedure for start symbol.

Void AC {

    choose an A - production,  $A \rightarrow x_1 x_2 x_3 \dots x_k$

    for ( $i=1$  to  $k$ )

        if ( $x_i$  is a non-terminal)

            call procedure  $x_i()$ ;

        else if ( $x_i$  equals the current input symbol a)

            advance the input to the next symbol;

        else

            error;

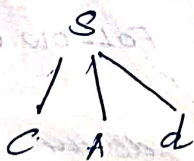
Procedure for a non-terminal in top-down parser

ex. Let grammar G be

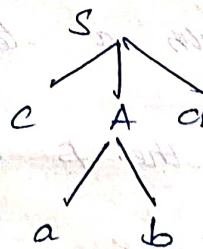
$S \rightarrow CAD$

$A \rightarrow ab|a$

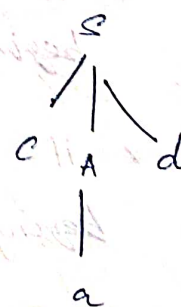
$w = cad$



$\Rightarrow$

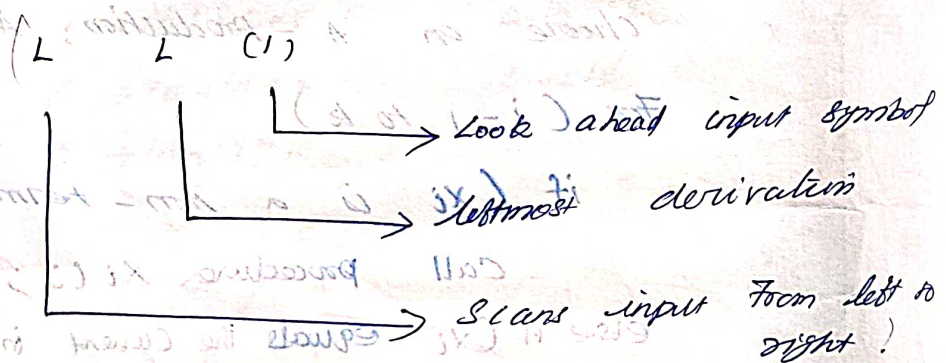


$\Rightarrow$



## PREDICTIVE PARSER / LL(1) PARSER

- Predictive parser are top-down parser.
- It is a type of recursive descent parser but with no backtracking.
- It can be implemented non-recursively by using stack data structure.
- They can also be termed as LL(1) parser as it is constructed for a class of grammar called LL(1).



A Grammar  $G$  is LL(1) if there are two distinct productions  $A \rightarrow \alpha \mid \beta$  with the following conditions hold:

- \* For no terminal  $a$  and  $B$  derive strings beginning with  $a$ .
- \* At most one of  $\alpha$  and  $\beta$  can derive empty string.
- \* If  $\beta \xrightarrow{*} \epsilon$  then  $\alpha$  does not derive any string beginning with a terminal in FOLLOW(A).
- \* If  $\alpha \xrightarrow{*} \epsilon$  then  $\beta$  does not derive any string beginning with a terminal in FOLLOW(A).



(LL(1) parser is designed using stack data structure explicitly to hold grammar symbols.

— Left-recursion is eliminated

— Common prefixes are also eliminated (left-factoring).

### Eliminating left-recursion

A grammar is left recursion if it has a production of the form  $A \rightarrow A\alpha$ , for some string  $\alpha$ .

To eliminate left-recursion for the production

$$A \rightarrow A\alpha | B$$

Rule:

$$\begin{aligned} A &\rightarrow BA' \\ A' &\rightarrow \alpha A' | \epsilon \end{aligned}$$

$$\begin{aligned} A &\rightarrow BA' \\ A' &\rightarrow \alpha A' | \epsilon \end{aligned}$$

ex:

$$A \rightarrow A\alpha_1 | A\alpha_2 | \dots | B_1 | B_2 | \dots | B_m$$

solution

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

$$A' \rightarrow \alpha_1 A' | \alpha_2 A' | \dots | \epsilon$$

ex:

$$S \rightarrow Aa | b$$

$$A \rightarrow Ac | Sd | \epsilon$$

$$\begin{aligned} A &\rightarrow BA' \\ A &\rightarrow bA' \\ A' &\rightarrow \alpha A' | \epsilon \end{aligned}$$

Sol:

Substitute the productions of S in A,

$$A \rightarrow Ac | Aad | bd | \epsilon$$

After eliminating left-recursion

$$S \rightarrow Aa | b$$

$$A \rightarrow bdA' | A'$$

$$A' \rightarrow cA' | adA' | \epsilon$$

$$\begin{aligned} S &\rightarrow Aa | b \\ A &\rightarrow bdA' | A' \\ A' &\rightarrow cA' | adA' | \epsilon \end{aligned}$$

## Left factoring

This can be done through rewriting the production until enough of the input has been seen

Rule

$$\begin{aligned} A &\rightarrow \alpha B_1 \mid \alpha B_2 \\ A &\rightarrow \alpha A' \\ A' &\rightarrow B_1 \mid B_2 \end{aligned}$$

ex:

$$A \rightarrow \alpha B_1 \mid \alpha B_2 \mid \dots \mid \alpha B_m \mid \gamma$$

Sol:

$$A \rightarrow \alpha A' \mid \gamma$$

$$A' \rightarrow B_1 \mid B_2 \mid \dots \mid B_m$$

ex:

$$S \rightarrow \frac{iEtS}{E \rightarrow b} \mid \frac{iEtSs}{E \rightarrow b} \mid a$$

Sol:  $S \rightarrow iEtSS' \mid a$

$$S' \rightarrow es \mid \epsilon$$

$$E \rightarrow b$$

## Computation of First

Rules:

To compute  $FIRST(X)$ , where  $X$  is a grammar symbol.

- If  $X$  is a terminal, then  $FIRST(X) = \{X\}$
- If  $X \rightarrow \epsilon$  is a production, then add  $\epsilon$  to  $FIRST(X)$
- if  $X$  is a non-terminal and  $X \rightarrow Y_1 Y_2 \dots Y_k$  is a production, then add  $FIRST(Y_1)$  to  $FIRST(X)$ . If  $Y_1$  derives  $\epsilon$ , then add  $FIRST(Y_2)$  to  $FIRST(X)$ .



## Computation of FOLLOW

(51)

- For the FOLLOW (Start Symbol) Place \$, where \$ is the input end marker.
- If there is a production  $A \rightarrow \alpha BB$ , then everything in  $FIRST(B)$  except  $\epsilon$  is in  $FOLLOW(B)$
- If there is a production  $A \rightarrow \alpha B$ , or a production  $A \rightarrow \alpha BB$  where  $FIRST(B)$  contains  $\epsilon$ , then everything in  $FOLLOW(A)$  is in  $FOLLOW(B)$ .

## Construction of parsing table

Input Grammar  $G$

Output parsing table  $M$

Method For each production  $A \rightarrow \alpha$ , do the following:

1. For each terminal  $a$  in  $FIRST(\alpha)$ , add  $A \rightarrow \alpha$  to  $M[A, a]$
2. If  $\epsilon$  is in  $FIRST(\alpha)$ , then for each terminal  $b$  in  $FOLLOW(A)$ , add  $A \rightarrow \alpha$  to  $M[A, b]$ .
3. If  $\epsilon$  is in  $FIRST(\alpha)$  and \$ is in  $FOLLOW(A)$ , add  $A \rightarrow \alpha$  to  $M[A, \$]$ .
4. If no production is found in  $M[A, a]$  then set error to  $M[A, a]$ .

## Parsing of input

- \* Stack - holds sequence of grammar symbols with \$ on the bottom of stack
- \* Input buffer - contains the input to be parsed with \$ as end marker for the string.
- \* Parsing table.

Alg:

Input: A string  $w$  and parsing table  $M$  for Grammar  $G$

Output: If  $w$  is in  $L(G)$  then Success; Otherwise error

Method

Let  $a$  be the first symbol of  $w$ ; let  $x$  be the top of stack symbol; while  $(x \neq \$)$

if  $(x = a)$  Pop the stack and let  $a$  be the next symbol of  $w$ ;

else if  $(x \text{ is a terminal})$  error();

else if  $(M[x, a] \text{ is an error entry})$  error();

else if  $(M[x, a] = X \rightarrow Y_1 Y_2 \dots Y_k)$  {

Output the production  $X \rightarrow Y_1 Y_2 \dots Y_k$ ;

Pop the stack;

Push  $Y_k, Y_{k-1}, \dots, Y_1$  onto stack with  $Y_1$  on the top; }

Let  $x$  be the top stack symbol; }

Non - recursive Predictive Parser

- This prevents implicit recursive calls.

- It can also be termed as table-derived Predictive parser.

Components

\* Input buffer - holds input string to be parsed.

\* Stack - holds sequence of Grammar symbols.

\* predictive parsing algorithm - contains steps to parse the input string, controls the process.





57

Step 3: Compute FIRST

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

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

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

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

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

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

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

↓

$$E \rightarrow TE', \text{ FIRST}(E) = \text{FIRST}(T)$$

↓

$$T \rightarrow FT', \text{ FIRST}(T) = \text{FIRST}(F)$$

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

$F \rightarrow (E), \text{ FIRST}(F)$   
 $F \rightarrow \text{id}, \text{ FIRST}(F)$

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

$$\text{Hence, FIRST}(E) = \text{FIRST}(T) = \text{FIRST}(F) = \{\text{c}, \text{id}\}$$

Rule: If a non-terminal derives  $E$ , then  $E$  is in  $\text{FIRST}(\text{non-terminal})$   
 If  $A \rightarrow E$ , then  $\text{FIRST}(A) = \{E\}$

$\text{FIRST}(E') = \{+, E\} \Rightarrow E' \rightarrow +TE'$  is a production starting with a terminal '+',  $\text{FIRST}(E')$  include '+'

$\text{FIRST}(T') = \{*, E\} \Rightarrow T' \rightarrow *FT'$  is a production starting with terminal '\*',  $\text{FIRST}(T')$  include '\*'

Step 4: Compute FOLLOW

Rule: Place  $\phi$  in  $\text{FOLLOW}(\text{start symbol})$

$$\text{FOLLOW}(E) = \{\phi, \text{d}\}$$

Since  $F \rightarrow (E)$  is a production where  $E$  is followed by a terminal '(',  $\text{FOLLOW}[E]$  contain '('

Production  $F \rightarrow (E)$   
 $\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$   
 $A \quad \alpha \quad B \quad \beta$

$$\text{FOLLOW}(E) = \text{FIRST}(\beta) = \{\text{d}\}$$

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



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

Production

$$E \rightarrow T E' \quad (55)$$

$$\begin{array}{ccc} \downarrow & \downarrow & \downarrow \\ A & \alpha & B \end{array}$$

(i.e)

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

$$\text{Follow}(T) = \{+, \epsilon, \$\}$$

Production

$$E' \rightarrow + T E'$$

$$\begin{array}{cccc} \downarrow & \downarrow & \downarrow & \downarrow \\ A & \alpha & B & P \end{array}$$

(i.e)

$$\text{Follow}(T) = \text{FIRST}(P) = \text{FIRST}(E')$$

$$= \{+, \epsilon\}$$

(i.e)

$$\text{Follow}(B) = \text{Follow}(A) \cup \{\text{FIRST}(P) - \epsilon\}$$

$$\text{Follow}(T) = \text{Follow}(E') \cup \{\text{FIRST}(E') - \epsilon\}$$

$$= \{\$, \epsilon\} \cup \{+, \epsilon\} - \epsilon$$

$$= \{\$, \epsilon, +\}$$

Rules :

1)  $\$$

2) if  $A \rightarrow \alpha P \beta$

3)  $A \rightarrow \alpha B$

$$= \{\epsilon\}$$

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

Production

$$T \rightarrow F T'$$

$$\begin{array}{ccc} \downarrow & \downarrow & \downarrow \\ A & \alpha & B \end{array}$$

$$\text{ie } \text{Follow}(B) = \text{Follow}(A)$$

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

$$\text{Follow}(F) = \{+, \epsilon, \$, *\}$$

Production

$$T' \rightarrow * F T'$$

$$\begin{array}{cccc} \downarrow & \downarrow & \downarrow & \downarrow \\ A & \alpha & B & P \end{array}$$

$$\text{ie. } \text{Follow}(F) = \text{FIRST}(P) = \text{FIRST}(T') = \{*, \epsilon\}$$

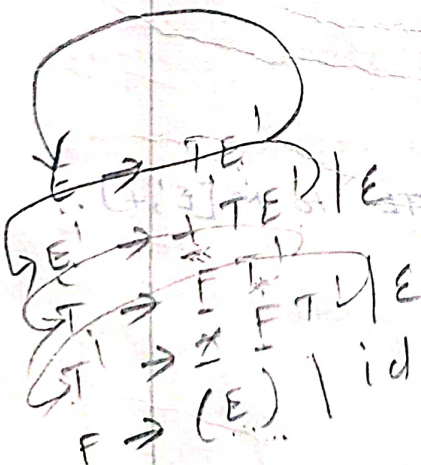
Since  $\epsilon$  is in  $\text{FIRST}(T')$ , everything in  $\text{Follow}(T')$  is in  $\text{Follow}(F)$ .

$$\text{ie. } \text{Follow}(B) = \text{Follow}(A) \cup \{\text{FIRST}(P) - \epsilon\}$$

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

$$= \{\$, \epsilon, +\} \cup \{*, \epsilon\} - \epsilon$$

$$= \{\$, \epsilon, +, *\}$$



Step 5. Construct parsing table.

Non-terminal	Input symbol				
	id	+	*	(	)
E	$E \rightarrow TE'$			$E \rightarrow TE'$	
E'		$E' \rightarrow +TE'$			$E' \rightarrow \epsilon$
T	$T \rightarrow FT'$			$T \rightarrow FT'$	
T'		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$
F	$F \rightarrow id$			$F \rightarrow (E)$	

To make over parsing table entry, analyze each and every production.

For the production

$$E \rightarrow T E'$$

$$\begin{matrix} \downarrow & \downarrow \\ A & \alpha \end{matrix}$$

$$FIRST(TE') = FIRST(T) = \{(, id)\}$$

add the production  $E \rightarrow TE'$  to  $M[E, (]$  and  $M[E, id]$

For the production,

$$E' \rightarrow + TE'$$

$$\begin{matrix} \downarrow & \downarrow \\ A & \alpha \end{matrix}$$

$$FIRST(+ TE') = FIRST(+) = \{+\}$$

So, add the production  $E' \rightarrow + TE'$  to  $M[E', +]$

For the production

$$E' \rightarrow \epsilon$$

$$\begin{matrix} \downarrow & \downarrow \\ A & \alpha \end{matrix}$$

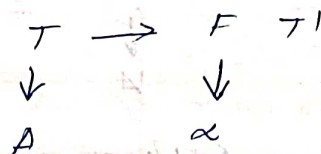
$$FIRST(\epsilon) = \{\epsilon\}$$

Since,  $FIRST(\alpha)$ , i.e.  $FIRST(\epsilon)$  contains  $\epsilon$ , add the production to  $M[A, FOLLOW(A)]$ ,  $FOLLOW(E') = \{), \$\}$

So, add the production  $E' \rightarrow \epsilon$  to  $M[E', )]$  and  $M[E', \$]$



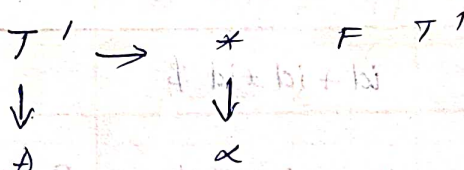
For the production,



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

add the production  $T \rightarrow FT'$  to  $M[T, (]$  and  $M[T, id]$ .

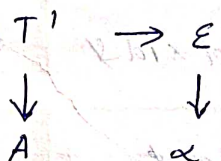
For the production



$$\text{FIRST}(FT') = \text{FIRST}(*) = \{ *, id \}$$

add the production  $T \rightarrow FT'$  to  $M[T, (]$  and  $M[T, id]$   
 $T' \rightarrow * FT' + \alpha M[T, *]$

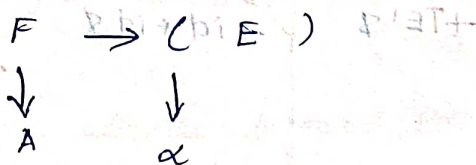
For the production



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

add the production  $T' \rightarrow \epsilon$  to  $M[T', +]$ ,  $M[T', )]$  and  $M[T', \$]$ .

For the production



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

add the production  $F \rightarrow (E)$  to  $M[F, (]$ .

For the production

$$\begin{array}{ccc} F & \rightarrow & id \\ \downarrow & & \downarrow \\ A & & \alpha \end{array}$$

$FIRST(id) = \{id\}$

add the production  $F \rightarrow id$  to  $M[F, id]$

Step 6: Parse the given input

$w = id + id * id$

$\begin{array}{|c|} \hline id \\ \hline + \\ \hline id \\ \hline * \\ \hline id \\ \hline \$ \end{array}$

Stack	Input	Action
$E \$$	$id + id * id \$$	
$TE \$$	$id + id * id \$$	Replace $E$ by its production $E \rightarrow TE'$ in $M[E, id]$
$FT'E \$$	$id + id * id \$$	Replace $T$ by its production $T \rightarrow FT'$ in $M[T, id]$
$id + TE' \$$	$id + id * id \$$	Replace $F$ by its production $F \rightarrow id$ in $M[F, id]$
$id + T'E' \$$	$id + id * id \$$	Match $id$
$T'E' \$$	$+ id * id \$$	
$E'E' \$$	$+ id * id \$$	Replace $T'$ by its production $T' \rightarrow E$ in $M[T', +]$
$E' \$$	$+ id * id \$$	
$+TE' \$$	$+ id * id \$$	Replace $E'$ by its production $E' \rightarrow +TE'$ in $M[E', +]$
$+TE' \$$	$+ id * id \$$	Match $+$
$TE' \$$	$id * id \$$	

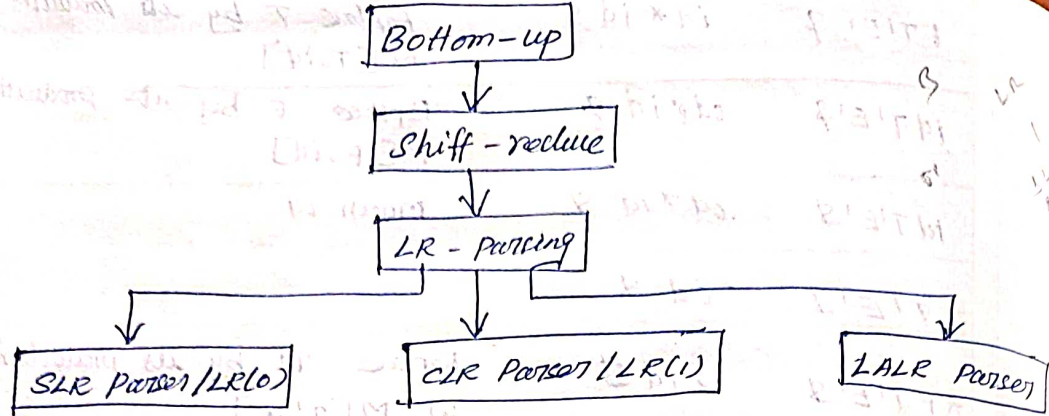


$FT E \$$	$id * id \$$	Replace $T$ by its production $T \rightarrow FT'$ in $M[T, id]$ (59)
$idT E \$$	$id * id \$$	Replace $F$ by its production $F \rightarrow id$ in $M[F, id]$
$idT E \$$	$id * id \$$	Match $id$
$T E \$$	$* id \$$	
$*FT E \$$	$* id \$$	Replace $T'$ by its production $T' \rightarrow *FT'$ in $M[T', *]$
$*FT E \$$	$* id \$$	Match $*$
$FT E \$$	$id \$$	
$idT E \$$	$id \$$	Replace $F$ by its production $F \rightarrow id$ in $M[F, id]$
$idT E \$$	$id \$$	Match $id$
$T E \$$	$\$$	
$E E \$$	$\$$	Replace $F$ by its production $T \rightarrow E$ in $M[T, \$]$
$E E \$$	$\$$	
$E\$$	$\$$	Replace $E'$ by its production $E' \rightarrow E$ in $M[E', \$]$
$\$$	$\$$	Successful parsing i.e. accept the string.

### BOTTOM - UP PARSING

- \* Bottom-up parsers construct parse trees starting from the leaves and work up to the root.
- \* Bottom-up syntax analysis is also termed as shift reduce parsing.
- \* The common method of shift-reduce parsing is called LR parsing.
- \* Operator precedence parsing is an easy-to-implement shift-reduce parser.





## Classification of bottom-up parsing

### Handles

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

ex:

grammar

- $S \rightarrow aABe$
- $A \rightarrow Abc \mid b$
- $B \rightarrow d$

The sentence  $abbcd$  can be reduced to  $S$  by the following steps.

$abbcd$   
 $aAbcd$   
 $aAde$   
 $aABe$   
 $S$

These reductions trace out the following  $r_m$  derivation in reverse.

$S \xrightarrow{r_m} aABe \xrightarrow{r_m} aAde \xrightarrow{r_m} aAbcd \xrightarrow{r_m} abbcd$

### Handle pruning

\* If  $A \rightarrow B$  is a production then reducing  $B$  to  $A$  by the given production is called handle pruning i.e. removing the children of  $A$  from the parse tree.



\* A rightmost derivation in reverse can be obtained by handle pruning.

ex:

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow (E)$$

$$E \rightarrow id$$

Right Sentential Form	Handle	Reduction production
$id_1 + id_2 * id_3$	$id_1$	$E \rightarrow id$
$E + id_2 * id_3$	$id_2$	$E \rightarrow id$
$E + E * id_3$	$id_3$	$E \rightarrow id$
$E + E * E$	$E * E$	$E \rightarrow E * E$
$E + E$	$E + E$	$E \rightarrow E + E$
$E$		

### Shift - reduce parsing

(i) shift-reduce parsing is a bottom-up parsing that reduces a string w to the start symbol of grammar

(ii) It scans and parses the input text in one forward pass without backtracking

### \* Stack implementation of shift-reduce parsing

\* Locating the substring to be reduced in a right sentential form.

\* Determining which production to choose in case there is more than one production with that substring on the right side.

②

## \* Implementation of shift reduce parser

- stack is used to hold grammar symbols
- An input buffer is used to hold the string to be parsed.

ex: consider grammar

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow (E)$$

$$E \rightarrow id$$

① Shift

② Reduce

③ Accept

④ Reject



Stack

input

Action

\$  
\$id<sub>1</sub>

id<sub>1</sub> + id<sub>1</sub> \* id<sub>3</sub> \$  
+ id<sub>2</sub> \* id<sub>3</sub> \$

shift id  
reduce by  $E \rightarrow id$

\$E

+ id<sub>2</sub> \* id<sub>3</sub> \$

shift +

\$E +

id<sub>2</sub> \* id<sub>3</sub> \$

shift id

\$E + E

\* id<sub>3</sub> \$

shift

\$E + E \*

id<sub>3</sub> \$

shift

\$E + E \* id<sub>3</sub>

\$

reduce by  $E \rightarrow id$

\$E + E \* E

\$

reduce by  $E \rightarrow E * E$

\$E + E

\$

reduce by  $E \rightarrow E + E$

\$E

\$

accept

## \* Viable Prefixes

The set of prefixes of right sentential forms that can be appears on the stack of a shift - reduce parser are called Viable prefixes



(63)

\* Conflicts during shift-reduces parsing

Shift-reduce parsing cannot be used in Context Free Grammar. bottom up Parsers

LR PARSERS

LR(0) or SLR(1)	LR(1) or CLR	LALR
-----------------------	--------------------	------

LR PARSERS are used to parse the large class of Context Free Grammars.

This technique is called LR(k) parsing

\* L is left-to-right scanning of the input

\* R is for constructing a rightmost derivation in reverse

\* k is the number of input symbols of lookahead that are used in making parsing decisions.

0000

0000A

Algorithm an LR Parser

\* SLR(1) - Simple LR

- \* works on smallest class of grammar.
- \* Few number of states hence very small table
- \* Simple and Fast construction

\* LR(1) - LR parser

- \* It called Canonical LR parser
- \* works on complete set of LR(1) Grammar.
- \* Generates large table and large number of states
- \* Slow construction

\* LALR(1) - Lookahead LR parser

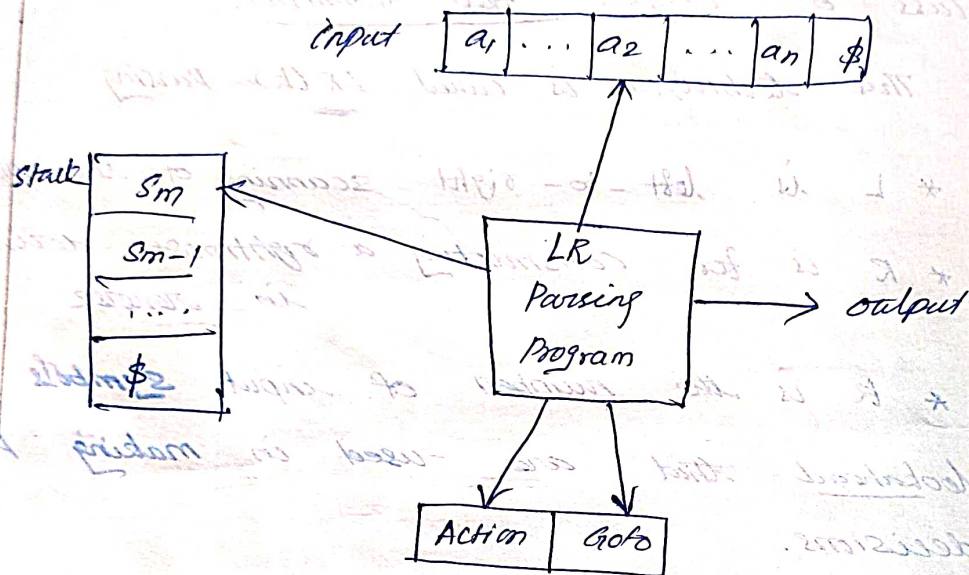
- \* works on intermediate size of grammar
- \* Number of states are same as in SLR(1).

## Model of LR Parser

LR Parser consists of an input, an output, a stack, a driver program and a parsing table that has two functions

1. Action

2. Goto



## Model of an LR Parser

### ACTION

This function takes as arguments a state  $i$  and a terminal  $a$  (or  $\$$ , the input end marker).

The value of  $ACTION[i, a]$  can have one of the four forms:

- (i) shift  $j$ , where  $j$  is a state
- (ii) reduce by a grammar production  $A \rightarrow \beta$
- (iii) accept
- (iv) error

### GOTO

This function takes a state and grammar symbol as arguments and produces a state.

If  $GOTO[i, A] = j$ , the Goto also maps a state  $i$  and non-terminal  $A$  to state  $j$ .



## \* LR Parsing Algorithm

(Refer book 164 Yedee Publication)

(65)

### \* LR(0) Items

An LR(0) item of a grammar  $G$  is a production of  $G$  with a dot at some position of the body.

ex.  $A \rightarrow \cdot XYZ$

$A \rightarrow X \cdot YZ$

$A \rightarrow XY \cdot Z$

$A \rightarrow XYZ \cdot$

One collection of set of LR(0) items. Called Canonical LR(0) collection, provides finite automaton that is used to make parsing decisions.

Such an automaton is called an LR(0) automaton.

### \* LR(0) parser / SLR(1) parser

An LR(0) parser is a shift-reduce parser that uses zero tokens of lookahead to determine what action to take (hence the 0.)

This means that in any configuration of the parser,

The parser must have an unambiguous action to choose - either it shifts a specific symbol or applies a specific reduction.

If there are ever two or more choices to make, the parser fails and the grammar is not LR(0).



(6)

Example:Construct LR(0) item for the grammarLR(0) or SLR(1) $E \rightarrow E + T$  $E \rightarrow T$  $T \rightarrow id$  $T \rightarrow (E)$  $Follow(E) = \{ +, \$ \}$  $Follow(T) = \{ +, \$ \}$ Solution:Step 1: Number the productions1.  $E \rightarrow E + T$ 2.  $E \rightarrow T$ 3.  $T \rightarrow id$ 4.  $T \rightarrow (E)$ Step 2: Construct augmented grammar by introducing a new start symbol  $E'$  and a new start production  $E' \rightarrow E$  $E' \rightarrow E$  $E \rightarrow E + T$  $E \rightarrow T$  $T \rightarrow id$  $T \rightarrow (E)$ Step 3: Find closure  $(E' \rightarrow \bullet E)$  $I_0$  $E' \rightarrow \bullet E$  $E \rightarrow \bullet E + T$  $E \rightarrow \bullet T$  $T \rightarrow \bullet id$  $T \rightarrow \bullet (E)$ 

List down symbol next to dot in each production

 $V = \{ E, T \}$  $T = \{ id, ( \}$



Step 4: Find goto for itemset  $I_0$

$goto(I_0, E)$	$E' \rightarrow E \bullet$ <i>acc</i> $E \rightarrow E \bullet + T$	$I_1$
$goto(I_0, T)$	$E \rightarrow T \bullet$	$I_2$
$goto(I_0, id)$	$T \rightarrow id \bullet$	$I_3$ <i>state</i>
$goto(I_0, ($	$T \rightarrow ( \bullet E )$ $E \rightarrow \bullet E + T$ $E \rightarrow \bullet T$ $T \rightarrow \bullet id$ $T \rightarrow \bullet ( E )$	$I_4$

Step 5: Find goto for itemset  $I_1$

$goto(I_1, +)$	$E \rightarrow E + \bullet T$ $T \rightarrow \bullet id$ $T \rightarrow \bullet ( E )$	$I_5$
----------------	--	-------

Step 6 Find goto for itemset  $I_4$

$goto(I_4, E)$	$T \rightarrow ( E \bullet )$ $E \rightarrow E \bullet + T$	$I_6$
$goto(I_4, T)$	$E \rightarrow T \bullet$	$I_2$
$goto(I_4, id)$	$T \rightarrow id \bullet$	$I_3$
$goto(I_4, ($	$T \rightarrow ( \bullet E )$ $E \rightarrow \bullet E + T$ $E \rightarrow \bullet T$ $T \rightarrow \bullet id$ $T \rightarrow \bullet ( E )$	$I_4$

(68)

Step 7 Find goto for itemset I5

goto (I5, T)	$E \rightarrow E + T \cdot$	I7
goto (I5, id)	$T \rightarrow id \cdot$	I3
goto (I5, (	$T \rightarrow (\cdot E)$ $E \rightarrow \cdot E + T$ $E \rightarrow \cdot T$ $T \rightarrow \cdot id$ $T \rightarrow \cdot (E)$	I4

Step 8: Find goto for itemset I6

goto (I6, ))	$T \rightarrow (E) \cdot$	I8
goto (I6, +)	$E \rightarrow E + \cdot T$ $T \rightarrow \cdot id$ $T \rightarrow \cdot (E)$	I5

Step 9: Find goto for itemset I6

goto (I6, ))	$T \rightarrow (E) \cdot$	I8
goto (I6, +)	$E \rightarrow E + \cdot T$ $T \rightarrow \cdot id$ $T \rightarrow \cdot (E)$	I5

Step 9: The LR(0) table of Grammar G1

\* For a shift, the token to be shifted determines the next state.

\* For a reduce, the state on top of stack specifies the production to be used.



State	Action					Goto	
	+	id	(	)	\$	E	T
0	+	$s_2$	$s_4$			1	2
1	$s_5$				accept		
2	$r_2$			$r_2$	$r_2$		
3	$r_3$			$r_3$	$r_3$		
4		$s_3$	$s_4$			6	2
5		$s_3$	$s_4$				7
6	$s_5$			$s_8$			
7	$r_1$			$r_1$	$r_1$		
8	$r_4$			$r_4$	$r_4$		

### Non-terminals

$$\text{goto}(I_0, E) = I_1 \quad \text{goto}(0, E) = 1$$

$$\text{goto}(I_0, T) = I_2 \quad \therefore \text{goto}(0, T) = 2$$

$$\text{goto}(I_4, E) = I_6 \quad \text{goto}(4, E) = 6$$

$$\text{goto}(I_4, T) = I_2 \quad \therefore \text{goto}(4, T) = 2$$

$$\text{goto}(I_5, T) = I_7 \quad \therefore \text{goto}(5, T) = 7$$

$$(I_1, E) = E \rightarrow E \rightarrow \text{reduce } r_1 = \text{follow}(E)$$

$$(I_0, T) = E \rightarrow T \rightarrow \text{follow}(E) = \{ \$, +, ) \}$$

### Terminals

$$\text{goto}(I_0, id) = I_3$$

$$\text{goto}(I_0, ( ) = I_4$$

$$\text{goto}(I_1, + ) = I_5$$

$$\text{goto}(I_4, id) = I_3$$

$$\text{goto}(I_4, ( ) = I_4$$

$$\text{goto}(I_5, id) = I_5$$

$$\text{goto}(I_6, + ) = I_5$$

$$\text{goto}(I_6, ( ) = I_8$$

$$\text{action}(0, id) = s_3$$

$$\text{action}(0, ( ) = s_4$$

$$\text{action}(1, + ) = s_5$$

$$\text{action}(4, id) = s_3$$

$$\text{action}(4, ( ) = s_4$$

$$\text{action}(5, id) = s_5$$

$$\text{action}(6, + ) = s_5$$

$$\text{action}(6, ( ) = s_8$$



# Computation of FOLLOW

For the inclusion

$$\begin{array}{cccc} E & \rightarrow & E & T & E \\ \downarrow & & \downarrow & \downarrow & \downarrow \\ A & & \alpha & B & B \end{array}$$

Since  $FIRST(B)$  contains  $\epsilon$ ,  $FOLLOW(B) = FOLLOW(A)$   
i.e.,  $FOLLOW(T) = FOLLOW(E) = \{\$, +, \}$

Computing FOLLOW(T)

Step 10: parser table on input id+id+id.

Stack	symbols	input	Action
0	$\epsilon$	id+id+id\$	$action[0, id] = s_3$ , shift 3, Push id then 3
03	id	id+id+id\$	$action[3, +] = r_3$ , i.e. 3rd production reduce by $T \rightarrow id$ Pop $[B]$ symbol from stack i.e. 1.
0	$T$	+id+id\$	$goto[0, T] = 2$
02	$T$	+id+id\$	$action[2, +] = r_2$ , i.e. 2nd production reduce by $E \rightarrow T$ Pop $[B]$ symbol from stack i.e. 1.
0	$E$	+id+id\$	$goto[0, E] = 1$ , push 1
01	$E$	+id+id\$	$action[1, +] = s_5$ , shift 5, Push + then 5
015	$E +$	id+id\$	$action[5, id] = s_3$ shift 3, Push id then 3

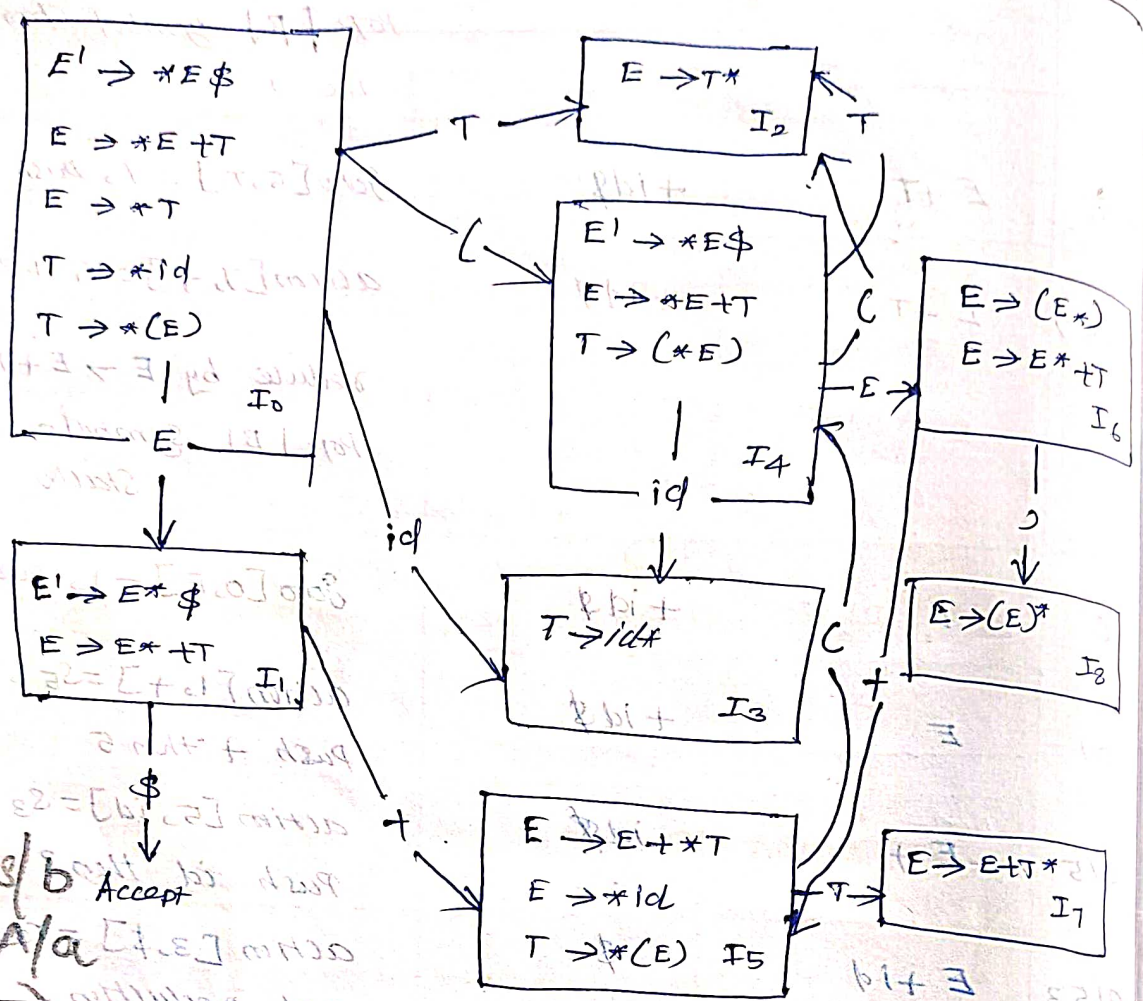


0153	$E + id$ DE+10153	+id\$	action[3,+] = $r_3$ , i.e., 3rd production reduce by $T \rightarrow id$ Pop  B  symbol from stack i.e. 1
015	$E + T$ DE+115	+id\$	goto[5,T] = 7, push 7
0157	$E + T$ DE+1157	+id\$	action[7,+] = $r_1$ , i.e. 1st production reduce by $E \rightarrow E + T$ Pop  B  symbols from stack i.e. 3
0	$E$ DE+1	+id\$	goto[0,E] = 1, push 1
01	$E$ DE+1	+id\$	action[1,+] = $s_5$ , shift 5 push + then 5
015	$E + T$ DE+115	+id\$	action[5,id] = $s_3$ , shift 3 push id then 3
0153	$E + id$ DE+10153	+id\$	action[3,\$] = $r_3$ , i.e. 3rd production, reduce by $T \rightarrow id$ Pop  B  symbol from stack i.e. 1
015	$E + T$ DE+115	\$	goto[5,T] = 7, push 7
0157	$E + T$ DE+1157	\$	action[7,\$] = $r_1$ , i.e. 1st production reduce by $E \rightarrow E + T$ Pop  B  symbols from stack i.e., 3
0	$E$ DE+1	\$	goto[0,E] = 1, push 1
01	$E$ DE+1	\$	action[1,\$] = accept Parser accept the input and halts.



2

# Step 11: Transition diagram of DFA for the grammar



HW

$S \rightarrow AS/b$  Accept

$A \rightarrow SA/a$

abab = w

## LR(1) Parser / Canonical LR (CLR)

\* Even more powerful than SLR(1) is the LR(1) parsing method.

\* LR(1) includes LR(0) items and a look ahead token in itemsets.

\* An LR(1) item consists

→ Grammar production rule

→ Right-hand position represented by the dot

→ Lookahead token

\* An LR(1) state is a set of LR(1) items.



## Example

(43)

Construct LR(1) items for the Grammar

$$\begin{aligned} S &\rightarrow CC \\ C &\rightarrow cCld \end{aligned}$$

Unit - 2

Solution:

Step 1: Number the productions

1.  $S \rightarrow CC$
2.  $C \rightarrow cC$
3.  $C \rightarrow d$

→ LR(0) -  
→ SLR -  
→ CLR  
→ LALR

Step 2: Construct augmented grammar by introducing

a production of the form

$S' \rightarrow S$  where  $S$  is the start symbol

$$S' \rightarrow S$$

$$S \rightarrow CC$$

$$C \rightarrow cC$$

$$C \rightarrow d$$

Step 3: Find closure ( $S' \rightarrow \bullet S, \$$ )

$I_0$

$$\begin{aligned} S' &\rightarrow \bullet S, \$ \\ S &\rightarrow \bullet C, \$ \\ C &\rightarrow \bullet c, Cld \\ C &\rightarrow \bullet d, cld \end{aligned}$$

List down symbol next to dot in each production

$$V = \{S, c\}$$

$$T = \{C, d\}$$

Note:

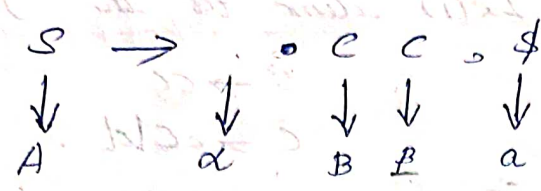
If a production is of the form  $A \rightarrow \alpha \bullet B \beta$ , then its lookahead is computed by  $FIRST(\beta \alpha)$ .

$$\begin{array}{c} S' \rightarrow \epsilon \bullet S \epsilon, \$ \\ \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \\ A \quad \alpha \quad B \quad \beta \quad a \end{array}$$

$$FIRST(\beta \alpha) = FIRST(\epsilon \$) = \{\$ \}$$

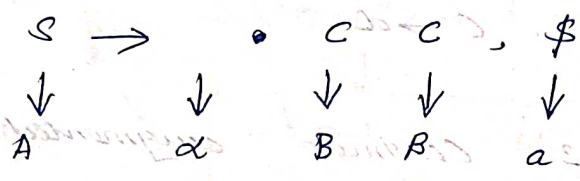
94

Production  $C \rightarrow \cdot c C, \$$



$$FIRST(Ba) = FIRST(c\$) = FIRST(c) = \{c, d\}$$

Production  $C \rightarrow \cdot d$



$$FIRST(Ba) = FIRST(c\$) = FIRST(c) = \{c, d\}$$

Step 4: Find goto for itemset  $I_0$

goto( $I_0, S$ )	$S' \rightarrow S \cdot \$$	$I_1$
goto( $I_0, c$ )	$S \rightarrow c \cdot C, \$$ $C \rightarrow \cdot c C, \$$ $C \rightarrow \cdot d, \$$	$I_2$ List down symbols next to dot in each production $V = \{c\}$ $T = \{c, d\}$
goto( $I_0, d$ )	$c \rightarrow c \cdot c, cld$ $c \rightarrow \cdot c c, cld$ $c \rightarrow \cdot d, cld$	$I_3$ List down symbols next to dot in each production $V = \{c\}$ $T = \{c, d\}$
goto( $I_0, d$ )	$c \rightarrow d \cdot, cld$	$I_4$

Step 5: Find goto for itemset  $I_2$

goto( $I_2, c$ )	$S \rightarrow c c \cdot \$$	$I_5$
goto( $I_2, c$ )	$C \rightarrow c \cdot C, \$$ $C \rightarrow \cdot c C, \$$ $C \rightarrow \cdot d, \$$	$I_6$
goto( $I_2, d$ )	$c \rightarrow d \cdot, cld$	$I_7$



Step 6: Find goto for item sets  $I_3$

goto( $I_3, c$ )	$C \rightarrow cC, cld$	$I_8$
goto( $I_3, c$ )	$C \rightarrow c \cdot C, cld$ $C \rightarrow \cdot cC, cld$ $C \rightarrow \cdot d, cld$	$I_3$
goto( $I_3, d$ )	$C \rightarrow d \cdot, cld$	$I_4$

Step 7: Find goto for item set  $I_6$

goto( $I_6, c$ )	$C \rightarrow cC, \$$	$I_9$
goto( $I_6, c$ )	$C \rightarrow c \cdot C, \$$ $C \rightarrow \cdot cC, \$$ $C \rightarrow \cdot d, \$$	$I_6$
goto( $I_6, d$ )	$C \rightarrow d \cdot, \$$	$I_7$

Step 8: All itemsets are processed, i.e., no more new item set.

Step 9: Construct parsing table

States	Action			Goto	
	c	d	\$ or	S	C
0	$S_3$	$S_4$		1	2
1			accept		
2	$S_6$	$S_7$			5
3	$S_3$	$S_4$			8
4	$r_3$	$r_2$			
5					
6	$S_6$	$S_7$			9
7			$r_3$		
8	$r_2$	$r_2$			
9			$r_2$		

(46)

Non-terminals

$$\text{goto}(I_0, s) = I_1 \quad \therefore \text{goto}(0, s) = 1$$

$$\text{goto}(I_0, c) = I_2 \quad \therefore \text{goto}(0, c) = 2$$

$$\text{goto}(I_2, c) = I_5 \quad \therefore \text{goto}(2, c) = 5$$

$$\text{goto}(I_3, c) = I_8 \quad \therefore \text{goto}(3, c) = 8$$

$$\text{goto}(I_6, c) = I_9 \quad \therefore \text{goto}(6, c) = 9$$

Terminals

$$\text{goto}(I_0, c) = I_3 \quad \therefore \text{action}(0, c) = s_3$$

$$\text{goto}(I_0, d) = I_4 \quad \therefore \text{action}(0, d) = s_4$$

$$\text{goto}(I_2, c) = I_6 \quad \therefore \text{action}(2, c) = s_6$$

$$\text{goto}(I_2, d) = I_7 \quad \therefore \text{action}(2, d) = s_7$$

$$\text{goto}(I_3, c) = I_3 \quad \therefore \text{action}(3, c) = s_3$$

$$\text{goto}(I_3, d) = I_4 \quad \therefore \text{action}(3, d) = s_4$$

$$\text{goto}(I_6, c) = I_6 \quad \therefore \text{action}(6, c) = s_6$$

$$\text{goto}(I_6, d) = I_7 \quad \therefore \text{action}(6, d) = s_7$$

HW ①

S → SS

S → a

S → ε

②

S → AA

A → Aa/b

OSI/INTRODUCTION TO LALR PARSER

87

\* LALR stands for lookahead LR parser

\* This is the extension of LR(0) items, by introducing the one symbols of lookahead on the input.

\* It supports large class of grammars.

\* The number of states in LALR parser is lesser than that of LR(1) parser.

\* LALR is preferable as it can be used with reduce memory

\* Most syntactic constructs of programming language can be easily converted to LALR.



## Steps to construct LR(1) parsing table

- \* Generate LR(1) items
- \* Find the items that have same set first component (core) and merge those sets into one.
- \* Merge the goto's of combined itemset
- \* Revise the parsing table of LR(1) parser by replacing states and goto's with combined states and combined goto's respectively.

ex: For the Grammar

$$1. S \rightarrow CC$$

$$2. C \rightarrow cC \mid d$$

Solution:

Step 1: Number the productions

$$1. S \rightarrow CC$$

$$2. C \rightarrow cC$$

$$3. C \rightarrow d$$

Step 2: Construct augmented grammar by introducing a production of the form  $S' \rightarrow S$  where  $S'$  is the start symbol

$$S' \rightarrow S$$

$$C \rightarrow cC$$

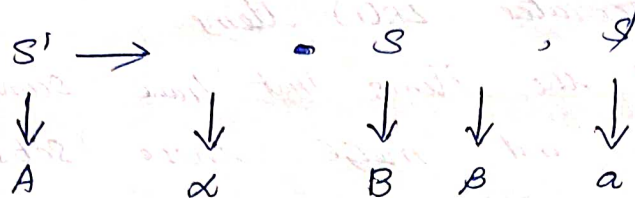
$$C \rightarrow d$$

Step 3: Find closure ( $S' \rightarrow \cdot S, \$$ )

$I_0$	$S' \rightarrow \cdot S, \$$ $S \rightarrow \cdot CC, \$$ $C \rightarrow \cdot cC, cId$ $C \rightarrow \cdot d, cId$	List down symbols next to dot in each module $V = \{S, C\}$ $T = \{c, d\}$
-------	---	--

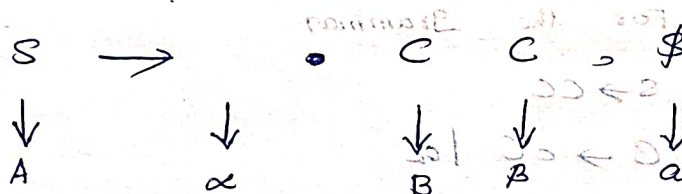
19

Production  $S \rightarrow \cdot CC$



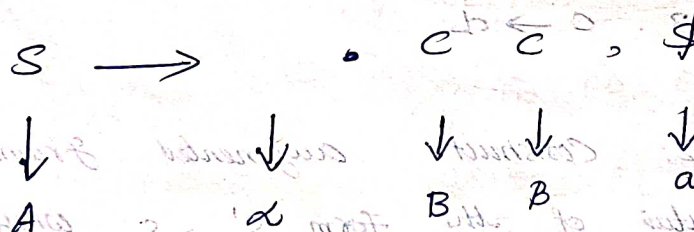
$$FIRST(B\alpha) = FIRST(\epsilon \$) = \{\$ \}$$

Production  $C \rightarrow \cdot cC$



$$FIRST(B\alpha) = FIRST(C\$) = FIRST(CC) = \{C, d\}$$

Production  $C \rightarrow \cdot d$



$$FIRST(B\alpha) = FIRST(C\$) = FIRST(C) = \{C, d\}$$

Step 4: Find goto for itemset  $I_0$

goto ( $I_0, S$ )	$S' \rightarrow S \cdot , \$$	$I_1$
goto ( $I_0, C$ )	$S \rightarrow C \cdot C, \$$ $C \rightarrow \cdot cC, \$$ $C \rightarrow \cdot d, \$$	$I_2$ List down symbols next to dot in each production $V = \{c\} \quad T = \{C, d\}$
goto ( $I_0, c$ )	$c \rightarrow c \cdot c, c   d$ $C \rightarrow \cdot cC, c   d$ $C \rightarrow \cdot d, c   d$	$I_3$ $V = \{c\} \quad T = \{C, d\}$
goto ( $I_0, d$ )	$c \rightarrow c \cdot d, c   d$	$I_4$



Step 5: Find goto for itemset  $I_2$

goto( $I_2, c$ )	$C \rightarrow cC, \$$	$I_5$
goto( $I_2, c$ )	$C \rightarrow c \cdot C, \$$ $C \rightarrow \cdot cC, \$$ $C \rightarrow \cdot d, \$$	$I_6$
goto( $I_2, d$ )	$C \rightarrow d \cdot, \$$	$I_7$

Step 6: Find goto for itemset  $I_3$

goto( $I_3, c$ )	$C \rightarrow cC, cld$	$I_8$
goto( $I_3, c$ )	$C \rightarrow c \cdot C, cld$ $C \rightarrow \cdot cC, cld$ $C \rightarrow \cdot d, cld$	$I_3$
goto( $I_3, d$ )	$C \rightarrow d \cdot, cld$	$I_4$

Step 7: Find goto for itemsets  $I_6$

goto( $I_6, c$ )	$C \rightarrow cC, \$$	$I_9$
goto( $I_6, c$ )	$C \rightarrow c \cdot C, \$$ $C \rightarrow \cdot cC, \$$ $C \rightarrow \cdot d, \$$	$I_6$
goto( $I_6, d$ )	$C \rightarrow d \cdot, \$$	$I_7$

Step 8: All itemsets are processed, no more new itemsets

8

Step 9: Construct parsing table

States	Action			Goto	
	C	d	\$	S	C
0	S <sub>3</sub>	S <sub>4</sub>		1	2
1			accept		
2	S <sub>6</sub>	S <sub>7</sub>			5
3	S <sub>3</sub>	S <sub>4</sub>			8
4	r <sub>3</sub>	r <sub>3</sub>			
5			r <sub>1</sub>		
6	S <sub>6</sub>	S <sub>7</sub>			9
7			r <sub>3</sub>		
8	r <sub>2</sub>	r <sub>2</sub>			
9			r <sub>2</sub>		

Non - terminals

$$\text{goto}(I_0, S) = I_1$$

$$\text{goto}(I_0, C) = I_2$$

$$\text{goto}(I_2, C) = I_5$$

$$\text{goto}(I_3, C) = I_8$$

$$\text{goto}(I_6, C) = I_9$$

$$\therefore \text{goto}(0, S) = 1$$

$$\therefore \text{goto}(0, C) = 2$$

$$\therefore \text{goto}(2, C) = 5$$

$$\therefore \text{goto}(3, C) = 8$$

$$\therefore \text{goto}(6, C) = 9$$

Terminals

$$\text{goto}(I_0, c) = I_3$$

$$\text{goto}(I_0, d) = I_4$$

$$\text{goto}(I_2, c) = I_6$$

$$\text{goto}(I_2, d) = I_7$$

$$\text{goto}(I_3, c) = I_3$$

$$\text{goto}(I_3, d) = I_4$$

$$\therefore \text{action}(0, c) = S_3$$

$$\therefore \text{action}(0, d) = S_4$$

$$\therefore \text{action}(2, c) = S_6$$

$$\therefore \text{action}(2, d) = S_7$$

$$\therefore \text{action}(3, c) = S_3$$

$$\therefore \text{action}(3, d) = S_4$$



$$\text{goto } (I_6, c) = I_6$$

$$\therefore \text{action } (6, c) = S_6 \quad (8)$$

$$\text{goto } (I_6, d) = I_7$$

$$\therefore \text{action } (6, d) = S_7$$

In LR(1) items

Item sets  $I_3$  and  $I_6$  have same core with different lookaheads, so merge  $I_3$  and  $I_6$  to form  $I_{36}$ .

Item sets  $I_4$  and  $I_7$  have same core with different lookaheads. So merge  $I_4$  and  $I_7$  to form  $I_{47}$ .

Item set  $I_8$  and  $I_9$  have same core with different lookaheads. So merge  $I_8$  and  $I_9$  to form  $I_{89}$ .

$I_3$ $C \rightarrow c \cdot C, cld$ $C \rightarrow \cdot cC, cld$ $C \rightarrow \cdot d, cld$	$I_6$ $C \rightarrow c \cdot C, \$$ $C \rightarrow \cdot cC, \$$ $C \rightarrow \cdot d, \$$	$I_{36}$ $C \rightarrow c \cdot C, cld   \$$ $C \rightarrow \cdot cC, cld   \$$ $C \rightarrow \cdot d, cld   \$$
$I_4$ $C \rightarrow d \cdot, cld$	$I_7$ $C \rightarrow d \cdot, \$$	$I_{47}$ $C \rightarrow d \cdot, cld   \$$
$I_8$ $C \rightarrow cC \cdot, cld$	$I_9$ $C \rightarrow cC \cdot, \$$	$I_{89}$ $C \rightarrow cC \cdot, cld   \$$

Parsing table

States	Action			Goto	
	c	cl	\$	non-terminal	terminal
0	$S_{36}$	$S_{47}$		1	2
1			accept		
2	$S_{36}$	$S_{47}$			5
36	$S_{36}$	$S_{47}$			89
47	$r_3$	$r_3$	$r_3$		
5					



(9)

# OPERATOR PRECEDENCE PARSER

## Operator Grammar

- No  $\epsilon$  in its right side of production
- No two adjacent non-terminals.

Operator precedence parser can parse only small class of grammars.

## Operator precedence relations

Operator precedence relation describes relation between pair of terminals which guides the selection of handles.

- Associativity and precedences of operators are known.
- unambiguous grammar is constructed for the language that reflects correct associativity and precedence in its parse tree.

operator precedence relation and its meaning

Relation	Meaning
$<$	a yields precedence to b
$>$	a takes precedence over b
$=$	a has same precedence as b

Algorithm as follows page no 193 :

## Precedence Functions.

precedence functions are used to map terminals to integers.



It eliminates the need for storing the table of precedence relations by encoding into functions.

For symbols  $a$  and  $b$

- \*  $f(a) < g(b)$  whenever  $a < b$
- \*  $f(a) > g(b)$  whenever  $a > b$
- \*  $f(a) = g(b)$  whenever  $a = b$

ex:

$E \rightarrow E + E \mid E * E \mid id$

input string  $id + id * id$

Step 1:

Construct operator precedence relation table

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

Note:

- Always  $id$  takes precedence over other terminals
- $\$$  yields precedence to other terminals
- Other terminals have precedence relations based on associativity and precedence.

Step 2: Parsing table

Stack	input	Action
\$	$id + id * id \$$	$a = \$, b = id \Rightarrow a < b$ so push $b(id)$
$id \$$	$+ id * id \$$	$a = id, b = + \Rightarrow a > b$ so pop from stack

Ⓞ	\$	+id*id\$	Compare element in stack and the most recently popped element $\Rightarrow$ $\$ < id$ , so stop pop
	\$	+id*id\$	$a = \$$ , $b = + \Rightarrow a < b$ so push $b (+)$
	+ \$	id*id\$	$a = +$ , $b = id \Rightarrow a < b$ so Push $b (id)$
	id + \$	*id\$	$a = id$ , $b = * \Rightarrow a > b$ so Pop from stack ( $id$ )
	+ \$	*id\$	$a = +$ , $b = id \Rightarrow a < b$ stop Pop
	+ \$	*id\$	$a = +$ , $b = * \Rightarrow a < b$ push $b (*)$
	* + \$	id \$	$a = *$ , $b = id \Rightarrow a < b$ push $b (id)$
	id * + \$	\$	$a = id$ , $b = \$ \Rightarrow a > b$ pop from stack ( $id$ )
	* + \$	\$	$a = *$ , $b = id \Rightarrow a < b$ stop pop
	* + \$	\$	$a = *$ , $b = \$ \Rightarrow a > b$ pop from stack ( $*$ )
	+ \$	\$	$a = +$ , $b = * \Rightarrow a < b$ stop pop
	+ \$	\$	$a = +$ , $b = \$ \Rightarrow a > b$ pop from stack ( $+$ )
	\$	\$	$a = \$$ , $b = + \Rightarrow a > b$ stop pop
	\$	\$	$a = \$$ , $b = \$ \Rightarrow$ HALT

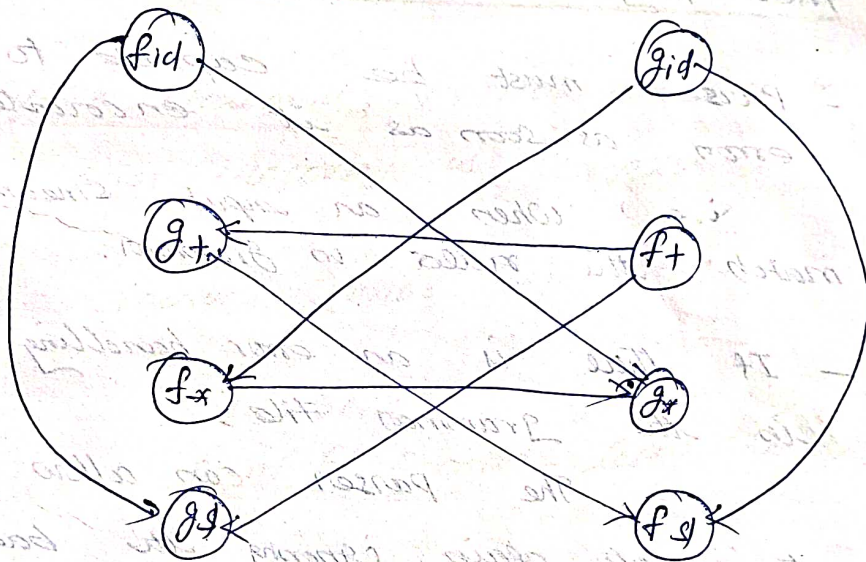


## For precedence Function.

(85)

1. Create symbol  $f_a$  and  $g_a$  for each terminal  $a$ , i.e.,  $f_{id}$ ,  $f_+$ ,  $f_*$ ,  $f_{\$}$ ,  $g_{id}$ ,  $g_+$ ,  $g_*$  and  $g_{\$}$ .
2. Create directed graph
3. If there are no cycles so precedence Function exists
4. Precedence function are given as follow.

The path of  $f_a$  is given by the grammar of edges from  $f_a$  to either  $f_{\$}$  or  $g_{\$}$ .



Path of  $f_+$  = 2  $\Rightarrow f_+ \rightarrow g_+ \rightarrow f_{\$}$

Path of  $f_*$  = 4  $\Rightarrow f_* \rightarrow g_* \rightarrow f_+ \rightarrow g_+ \rightarrow f_{\$}$

Path of  $f_{id}$  = 4  $\Rightarrow f_{id} \rightarrow g_* \rightarrow f_+ \rightarrow g_+ \rightarrow f_{\$}$

Path of  $f_{\$}$  = 0  $\Rightarrow f_{\$}$

Path of  $g_+$  = 1  $\Rightarrow g_+ \rightarrow f_{\$}$

Path of  $g_*$  = 3  $\Rightarrow g_* \rightarrow f_+ \rightarrow g_+ \rightarrow f_{\$}$

Path of  $g_{id}$  = 5  $\Rightarrow g_{id} \rightarrow f_* \rightarrow g_* \rightarrow f_+ \rightarrow g_+ \rightarrow f_{\$}$

Path of  $g_{\$}$  = 0  $\Rightarrow g_{\$}$

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

80

## Error Handling and Error Recovery in Syntax Analyzer

- An efficient program should not terminate on an parse error
- It must recover to parse the rest of the input and check the subsequent errors.
- For one line input, the routine `YYParse()` can be made to return 1 on error and then calls `YYParse()` again.

81

### YACC program error handling

- Parser must be capable to detecting the error as soon as it encounters.  
i.e.) When an input stream does not match the rules in grammar.
- If there is an error handling subroutine in the grammar file,  
The parser can allow for entering the data again, ignoring the bad data or initiating a cleanup and recovery action.
- When the parser finds an error, it may need to reclaim parse tree storage, delete or alter symbol table entries and set switches to avoid generating further output.
- Error handling routines are used to restart the parser to continue its process even after the occurrence of error.
- Tokens following the error get discarded to restart the parser.



e.g. state : error ;

(27)

### providing for error correction

The input errors can be corrected by entering a line in the data stream again.

```
input : error '\n'
```

```
{
```

```
printf ("Reenter last line :");
```

```
}
```

```
input
```

```
{
```

```
$$ = $4;
```

```
};
```

The yacc statement, yyerror is used to indicate that error recovery is complete.

```
input : error '\n'
```

```
{
```

```
yyerror;
```

```
printf ("Reenter last line :");
```

```
}
```

```
input
```

```
{
```

```
$$ = $4;
```

```
};
```

### clearing the lookahead token

- when an error occurs, the lookahead token becomes the token at which the error was detected.
- The lookahead token must be changed if the error recovery action includes code to find the correct place to start processing again.

(23)  
- To clear the lookahead token, the error recovery action issues the following statements.

YYclearin

To assist in error handling, macros can be placed in yacc action.

YYERROR	causes the parser to initiate error handling.
YYABORT	causes the parser to return with a value of 1.
YYACCEPT	causes the parser to return with a value of 0.
YYRECOVERING()	Returns a value of 1 if a syntax error has been detected and the parser has not yet fully recovered.

Macros for error handling.

YACC

- \* YACC stands for Yet Another Compiler-Compiler.
- \* It is a LALR parser generator.

Design of syntax analyzer / parser generator Yacc

Translator  $\rightarrow$  Yacc compiler  $\rightarrow$  y.tab.c

y.tab.c  $\rightarrow$  C compiler  $\rightarrow$  a.out

input  $\rightarrow$  a.out  $\rightarrow$  output

Creating translator using yacc



(89)

- YACC Specification of the translator (transl<sup>td</sup>.y) is given to YACC compiler which transform it into y.tab.c that represent LALR Parser in C program.

- The transformation of YACC specification into C program is done by the UNIX Command 'yacc yaccfile',

ie. yacc translate.y which uses the LALR method.

- It contains a routine yyparse() which returns 0 if the program is correct, non-zero otherwise.

- The object program a.out is obtained by the compilation of y.tab.c with ly library using the Command cc y.tab.c.ly.

- It also works with lex. YACC calls yylex() to get next token.

- YACC and lex must agree on the values for each token.

e.g. yyerror(str)

char \*str;

```
{ printf("yy-error: %s at line %d\n", str, yyline);
```

```
}
```

```
main()
```

```
{ if (!yyparse())
```

```
{ printf("accept\n");
```

```
}
```

```
else printf("reject\n");
```

```
}
```

# YACC File format

declarations

% %

translation rules

% %

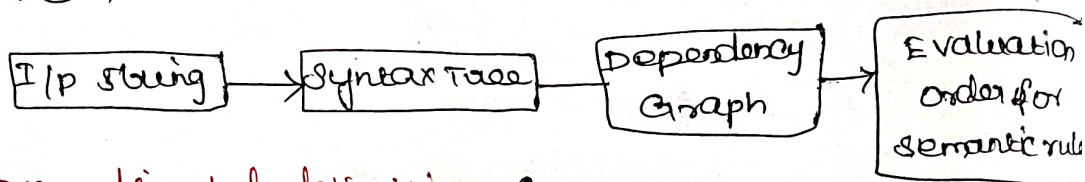
supporting C-routines



## Unit-III Syntax Directed Translation and Intermediate Code Generation

### ① Syntax directed Definitions: - augmented CFG is gen

It is a kind of abstract specification. The conceptual view of syntax directed translation, can be ,



Syntax -directed definition is a generalization of Context free grammar in which each grammar production  $x \rightarrow \alpha$  is associated with it a set of semantic rules of the form  $a := f(b_1, b_2, \dots, b_k)$  where  $a$  is an attribute obtained from the function  $f$ .

attribute: can be string, number, type,

$x \rightarrow \alpha$  be a CFG

$a := f(b_1, b_2, \dots, b_k)$  where  $a$  is attribute,

Two types attribute:

① Synthesized attribute: 'a' is called synthesized attribute of  $x$  and  $b_1, b_2, \dots, b_k$  are attributes belonging to the production symbols.

The value of synthesized attribute at a node is computed from the values of attributes at the children of that node.

② Inherited attribute: 'a' is called inherited attribute of one of the grammar symbol on the right side of the production ( $\alpha$ ) &  $b_1, b_2, \dots, b_k$  are belonging to either  $x$  or  $\alpha$ .

value is computed siblings & parent of the

## Synthesized attribute:

92

Eg:

$$S \rightarrow EN$$

$$E \rightarrow E + T$$

$$E \rightarrow E - T$$

$$E \rightarrow T$$

$$T \rightarrow T * F$$

$$T \rightarrow T / F$$

$$T \rightarrow F$$

$$F \rightarrow (E)$$

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

$$N \rightarrow ;$$

Sol:

production rule	Semantic actions
$S \rightarrow EN$	$\text{print}(E.\text{val})$
$E \rightarrow E_1 + T$	$E.\text{val} = E_1.\text{val} + T.\text{val}$
$E \rightarrow E_1 - T$	$E.\text{val} = E_1.\text{val} - T.\text{val}$
$E \rightarrow T$	$E.\text{val} := T.\text{val}$
$T \rightarrow T_1 * F$	$T.\text{val} := T_1.\text{val} * F.\text{val}$
$T \rightarrow T_1 / F$	$T.\text{val} := T_1.\text{val} / F.\text{val}$
$T \rightarrow F$	$T.\text{val} = F.\text{val}$
$F \rightarrow (E)$	$F.\text{val} := E.\text{val}$
$F \rightarrow \text{digit}$	$F.\text{val} := \text{digit}.\text{lexval}$
$N \rightarrow ;$	Can be ignored by lexical analyzer as ; is terminating symbol.

To compute S-attribute definition:

① write the syntax directed def. using the appropriate semantic actions for corresponding production rule.

② Annotated parse tree is generated & attribute values are computed. (bottom up manner)

③ The value is obtained at the root node is supposed to be the final o/p.

\_\_\_\_\_ x \_\_\_\_\_ x \_\_\_\_\_

Eg: Construct parse tree, syntax tree & annotated parse tree for the i/p string is  $5 * 6 + 7$ .

Sol:

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

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

$$F.\text{val} = 5$$

$$T.\text{val} = T_1.\text{val} * F.\text{val}$$

$$= 5 * 6$$

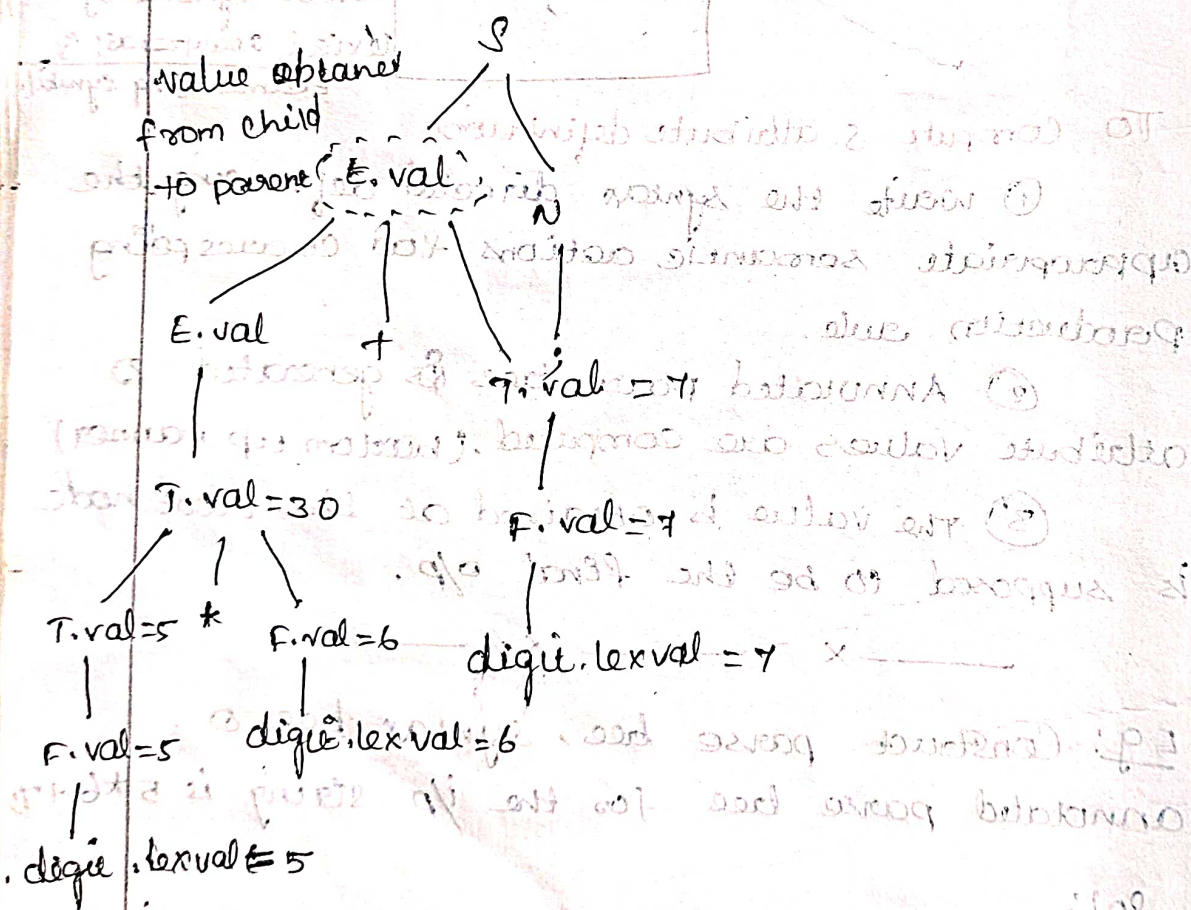
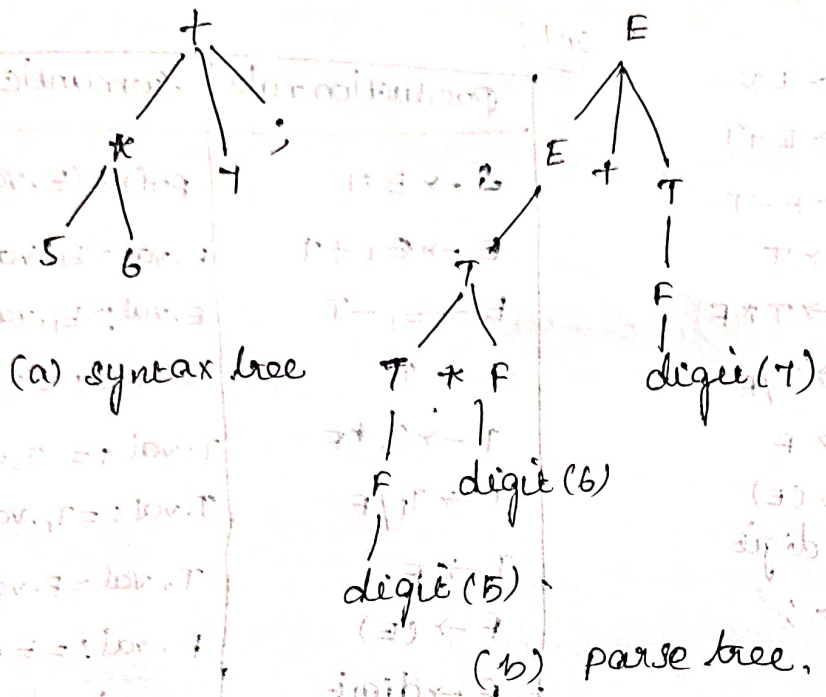
$$T.\text{val} = 30$$

$$E.\text{val} = E_1.\text{val} + T.\text{val}$$

$$= 30 + 7$$

$$E.\text{val} = 37$$

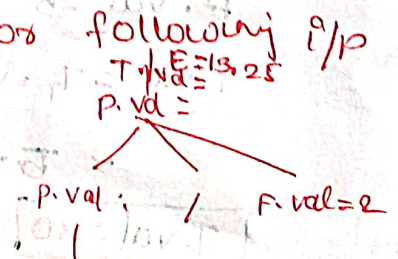




(c): Annotated parse tree.

2) Construct a decorated parse tree according to the syntax directed definition for following i/p

stmt:  $(4 + 7.5 * 3) / 2$



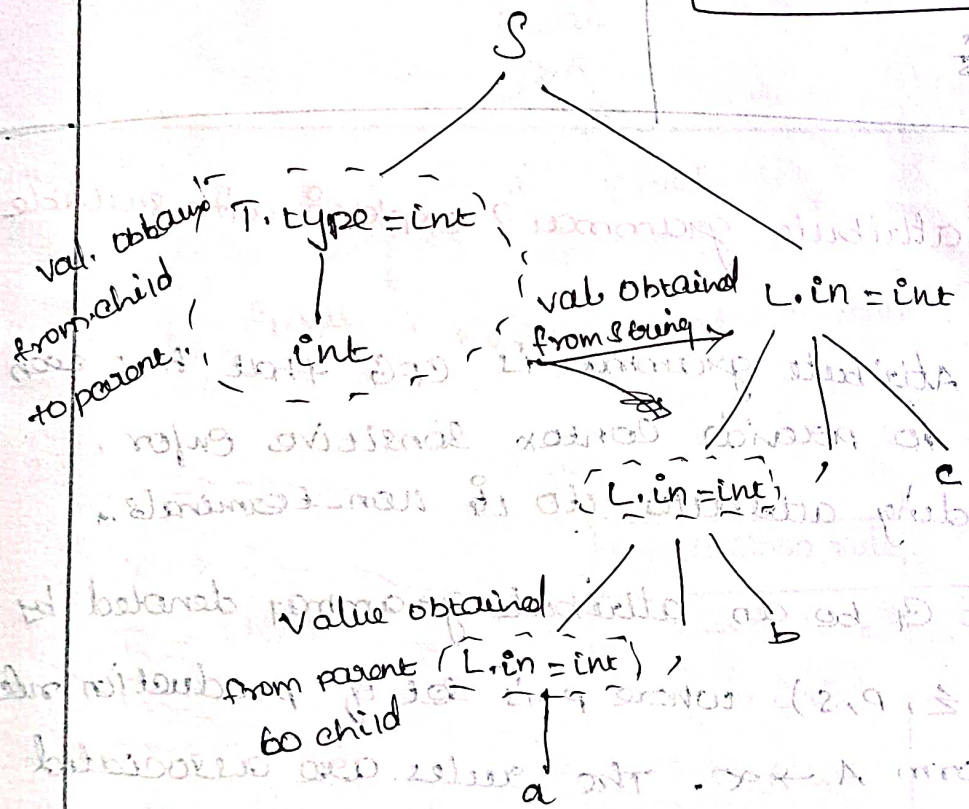
## ② Inherited attribute:

The value of inherited attribute at a node in a parse tree is defined using the attribute values at the parent or siblings.

$S \rightarrow TL$   
 $T \rightarrow \text{int}$   
 $T \rightarrow \text{float}$   
 $T \rightarrow \text{char}$   
 $T \rightarrow \text{double}$   
 $L \rightarrow L_1, \text{id}$   
 $L \rightarrow \text{id}$

production Rule	Semantic action
$S \rightarrow TL$	$L.in = T.type$
$T \rightarrow \text{int}$	$T.type := \text{Integer}$
$T \rightarrow \text{float}$	$T.type := \text{float}$
$T \rightarrow \text{char}$	$T.type := \text{char}$
$T \rightarrow \text{double}$	$T.type := \text{double}$
$L \rightarrow L_1, \text{id}$	$L_1.in = L.in$ $Enter.type(id, Enter, L.in)$
$L \rightarrow \text{id}$	$Enter.type(id, Enter, L.in)$

Annotated parse tree

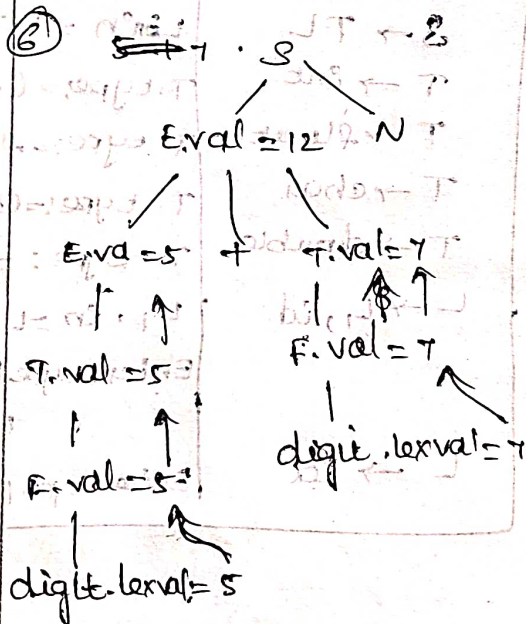


Synthesized translation	Inherited Translation
<ol style="list-style-type: none"> <li>① Definition</li> <li>② Syn. attributed are computed</li> <li>③ using values of children</li> </ol>	<ol style="list-style-type: none"> <li>① Definition</li> <li>② Inherited attribute are computed.</li> <li>③ " using values of parent &amp; siblings node</li> </ol>



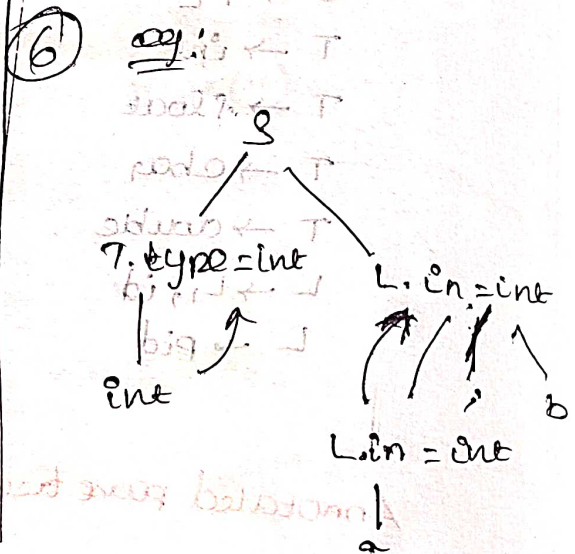
④ Info. is passed up in syntax tree

⑤ simple to compute the value



④ Info. is passed down in syntax tree.

⑤ Complex to compute the value.



What is attribute grammar? explain with suitable example.

Attribute grammar is CFG that has been extended to provide context sensitive error, by appending attributes to its non-terminals.

[ Let  $G$  be an attribute grammar denoted by  $G = (N, \Sigma, P, S)$  where  $P$  is set of production rules in the form  $A \rightarrow \alpha$ . The rules are associated with set of semantic rules of the form  $b_i = f(a_1, a_2, a_3, \dots, a_n)$  where  $f$  denotes functionality.

Two types:

① synthesized attribute

② Inherited attribute.

### 3) Dependency graph:

A directed graph that represents the interdependencies b/w synthesized & inherited attributes at nodes in the parse tree is called dependency graph.

$$E \rightarrow E_1 + E_2$$

$$E \rightarrow E_1 * E_2$$

Sol:

production rule      semantic rule

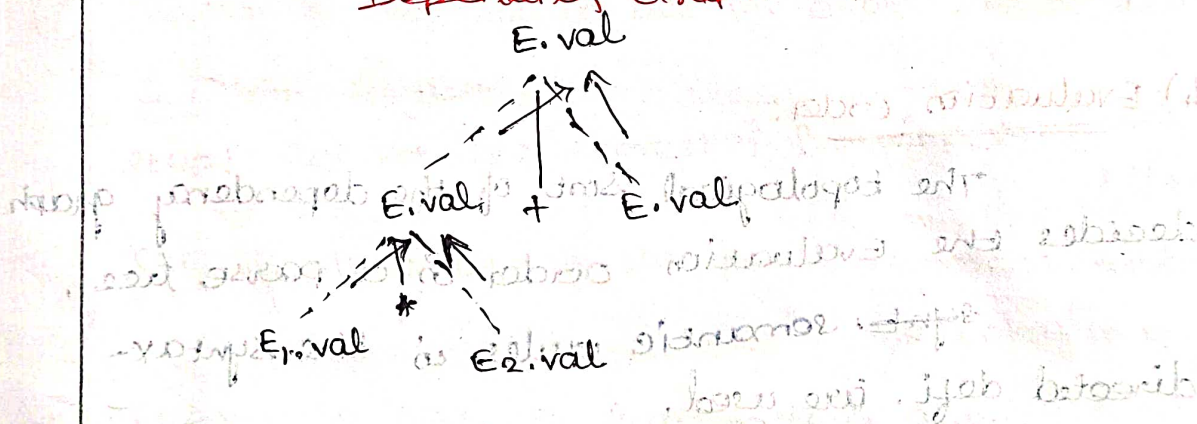
$$E \rightarrow E_1 + E_2$$

$$E.val := E_1.val + E_2.val$$

$$E \rightarrow E_1 * E_2$$

$$E.val := E_1.val \times E_2.val$$

Dependency Graph.



$$2) S \rightarrow T \text{ list}$$

$$T \rightarrow \text{int}$$

$$T \rightarrow \text{float}$$

$$T \rightarrow \text{char}$$

$$T \rightarrow \text{double}$$

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

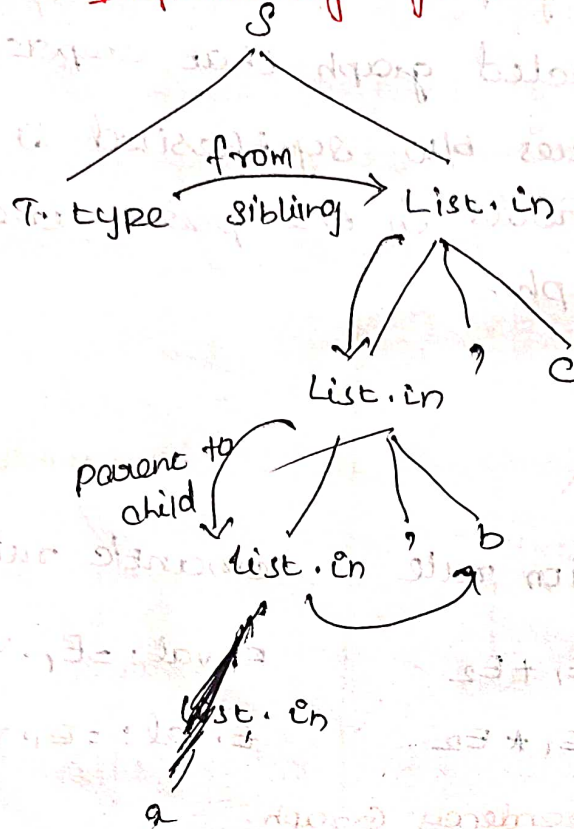
$$\text{list} \rightarrow \text{id}$$

Sol:

production rule	semantic rule
$S \rightarrow T \text{ list}$	$\text{List.in} := T.\text{type}$
$T \rightarrow \text{int}$	$T.\text{type} := \text{integer}$
$T \rightarrow \text{float}$	$T.\text{type} := \text{float}$
$T \rightarrow \text{char}$	$T.\text{type} := \text{char}$
$T \rightarrow \text{double}$	$T.\text{type} := \text{double}$
$\text{list} \rightarrow \text{list}, \text{id}$	$\text{List}_1.\text{in} := \text{List}_2.\text{in}$
	$\text{Enter-type}(\text{id}, \text{entry}, \text{list.in})$
$\text{list} \rightarrow \text{id}$	$\text{Enter-type}(\text{id}, \text{entry}, \text{list.in})$

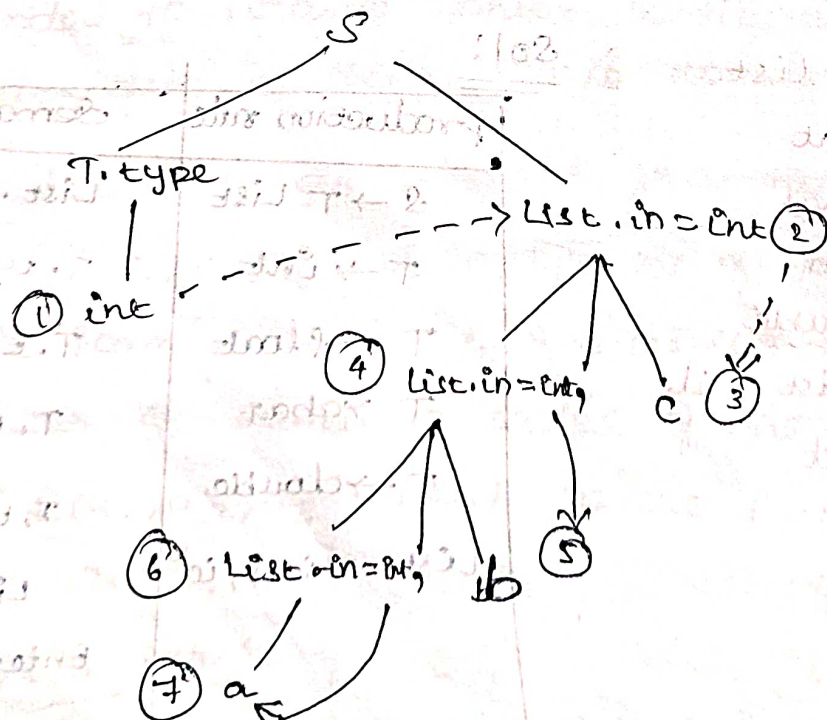


# Dependency graph:



## 4) Evaluation order:

The topological sort of the dependency graph decides the evaluation order in a parse tree, synt. semantic rules, in the syntax-directed def. are used.



Evaluation order.

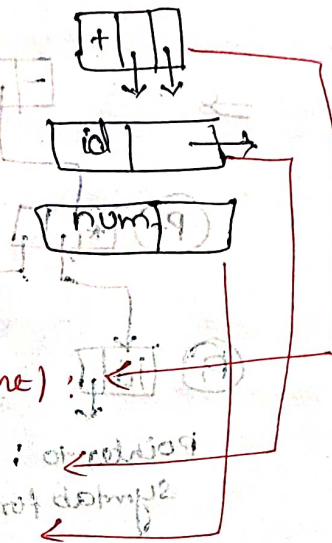
2)

## Construction of Syntax Tree Expression:

99

Grammar:

$E \rightarrow E + T$   
 $E \rightarrow E - T$   
 $E \rightarrow E * T$   
 $E \rightarrow T$   
 $T \rightarrow id$   
 $T \rightarrow num$



- ① mkenode (op, left, right)
- ② mkleaf (id, entry)
- ③ mkleaf (num, val)

### ① Construct the syntax tree for expression $x * y - 5 + z$

Sol:

Step 1: Convert the expression from infix to postfix

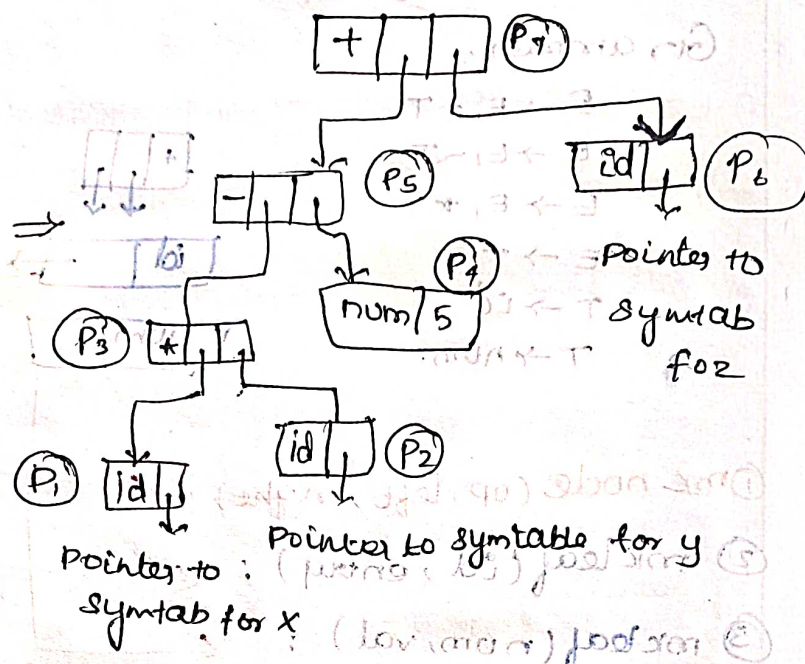
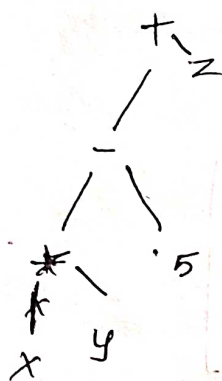
Step 2: Make use of the fun. mkenode(), mkleaf(id, entry)

Step 3: The seq. of fun calls is given.

postfix expression  $xy * 5 - z +$

Symbol	Operation
x	$P_1 = \text{mkleaf}(\text{id}, \text{ptr to entry } x)$
y	$P_2 = \text{mkleaf}(\text{id}, \text{ptr to entry } y)$
*	$P_3 = \text{mkenode}(*, P_1, P_2)$
5	$P_4 = \text{mkleaf}(\text{num}, 5)$
-	$P_5 = \text{mkenode}(-, P_3, P_4)$
z	$P_6 = \text{mkleaf}(\text{id}, \text{ptr to entry } z)$
+	$P_7 = \text{mkenode}(+, P_5, P_6)$





Syntax directed definition for grammar,

### production Rule

Semantic operation

$$E \rightarrow E_1 + T$$
$$E.npt_1 := \text{middle}('+', E_1.npt_1,$$
$$E \rightarrow E, \neg J$$

$E, \text{np}^2 = \text{mic mode } (1^-, E, \text{np}^2)$

$$E \rightarrow E_1 * T$$
$$E.npb[i] = mcode(i)('t')$$

$E_1, n p b_1, T, n p b_1)$

$$E_{\text{TOT}} = T \cdot \ln \left( \frac{1}{\Omega} \right) = T \cdot \ln \left( \frac{1}{\Omega_{\text{TOT}}} \right) = T \cdot \ln \left( \frac{1}{\Omega_{\text{A}} \cdot \Omega_{\text{B}}} \right) = T \cdot \ln \left( \frac{1}{\Omega_{\text{A}}} \right) + T \cdot \ln \left( \frac{1}{\Omega_{\text{B}}} \right) = E_{\text{A}} + E_{\text{B}}$$

$\tau \rightarrow \text{id}$  in  $\text{box}$   $\vdash \text{E} : \text{npt} : = \text{mleaf} (\text{id}, \text{id}, \text{pr}_2, \text{empty})$

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

(3)  $\text{H}_2\text{O}$  = 49

(179)  $\text{H}_2\text{O} + \text{CO}_2 = \text{H}_2\text{CO}_3$

(Spend of 10, 15, 20, 25, 30, 35, 40, 45, 50, 55, 60, 65, 70, 75, 80, 85, 90, 95, 100)

(39) 1919 - 1920 - 1921





## S-attributes on the parser stack!

- ① A translator for S-attributes def. is implemented using LR parser generator.
- ② A bottom up method is used to parse the i/p str.
- ③ A parser stack is used to hold the values of S-attribute.

production:  $x \rightarrow ABC$

state	value
A	A.a
B	B.b
C	C.c

TOP  $\rightarrow$

state	value
X	X.x

TOP  $\rightarrow$

Bottom up Evaluation of

Before reduction  
The TOP symbol on the stack is pointed by pointer

Prod. rule	Sem. action
$x \rightarrow ABC$	$x.x = f(A.a, B.b, C.c)$

- ④ After reduction, TOP is decremented by 2 the state containing X is placed at the top of state [TOP] & value of S-attribute X.x is put in value [TOP]

- ⑤ If the symbol has no attribute then the corresponding entry in the value array will be kept undefined.

Eg:

Construct a syntax directed translation scheme that translates arithmetic exp. from infix to postfix notation. Using ~~arithmetic~~ semantic rules for each of the grammar symbols & semantic rules, evaluate input:  $3 * 4 + 5 * 2$ .

$$E \rightarrow E + T$$

$$E \rightarrow E - T$$

$$E \rightarrow T$$

$$T \rightarrow id$$

Sol:

Semantic rules for evaluation of expr. for given grammar

$$L \rightarrow E \quad \{ \text{print}(E.val) \}$$

$$E \rightarrow E + T \quad \{ E_1.val = E_1.val + T.val \}$$

$$E \rightarrow T \quad \{ E.val = T.val \}$$

$$T \rightarrow T * F \quad \{ T.val = T_1.val * F.val \}$$

$$T \rightarrow F \quad \{ T.val = F.val \}$$

$$F \rightarrow id \quad \{ F.val = id.lexval \}$$

Symbol stack	value stack	I/P string	Syntax action	Semantic rule
\$	\$	3 * 4 + 5 * 2	shift	
\$ id	\$ 3	* 4 + 5 * 2	Reduce $F \rightarrow id$	$F.val = id.val$
\$ F	\$ 3	* 4 + 5 * 2	Reduce $T \rightarrow F$	$T.val = F.val$
\$ T	\$ 3	* 4 + 5 * 2	shift	
\$ T *	\$ 3 *	+ 5 * 2	shift	
\$ T * id	\$ 3 * 4	+ 5 * 2	Reduce $F \rightarrow id$	$F.val = id.val$



\$ T * F	\$ 13 * 4	+ 5 * 2	reduce $F \rightarrow T * F$	$T.val = T.val$ $F.val$
\$ T	\$ 12	+ 5 * 2	reduce $E \rightarrow T$	$E.val = T.val$
\$ E * F	\$ 12	+ 5 * 2	shift	
\$ E + id	\$ 12 + id	5 * 2	shift	
\$ E + id	\$ 12 + 5	* 2	reduce $F \rightarrow id$	$F.val = id.val$
\$ E + F	\$ 12 + 5	* 2	reduce $T \rightarrow F$	$T.val = F.val$
\$ E + T	\$ 12 + 5	* 2	shift	
\$ E + T *	\$ 12 + 5 *	2	shift	
\$ E + T * id	\$ 12 + 5 * 2		reduce $F \rightarrow id$	$F.val = id.val$
\$ E + T * F	\$ 12 + 5 * 2		reduce $T \rightarrow T * F$	$T.val = T.val + F.val$
\$ E + T	\$ 12 + 10		reduce $E \rightarrow E + T$	$E.val = E.val + T.val$
\$ E	\$ 22		reduce $L \rightarrow E$	Print (E.val)

### L-attributed definition

SDP can be defined as the L-attributed for the production rule  $A \rightarrow x_1 x_2 \dots x_n$  where the embedded attribute  $x_k$  is such that  $1 \leq k \leq n$ . The production  $A \rightarrow x_1 x_2 \dots x_n$  is such that.



① It depends upon the attributes of the symbol

$x_1, x_2, \dots, x_{j-1}$  to left of  $x_j$ .

② It also depends upon the inherited attribute A.  
check whether the given SDD is L-attributed or not.

$$\begin{aligned}
 A \rightarrow PQ & \quad P.in := p(A.in) \\
 & \quad Q.in := q(P.sy) \\
 & \quad A.sy := f(Q.sy) \\
 A \rightarrow xy & \quad y.in := y(A.in) \\
 & \quad x.in := x(y.sy) \\
 & \quad A.sy := f(x.sy)
 \end{aligned}$$

Sol:

The attr.  $in$  and  $sy$  represents the inherited and S-attr respectively. The given SDD is not L-attr.

prod. rule	sem. action	obs of attri
$A \rightarrow PQ$	$P.in = p(A.in)$ $Q.in := q(P.sy)$ $A.sy := f(Q.sy)$	<div style="border-left: 1px solid black; padding-left: 10px;"> <math>\hookrightarrow</math>  L-attribute </div>
$A \rightarrow xy$	$y.in := y(A.in)$ $x.in := x(y.sy)$ $A.sy := f(x.sy)$	<div style="border-left: 1px solid black; padding-left: 10px;"> L-attribute  not-L-attri  L-attribute. </div>

Definition  $x.in := x(y.sy)$ .

sem. action suggests the value  $x.in$  depends upon value  $y.sy$ . violates the 1st rule of L-attr defn.

Thus  $x.in := x(y.sy)$  is not L-attribute.



# S-attributed grammar, L-attributed grammar.

① Grammar having no inherited attr, but only synthesized attr  
 attr to be evaluated in one left-to-right traverse of abstract syntax tree.

Grammar in which inherited attr can be evaluated.

Incorporated in both

② Both Bottomup & Top-down parsing  
 Incorporated in using top-down parsing

③ The YACC family can be broadly considered as S-attributed grammar  
 special type of compiler narrow compiler are based on some of the L-attr. grammar.

④ S-attr grammar can be L-attr grammar  
 L-attr grammar can not be S-attr

⑤ Eq:  $(p \cdot A)q = (p \cdot A)q$   
 $(p \cdot A)q = (p \cdot A)q$   
 $(p \cdot A)q = (p \cdot A)q$

L-attr grammar can not be S-attr  
 Eq:  $q \leftarrow A$   
 Eq:

## Translation Scheme:

During the process of parsing the evaluation of attribute takes place by consulting the semantic action enclosed in  $\{ \}$  at the sight of the grammar symbol. This process of eval. of code fragment semantic actions from the syntax directed def is called Syntax directed translation. Thus the SDD can be done by Syntax directed Translation scheme.

$\Rightarrow$  A translation scheme generates the T/p (107)  
 by executing the sem. actions in order manner  
 $\Rightarrow$  This processing is using depth first traversal

Give the translation scheme that converts infix to postfix form for the following grammar. Also generates the annotated parse tree T/p string  $2+6+1$ .

$E \rightarrow E+T \quad \{ \text{print}(' + ') \}$   
 $E \rightarrow T$

$T \rightarrow 0 \quad \{ \text{print}(' 0 ') \}$

$T \rightarrow 1 \quad \{ \dots \}$

$T \rightarrow 2$

$T \rightarrow 3$

$T \rightarrow 4$

$T \rightarrow 5$

$T \rightarrow 6$

$T \rightarrow 7$

$T \rightarrow 8$

$T \rightarrow 9 \quad \{ \text{print}(' 9 ') \}$

Sol: Let us convert this grammar into non left recursive form.

$E \rightarrow TP$

$T \rightarrow 0|1|2|3|4|5|6|7|8|9$

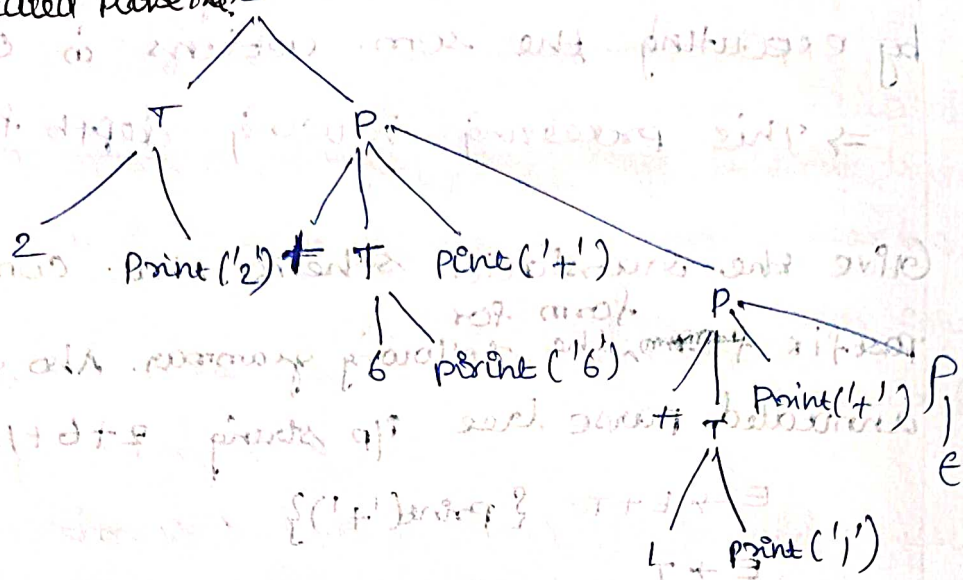
$P \rightarrow +TP | \epsilon$

Translation scheme for this grammar!

productions	semantic actions.
$E \rightarrow TP$ $T \rightarrow 0$ $T \rightarrow 1$ $T \rightarrow 2$ $T \rightarrow 3$ $\vdots$ $T \rightarrow 9$ $P \rightarrow +TP   \epsilon$	$\{ \text{print}(' 0 ') \}$ $\{ \text{print}(' 1 ') \}$ $\{ \text{print}(' 2 ') \}$ $\{ \text{print}(' 3 ') \}$ $\vdots$ $\{ \text{print}(' 9 ') \}$ $\{ \text{print}(' + ') \} P   \epsilon$



# Annotated parse tree E

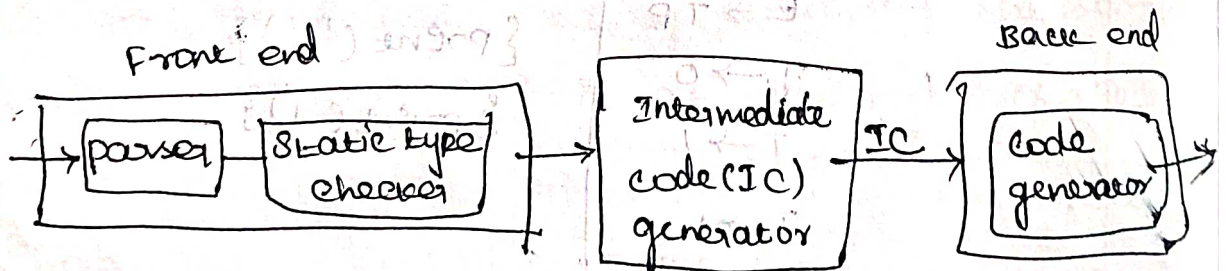


## Intermediate languages:

### Introduction:

Compiler converts the source pgm into machine pgm. This activity can be done directly, but it is not always possible to generate such machine code directly in one pass.

Then typically compilers generate an easy to represent form of a source language which is called Intermediate language. The generation of an intermediate language leads to efficient code generation.



## Properties:

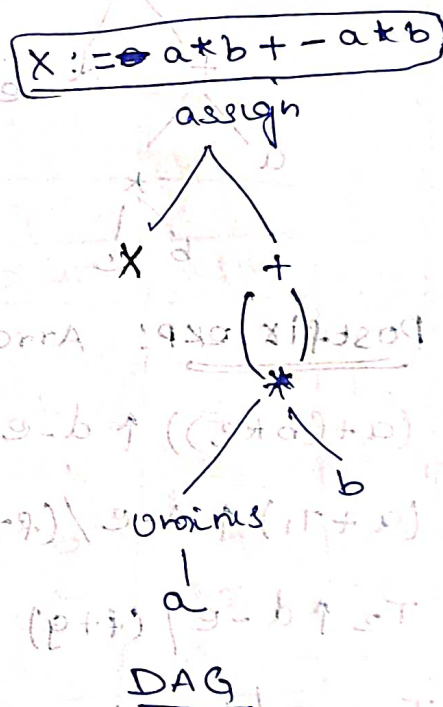
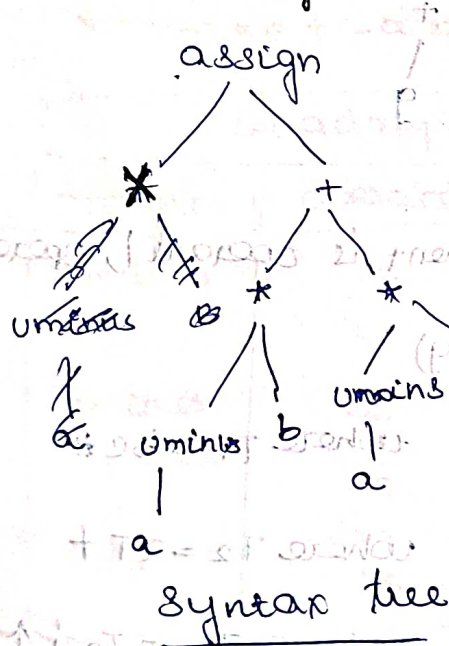
(109)

- ① The IC is an easy form of source language which can be generated eff. by the compiler.
- ② The gener. of IC should lead to eff. code gener.
- ③ The Intermediate language should act as "effective mediator" b/w front & back end.
- ④ I.L should be flexible enough so that optimized code can be generated.

## Forms of IC:

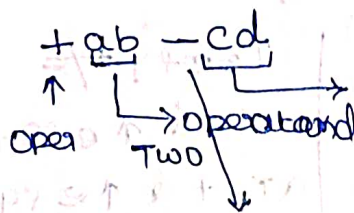
Three types:

① Abstract syntax tree:



② Polish notation: The Polish notation is also called as Prefix notation in which the operator occurs first then operands are arranged.

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





10

There is a reverse polish notation which can be using postfix or posin representation.

$$X := -a * b + -a * b$$

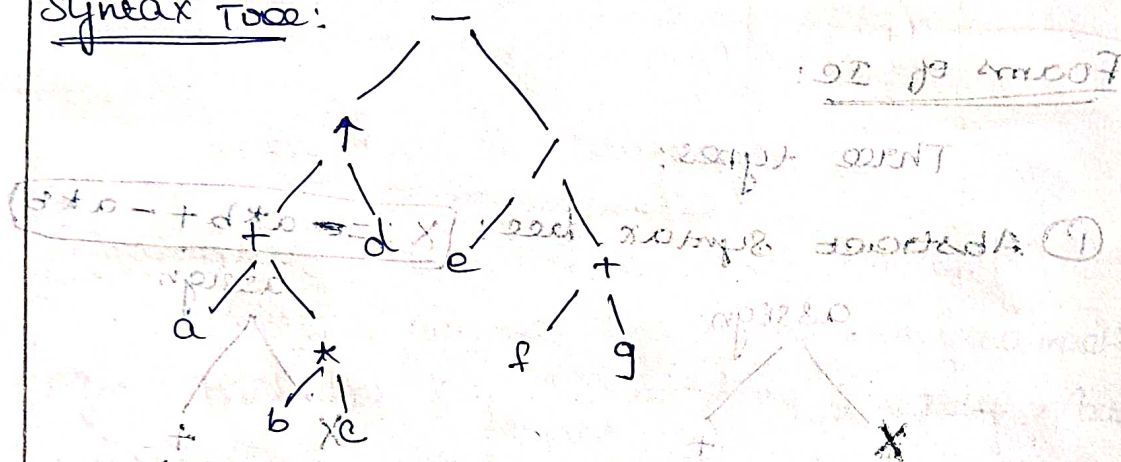
Postfix:

$$X a - b * a - b * + ; =$$

① Construct syntax tree & postfix notation for the following exp.

$$(a + (b * c)) \uparrow d - e / (f + g)$$

Syntax Tree:



Postfix exp: Arrangement is operand 1, operand 2, operator

$$(a + (b * c)) \uparrow d - e / (f + g)$$

$$(a + T_1) \uparrow d - e / (f + g)$$

$$\text{where } T_1 = b * c$$

$$T_2 \uparrow d - e / (f + g)$$

$$\text{where } T_2 = a * T_1$$

$$T_3 = e / (f + g)$$

$$\text{where } T_3 = T_2 d \uparrow$$

$$T_3 = e / T_4$$

$$\text{where } T_4 = f + g$$

$$T_3 = T_5$$

$$\text{where } T_5 = e * T_4 / (f + g)$$

$$T_6$$

$$\text{where } T_6 = T_3 T_5 -$$

backward substituting the values:

$$T_6$$

$$T_3 T_5 -$$

$$T_3 e T_4 / -$$

$$T_3 e f g + / -$$

$$T_3 e f g + / -$$

$$T_2 d \uparrow e f g + / -$$

$$a T_1 + d \uparrow e f g + / -$$

$$a b c + d \uparrow e f g + / -$$

### ③ Three address code: (11)

Abstract form of IC that can be implemented as a record with the address fields.

Representations used three address code such as

- (11) ① Quadruples ② Triples ③ Indirect triples.

#### Quadruple representation:

(11) Four fields - op, arg1, arg2, result.

(11) op - field is represent internal code for operators.

(11) arg1, arg2 - Two operands used

(11) result - result to store of an expr.

I/p statement:

$$x := -a * b + -a * b$$

Quadruples				
	Operator	operand1	operand2	Result
(0)	Uminus	a		t1
(1)	*	t1	b	t2
(2)	Uminus	a		t3
(3)	*	t3	b	t4
(4)	+	t2	t4	t5
(5)	:=	t5		x

$$t1 := \text{Uminus } a$$

$$t2 := t1 * b$$

$$t3 := -a$$

$$t4 := t3 * b$$

$$t5 := t2 + t4$$

$$x := t5$$

Triples: Temporary variable is avoided by referencing

the pointers in the symbol table.

$$x := -a * b + -a * b.$$

Number	operator	operand1	operand2
(0)	Uminus	a	
(1)	*	(0)	b
(2)	Uminus	a	
(3)	*	(2)	b
(4)	+	(1)	(3)
(5)	:=	x	(4)



Indirect triples - listing pointers are used instead of using statements.

Num	operator	operand1	operand2	Statement	
(0)	ominus	a		(0)	(11)
(1)	*	(11)	b	(1)	(12)
(2)	ominus	a		(2)	(13)
(3)	*	(13)	b	(3)	(14)
(4)	+	(12)	(4)	(4)	(15)
(5)	:=	x	(15)	(5)	(16)

Quadruple	Triple	Indirect triple
-----------	--------	-----------------

① Format of quadruple (OP, operand1, operand2, result)  
is (OP, op<sub>1</sub>, operand, result)

but it makes use of two tables.

② Temporary var. is used from access of symbol table.

pointer is used. To access from symbol table.

Computations is made separately & stored saves some <sup>amt of</sup> space compared with quadruple triples.

③ This representation is beneficial for code optimization.

The use of pointers allows to access the symbol table entries quickly.

During the code optimization the table of inserting, deleting the inst. is involved. Easy for quadruples & difficult for triples.

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

Three address Code:

$$\begin{aligned} t_1 &:= a * b & t_5 &:= a + b \\ t_2 &:= 0 \text{ minus } t_1 & t_6 &:= t_5 + t_2 \\ t_3 &:= c + d & t_7 &:= t_4 - t_6 \\ t_4 &:= t_2 + t_3 \end{aligned}$$

Quadruple:

Location	operator	operand1	operand2	Result
(1)	*	a	b	t1
(2)	0 minus	t1		t2
(3)	+	c	d	t3
(4)	+	t2	t3	t4
(5)	+	a	b	t5
(6)	+	t5	t3	t6
(7)	-	t4	t6	t7

Triple:

Location	operator	operand1	operand2
(1)	*	a	b
(2)	0 minus	(1)	
(3)	+	c	d
(4)	+	(2)	(3)
(5)	+	a	b
(6)	+	(5)	(3)
(7)	-	(4)	(6)

Indirect triple:

Loc.	operator	operand1	operand2	Stmt
(1)	*	a	b	(11)
(2)	0 minus	(11)		(12)
(3)	+	c	d	(13)
(4)	+	(12)	(13)	(14)
(5)	+	a	b	(15)
(6)	+	(15)	(13)	(16)
(7)	-	(14)	(16)	(17)



(14)

# Types of Three address code: (i + o) - 123

Language Construct	IC Form	Meaning
Assign. Stmt	$x := y \text{ op } z$	Binary operation is performed using operator 'op'
Assign. Stmt	$x := \text{op } y$	Unary ops. is performed. op is unary operator.
Copy Stmt	$x := y$	value of y is assigned to x.
Unconditional Jump	goto L	The control flow goes to stmt labeled by L.
Conditional Jump	if x rel op y goto L	relational operators such as $<, =, >$ . if x rel op is true then it executes goto L stmt.
procedure calls	Param $x_1$ Param $x_2$ ... Param $x_n$ call P, n return y	$x_1, x_2, \dots, x_n$ are used as parameters to the procedure P. return value of y.
Array Stmt	$x := y[i]$ $x[i] := y$	i-th index of array y is assigned to x. value of y is assigned to index i of array x.
Address & pointer assignment	$x := \&y$ $x := *y$ $*x := y$	address of y stores to x pointer of x stores x

Declaration:

In declarative stmt the data items along with their data types are declared.

$S \rightarrow D$	$\{ \text{offset} := 0 \}$
$D \rightarrow \text{id} : T$	$\{ \text{enter\_tab}(\text{id.name}, T.\text{type}, \text{offset});$ $\{ \text{offset} := \text{offset} + T.\text{width} \} \}$
$T \rightarrow \text{Integer}$	$\{ T.\text{type} := \text{integer};$ $T.\text{width} := 4 \}$
$T \rightarrow \text{real}$	$\{ T.\text{type} := \text{real};$ $T.\text{width} := 8 \}$
$T \rightarrow \text{array}[\text{num}]$ of $T_1$	$\{ T.\text{type} := \text{array}(\text{num.val}, T_1.\text{type});$ $T.\text{width} := \text{num.val} \times T_1.\text{width} \}$
$T \rightarrow *T_1$	$\{ T.\text{type} := \text{pointer}(T_1.\text{type});$ $T.\text{width} := 4 \}$

$D \rightarrow \text{id} : T$  is a declarative stmt for id declaration.

Assignment Stmt:

Mainly deals with expressions. The exp. can be type of integer, real, array & record.

- ① Obtain the translation scheme for obtaining the tree address code for the grammar

$S \rightarrow \text{id} := E$

$E \rightarrow E_1 + E_2$

$E \rightarrow E_1 * E_2$

$E \rightarrow -E_1$

$E \rightarrow (E_1)$

$E \rightarrow \text{id}$



(116)

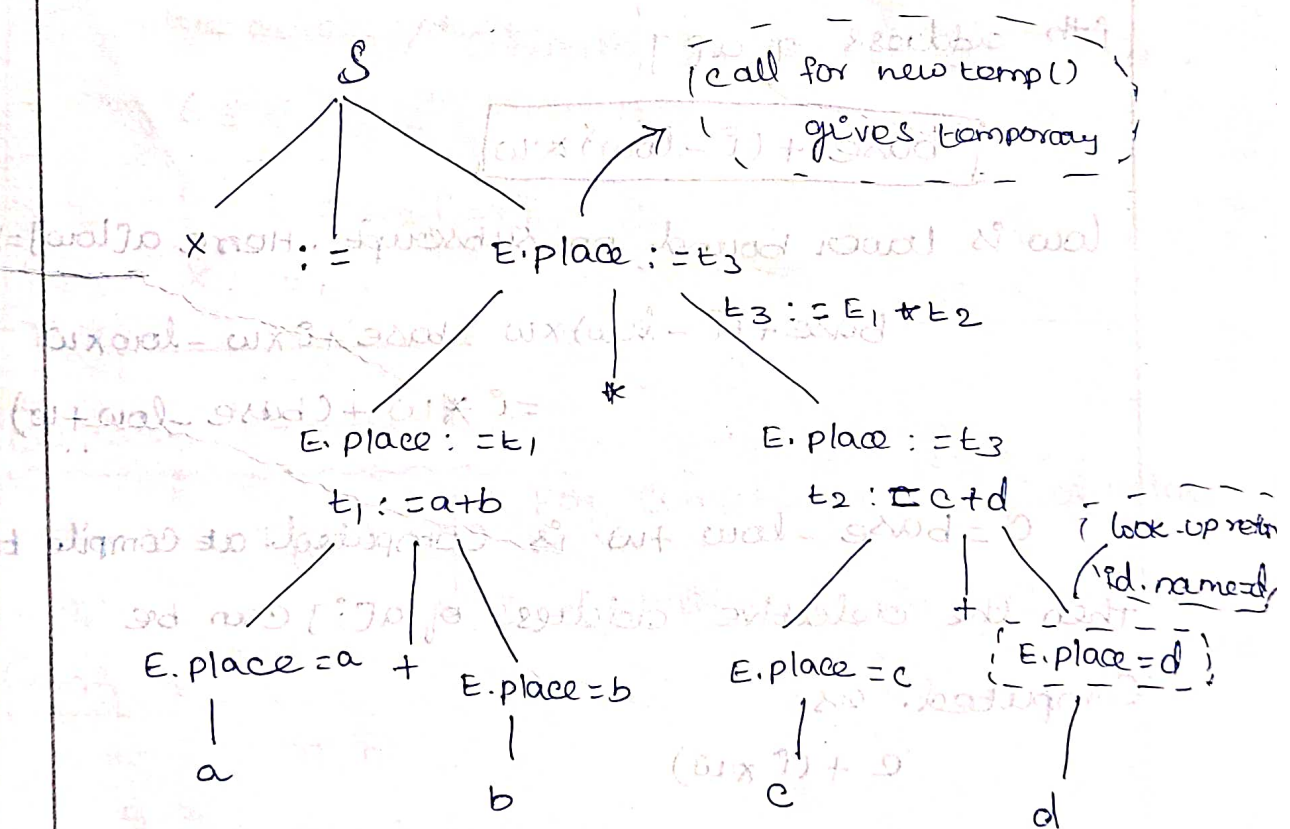
Sol:

prod. Rule	Semantic Actions.
$S \rightarrow id := E$	<pre> { id.entry := look-up(id.name);   if id.entry ≠ nil then     append(id.entry := 'E.place);   else error; /* id not decl */ }</pre>
$E \rightarrow E_1 + E_2$	<pre> {   E.place := newtemp();   append(E.place := 'E<sub>1</sub>.place' +         'E<sub>2</sub>.place') }</pre>
$E \rightarrow E_1 * E_2$	<pre> {   E.place := newtemp();   append(E.place := 'E<sub>1</sub>.place'         '* E<sub>2</sub>.place') }</pre>
$E \rightarrow -E_1$	<pre> {   E.place := newtemp();   append(E.place := 'uminus'         'E<sub>1</sub>.place') }</pre>
$E \rightarrow (E_1)$	<pre> {   E.place := E<sub>1</sub>.place }</pre>
$E \rightarrow id$	<pre> {   id.entry := look(id.name);   if id.entry ≠ nil then     append(id.entry := 'E.place)   else error; }</pre>

Consider the assignment stmt:  $X := (a+b) * (c+d)$

produ. rule	sem. Action	O/p
$E \rightarrow id$	$E.place := a$	
$E \rightarrow id$	$E.place := b$	
$E \rightarrow E_1 + E_2$	$E.place := t_1$	$t_1 := a + b$
$E \rightarrow id$	$E.place := c$	
$E \rightarrow id$	$E.place := d$	
$E \rightarrow E_1 + E_2$	$E.place := t_2$	$t_2 := c + d$
$E \rightarrow E_1 * E_2$	$E.place := t_3$	$t_3 := (a + b) * (c + d)$
$S \rightarrow id := E$		$X := t_3$

Annotated parse tree: for generation of 3 address code





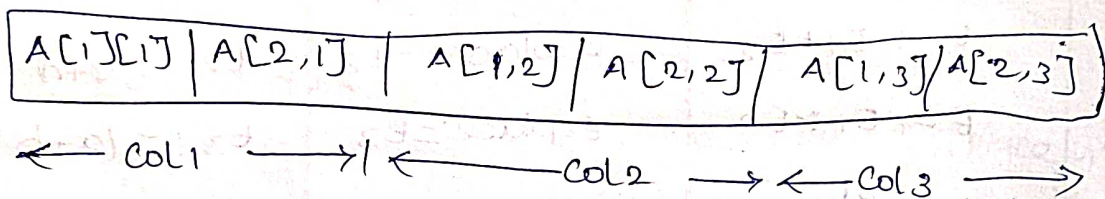
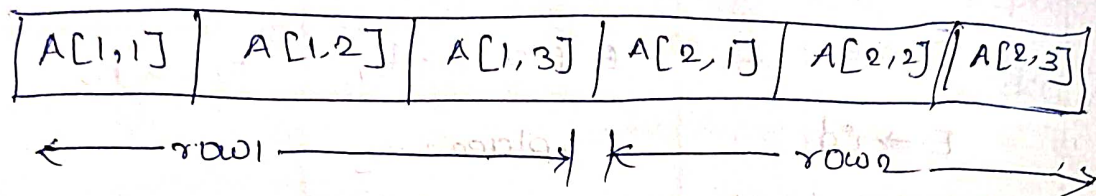
(13)

Array: Contiguous storage of elements.

2 Representations:

① Row major repres.

② Column major repres.



Let base is the address at  $a[]$  and  $w$  is the width of element (req. many units) then to compute  $i^{\text{th}}$  address of  $a[]$

$$\boxed{\text{base} + (i - \text{low}) \times w}$$

low is lower bound on subscript. Here  $a[\text{low}] = \text{base}$

$$\begin{aligned} \text{base} + (i - \text{low}) \times w &= \text{base} + i \times w - \text{low} \times w \\ &= i \times w + (\text{base} - \text{low} \times w) \end{aligned}$$

$C = \text{base} - \text{low} \times w$  is computed at compile time.

Then the relative address of  $a[i]$  can be computed as

$$C + (i \times w)$$

Generate the three address code for expression (119)

$x := A[i, j]$  for an array  $10 \times 20$ . Assume  $low_1 = 1$  and  $low_2 = 1$

Sol: Give that

$$low_1 = 1 \text{ and } low_2 = 1$$

$$n_1 = 10, n_2 = 20.$$

$$A[i, j] = ((i \times n_2) + j) \times w + (\text{base} - ((low_1 \times n_2) + low_2) \times w)$$

$$A[i, j] = ((i \times 20) + j) \times 4 + (\text{base} - ((1 \times 20) + 1) \times 4)$$

$$A[i, j] = 4 \times (20i + j) + (\text{base} - 84)$$

Three address code for expr.

$$t_1 := i \times 20$$

$$t_1 := t_1 + j$$

$$t_2 := c \quad /* \text{Computation of } c = \text{base} - 84 */$$

$$t_3 := 4 * t_1$$

$$t_4 := t_2 [t_3]$$

$$x := t_4$$

Boolean Expression:

2 Types  $\rightarrow$  For computing the logical values  
 $\rightarrow$  In conditional exp. using if-then-else or while-do.

Grammar:

$$E \rightarrow E \text{ OR } E$$

$$E \rightarrow E \text{ AND } E$$

$$E \rightarrow \text{NOT } E$$

$$E \rightarrow (E)$$

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

$$E \rightarrow \text{TRUE}$$

$$E \rightarrow \text{FALSE}$$

(\*) Relop denoted by  $\leq, \geq, \neq, <, >$

(\*) OR, AND - left associate.

(\*) Highest precedence: NOT  
AND  
OR



## Numerical Representation:

Translation scheme for generation of 3 address code for Boolean Expr.

Production Rule	Semantic Rule
$E \rightarrow E_1 \text{ OR } E_2$	$\{$ $E.\text{place} := \text{newtemp}()$ $\text{append}(E.\text{place} := 'E_1.\text{place}' \text{ OR } E_2.\text{place})$ $\}$
$E \rightarrow E_1 \text{ AND } E_2$	$\{$ $E.\text{place} := \text{newtemp}()$ $\text{append}(E.\text{place} := 'E_1.\text{place}$ $' \text{ AND } ' E_2.\text{place})$ $\}$
$E \rightarrow \text{NOT } E_1$	$\{$ $E.\text{place} := \text{newtemp}()$ $\text{append}(E.\text{place} := 'NOT' E_1.\text{place})$ $\}$
$E \rightarrow (E_1)$	$\{$ $E.\text{place} := E_1.\text{place}$ $\}$
$E \rightarrow \text{id}_1 \text{ relop id}_2$	$\{$ $E.\text{place} := \text{newtemp}()$ $\text{append}('E_1.\text{place relop op}$ $E_2.\text{place 'goto' next-state + 3};$ $\text{append}(E.\text{place} := '0');$ $\text{append}('goto' \text{ next-state} + 2);$ $\text{append}(E.\text{place} := '1');$ $\}$
$E \rightarrow \text{TRUE}$	$\{$ $E.\text{place} := \text{newtemp}()$ $\text{append}(E.\text{place} := '1')$ $\}$
$E \rightarrow \text{FALSE}$	$\{$ $E.\text{place} := \text{newtemp}()$ $\text{append}(E.\text{place} := '0')$ $\}$

$P > Q$  AND  $r < s$  OR  $u > v$

100: if  $P > Q$  goto 103

574

101:  $t_1 := 0$

102: goto 104

103:  $t_1 := 1$

104: if  $r < s$  goto 107

105:  $t_2 := 0$

106: goto 108

107:  $t_2 := 1$

108: if  $u > v$  goto 111

109:  $t_3 := 0$

110: goto 112

111:  $t_3 := 1$

112:  $t_4 := t_1$  AND  $t_2$

113:  $t_5 := t_4$  OR  $t_3$

(\*) This method of evaluation is called "Short-circuit".

### Flow of Control Stmt:

Stmts are ① If-then-else

② while-do

$S \rightarrow$  if  $E$  then  $S_1$

| if  $E$  then  $S_1$  else  $S_2$

| while  $E$  do  $S_1$

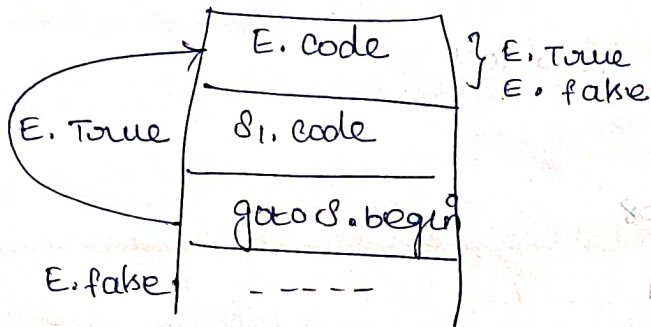
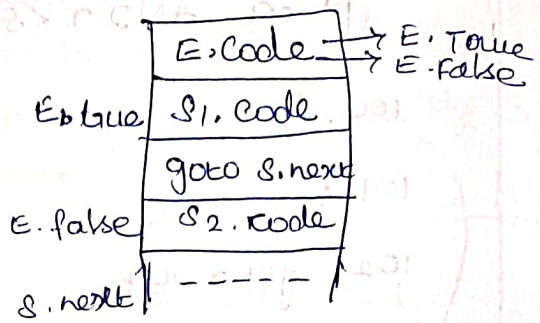
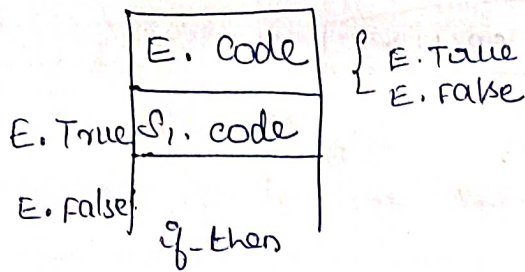
generating three address code -

→ To generate new symbolic label the function new-label() is used.

→ With the Expr,  $E.true$  &  $E.false$  are the labels associated.



122



while do

Short circuit code is:

$S \rightarrow \text{if } E \text{ then } S_1$

$E.\text{true} := \text{new\_label}()$

$E.\text{false} := S.\text{next}$

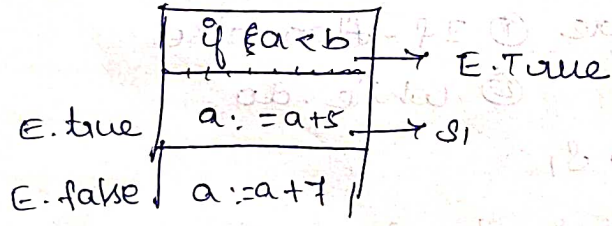
$S_1.\text{next} := S.\text{next}$

$S.\text{code} := E.\text{code} \parallel \text{gen\_code}(E.\text{true})$

$(E.\text{true}) \parallel S_1.\text{code}$

|| - Concatenate string

if  $a < b$  then  $a = a + 5$  else  $a = a + 7$



Three address code

100 : if  $a < b$  goto L1

101 : goto 103

102 : L1 :  $a = a + 5$  /  $E.\text{True}$

103 :  $a := a + 7$

✱

Construct 3 address code for:

(123)

if  $((a < b) \text{ and } (c < d) \text{ or } (a > d))$  then

$z = x + y * z$

else  $z = z + 1$

Sol: Three address code is

100: if  $a < b$  goto 102

101: goto 110

102: if  $c < d$  goto 106

103: goto 104

104: if  $a > d$  goto 106

105: goto 110

106:  $t_1 := x + y$

107:  $t_2 := t_1 * z$

108:  $z := t_2$

109: goto 112

110:  $t_3 := z + 1$

111:  $z := t_3$

112: STOP.



## Concept of Backpatching:

Backpatching is the activity of filling up unspecified info. of labels using appropriate semantic actions in during the code generation process.

will use

- Boolean Expression
- Flow of Control starts

Following fun are used:

→  $mklist()$

→  $merge\_list(P_1, P_2)$  | Concatenates  $P_1, P_2$

→  $backpatch(P, i)$  | insert  $i$  as target label for the stmt pointed by pointer  $P$ .

## Backpatching using Boolean Expression:

grammar for BE:

$E \rightarrow E_1 \text{ OR } M \ E_2$

$E \rightarrow E_1 \text{ AND } M \ E_2$

$E \rightarrow \text{NOT } E_1$

$E \rightarrow (E_1)$

$E \rightarrow id_1 \text{ relop } id_2$

$E \rightarrow \text{TRUE}$

$E \rightarrow \text{FALSE}$

$M \rightarrow E$

$M$  is nonterminal inserted as max. cont. mid



(12)

$E.Tlist$  &  $E.Flist$  are used to generate jumping code.  
(\*) attribute 'state' will be associated with the  $M$  and that is used to record the number (address) of stmt. denoted as  $M.state$ . The 'nextstate' will point to next stmt.

Production Rule	Semantic action.
$E \rightarrow E_1 \text{ OR } M E_2$	$\{$ $backpatch(E_1.Flist, M.state);$ $E.Tlist := merge(E_1.Tlist, E_2.Tlist);$ $E.Flist := E_2.Flist;$ $\}$
$E \rightarrow E_1 \text{ AND } M E_2$	$\{$ $backpatch(E_1.Tlist, M.state);$ $E.Tlist := E_2.Tlist;$ $E.Flist := merge(E_1.Flist, E_2.Flist);$ $\}$
$E \rightarrow \text{NOT } E_1$	$\{$ $E.Tlist := E_1.Flist;$ $E.Flist := E_1.Tlist;$ $\}$
$E \rightarrow (E_1)$	$\{$ $E.Tlist := E_1.Tlist;$ $E.Flist := E_1.Flist;$ $\}$
$E \rightarrow id_1 \text{ relop } id_2$	$\{$ $E.Tlist := mclist(nextstate);$ $E.Flist := mclist(nextstate+1);$ append('if' $id_1$ , place, relop, 'or' $id_2$ , place, 'goto -') append('goto -'); $\}$
$E \rightarrow \text{True}$	$\{$ $E.Tlist := mclist(nextstate);$ append('goto -'); $\}$

$E \rightarrow \text{FALSE}$	$\{$ $E.Flist := \text{mklist}(\text{nextstate});$ $\text{append}('goto -');$ $\}$
$M \rightarrow E$	$\{$ $M.state := \text{nextstate};$ $\}$

(125)

(\*) Using Backpatching, generate an IC for following EXP.

EXP.  $A < B$  OR  $C < D$  AND  $P < Q$

Sol: scan it from left to right:

$A < B$  matches with rule  $E \rightarrow id, \text{rel op } id_2$ .

Three address code generated as

```

100 if A < B goto }
1001 goto -

```

similarly, for remaining part of EXP.

```

102 if C < D goto -
103 goto -
104 if P < Q goto -
105 goto -

```

Hence

```

102 if C < D goto -
103 goto -
104 if P < Q goto -
105 goto -

```

backpatch( $E.Tlist$ , 104)  
 $\rightarrow M.state = \text{nextstate} = \{103\}$

$E.Tlist := \{104\}$

$E.Flist := \{103, 105\}$

```

102 if C < D goto 104
103 goto -
104 if P < Q goto -
105 goto -

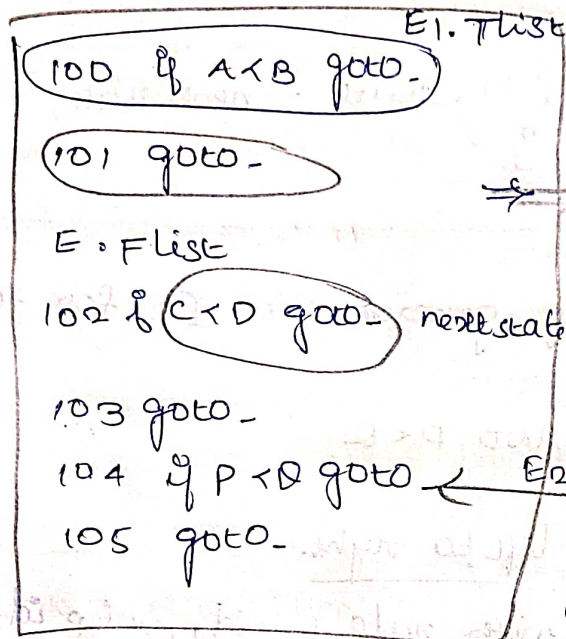
```

XXXX



(26)

Now Consider,

 $A < B$  OR  $C < D$  AND  $P < Q$ 


backpatch(E1.Flist, 102)  
 $E.Tlist = \{100, 101\}$

At this line  $C < D$  AND  $P < Q$  decided to be true.

That means if  $C < D$  is true and if  $P < Q$  is also true then  $C < D$  AND  $P < Q = \text{TRUE}$ .

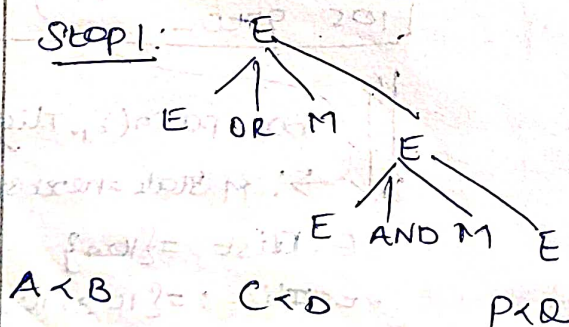
Hence  $E2.Tlist = \{104\}$

Finally 3 address code with backpatching:

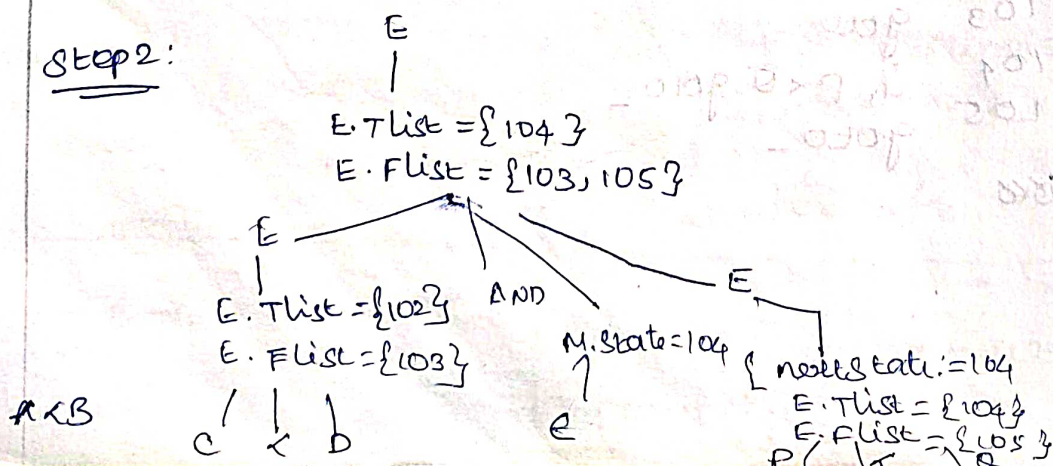
```

100 if A < B goto -
101 goto 102
102 if C < D goto 104
103 goto -
104 if P < Q goto -
105 goto -
  
```

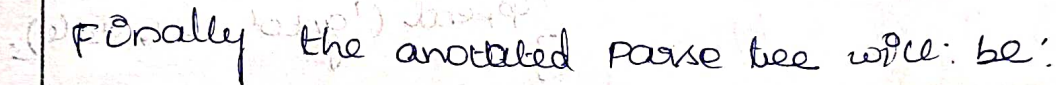
Annotated parse tree:



Step 2:



127.



\_\_\_\_\_ x \_\_\_\_\_

$S \rightarrow \text{if } E \text{ then } S$   
 $S \rightarrow \text{if } E \text{ then } S \text{ else } S$   
 $S \rightarrow \text{while } E \text{ then do } S$   
 $S \rightarrow \text{begin } L \text{ end}$   
 $S \rightarrow A$   
 $L \rightarrow L; S$   
 $L \rightarrow S$

S - stmt  
E - Boolean exp  
L - stmt list  
A  $\Rightarrow$  Assignment  
str

$S \rightarrow q \text{ E then } M_1 S_1$   
 $S \rightarrow \bar{q} \text{ E then } M_1 S_1 \text{ N else } M_2 S_2$   
 $S \rightarrow \text{while } M_1 \text{ E then do } M_2 \text{ S}$   
 $S \rightarrow \text{begin } L \text{ end}$   
 $S \rightarrow A$

$$\begin{aligned} L &\rightarrow L; MS \\ L &\rightarrow S \\ M &\rightarrow e \\ N &\rightarrow e \end{aligned}$$



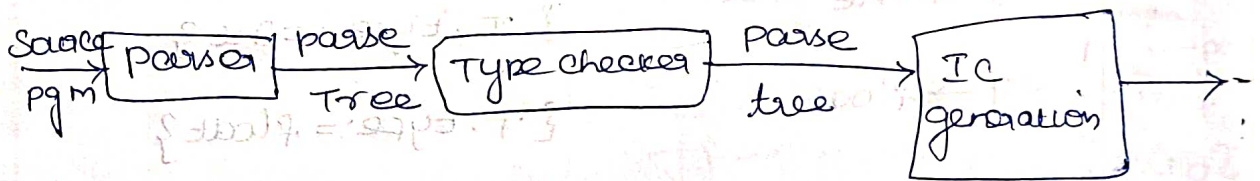
prod. rule	semantic action
$S \rightarrow \epsilon$ E then $M, S_1$	<pre> { backpatch(E.Tlist, M.state);   S.next := merge(E.Flist, S1.next); } </pre>
$S \rightarrow \epsilon$ E then $M, S_1$ N else $M, S_2$	<pre> { backpatch(E.Tlist, M1.state);   backpatch(E.Flist, M2.state);   S.next := merge(S1.next,     merge(N.next, S2.next)); } </pre>
$S \rightarrow \text{while } M, E \text{ then}$ do $M, S$	<pre> { backpatch(S1.next, M1.state);   backpatch(E.Tlist, M2.state);   S.next = E.Flist;   append('goto', M1.state); } </pre>
$S \rightarrow \text{begin } L \text{ end}$	<pre> { S.next := L.next; } </pre>
$S \rightarrow A$	<pre> { S.next := <del>S.next</del> NULL; } </pre>
$L \rightarrow L', M, S$	<pre> { backpatch(L1.next,   M.state);   L.next = S.next; } </pre>
$L \rightarrow S$	<pre> { L.next := S.next; } </pre>
$M \rightarrow E$	<pre> { M.state := next state; } </pre>
$N \rightarrow E$	<pre> { N.next := mlist(next state);   append('goto'); } </pre>

## Type Checking:

### Type System:

Type analysis and type checking is important activity done in sem. analysis phase. need for type checking is -

- ① To detect the errors arising in the exprs due to incompatible operand.
- ② To generate intermediate code for exp. & str.



### Role of type checker

Type Expression: - int, char, float, double, enum are type expr.

typedef int \*INT\_PTR

- ① Array : array (1, T) Ex: int arr[20];
- ② Product :
- ③ Struct : struct stud { char name[10]; float marks;
- ④ pointers ;  
float \*xyz; struct stud student[10];  
Pointers (float);
- ⑤ Function - int sum(int a, int b)



## Syntax directed Translation and Intermediate code generation.

### Design of predictive Translator:

⇒ During the process of parsing the evaluation of attribute takes place by consulting the semantic action enclosed in  $\{ \}$  at the right of the grammar symbol. This process of execution of code fragment semantic actions form the Syntax directed definition is called syntax-directed translation.

⇒ A Translator scheme generates the o/p by executing the semantic actions in an ordered manner.

⇒ This processing is using depth first traversal.

Give the translation scheme that converts infix to postfix form for the following grammar. also generate the annotated parse tree for string  $2+6+1$ .

$E \rightarrow E + T \quad \{ \text{print}(' + ') \}$

$E \rightarrow T$

$T \rightarrow 0 \quad \{ \text{print}(' 0 ') \}$

$T \rightarrow 1 \quad \{ \text{print}(' 1 ') \}$

$T \rightarrow 2 \quad \{ \text{print}(' 2 ') \}$

$T \rightarrow 3 \quad \{ \text{print}(' 3 ') \}$

$T \rightarrow 4$

$T \rightarrow 5$

$T \rightarrow 6$

$T \rightarrow 7$

$T \rightarrow 8$

$T \rightarrow 9 \quad \{ \text{print}(' 9 ') \}$

12/12/20  
12/12/20  
12/12/20  
12/12/20

Solution: Let us convert this grammar into non left recursive

$$E \rightarrow TP$$

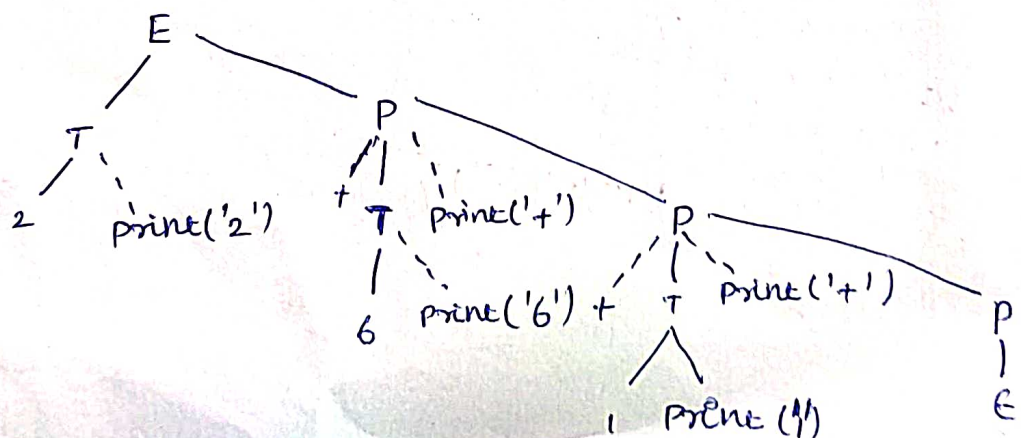
$$T \rightarrow 0|1|2|3|4|5|6|7|8|9$$

$$P \rightarrow +TP|\epsilon$$

Translation scheme:

Productions	Semantic Actions.
$E \rightarrow TP$	
$T \rightarrow 0$	{ print('0') }
$T \rightarrow 1$	{ print('1') }
$T \rightarrow 2$	{ print('2') }
$T \rightarrow 3$	{ print('3') }
$T \rightarrow 4$	{ print('4') }
$T \rightarrow 5$	{ print('5') }
$T \rightarrow 6$	{ print('6') }
$T \rightarrow 7$	{ print('7') }
$T \rightarrow 8$	{ print('8') }
$T \rightarrow 9$	{ print('9') }
$P \rightarrow +TP \epsilon$	{ print('+') } P   $\epsilon$

Annotated parse tree





## Guideline for Designing the translation scheme:

(2)

Production Rule	Semantic Action
$E \rightarrow E_1 * T$	$\{ E.val := E_1.val * T.val \}$

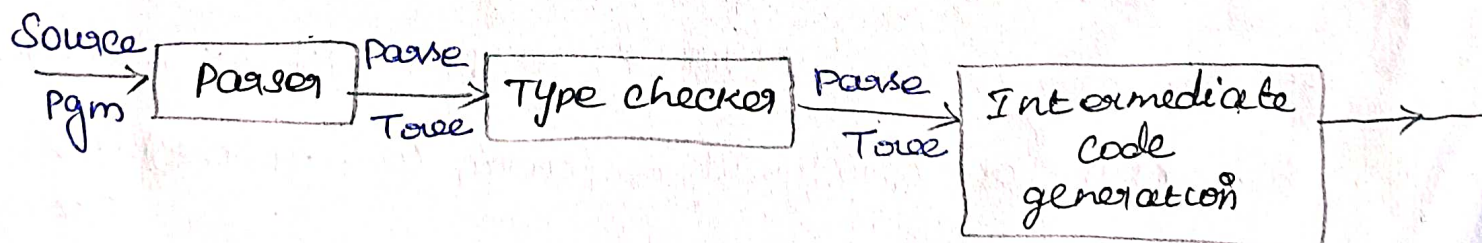
val - synthesized attribute.

- ①  $A \rightarrow x y$  then computation of  $x.in$  must be done <sup>always</sup> before  $y$ .
- ②  $A \rightarrow B_1 B_2$  the semantic action  $\{ B_1.s := B_2.s \}$  is invalid as computation of attribute for  $B_1$  is based on computation of attribute of  $B_2$ .
- ③ Computation of synthesized attribute for non-terminal on the left, first do all the computations of the attributes that appear in non-terminal.

Type Systems: - Important activity done in the semantic Analysis Phase

Need:

- ① To detect the errors arising in the expression due to incompatible operand.
  - ② To generate intermediate code for expressions and stmts.
- Supports two types of data types.
- ① Basic data types - Integer, character, real, Boolean, ...
  - ② Constructed data types - Arrays, record, set, pointers, ...



Role of type checker.



## Type Expression:

① Basic type is called type expression. Hence int, float, double, enum are type expressions.

② While performing type checking two special basic types are needed such as type\_error and void.

③ The type name is also type expression.

typedef int \*INT\_PTR

④ Type Constructors are also type expressions.  
eg: array, product, struct, product.

1) Array - array(I, T)    I - Index set, usually Integer

2) Product -  $T_1 \times T_2$

3) struct - struct stud {  
    char name[10];  
    float marks;  
};  
struct stud student[10];

4) Pointers - float \*xyz;  
    pointer(float)

5) Function - int sum (int a, int b)

LIST  $\rightarrow$  id    LIST.type = T.type

Here T.type = int hence LIST.type = int.

$\therefore$  LIST  $\rightarrow$  L, [num] -    L.type := { array (0, ..., num.val - 1), LIST.type }  
    L.type := { array (0, ..., 9), int } as num.val = 10

the entry for A[10] will be array (0, ..., 9, int) as type expression



## Specification of simple type checker:

(3)

### Type checking of Expression:

$E \rightarrow \text{literal} \quad \{ E.type := \text{char} \}$

$E \rightarrow \text{num} \quad \{ E.type := \text{int} \}$

Here  $E$  can be literal or num, the data types associated with them can be char or int respectively.

$E \rightarrow \text{id} \quad \{ E.type = \text{look-up}(\text{id}.entry) \}$

look-up function reads the symbol table for id entry and thereby it obtains the type of identifiers.

$E \rightarrow E_1 \text{ mod } E_2 \quad \{ E.type := \begin{cases} E_1.type & \text{if } E_1.type = \text{int and } E_2.type = \text{int} \\ \text{type-error} & \text{otherwise} \end{cases} \}$

then int

else

type-error

}

### Type checking of stmts:

Special type void is used.

$S \rightarrow \text{id} := E$

$S \rightarrow \text{if}(E) S_1$

$S \rightarrow \text{while}(E) S_1$

$S \rightarrow S_1 ; S_2$

production rule	Semantic rule.
$S \rightarrow \text{id} := E$	$\{ S.type := \begin{cases} \text{id.type} & \text{if } \text{id.type} = E.type \\ \text{void} & \text{otherwise} \end{cases} \}$
$S \rightarrow \text{if}(E) S_1$	$\{ S.type := \begin{cases} S_1.type & \text{if } E.type = \text{boolean} \\ \text{type-error} & \text{otherwise} \end{cases} \}$

$S \rightarrow \text{while}(E) S_1$	$\{ S.\text{type} := \text{if } E.\text{type} := \text{boolean}(\text{True}) \text{ then } S_1.\text{type}$ $\quad \text{else type-error}$ $\}$
$S \rightarrow S_1; S_2$	$\{ S.\text{type} := \text{if } S_1.\text{type} = \text{void and } S_2.\text{type} = \text{void then}$ $\quad \text{void}$ $\quad \text{else type-error}$ $\}$

### Equivalence of Type Expressions:

The Job of type checker is to find whether two type expressions are equivalent or not.

Two categories:

- $\Rightarrow$  Name equivalence
- $\Rightarrow$  Structural equivalence.

#### ① Structural Equivalence of Type Expressions:

When two expressions are the same basic type or formed by applying the same constructor to structurally equivalent types then those expressions are called structurally equivalent.

$S_1$	$S_2$	Equivalence	Reason.
char	char	$S_1$ is equivalent to $S_2$	similar basic types
pointer (char)	pointer (char)	$S_1$ is equivalent to $S_2$	similar constructor pointer to the char type.



Four bits  
encoding

Basic type	Encoding
boolean	0000
char	0001
Integer	0010
real	0011

Two bits  
encoding

Type constructor	Encoding
pointer	01
array	10
function	11

### Name Equivalence of Type Expressions:

In the name equivalence the type expressions are given the names. The two types expressions are said to be name equivalent if and only if they are identical.

typedef struct Node

{ int x;

}; node;

node \*first, \*second;

struct Node \*last1, \*last2;

first, second are name equivalent.

last1, last2 are name equivalent.

### Type conversions:

⇒ The process of converting one type to another.

eg. stmt  $f + i$  where  $f$  is a float type identifier and  $i$  is an integer type identifier and addition of one has to be done.

⇒ Computer wants to convert one type to another type.

## Two types of conversions:

### ① Explicit Conversion:

```
int xyz, p;
```

```
p = (float)xyz;
```

Identifies xyz is type-casted and this is how explicit conversion from int to float takes place.

### ② Implicit type conversion:

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

```
A[i] = 1;
```

} Takes 4.8 microseconds to execute

Explicit type conversion.

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

```
A[i] = 1.0
```

} 5.4 ms to execute.

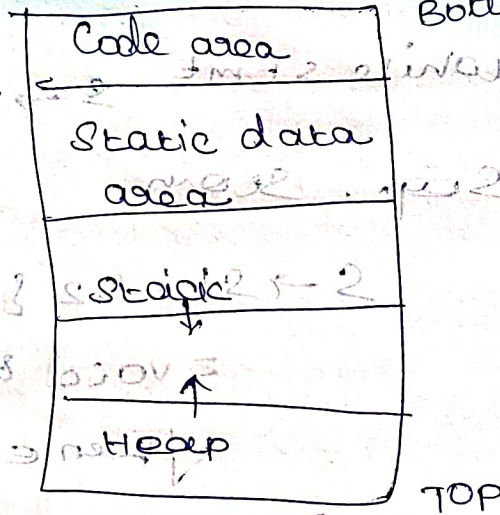


Unit - IVRun Time Environment andCode GenerationStorage Organization:

- N/D - 2022

- A/P - 2022

Bottom.



(X)

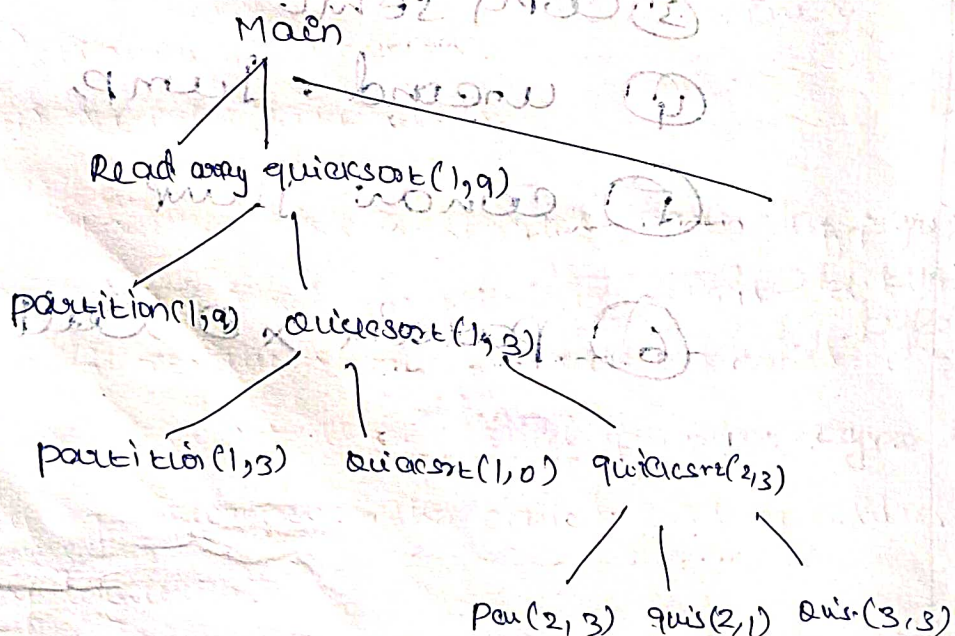
Strategies:

① Static Allocation

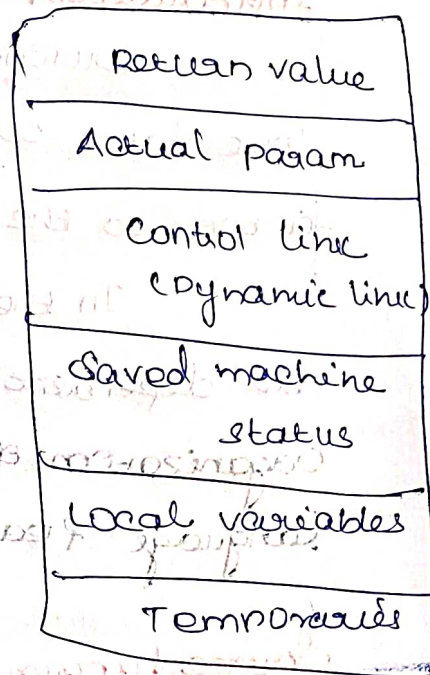
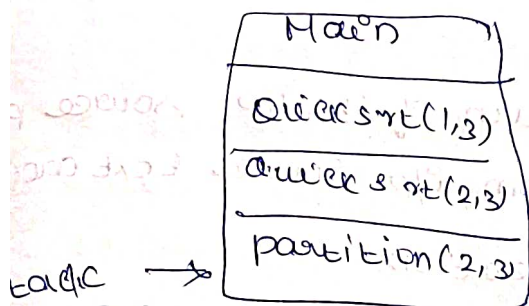
April / May, -2022

② Stack Allocation

③ Heap Allocation.

Storage Allocation space:Activation Tree: - graphical representation,





### Calling sequences:

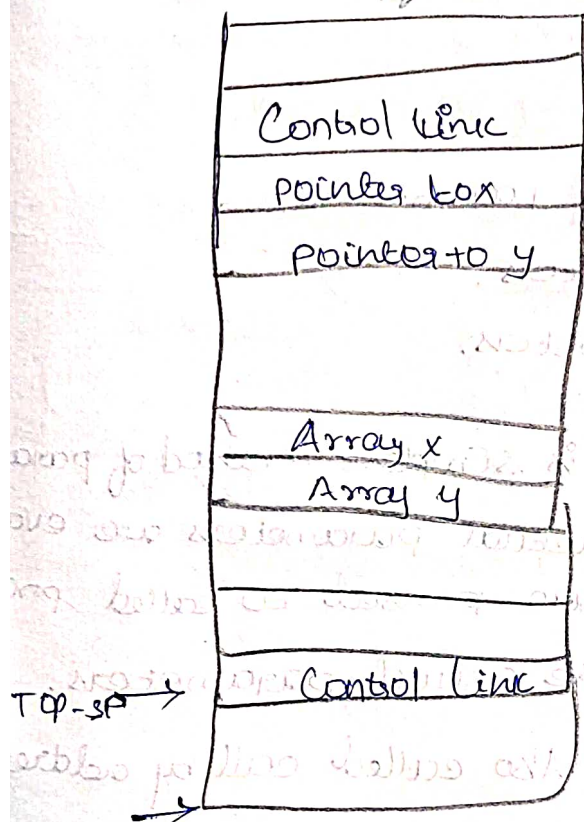
⇒ Series of calls to the procedures

⇒ Allocate Activation Record on the stack

⇒ return sequence is used to restore the state of machine.

⇒ caller is procedure that calls callee

Variable Length data:



Act.  
record for A

Array of A

Act.  
record for B

place  
available length array of B



## Runtime Environments:

### Introduction:

During execution of the i/p source pgm some data objects from i/p source text are imported factors to the code generation.

In the code generation these data objects are dependent upon the org. of mem. And the organization of memory is driven by the source language features.

### Source language Issues:

- ① Does the source language allow recursion?
- ② How the parameters are passed to the procedure?
- ③ Does the procedure defer nonlocal names?
- ④ Does the language support the mem allocation and deallocation dynamically?

### Parameter Passing:

Two types of Parameters:

- i) Formal parameters.
- ii) Actual parameters.

i) Call by value: This is simplest method of parameter passing. The actual parameters are evaluated and their r-values are passed to called procedure.

Example: C, C++ use actual parameters.

ii) Call by reference: Also called call by address or call by location.

The values of actual parameters can be changed.



Eg: Reference parameters in c++, pascal's var param

3) Copy restore: This method is hybrid b/w call by value and call by reference.

Eg: Also called copy-in-copy-out or values result.

4) call by name: This is less popular method of

parameter passing.

⇒ The locals names of called procedure and names of calling procedure are distinct.

Eg: ALGOL uses call by name method.

call by value	call by reference	copy restore	call by name.
1 50	1 50	1 50	1 50
1 50	50 50	50 1	Error.

### Symbol Table: - A/M-2022

(\*) The symbol table is a DS used by compiler to keep track of semantics of variable.

(\*) The symbol table is built in lexical and syntax analysis phase.

I-value and r-value:

The I and r prefixes come from left and right side assignment.

$a := i + 1$   
 $\uparrow \quad \uparrow$   
 I-value r-value



## Symbol Table Entries:

The items to be stored in symbol table are

- 1) variable names
- 2) constants.
- 3) procedure names
- 4) function names
- 5) literal constants & string
- 6) compiler generated temp.
- 7) Labels in source language

Compiler uses following types of information from symbol table:

- 1) Data type
- 2) Name
- 3) Declaring procedures
- 4) Offset in storage
- 5) If structure or record
- 6) No. and type of arguments passed,
- 7) Base address.

How to store names in symbol table?

Two types of name representation.

i) Fixed-length name:

Fixed space for each name is allocated.

Name	Attribute
calculate	
sum	
a	
b	



## 2) Variable Length name:

The amount of space required by string is used to store the names.

Name		Attribute
Starting Index	Length	
0	10	
10	4	
14	2	
16	2	

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
c	a	l	c	u	l	a	e	e	l	u	m	l	a	l	b	l	

## Symbol table Management:

Requirements:

- ① For quick insertion of identifiers and related info.
- ② For quick searching of identifiers.

### i) List Data structure for symbol table:

⇒ Linear Data structure list is a simplest kind of mechanism to implement the symbol table.

⇒ In this method an array is used to store names and associated information.

Name 1	Info 1
Name 2	Info 2
Name 3	Info 3
⋮	⋮
Name n	Info n

⇒ The advantages of list organization is that it takes minimum amount of space.



## 2) Self organizing List:

The symbol table implementation is using linked list. A link field is added to each record.

A pointer "First" is maintained to point to first record of the symbol table.

Name	Info
Name 1	Info 1
Name 2	Info 2
Name 3	Info 3
Name 4	Info 4

## 3) Hash Table:

Hashing is an important technique used to search the records of symbol table. This method is superior to list organization.

$$\text{position} = h(\text{name})$$

Name	Info	Hash line
Sum		
i		
j		
Avg		

Advantages of hashing is quick search & possible and the disadvantage is that hashing is complicated to implement.



# Dynamic Storage Allocation:

142

Two techniques.

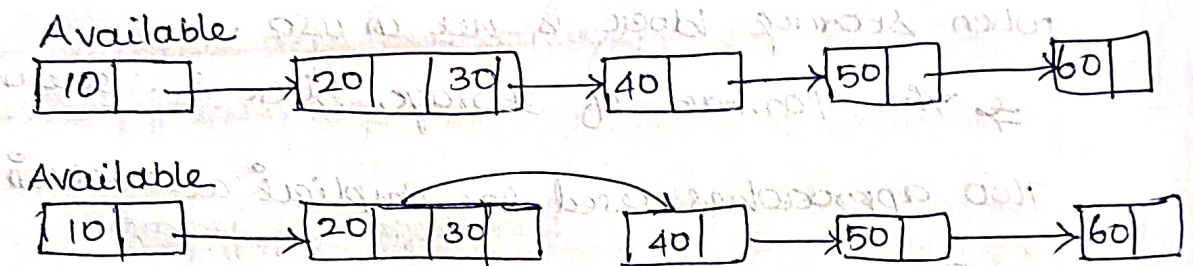
- ① Explicit allocation.
- ② Implicit allocation.

## Explicit Allocation:

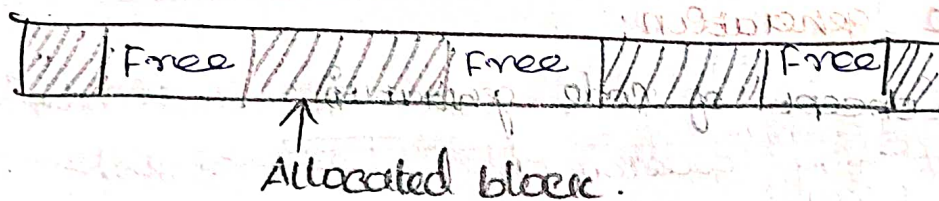
It can be done for fixed size and variable sized blocks.

In this techniques free list is used. Free list is a set of free blocks. This list can observed when we want to allocate memory.

The pointer which points to first block of memory is called "Available".



## Explicit Allocation of variable sized blocks:

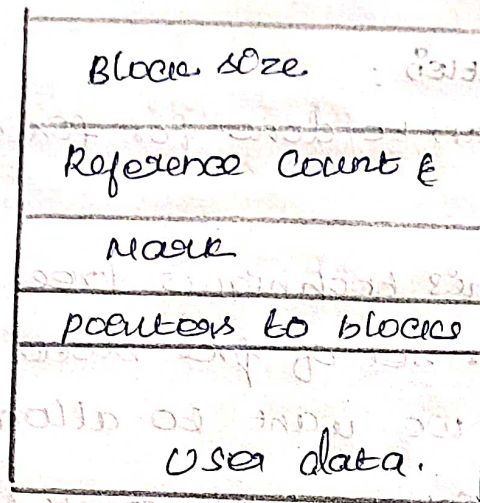


List of 7 blocks gets allocated and second, fourth and sixth block is deallocated then fragmentation occurs. Thus we get variable sized blocks that are available free.



## Implicit Allocation:

The implicit allocation is performed using user pgm and runtime packages.



⇒ The run time package is required to know when storage block is not in use.

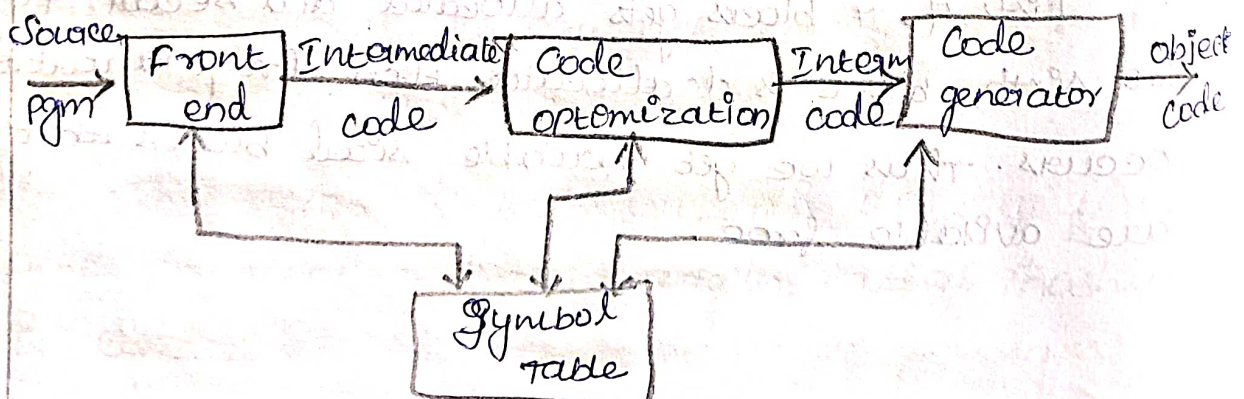
⇒ The format of storage block is shown,

Two approaches used for implicit allocation,

- ① Reference Count
- ② Marking techniques.

## Code generation:

### Concept of code generation:







## Issues in the design of a code generator: - A/M } 2022 N/O }

(145)

### ① I/P to the code generator:

(\*) The code generator phase takes IC as I/P. IC along with the symbol table info & is used to determine the runtime addresses of the data objects.

(\*) The IC may be represented using graphical representations such as syntax trees, or DAGs.

### Target Programs:

The O/P of code generator is target code.

The target code comes in three forms such as;

- 1) absolute machine language,
- 2) relocatable machine language,
- 3) assembly language.

### Memory Management:

(\*) Both the front end and code generator performs the task of mapping the names in the source prog to addresses to the data objects at run time.

(\*) The names in source prog, in these address code refer to the entries in the symbol table.

### Instruction Selection:

The uniformity and completeness of instruction set is an important factor for the code generator. The selection of instruction depends upon the instruction set of target machine.

eg:  $x := a + b$   
MOV a, R0.  
ADD b, R0  
MOV R0, x



(146)

## Register allocation

During register allocation, select appropriate set of variables that will reside in registers.

$$t_1 = a + b$$

MOV a, R0

$$t_2 = t_1 * c$$

ADD b, R0

$$t_1 = t_1 / d$$

MUL c, R0

DIV d, R0

MOV R0, t1

## Target Machine Description:

For designing the good code generator it is necessary to have prior knowledge of target machine and instruction set used for this target machine.

MOV - moves from source to destination

ADD - Add source to des.

SUB - subtracts source from des.

Eg:

MOV R1, M stores the contents of register R1 into memory location M

→ Indexed addressing mode the address offset C from the value of register R0 can be written as, MOV T(R1), M

MOV #5, R0 this instruction we can store the constant 5 into register R0



## (X) Cost of Instruction:

Instruction	Cost	Interpretation
MOV R <sub>0</sub> , R <sub>1</sub>	1	Cost of register mode $+1 = 0+1=1$
MOV R <sub>1</sub> , M	2	use of memory variable $+1 = 1+1=2$
SUB 5(R <sub>0</sub> ), *10(R <sub>1</sub> )	3	use of first constant + use of 2nd constant $+1=3$

## Design of a simple code generator: - A/M - 2018

N/D - 2022

Generating target code from three address statement:

$x := a + b;$

Corresponding target code:

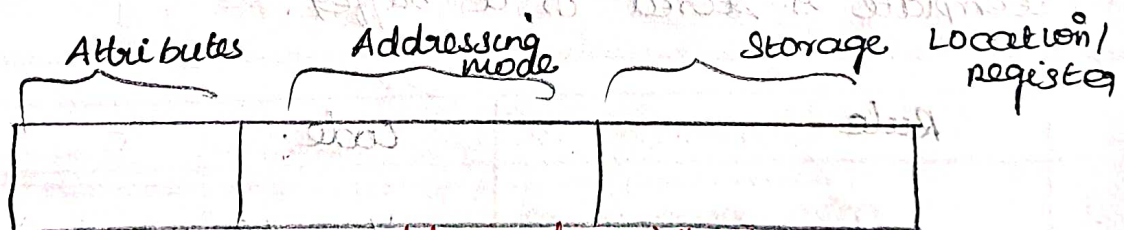
ADD b, R<sub>i</sub>

Here R<sub>i</sub> holds value of a  
Here cost = 2

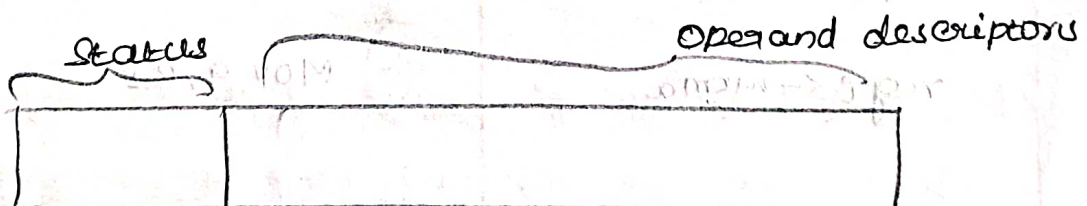
(or)

MOV b, R<sub>0</sub>  
ADD R<sub>1</sub>, R<sub>0</sub>

Here R<sub>0</sub> holds value of a  
Here cost = 2



Address descriptor



Register descriptor



Three address code.	Target code:
$t_1 = a + b$	MOV a, R0
$t_2 = b - c$	ADD b, R0
$t_3 = t_1 / t_2$	MOV b, R1
$t_4 = t_1 * t_2$	SUB c, R1
$t_5 = t_3 - t_4$	MOV R1, R2
$t_6 = t_5 + f$	DIV R0, R2
$x = t_6$	SUB R1, R0
	ADD f, R0
	MOV R0, X

### Optimal Code Generation for Expressions!

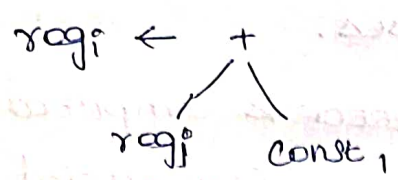
The code generator is a tree rewriting techniques which the instruction selection can be done automatically from a high-level specification of the target machine.

Draw the tree structure for complete expression

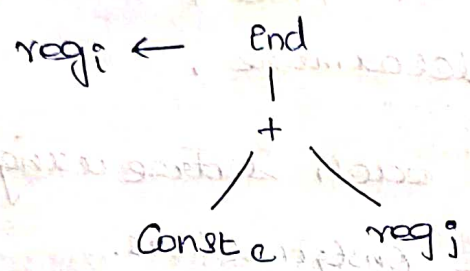
To traverse the tree in bottom up fashion and match the subtrees

The corresponding code for the matched template is stored in a buffer.

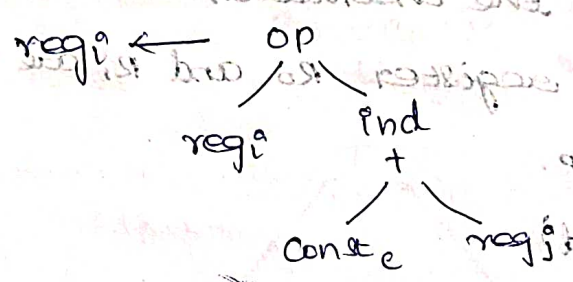
Rule	code.
$reg_0 \leftarrow const$	MOV #c, R0
$reg_0 \leftarrow mem_a$	MOV a, R0
$  \begin{array}{c}  mem \leftarrow \cdot \\  \swarrow \quad \searrow \\  mem_a \quad reg_0  \end{array}  $	MOV R0, a



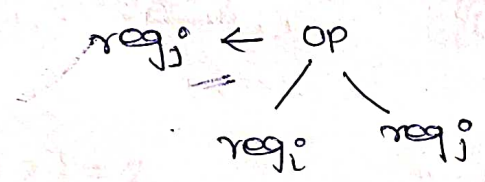
mov R<sub>j</sub>, R<sub>i</sub>



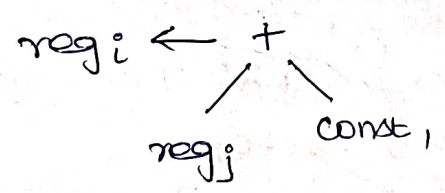
mov C(R<sub>j</sub>), R<sub>i</sub>



op C(R<sub>j</sub>), R<sub>i</sub>



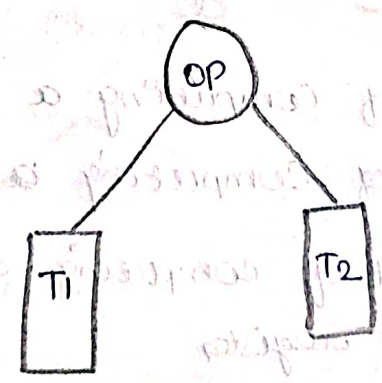
mov R<sub>j</sub>, R<sub>i</sub>  
op R<sub>j</sub>, R<sub>i</sub>



mov const, R<sub>j</sub>  
op R<sub>j</sub>, R<sub>i</sub>

## Dynamic programming code Generation:

Principle of dynamic programming



Tree T



150

Its working in three phases.

1<sup>st</sup> phase the cost vector is computed in the bottom up fashion for the constructed exp. T<sub>2</sub>

2<sup>nd</sup> phase of alg, the T<sub>tree</sub> + is computed using cost vector to determine.

3<sup>rd</sup> phase traverse each subtree using cost vectors and associated instructions.

Eq: By considering the instruction set each of one cost and two register R<sub>0</sub> and R<sub>1</sub> are available the cost vector for exp.

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

$$R_i = M$$

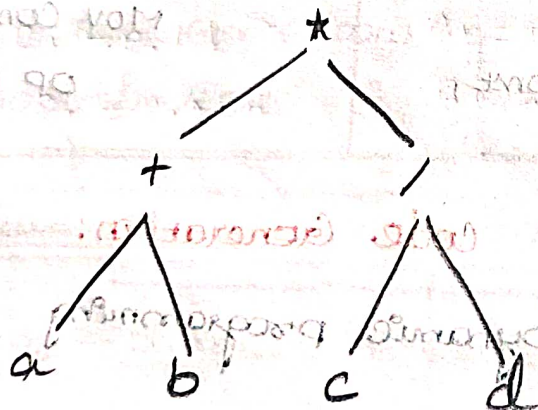
$$R_i = R_i \text{ OP } R_j$$

$$R_i = R_i \text{ OP } M$$

$$R_i = R_j$$

$$M = R_i$$

**Solution:** For node labelled as 'a'



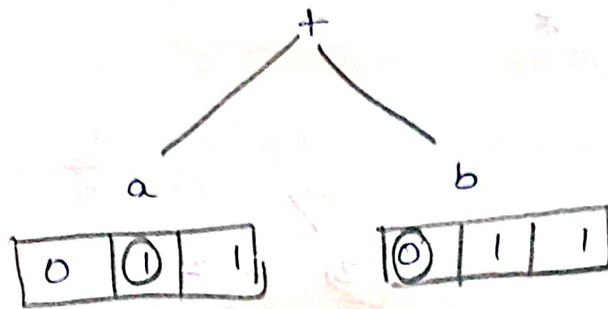
$c[0] = 0$  The cost of computing a in mem.

$c[1] = 1$  The cost of computing a in register

$c[2] = 1$  The cost of computing a into two available registers

$C[0] = 3$  This cost of Computing + Ento max (5)

$$R_0 = R_0 + M$$

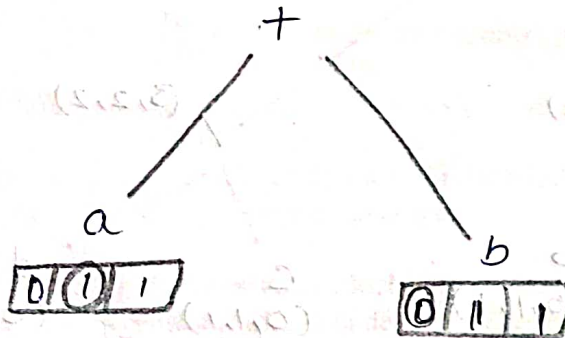


$$\text{Total cost} = 1 = 2 + 1 = 3$$

$$C[1] = 2$$

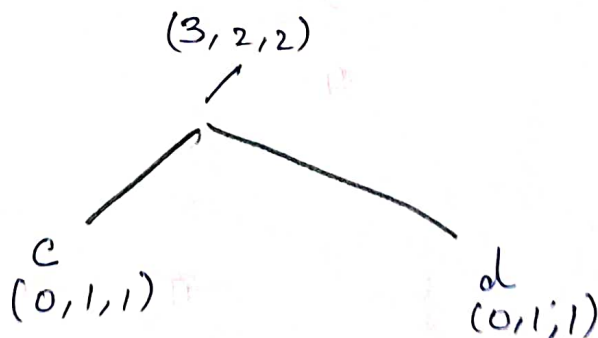
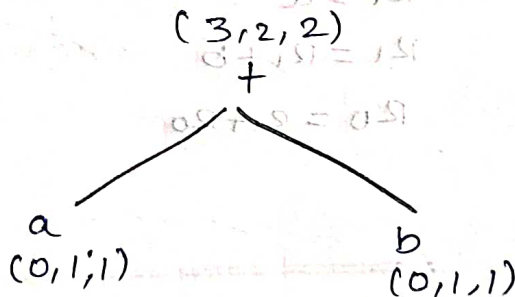
$$R_0 = R_0 + M$$

$$\Rightarrow 1 + 0 + 1 = 2$$



$$C[2] = 2$$

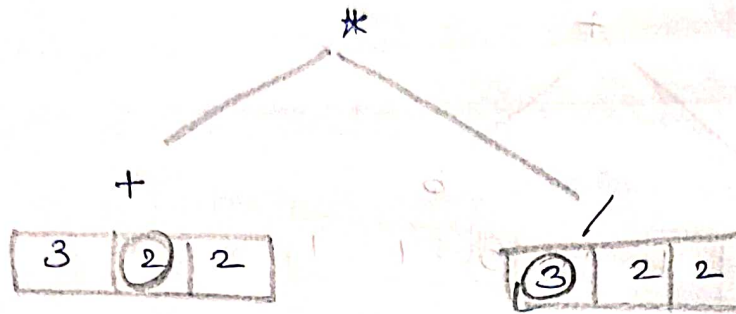
$$C[2] = 1 + 0 + 1 = 2$$



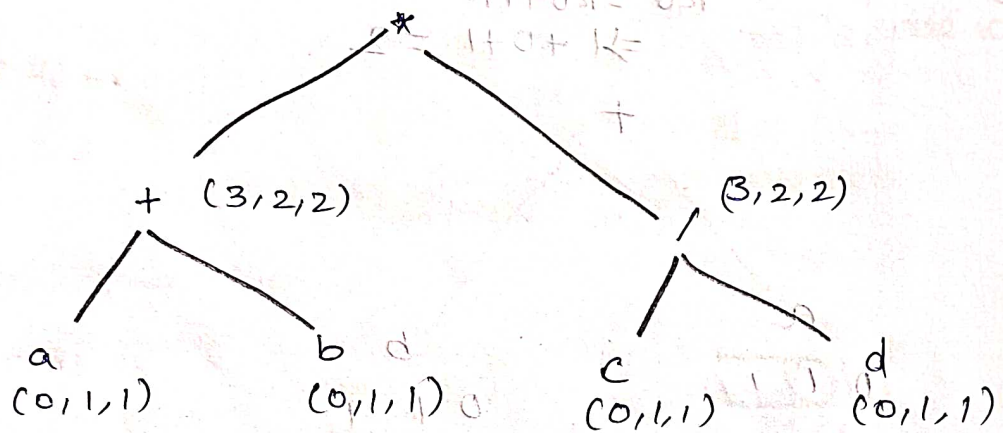


$$e[1] = 6$$

$$C[1] = 2 + 3 + 1 = 6$$



$(6, 6.5)$



$$R_0 = a \quad \text{or}$$

$$R_0 = R_0 + b$$

$$R_1 = \mathbb{C}$$

$$R_1 = R_1 / d$$

$$R_1 = R_1 \wedge R_0$$

$$R_0 = C$$

$$r_0 = r_0/d \cdot 0 + 1 = [1.0]$$

$$R_1 = a$$

$$R_1 = R_1 + b$$

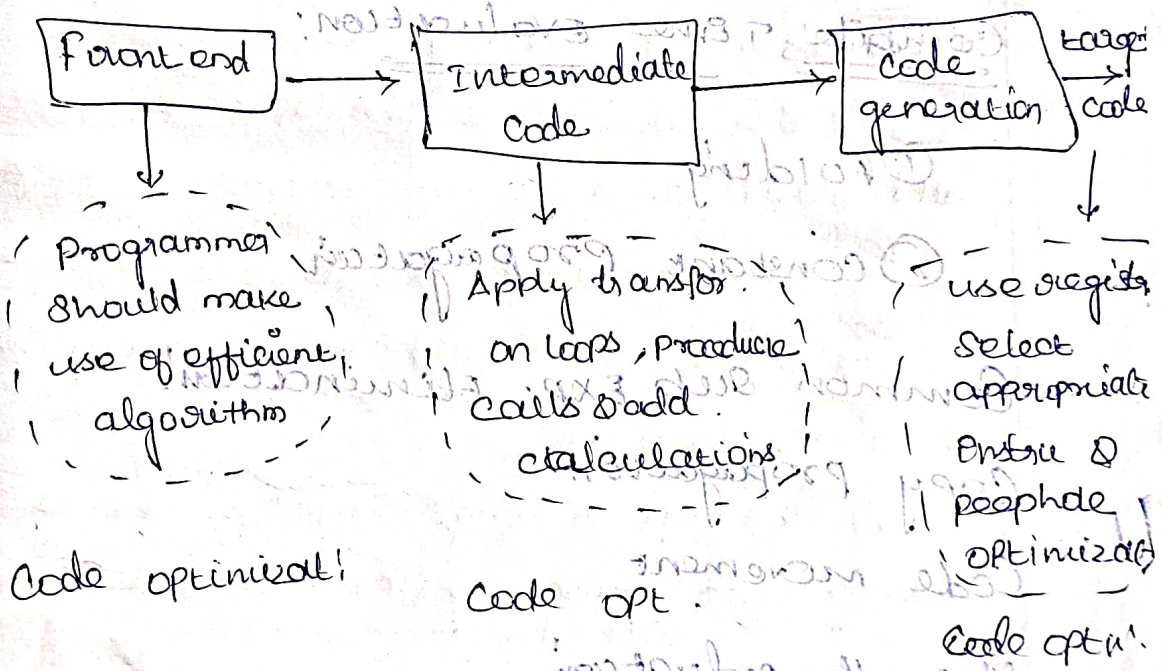
$$R_0 = R_1 + R_2$$

# Code Optimization

The Code optimization is a technique required to produce an efficient target code.

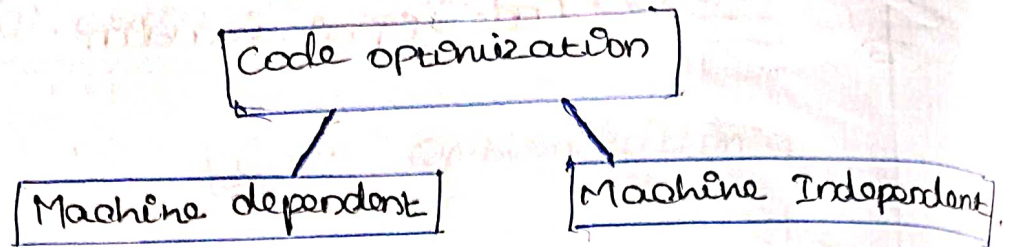
Two important issues:

- ① The semantic equivalence of the source pgm must not be changed.
- ② Pgm. efficiency must be achieved without changing the algorithm of the pgm.





## classification of optimization:



Machine dependent optimization is based on characteristics of the target machine for the instruction set used and addressing modes used for the instructions to produce the efficient target code.

Machine independent optimization is based on the characteristics of the programming languages for appropriate programming structures.

### Properties of optimizing compilers:

- ① The source code should be such that it should produce minimum amount of target code
- ② There should not be any unreachable code
- ③ Dead code should be completely removed from source language.
- ④ The optimizing compilers should apply following code improving on source languages.
  - i) Common subexp. elimination
  - ii) Dead code elimination
  - iii) Code movement
  - iv) Strength reduction.

## Principle source of optimization:

155

Optimization can be done locally or globally. If the transformation is applied on the same basic block then that kind of transformation is done locally otherwise transformation is done globally.

### Function preserving transformations:

- 1) Common subexpression elimination.
- 2) Copy propagation
- 3) dead-code elimination
- 4) Constant folding

### Compile time Evaluation:

#### 1) Folding:

In the folding technique the computation of constant is done at compile time instead of execution time.

length =  $(22/7) * d$

#### 2) Constant propagation:

In this technique of the value of variable is replaced and computation of an expression is done at the compilation time.

$$P = 3.14$$

$$r = 5$$

$$Area = P * r * r$$

$$Area = 3.14 * 5 * 5$$

#### 2) Common sub expression Elimination:

An expression appearing repeatedly in the pgm which is computed.



Sub expression eliminated!

$$t_1 = 4 * i$$

$$t_2 = a[t_1]$$

$$t_3 = 1 * j$$

$$t_4 = 1 * i$$

$$t_5 = n;$$

$$t_6 = b[t_4] + t_5$$

$$t_1 = 4 * i$$

$$t_2 = a[t_1]$$

$$t_3 = 4 * j$$

$$t_5 = n$$

$$t_6 = b[t_4] + t_5$$

### Copy propagation:

variable propagation means use of one variable instead of another.

Eg:  $x = pi;$

$\vdots$

$$area = x * r * r;$$

Optimization using variable propagation:

$$area = pi * r * r;$$

### Code move more:

① To reduce the size of the code.

② To reduce the frequency of exe. code.

for ( $i=0; i <= 10; i++$ )

$$x = y * 5;$$

$$k = (y * 5) + 50;$$

}

Optimized:  $z = y * 5$

for ( $i=0; i <= 10; i++$ )

$$x = z;$$

$$k = z + 50;$$

}

## Strength reduction:

The strength of certain operators is higher than others.

```
for (i=1; i<=50; i++)
```

```
{
```

```
...
```

```
    Count = i * 7;
```

```
...
```

```
}
```

The code can be replaced by using strength reduction:

```
temp = 7;
```

```
for (i=1; i<=50; i++)
```

```
{
```

```
...
```

```
    Count = temp;
```

```
    temp = temp + 7;
```

```
...
```

```
}
```

## Dead code elimination:

The variable is said to be live in a pgm if the value contained into it is used subsequently.

for eg:

```
i = j;
```

```
x = i + 10;
```

```
...
```

The optimization can be performed by eliminating the assignment stmt `i = j`. The assignment stmt is called dead assignment.

```
i = 0;
```

```
if (i == 1)
```

```
{
```

```
    a = x + 5;
```

```
}
```



## Peep-hole - Optimization:

Peep-hole optimization is a simple and effective technique for locally improving target code.

### Characteristics of peephole optimization:

#### ① Redundant Instruction elimination:

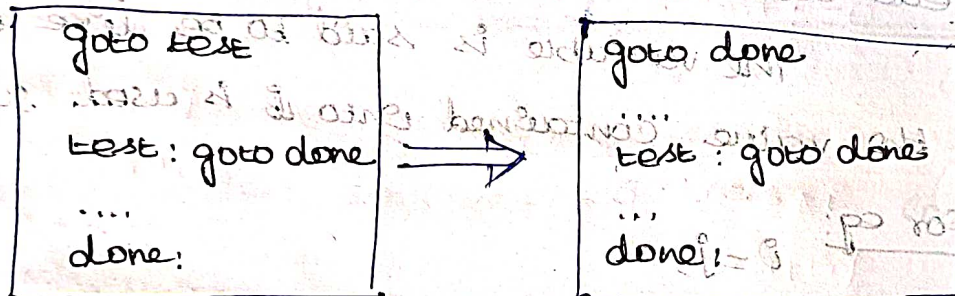
Redundant loads and stores can be eliminated in this type of transformations.

Eg: `MOV R0, x`

`MOV x, R0`

#### ② Flow of control optimization:

Unnecessary jumps or jumps can be eliminated.



#### ③ Algebraic simplification:

Peephole optimization is an effective technique for algebraic simplification.

$x := x + 0$

(or)

$x := x * 1$

#### ④ Reduction in strength:

If we can replace these instructions by cheaper instructions.

Eg:  $x^2$  is cheaper than  $x * x$ .

#### ⑤ Machine idioms:

The target instructions have equivalent machine instructions for performing some operations. Hence we can replace these target instructions by equivalent machine instructions.

#### 5) Basic blocks and flow graphs:

The basic block is a sequence of consecutive stmts in which flow of control enters at the beginning and leaves at the end.

Block 1:  
 $t_1 = a + 5$   
 $t_2 = t_1 + 7$   
 $t_3 = t_2 - 5$   
 $t_4 = t_1 + t_3$   
 $t_5 = t_2 + t_4$

Block 2:  
 $s_x = t_1 + 2$

Block 3:  
 $t_1 = t_2 + 7$

Block 4:  
 $s_x = t_3 + 8$

Block 5:  
 $t_1 = t_2 + 8$

Block 6:  
 $t_1 = t_2 + 8$

Algorithms for partitioning into blocks:

① First determine the leaders by using following rules

⇒ The 1<sup>st</sup> stmt is leader.

⇒ Any target stmt of conditional or unconditional goto is a leader

⇒ Any stmt that immediately follow a goto or unconditional goto is a leader.



(160)

2) The basic block is formed starting at the leader stmt & ending just before the next leader stmt.

Eg: Two vectors a and b length 10 and partition it into basic blocks.

**Solution:**

prod = 0;

i = 1;

do

prod = prod + a[i] \* b[i]

i = i + 1

while (i <= 10);

**Solution:** Equivalent three address code :-

1. prod := 0

2. i := 1

3. t1 := 4 \* i

4. t2 = a[t1]

5. t3 = 4 \* i

6. t4 = b[t3]

7. t5 = t2 \* t4

8. t6 = prod + t5

9. prod = t6

10. t7 = i + 1

11. i = t7

12. if i <= 10 goto (3)

**Block:1**

1. prod := 0

2. i := 1

**Block:2**

3. t1 = 4 \* i

4. t2 = a[t1]

5. t3 = 4 \* i

6. t4 = b[t3]

7. t5 = t2 \* t4

8. t6 = prod + t5

9. prod = t6

10. t7 = i + 1

11. i = t7

12. if i <= 10 goto (3)

# DAG Representation:

161

Construct DAG for

$sum = 0;$

for ( $i = 0; i < 10; i++$ )

$sum = sum + a[i]$

Solution:

1)  $sum = 0$

6)  $sum = t_3$

2)  $i = 0$

7)  $t_4 = i + 1$

3)  $t_1 = a * i$

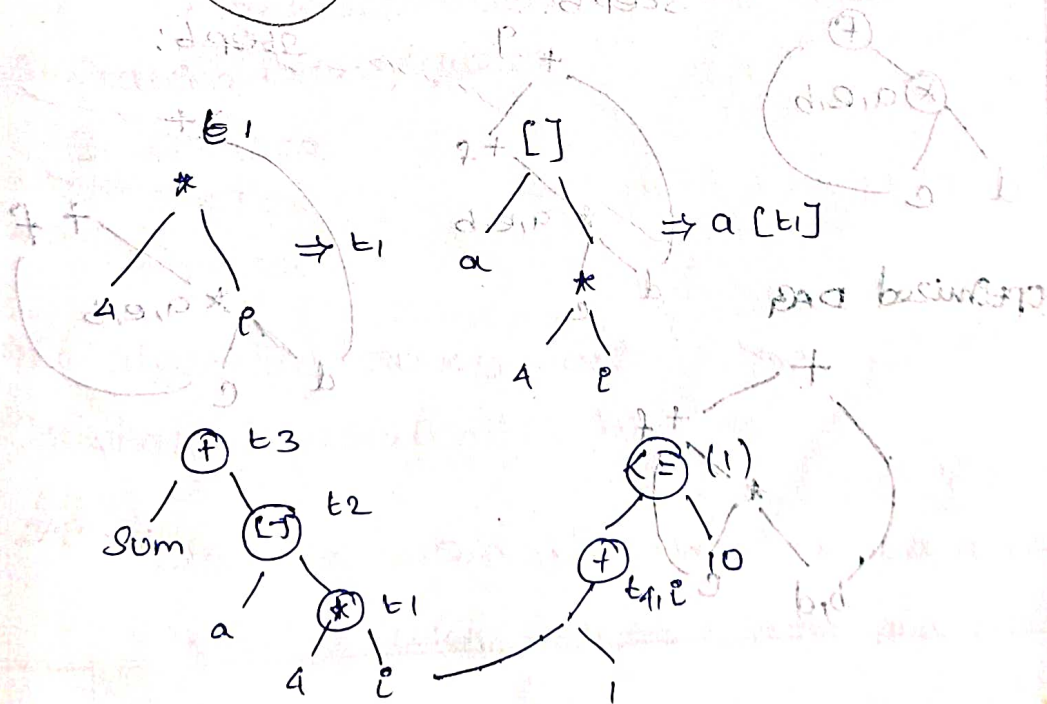
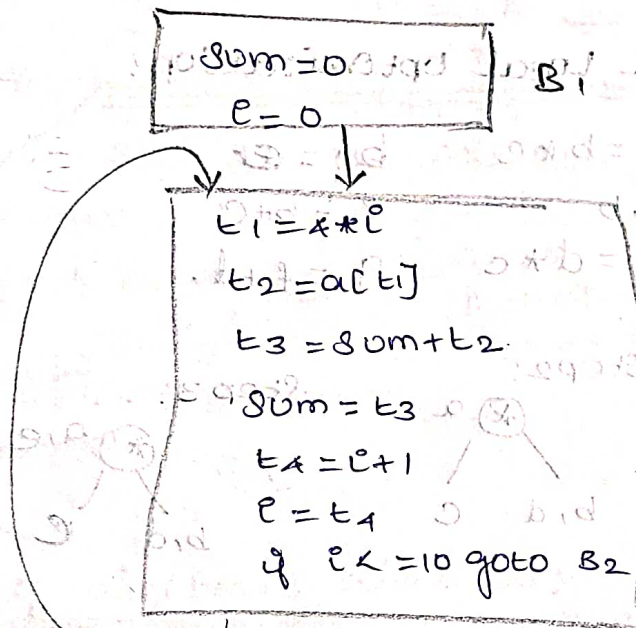
8)  $i = t_4$

4)  $t_2 = a[t_1]$

9) if  $i < 10$  goto (3)

5)  $t_3 = sum + t_2$

Block  $B_2$  for construction of DAG:





# Algorithm for Construction of DAG:

case (i)  $x \leftarrow y \text{ op } z$

Case ii)  $x := \text{op } y$

case iii)  $x := y$

## Applications of DAG:

- ① determining the common sub-expressions.
- ② determining which names are used inside the block and computed outside the block.
- ③ determining which stmts of the block could have their computed value outside the block.

## DAG based Local optimization:

$a := b * c$

$b := d$

$d := b$

$f := b + c$

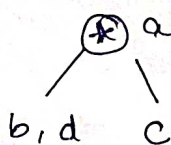
$e := d * c$

$g := f + d$

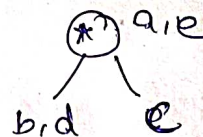
Step 1:



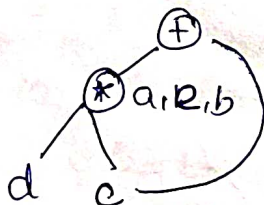
Step 2:



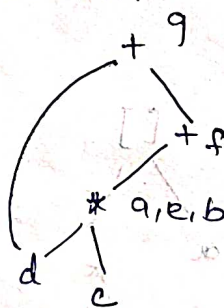
Step 3:



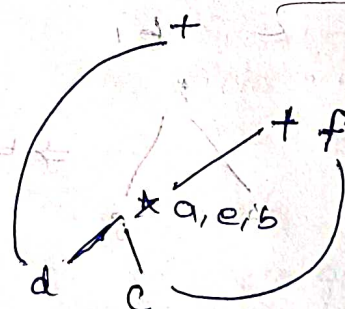
Step 4:



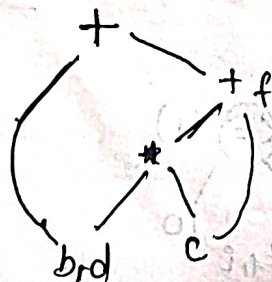
Step 5:



Step 6:



Optimized DAG



## 7) DAG - Optimization of basic block:

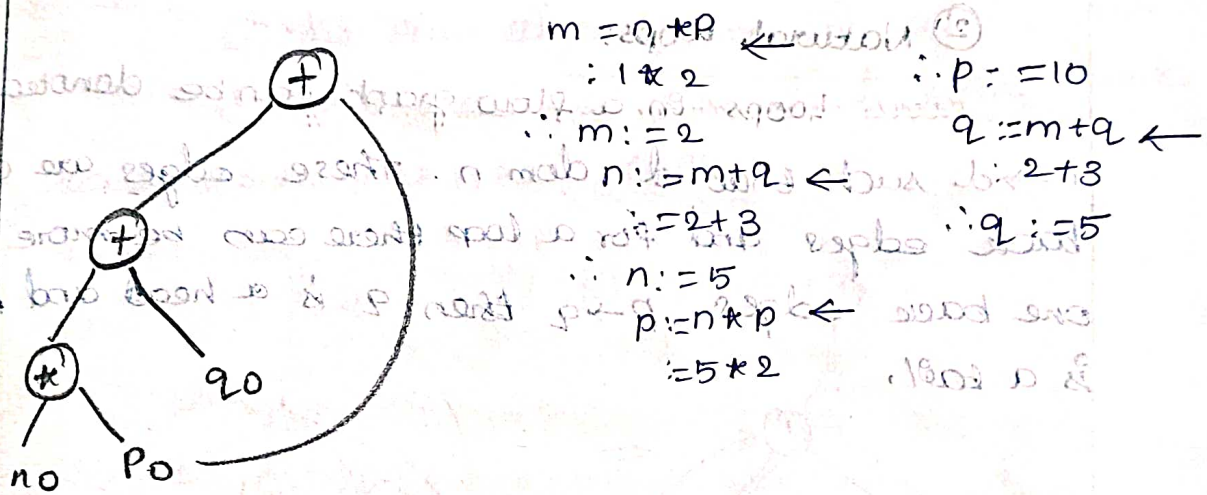
(163)

Two types of optimizations that can be done on basic blocks.

### ① Structure preserving transformation:

Structure preserving transp. can be applied by applying some principle techniques such as common sub expression elimination, variable and constant propagation, code movement, dead code elim.

Eg:  $m := n * p$  If we assume the values of  
 $n := m + q$   $n = 1$   
 $p := n * p$   $p = 2$  and  
 $q := m + q$   $q = 3$  then expressions becomes



### 2) Use of algebraic identities:

Algebraic identities are used on peephole optimization techniques.

Eg:  $a + 0 = a$   
 $a * 1 = a$   
 $a / 1 = a$

The algebraic transformation can be obtained using the strength reduction technique.

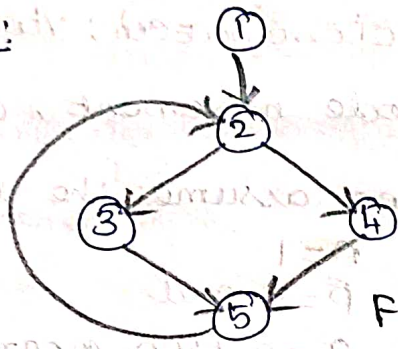
Eg: Instead of using  $2 * a$  we can use  $a + a$   
Instead of using  $a / 2$  we can use  $a * 0.5$



## 8) Loops in Flow Graph:

① Dominators: In a flow graph, a node  $d$  dominates  $n$  if every path to node  $n$  from initial node goes through  $d$  only. This can be denoted as ' $d \text{ dom } n$ '.

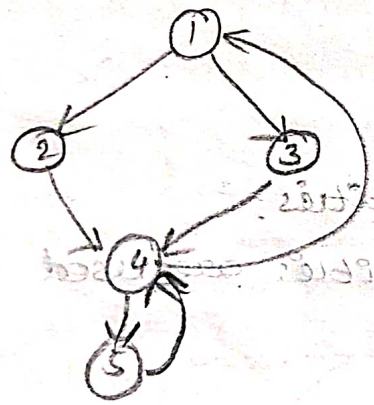
Eg:



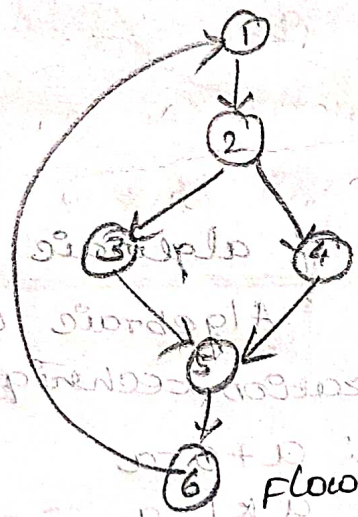
Flow graph.

## ② Natural loops:

Loops in a flow graph can be denoted by  $n \rightarrow d$  such that  $d \text{ dom } n$ . These edges are called back edges and for a loop there can be more than one back edges.  $p \rightarrow q$  then  $q$  is a head and  $p$  is a tail.



Flow graph with loops



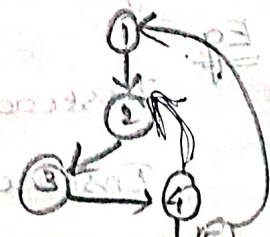
Flow graph with natural loop.

## 3) Inner loops:

The inner loop is a loop that contains no other loop.

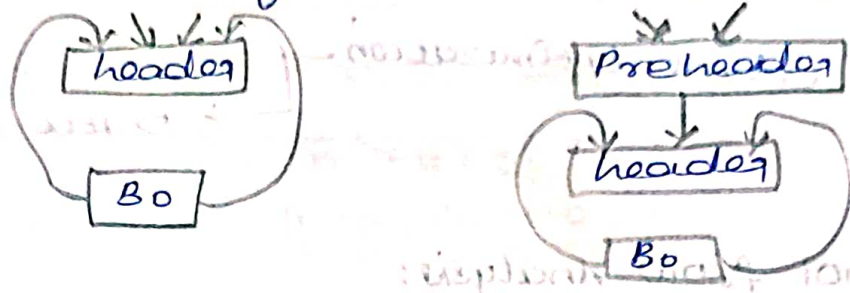
Here the inner loop is

$4 \rightarrow 2$  that means edge given by  $2-3-4$



#### 4) Pre-header:

The pre-header is a new block created such that successor of this block is the header block.

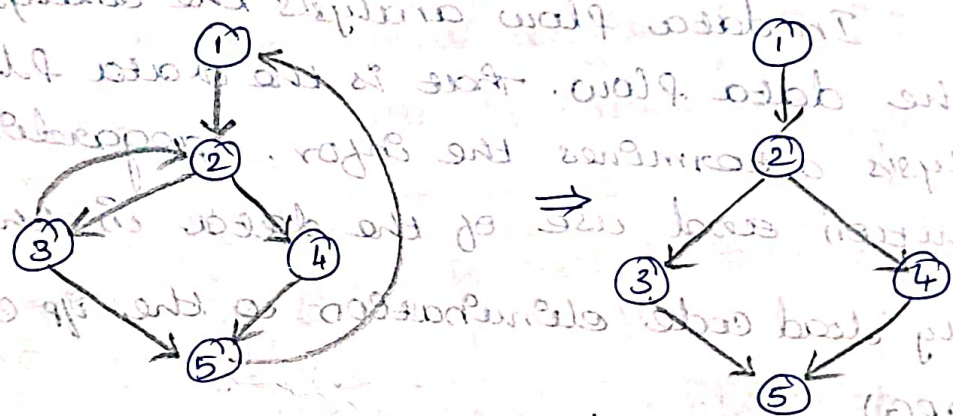


#### 5) reducible flow graph:

The reducible flow graph is a flow graph in which there are two types of edges forward edges and backward edges.

properties:

- ① The flow edge from an acyclic graph.
- ② The back edge are such edge whose head dominates their tail.



Flow graph.

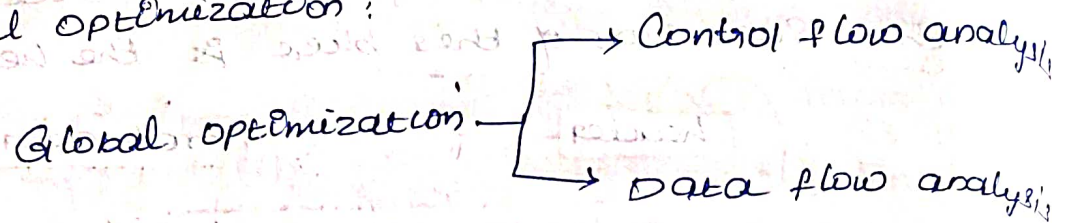
#### 9) Global data flow Analysis:-

The local optimization has a very restricted scope on the other hand the global optimization is applied over a broad scope such as procedure or function body.



(166)

There are two types of analysis performed for global optimization:



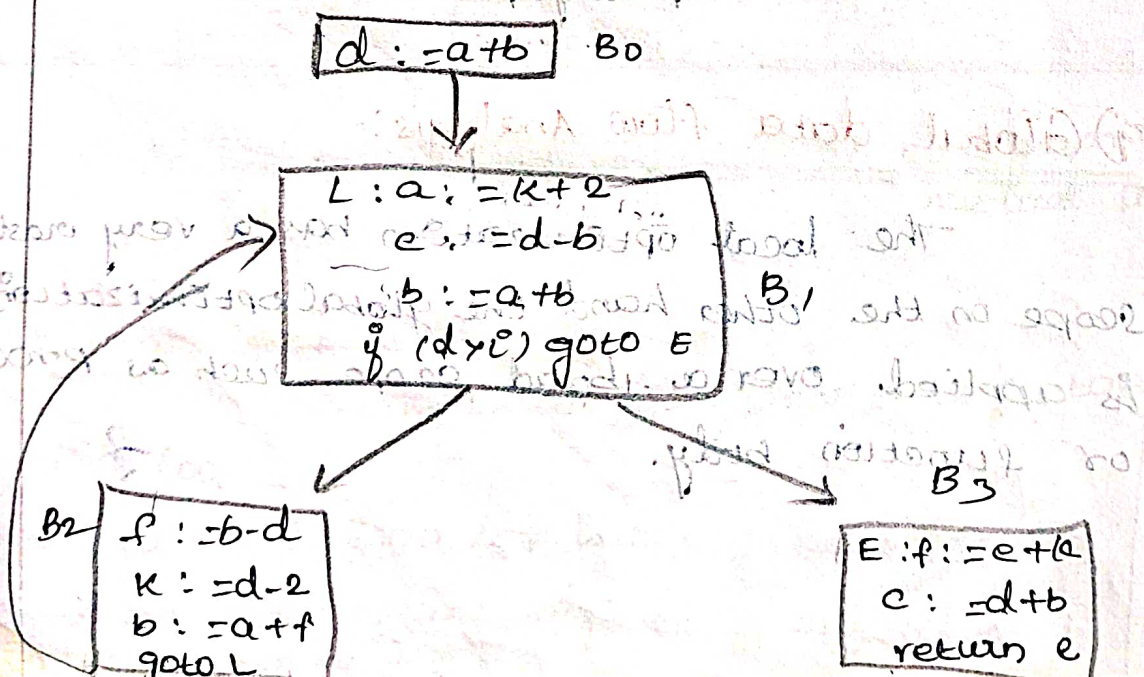
### Control flow Analysis:

The control flow analysis determines the info. regarding arrangement of graph nodes, presence of loops, nesting of loops and nodes visited before execution of a specific node.

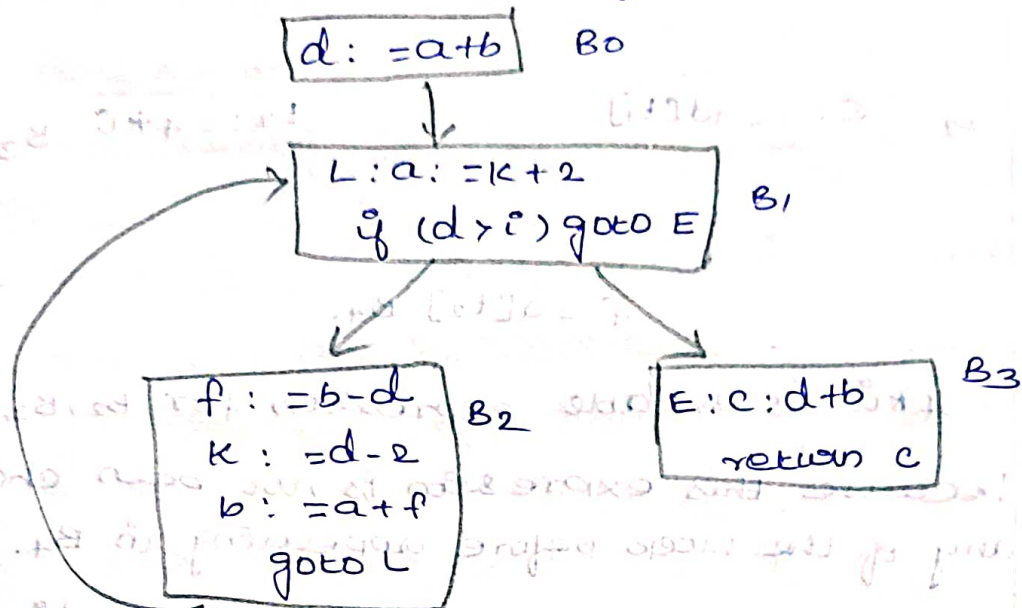
Thus Control flow analysis is made on the flow of control by carefully examining the program flow graph.

In data flow analysis the analysis is made on the data flow. That is the data flow analysis determines the info. regarding the definition and use of the data in the pgm.

Apply dead code elimination to the Cpp code given in (CFG).



dead code  $c := d - b$ , similarly  $B_3$  block (167)  
 $f := e + k$  has no meaning.

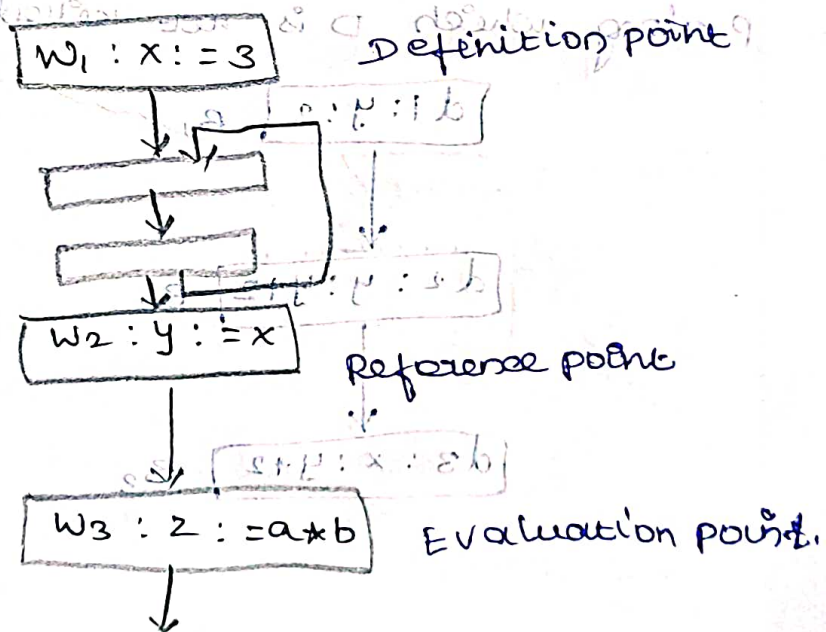


### Data flow properties:

A pgm point containing the definition is called "definition point".

A pgm point at which a reference to a data item is made is called "reference point".

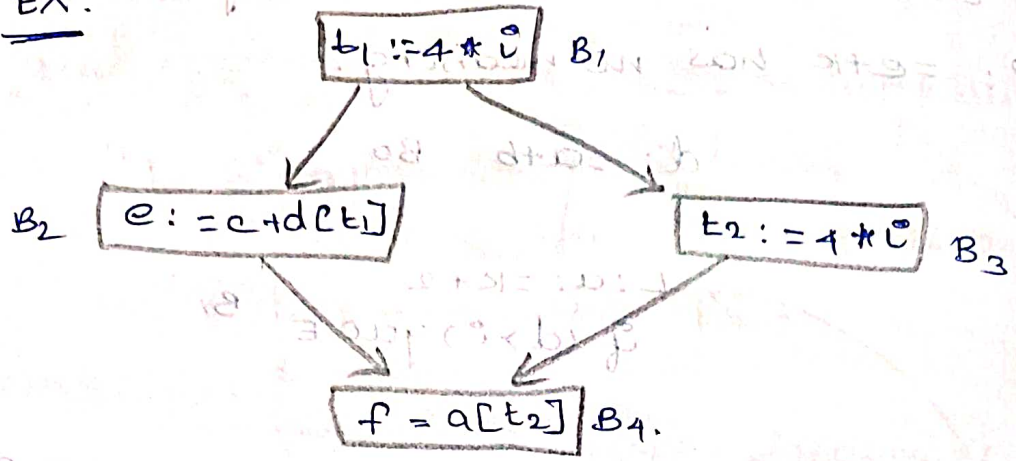
A pgm point at which some evaluating expression is given is called "evaluation point".





168

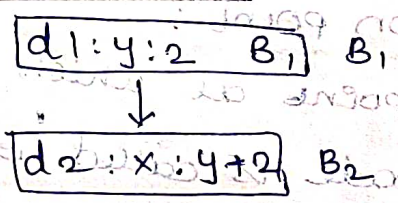
EX:



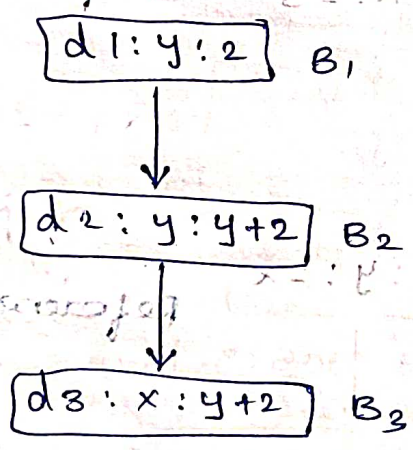
$4 * c$  is available expression for  $B_2, B_3, B_4$ , because this expression is not been changed by any of the block before appearing in  $B_4$ .

Adv: The use of available exp. is to eliminate common sub expressions.

## 2) Reaching definitions:



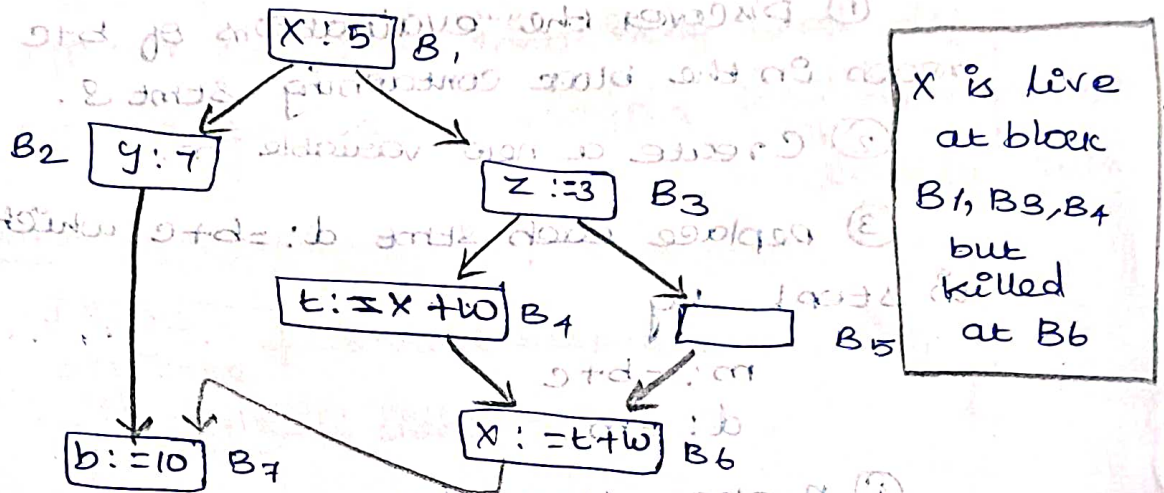
A definition  $D$  reaches at point  $P$  if there is a path from  $D$  to  $P$  if there is a path from  $D$  to  $P$  along which  $D$  is not killed.



Adv: Reaching definitions are used in constant and variable propagation. (69)

### 3) Live variable:

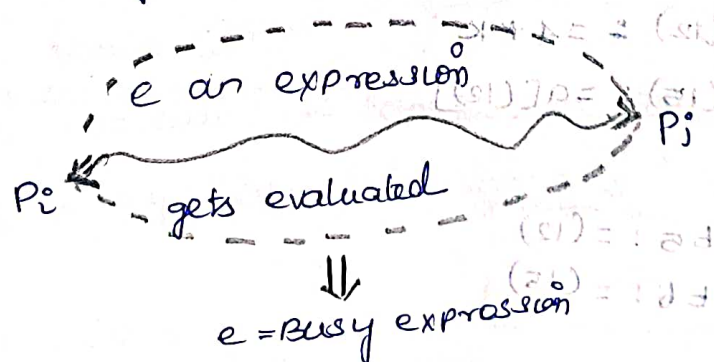
A variable  $x$  is live at some point  $p$ , if there is a path from  $p$  to the exit, along which the value of  $x$  is used before it is redefined.



Adv: ① Live variables are useful in reg. allocation.  
 ② Live variables are useful for dead code elimination.

### 4) Busy expression:

An expression  $e$  is said to be a busy expression along some path  $P_i \dots P_j$  if and only if an evaluation of  $e$  exists along some path  $P_i \dots P_j$  and no definition of any operand exists before its evaluation along the path.





(170)

## 11) Efficient Data Flow Algorithms: SSA

### Redundant Common Subexpression Elimination:

Alg: Algorithm for global common exp. elimin.

I/P: A flow graph with available expression.

O/P: A flow graph after eliminating common subexp.

- ① Discovers the evaluations of  $b+c$  that reach on the block containing stmt  $S$ .
- ② Create a new variable  $m$ .
- ③ Replace each stmt  $d := b+c$  which obtains in step 1. by  
 $m := b+c$   
 $d := m$
- ④ Replace stmt  $S$  by  $a := m$ .

Step 1:

$E1 := 4 * K$   
 $E2 := a[E1]$

$E5 := 4 * K$   
 $E6 := a[E5]$

Step 4:

$(12) := 4 * K$   
 $(15) := a[(12)]$

$E5 := (12)$   
 $E6 := (15)$

Step 2 & 3:

$m := 4 * K \rightarrow (12)$   
 $E1 := m \rightarrow (15)$   
 $E2 := a[E1]$

$E5 := m \rightarrow (12)$   
 $E6 := a[E5] \rightarrow (15)$

## Copy propagation:

(171)

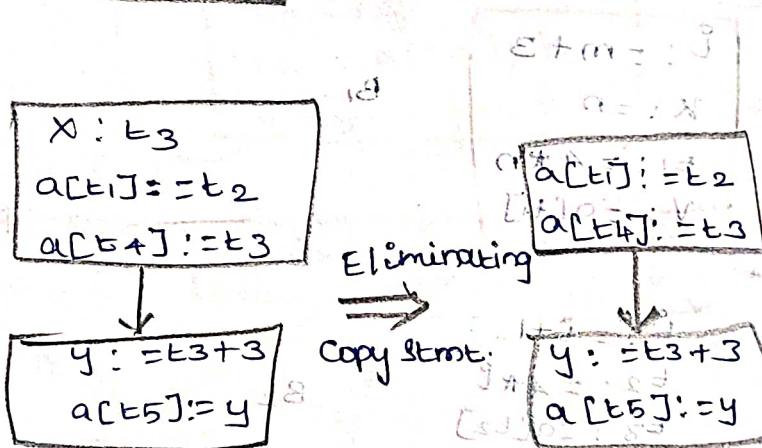
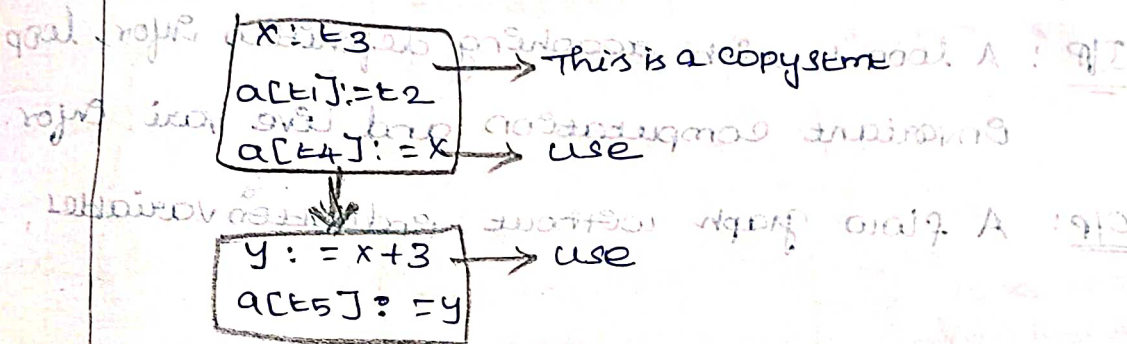
The assignment on the form  $a := b$  is called copy stmt. The idea behind the copy prop. is to ~~use~~ use  $b$  for  $a$  wherever possible after copy stmt  $a := b$ . Let us see the alg. for copy propagation.

Alg: copy propagation.

I/P: A flow graph containing used-def.

O/P: A graph after applying copy propagation transformation.

Step 1:



## Induction variable:

A variable  $i$  is called an induction variable of loop  $L$  if every time the variable  $i$  changes values.

$$a := i * b$$

$$a := b * i$$

$$a := i + b$$

$$a := b + i$$



172 a:  $= i * b$  then the triple for a is  $(i, b, 0)$ . we will understand this concept of writing triple with the help of block.

$t_1 := 0 + 1$   
 $t_2 := 4 * t_1$   
 $t_3 := a[t_2]$   
 $if\ t_3 < 10\ goto\ B_2$

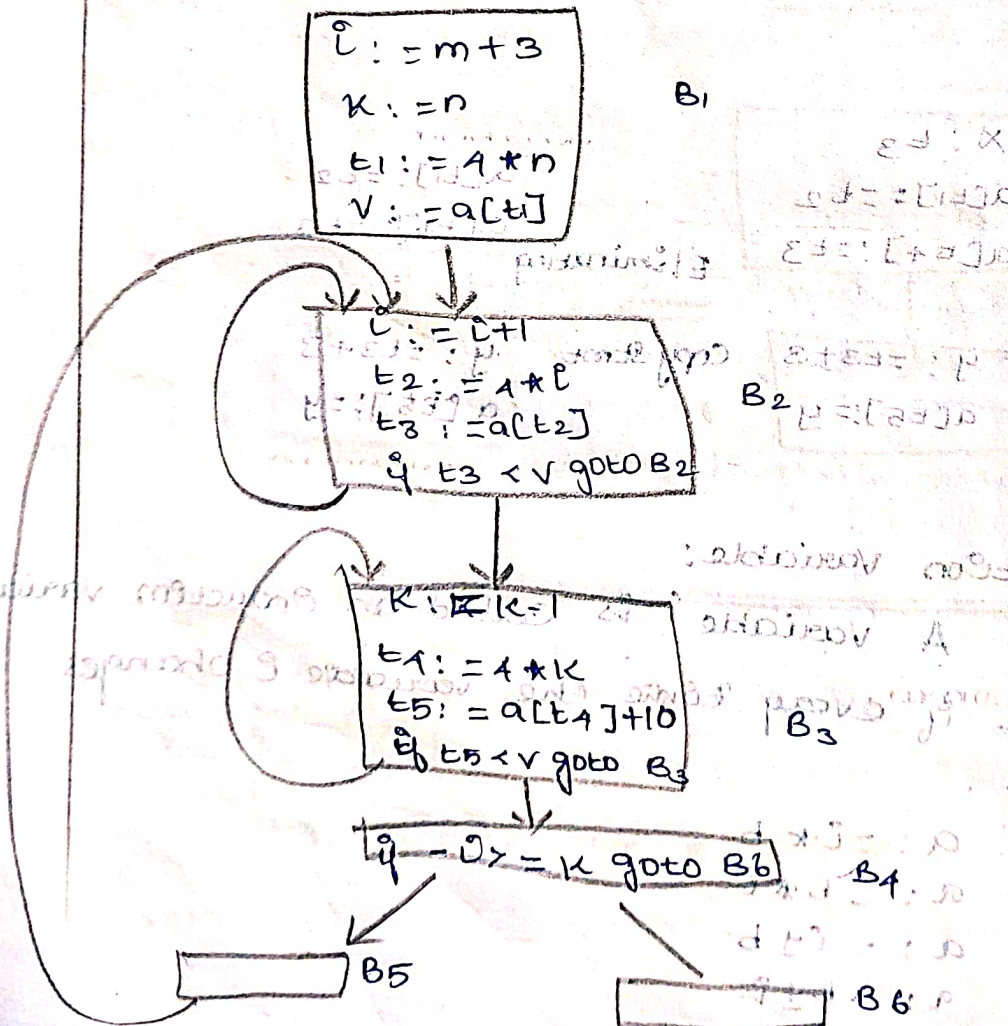
$(0, 4, 0)$

$\downarrow$                        $\downarrow$   
 Induction              Matching with  
 variable              constant b.

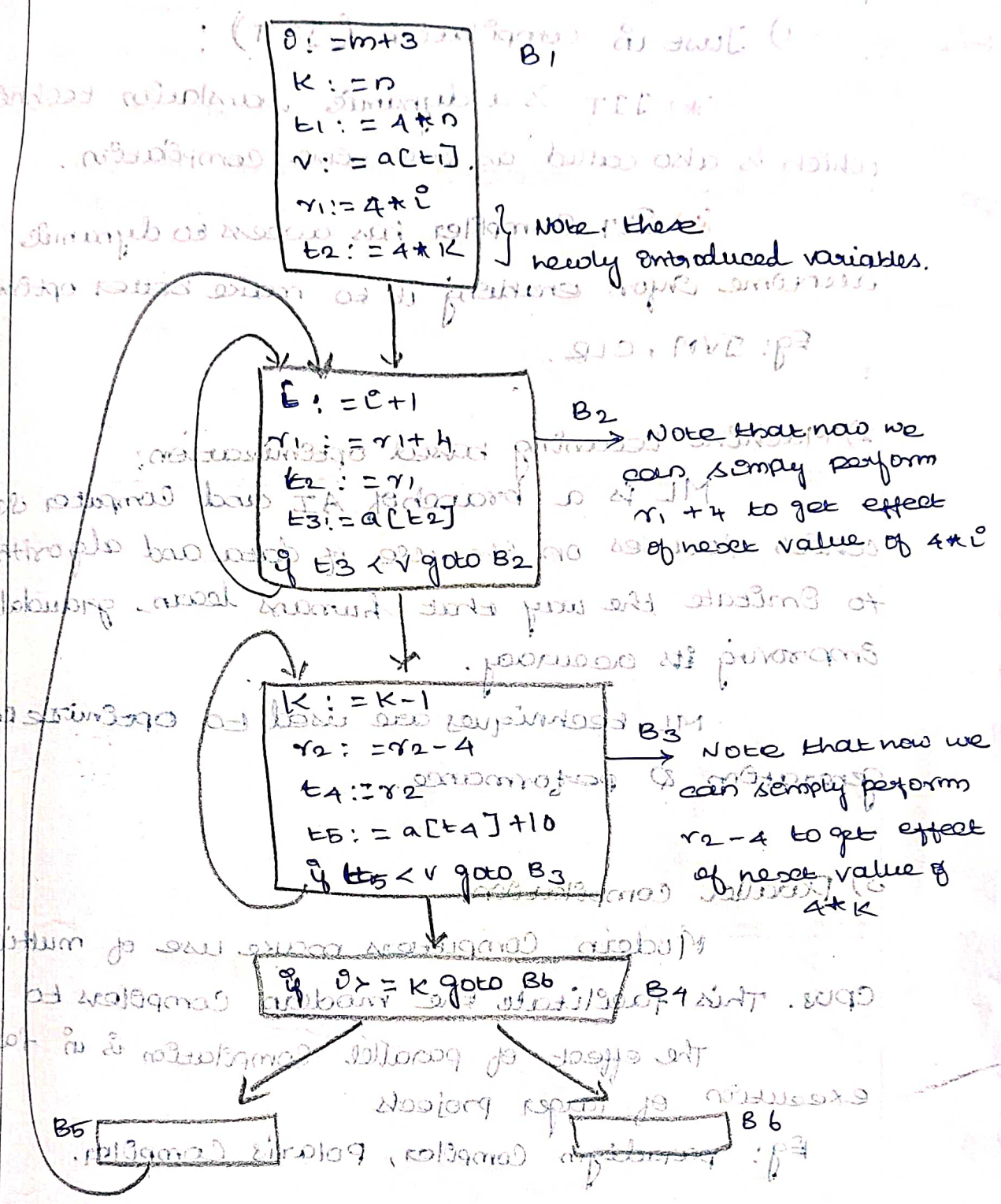
Alg: Elimination of induction variables

I/P: A loop  $L_i$  with reaching definition info, loop Invariant computation and live vari. info.

O/P: A flow graph without induction variables.



The flow can be rewritten as:



This allows the developer to use the same code development across diff. languages  
compiler tool to compile and optimize the  
code. The modern compiler language  
will be for writing the code.



(174)

## 12) Recent Trends in Compiler Design:

### 1) Just in Compilation (JIT) :

(\*) JIT is a dynamic translation technique which is also called as run time compilation.

(\*) JIT Compiler has access to dynamic runtime info. enabling it to make better optimization.  
Eg: JVM, CLR.

### 2) Machine Learning based optimization:

ML is a branch of AI and Computer Science which focuses on the use of data and algorithms to imitate the way that humans learn, gradually improving its accuracy.

ML techniques are used to optimize code generation & performance.

### 3) Parallel compilation:

Modern Computers make use of multicore CPUs. This facilitates the modern compilers to

The effect of parallel compilation is in fast execution of larger projects

Eg: paradigm compiler, Polaris compiler.

### 4) Language - Independent Compiler Infrastructure:

The support for multiple programming languages make the modern compilers language independent.

This allows the developers to use the same compilation tools to compile and optimize the code generation across diff. languages

5) Advanced compiler analysis and optimization: (175)

Today, we embrace new technology with great changes in architectural design, highly parallel processing, registers allocations, cache hierarchy and high processing speed demands for better compilers.

---