

unit - 1.

Introduction to Compilers.

The today's world depends up on programming languages, because all the softwares running on computers was written in some programming languages. But before a program can run, it first must be translated to a suitable form in which it can execute on computer.

The software that does this Translation is called as Compiler.

Language Pre processors.

Compiler is a program that can read a program in one language i.e; source language, and translate it into another language, known as target language.

An important role of the Compiler is to report any errors in source program that it detects during the Translation process.

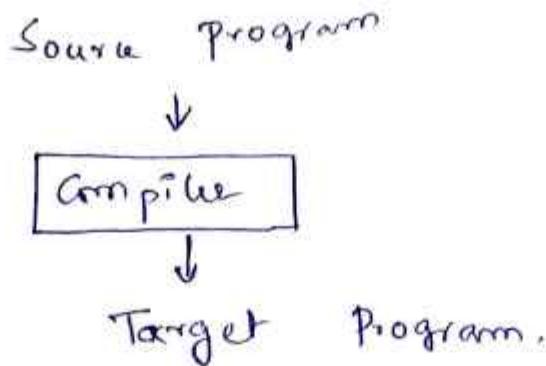
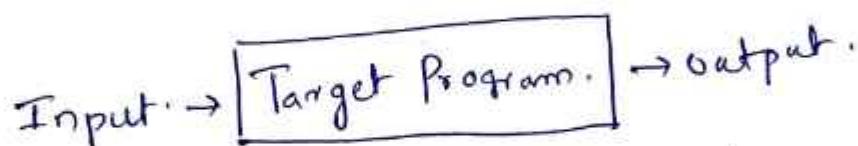
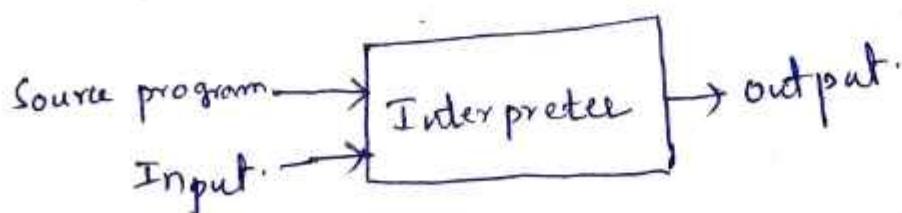


fig. compiler.

If target program is executable machine program, it can take the inputs from user and produce output.



Interpreter is another common kind of language pre processor, instead of producing target program, interpreter directly executes the operations specified in the source program, and on inputs provided by user.



Machine language Target program produced by interpreter.
a Compiler is much faster than interpreter.
But interpreter can give better error diagnosis

Then compile because it executes source program statement by statement.

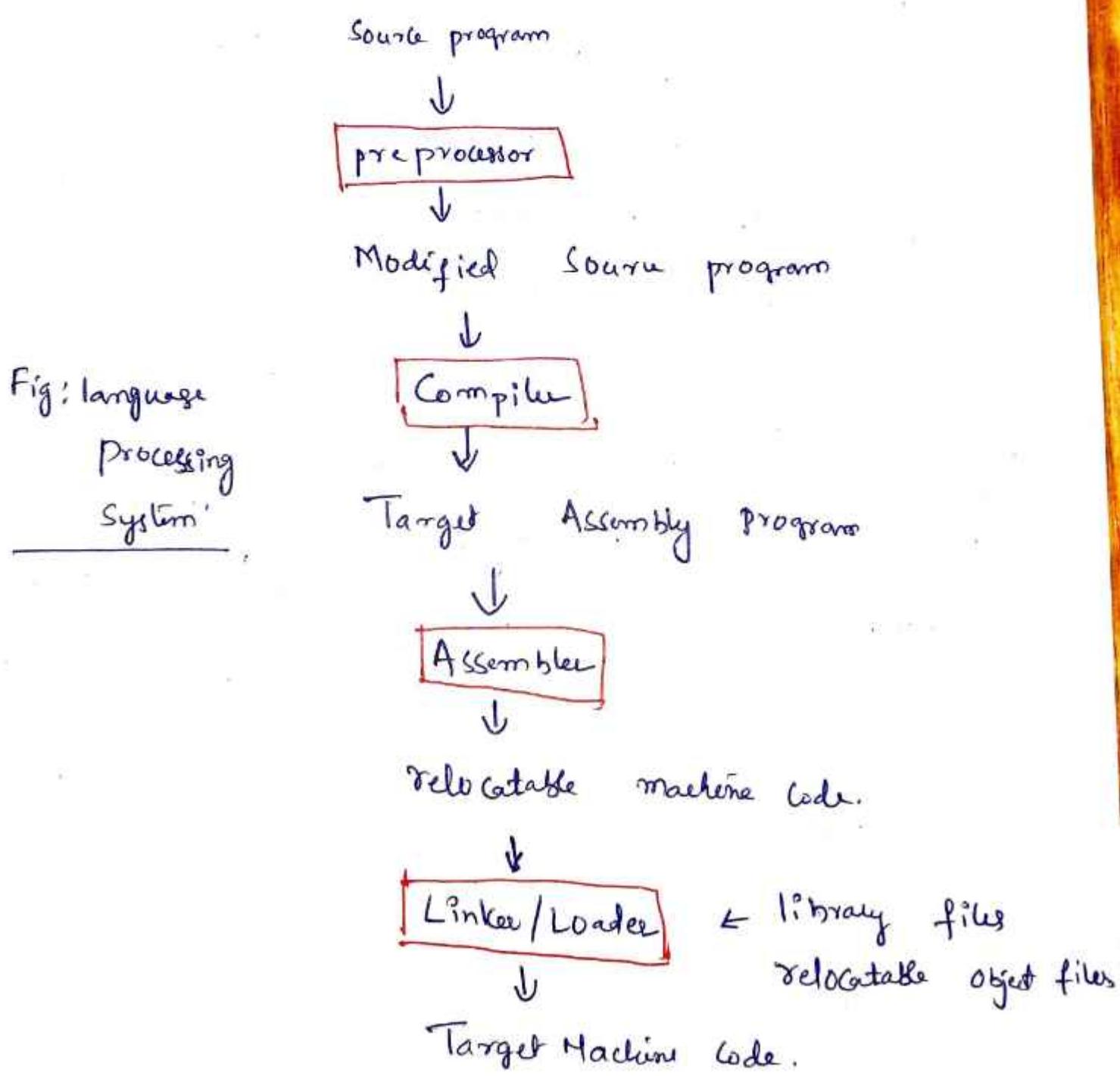
In addition to a compiler several other programs may be needed to create an executable target program.

1. A source program may be divided into several modules and stored in separate files. Preprocessor collects these files into one unit. Preprocessor may also expand short cuts, called as macros, into source language statements.

2. The modified source program is then fed to a Compiler. Compiler produces an Assembly language output. Because Assembly language program is easier to debug and to produce output.

3. The Assembly language is then processed by Assembler, that produces machine code as output.

4. Large programs are often compiled as pieces. The relocatable machine code has to be linked together. The linker links these files.
5. The loader puts together all executable object file into main memory for execution.



The Structure Of a Compiler

Compiler consists of two parts:

- i. Analysis.
- ii. Synthesis.

Analysis part breaks up source program into pieces and imposes a grammatical structure on them. It then uses this structure to create intermediate code representation of source program.

If the Analysis part detects that source program is syntactically or semantically weak, i.e., not following the rules of grammar, compiler provides informative messages to user in the format of errors, so that user can take corrective action.

This Analysis part collects the info about source program, and stores in a Symbol table, which is passed along with intermediate code to synthesis part.

The Synthesis part constructs desired code from intermediate form in target program and info in Symbol Table.

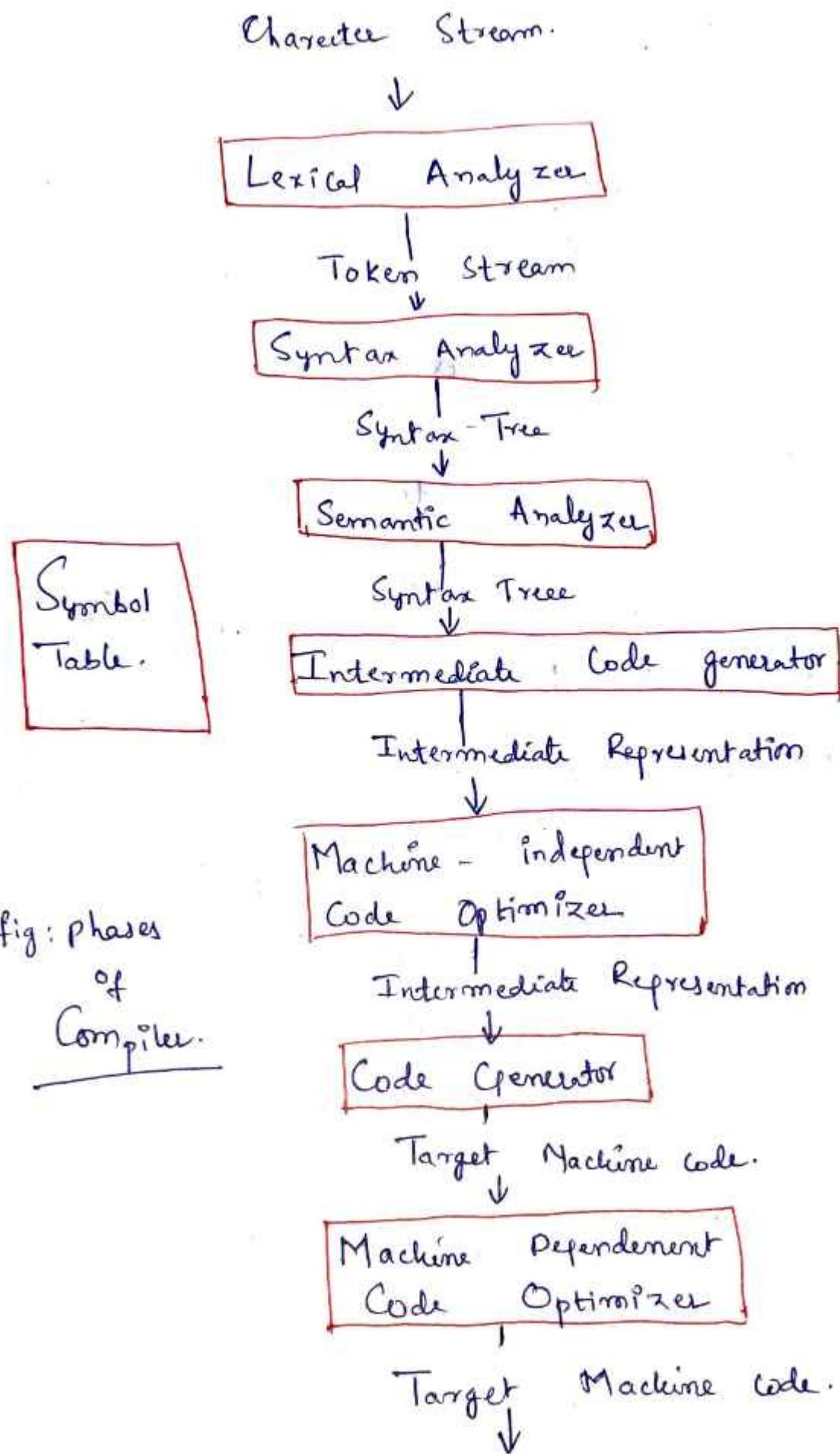
The Analysis part is often called as front end of the compiler, The Synthesis part is called as backend.

Compilation process operates as a sequence of phases, each of which transforms one representation to other. A typical decomposition of compiler into phases is shown below figure.

Several phases may be grouped together. The Symbol table which stores the entire source program info about the compiler is used by all phases of the compiler. Some compilers have optimization phase between the front end and backend.

The purpose of code optimization is to do some transformations in intermediate code so that backend can produce a better target program.

(4)



Lexical Analysis

The first phase of compiler is called as Lexical Analysis or scanning. The Lexical Analyzer reads the stream of characters in a source program, and groups them into meaningful sequences called as lexemes.

For each lexeme, the lexical analyzer produces as o/p a token of the form

$\langle \text{token-name}, \text{Attribute value} \rangle$

which is passed to next phase Syntax Analysis.

In token first component $\langle \text{token-name} \rangle$ is a symbol that is used during Syntax analysis, Attribute Value points to an entry in the Symbol table for this token. Info in Symbol table is needed for Semantic analysis and code generation.

For e.g. consider a source program contains assignment statement

position = initial + rate * 60.

(5)

1. position is a lexeme, that will be mapped to a token $\langle id, 1 \rangle$, where id stands for identifier, 1. points to symbol-table entry. This symbol-table entry holds info about identifier such as its name and type.
2. The assignment symbol '=' is a lexeme which is mapped to token $\langle = \rangle$, it will not have attribute value, so second component is omitted.
3. initial is a lexeme that is mapped to a token $\langle id, 2 \rangle$, where 2 points to symbol-table entry for initial.
4. + is a lexeme, mapped to token $\langle + \rangle$
5. rate is a lexeme, mapped to token $\langle id, 3 \rangle$, where 3 points to symbol-table entry for rate.
6. * is a lexeme, mapped to token $\langle * \rangle$
7. 60 is a lexeme, mapped to token $\langle 60 \rangle$

After Lexical Analysis The assignment statement can be represented as sequence of tokens as below

$$\boxed{\text{Position} = \text{initial} + \text{rate} * 60}$$

$$\boxed{<\text{id}, 1> <= > <\text{id}, 2> <+> <\text{id}, 3> * <60>}$$

Syntax Analysis.

The second phase of compiler is Syntax analysis or parsing. Parser uses first component of token produced by Lexical analyzer to create a tree like intermediate representation that depicts Grammatical structure of token stream.

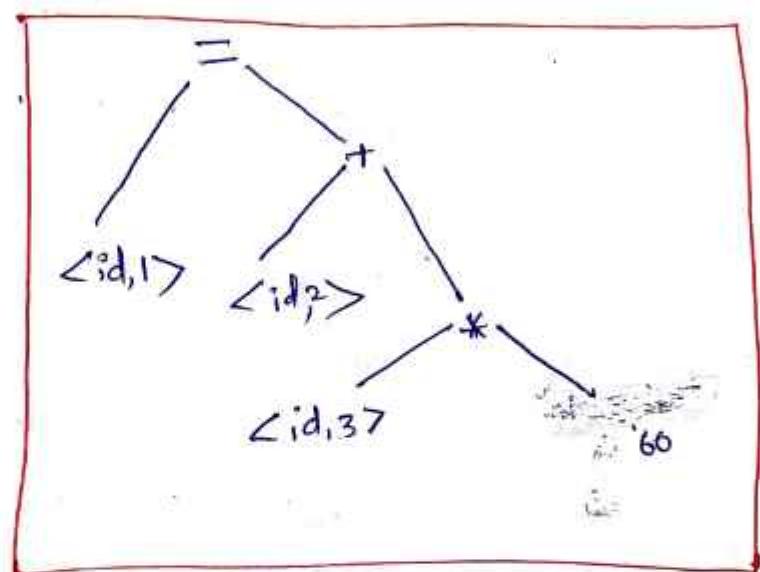
This tree structure is known as Syntax Tree.

The tree shows the order in which assignment is performed. For ex. in the statement $\text{position} = \text{initial} + \text{rate} * 60,$

The tree has a node $*$, with left child $<\text{id}, 3>$ as its left child, and 60 as its right child. The node $*$ tells that

We must first multiply rate by 60. ⑥

- The node + represent, we should add the result of multiplication, to the value of initial.
- The root '=' denotes the right child result should be stored in identifier position.



Syntax Tree.

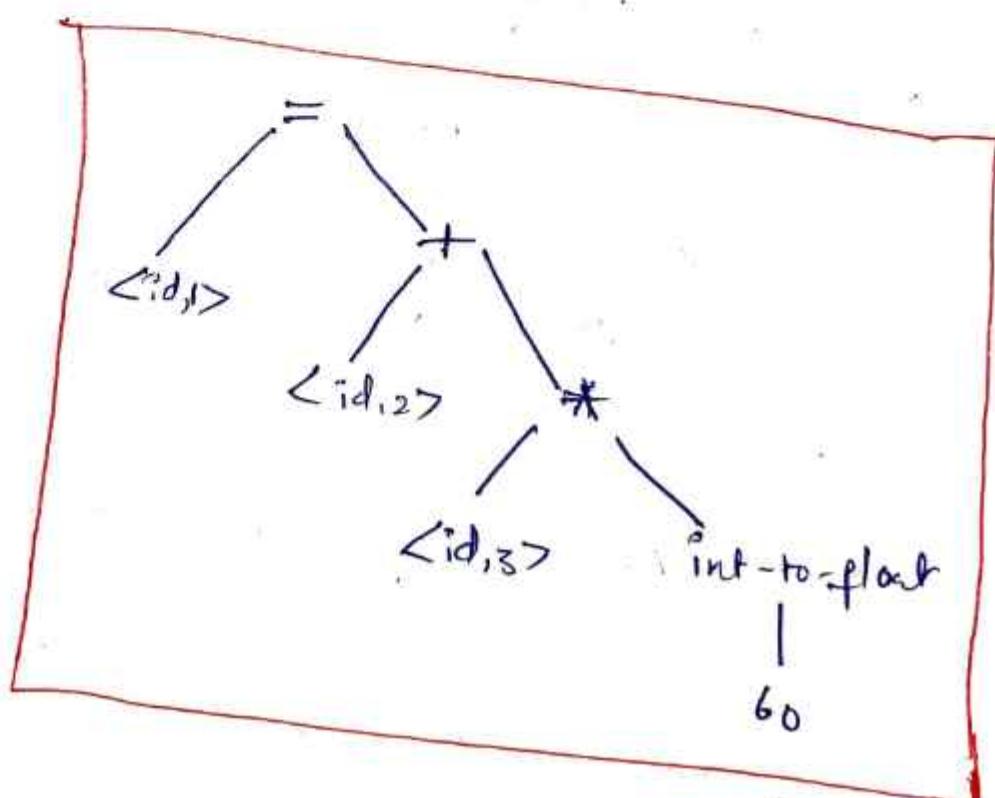
Semantic Analysis.

Semantic Analyzer uses the Syntax Tree and the info in symbol table to check the source program for semantic consistency with language definition. An important part of semantic analysis is type checking, compiler checks each operator has matching operands.

The language specification may permit type conversions known as coercion.

For ex., binary arithmetic operator may be applied to either two integers, or on two floating point numbers. If Operator is applied to floating point number and integer, the compiler will convert integer into a floating point number.

Such coercion appears in following syntax tree, we consider that position, initial and rate have been declared to floating point numbers, and 60 is an integer. The type checker in semantic analysis converts 60 to float.



Intermediate Code Generation.

While converting source program into machine code, a compiler may convert to one or more intermediate forms. Syntax Trees are a form of intermediate representations which are used in Syntax and Semantic phases.

After Syntax and semantic analysis, compiler will generate an explicit low-level (or) machine like ^{intermediate} representation. It will have two important properties:

- i. it should be easy to produce
- ii. it should be easy to Translate into the target machine code.

One format of intermediate code used is Three address code, which consists of Assembly like instructions with three operands per instruction. Each operand act like a register.

$$t_1 = \text{inttofloat}(60).$$

$$t_2 = id_3 * t_1$$

$$t_3 = id_2 + t_2$$

$$id_1 \rightarrow t_3.$$

In Three address code.

1. At most one operator should be present at right side of instⁿ.
2. To store the result temporary names are used at left side.
3. Some instructions may have fewer than 3. Operands.

Code Optimization.

Machine independent code optimization phase attempts to improve the intermediate code so that better target code will result. Better means faster, shorter, and which consumes less power.

In above intermediate code representation The optimizer can deduce the conversion of 60 from int to float, by replacing 60 by 60.0 t_3 is used only once, to transmit its value to id1, So optimizer can transform previous code as below.

$$t_1 = id_3 * 60.0$$

$$id_1 = id_2 + 1.$$

Code generation.

→ Code generator takes intermediate representation of source program and maps it into the target language.

→ If target language code is machine code of code registers or memory locations are selected for each variable used in the program.

→ An important task of code generation is assignment of registers to hold variables.

→ For e.g. previous optimized intermediate code , using Registers R₁ and R₂, can be translated to machine code as follows.

LDF R₂, id₃.

MULF R₂, R₂, #60.0

LDI R₁, id₂

ADD F R₁, R₁, R₂

STF id₁, R₁.

→ The first operand in each "inst" specifies a destination. F in each "inst" tells us that it deals with floating point numbers.

Symbol Table Management

- Compiler will record the variable names used in source program, and various info about attributes of each name.
- The Attributes may provide info about storage allocated for a name, type, scope, and in case of functions, The number of arguments, Type of arguments, parameter passing methods and the type returned.
- Symbol table is a data structure which contains a record for each variable, with fields for the attributes of that name.
- Symbol table is designed in such a way that Compiler will find the record for each name quickly and stores and retrieves data from that name quickly.

The Science of Building a Compiler.

- Compiler design is full of beautiful examples, where complicated real world problems are solved easily, mathematically. By using abstraction
- A compiler must accept all source programs that conforms to specification of programming languages. The programs can be large, possibly millions of lines of code.
- While converting such big programs to machine code the compiler must preserve meaning of the program.

Modelling in Compiler Design and Implementation.

Compiler design is a study of how to design a mathematical model and choose the right algorithm.

Some of most fundamental models are finite-state machines and regular expressions which are used to describe syntactic lexical units of programs (Keywords, identifiers etc). And Content-Free Grammars used to describe Syntactic Structure of programming languages.

The Science of Code Optimization.

- Optimization in Compiler design, attempts to produce code that is most efficient than obvious code.
- Optimization is complex because processor architecture becomes very complex.
- Optimization is important because, Parallel Computer needs substantial optimization. Otherwise their performance suffers by order of magnitude.

Programming Language Basics.

(10)

The static / Dynamic Distinction.

→ If a programming language uses a policy that allows a compiler to decide an issue, then we say language uses static policy, or that issue can be decided at compile time.

→ A policy that allows a decision to made when we execute a program is said to be a dynamic policy or to require a decision at runtime.

scope → scope of a variable x , is the region in which x is declared. A language uses static scope or lexical scope if we can say scope at time of declaration. In Dynamic scope, same use of x could refer to different declaration of x .

Ex:- Consider the use of term "static" in java,

to declare a variable, as follows.

```
public static int x;
```

Variable is a name of memory location used to hold data value.

Above statement declares `x` as a static, i.e., only one variable is created, where all objects access same memory location while accessing this variable.

If static is omitted, all objects refers to different memory locations.

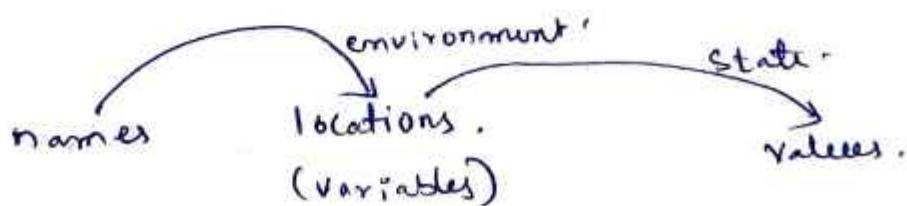
Environments and States

The association of names with location in memory, and then values can be described by two mapping which changes as the program runs.

1. Environment → Mapping from names to memory location is known as Environment.

Variable refers to location^(l-value). Environment is mapping b/w name and variable.

2. State → It is mapping from memory location to their values. State maps l-value to corresponding r-values in C terminology.



- Most bindings of names to locations is dynamic.
- The bindings of locations to values is generally dynamic as well as static.

for ex:-

```
#define ARRSIZE 1000.
```

binds the name after ARRSIZE to a value statically.

Static scope and block structure.

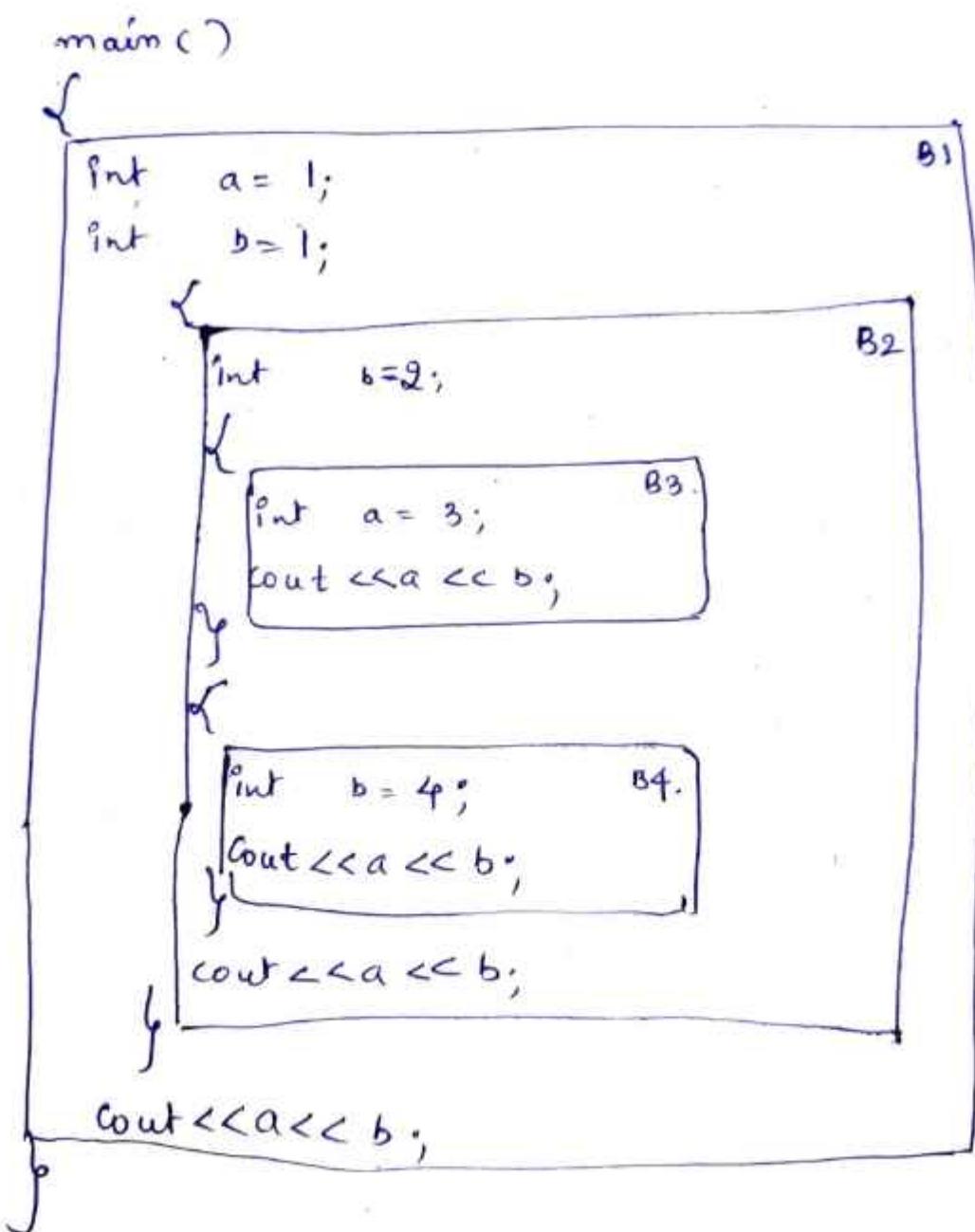
The scope rules for C are based on program structure. Scope of a variable is implicitly determined by where declaration appears in the program. In languages like C++, C# and java provides explicit control over scope through use of keywords like public, private and protected.

Static scope can be given by using blocks. Block is a grouping of declarations and statements. C uses braces '{' and '}' to denote a block. In some languages like Algo we use "begin" and "end" to denote a block.

Blocks can be nested inside other blocks. We say a variable declaration D, belongs to a block B, if it is declared in B. If the declaration D of a name x is in B, Then scope of x is B, as well as in any nested block. x will not have scope in

nested block if it is redeclared in inner block.

Eg:- Consider the program shown in below figure, which has four blocks. With several definitions of declaration initializers give variable to number of block to which it belongs.



(12)

For ex, Consider The Declaration of 'int a=1' in B₁, it Scope is entire block B₁, except the nested block of B₁, where 'a' is declared. B₂ is immediately nested under B₁, but does not contain redeclaration of a. B₃, which is nested under B₁, has redeclaration of a, so it is out of scope of declaration of a in B₁. B₄ does not contain redeclaration of a, so it is in scope of B₁'s declaration.

Consider print statement of B₄, since B₄ has redeclared b=4, so it prints a b value as 4, and it does not contain a, its surrounding block B₂ also nor contains a's declaration. B₂'s surrounding Block B₁, contains a's declaration. So a's value is printed as 1.

<u>Declaration</u>	<u>Scope</u>
int a=1	B ₁ - B ₃
int b=1	B ₁ - B ₂
int b=2	B ₂ - B ₄ .
int a=3	B ₃ .
int b=4	B ₄ .

Explicit Access Control.

If p is an object of class with a member x , then $p.x$ refers to access of object p to class member x . And if a class C with member x , its subclass C' can also access x , except if it redeclared x .

The Object Oriented languages like C++, Java provides explicit access control through keywords public, private, protected. These keywords provides encapsulation by restricting access.

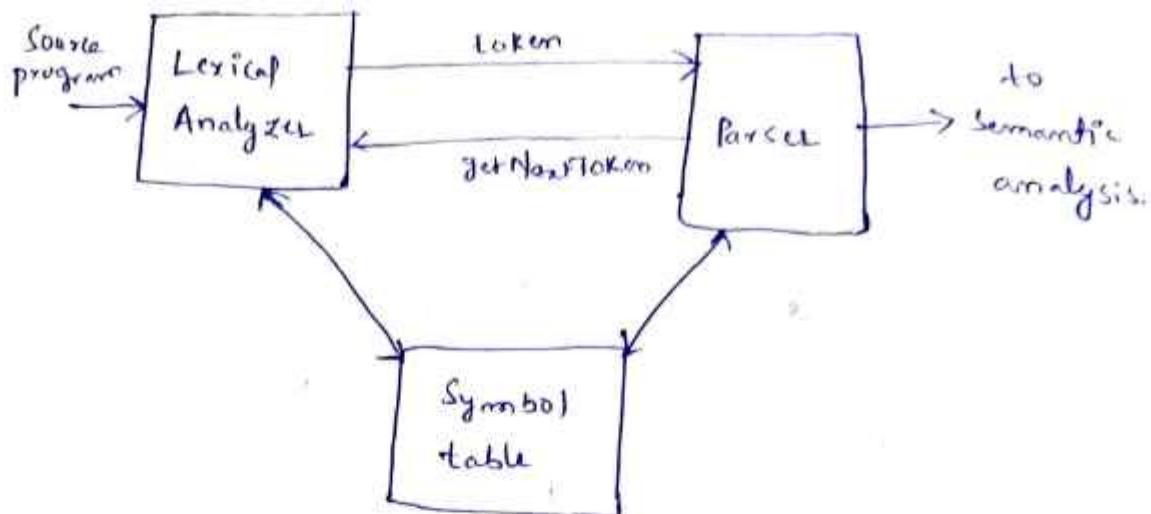
- private limits the access to the members declared within the class and "friend" classes.
- protected members can be accessed from subclasses.
- public members can be accessed from outside the class.

Chapter-2.

Lexical Analysis.Role of the Lexical Analyzer.

As it is the first phase compiler's main task is to read input characters of source program, group them into tokens sequentially from each lexeme. The stream of tokens is sent to build a parser for syntax analysis. Lexical analyzer also interacts with symbol table. When lexical analyzer discovers a lexeme as an identifier, it enters the lexeme into symbol table. These interactions are shown in following figure.

Commonly the interaction is started by parser. The call from parser causes the lexical analyzer "getNextToken" command to read characters from input until it identifies next lexeme, and produces token from it, which it returns to parser.



As the Lexical analyzer is part of compiler
it will also do some other tasks besides
identification of lexemes.

1. It will remove comments and whitespaces (blank, newLine, tab).
2. It will correlate error messages generated by compiler with source program. It will associate a line number in source program with error message.
3. If the source program uses macros, the expansion of macros will be performed by Lexical Analyzer.

Lexical Analyzer consists of two processes.

- a). Scanning → Simply reads input, does not concern with tokenization.
- b). Lexical Analysis → It is complex, generates tokens from output of scanner.

Lexical Analysis Versus Parsing.

The reasons for separating Lexical Analysis from parsing are as follows.

1. The design of Compiler will be simplified. For ex., if Comments and White Spaces are not removed in Lexical Analysis, The parser's task will become more complex.
2. Compiler's Efficiency is improved.
3. Compiler's portability is enhanced.

Tokens, Patterns, Lexemes.

→ Token is a pair consists of token name and optional attribute value. Token name represents kind of lexical unit. Token name works as input to parser.

→ pattern denotes type / form of source of a token. For ex., if token is a keyword, the pattern will be character that forms key word. For identify the pattern is more complex.

→ lexeme is a sequence of characters in source program that matches pattern for a token.

For ex., in below statement

```
printf("Total = %d\n", score);
```

printf, score matches the pattern for token id, and "Total=%d\n", is a lexeme matching literal.

Token	Informal Description	Sample lexeme
if	Character i,f	if.
else	Character e,l,s,e	else.
Comparision.	< or > or <= or >= or == or !=	<=, !=
id	letter followed by letters or digits	p1, score, d2.
literal	Anything that surrounded by " "s.	" hello".

Examples of tokens.

Attributes for tokens.

When more than one token lexeme matches a pattern, we need additional info about matched lexeme and pattern.

(15)

For ex., for pattern of token "number", both 0 and 1 matches. So it is important to know for code generator to know which lexeme is found in source program. Thus in many cases lexical analyzer sends token name and attribute value pair to parser that describes lexeme represented by token.

The most important example is the token id. When we need to associate several info with identifier like its lexeme, token, type, and where it can found in symbol table. Thus appropriate value for identifier is a pointer to symbol table entry for that identifier.

For ex., The id token names and associated attribute values for the Fortran statement:

$E = M * C^{**2}$.
are given as follows.

<id, pointer to symbol table entry for E>

<assign-op>

<id, pointer to symbol table entry for M>

<mult-op>

<id, pointer to symbol table entry for C>

<exp-op>

<number, integer value>.

Input Buffering -

The task of reading lexeme from source program can be speed up by using Buff'ers. It is a difficult task because we need to see next character, before we can see that we have right lexeme. For ex., we can not be sure that we've have seen end of identifier until we see a character that is not a letter or digit, and any thing which is not a part of lexeme id.

Buffer pairs.

Buffers are used to reduce the overhead involved in reading input characters. This scheme consists of two buffers. Each buffer size is same, say N , i.e., size of disk block e.g. 4096 bytes. Using a read command we can read N characters from one system call at a time, rather than using one system call per character. If fewer than N characters remain in file, special character "eof" marks the end of source file, which is different from characters in source program.

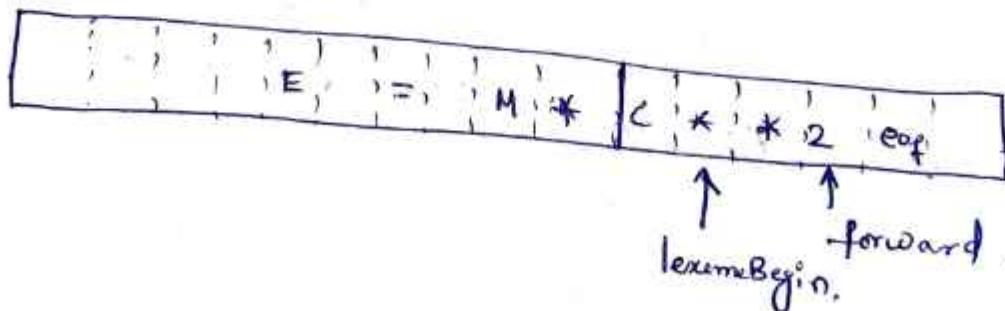
Two pointers to input are maintained

1. pointer "LexemeBegin" marks the beginning of current Lexeme, whose length we are attempting to determine.

2. Pointer "forward" scans until the pattern token is matched.

Once the lexeme is matched, forward is set to character at right end. Then after the lexeme is found, and passed to parser as token, "lexeme begin" is set to next character found immediately after the lexeme just found.

When "forward" pointer has reached end of the buffer, we must reload other input, and move forward to beginning of newly loaded buffer.

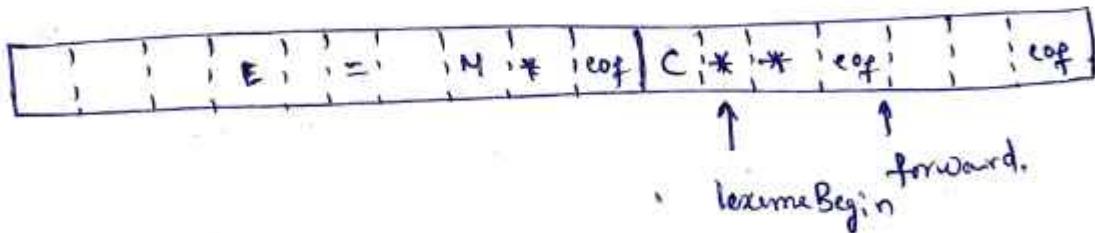


Sentinel.

For each character we read, we make two tests, one for end of the Buffer, and second one to determine, which character we read.

We can combine the buffer end test, and with test for new character read by using a special character known as Sentinel Character at the end. Sentinel is a special character which is not part of source program, and this is represented by eof.

eof is placed at the end of the buffer, as well as at the end of entire I/P. As shown in below figure.



Recognition of Tokens.

The input string is compared with patterns, and then tokens are generated for ex. Consider below grammar

$$\text{stmt} \rightarrow \text{if expr then stmt}$$

$$\begin{aligned} &| \text{if expr then stmt else stmt} \\ &| c \end{aligned}$$

$$\text{expr} \rightarrow \text{term} \cdot \text{relop term}$$

$$| \text{term}$$

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

$$| \text{number}$$

The above grammar is a simple fragment of branching and conditional statements. This syntax is similar to pascal.

$\Rightarrow \text{relop}$ is a comparison operator like " $=$ " for equals, and " $<>$ " for not equals in pascal.

Terminals of the grammar is which if, then, else, relop, id and number are named of tokens.

The patterns for these tokens are given by regular definitions as follows:

digit $\rightarrow [0-9]$

digits $\rightarrow \# \text{ digit}^*$

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

letter $\rightarrow [A-Za-z]$

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

if $\rightarrow \text{ if}$.

then $\rightarrow \text{ then}$.

else $\rightarrow \text{ else}$.

operator $\rightarrow < | > | <= | >= | = | <>$

The keywords if, then, else are the lexemes. They are not identifiers, but their lexemes matches patterns of identifiers.

Another job of lexical analyzer is to remove white spaces, by using token "ws" which is defined as follows.

ws $\rightarrow (\text{ blank} | \text{ tab} | \text{ newline})^*$

When any of them are recognized as "ws", we do not return it to parser, rather the next character after whitespace is returned. i.e; next token will be forwarded to parser.

Transition Diagrams

An intermediate step in lexical analysis is we first convert patterns into flowcharts called as "Transition Diagram".

Transition diagram have a collection of nodes or circles called as states. Each state represents a condition that could occur during the process of scanning the input for matching the lexeme to pattern. State is a summary between lexeme begin pointer and forward pointer.

Edges directs from one state of Transition diagram to other. Every edge is labeled by one or more symbols.

If we are in state 's', and next input symbol is 'a', we look for an edge out of s labeled by a. If such an edge is found, we advance forward pointer, and enter into the state where edge leads. If every state has only one edge The transition diagram is deterministic otherwise Non deterministic.

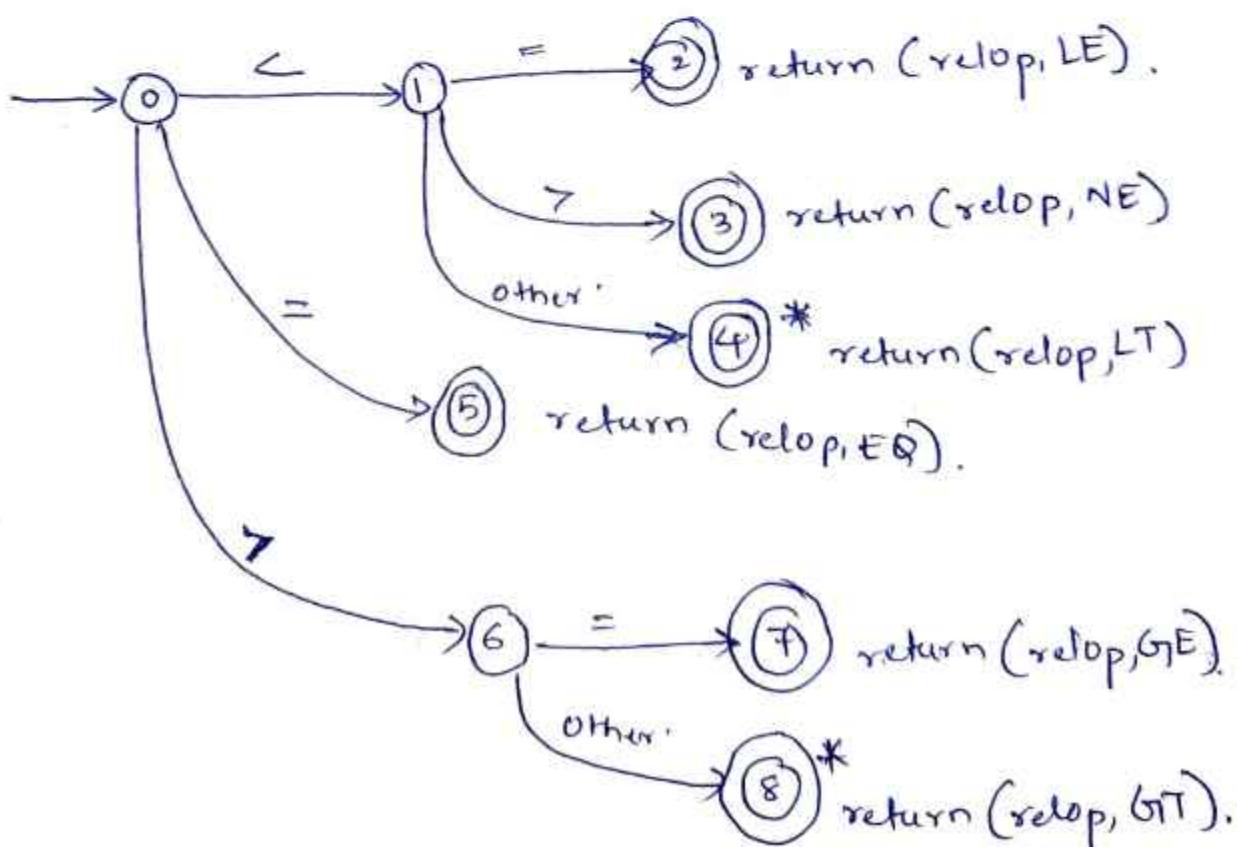
Some important conventions about Transition diagram are -

1. Some states are called as accepting or final states. These states indicates that some token has been accepted. The accepting / final state is represented by double circle. If any action to be taken after this state price returning value to the parser, a token and attribute will be attached to accepting state.
2. If it is necessary to move backward to one position, we should place forward pointer to * near accepting state.
3. One state is designated as "start" state, "initial state", which is indicated by label start entering from nowhere.

Ex:- Following diagram shows a transition figure that recognizer lexemes matching the token "relOp" (relational operators).

→ Take begin at state 0.
 → If input symbol is <, then among the lexemes that matches the pattern for relOp we look for <, <>, <=.

- we therefore goto state 1, and looks for next character.
- If next character is =, we recognize \leq and enter into state 2. and return the token relOp with attribute \leq .
- If in state 1, the next character is $>$, Then we have become \neq , and enter into state 3, and returns not-equals operators token.
- On Any other character, The lexeme is \neq , and we enter into state 4. and returns that info. State 4 has * to indicate that we must retract the input one position.



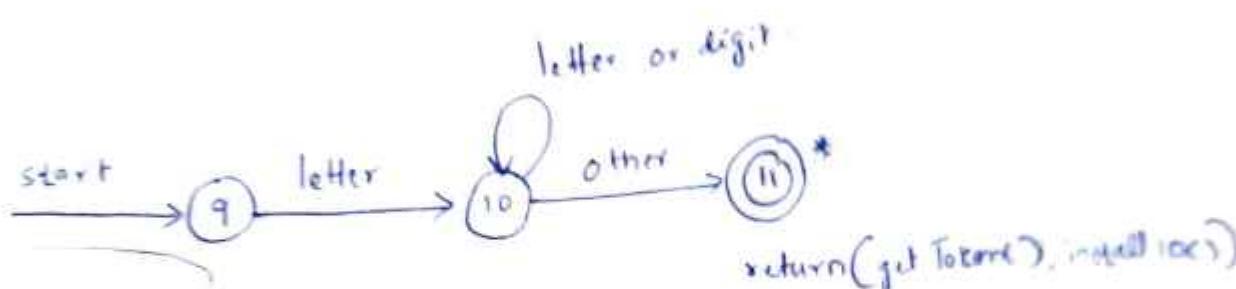
Recognition of Reserved words and Identifiers

Recognizing keywords and identifiers

Some what different from recognizing operators

Keywords like if and else are reserved.

They are not identifiers. Following Transition diagram identifies keywords if, then, else etc.



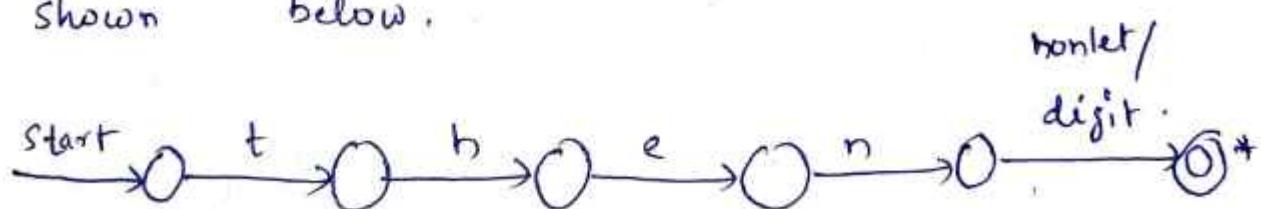
There are two ways, to recognize reserved words, that looks like identifiers

tell the reserved words in symbol table initially, A field of the symbol tells which token they represent. They are not identifiers.

When we find an identifier, installID(), places the identifier in symbol table. If it is not there already, and returns a pointer to the symbol table entry found.

- Any identifier not in the symbol table during lexical Analyzer can not be a Reserved Word., so its token is id.
- getToken(), examines the symbol table entry for lexeme found, and returns the token name either id or one of the keyword tokens that initially installed in the table.

2. Draw a separate transition diagram for each keyword, an one for keyword "then" is shown below.



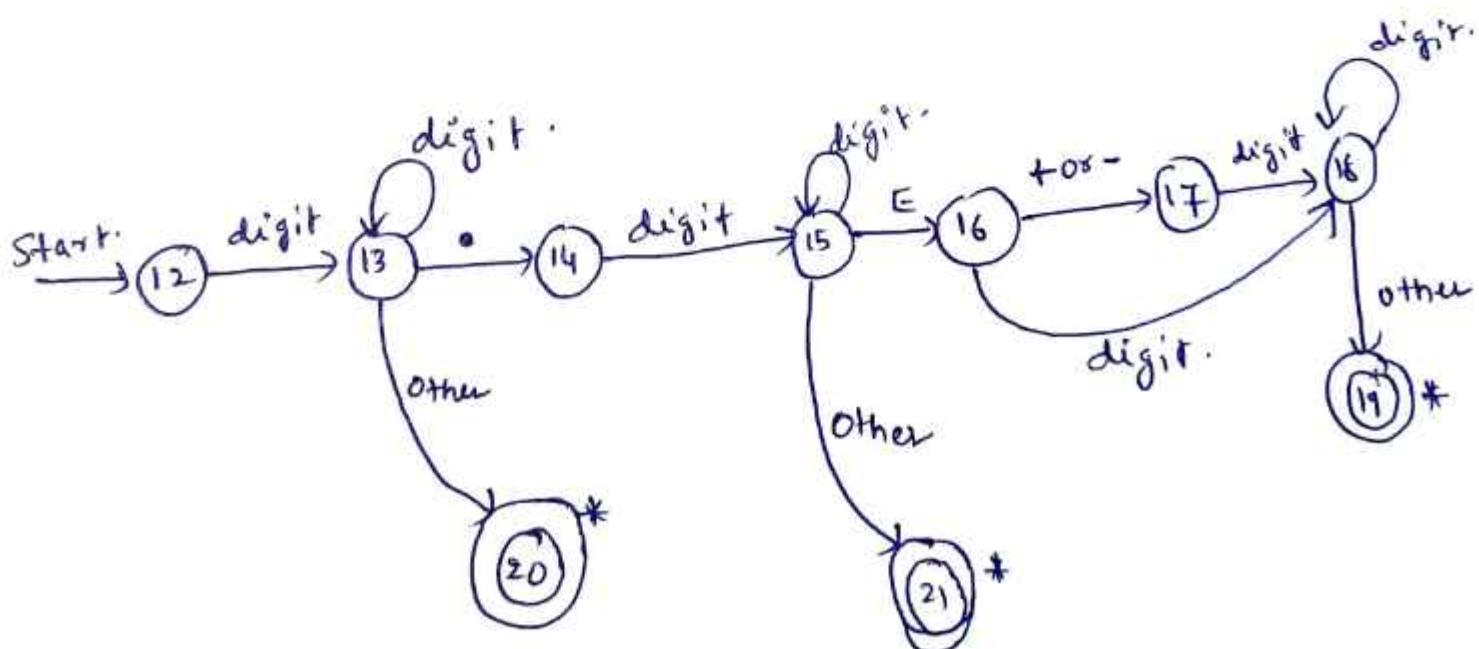
→ Transition diagram consists of states representing successive letters of keywords seen. followed by test for "non-letter-or-digit".

(21)

i.e., Any character that can not be a continuation of identifier. It is necessary to check the end of identifier, or else we return "then", where the correct token was "id" for lexeme "thenent value". If this approach is used, we should prioritize tokens so that reserved-word tokens are recognized in preference to id, when lexeme matches both patterns.

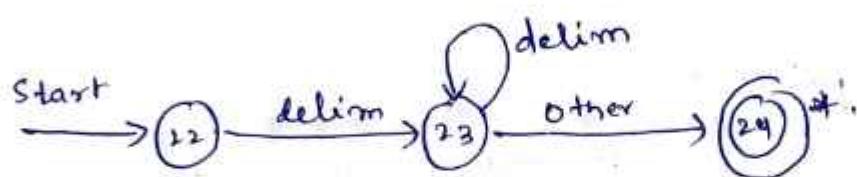
→ This approach is not used in our example, so the states in Transition diagram are numbered.

Transition Diagram for Unsigned Number.



A Transition Diagram for White Space.

In following diagram we look for one or more White Spaces, represented by delim. These delim characters could be blank, tab, newline and characters which are nor considered to be part of any token.



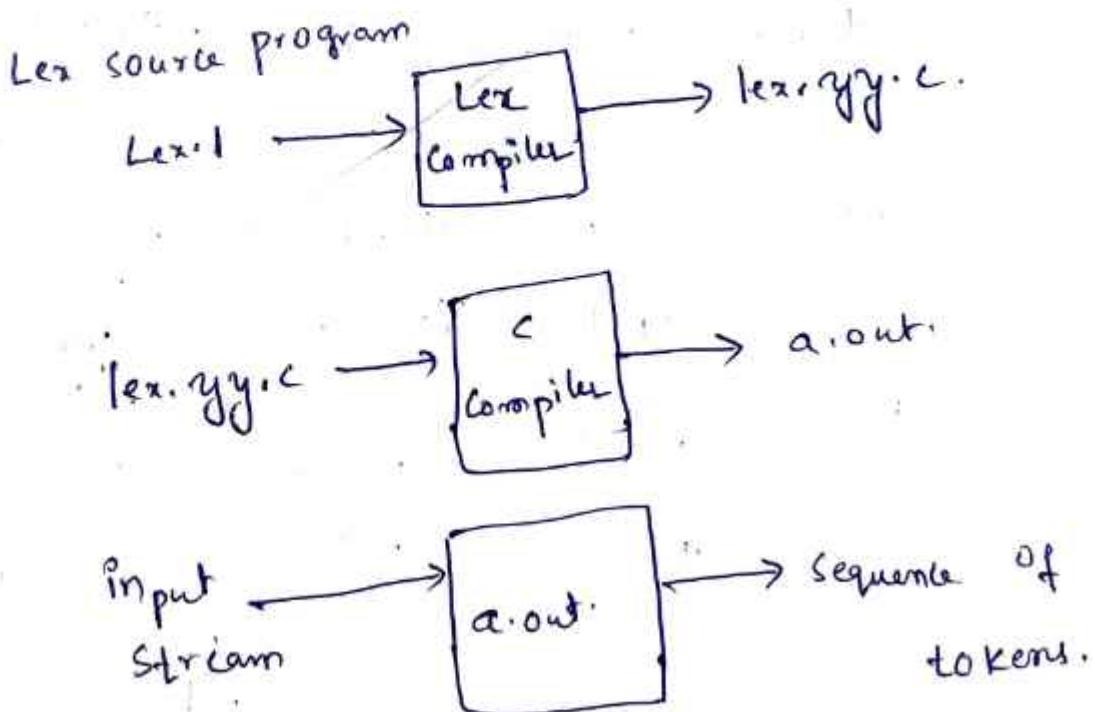
The Lexical - Analyzer Generator

Lex or Flex, allows to specify patterns for lexical analyzer, which describes patterns for tokens.

The input notation for Lex tool is known as the Lex language, and the tool is known as Lex Compiler. Lex Compiler transforms input patterns into transition diagram and generates code in file called as lex.yy.c.

Use of Lex

Following figure shows how Lex is used.



- An input file lex.l, which is written in Lex language describes lexical Analyzer to be generated.
- Lex compiler transforms Lex.l to a C program, which is named as lex.yy.c.
- This C file is compiled by C compiler, to a file a.out.
- This a.out file takes the input Stream of source program and generates sequence of tokens.
- a.out is a function in C, that returns an integer, which is code for token names, The Attribute value, i.e; either a numeric value, or pointer to symbol table entry is placed in a global variable "yyval." which is shared between Lexical Analyzer and parser. So a.out returns attribute Token name and attribute Value.

Structure of Lex Programs.

Lex program has following pattern

declaration

%%

Transition rules.

%%

Auxiliary Functions.

- Declaration section consists of declaration of variables and constants.
- The Transition rules will have the form
 - pattern of action}
- pattern uses definitions of declarations in declaration section. The Actions fragments of code written in C.
- The Third section holds additional functions used in Actions.

Lexical Analyzer created by Len, when called by parser, reads input string one character at a time, until longest prefix matches any one of the pattern P_i , it then executes Action A_i . A_i returns a single Value Token Name to parser, and a shared Variable yyval passes additional info about lexeme found, if needed.

If P_i describes white spaces, lexical Analyzer proceeds to find additional lexemes which matches to a pattern, and corresponding Action is returned to parser.

Example: Following is a lex program that recognizes tokens of below figure and returns Tokens Found.

Lexeme	TokenName	AttributeValue
Any ws	—	—
if	if	—
Then	Then	—
else	else	—
Any id	id	pointer to table entry
Any Number	Number	pointer to table entry.

$<$	relop	LT
\leq	relop	LE
$=$	relop	EQ
$>$	relop	GT
\geq	relop	GE

Finite Automata.

Lex is a tool, which converts input program to lexical Analyzer. The heart of transition is known as Finite Automata.

Finite automata are graphs, like Transition diagrams, with following differences.

1. Finite Automata are recognizers. They say "yes" or "no" about each possible string.
2. Finite Automata is of two types.

i. Non Deterministic Finite Automata (NFA).

A symbol can label many edges out of same state, and ϵ is an empty state. String, is a possible label.

ii. Deterministic Finite Automata (DFA)

Each state for each symbol as input, has exactly one edge with symbol leaving that state.

Both DFA, and NFA can recognize some languages known as regular languages.

Non Deterministic finite Automata (NFA)

NFA consists of

1. A finite set of states S .
2. A set of input symbols Σ , known as input alphabet., ϵ is an empty input string, which is never a member of Σ .
3. A Transition function δ ; which gives for each state, and symbol in Σ like a set of next states.
4. A start state s_0 , from S . (or initial state).
5. A set of states F , known as accepting states (or final states).

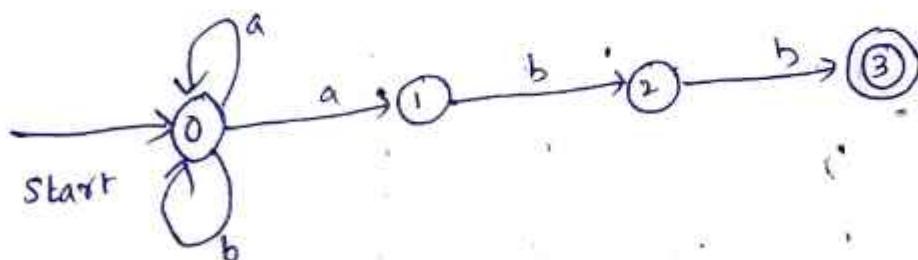
We can represent NFA or DFA by using transition graph, where nodes are states, labeled edges represent transition function.

This graph is same like Transition

diagram but with following differences.

- The same symbol, can label edges from one state to several different states.
- An edge may represent by ϵ , the empty string, in addition to symbols of input string.

Example:- An NFA, recognizing the language of regular expression $(a|b)^*abb$. can be given as follows. This example describes all strings of a's and b's ends with abb



This NFA, accepts all strings ends with abb.

Transition Tables.

NFA can also be represented by Transition table, whose rows corresponds to states, and columns corresponds to input symbol and ϵ . The entry in the table for a given state and input is the target state. If the transition function has no info about state input pair, we put \emptyset in the table.

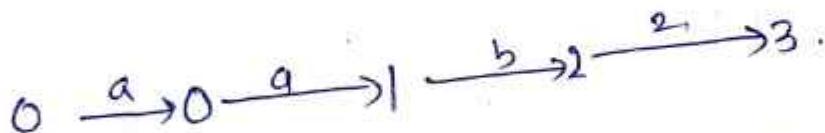
The Transition table for previous NFA can be given as follows.

State	a	b	c
0	$\{q_0, q_1\}$	$\{q_0\}$	\emptyset
1	\emptyset	$\{q_2\}$	\emptyset
2	\emptyset	$\{q_3\}$	q_1
3	\emptyset	\emptyset	\emptyset

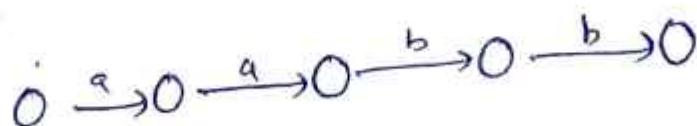
Accepting input strings.

→ NFA accepts input string x , if only there is some path in the transition graph from the start state to one of the accepting states. Labels with ϵ along the path are ignored.

Ex:- The string aabb, of previous NFA is accepted, The path is from state 0 to 3 demonstrating the fact as follows.

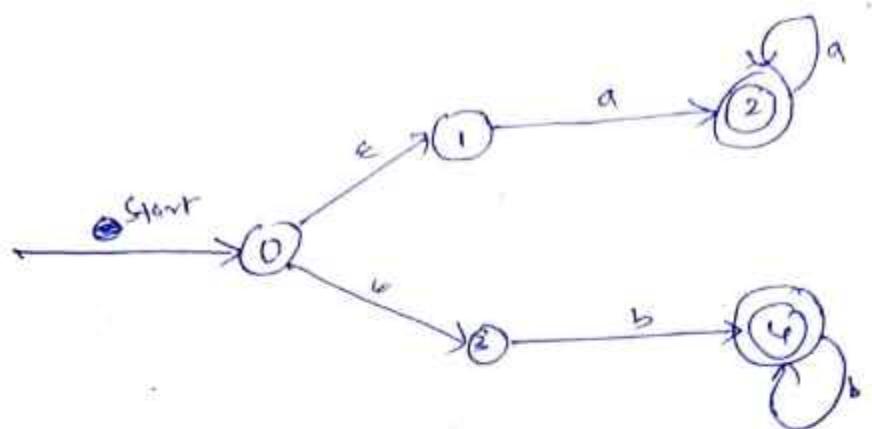


Several paths for some string may lead to different states for all



This path leads to state 0, which is not final state, so string on this path is not accepted.

Ex:- The following NFA accepts $L(a^*b^*)$.



The string "aaa" is accepted because of path

$$0 \xrightarrow{\epsilon} 1 \xrightarrow{a} 2 \xrightarrow{a} 3 \xrightarrow{a} 4$$

Deterministic finite Automata (DFA)

DFA is special case of NFA, where

1. There are no moves on input ϵ , and
2. for each state s , and input symbol a ,
there is exactly one edge labeled by $'a'$ out of s .

In Transition Table of DFA, each entry is single state, so states are represented by $\{s\}$ without using curly braces, which are used to represent sets.

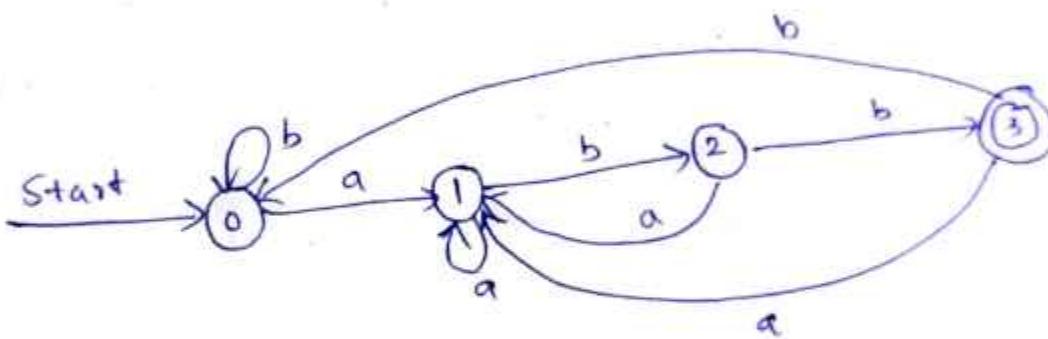
(14)

Every regular expression and DFA can be converted to DFA accepting same language.

DFA is used to build or simulate lexical Analyzers.

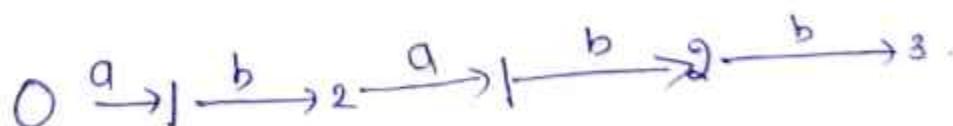
Ex:- Following figure represents DFA accepting

language $(a|b)^*abb$, same as previous NFA



Given string 'a babb' is accepted along the

path



Regular Expression to Automata.

Regular Expressions is a choice notation to describe lexical Analyzer and other pattern processing software.

Implementation of SW requires simulation of a DFA. Because NFA has choice of input symbol or in ϵ , NFA's simulation is less straightforward than DFA. Thus it is often important to convert NFA to DFA that accepts same language.

Conversion of an NFA to DFA.

→ Afa conversion Each state of DFA, corresponds to a set of ^{NFA} states. This technique is known as subset construction. This can be given by using following Algorithm.

Algorithm: Subset construction of a DFA from NFA.

Input : An NFA N .

Output : A DFA D , accepting same Language N .

Method :- This algorithm constructs a Transition table D_{tran} for DFA.

→ Each state of DFA, also contains set of states of NFA, which has moved on a given input string. First we need to deal with null transitions of NFA.

Following figure shows different functions that describes Computations on states of NFA. s denotes single state of NFA and $T \subseteq S$ set of states of NFA.

Operation

ϵ -closure(s)

Description

→ Set of NFA states reachable on ϵ -transitions alone.

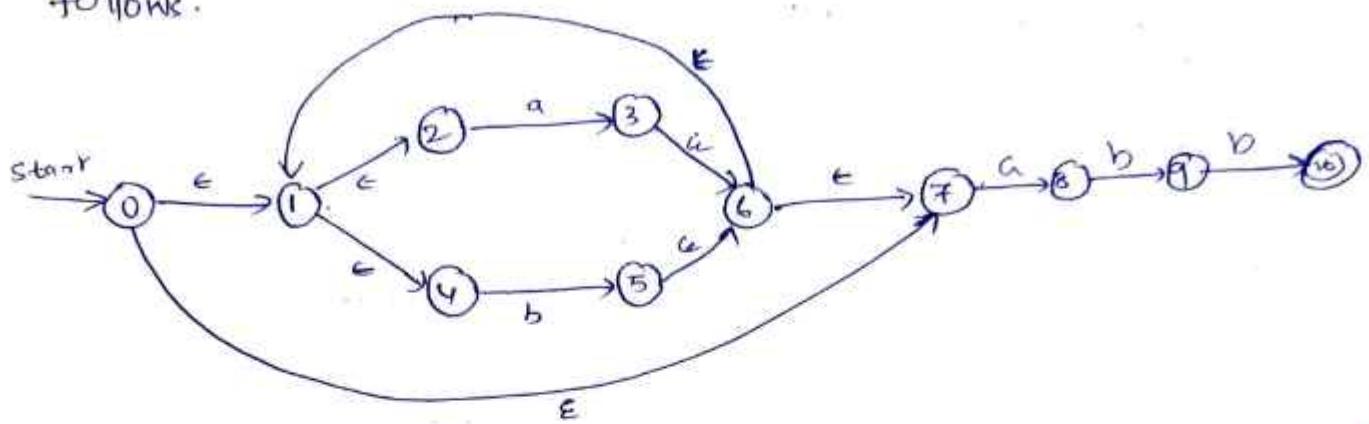
ϵ -closure(T)

→ Set of NFA states, reachable from set of states in set T on ϵ Transitions.

move(T, a)

→ Set of states $\overset{from}{\in} T$ on a ,

→ Following figure shows NFA, accepting $(a/b)^*abb$. It can be converted to DFA as follows.



→ ϵ -closure(0) is $\{0, 1, 2, 4, 7\}$, These are the states which are reachable from 0, on ϵ -transition. State 0, can be reached from itself using ϵ -transition. These set of states are equivalent to a set A, in DFA.

→ The input alphabet is $\{a, b\}$, we calculate/find

ϵ -closure (move(A, a)) i.e; $D_{trans}[A, a]$, and

ϵ -closure (move(A, b)) i.e; $D_{trans}[A, b]$.

→ In state A, only 2 and 7, have Transition on a to 3 and 8 respectively. Thus

$$Mov(A, a) = \{3, 8\}, \text{ and } \epsilon\text{-closure}(\{3, 8\}) =$$

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

So,

$$D_{trans} [A, a] = \epsilon\text{-closure}(\text{move}(A, a)) = \epsilon\text{-closure}(\{3, 8\}) \\ = \{1, 2, 3, 4, 6, 7, 8\},$$

Call this state as B.

$$D_{trans} [A, b] = B.$$

\Rightarrow Now $D_{trans} [A, b]$, among states of A, only 4 has transition on b to 5. So

$$D_{trans} [A, b] = \epsilon\text{-closure}(\{5\}) = \{1, 2, 4, 5, 6, 7\}.$$

Call this state as C.

$$D_{trans} [A, b] = C.$$

$\Rightarrow D_{trans} [B, a] = \epsilon\text{-closure}(\{3, 8\})$, which is similar to B itself.

$$\text{so } D_{trans} [B, a] = B$$

$$\Rightarrow D_{trans} [B, b] = \epsilon\text{-closure}(\{5, 9\}) = \{1, 2, 4, 5, 6, 7, 9\}$$

$$\text{let } D_{trans} [B, b] = D.$$

$$D_{trans} [C, a] = \leftarrow\text{-closure}(\text{move}(C, a)) = \leftarrow\text{-closure}\{3, 8\}$$

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

$$D_{trans} [C, a] = B.$$

$$D_{trans} [C, b] = \leftarrow\text{-closure}(\text{move}(C, b)) = \leftarrow\text{-closure}\{5, 6\}$$

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

$$D_{trans} [D, a] = \leftarrow\text{-closure}(\text{move}(D, a)) = \leftarrow\text{-closure}\{3, 8\}$$

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

$$= B.$$

$$D_{trans} [D, b] = \leftarrow\text{-closure}(\text{move}(D, b)) = \leftarrow\text{-closure}\{5, 10\}$$

$$= \{1, 2, 4, 5, 6, 7, 10\} = E.$$

$$D_{trans} [E, a] = \leftarrow\text{-closure}(\text{move}(E, a)) = \leftarrow\text{-closure}\{3, 8\}$$

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

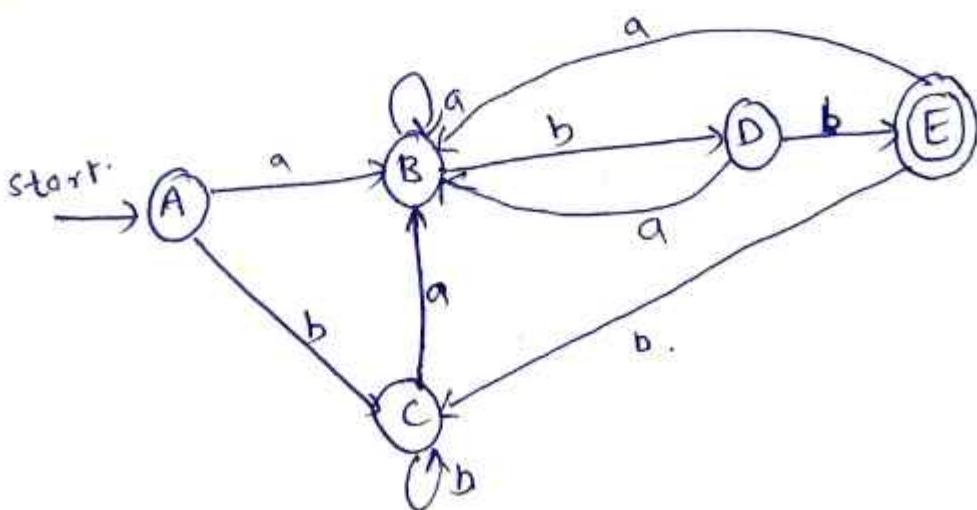
$$D_{trans} [E, b] = \leftarrow\text{-closure}(\text{move}(E, b)) = \leftarrow\text{-closure}\{5\}$$

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

State A is the start state, and State E, which contains state 10 of NFA, is only accepting state. The final DFA Trans can be given as.

NFA State.	DFA State	a	b.
$\{0, 1, 2, 4, 7\}$	A	B	C.
$\{1, 2, 3, 4, 6, 7, 8\}$	B	B	D.
$\{1, 2, 4, 5, 6, 7\}$	C	B	C
$\{1, 2, 4, 5, 6, 7, 9\}$	D	B	E
$\{1, 2, 4, 5, 6, 7, 10\}$	E	B	C

The DFA for above NFA, using subset construction is given as follows.



Design of a Lexical-Analyzer Generator.

Following figure shows the architecture of lexical analyzer generated by Lex. The program uses either finite automaton which can be either deterministic or non deterministic.

Following is a lex program which is turned into a Transition table and actions which will be used by Finite Automata

Simulators.

The components include

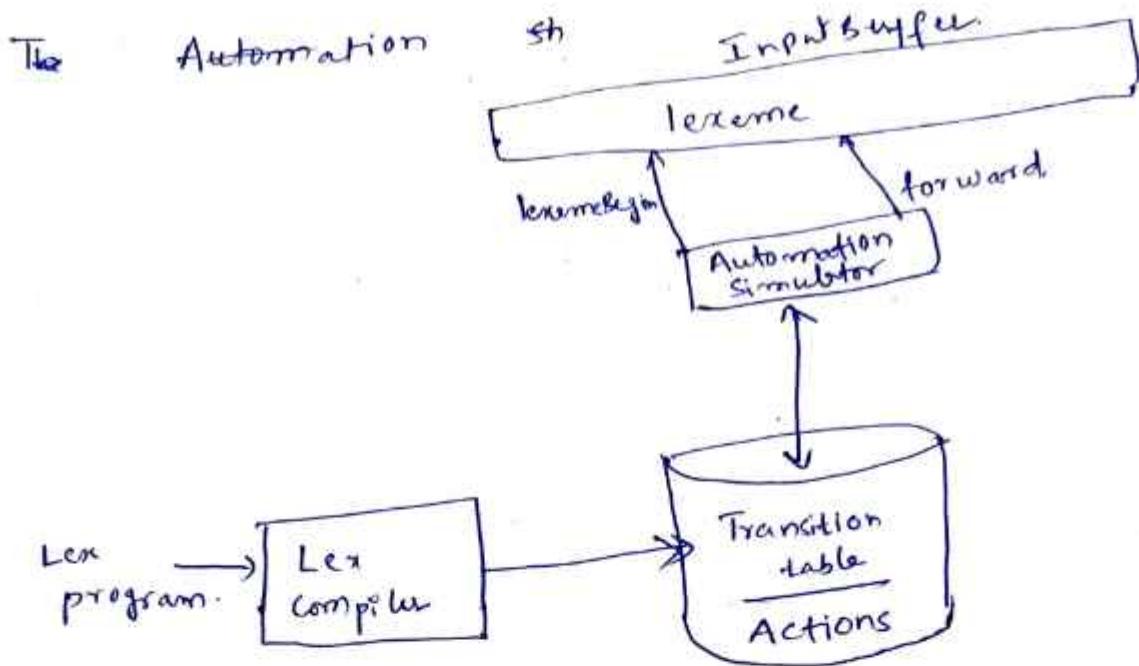
1. A Transition table for automation.

2. Functions passed directly through Lex

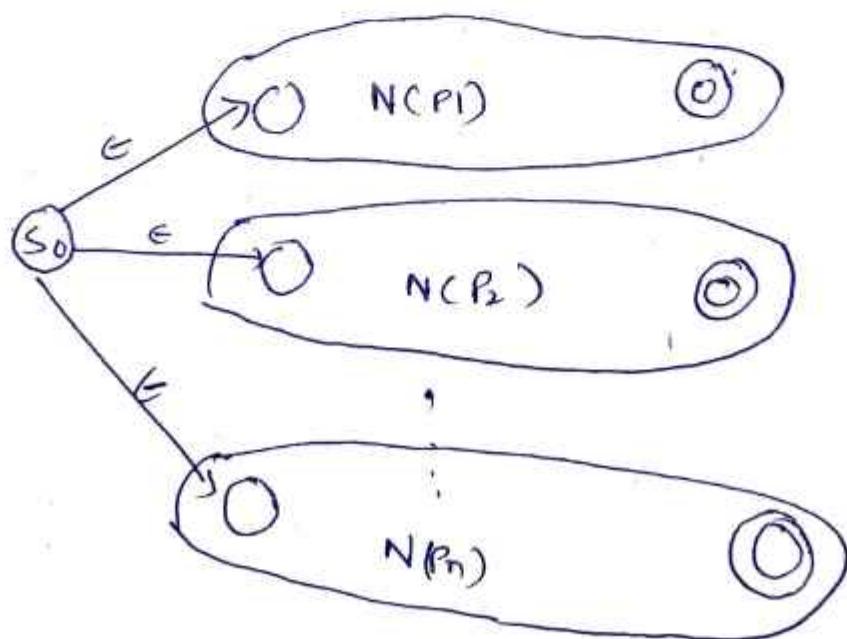
to the output.

3. The actions from the input program, which are to be invoked at the appropriate time by simulator.

To construct our automation, we begin by taking each regular-expression pattern in Lex program and converting it to NFA.



The Automation will recognize all lexemes matching any of the patterns in source program with ϵ -transitions. to from each start states of the NFA's N_i of pattern P_i as shown in below.



(32)

Example:

Consider following patterns.

a δ action A_1 for pattern $P_1 \}$

abb δ action A_2 for pattern $P_2 \}$

a^*b^* δ action A_3 for pattern $P_3 \}$

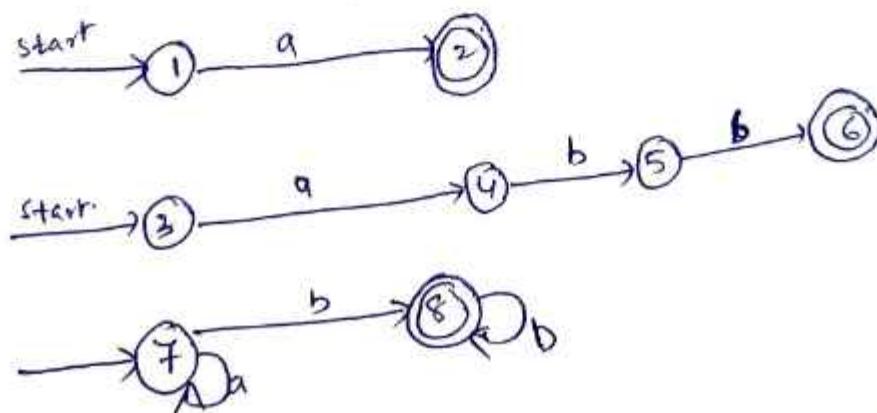
There are some conflicts in this example.

The pattern abb match P_2 and P_3 but we consider it as lexeme for P_1 , because it is listed first.

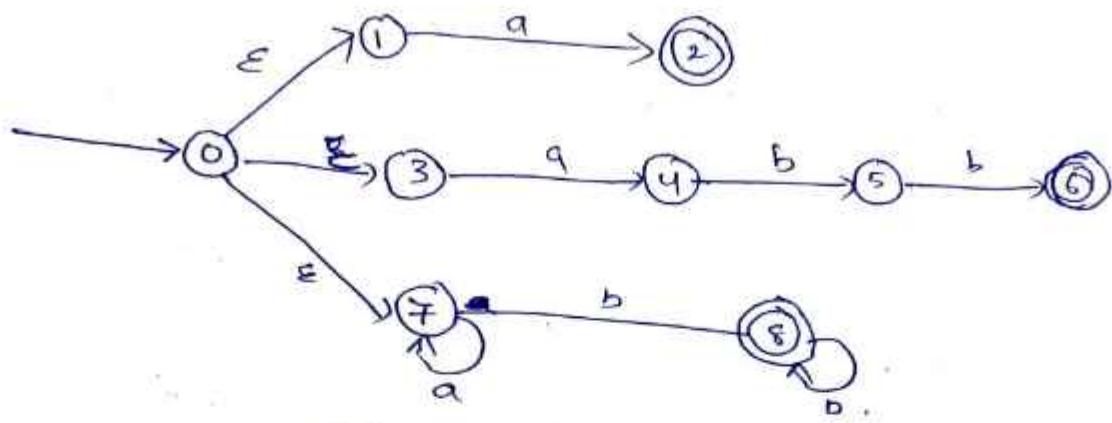
Following figure shows NFA's that recognize

three patterns. The last figure shows

combination of all above three NFA's.



NFA's for a, abb, and a^*b^* .



Combined NFA.

Pattern

Matching Based On NFA's.

Lexical Analyzer which stimulates NFA, reads the input from the beginning referred to as "Lexeme Begin", As it moves forward pointer, it calculates the set of states it is in each point.

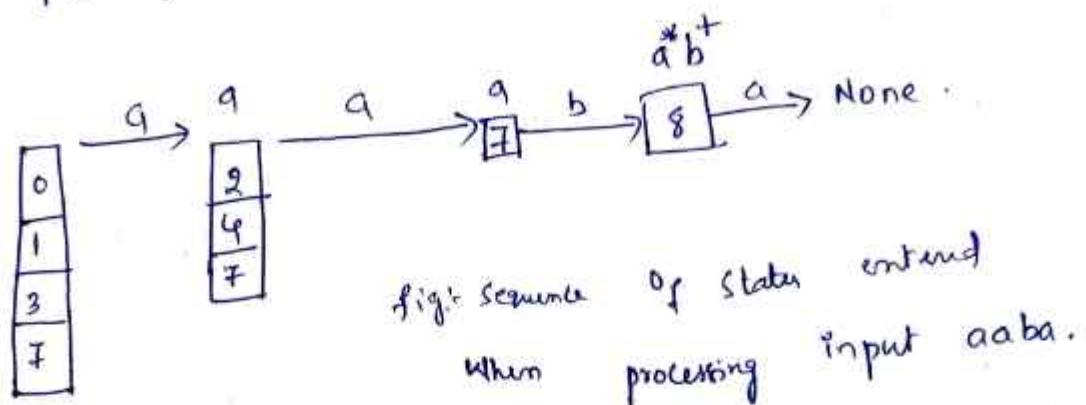
(32)

When NFA reaches a point on input, where there are no next states, the next set of states will be empty. Then by seeing accepting states in transition we decide the longest prefix that is a lexeme matching the pattern.

We look backward in sequence of states until we find one or more accepting states. If there are more than one accepting state we take the state with earliest p_i (pattern) in the list from Lexical Analyzer, move the forward pointer back to the end of lexeme, and performs the action A_i associated with pattern p_i .

Ex:- Consider patterns of given example, and input begins with "aab". The below figure shows set of states of NFA of previous figure, we enter state $\{0, 1, 3, 7\}$ on closure a which is from initial state. After reading fourth i/p symbol, T we are in empty set of states, since

There are no transitions out of state 8
on input a,



Thus we backup from that position, looking for accepting states. We notice in above figure after reading first input symbol a, we are in state 2, which indicates pattern "a" is matched, but after reading aab, we are in state 8, which indicates $a^* b^+$ has been matched. Thus prefix aab, is the longest prefix up to an accepting state, we therefore select "aab" as lexeme, and execute Action A₃, which returns the parser that token whose pattern $P_3 = a^* b^+$ has been found.

Optimization of DFA-Based Pattern Matchers.

There are three algorithms that can be used to implement and optimize pattern matchers from Regular Expressions.

1. First algorithm is used in LEX Compiler, it constructs DFA directly from a regular expression, without constructing intermediate NFA.
2. Second algorithm minimizes the states of DFA, by combining the states that have the same future behavior.
3. The third algorithm produces more compact representations of transition tables than the standard, two dimensional table.

Important states of NFA

An NFA state is known as important if it has non- ϵ out transitions. The subset construction uses only the important states in a set T , when it computes ϵ -closure ($\text{move}(T, a)$). The set of states reachable from T on a .

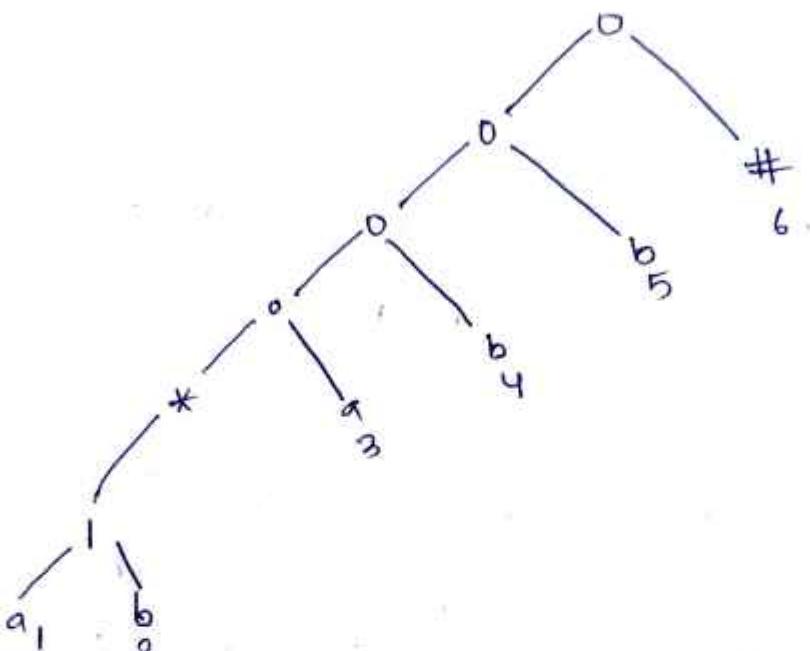
During Substitution construction, two sets of NFA states can be identified

- 1. Have the same important states ; and.
- 2. Either both have accepting states or neither does.

Each important state corresponds to a particular Operand in regular expression. The constructed NFA, has only one Accepting state but does not have out-transitions. By concatenating right end marker $\#$, to a regular expression(s) we give accepting state for r or on $\#$ transition of $\#$, making it important state of NFA for $(r)\#$.

The important states of NFA, corresponds to positions in the regular expression that holds symbols of alphabets. Regular expression can be represented by Syntax Tree, where leaves denote Operands, Operators. Interior nodes are cat-node (dot), or-node (\sqcap), or star-node ($*$) interior nodes denotes

Following figure shows Syntax Tree for regular expression $(a|b)^*abb\#\#$



Leaves in syntax tree are labeled by ϵ , or by alphabet symbol, if the leaf is not labeled by ϵ , we attach an integer, which refers to position of leaf / symbol in regular expression. One symbol can have many positions, for ex, a has position 1 and 3. positions corresponds to important states of NFA.

Functions Computed from the Syntax Tree.

To construct DFA, directly from regular expression, we construct Syntax Tree, and then compute four functions nullable, firstpos, lastpos and

followpos, which are defined as follows

1. nullable(n) \Rightarrow This function is true, in a Syntax Tree for node n , if the Subexpression represented by n has ϵ in its language. made null.
The Subexpression can be
2. Firstpos(n) \Rightarrow is set of positions in subtree rooted at n , that corresponds to first symbol.
3. lastSymbol(n) \Rightarrow set of positions in the subtree rooted at n that corresponds to the last symbol of at least one string.
4. Followpos(\emptyset) \Rightarrow for a position p , there is a set of positions q in entire Syntax Tree, such that there is a string $x = a_1 a_2 \dots a_n$ in $L((\cdot)^{\#})$, where a_i is matched to position p , and a_{i+1} to position q .

Example.. Consider the cat-node n , of previous
Syntax Tree That corresponds to expression
 $(a|b)^*a$.

Nullable(n) \Rightarrow Nullable is false, for this node
as it generates all the strings of a's and b's.
ending with a. On the other hand the
star node, below, it can be nullable, it
may generate along with strings of
a's and b's.

firstpos(n) $\Rightarrow \{1, 2, 3\}$, In strings generated like
'aq', The first position corresponds to 1,
string 'ba', The first position corresponds to 2,
And string with only 'a', first position
corresponds to 3.

lastpos(n) $\Rightarrow \{3\}$, For all strings, generated from
expression of node n, The last position is for
'a', and it is 3.

Computing nullable, firstpos, and lastpos.

We can compute nullable, firstpos and lastpos depending upon the height of tree. The following table represents the inductive rules for nullable and firstpos. The rules for lastpos are same as for firstpos.

Node (n)	Nullable(n)	firstpos(n)
A leaf node labeled ϵ	true	\emptyset
A leaf with position i	false	$\alpha_i \beta$.
An or-node $n = c_1 c_2$	nullable(c_1) or nullable(c_2)	firstpos(c_1) \cup firstpos(c_2)
A cat-node $n = c_1 c_2$	nullable(c_1) and nullable(c_2)	if (nullable(c_1)) firstpos(c_1) \cup firstpos(c_2) else firstpos(c_1). firstpos(c_1).
A Star Node $n = c_1 *$	true	

Example :- From all the nodes of previous Syntax Tree only Star-node is nullable, and non of the leaves are nullable because they corresponds to non-ε operands.

→ The 'or' node is not nullable, because neither of its children or null.

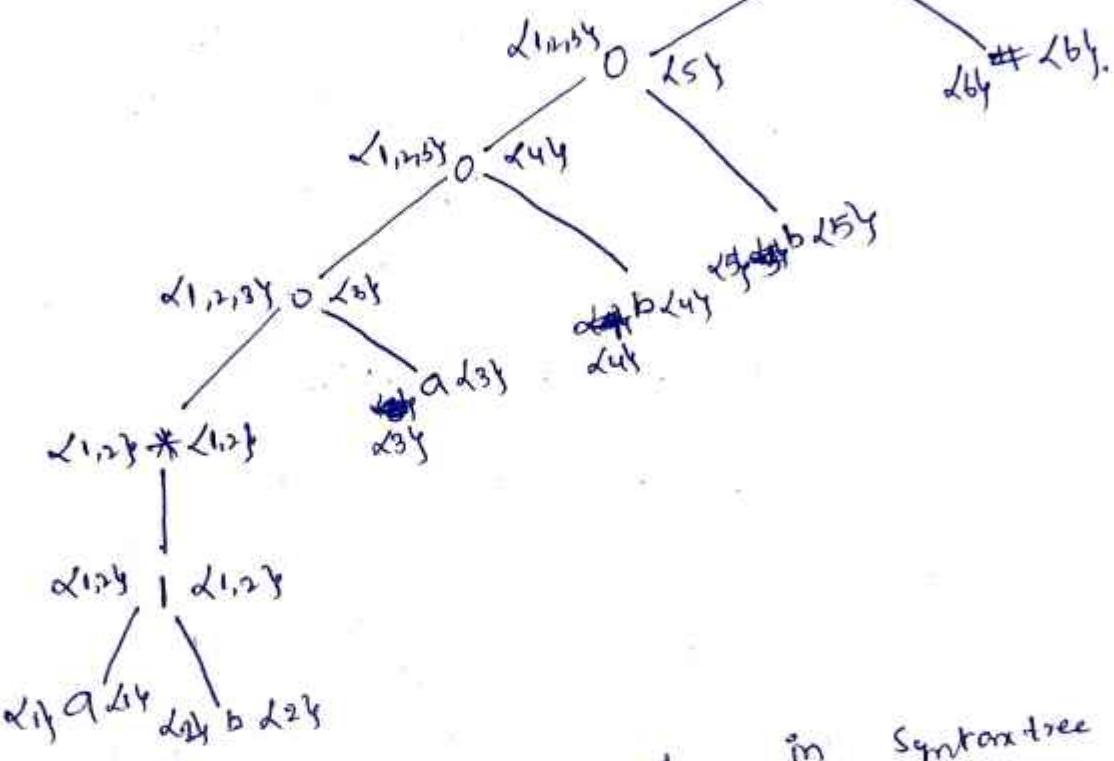
⇒ The Star . node is nullable because every star node is nullable.

⇒ The cat . node is ^{not} nullable if atleast one child is non nullable.

The Computation of firstpos and lastpos for each node is shown in below figure, with firstpos(n) to the left of node n , and lastpos(n) to its right. Each leaf has only itself for firstpos and lastpos, as the rule for non-ε leaves. For the or-node, we take the union of firstpos at the children and do the same for lastpos. The rule for the star-node says that we take value of firstpos or lastpos at any one child. of that node.

For ex, consider lowest cat-node, consider it as n . Its left child is nullable, and so its first pos is $\{1,2\}$, and first pos of right child is $\{3\}$'s union will give i.e; $\{1,2,3\}$ will give

n's firstpos. The rule for lastpos β is
 same like firstpos, with children interchanged. That is
 to compute lastpos(β), we must check whether
 its right child is nullable, which is not.
 Therefore lastpos(β) is last pos of $\alpha \in \{a|b\}^*$.



First pos and last pos for nodes in syntax tree
 for $(a|b)^*abb\#$

Followpos.

Follow pos can be applied to only. cat node(\bullet), and closure/start node ($*$). The rules for finding followpos are for a node n are.

if n is \bullet then

$$\text{Firstpos}(c_2) \rightarrow \text{lastpos}(c_1)$$

i.e; copy firstpos c_2 to last pos of c_1

if n is $*$ then

$$\text{Firstpos}(n) \rightarrow \text{lastpos}(n).$$

i.e; copy firstpos of n i.e; $*$ to last pos of $*$

* The final follow pos according to above two rules is shown in below table. First we find followpos for \bullet and then for $*$.

Position	n	$\text{followpos}(n)$
1		$\{1, 2, 3\}$
2		$\{1, 2, 3\}$
3		$\{4\}$
4		$\{5\}$
5		$\{6\}$
6		\emptyset

Converting the Regular Expression to DFA.

We now put all the steps of our example to construct DFA, for regular expression $r = (a/b)^*abb$. We observed that nullable is true only for * node.

The value of firstpos for the root is $\{1, 2, 3\}$, so this is the start node of D, consider this state as A. i.e; $A = \{1, 2, 3\}$, now we find $D_{trans}[A, a]$, and $D_{trans}[A, b]$. From the tree, among the positions of A, 1 and 3 corresponds to a, and 2 corresponds to b.

So

$$\begin{aligned}D_{trans}[A, a] &= \text{followpos}(1) \cup \text{followpos}(3) \\&= \{1, 2, 3\} \cup \{4\} \\&= \{1, 2, 3, 4\} \Rightarrow \text{let this state } B.\end{aligned}$$

$$\begin{aligned}D_{trans}[A, b] &= \text{followpos}(2) \\&= \{1, 2, 3\} \Rightarrow \text{it is same as } A.\end{aligned}$$

Similarly,

$$\begin{aligned}D_{trans}[B, a] &= \text{followpos}(1) \cup \text{followpos}(3) \\&= \{1, 2, 3, 4\} = B.\end{aligned}$$

$$\begin{aligned}D_{trans}[B, b] &= \text{followpos}(2) \cup \text{followpos}(4) \\&= \{1, 2, 3\} \cup \{5\} \\&= \{1, 2, 3, 5\} = C.\end{aligned}$$

(40)

$$D_{trans}[C, a] = \text{followpos}(1) \cup \text{followpos}(3)$$

$$= \{1, 2, 3, 4\} = B$$

$$D_{trans}[C, b] = \text{followpos}(2) \cup \text{followpos}(4) \cup \text{followpos}(5)$$

$$= \{1, 2, 3\} \cup \{5\} \cup \{6\}$$

$$= \{1, 2, 3, 5, 6\} = D.$$

\Rightarrow it is the final state as it has position of end marker.

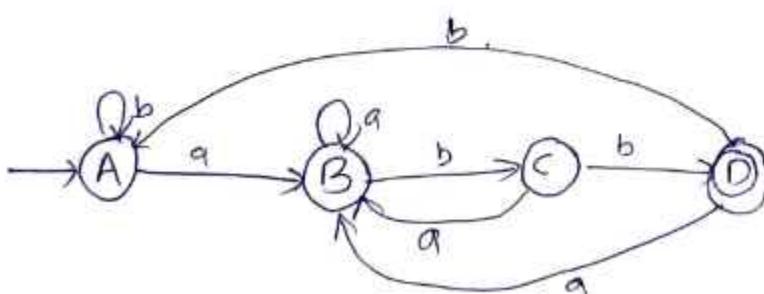
$$D_{trans}[D, a] = \text{followpos}(1) \cup \text{followpos}(3)$$

$$= \{1, 2, 3, 4\} = B.$$

$$D_{trans}[D, b] = \text{followpos}(2)$$

$$= \{1, 2, 3\} = A.$$

The DFA can be given as.

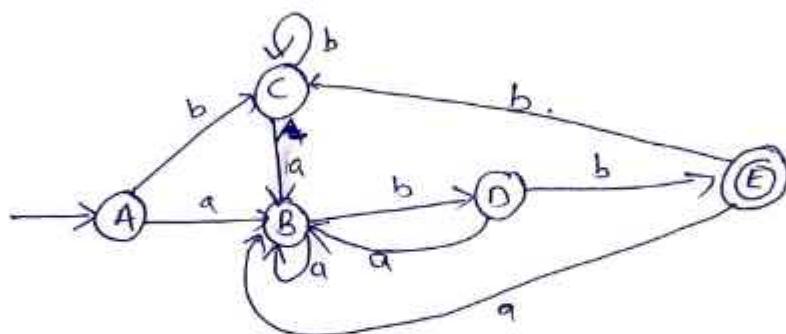


Minimizing The number of States of a DFA.

There can be many DFA's that recognize same language. If we implement a Lexical Analyzer using DFA, we prefer the DFA with few states. The Minimization can be done as follows.

Method:-

Consider below DFA



Step 1:- Separate states into two groups according to accepting states $\Rightarrow \{A, B, C, D\}, \{E\}$.

Step 2:- Consider both groups and inputs a, b.
E can not be split, because it has only one state $\{E\}$ remains in final group.

Step 3:- In group $\{A, B, C, D\}$, on input a, each state goes to B. On input b, A, B, C so we split members of $\{A, B, C, D\}$, D goes to $\{E\}$, $\{A, B, C, D\}$ as $\{A, B, C\}, \{D\}$.
The final group containing ~~$\{B\}$~~ $\{D\}, \{E\}$.

Step 4 :- Now in Group $\{A, B, C\}$, A and C all states goes to members of $\{A, B, C\}$. On input a , and on b , A and C goes to members of Group, B goes to D . So we split $\{A, C\}, \{B\}$.

Step 5 :- $\{A, C\}$, can't be splitted, because both goes to same state on input a and b .

The final Group consists of $\{A, C\}, \{B\}, \{D\}, \{E\}$

The minimum DFA, consists of four states. And we consider A as representative of $\{A, C\}$.

The initial state is A , and final state is E .

* The Transition Table ^{and DFA.} can be given as

follows.

state	$\frac{a}{B}$	$\frac{b}{A}$
A	B	A, C
B	B	D
D	B	E
E	B	A

```

graph LR
    Start(( )) --> A((A))
    A -- a --> B((B))
    A -- b --> C((C))
    B -- a --> B
    B -- b --> D((D))
    D -- a --> D
    D -- b --> E(((E)))
    C -- a --> A
    style Start fill:none,stroke:none
    style C fill:none,stroke:none
    style E fill:double
  
```

SYNTAX ANALYSIS = (ROLE OF SYNTACTIC ANALYZER) :-

- Syntax Analyzer is also called as the parser, the parser obtains the stream of tokens from the lexical analyzer.
- And, verifies the tokens syntactically, if the tokens are syntactically correct then it will generate a parse tree.

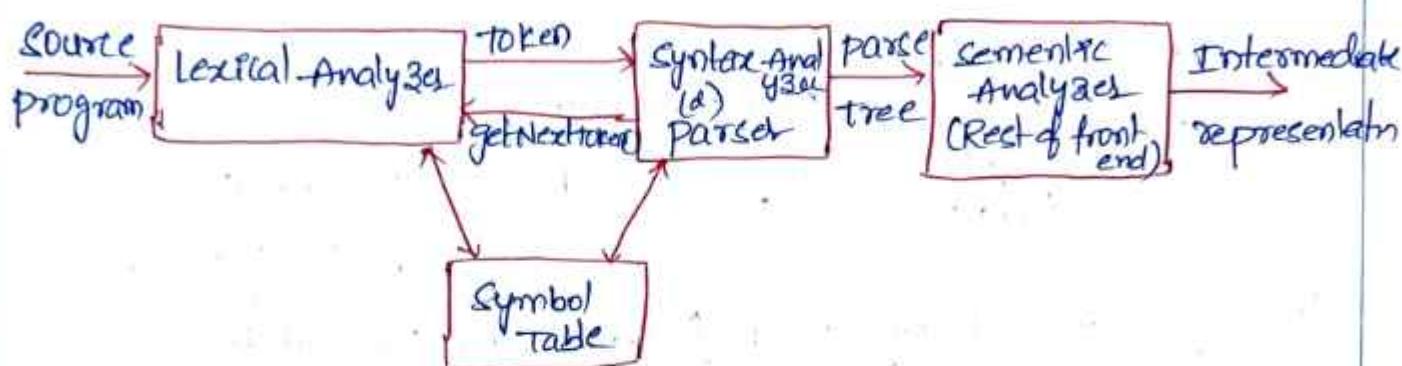
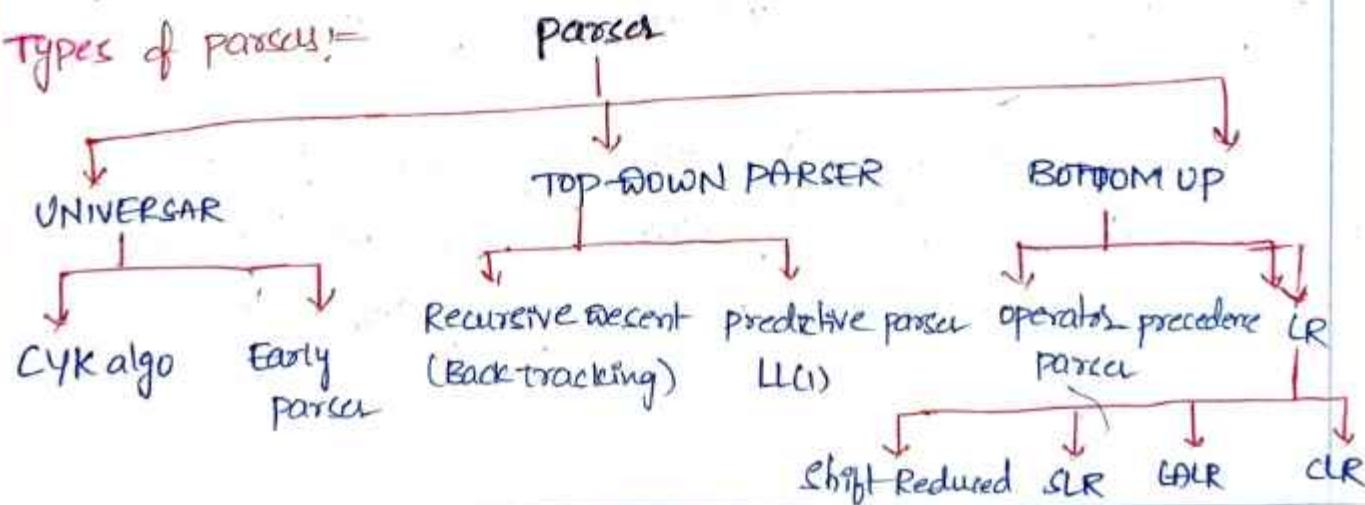


Fig:- position of parser in compiler model.

- The main function of parser is to check the syntactic structure of ip code, if the given ip code is out of syntactic errors then it will generate a parse tree.
- If ip code is consists of syntax errors, then the parser reports an errors to user, now, the user will correct it. It is responsibility of the user to correct those errors.
- Symbol Table communicate with both Lexical Analyzer & Syntax Analyzer for the token information.

Types of parsers:-



- In top-down parser construct the parse tree from top (root) to the bottom (leaves)
- Bottom-up parser construct the parse tree from leaves to root (starts from ϵ symbol and reduced to starting state)

Representative Grammars:

Associativity and precedence are captured in the following grammar

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

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

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

$E \rightarrow$ Expression

$T \rightarrow$ Terms, $F \rightarrow$ Factors that can be either parenthesized expression

These grammars belong to LR grammars that are suitable for bottom-up parsing. This grammar cannot be used for top-down parsing

The following non-left-recursive grammar will be used for top-down parsing

$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

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

Syntax Error Handling:

Programming errors can occur at many different levels.

Lexical Errors:

→ Include misspelling of identifiers, keywords, & operators

e.g.: ^{Identifier} ellipsesize instead of ellipsesize and missing quotes around text intended as a string. (e.g. "String");

Syntax Errors :- Misplaced semicolon(;), extra or missing braces {}, }.

Eg :- appearance of case statement without an enclosing switch.

Semantic Errors :- Eg :- $(a + (b *)) \rightarrow$ missing parenthesis.

→ type mismatch b/w operators and operands.

Eg :- $2 + a[i] \rightarrow$ type mismatch.

② The return of a value in Java method with result type void.

logical Errors :-

→ errors can be anything from incorrect reasoning on the part of programming.

Eg :- In "c" - the assignment operator instead of comparison operator.

if ($C = A$) X if ($C == A$) ✓

Errors Handles in parser :-

- should report the presence of errors clearly and accurately.
- should recover from each error quickly enough to detect the subsequent errors.
- should not slow down the processing of remaining program.

Error Recovery Strategies :-

(i) Panic Mode Recovery :-

Eg :- int a, b; }
printf(" ");

→ The parser discards I/P symbols at a time until one of a designated set of synchronizing tokens is found.

→ The synchronizing tokens are delimiters ; ; } }

→ It is simple to implement & does not goto & loop

(ii) phrase level Recovery :-

→ The parser performs local correction on remaining I/P when the error is discovered.

→ The parser replaces the prefix of the remaining I/P by some strings that allows the parser to carry on its execution

Eq: = replace a (,) comma by semicolon (;), delete an extra semicolon, insert a missing semicolon (;

(i) $\text{print}(" ")$, $\rightarrow \text{printf}(" ")$;

Disadvantage: = Error correction is difficult when actual errors occur before the point of detection

(ii) Form production.

\rightarrow common errors that can be encountered, we can augment the grammars for the language with productions that generate erroneous constructs. \rightarrow

\rightarrow use a new grammar for the parser.

(iv) Global correction:

\rightarrow The aim is to make some changes while converting incorrect i/p string to a valid string

- Given an incorrect i/p x find a parse tree for a related string w (^{& Grammar} using the given Grammar) such that no of changes (insertion/deletion) required to transform x to w is minimum
- Too costly to implement

(2) CONTEXT FREE GRAMMAR:

Context Free Grammars consists of 4-tuples

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

- $V \rightarrow$ Set of variables (or) Non-terminals
- $T \rightarrow$ Set of terminals.
- $P \rightarrow$ Set of Production Rules
- $S \rightarrow$ Start symbol.

(i) Non-terminal:

\rightarrow denotes set of strings.

- Uppercase letters early in alphabet A, B, C, ...
- $S \rightarrow$ start symbol
- E, T, F (Expression - E, Term - T, Factor - F)

(ii) Terminal:

\rightarrow Terminals are the basic symbols from which strings are formed

- TOKEN-name is also called as terminal
- The following are terminal symbols.
 - (a) lowercase letters in the alphabets such as a, b, c,
 - (b) operator symbols such as +, *, ...
 - (c) punctuations (" { ", " } ", " [", "] ", " , ", " . . . ") & digit 0, 1, ..., 9

(iii) Production:

consists of

- (a) A Non-terminal called head or left side of the production.
This production defines some of the symbols denoted by the head
- (b) body or right side \rightarrow consisting of zero or more terminals
and Non-terminals.
- \Rightarrow uppercase letters late in alphabet X, Y, Z represent Grammar
symbols either do terminals or Non-terminals.

$\rightarrow \alpha, \beta, \gamma \dots$ represent string of grammar symbols

\rightarrow A set of productions $A \rightarrow \alpha_1, A \rightarrow \alpha_2, \dots, A \rightarrow \alpha_K$ also written as
Product head $A \rightarrow \alpha_1 | \alpha_2 | \dots | \alpha_K$
Non-Terminal either terminal α_i or Nonterminal (NVT)*

Eg: Consider the grammar

$$E \rightarrow E+E \quad E \rightarrow E*E \quad E \rightarrow I \quad I \rightarrow a$$

$$CFG = (\{E, I\}, \{+, *, a\}, P, S)$$

\Rightarrow derive

Derivation: Derivation is process of applying a sequence of production rules in order to derive a string.

\rightarrow deriving an i/p string from the start symbol of Grammars in one or more steps by replacing the head of the production by its body of the production.

There are two types of derivation.

(i) Left most derivation (LMD)

(ii) Right most derivation (RMD)

Left most derivation: In each step

We have to expand left-most Non-terminal by one of its production body.

Eg: $E \rightarrow E+E | E*E | -E | (E) | id$

Derive a string $id + id * id$

LMD: $E \xrightarrow{\text{LMD}} E+E \quad (E \rightarrow E+E)$

$E \xrightarrow{\text{Lm}} id + E \quad (E \rightarrow id)$

$E \xrightarrow{\text{Lm}} id + E*E \quad (E \rightarrow E*E)$

$E \xrightarrow{\text{Lm}} id + id * E \quad (E \rightarrow id)$

$E \xrightarrow{\text{Lm}} id + id * id \quad (E \rightarrow id)$

Right most derivation:

\rightarrow Here, we have to expand the Right-most non-terminal by

RMD:

$E \xrightarrow{\text{Rm}} E+E \quad (E \rightarrow E+E)$

$E \xrightarrow{\text{Rm}} E+E*E \quad (E \rightarrow E*E)$

$E \xrightarrow{\text{Rm}} E+E*id \quad (E \rightarrow id)$

$E \xrightarrow{\text{Rm}} E+id*id \quad (E \rightarrow id)$

$E \xrightarrow{\text{Rm}} id + id * id \quad (E \rightarrow id)$

Parse tree:

→ A Graphical representation of derivation is called parse-tree.

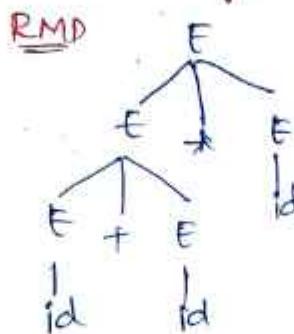
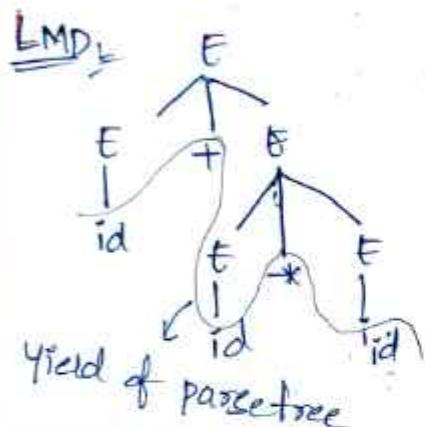
These types of nodes in parse tree

(i) Interior node → are Non terminals

(ii) children node → terminals

→ In parse tree the root node must be start symbol

construct parse tree for i/p string $w = id + id * id$



Ambiguity:

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

→ If a grammar produces more than one parse-tree, then that grammar is called as ambiguous.

(i) more than one LMD (different)

(ii) more than one RMD (different).

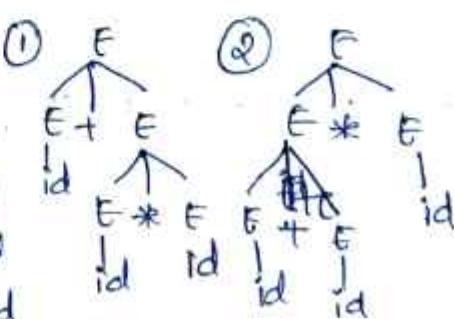
e.g. $E \rightarrow E+E \mid E * E \mid (E) \mid id$ string $w = id + id * id$

LMD₁:

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

RMD₁:

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

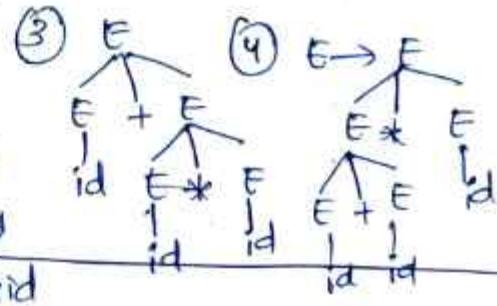


LMD₂:

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

RMD₂:

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



→ For the above string we got 2 LMD & RMD and 2 parse trees

So, the given grammar is ambiguous

→ Top down parser can't handle ambiguous grammar, so we need to convert this grammar into unambiguous grammar.

* While converting ambiguous to unambiguous:

→ The grammar should follow associativity and precedence rules

* , /, +, - → left association (when an ^{operand} operator has ^{operator} operand on both sides, the operand should associate with left-side operator)

1, *, +, - → precedence order
1 2 3 4

② WRITING A GRAMMAR:

→ Grammar is used to describe the syntax of programming language.

Lexical versus Syntactic Analysis:

- Separating syntactic structure into smaller and manageable components of lexical and non-lexical parts.
- Lexical rules are simple to understand than grammar.
- We need notations to describe the grammar.
- RE are used to Realtime Examples to make them easy to understand.
- Grammars are used for describing nested structure such as balanced parenthesis, matching begin-end, if-then-else.

Eliminating Ambiguity:

Properties of CFG

1. Ambiguous
2. Left Recursive
3. Left Factoring.

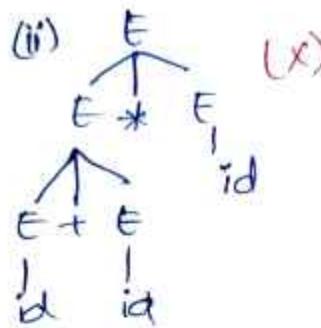
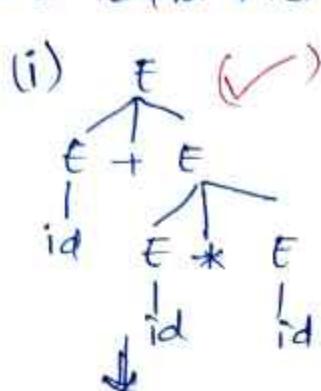
Elimination of ambiguity:-

(Consider CFG) :-

$$\text{Eg:- } E \rightarrow E+E \mid E * E \mid (E) \mid id$$

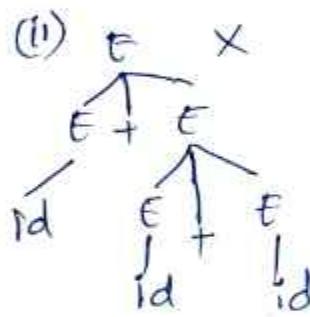
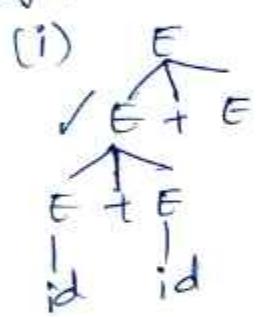
$$w = id + id * id$$

(There are 2 parse trees for string)



→ This is the valid parse tree, because the "*" operator is appeared at bottom of parse, so it is evaluated first.

$$\text{Eg2:- } w = id + id + id$$



→ two factors that are needed to be taken

(1) precedence

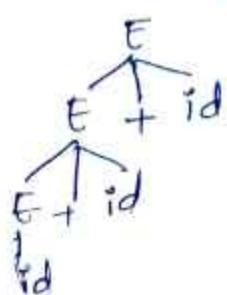
(2) left associativity → to ensure left associativity, we need to convert the grammar into left recursive RHS

→ left recursive $E \rightarrow E+E$ (the leftmost symbol is equal to the LHS)

* In left recursive, the parse tree can be grown on left side only

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

$$id + id + id \Rightarrow (id + id) + id \text{ (left associative)}$$



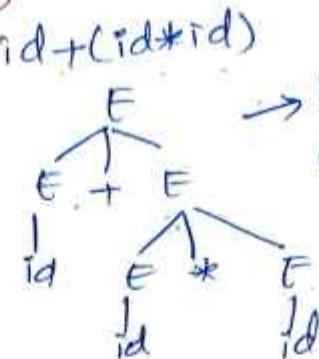
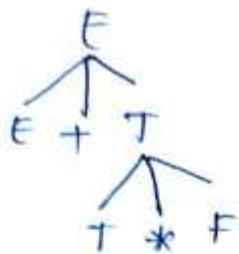
left associativity can be achieved

$$(id + id) + id$$

↓
This expression is evaluated first, so that it is valid parse tree.

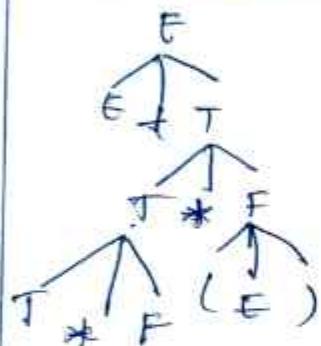
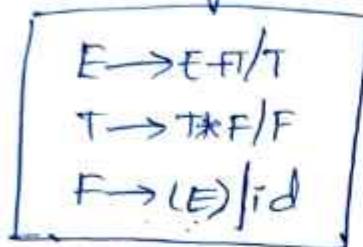
$$\begin{array}{l} E \rightarrow E + T \\ T \rightarrow T * F \end{array}$$

(precedence)

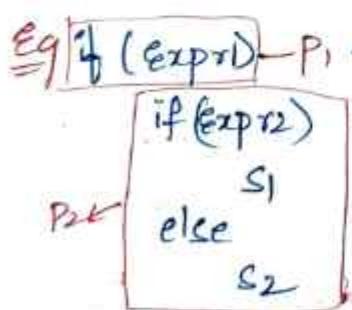


→ highest precedence operation is executing first.

The Final ~~gen~~ unambiguous Grammars



Dangling else: In a program if there are more than one if-
stmts then else part is matched with wrong if
stmt, that will lead to wrong results, that is called dangling
else.



-) if P_1 is true Then it will goto P_2 and
 - check the cond1 and condition, if it also true then it will execute S_1
 - if and condition is false it will execute else part i.e S_2

other stnt

(ii) if $P_i \rightarrow$ is false it has to execute the other statement

(ii) But here if P_1 - False then it is executing else part i.e statement
it is matching with wrong else.

- Consider the following dangling else grammars.

$\text{stmt} \rightarrow \text{if Expr then stmt}$

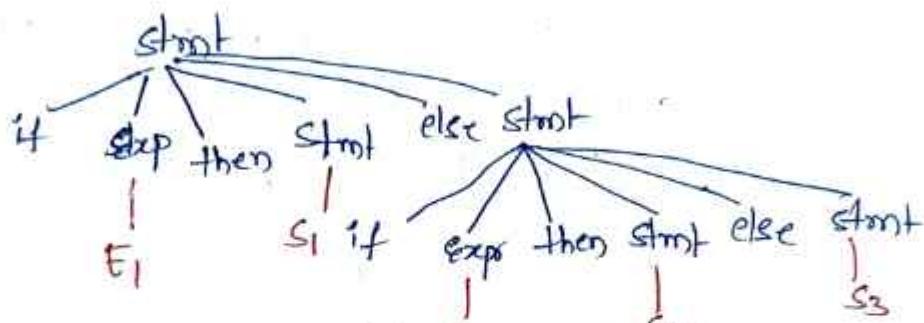
$\text{if Expr then stmt else stmt}$

| other

- The compound conditional stmt for the above grammars is

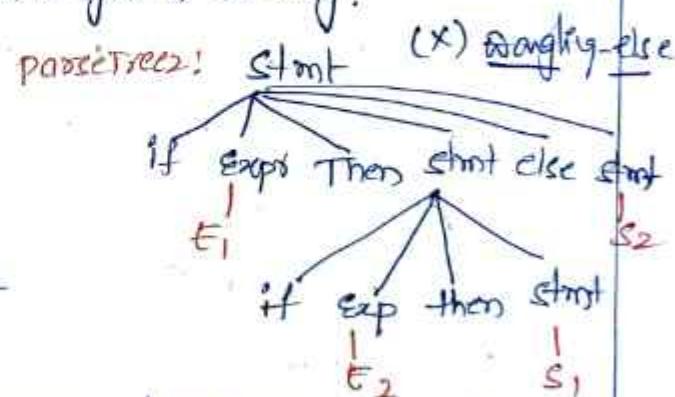
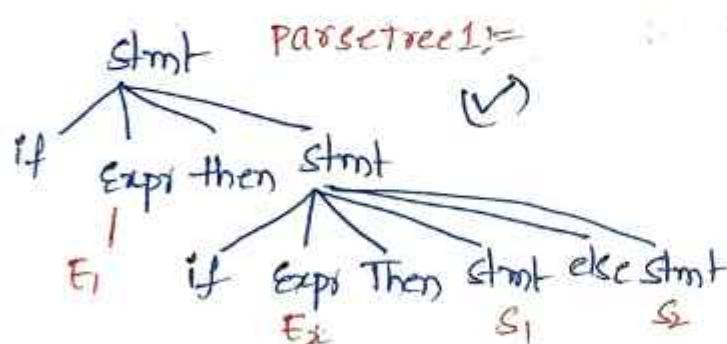
$\text{if } E_1 \text{ then } S_1 [\underline{\text{else if }} E_2 \text{ then } S_2 \text{ else } S_3]$

- The parse for this stmt



$\text{Ex i/p string} = \text{if } E_1 \text{ then } [\text{if } E_2 \text{ then } S_1 \text{ else } S_2]$

there exist two parse trees for the given string.



→ In these two parse trees, in programming language, the 1st parse tree is considered because, match the else statement with closest then. $\text{if } E_1 \text{ then } [\text{if } E_2 \text{ then } S_1 \text{ else } S_2] \rightarrow \text{matched stmt}$

→ So, the given Grammars is ambiguous Grammars.

The unambiguous grammar will be.

$\text{stmt} \rightarrow \text{matched stmt} \mid \text{open stmt}$

matched stmt \rightarrow perfect if
open stmt \rightarrow simple if

matched stmt \rightarrow if Expr then matched stmt else matched stmt / other

open stmt \rightarrow if Expr then stmt

| if Expr then matched stmt else open stmt.

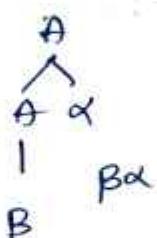
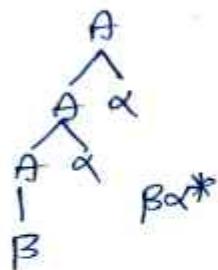
(ii) Elimination of left Recursion:

- A grammar is said to be left recursive if it has a non-terminal A such that there is a derivation $A^+ \Rightarrow A\alpha$ for some string α .
- Top-down parsing methods cannot handle this grammar, so we need to eliminate left recursion.

$A^+ \Rightarrow A\alpha / \beta$ (left recursive) $\alpha, \beta \in (VUT)^*$

[if LHS is equal to the leftmost Non-terminal of RHS]

The string derived from the above grammar.



string(w) = $\beta\alpha^*$
(any production should not contain * symbol)

production $A \rightarrow A\alpha / \beta$ could be replaced by the following non-left-recursive productions.

$$\begin{aligned} A &\rightarrow \beta A' \\ A' &\rightarrow \alpha A' / \epsilon \end{aligned}$$

Eg:- $E \rightarrow E + T / T$

$T \rightarrow T * F / F$

$F \rightarrow (E) / id$ (left recursive grammar)

so, we need to eliminate left-recursion from the above grammar

$$(i) \quad E \rightarrow E + T / T$$

$A \quad A \quad \alpha \quad \beta$

$$\begin{aligned} A &= E \\ \alpha &= +T \\ \beta &= T \end{aligned}$$

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow +TE' / \epsilon \end{aligned}$$

(7)

$$\text{(i) } T \rightarrow T * F / F$$

$$A \quad A \alpha / B$$

$$T \rightarrow FT'$$

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

$$\text{Eg2: } S \rightarrow \underline{SOS} \underline{S} / \underline{\epsilon}$$

$$A \quad A \alpha \quad B$$

$$A \rightarrow A \alpha / B$$

$$S \rightarrow \cancel{AOS} OIS'$$

$$S' \rightarrow OSIS' / \epsilon$$

$$\text{(ii) } F \rightarrow (E) / \text{id} \text{ (Non left recursive)}$$

After eliminating the left recursion
from the grammar

$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

$$F \rightarrow (E) / \text{id} / \epsilon$$

$$\text{Eg3: } S \rightarrow (L) / \epsilon - \text{No left recursion}$$

$$L \rightarrow L_1 S / S$$

* Elimination of left recursion for multiple production.

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

$$\boxed{A \rightarrow B_1 A' | B_2 A'}$$

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

$$\text{Eg1: } \text{Expr} \rightarrow \text{Expr} \underline{+ Expr} | \underline{Expr * Expr} / \text{id}$$

$$A \quad A \quad \alpha_1 \quad A \quad \alpha_2 \quad B_1$$

$$\text{Expr} \rightarrow \text{id Expr}$$

$$\text{Expr} \rightarrow + \text{Expr Expr} | * \text{Expr Expr} / \epsilon$$

$$\text{Eg2: } S \rightarrow Sx | S \underline{sb} | xS | a$$

$$A \quad A \alpha_1 \quad A \alpha_2 \quad B_1 \quad B_2$$

$$S \rightarrow xS | aS'$$

$$S' \rightarrow xs' | sbss' / \epsilon$$

$$\text{Eg3: } S \rightarrow Aa | b$$

$$A \rightarrow Ac | Sd | f$$

$$S \rightarrow Aa | b \rightarrow \text{NO left recursion}$$

$$A \rightarrow Ac | Sd | f$$

(1) Elimination of left factoring:

→ A grammar contains production rule in the form

$$A \rightarrow \alpha B_1 | \alpha B_2 | \dots | r_1 | r_2$$

↓

Then that grammar contains left factoring.

→ TDOP parsers can't handle left factoring the grammar contains left factoring.

→ we can eliminate the left factoring by replacing with the following production.

$$\begin{array}{l} A \rightarrow \alpha A' | r_1 | r_2 \\ A' \rightarrow B_1 | B_2 \end{array}$$

Ex1: $S \rightarrow iEts | iEtSs | a$ (it contains left factoring)
 $E \rightarrow b$

Here, $A = S$, $\alpha = iEts$ $B = E$ $B_1 = es$ $r_1 = a$

$$\begin{array}{l} B \rightarrow bB | b \\ A \rightarrow bB | b \\ B_1 \rightarrow B | \epsilon \end{array}$$

$$S \rightarrow iEtss | a$$

$$S \rightarrow \epsilon | es$$

Ex2: $A \rightarrow \frac{\alpha}{\alpha B_1} \frac{\alpha}{\alpha B_2} \frac{\alpha}{\alpha}$

$$\left(\begin{array}{l} B \rightarrow bB | b \\ A \rightarrow \frac{\alpha AB}{\alpha B_1} \frac{\alpha AB}{\alpha B_2} \frac{\alpha}{\alpha B_3} \\ A \rightarrow \alpha A' \\ A' \rightarrow \alpha B | \alpha \epsilon \end{array} \right)$$

Ex3: $X \rightarrow \frac{x+x}{\alpha B_1} \frac{x+x}{\alpha B_2} \frac{x+x}{\alpha}$

$$\begin{array}{l} X \rightarrow xx^1 | D \\ x^1 \rightarrow +x | *x \end{array}$$

Ex4: $E \rightarrow T+E/T$
 $T \rightarrow int / int * T / (E)$

$$T \rightarrow \frac{int}{\alpha B_1} \frac{int * T}{\alpha B_2} / (E)$$

$$\begin{array}{l} T \rightarrow T+E/T \\ \alpha B_1 = \alpha B_2 = \epsilon \\ E \rightarrow TE^1 \\ E^1 \rightarrow +E/\epsilon \end{array}$$

$$T \rightarrow int T^1 / (E)$$

$$T^1 \rightarrow \epsilon / *T^1$$

Top Down Parsing:

→ In Top-down parsing, The parse tree is constructed from root node to child node.

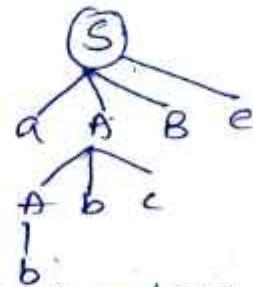
→ Top-down parser uses left-most derivation to derive the input string from the grammar.

→ TDP is starting from start symbol of grammar and reaching the i/p string

Eg:- Grammas is, derive the i/p string w=abbcde,

$$\begin{aligned} S &\rightarrow aABe \\ A &\rightarrow Abc/b \\ B &\rightarrow d \end{aligned}$$

LMD:- $S \rightarrow aABe$
 $\quad\quad\quad \underline{a}AbcBe$
 $\quad\quad\quad abbcBe$
 $\quad\quad\quad abbcde$



→ TDP constructed for the grammar if it is free from
 • ambiguity • left recursion.

→ The problem with Top parser is if we have more than one alternative for production, which alternative we should use. $\left\{ \begin{array}{l} E \rightarrow E+F/T \rightarrow T+F/F \\ F \rightarrow (E)/id \end{array} \right.$
 ↴ ambiguous, left recursive

Eg:- $E \rightarrow TE'$ → Unambiguous grammar

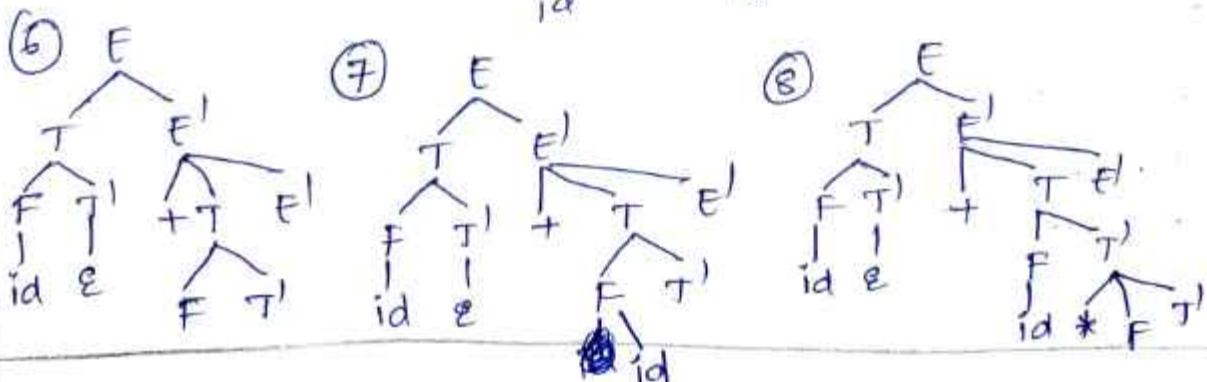
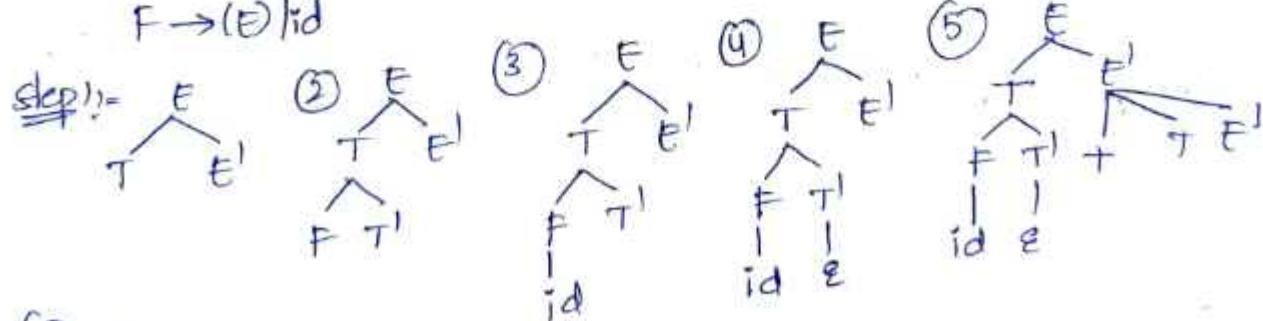
$$E' \rightarrow +TE'$$

$$T \rightarrow FT'$$

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

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

↳ Derive on i/p string $w=id+id*id$



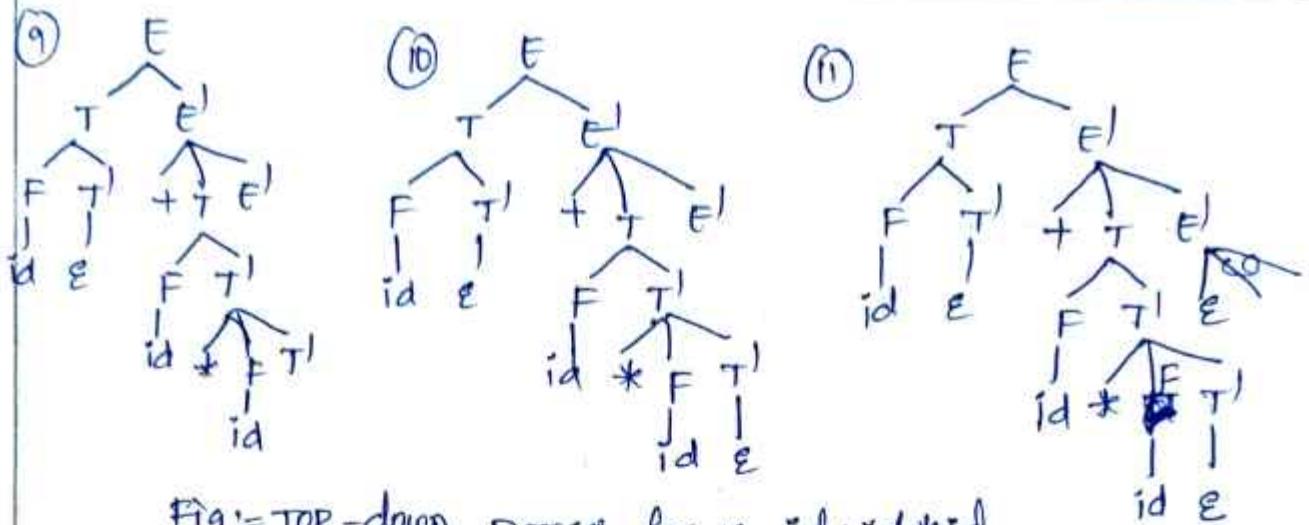


Fig:=TOP-down parses for $w=id+id*id$

i) Recursive-descent Parsing:= (with back tracking)

→ The construction of parse tree starts from root & proceed to child node.

→ Recursion= A function which is called by itself.

steps for constructing recursive descent parser=

- (i) If the i/p is a non-terminal, then call corresponding function
- (ii) If the i/p is terminal, then compare terminal with i/p symbol.
if both are same, then increment input pointer
- (iii) If a Non-terminal contains more than one production then
all the production code should be written in the corresponding
function.

$$\text{Eq1: } E \rightarrow iE \quad E \rightarrow +iE' \mid \epsilon$$

(In the given grammar we have 2 Nonterminals E, E' , so define
two functions, same as C-lang functions)

$E()$

{

```
if (input == 'i')
    input++;
    EPRIME();
}
```

EPRIME();

```

{
    if (input == '+')
    {
        if (input == 'i')
            input++;
        EPRIM();
    }
    else
        return;
}

```

i/p string: i+i

i/p string

~~Eg2:~~ Grammatical E → TE¹ E¹ → +TE¹/ε
 T → FT¹ FT¹ → *FT¹/ε F → (E)/ε

```

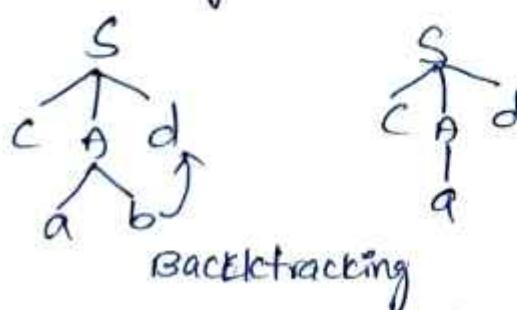
E()
{
    T();
    EPRIME();
}
EPRIME();
{
    if (input == '+')
    {
        input++;
        T();
        EPRIME();
    }
    else
        return;
}
T()
{
    if (input == '(')
    {
        input++;
        F();
        TPRIME();
    }
    else
        return;
}
F()
{
    if (input == ')')
        input++;
    else if (input == '*')
        input++;
    else if (input == 'i')
        input++;
    else if (input == 'd')
        input++;
}

```

~~Eg3:~~ S → CAD

A → ab/a

w = i/p string w = cad



(If None of the product for is satisfying the condition
 then we should check another possibility for "S")

```

input++
EPRIME();
{
    EPRIME()
    {
        if (input == '+')
        {
            input++;
            if (input == 'i')
                input++;
            EPRIME();
            else
                return;
        }
    }
}

```

$i/p = id + id$

eg: $\frac{E}{T} \rightarrow T E$

$T \rightarrow T E$

$T \rightarrow F T$

$F \rightarrow (E) / id$

$E \rightarrow T$

$T \rightarrow EPRIME()$

$EPRIME()$

$\frac{E}{T} \rightarrow T E$

$T \rightarrow F T$

$F \rightarrow (E) / id$

$E \rightarrow T$

$T \rightarrow EPRIME()$

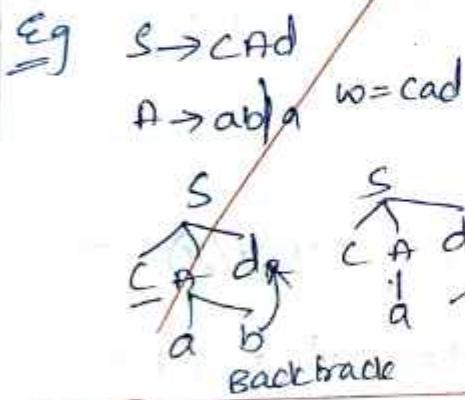
$EPRIME()$

```

    T()
    EPRIME();
    {
        else
            return;
    }
    T()
    {
        F()
        TPRIME();
    }
    TPRIME()
    {
        if (input == '+')
        {
            input++;
            F()
            TPRIME();
        }
        else
            return;
    }
}

```

(8) F()
 {
 if (input == 'i')
 {
 input++;
 E();
 }
 if (input == 'l')
 input++;
 else
 if (input == 'd')
 input++;
 }
}



(\because If none of the product for A is satisfying the condition, then we should check another possibility for S)

(ii) predictive Parsing: (a) LL(1) & Non recursive parser :-

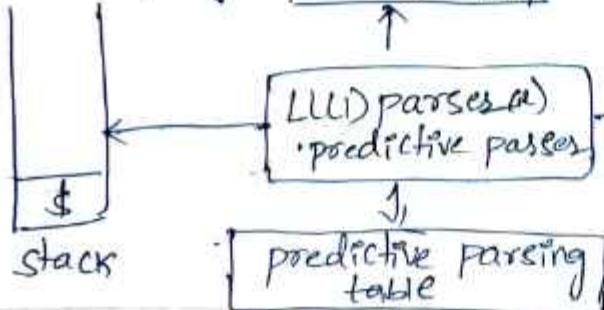
LL(1) - means

- It scans the input from left to Right by using left most derivation.

• $LL(1) \rightarrow 1$ - \Rightarrow look ahead - at a time we are reading 1 character at a time

i/p string \hookrightarrow

	\$	i/p buffer
--	----	------------



Steps for constructing

LL(1) parser:-

- ① Elimination of left recursion
- ② Elimination of left factoring
- ③ Calculation of FIRST() & FOLLOW
- ④ Construction of parsing table.

LL(1) parser \rightarrow parsing algorithm

LL(1) parsing table - Data structure which is constructed from the LL(1) grammar.

Stack = Data structure used to store the grammar symbols

- Always the bottom of the stack is \$
 - \rightarrow The end of ip buffer is \$.
 - \rightarrow after \$, the starting symbol of ip is pushed to the stack.
- To construct the parsing table we should compute two functions. (1) FIRST() (2) FOLLOW()

FIRST(x):

- if x is a terminal then $\text{FIRST}(x) = \{x\}$
- if x is a Non-terminal & $x \rightarrow y_1 y_2 y_3 \dots y_n$ is a production for x.

$$\text{FIRST}(x) = \text{FIRST}(y_1) \cup$$

if $\text{FIRST}(y_1)$ contains ϵ add $\text{FIRST}(y_2)$ to $\text{FIRST}(x)$

if all $\text{FIRST}(y_1, y_2, \dots, y_n)$ contains ϵ then add $\underline{\epsilon}$ to $\text{FIRST}(x)$

FOLLOW(B):

\rightarrow Always the follow() of starting symbol is \$.

\rightarrow If $A \rightarrow \alpha B \beta$

$$\text{Follow}(B) = \text{FIRST}(\beta) \cup$$

if $\text{FIRST}(\beta)$ contain ϵ or $A \rightarrow \alpha B$ then
add $\text{Follow}(A)$ to $\text{Follow}(B)$

Eg! Given Grammar is

$$E \rightarrow E + T \mid T \quad \text{left recursion} \quad \left. \begin{array}{l} \text{eliminate left} \\ \text{recursion} \end{array} \right\} \text{from these 2 productions}$$
$$T \rightarrow T * F \mid F \quad \text{left recursion}$$
$$F \rightarrow (F) \mid \text{id}$$

$E \rightarrow TE'$
 $E' \rightarrow +TE' / \epsilon$
 $T \rightarrow FT'$
 $T' \rightarrow *FT' / \epsilon$
 $F \rightarrow (E) / id$

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

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

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

- $\text{FIRST}(T') = \{ *, \epsilon \}$

$\rightarrow \text{Follow}(E) = \{ +, \$,) \}$

$\downarrow \quad F \rightarrow (E) / id$

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

Follow(E) = $\{ \text{Follow}(E) + \text{Follow}(E') \}$

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

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

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

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

Follow(T') $\Rightarrow \{ \text{Follow}(T) + \text{Follow}(T') \}$

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

Follow(F) = $\text{Follow}(T') + \text{FIRST}(T) + \text{Follow}(T)$

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

(iii)
predictive parsing table :-

	id	+	*	()	\$
E	$E \rightarrow TE^1$			$E \rightarrow TE^1$		
E^1		$E^1 \rightarrow +TE^1$			$E^1 \rightarrow E$	$E^1 \rightarrow \epsilon$
T	$T \rightarrow FT^1$				$T \rightarrow FT^1$	
T^1		$T^1 \rightarrow E$	$T^1 \rightarrow *FT^1$		$T^1 \rightarrow E$	$T^1 \rightarrow \epsilon$
F	$F \rightarrow id$				$F \rightarrow (E)$	

$W = id + id$

Stack	Input	Action	Comments
$E\$$	$id + id \$$	$E \rightarrow TE^1$	The S/P is parsed
$TE^1 \$$	$id + id \$$	$T \rightarrow FT^1$	$eg_1: S \rightarrow (L) / q$ $L \rightarrow L, S / S$
$FT^1 \$$	$id + id \$$	$F \rightarrow id$	$W = (a)$
$id T^1 E^1 \$$	$id + id \$$		$eg_2: S \rightarrow AAB BAA \epsilon$
$T^1 E^1 \$$	$+ id \$$	$T^1 \rightarrow \epsilon$	$A \rightarrow aAB \epsilon$
$\epsilon E^1 \$$	$+ id \$$		$B \rightarrow bB \epsilon$
$E^1 \$$	$+ rd \$$	$E^1 \rightarrow +TE^1$	$W = "aabbb"$
$*TE^1 \$$	$*id \$$	$T \rightarrow FT^1$	$S \rightarrow RETS RETS cS q$
$FT^1 E^1 \$$	$id \$$	$F \rightarrow id$	$E \rightarrow B \text{ (Non-LL(1))}$
$id T^1 E^1 \$$	$id \$$	$T^1 \rightarrow \epsilon$	$eg_3: S \rightarrow AB CBR Bq$
$\epsilon E^1 \$$	$\$$	$E^1 \rightarrow \epsilon$	$A \rightarrow da BC$
$\$$	$\$$	accepted	$B \rightarrow g k$ $C \rightarrow h k$

Bottom-up Parsers :-

→ Bottom-up parsers construct parse tree from child (bottom) node to root node (top).

→ Bottom up parsers use RMDP in Reverse order.

→ In bottom-up parser the ~~part of~~ "reducing" I/P string w to start symbol of grammar

$$\text{Eq1: } S \rightarrow AABE$$

$$A \rightarrow Abc/b$$

$$B \rightarrow d$$

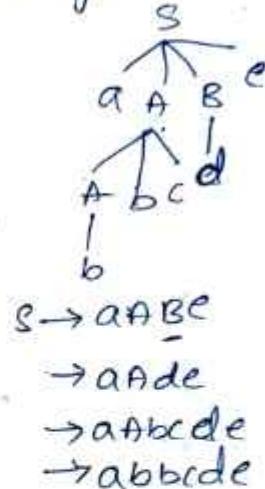
Input string $w = abbcde$

$$a\cancel{A}bcde \quad (A \rightarrow b)$$

$$a\cancel{A}de \quad (A \rightarrow Abc)$$

$$a\cancel{A}Be \quad (B \rightarrow d)$$

$$S \quad (C \rightarrow a\cancel{A}Be)$$



$$\text{Eq2: } E \rightarrow ETT/T$$

$$E \rightarrow T*T/F/F$$

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

$$w = id * id$$

$$id * id \quad (F \rightarrow id)$$

$$T * id \quad (T \rightarrow id)$$

$$T * id \quad (\cancel{T \rightarrow T})$$

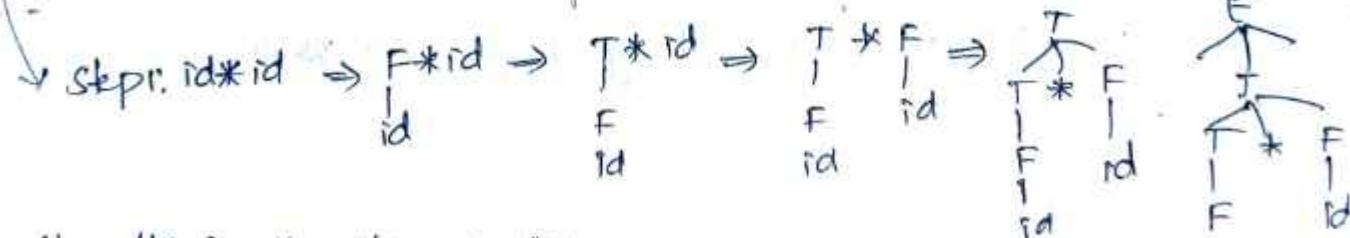
$$\cancel{E \rightarrow id} \quad T * id \quad (F \rightarrow id)$$

$$T * F \quad (\cancel{T \rightarrow T} * F)$$

$$T \quad (E \rightarrow ST)$$

$$E$$

→ The problem with Bottom-up parsers is when to reduce the ~~or~~ to production.



Handle & Handle pruning :-

Handle: Handle is a substring, which matches with right side of production

→ if Handle matches with the right-side of production, then it is replaced with left hand side Non-terminal

Eg1: $S \rightarrow aABC$ $w = abcd e$
 $A \rightarrow abc/b$
 $B \rightarrow d$

Right sentential form	Handle	Reducing production
<u>abbcde</u>	b	$B \rightarrow b$
a <u>bcde</u>	abc	$A \rightarrow abc$
aAde	d	$B \rightarrow d$
aABC	aABC	$S \rightarrow aABC$
s		

Handles during parse of $id_1 * id_2$

Eg2: $E \rightarrow E + T/T$ $w = id_1 * id_2$
 $T \rightarrow T * F/F$
 $F \rightarrow (E)id$

Right sentential form	Handle	Reducing production
<u>id_1 * id</u>	$E id$	$F \rightarrow id$
$F * id$	F	$F \rightarrow F$
$T * id$	id	$T/F \rightarrow id$
$T * F$	TF	$T \rightarrow T * F$
T	T	$E \rightarrow T$

→ Handle pruning is obtained by Right most derivation in reverse order.

(f) Shift-Reduce parser:

→ Shift-Reduce parser use two data structures ① stack ② input buffer

① Stack - is used to store the grammar symbol

② Input buffer - holds the i/p string to be parsed

Stack	Input string
\$	w\$

Actions of shift-Reduce parser:

1. Shift — shift the next I/p symbol onto the top of the stack.
2. Reduce — If the top of stack is matched with Right side of production, then it is reduced to corresponding non-terminal.
3. Accept — successful completion of parsing
4. Error — Discovers a ^{syntax} error and call error recovery methods

Eg :- $E \rightarrow E + T / T$ $w = id * id$
 $T \rightarrow T * F / F$
 $F \rightarrow (E) / id$

Fig :- configuration of shift-Reduce parser on input id

Stack	Input Buffer	Action
\$	<u>id * id \$</u>	shift
\$ id	* id \$	Reduce by $F \rightarrow id$
\$ F	* id \$	Reduce $F \rightarrow F$
\$ T	* id \$	shift
\$ T *	id \$	shift
\$ T * id	\$	Reduce $F \rightarrow id$
\$ T * F	\$	Reduce $T \rightarrow T * F$
\$ T	\$	Reduce $E \rightarrow T$
\$ E	\$	Accepted

Conflicts during shift-Reduce parsing

- ① Shift-Reduce conflict : cannot decide whether to shift or to reduce
- ② Reduce-reduce conflicts if the ^{two} productions have same handle (substring),

4

Shift-Reduce

Eg 2: $S \rightarrow (L) \cdot a$

$L \rightarrow L, S / S$ parse the string $(a, (a, a))$ using SR-Power

Stack	I/P Buffer	parsing Action
\$	$((a, (a, a)) \cdot a)$	shift
$\$ (\cdot a$	$a, (a, a)) \cdot a$	shift
$\$ (a$	$, (a, a)) \cdot a$	Reduce $S \rightarrow a$
$\$ (S$	$, (a, a)) \cdot a$	$L \rightarrow S$ Reduce
$\$ (L$	$, (a, a)) \cdot a$	shift
$\$ (L$	$(a, a)) \cdot a$	shift
$\$ (L C$	$a, a)) \cdot a$	shift
$\$ (L, a$	$, a)) \cdot a$	Reduce $S \rightarrow a$
$\$ (L, S$	$, a)) \cdot a$	shift Reduce $L \rightarrow S$
$\$ (L, L$	$, a)) \cdot a$	shift
$\$ (L, (L$	$, a)) \cdot a$	shift
$\$ (L, (L, a$	$) \cdot a$	Reduce $S \rightarrow a$
$\$ (L, (L, S$	$) \cdot a$	Reduce $L \rightarrow L, S$
$\$ (L, (L$	$) \cdot a$	shift
$\$ (L, (L$	$) \cdot a$	Reduce $S \rightarrow (L)$
$\$ (L, S$	$) \cdot a$	Reduce $L \rightarrow L, S$
$\$ (L$	$) \cdot a$	shift
$\$ (L$	$\cdot a$	Reduce $S \rightarrow (L)$
$\$ S$	$\cdot a$	Accept

Eg 3:- $S \rightarrow TL;$
 $T \rightarrow \text{int} / \text{float}$
 $L \rightarrow L, \text{id} / \text{id}$

w = int id id ④ s → os0 / |s1|2

$$\omega = 10201$$

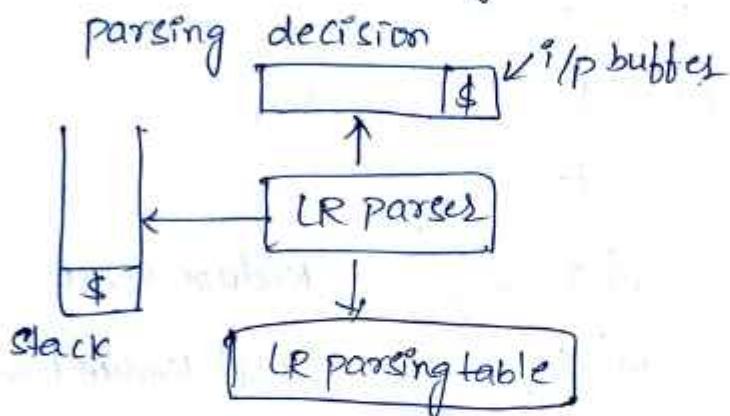
Introduction to LR Parsing := LR parser is a type of Bottom-up parser

→ LR(λ) parses some i/p from left to right.

→ It uses Rightmost Derivation in Reverse ^{order} ~~order~~.

desire ~~the~~ to reduce ip string to start symbol of grammar

$K \rightarrow$ No. of lookahead symbols that are used to make



Stack - Data structure used to store grammatical symbols.

i/p - " " " " " i/p string to be parsed

LR pause → uses algorithm to make decisions.

LR Parsing table - is constructed by using LR(0) items.

- These table uses two functions (ϕ)

Benefits of LR(k) parsing:

→ most generic Non-Backtracking shift Reduced parsing
Technique.

→ These parsers can recognize all programming languages for which CFG can be written.

→ They are capable of detecting syntactic errors as soon as possible while scanning of i/p.

Simple LR Parsing:

(15)

$$\text{Eg: } E \rightarrow E + T$$

$$E \rightarrow T$$

$$T \rightarrow T * F$$

$$T \rightarrow F$$

$$F \rightarrow (E)$$

$$F \rightarrow \text{id}$$

① Augmented Grammars

$$E' \rightarrow E \quad E' \rightarrow F$$

$$E \rightarrow E + T$$

$$E \rightarrow T$$

$$T \rightarrow T * F$$

$$T \rightarrow F$$

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

Steps to construct SLR parser

① Introduce augmented grammar

② calculate canonical collection of LR(0) items.

③ construct

SLR parsing table by using (i) Goto (ii) Action function

$E' \rightarrow \cdot E$	$(I_1, +) - I_6$
$E \rightarrow \cdot E + T$	$E \rightarrow E + \cdot T$
$E \rightarrow \cdot T$	$T \rightarrow \cdot F$
$T \rightarrow \cdot T * F$	$T \rightarrow \cdot T * F$
$T \rightarrow \cdot F$	$F \rightarrow \cdot (E)$
$F \rightarrow \cdot (E)$	$F \rightarrow \cdot \text{id}$
$F \rightarrow \cdot \text{id}$	

Goto (I_0, E) - I_1

$$E' \rightarrow E \cdot$$

$$E \rightarrow E \cdot T$$

$I_0, T) - I_2$

$$E \rightarrow T \cdot$$

$$T \rightarrow T \cdot * F$$

$I_0, F) - I_3$

$$T \rightarrow F \cdot$$

$I_0, () - I_4$

$$F \rightarrow (\cdot E)$$

$$E \rightarrow \cdot E + T$$

$$T \rightarrow \cdot F$$

$$E \rightarrow \cdot T$$

$$F \rightarrow \cdot \text{id}$$

$I_0, \text{id}) - I_5$

$$F \rightarrow \text{id} \cdot$$

$(I_2, *) - I_7$
$T \rightarrow T * \cdot F$
$F \rightarrow \cdot (E)$
$F \rightarrow \cdot \text{id}$

$(I_4, E) - I_8$
$F \rightarrow (E \cdot)$
$E \rightarrow \cdot \text{id}$

$(I_4, T) - I_2$
$E \rightarrow T \cdot$
$E \rightarrow \cdot E + T$
$T \rightarrow T \cdot * F$

$(I_4, F) - I_3$
$T \rightarrow F \cdot$

$(I_4, +) - I_6$
$E \rightarrow E + \cdot T$
$T \rightarrow \cdot T * F$
$T \rightarrow \cdot F$
$F \rightarrow \cdot (E)$

$(I_6, () - I_4$
$F \rightarrow (\cdot E)$
$F \rightarrow \cdot \text{id}$

$(I_6, T) - I_1$
$E \rightarrow E + T \cdot$
$T \rightarrow T \cdot * F$
$T \rightarrow F \cdot$

$(I_6, F) - I_3$
$T \rightarrow F \cdot$

$(I_6, id) - I_5$
$F \rightarrow \cdot \text{id}$

$(I_7, F) - I_{10}$
$T \rightarrow T * F \cdot$

$(I_7, C) - I_4$
$F \rightarrow (\cdot E)$
$E \rightarrow \cdot E + T$
$T \rightarrow \cdot T$
$T \rightarrow \cdot T * F$

$(I_7, id) - I_5$
$F \rightarrow \cdot \text{id}$

$(I_8,)) - I_{11}$
$F \rightarrow (\epsilon)$

$(I_8, +) - I_6$
$E \rightarrow E + \cdot T$
$T \rightarrow \cdot T * F$
$T \rightarrow \cdot F$
$F \rightarrow \cdot (E)$

$(A_1, *) - \text{I}_1$

13

$T \rightarrow T * F$

$F \rightarrow \cdot (E)$

$F \rightarrow \cdot id$

STATE	ACTION						GOTO		
	id	$+$	$*$	$($	$)$	$\$$	E	T	F
0	s_5				s_4		1	2	3
1							Accept		
2		r_2	s_7			r_2	r_2		
3		r_4	r_4			r_4	r_4		
4	s_5				s_4		8	2	3
5	s_5	r_6	r_6			r_6	r_6		
6	s_5				s_4			9	3
7	s_5				s_4				10
8		s_6				s_{11}			
9		r_1	s_7			r_1	r_1		
10		r_3	r_3			r_3	r_3		
11		r_5	r_5			r_5	r_5		

1. $I_1: E \xrightarrow{*} E$.

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

2. $I_2: E \rightarrow T$.

$E \xrightarrow{*} \cdot E$

$\text{Follow}(E) = \{+,), \$\}$

3. $I_3: T \rightarrow F, (r_4)$

1 $E \rightarrow EFT$

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

4. $I_5: F \rightarrow id$.

2 $E \rightarrow T (r_2)$
3 $T \rightarrow T * F$

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

5. $I_9: E \rightarrow E + T$.

4 $T \rightarrow F (r_3)$

$(2, t) (2, +) (2, \$) \rightarrow r_2$

$I_{10}: T \rightarrow T * F$.

5 $F \rightarrow (E)$

$(3, \$) (3, +) (3, \$) (3, *) \rightarrow r_3$

$I_{11}: F \rightarrow (E)$.

6 $F \rightarrow id$

$(5, \$) (5, *) (5, \$) (5, *) \rightarrow r_6$

$I_1: S^* \xrightarrow{*} S$

$$w = id * id + id$$

STACK	INPUT	ACTION
\$0	id * id + id \$	shift 5
\$0 id\$	* id + id \$	reduce 6 $E \rightarrow id$
\$0 F3	* id + id \$	reduce 4 $T \rightarrow F$
\$0 T2	* id + id \$	S7
\$0 T2 * 7	id + id \$	S5
\$0 T2 * 7 id\$	+ id \$	r6 $F \rightarrow id$
\$0 T2 * 7 F10	+ id \$	r3 $T \rightarrow T * F$
\$0 T2	+ id \$	r2 $E \rightarrow T$
\$0 E1	+ id \$	S6
\$0 E1 + 6	id \$	S5 to (E)
\$0 E1 + 6 id \$	\$	r6 $F \rightarrow id$
\$0 E1 + 6 F3	\$	r2 $T \rightarrow F$
\$0 E1 + 6 T9	\$	r1 $E \rightarrow E + T$
\$0 E1	\$	acceptance

Eg:- 2 $S \rightarrow AS / b$

$$A \rightarrow SA / a \quad w = ababab$$

→ shift → while performing shifting first we have to shift input symbol on to the top of the stack & then shift the state number.

→ reduce → while performing reduce operation. (if $F \rightarrow id$ are reducing) If Right hand side contains one symbol we need to pop two symbols from top of the stack.

(ii) TOP LR(0) parser:

steps for construct LR(0) parser

(1) Introduce augmented grammars

(2) calculate canonical collection of LR(0) items.

(3) construct LR(0) parsing table by using (i) Goto (ii) action functions.

w = string = aaabb.

$$\text{e.g. } S \rightarrow AA$$

$$A \rightarrow aA/b$$

Augmented grammar

$$S' \rightarrow S$$

$$S \rightarrow AA$$

Closure =
LR(0) items - I_0

$$S' \rightarrow S.$$

$$S \rightarrow .AA$$

$$A \rightarrow .aA/b$$

Goto
 $(I_0, S) \rightarrow I_1$

$$S' \rightarrow S.$$

Goto $(I_0, A) - I_2$

$$S \rightarrow A \cdot A$$

$$A \rightarrow .aA/b$$

Goto $(I_0, a) - I_3$

$$A \rightarrow a \cdot A$$

$$A \rightarrow .aA/b$$

$(I_2, A) \rightarrow I_5$

$$S \rightarrow AA.$$

$(I_2, a) - I_3$

$$A \rightarrow a \cdot A$$

$$A \rightarrow .aA/b$$

$(I_2, b) - I_4$

$$A \rightarrow A \cdot A$$

$$A \rightarrow b.$$

$(I_3, a) - I_3$

$$A \rightarrow a \cdot A$$

$$A \rightarrow .aA/b$$

$(I_3, A) - I_6$

$$A \rightarrow aA.$$

$(I_3, b) - I_4$

$$A \rightarrow b.$$

LR(0) Parsing Table

	Action		Goto			
	a	b	\$	A	S	
0	S_3	S_4		2	1	
1			Accept			
2	S_3	S_4		5		
3	S_3	S_4		6		
4	r_3	r_3	r_3			
5	r_1	r_1	r_1			
6	r_2	r_2	r_2			

wrote down the states that contains "•" at the end of the production.

$$I_1 \rightarrow S' \rightarrow S.$$

$$I_4 \rightarrow A \rightarrow b.$$

$$I_6 \rightarrow A \rightarrow aA.$$

$$I_5 \rightarrow S \rightarrow AA.$$

Given Grammar

$$S \rightarrow AA \quad (1)$$

$$A \rightarrow aA \quad (2)$$

$$A \rightarrow b \quad (3)$$

→ The given string $w = aa\ bb \Rightarrow w = aabb\$$

Stack	Input	Action
\$0	aabb\$	shift 3
\$0a3	abb\$	shift 3
\$0a3a3	b\$	$\tau_1 \rightarrow AA$ shift 4
\$0a3a3b4	b\$	reduce τ_3 $A \rightarrow b$
\$0a3a3b4	\$	shift 4

→ while performing shifting 1st we have to shift i/p symbol onto the top of the stack & then state Number

Stack	Input	Action
\$0	aabb\$	S_3
\$0a3	abb\$	S_3
\$0a3a3	bb\$	S_4
\$0a3a3b4	b\$	reduce τ_3 $A \rightarrow b$
\$0a3a3A6	b\$	$\tau_2 A \rightarrow AA$
\$0a3A6	b\$	$\tau_2 A \rightarrow AA$
\$0A2	b\$	S_4
\$0A2B4	\$	$\tau_3 A \rightarrow b$
\$0A2A5	\$	$\tau_1 S \rightarrow AA$
\$0S1	\$	Accepted.

String $w = aabb$ is parsed through the grammar by using LR(0) parser.

Canonical LR parsing:-

Step 1: Augmented grammars
 Eg.: $S \rightarrow CC$ $S \rightarrow S$ $S \rightarrow CC$ $w = add$
 $C \rightarrow ac$ $C \rightarrow ac$ $C \rightarrow ac$
 $C \rightarrow d$ $C \rightarrow d$ $C \rightarrow d$

Step 2: LR(1) items = $LR(0) + \text{lookahead}$

$I_0: S^l \rightarrow \cdot S, \$$ $S \rightarrow \cdot CC, \$$ $C \rightarrow \cdot ac, ald$ $C \rightarrow \cdot d, ald$	$(I_2, d) - I_1$ $C \rightarrow d \cdot, \$$
$(I_0, S) - I_1$ $S^l \rightarrow S \cdot, \$$	$(I_3, c) - I_2$ $C \rightarrow ac \cdot, ald$
$(I_0, C) - I_2$ $S \rightarrow C \cdot C, \$$ $C \rightarrow \cdot ac, \$$ $C \rightarrow \cdot d, \$$	$(I_3, a) - I_3$ $C \rightarrow a \cdot c, ald$ $C \rightarrow \cdot ac, ald$ $C \rightarrow \cdot d, ald$ $(I_3, d) - I_4$ $C \rightarrow d \cdot, ald$
$(I_0, a) - I_3$ $C \rightarrow a \cdot c, ald$ $C \rightarrow \cdot ac, ald$ $C \rightarrow \cdot d, ald$	$(I_6, C) - I_5$ $S \rightarrow ac \cdot, \$$
$(I_0, d) - I_4$ $C \rightarrow d \cdot, ald$	$(I_6, a) - I_6$ $C \rightarrow a \cdot c, \$$ $C \rightarrow \cdot ac, \$$ $C \rightarrow \cdot d, \$$
$(I_2, C) - I_5$ $S \rightarrow CC \cdot, \$$	$(I_6, d) - I_7$ $C \rightarrow d \cdot, \$$
$(I_2, a) - I_6$ $S \rightarrow a \cdot c, \\$ $C \rightarrow \cdot ac, \$$ $C \rightarrow \cdot d, \$$	

Step 3: construction of CLR parsing Table.

Action			Goto	Goto
a	d	\$	s	c
0 S_3	S_4		1	2
1				
2 S_6	S_7			5
3 S_3	S_4			8
4 $r_3 \rightarrow d$	$r_3 \rightarrow d$			
5				
6 S_6	S_7			
7				9
8 $r_2 \rightarrow ac$	$r_2 \rightarrow ac$			
9				

Accepted

1: $S \rightarrow cc$
 2: $c \rightarrow ac$
 3: $c \rightarrow d$
 4: $s \rightarrow s, \$$
 5: $c \rightarrow d, ald$
 6: $s \rightarrow cc, \$$
 7: $c \rightarrow d, \$$
 8: $c \rightarrow ac, ald$
 9: $s \rightarrow ac, \$$

LALR Table: $I_3 = I_6$

$I_4 = I_7$

$I_8 = I_9$

Action			Goto
a	d	\$	s c
0 S_3, S_6	S_4, S_7		1 2
1 S_3, S_6	S_4, S_7	ACC	ACC
2 S_3, S_6	S_4, S_7		5
36 S_3, S_6	S_4, S_7		89
47 r_3	r_3	r_3	
5			
89 r_2	r_2	r_2	

LALR Parsing Table

w = add

stack	input buffer	Action
\$0	add \$	S_3, S_6
\$0936	dd \$	S_4, S_7
\$0936d47	d \$	$r_3 \rightarrow d$
\$0936c89	d \$	$r_2 \rightarrow ac$
\$0C2	d \$	S_7
\$092d49	\$	$r_3 \rightarrow do$
\$092d	\$	
\$0C2L5	\$	$r_1, S \rightarrow cc$
\$0S1	\$	Accepted



→ In LALR parser combine the ~~state~~ same states that are having different lookahead symbols.

Canonical LR parser = (CLR parser)

(15) (16)

Step 1: Augmented grammar

$S \rightarrow CC$
 $C \rightarrow cC$

$C \rightarrow d$
 $S \rightarrow CC$

$C \rightarrow cC$

Step 2: Calculating LR(0) items

$I_0 \rightarrow S \cdot, \$$
 $(I_0, S) - I_1$

$S \rightarrow C \cdot C, \$$
 $S \rightarrow C \cdot C, \$$

$C \rightarrow C \cdot C, c/d$
 $(I_0, C) - I_2$

$C \rightarrow C \cdot C, c/d$
 $S \rightarrow C \cdot C, \$$

$C \rightarrow C \cdot C, \$$
 $C \rightarrow C \cdot C, \$$

$C \rightarrow C \cdot C, \$$
 $C \rightarrow C \cdot C, \$$

(I_0, S)

$(I_1, C) - I_3$

$S \rightarrow C \cdot C, \$$

$(I_2, C) - I_4$

$C \rightarrow C \cdot C, \$$

$C \rightarrow C \cdot C, \$$

$C \rightarrow C \cdot C, \$$

$(I_3, C) - I_5$

$C \rightarrow C \cdot C, c/d$

$C \rightarrow C \cdot C, c/d$

$C \rightarrow C \cdot C, c/d$

construction of CLR Parsing tables =

	Action			Goto		S.	C.
	c	d	\$	1	2		
0	S_3	S_4					
1							
2	S_6	S_7					
3	S_3	S_4					
4	r_3	r_3					
5			r_1				
6	S_6	S_7					
7			r_3				
8	r_1	r_1					
9							

Parse the input string $w = dd\$$

state	input	string
\$0	dd\$	shift y
\$0d4	d\$	reduce $3 \rightarrow d$
\$0c2	d\$	shift z
\$0c2d7	\$	reduce $3 \rightarrow d$
\$0c2c5	\$	reduce $8, S \rightarrow cc$
\$0s1	\$	Accepted

↳ So, the given input string is accepted by the grammar.

↳ LALR parser = (Look Ahead LR parser)

Eg/1: steps to construct LALR parser

- (1) Introduce Augmented grammars
- (2) calculate the (LRL) items
- (3) construction of LALR parse table.
- (4) Parsing of given I/p string.

Eg/1: construct LALR parser for

$$S \rightarrow cc$$

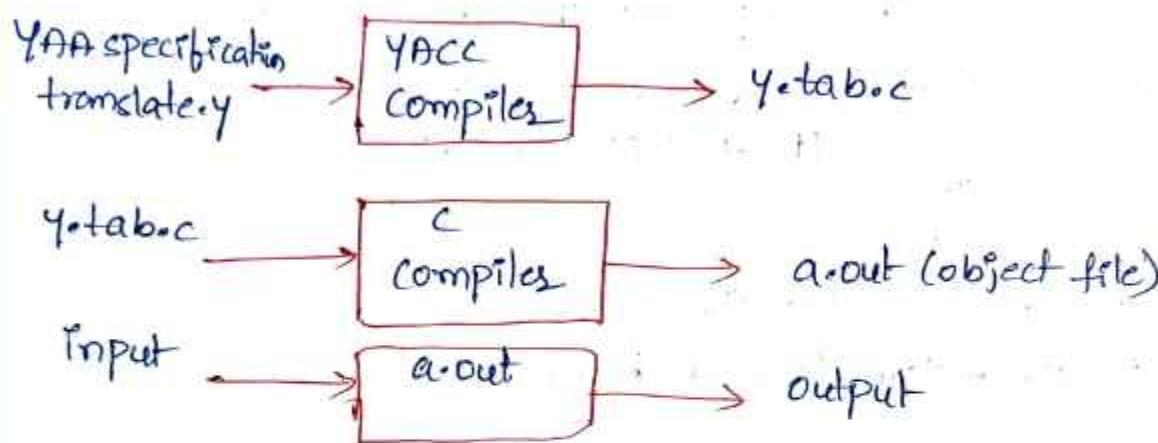
$$C \rightarrow ac$$

$$C \rightarrow d$$

and also parse the string $w = add$.

Parsers generators:-

- we use LALR parsers generators YACC.
- YACC → "yet another compiler compiler"
- YACC is a tool to generate parsers, YACC accepts an tokens as I/P and produces parse tree as O/P.

The parsers generators YACC:-

→ 'translate.y' consists of YACC specification

Structure of YACC program:-

it has three parts.

declarations

%.

translation rules

%.

Auxiliary functions

Declarations:-

→ used to declare the C variables and constants, & header files are also specified here.

Syntax: = %d eg: = %d

%}

int a, b;

Const int a=20;

#include<stdio.h>

%}

(i) Translation Rules:-

→ Translation Rules are enclosed between "%{" & "%}"

Syntax:-

head → body₁ | body₂ | body₃ ...
Eg :- C → aa | bb

head : body₁ {semantic action}

| body₂ {semantic action}

| body₃ {semantic action}

→ In translation Rules we use two symbols $\$\$$, $\$i$, $\$i \rightarrow$

(ii) Auxiliary function:- $\$i \rightarrow$ represent the i th symbol of body.

Eg :- E → E + T (If we want to access E, we have to use
 $\$1, + \rightarrow \$2, T \rightarrow \$3$)

(iii) Auxiliary function:-

→ used to define "C" function.

→ yylex is ~~use~~ function is used here.

Eg :- YACC specification (program) of simple desk calculator.

$\% E = E \rightarrow E + E / T$

$T \rightarrow T * F / F$

$F \rightarrow (E) / id.$

Declaration part:-

%{

#include<ctype.h>

%}

% token weight → specification of RE

% - % → it specifies I/P to desk calculator is an expression by In

lines: Expr '%n' & printf("%d\n", \$1); }

Expr: Expr '+' term { \$\$ = \$1 + \$3; } ;
| term

term: term '*' factor { \$\$ = \$1 * \$3; } ;
| factor

(2) (5)

```
factor : ('Expr') { $$ = $2; }  
| DIGIT;
```

1.1.

```
yylex()  
{  
    int c;  
    c = getchar();  
    if (isdigit(c))  
    {  
        yyval = c - '0';  
        return DIGIT;  
    }  
    return c;  
}
```

yylex → is a lexical analyzer function
which takes our source program
as input and produces token as output

1

2

3

1) Using Ambiguous Grammars

- For language constructs like expressions, an ambiguous grammar provides a shorter, more natural specification than any equivalent unambiguous grammar.
- Another use of ambiguous grammar is in isolating commonly occurring syntactic constructs for special-case optimization, we can add new productions to grammar.
- Sometimes ambiguity rules allow only one parse tree, in such case it is unambiguous, it is possible to design an LR parser that allows same ambiguity resolving choices.

Precedence & Associativity to Resolve conflicts:-

- Consider ambiguous grammar with operators '+' & '*'

$$E \rightarrow E+E \mid E*E \mid (E) \mid id \quad (1)$$
- $E \rightarrow E+T, T \rightarrow T*F$, generates same language gives lower precedence to '*' than '+', makes left associative. [use ambiguous grammar]
- But first we change associativity and precedence of operators + & * without disturbing above grammar.
- Second, the parser for the unambiguous grammar will spend a substantial fraction of its time reducing by the productions $E \rightarrow T \& T \rightarrow F$
- $E^l \rightarrow E$, (1) is ambiguous there will be parsing-action conflicts when we try to produce LR parsing table.

$I_0: E^{\dagger} \rightarrow E$
 $E \rightarrow E+E$
 $E \rightarrow \cdot E+E$
 $E \rightarrow \cdot (E)$
 $E \rightarrow \cdot id$

$I_1: E^{\dagger} \rightarrow E^{\cdot}$
 $E \rightarrow E^{\cdot}+E$

$E \rightarrow E \cdot + E$

$I_2: E \rightarrow (\cdot E)$
 $E \rightarrow \cdot E+E$
 $E \rightarrow \cdot E \times E$
 $E \rightarrow \cdot (E)$
 $E \rightarrow \cdot id$

$I_3: E \rightarrow id$
 $I_4: E \rightarrow E+E$
 $E \rightarrow \cdot E+E$
 $E \rightarrow \cdot E \times E$
 $E \rightarrow \cdot (E)$
 $E \rightarrow \cdot id$

$I_5: E \rightarrow E \times \cdot E$
 $E \rightarrow \cdot E+E$
 $E \rightarrow \cdot E \times E$
 $E \rightarrow \cdot (E)$
 $E \rightarrow \cdot id$

$I_6: E \rightarrow (E^{\cdot})$
 $E \rightarrow E \cdot + E$
 $E \rightarrow E \cdot \times E$
 $I_7: E \rightarrow E+E$
 $E \rightarrow E \cdot + E$
 $E \rightarrow E \cdot \times E$

$I_8: E \rightarrow E \times E$
 $E \rightarrow E \cdot + E$
 $E \rightarrow E \cdot \times E$

$I_9: E \rightarrow (E) \cdot$

conflict occurs in $I_7 \& I_8$, $E \rightarrow E+E^{\cdot}$, $E \rightarrow E \times E^{\cdot}$

PREFIX

$E+E$

STACK

0147

INPUT

*id\$

If ilp is $id+id$.

State	ACTION						GOTO
	id	+	*	()	\$	
0	s_3			s_2			1
1		s_4	s_5			accept	
2	s_3			s_2			6
3		r_4	r_4		r_4		
4	s_3			s_2			7
5	s_3			s_2			8
6		s_4	s_5		s_9		
7		r_1	s_5		r_1	r_1	
8		r_2	r_2		r_2	r_2	
9		r_3	r_3		r_3	r_3	

Fig - parsing table for grammar

"Dangling- Else" Ambiguity:-

stmt → if expr then stmt else stmt
if, expr then stmt

Consider other

Consider the grammar, an abstraction of this grammar where 'i' stands for if expr then, e stands for else and 'a' stands for "all other production".

$s' \rightarrow s$ (augmented grammar)

s → ises / is / a

$$I_0: S^1 \rightarrow S \quad I_2: S \rightarrow i_* S_{\leq 2}$$

$s \rightarrow \cdot ises$

$S \rightarrow .$ is

$s \rightarrow a$

$$T_1: S^1 \rightarrow S.$$

$\Gamma_3 : S \rightarrow a$.

T₄: S → ls, es

I5: 5 → ise. 5

$s \rightarrow \cdot iscs$

$$S \rightarrow \cdot ls$$

$s \rightarrow \cdot a$

16: $S \rightarrow \text{ises}$.

Fig:- LR(0) states for augmented grammars

if expr then stmt — (2)

Ambiguity arises in LR shift/reduce conflict

$s \rightarrow is$ calls for a shift of e

$\text{FOLLOW}(S) = \{e, \$\}$, $S \rightarrow i\$$ call for reduction by ~~several~~

S → is on tip 'e'

In (2) should we shift else onto stack or reduce if expr then Stmt.) we should shift else because it is associated with previous then.

SLR parsing table is constructed.

STATE	ACTION				GOTO
	i	e	a	\$	
0	s ₂		s ₃		1
1				accepted	
2	s ₂		s ₃		4
3		r ₃		r ₃	
4		s ₅		r ₂	
5	s ₂		s ₃		6
6		r ₁		r ₁	

input is iiaea, At line (5) state 4 selects the shift action on input e, whereas at line (9) state 4 calls for reduction by $S \rightarrow iS$ on input \$.

STACK	SYMBOLS	INPUT	ACTION
(1) 0		iiaea \$	shift
(2) 02	i	iae a \$	shift
(3) 022	ii	aea \$	shift
(4) 0223	ii a	ea \$	shift
(5) 0224	ii's	ea \$	reduce by $S \rightarrow a$
(6) 02245	ii'se	a \$	shift
(7) 022453	ii'ea	\$	reduce by $S \rightarrow a$
(8) 022456	ii'ea s	\$	reduce by $S \rightarrow i's$
(9) 024	iS	\$	reduce by $S \rightarrow iS$
(10) 01	S	\$	accept

Fig:- parsing actions on input iiaea

Syntax Directed Definition (SDN) → Semantic

→ SDD is completed in Economic Analysis phase & ICG phase were SDD's.

SDD is Context Free Grammars with Semantic Rules and attributes

SDD = CFG + semantic rules

\rightarrow SDD is also called as "attribute grammars".
 \rightarrow After,

- Attributes are associated with grammar symbols and semantic rules.
- Rules are associated with productions.

→ If 'x' is a symbol and 'a' is one of attribute then $x.a$ denotes value at node 'x'.

→ Attributes may be numbers, strings, data types & references etc

Production	Semantic Rules.
$E \rightarrow E + T$	$E \cdot \text{Val} = E \cdot \text{Val} + T \cdot \text{Val}$ (Here $E \rightarrow$ is grammar symbol)
$E \rightarrow T$	$E \cdot \text{Val} = T \cdot \text{Val}$ $\text{Val} \rightarrow$ is attribute of $(+)$
$D \rightarrow TL$	$Linh = T \cdot \text{type}$.

↳ Semantic Rules have two Notations. (i) SDD (ii) SDT

Types of attribute:-

i) Synthesized Attribute = If a node takes value from its children node, then it is called as synthesized attribute.

(ii) Inherited Attributes: If a node takes value from its sibling or parent then it is called as inherited attributes. D \rightarrow TL

Production Semiotic Rules Categories

Eg:- $A \rightarrow BCD$ $\rightarrow C$ is inherited attribute

$c.i = A.i \rightarrow$ Here c is taking value from its parent A .

$C.i = B.i \rightarrow C$ is taking value from its sibling B .

$C \cdot i = D \cdot i \rightarrow C = D$ if and only if $i \neq 0$.

Production	Semantic Rules
$L \rightarrow E_n$	$E\text{-val} = E\text{-val}$
$E \rightarrow E_1 + T$	$E\text{-val} = E_1\text{-val} + T\text{-val}$
$E \rightarrow T$	$E\text{-val} = T\text{-val}$
$T \rightarrow T_1 * F$	$T\text{-val} = T_1\text{-val} * F\text{-val}$
$T \rightarrow F$	$T\text{-val} = F\text{-val}$
$F \rightarrow (E)$	$F\text{-val} = E\text{-val}$
$F \rightarrow \text{digit}$	$F\text{-val} = \text{digit.lexval}$

Fig: SPP for a simple desk calculator

Evaluating an SPP at Nodec of parse tree :-

Annotated parse tree:-

- A parse tree which contain values at each node is known as annotated parse tree.
- A parse tree is constructed in order to evaluate the attribute value at each node of parse tree.
- If an attributed is synthesized.
 - 1st evaluate the val attribute at all the children node.
 - evaluate the value attribute at parent node.
 - synthesized attributes, attributes are evaluated in bottom-up manner.

Production

$$A \rightarrow B$$

Semantic Rules

$$A.i = B.i$$

$$B.i = A.i + 1$$

These rules are circular, it is impossible to evaluate A.i without first evaluating B.i at some node.

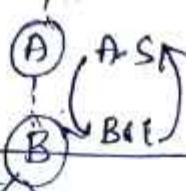


Fig: The circular dependency of A.i and B.i on one another

Eg1: Construct an SDD for simple weak calculator. Grammars and construct parse tree for $3 * 5 + 4n$.

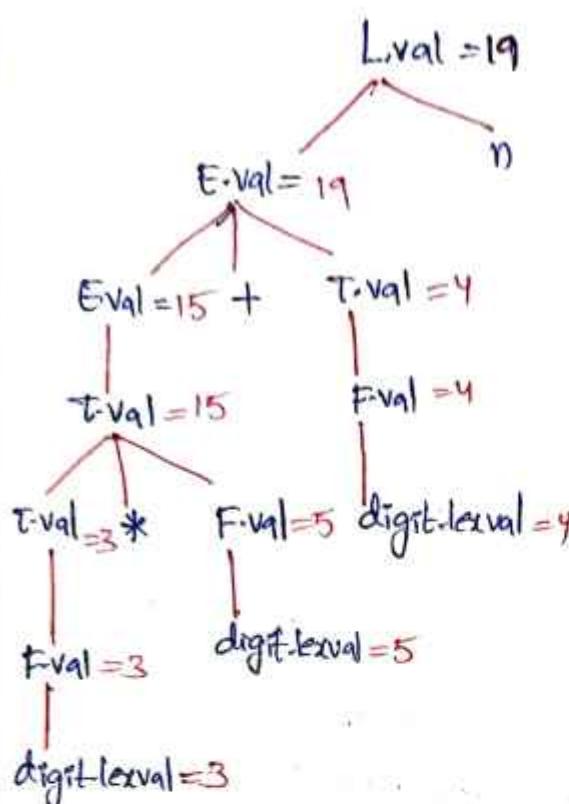


Fig: Annotated parse tree for $3 * 5 + 4n$

Grammars

production	semantic rule
$L \rightarrow E_n$	$L.val = E.val$
$E \rightarrow E+T$	$E.val = E.val + T.val$
$E \rightarrow T$	$E.val = T.val$
$T \rightarrow T_* F$	$T.val = T.val * F.val$
$T \rightarrow F$	$T.val = F.val$
$F \rightarrow (E)$	$F.val = E.val$
$F \rightarrow \text{digit}$	$F.val = \text{digit.level}$

Fig: SDD for simple weak calculator

Eg2: construct annotated parse tree for $a+b+c$.

- ↳ Annotated parse tree contains values at each node.
 ↳ To construct annotated parse tree, we have to perform top-down left to right traversing, if there is reduction execute corresponding action.

Eg2: production semantic rules

$$D \rightarrow TL \quad L.inh = T.type \rightarrow \text{list inheritance type } T$$

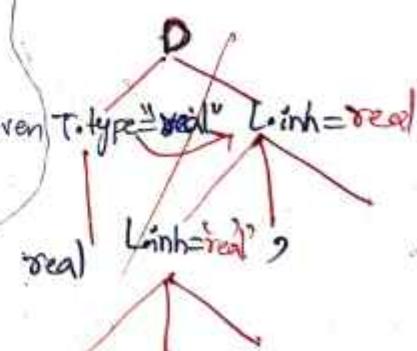
$$T \rightarrow \text{int} \quad T.type = "integers"$$

$$T \rightarrow \text{real} \quad T.type = "real"$$

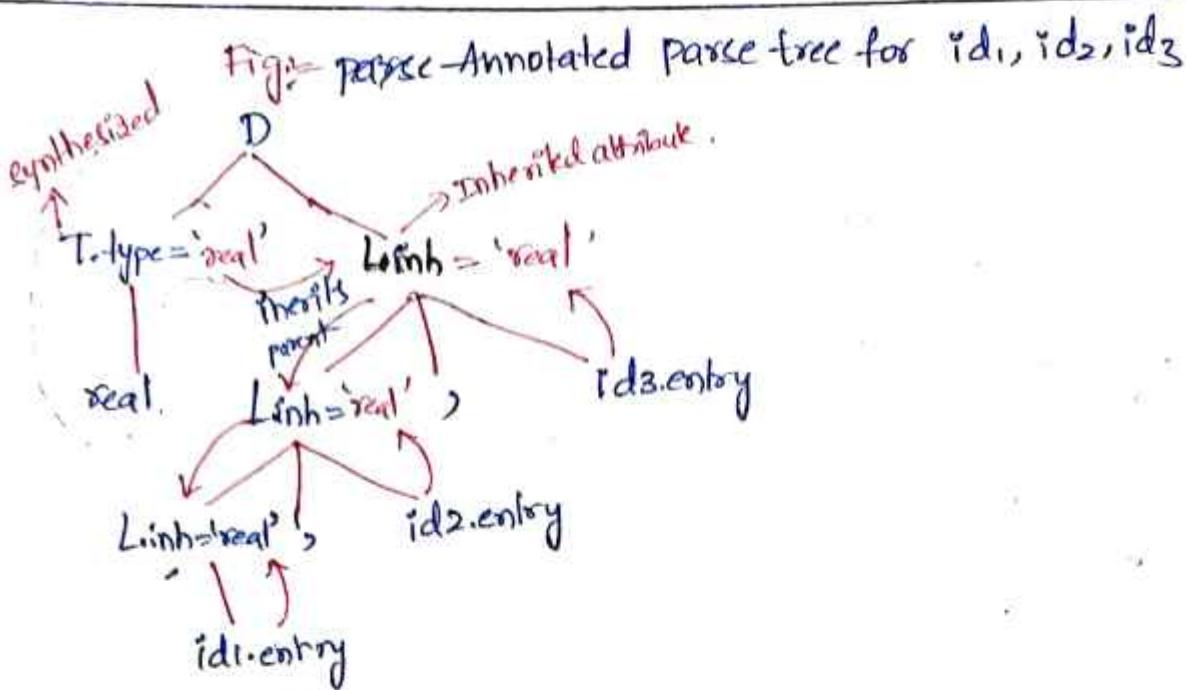
$$L \rightarrow \text{id}, id \quad L.inh = L.inh \rightarrow \text{here id needs to be given } T.type \\ \text{addtype(id.entry, L.inh)} \\ \text{addtype(id.entry, L.inh)}$$

$$L \rightarrow id$$

parse tree:



In dependency graph we have an edge directed edges



② Evaluating order SDD's:

→ Dependency graph is used to determine an evaluation order for the attribute given

Dependency Graph:= A directed graph that represents the interdependency b/w synthesized and inherited attribute at node in the parse tree

→ Dependency Graph represents the flow of information among the attributes in parse tree.

→ used to determine the evaluation order for attribute in a parse tree. (which semantic action should execute first)

→ An Annotated parse tree shows the value of attributes, a dependency graph determines how those values can be computed.

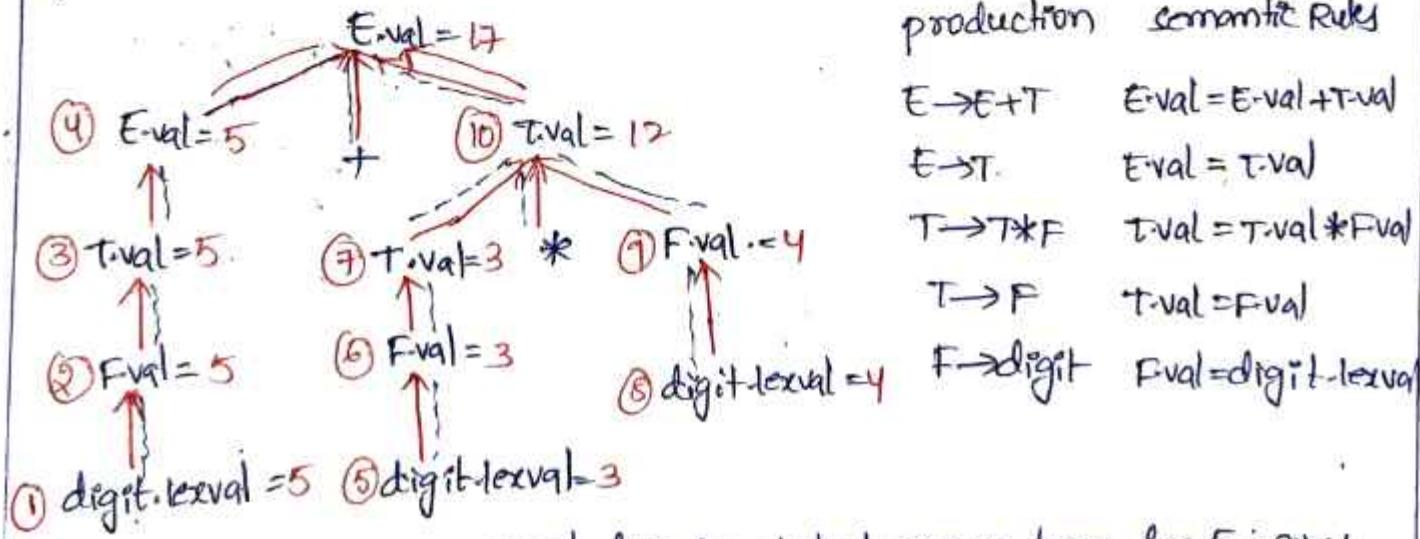
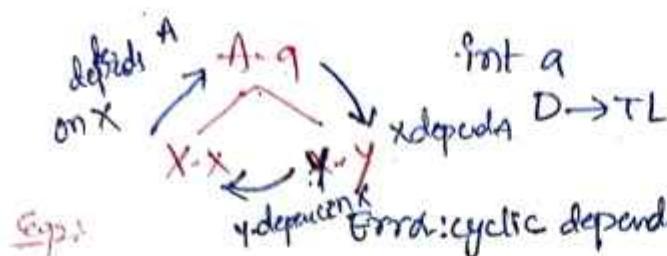


Fig: Dependency graph for Annotated parse tree for $5+3*4$

(3)

→ Edges in dependency graph show the interdependency b/w synthesized & inherited attributes at nodes in the parse tree.



→ Dependency graph cannot be cyclic
In DG - there is a edge backwards from the dependent-node to originating node.

Production Semantic Rules

$$T \rightarrow PT \quad T.inh = F.val$$

$$T.val = T.syn$$

$$T.inh = T.inh * F.val$$

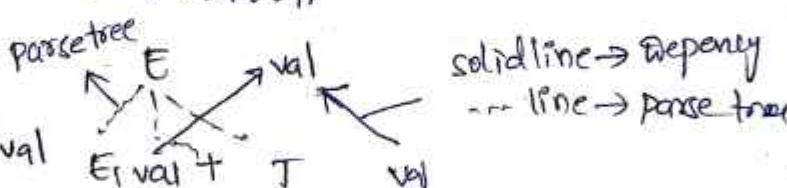
$$T.syn = T.syn$$

$$T.inh = T.inh$$

$$F \rightarrow digit \quad F.val = digit.lexval$$

Production Semantic Rule

$$E \rightarrow ET \quad E.val = E.val + T.val$$



Types of SDD's -

S-attributed Definition = $E.val$ is synthesized from $E.val$ and $T.val$.

* A SDD that use only synthesized attributes is called as S-attributed SDD.

$$\text{Eq} \quad A \rightarrow BCD$$

$$A.S = B.S$$

$$A.S = C.S$$

$$A.S = D.S$$

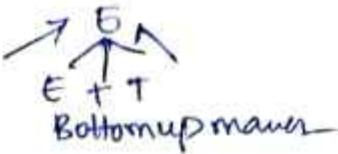
<u>Production</u>	<u>Semantic Rules</u>	<u>Eq</u> = <u>SDD for</u> <u>S-attributed</u> <u>Definition</u>
$L \rightarrow En$	$L.val = E.val$	
$E \rightarrow ET$	$E.val = E.val + T.val$	
$E \rightarrow T$	$E.val = T.val$	
$T \rightarrow T * F$	$T.val = T.val * F.val$	
$T \rightarrow F$	$T.val = F.val$	
$F \rightarrow (E)$	$F.val = E.val$	
$F \rightarrow digit$	$F.val = digit.lexval$	
$E \rightarrow ET \quad \{ \}$	$E.val = E.val + T.val$	$E.val = E.val + T.val$

* Semantic actions are placed at the end of the production.

$$\text{Eq} \quad E \rightarrow ET \quad \{ \quad E.val = E.val + T.val \quad \}$$

- It is also called as postfix "SDD"
- Attributes are evaluated with Bottom-up parsing

L-attributed SDD:



* A SDD that use both synthesis attributes & inherited attributes is called as L-attributed SDD.

* In L-attributed SDD, Inherited attributed is restricted to inherit from parent or left sibling only.

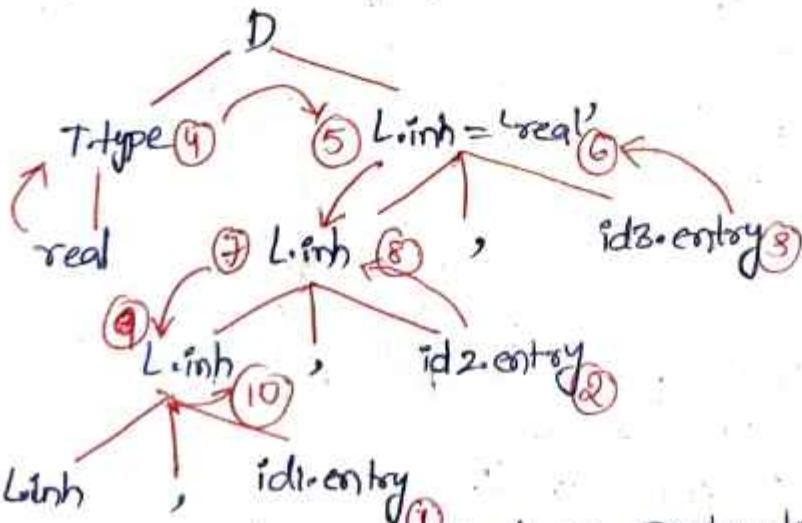
Eg: $A \rightarrow XYZ \{ y.S = A.S, y.S = X.S, y.S = Z.X \}$ production semantic action

* semantic actions are placed anywhere on R.H.S. eq= $E \rightarrow E + T \{ \} \rightarrow \{ \} E + T$

* evaluated attributes are evaluated by traversing parse tree depthfirst, left to right order

Production	Semantic Rules
$D \rightarrow TL$	$L.inh = T.type$
$T \rightarrow int$	$T.type = integer$
$T \rightarrow float$	$T.type = float$
$L \rightarrow L, id$	$L.inh = L.inh.addtype(id.entry, L.inh)$
$L \rightarrow id$	$L.inh = addtype(id.entry, L.inh)$

Eg:- SDD for simple type declaration for L-attributed SDD



Eg:- dependency graph for a declaration id, id_1, id_2

(3) Application of syntax directed Translation:

6

Syntax Directed Translation (SDT):

↳ SDT is a CFG together with semantic actions.

↳ Semantic actions are enclosed in '{', '}' → eq: $A \rightarrow ABC \rightarrow \{ABC\}$

↳ Semantic actions are placed at any where, on RHS of producers

↳ Semantic actions specifies in which order the expression is executed.

Eq 1: Production

$A \rightarrow B + C$

semantic action

{printf ('+') ; }

Eq 2: production : semantic action

$E \rightarrow E + T \quad \{ \text{printf} ('+') ; \} \quad ①$

$E \rightarrow T \quad \{ \quad \} \quad ②$

$E \rightarrow T * F \quad \{ \text{printf} ('*') ; \} \quad ③$

$T \rightarrow F \quad \{ \quad \} \quad ④$

$F \rightarrow \text{num} \quad \{ \text{printf} (\text{num-val}) ; \} \quad ⑤$

Fig 1: SDT for evaluation of expression.

Applications of SDT:-

Construction of syntax Tree:-

→ Syntax tree is an intermediate representation

→ The nodes in syntax tree is implemented by objects with suitable no. of fields

→ Each field will have an op-field that is the label of the node.

eq: $A \rightarrow ABC \rightarrow \{ABC\}$

$\rightarrow A \{ \quad \} BC$

$\rightarrow \{ \quad \} ABC$

Eq 2: production : semantic action

$E \rightarrow E + T \quad \{ E\text{-val} = E\text{-val} + T\text{-val} \}$

$T \quad \{ T\text{-val} = T\text{-val} \}$

$T \rightarrow T * F \quad \{ T\text{-val} = T\text{-val} * F\text{-val} \}$

$F \quad \{ T\text{-val} = F\text{-val} \}$

$T \rightarrow \text{num} \quad \{ E\text{-val} = \text{num-val} \}$

$2 + 3 * 4$

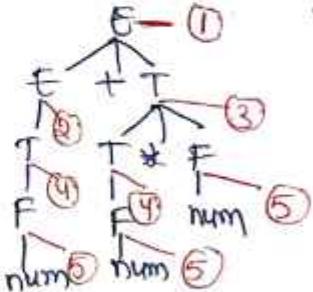


Fig 2: SDT

We can construct the syntaxtree by using the following functions.

(1) $\text{mknode}(\text{op}, \text{left}, \text{right})$

(2) $\text{mkleaf}(\text{id}, \text{entry} \rightarrow \text{symbol-table})$

(3) $\text{mkleaf}(\text{num}, \text{value})$

→ If the node is leaf, an additional field hold the lexical value for the leaf

$\text{leaf}(\text{OP}, \text{val})$ - create leaf object

→ If Node is an operator, Create an object with first field op and k additional for the k children $c_1 \dots c_2$

Eg:- Construct the syntaxtree for the following grammar for the expression $x * y - s + z$.

Step 1: Construct SDD for the given grammar

production Semantic Rule

$E \rightarrow E_1 + T$	$E\text{-node} = \text{mknode} (+', E_1\text{-node}, T\text{-node})$ (1)	$E\text{-node} = \text{newleaf}$ $(+, E_1\text{-node}, T\text{-node})$
$E \rightarrow E_1 - T$	$E\text{-node} = \text{mknode} (-', E_1\text{-node}, T\text{-node})$	
$E \rightarrow T$	$E\text{-node} = T\text{-node}$	
$T \rightarrow (E)$	$T\text{-node} = E\text{-node}$	
$T \rightarrow id$	$T\text{-node} = \text{mkleaf}(\text{id}, \text{id}\text{-entry})$ (2)	$T\text{-node} = \text{newleaf}(\text{id}, \text{id}\text{-entry})$
$T \rightarrow num$	$T\text{-node} = \text{mkleaf}(\text{num}, \text{num}\text{-val})$	

Step 2: - symbol

operation

x $P_1 = \text{mkleaf}(\text{id}, \text{entry}-x)$

y $P_2 = \text{mkleaf}(\text{id}, \text{entry}-y)$

$*$ $P_3 = \text{mknode}(*, P_1, P_2)$

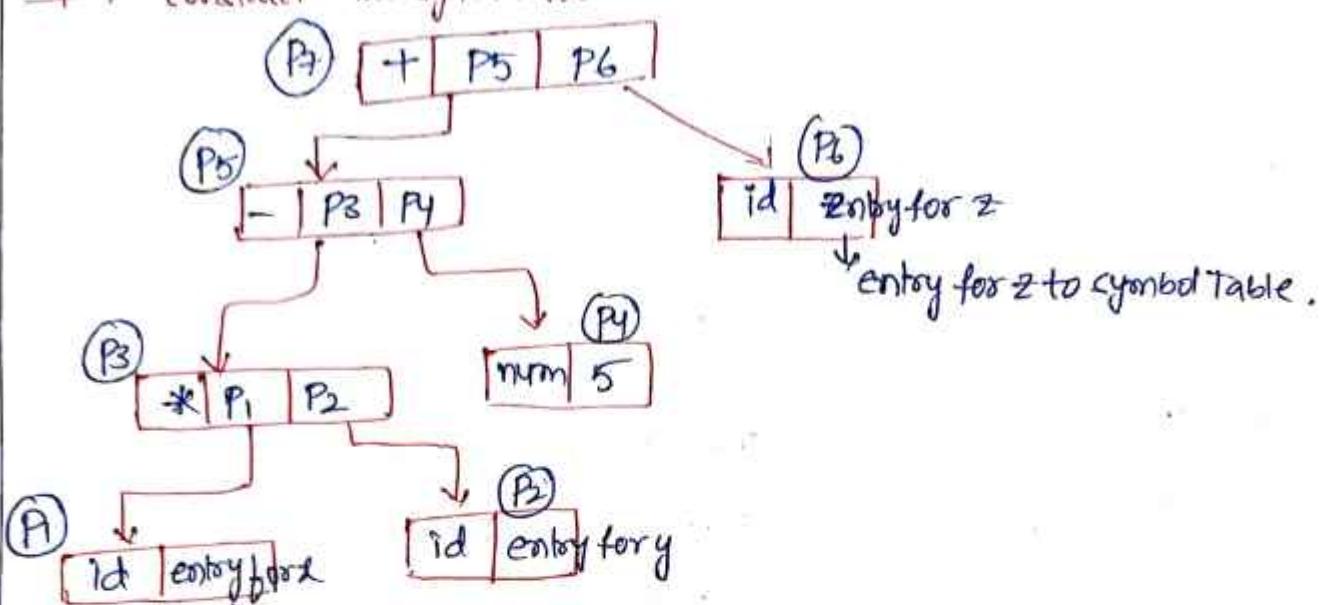
5 $P_4 = \text{mkleaf}(\text{num}, 5)$

$-$ $P_5 = \text{mknode}(-, P_3, P_4)$

$3 \quad P_6 = \text{mkleaf}(\text{id}, \text{entry}-3)$

$+ \quad P_7 = \text{mknoden}(+, P_5, P_6)$

Step 3:- construct the syntax tree.



Ex 2:- $a - 4 + c$

symbol operation

$a \quad P_1 = \text{mkleaf}(\text{id}, \text{entry}-a)$

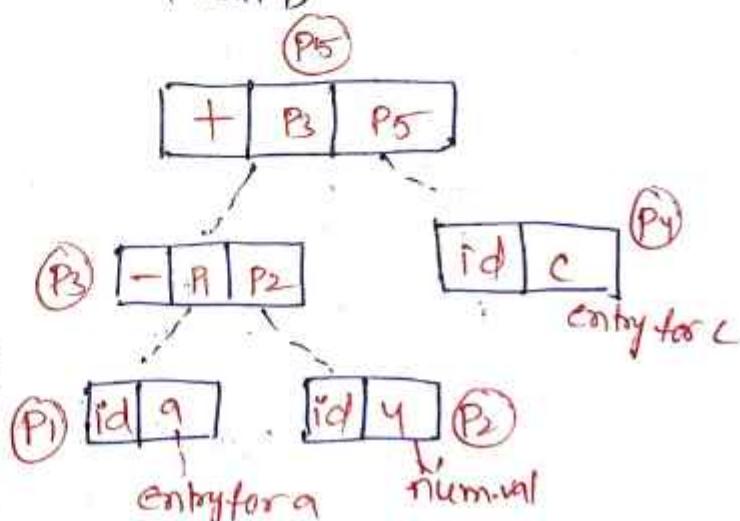
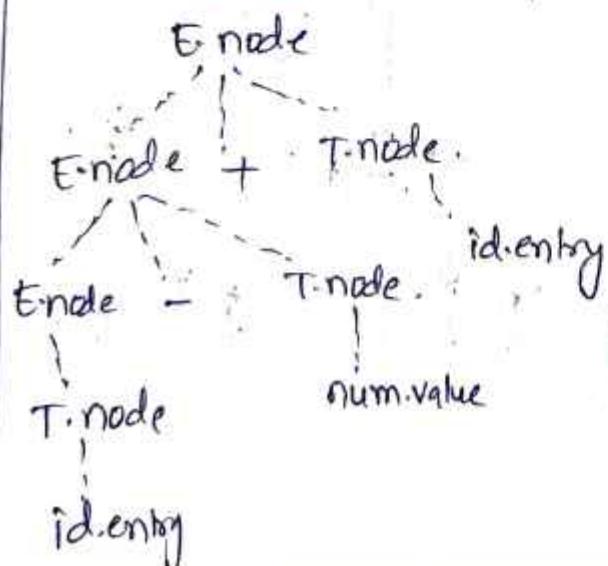
$4 \quad P_2 = \text{mkleaf}(\text{num}, 4)$

$- \quad P_3 = \text{mknoden}(-, P_1, P_2)$

$c \quad P_4 = \text{mkleaf}(\text{id}, \text{entry}-c)$

$+ \quad P_5 = \text{mknoden}(+, P_3, P_4)$

Constructing Syntax tree:-



Eg1: production

$E \rightarrow E^1$

$E^1 \rightarrow + T E^1$

$E^1 \rightarrow - T E^1$

$E^1 \rightarrow \epsilon$

$T \rightarrow (E)$

$T \rightarrow id$

$T \rightarrow num$

Semantic Rules

E-node = El-syn

T-inh = T-node

$E^1\text{-inh} = \text{newnode} (+, E^1\text{-inh}, T\text{-node})$

$E^1\text{-syn} = E^1\text{-syn}$

$E^1\text{-inh} = \text{newnode} (-, E^1\text{-inh}, T\text{-node})$

$E^1\text{-syn} = E^1\text{-syn}$

$E^1\text{-syn} = E^1\text{-inh}$

T-node = E-node

T-node = newleaf (id, entry-id)

T-node = newleaf (num, num-val)

Eg1: $a - 4 + c$

symbol operation

a $P_1 = \text{mkleaf}(\text{id}, \text{entry-}a)$

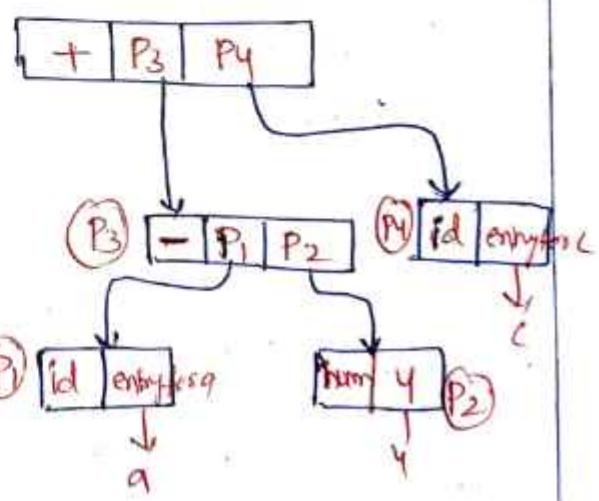
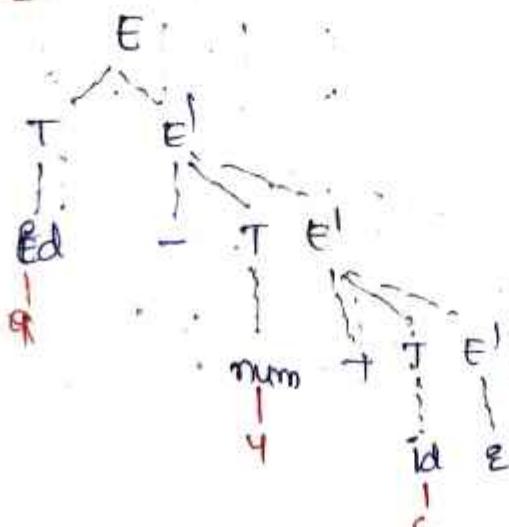
4 $P_2 = \text{mkleaf}(\text{num}, 4)$

- $P_3 = ('-', P_1, P_2)$

c $P_4 = \text{mkleaf}(\text{id}, \text{entry-}c)$

+ $P_5 = \text{mknode} ('+', P_3, P_4)$

Syntax tree:



Syntax-Directed Translation schemes:-

- (3) Syntax-Directed Translation schemes:-
- A SDT is a context free Grammar combined with semantic actions.
- Semantic actions are also called as program fragments.
- Semantic actions are embedded with production body.
- Any SDT can be implemented by constructing parse tree and then performing the actions in a left-to-right depth-first order.

SDT's are used to implement two important classes of SDD's

- 1) The underlying grammar is LR-parsable, & the SDD is L-attributed.
- 2) The underlying grammar is LL-parsable, and the SDD is L-attributed.

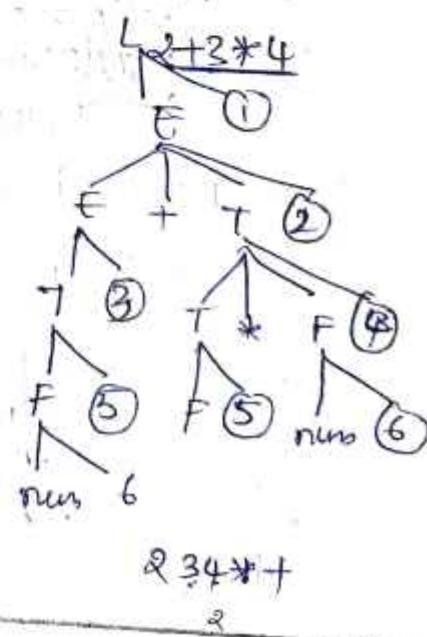
Postfix Translation schemes:-

- It is used to convert infix expression to postfix expression.
- Infix expression - operator appears b/w operands
- Postfix expression - operator appears after the operands.

Eg:-

Production	semantic actions
$L \rightarrow E_n$	{ print(Eval); }①
$E \rightarrow E + T$	{ print('+'); }②
$E \rightarrow T$	{ }③
$T \rightarrow T * F$	{ print('*'); }④
$T \rightarrow F$	{ }⑤
$F \rightarrow \text{num}$	{ print(numval); }⑥

fig:- SDT for desc calculator



2nd method to convert the postfix infix expression to postfix expression

$$\text{Eg:- } E \rightarrow E + T \mid T \quad \text{for } = 1 + 2 + 3$$

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

(here we are having only operation, so that we need to check the left recursion in the grammar)

→ if Grammar contains left recursion we need to convert the eliminate left recursion from the grammar

production LDT

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

$$E \rightarrow T \quad \{ \}$$

$$T \rightarrow \text{num} \quad \{ \text{print num-value} \}$$

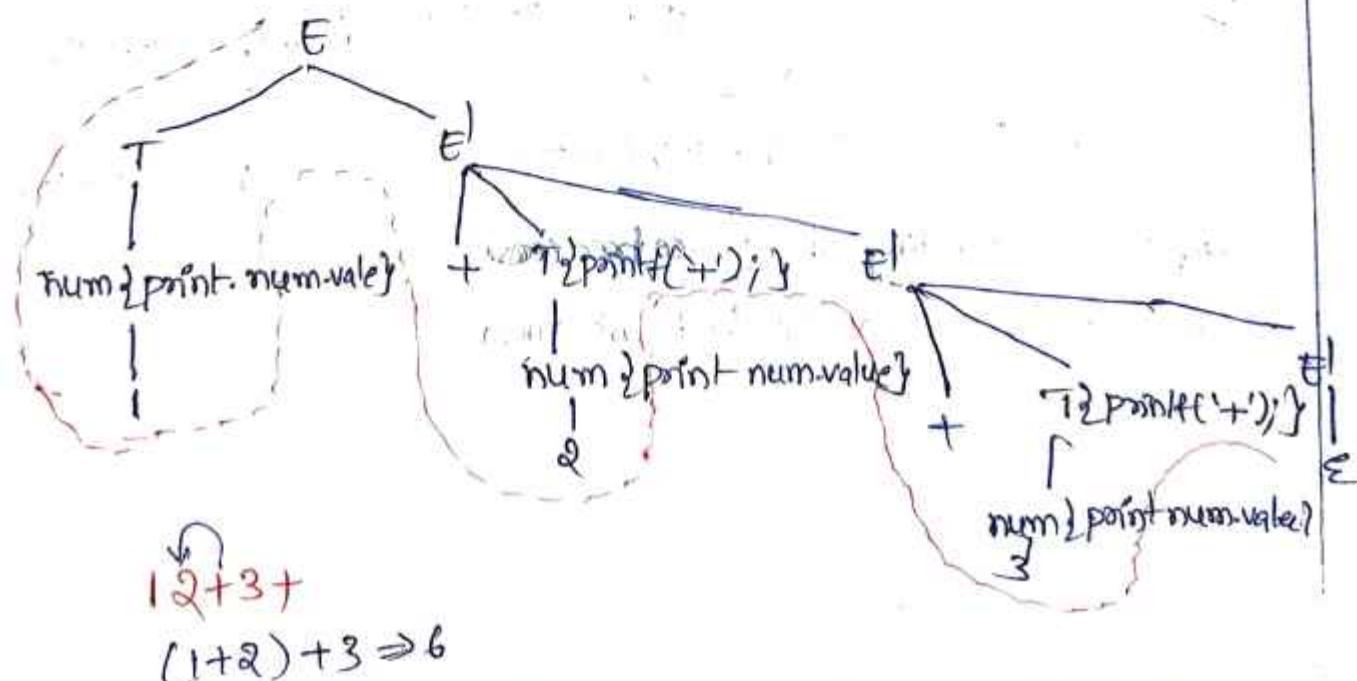
LDT for given grammar

so, the given Grammar contains left recursion

$$\begin{array}{c} E \rightarrow E + T \mid T \\ \cancel{\downarrow} \quad \cancel{\downarrow} \\ T \rightarrow \text{num} \end{array} \quad \frac{E \rightarrow E + T \{ \text{printf}('+''); \} \mid T}{\cancel{\downarrow} \quad \cancel{\downarrow} \quad \cancel{\alpha}} / \beta$$

$$\left. \begin{array}{c} \text{The grammar} \\ \text{after eliminating} \\ \text{the left recursion} \end{array} \right\} \begin{array}{l} E \rightarrow TE' \\ E' \rightarrow +T \{ \text{printf}('+''); \} E' \\ E' \rightarrow \epsilon \\ T \rightarrow \text{num} \quad \{ \text{print num-value} \} \end{array}$$

$$\begin{array}{l} A \rightarrow \alpha x \mid \beta \\ A \rightarrow \beta A \\ A \rightarrow \alpha A \mid \epsilon \end{array}$$



Q1 = Give translation scheme that convert infix expr to postfix expression for the following grammar and also generate annotated parse for input string 2+6+1

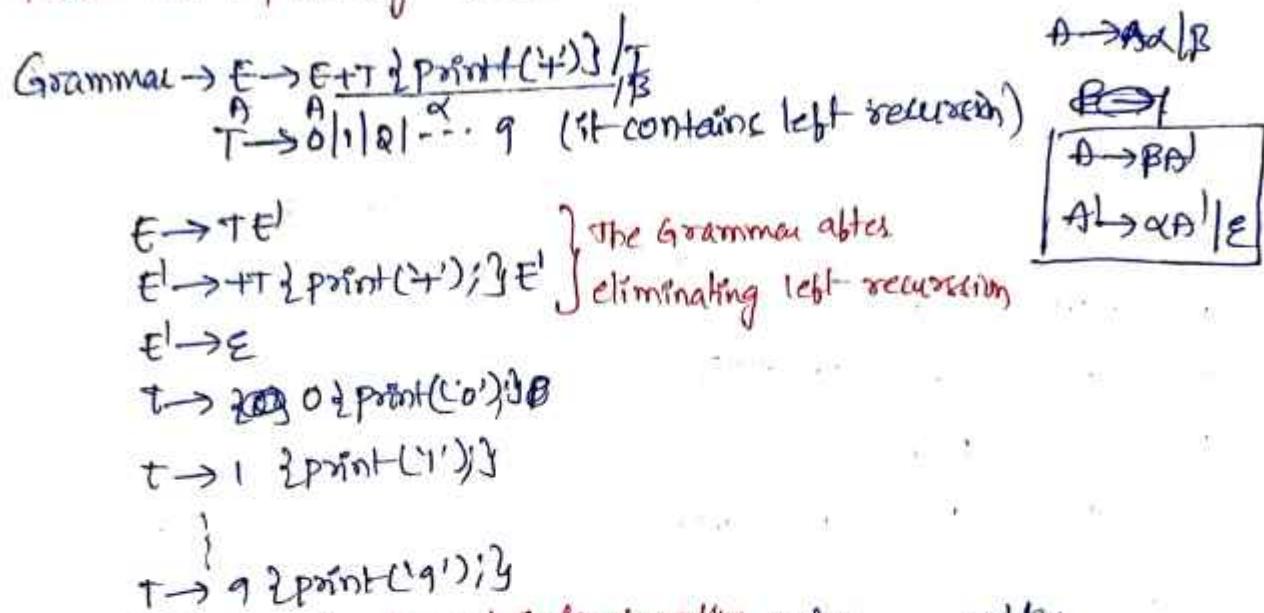
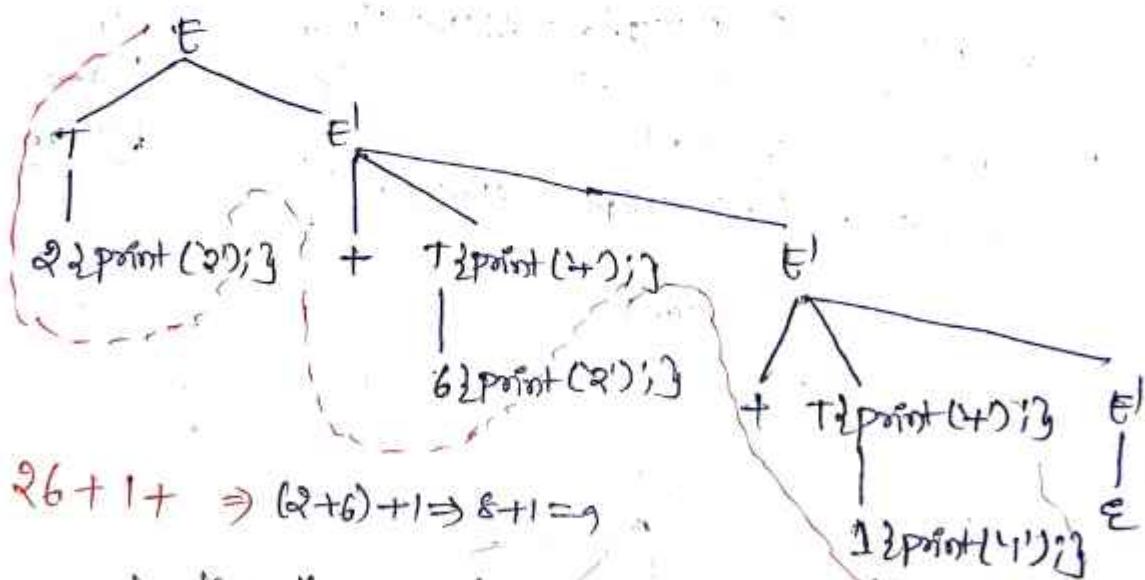


Fig: SDT to convert infix to postfix infix postfix

Annotated Parse tree = The given i/p string $2+6+1 \Rightarrow 26+1+$



→ After constructing the parse tree, traverse the parse from top-down and left to right manner

$L \rightarrow \epsilon_n \{ \text{print}(\text{val}); \}$

$E \rightarrow E + T \{ E\text{-val} = E\text{-val} + T\text{-val}; \}$

$E \rightarrow T \{ E\text{-val} = T\text{-val}; \}$

$T \rightarrow T_1 * F \{ T\text{-val} = T_1\text{-val} * F\text{-val}; \}$

$T \rightarrow F \{ T\text{-val} = F\text{-val}; \}$

$F \rightarrow (E) \{ F\text{-val} = E\text{-val}; \}$

$F \rightarrow \text{digit} \{ F\text{-val} = \text{digit-lexval}; \}$

Fig :- postfix SDT's for implementing the desc calculator

parser - stack implementation of postfix SDT's :-

→ postfix SDT's can be implemented during LR parsing by executing the actions when the reductions occur.

→ The grammar symbols of each grammar

→ The attributes of each grammar symbols can be placed in stack during the parsing.

→ The parser stack contain records with fields for a grammar symbol.

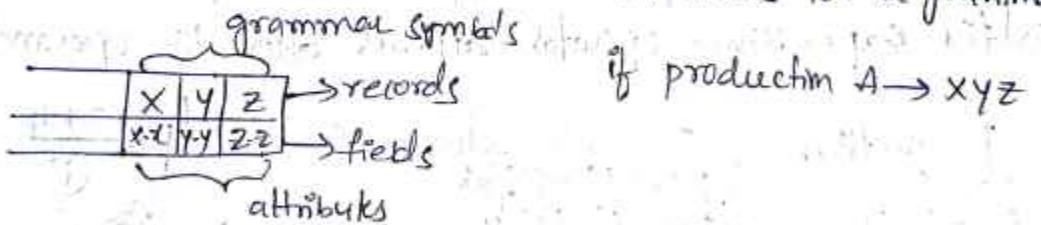
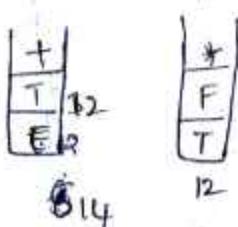
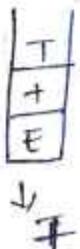


fig :- parser stack with a field for synthesized attributes

or $E \rightarrow E + T \{ \text{print}(+); \}$ or $E \rightarrow E + T \cdot \quad 2+3*4$



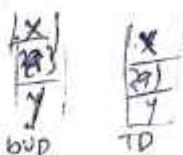
production	Actions.
$L \rightarrow E_n$	{ print (stack [top-1], val); top=top-1 }
$E \rightarrow E + T$	{ stack [top+2].val = stack [top-2].val + stack [top].val; top=top-2; }
$E \rightarrow T$	
$T \rightarrow T * F$	{ stack [top-2].val = stack [top-2].val * stack [top].val; top=top-2; }
$T \rightarrow F$	
$F \rightarrow (E)$	{ stack [top-2].val = stack [top-1].val; top=top-2; }
$F \rightarrow \text{digit}$	

Fig: Implementing the desc calculator on bottom-up-parsing stack.

SDT's with Actions Inside productions:-

→ An actions may be placed at any position ~~or~~ right within the body of production.

e.g.: $B \rightarrow x \{ a \} y$.



The action "a" is executes after recognizing the ~~or~~ "x".

- If the parse is bottom-up then we perform action "a" when "x" is appears on the top of the stack.
- If the parse is top-down, we perform a just before expanding the "y"

Any SCDT can be implemented as follows.

- 1) Ignoring the actions, parse the P/P & produce a parse tree as a result.
- 2) examine each node and add additional node for corresponding action
- 3) perform the preorder traversal of the tree, and as soon as a node is labeled by action is visited, perform that action.

b

Eg → parse tree for expression $3 * 5 + 4$ with actions itself
we get if we visit the nodes in preorder, we get the prefix form of expression $+ * 3 5 4$

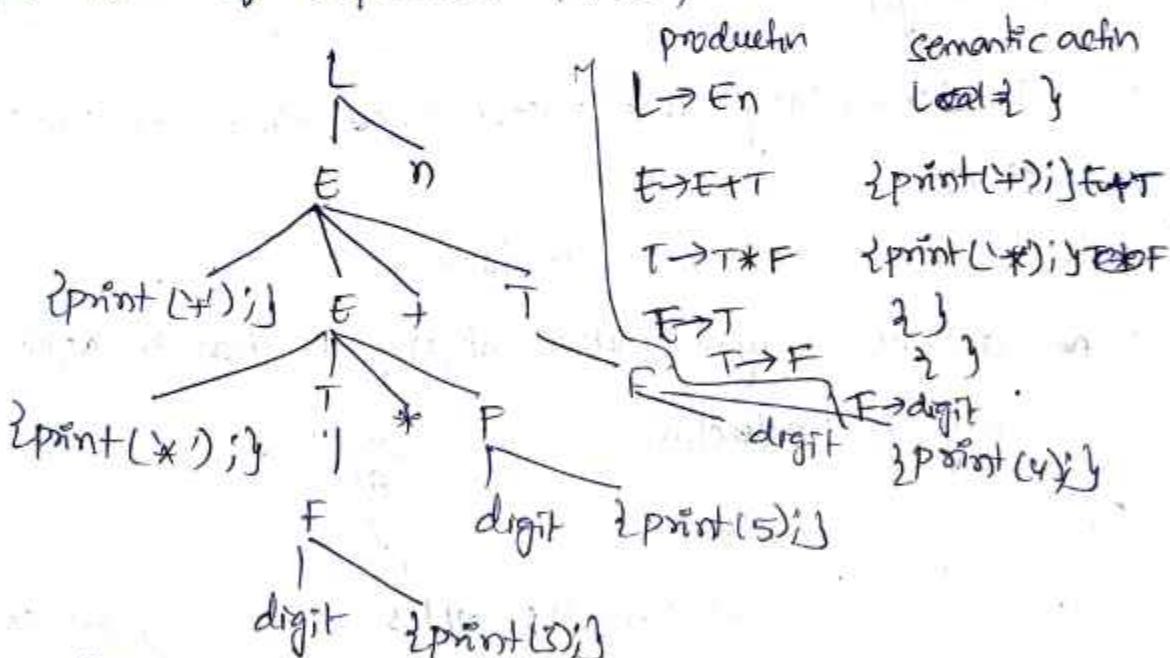


Fig: Parse tree with actions embedded.

Intermediate scale generations

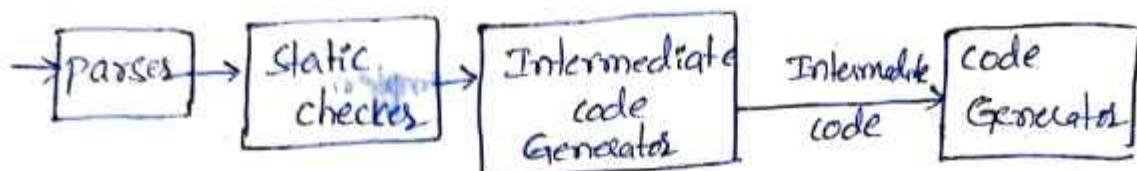


Fig: logical structure of a compiler's Frontend.

- Intermediate code is used to translate the sourcecode into the machine-code.
 - In the above figure parsing, static checking and Intermediate-code generation are done sequentially.
 - Static type checking includes type checking, which ensures that operators are applied to compatible operand.
 - ICG receives from its predecessor phase & semantic analysis phase
 - It takes i/p in the form of an annotated syntax tree.
 - In process of translating a source program into target code compiler may construct a sequence of intermediate representation



- Syntax trees are high level representation
 - A low level representation is suitable for machine-dependent tasks like register allocation and instruction selection.

variants of syntax tree :-

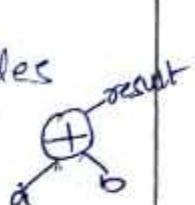
Directed Acyclic Graph (DAG) for expressions:-

→ DAG is a datastructure used for implementing transformations on basic blocks.

→ same nodes in

→ DAG represent the structure of a basic block.

- In DAG internal nodes represent ^{operators} and leaf nodes represent identifiers, constants.
- Internal nodes represent the result of expression



→ The only difference b/w syntax tree and DAG is, in DAG a node has more than one parent.

Applications of DAG:-

* Determining the common subexpression.

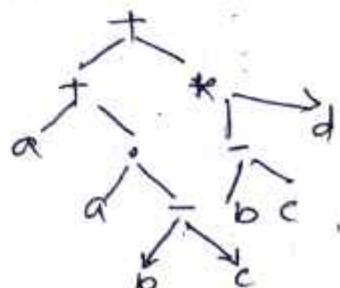
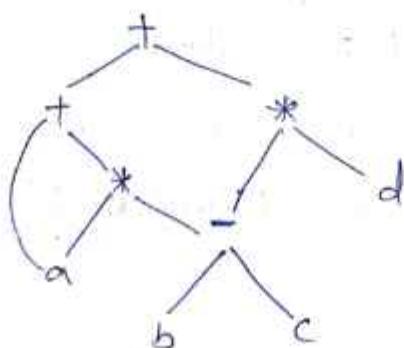
* Determining which names are used inside the block and computed outside the block.

* Determining which statement of the block could have their computed value outside the block.

• by Eliminating common subexpressions it simplify the code.

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

Syntax tree



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

production	semantic Rule
$E \rightarrow E_1 + T$	$E\text{-node} = \text{new Node}(' + ', E_1\text{-node}, T\text{-node})$
$E \rightarrow E_1 - T$	$E\text{-node} = \text{newNode}(' - ', E_1\text{-node}, T\text{-node})$
$E \rightarrow T$	$E\text{-node} = T\text{-node}$
$T \rightarrow (E)$	$T\text{-node} = E\text{-node}$
$T \rightarrow \text{id}$	$T\text{-node} = \text{new leaf}(\text{id}, \text{id}, \text{entry})$
$T \rightarrow \text{num}$	$T\text{-node} = \text{new leaf}(\text{num}, \text{num}, \text{val})$

Fig: SDD for to produce Syntax tree & WAG's

$$1) P_1 = \text{leaf}(\text{id}, \text{entry-}a)$$

$$P_2 = \text{leaf}(\text{id}, \text{entry-}a) = p_1 \quad \text{Eq 2: construct WAG for}$$

$$P_3 = \text{leaf}(\text{id}, \text{entry-}b)$$

$$\begin{aligned} 1) a &= b+c \\ 2) b &= a-d \end{aligned}$$

$$P_4 = \text{leaf}(\text{id}, \text{entry-}c)$$

$$P_5 = \text{Node}(' + ', P_3, P_4)$$

$$P_6 = \text{Node}(' + ', P_1, P_5)$$

$$P_7 = \text{Node}(' + ', P_1, P_6)$$

$$P_8 = \text{leaf}(\text{id}, \text{entry-}b) = P_3$$

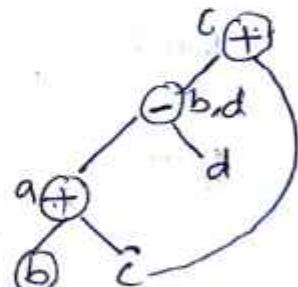
$$P_9 = \text{leaf}(\text{id}, \text{entry-}c) = P_4$$

$$P_{10} = \text{Node}(' - ', P_3, P_4) = P_5$$

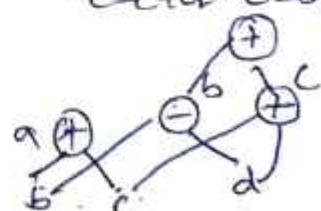
$$P_{11} = \text{leaf}(\text{id}, \text{entry-}d)$$

$$P_{12} = \text{Node}(' * ', P_5, P_{11})$$

$$P_{13} = \text{Node}(' + ', P_7, P_{12})$$



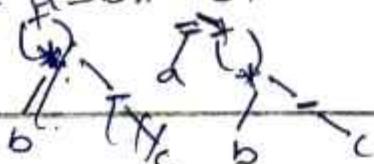
$$\begin{aligned} \text{Eq 3: } 1) a &= b+c \\ 2) b &= b-d \\ 3) c &= c+d \\ 4) e &= b+c \end{aligned}$$



The value number Fig: steps for constructing the WAG.

$$\text{Eq 2: } a = b * -c + b * -c$$

$$\text{Eq 3: } a = (a * b + c) - (a * b + c)$$



The Value Number methods for constructing WAG's

- Nodes of syntax tree or WAG are stored in array of records.
- Each row of array represent one record.(node)
- In each record first field is operation code, indicating the label of the node.
- leaves has the leaves one additional field which holds the lexical value.
- Interior nodes have two additional fields indicating left and right children.

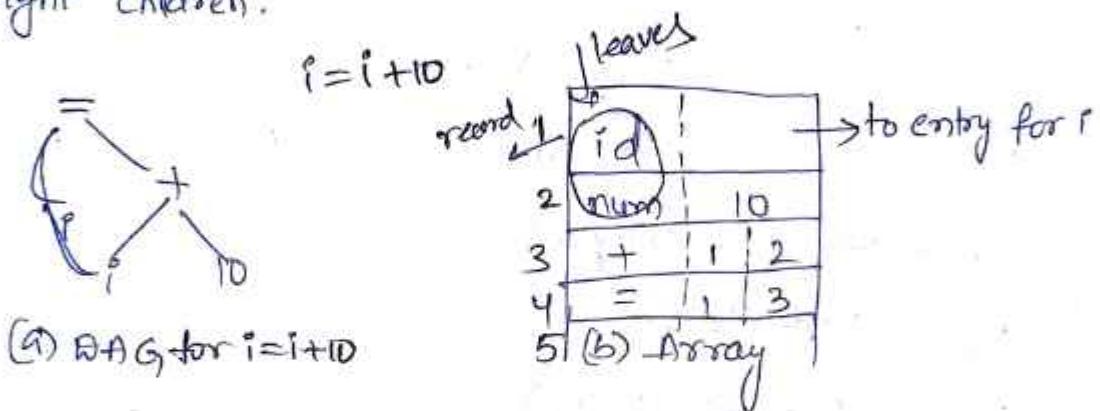


Fig: Nodes of a WAG for $i = i + 10$ allocated in an array.

- The array index is used for reference a node rather than a pointer.
- Initially the array is empty.
- First it searches for $\langle id, \text{lexical} \rangle$, if it is not there, we will make new record and so on.
- If already, record is present, it is just used for further records.
- In the array, we refer to node by giving integer index of the record, for that node within the array this integer is called Value Number. $\langle \text{op, value-num of left, value-num of right} \rangle$ is also called as signature of node.

Drawbacks:

- Search options: It takes time for every New/old record to search.
- To overcome this we are using hash functions, in which the nodes are put into "buckets" (hash table)
 - These buckets have only few nodes.
 - hashtable is a data structure that supports the dictionaries.
 - Dictionaries are used to insert & delete elements of a set.
 - Dictionaries are used to determine whether a given element is currently in the set. This is done.
 - It search the elements in less time and independent of the size of set.

To construct a has table for node.. of a DAG, use hashfunction
"h" is used that computes the index of bucket.

- ~~hashtable~~ The bucket index $h \in [0, l]$
- bucket can be implemented as linked list.
- An array indexed by has value, hold the bucket headers, each of which point to first cell of list

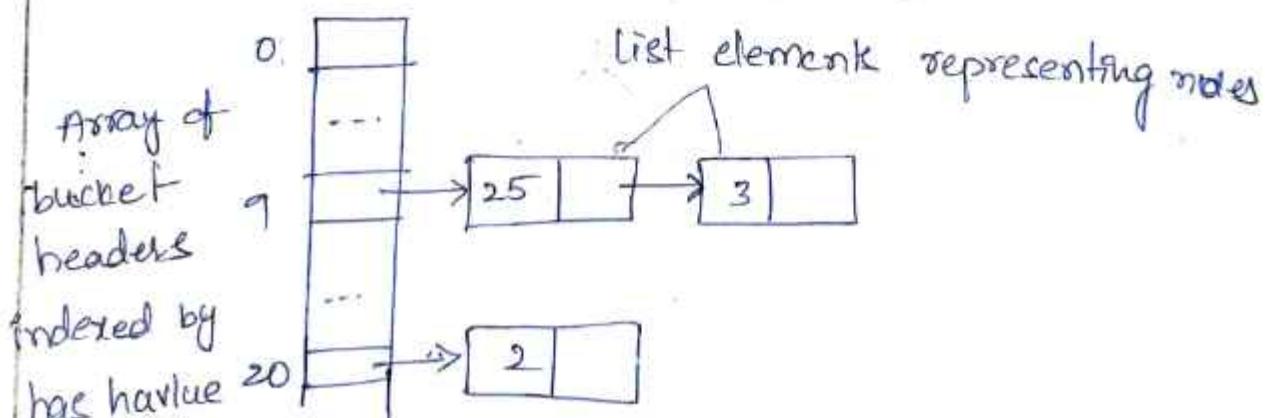


Fig: Data structure for searching buckets

8) Three Address code:

Intermediate code is three types

- ① Syntax tree representation
- ② postfix notation
- ③ Three address code.

In three-address code

- ① Each instruction should contain atmost 3 addresses
- ② Each instruction should contain 1 opertaor on RHS.

Eg:- source language expression $x+y+z$ ie converted into sequence of 3-address instructions.

$$\begin{aligned} t_1 &= y * z \\ t_2 &= x + t_1 \end{aligned}$$

$t_1, t_2 \rightarrow$ compiler generated temporary variables.

→ 3-address code is linearized representation of Syntax tree or RAG

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

$$\begin{aligned} t_1 &= b - c \\ t_2 &= a * t_1 \\ t_3 &= a + t_2 \\ t_4 &= t_3 * d \\ t_5 &= t_2 + t_4 \end{aligned}$$

Fig.(a) Three address code

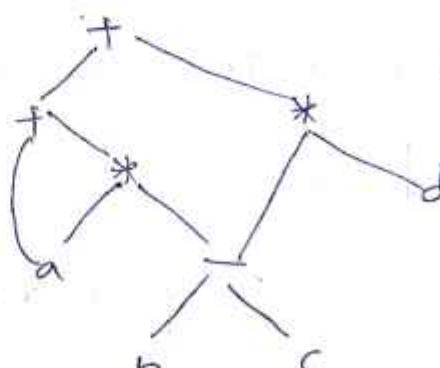


Fig.(b) = RAG

→ Three - address code is represented in 3 ways

- ① Quadruple
- ② Triple
- ③ Indirect Triples

Types of 3-address code:
Addressees and Instructions:

→ 3-address code can be implemented by using records with fields for the addresses

→ records are called quadruples and triples.

The address can be one of the following:

- A name → Source program names are addresses in 3-address code and names are replaced by pointer to its symbol table entry
- A constant
- compiler generated temporary variables.

List of the common three-address instruction forms:

(i) Assignment instruction $\rightarrow x = y \text{ op } z$, where x,y,z - addresses

(ii) $\rightarrow x = \text{op } y$ where op - is unary operation ($- , \sim$)

(iii) copy instruction $\rightarrow x = y$

(iv) unconditional jump $\rightarrow \text{goto } L$

(v) conditional jump $\rightarrow \text{if } x \text{ goto } L$
 $\rightarrow \text{if false } x \text{ goto } 'L'$.
 $\rightarrow \text{if } x \neq 0 \text{ op } y \text{ goto } L$

(vi) procedure call $\rightarrow y = \text{call } p, n$
 $\quad \quad \quad \text{return } y$

p → is the address of starting of line of procedure p

n → argument address.

y → return value

(vii) Indexed copy instruction $\rightarrow x = y[i] \quad \left. \begin{array}{l} x, y, i \text{ are the variables} \\ x[i] = y \end{array} \right\}$

(viii) Address and pointer assignment $x = & y$
 $x = * y$
 $* x = y$

Three address code

Quadruples:

3-address code is implemented as objects (as) records with fields for the operators and operand.

→ Quadruple has 4-fields (1) op (i) arg1 (ii) arg2 @result.

Ex: * instruction like $x=y$ or $z=-y$ do not use arg2

* operator like param use neither arg2 nor result

* conditional & unconditional jumps put the target label in result.

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

$$t_1 = -c$$

$$t_2 = b * t_1$$

$$t_3 = -c$$

$$t_4 = b * t_3$$

$$t_5 = t_2 + t_4 \quad \left. \begin{array}{l} t_3 = b * t_1 \\ t_4 = t_2 + t_3 \end{array} \right\}$$

$$a = t_5$$

(a) Three address code

	OP	arg1	arg2	result
(0)	-	c		t ₁
(1)	*	b	t ₁	t ₂
(2)	-	c		t ₃
(3)	*	b	t ₃	t ₄
(4)	+	t ₂	t ₄	t ₅
(5)	=	t ₅		a

(b) Quadruples

→ The disadvantage of Quadruples is too many temporary variables are needed, it require more amount of memory.

→ In order to overcome we are using Triples.

Triples: Triples has only three fields (1) op (2) arg₁ (3) arg₂

$$\text{Eq.} : a = b * -c + b * -c$$

	OP	-Arg1	Arg2
(0)	-	c	
(1)	*	b	(0)
(2)	-	c	
(3)	*	b	(2)
(4)	+	(1)	(3)
5	=	a	(4)

$$\begin{aligned}
 t_1 &= -c \\
 t_2 &= b * t_1 \\
 t_3 &= -c \\
 t_4 &= b * t_3 \\
 t_5 &= t_2 + t_4
 \end{aligned}$$

(b) **Triples** representation of $a = b * -c + b * -c;$

→ using triples we refer results of an operation by its position, rather than by an explicit temporary variable.

Indirect Triples:

Indirect triples consist of listing of pointers to triples, rather than a listing of triples themselves.

→ with triples the result of an operation is referred to by its position, so moving an instruction may require us to change all references to that result. This problem does not occur with indirect triples.

instruction

35	(0)
36	(1)
37	(2)
38	(3)
39	(4)
40	(5)
...	

	OP	arg1	arg2
0	-	c	
1	*	b	(0)
2	-	c	
3	*	b	(2)
4	+	(1)	(3)
5	=	a	(4)

Fig: Indirect triple representation of three address code

static single assignment form (SSA)

SSA is a special case of 3-address code. SSA is an intermediate representation that facilitates certain code optimizations. → In SSA each assignment to a variable should be specified with distinct names.

$$q := p = a + b$$

$$P_1 = a + b$$

$$q = p - c$$

$$q_1 = P_1 - c$$

$$p = q * d$$

$$P_2 = q_1 * d$$

$$p = c - d$$

$$P_3 = c - P_2$$

$$q = p + q$$

$$q_2 = P_3 + q_1$$

(a) Three-address code

(b) static single assignment form.

Ex: Intermediate program in three-address code SSA

→ The least no. of temporary variables required to create 3-address code in SSA.

* A variable can only be initialized one in L-H-S

* A variable which is initialized in L-H-S could only used R-H-S

⑨

Control Flow:

Simple if, if-else, else-if, switch, for, while, do-while.
The translation of statements such as if-else-statement and
while-statement is tied with translation of Boolean Expressions.

→ Boolean Expressions are used to.

- i) change the flow of control → Boolean expo are used as conditional expressions that alter the flow of control.
- ii) compute the logical values → for e.g. if (E) S,

Boolean Expressions =

→ A Boolean expression can represent true & false as value.

Boolean Expressions are composed of boolean operators

& $\&$, ||, and !

→ Boolean Expressions are generated by the following grammar

$$B \rightarrow B \& B \mid B \& B \mid !B \mid (B) \mid E \text{ rel } E \mid \text{true} \mid \text{false}$$

→ AND (or) OR are left associative

→ "NOT" has higher precedence than AND & or

short circuit code = (Jumping code)

In short-circuit code, the boolean operators &&, || and ! are translate into jumps.

→ In short-circuit code the and argument is evaluated only if 1st argument does not suffice to determine the value of expression.

e.g:- if ($x < 100 \& x > 200 \&& x != y$) $x = 0;$

In this translation the BE is true if control reaches label l_2 . If the expression is false, control immediately to l_1 , skipping l_2 and the assignment $x = 0$.

if $x < 100$ goto L₂
 if false $x > 200$ goto L₁
 if false $x_1 = y$ goto L₁
 L₂: $x = 0$

$\Leftarrow q_2: (x == y \text{ || } y == z) \dots$

fig: Jumping code.

Flow of control statements:

→ b Translation of Boolean expressions control - statements into three - address code.

Grammars

$S \rightarrow \text{if } (B) S_1$	$(\delta) S \rightarrow \text{if } (B) \text{then } S_1$
$S \rightarrow \text{if } (B) S_1 \text{ else } S_2$	↳ condition (δ) Boolean expression
$S \rightarrow \text{while } (B) S_1$	

- Grammars for simple if, if-else, while statements

① $S \rightarrow \text{if } (B) \text{ then } S_1$ (B) is evaluated 1st

code for simple if:

semantic Rule:

B.code	B.true . B.false	B.true = newlabel() B.false = S₁.next = S.next B.false = S ₁ .next = S.next S.code = B.code label(B.true) S ₁ .code
B.true	S ₁ .code	Intermediate code: S.code = B.code label(B.true) S ₁ .code
B.false	S.next	

fig: SDD for simple if - statement

if () {
 \downarrow Label, newlabel() function produce three address code
 \downarrow for B.true.
 \downarrow
 S}

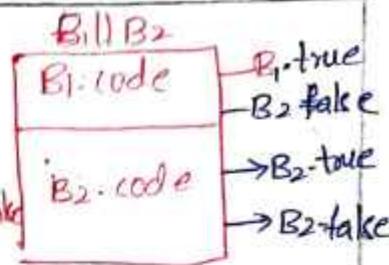
control flow translation of Boolean expression (d)
Three address code for Boolean expression (d) CDR (d) CDT for Boolean expression

production

Semantic Rules

$$B \rightarrow B_1 \parallel B_2$$

$\begin{cases} B_1.\text{true} = B.\text{true}; \\ B_1.\text{false} = \text{newlabel}(); \\ B_2.\text{true} = B.\text{true}; \\ B_2.\text{false} = B.\text{false}; \end{cases}$

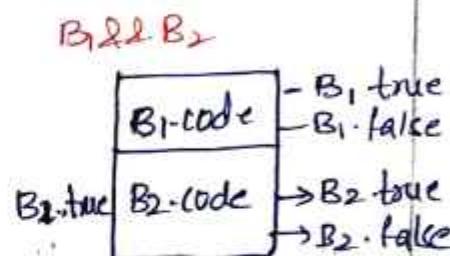


Intermediate code:

$$B.\text{code} = B_1.\text{code} \parallel \text{label}(B_1.\text{false}) \parallel B_2.\text{code}$$

$$B \rightarrow B_1 \& \& B_2$$

$\begin{cases} B_1.\text{true} = \text{newlabel}(); \\ B_1.\text{false} = B.\text{false}; \\ B_2.\text{true} = B.\text{true}; \\ B_2.\text{false} = B.\text{false}; \end{cases}$



$$B.\text{code} = B_1.\text{code} \parallel (\text{label}(B_1.\text{true})) \parallel B_2.\text{code}$$

$$B \rightarrow !B_1$$

$\begin{cases} B_1.\text{true} = B.\text{false}; \\ B_1.\text{false} = B.\text{true}; \\ B.\text{code} = B_1.\text{code} \end{cases}$

$$B \rightarrow \text{true}$$

$B.\text{code} = \text{gen}(\text{'goto'} B.\text{true});$

$$B \rightarrow \text{false}$$

$B.\text{code} = \text{gen}(\text{'goto'} B.\text{false});$

$$B \rightarrow E_1 \text{ relop } E_2$$

$B.\text{code} = E_1.\text{code} \parallel E_2.\text{code}$

$\parallel \text{gen}(\text{'if'} E_1 \text{ relop } E_2 \text{ 'goto' } B.\text{true})$

$\parallel \text{gen}(\text{'if'} E_1 \text{ relop } E_2 \text{ goto })$

$\parallel \text{gen}(\text{goto } E.\text{false})$

$\begin{cases} E_1.\text{true} = E.\text{true}; \\ E.\text{false} = E.\text{false}; \end{cases}$

$E.\text{code} = E_1.\text{code};$

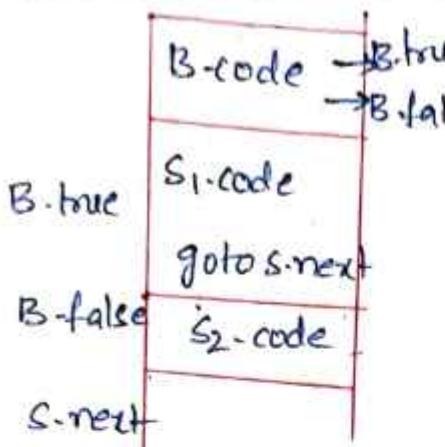
Ex. If back

if a < b goto E.true
goto E.false

$$E \rightarrow (E_1)$$

$S \rightarrow \text{if } (B) \text{ then } S_1 \text{ else } S_2$

code for if-else:



Semantic rules for if-else (contd.)

$B.\text{true} = \text{newlabel}()$

$(B.\text{false} = \text{newlabel}())$

$S_1.\text{next} = S.\text{next}$

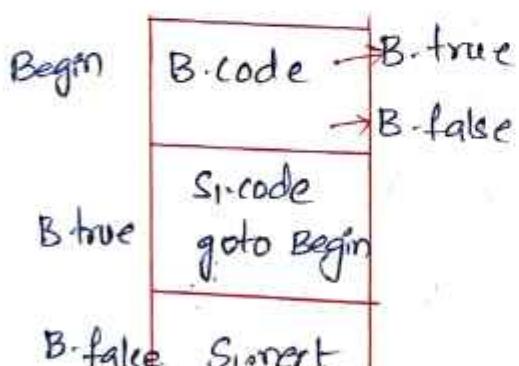
$S_2.\text{next} = S.\text{next}$

~~Scancode Three addresses code~~
Intermediate code

$s.\text{code} = B.\text{code} \parallel \text{label}(B.\text{true}) \parallel S_1.\text{code}$
 $\parallel \text{gen('goto' S.next)}$
 $\parallel \text{label}(B.\text{false}) \parallel S_2.\text{code}$

(iii) while (B) then S_i

code for while



Semantic Rules

Begin = newlabel()

B.true = newlabel()

B.next = begin

B.false = S.next

Intermediate code

$s.\text{code} = \text{label}(Begin) \parallel B.\text{code}$

$\parallel \text{label}(B.\text{true}) \parallel S_i.\text{code}$

production

$p \rightarrow s$

Semantic Rules

$s.\text{next} = \text{newlabel}()$

$p.\text{code} = s.\text{code} \parallel \text{label}(s.\text{next})$

$s.\text{code} = \text{assign_code}$

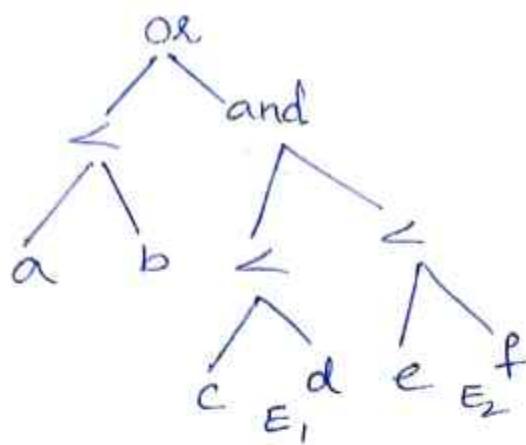
$s \rightarrow \text{assign}$

$s \rightarrow s_1 s_2$

$s_1.\text{next} = \text{newlabel}() \quad s_2.\text{next} = s.\text{next}$

$s.\text{code} = s_1.\text{code} \parallel \text{label}(s_1.\text{next}) \parallel s_2.\text{code}$

eg:- ① $a < b \text{ or } c < d \text{ and } e < f$



if $a < b$ goto $E \cdot \text{true}$

goto E_1

E_1 : if $c < d$ goto E_2
goto $E \cdot \text{false}$

E_2 : if $e < f$ $\cdot E \cdot \text{true}$
goto $E \cdot \text{false}$

eg:- ② if $(x < 100 \text{ || } x > 200 \text{ ff } x \neq y), x = 0$

if $x < 100$ goto L_2

goto L_3

L_3 : if $x > 200$ goto L_4

goto L_1

L_4 : if $x \neq y$ goto L_2

goto L_1

L_2 : $x = 0$

L_1 :

Type and Declarations:-

- * Type checking uses logical rules to decide about the behaviour of program at runtime.
- * It also ensures that types operand match type expected by the operators
Eg:- " && " operation Java expect its two operands to be boolean
int * float → type error

→ Determine the storage needed

Translation Application:

Compiler translates a type of name into storage

Compiler also determines the amount of storage required to store the type name at run time.

Type Expression:-

Type expression is either a basic type or formed by applying an operator called type constructor to a type expression.

→ T-E are used to represent the structure of type,

→ T-E are primitive datatypes.

→ Type name :- is a Type expression Eg: = ~~not~~ type def abc int.

Type expressions are of two types.

(i) Basic type :- Basic type for language are int, real, boolean, char

float, and void. A special type, type-error is used to indicate type errors.

Eg: = int x; Eg: = type abc int;
 int a; is a=b; → depends on language
 abc b;

(ii) Type constructor (or) Type Name :-

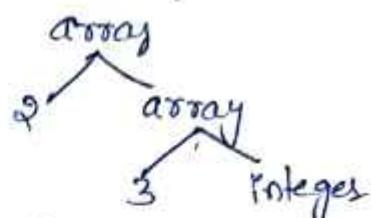
↳ Type constructor applied to list of type expressions.

↳ Types are formed applying an operator called type constructor to type expression.

(1) Arrays:= Arrays are specified as array (I, T) where I \rightarrow Int (s) range of integers
 T \rightarrow Type.

Eg:= In "C" Declaration "int a[100]" identifies type of "a" to be
 array(100,integer) \downarrow int(100) a;

Eg2:= T.E for int [2][3], "a array with 3-integers".
 obj name \downarrow array (2, array (3, integer))
 how many object
 Type Expression



Eg3:= int a(100), b(50);

$a = b$ X type error.

Type Expression for int[2][3]

\rightarrow A type expression can be formed by applying array constructor to a number & type expression.

(11) Record:= A record is a data structure with named fields.

\rightarrow A type expression can be formed by applying a record type constructor to the field name and their types.

Eg:= Struct st

```

    {
        int s;
        float f;
    };
    Struct st s1;
    record ((int, float))
    record (s1, integer)
    record (s1, float)
  
```

Eg2:= Named records are product with named elements
 for record structure with 2 named fields
 length (an integer) and word (of type array (char, char))
 the record is of type.

record (length * Int) X (word X array (char, char))
 struct
 {
 int length;
 char word[10];
 }

\rightarrow Type Expression may contain variable whose values are T.E Eg: int a=5;
product S and T are 2 T.E then their cartesian product S*T is a
 type expression Eg:= int * int.

Function:= Function maps a collection of types to another represented by D \rightarrow R,
 where D is domain R is range of function.

Eg:= int f1 (int x, char y, float z)
 {
 : return m;
 }

Domain = (int X char X float)
 Range -

O/P: int

\rightarrow T.E "int * int \rightarrow char" represent a function that takes 2 integers & returns a char value

Type equivalence

42

→ Two types are said to be equivalent if and only if an operand of one type in an expression is substituted for one of the other types, without type conversion.

Type equivalence are of two types.

1) Name equivalence :=

The two type expression are said to be name equivalence if they have same name or label.

value var1, var2;
total var3, var4;

Int x;

↳ node;

```
node * first, second;
```

struct node *last1, *last2;

→ In the above eg's, var1 and var2 are struct node *last1, *last2;
because their types are same.

→ var3 & var4 also Name equivalence

→ but var & vars are not name equivalent because their types are different

structural equivalence

→ If two expression are the basic type (a)

- Formed by applying the same constructors to structurally types equivalent types then those expression are called structurally equivalent

- (i) It checks the structure of type
- (ii) determines equivalence by whether they have same constructs applied to structurally equivalent types

Eg: type array (I_1, T_1) and array (I_2, T_2) structurally equivalent if $I_1 = I_2 \& T_1 = T_2$

$I_i \rightarrow$ Index of array
 $T_i \rightarrow$ Type of

array $a[100], b[50]$
array $a[100], b[100]$
↓
structurally equivalent

e.g1: type def int value
typedef int number

x : array (50, int)

y : array (100, int)

<u>S₁</u>	<u>S₂</u>	Equivalence	<u>Reason</u>	<u>Reason</u>
char	char	S_1 is equivalent to S_2	similar basic type	
pointer (char)	pointer (char)	S_1 is equivalent to S_2		similar constructs pointers to the char type.

Declarations :=

D → T id ; D E

T → B C | record '{D}'

B → int / float

C → ε / [num] C

D → Sequence of declarations.

T → basic & array and record types

B C → 'component' - generates zero or more integers within the brackets.

- Array type consists of basic type specified "B", followed by array component C.

Eg:- int [10][11]

- Record type is sequence of declaration for field of the record all surrounded by curly braces

record {int, a}

Storage layout for local names:-

→ compiler converts the typenames into the storage.

→ and determines the amount of storage needed to store the typename at runtime.

→ at compile time we can use these amount to assign a ^{type} name to relative address.

$$\text{relative address} = \text{offset} + \text{program counter}$$

→ Relative & types are saved in symbol table entry for type name.

→ Data of varying length such as string or whose size cannot be determined until runtime such as dynamic arrays.

→ The width of a type is no. of storage units needed for objects of that type.

SDT computes types and their widths for basic and array types.

$T \rightarrow B \quad \{ t = B.type; w = B.width; \}$

$C \quad \{ t.type = C.type; T.width = C.width; \}$

$B \rightarrow \text{int} \quad \{ B.type = \text{integer}; B.width = 4; \}$

$B \rightarrow \text{float} \quad \{ B.type = \text{float}; B.width = 8; \}$

$C \rightarrow \epsilon \quad \{ C.type = t; C.width = w; \}$

$C \rightarrow [\text{num}]^{\text{t}}, \{ C.type = \text{array}(\text{num}.value, C_1.type); C.width = \text{num}.value \times C_1.width; \}$

Fig: SDT for computing their types & widths

→ These declarations are represented with DAG & parse tree

Ex: parse tree for $\text{int}[2][3]$

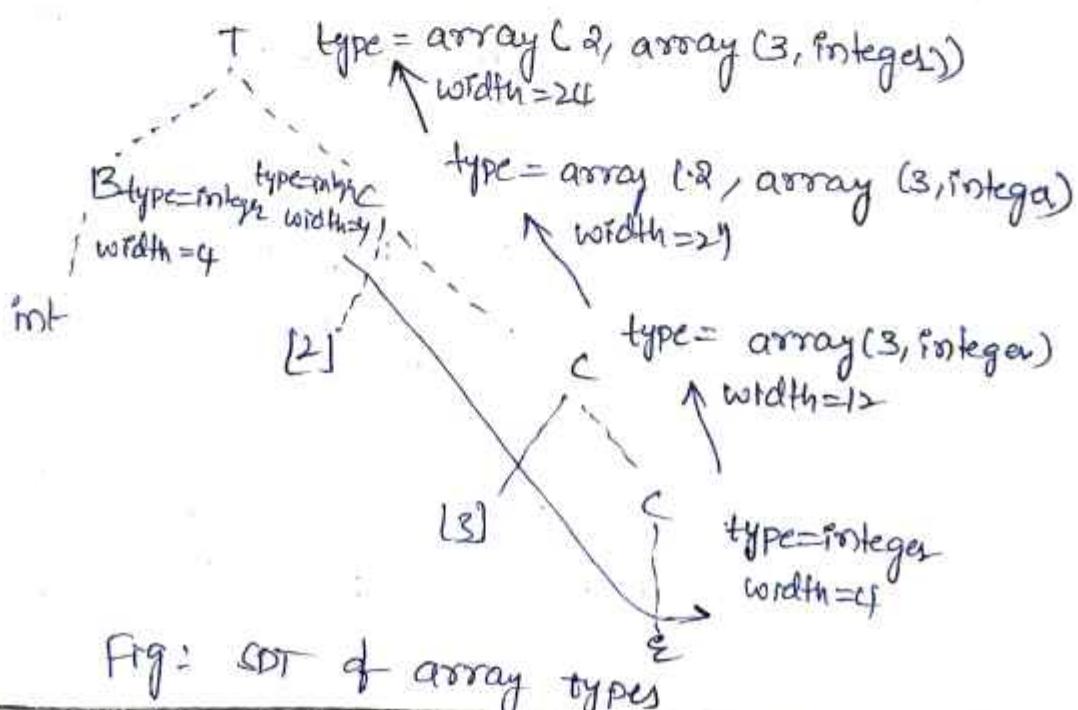


Fig: SDT of array types

Sequence of declarations:

6
24

→ In procedures all the declarations are passed at a time.

→ all declarations in single procedure to be passed as a group.

$P \rightarrow_D \{ \text{offset} = 0; \}$

$D \rightarrow e$

offset — is variable to keep track of the next available relative address.

$D \rightarrow T \text{ id} ; D$ - creates a symbol-table entry for
executing top-put (id.lexeme, Ttype, offset)

top - The current symbol table

~~top-put~~ → creates a symbol table entry for
Id-lexeme with type & relative address.

Field in Records & classes, =

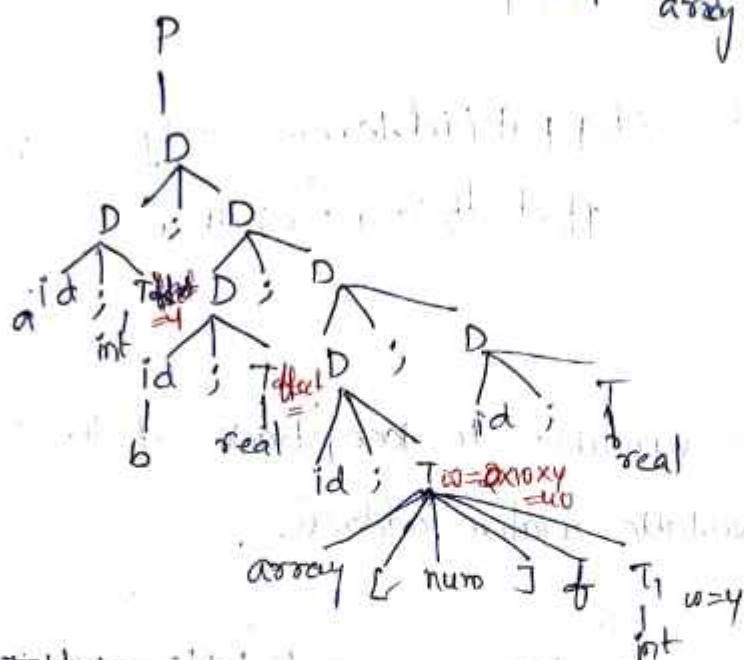
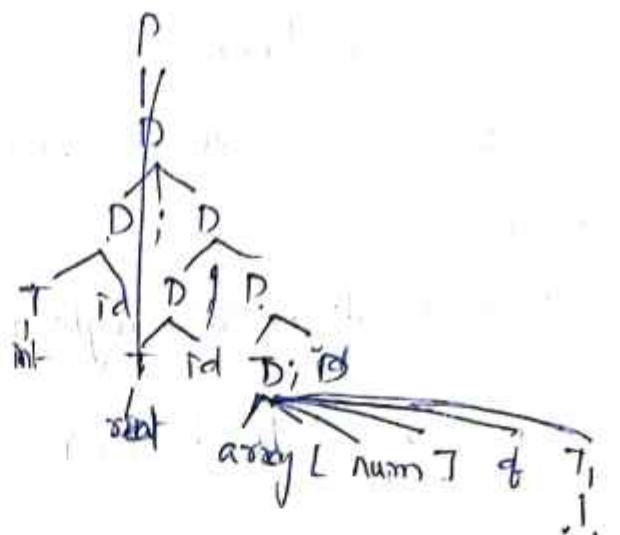
$t \rightarrow \text{record } \langle \text{y} \rangle$

The field in this record is type specified by the sequence of declaration grouped by D.

- The field names within record must be distinct
 - The offset or relative address for fieldname is relative to the data area for that record

Eqn =

```
a:int;
b:real;
c:array[10] of integers;
d:real;
}
```



Symbol Table:

Name	Type	offset	
a	integers	0 [0-3]	offset = 0 + 4 = 4
b	real	4	offset = 0 + 4 + 8 = 12
c	array(10,int)	12	offset = 12 + 40 = 52
d	real	52	

Type checking:-

- Compiler checks whether the program is following type rules or not.
- information about data types is maintained & computed by compiler.
- Type checker is a module of a compiler devoted to typechecking tasks.
- To do typechecking a compiler needs to assign a TE to each component of source program.
- Compiler determines TE conform to collection of logical rules that is called type system for the source language.
- Typechecking catch the errors in program.
- Assign types to values.
- single situation :- check types of objects & report a type error in case of a violation.
- more complex :- Incorrect types may be corrected (type Coercing).

Static

- Type checking done at compile time.
- properties can be verified before program run.
- can catch many common errors.
- Desirable when faster execution importance

Eg: pascal, c-type

- Type checking have been used to ~~genes~~ improve the security of system.

Rules for Type checking:-

Type checking has two forms

- Synthesis
- Interference

Dynamic

- perform during program execution.
- permits programmer to be less concern with c, Pascal strongly.
- Mandatory in some situations such as array, bounds check.
- more robust and clearer code design.

2) Type Synthesis :-

- It derives the expressions from the types of its subexpressions.
- It must be declared before they are used.
- Ex:- The type of $E_1 + E_2$ is defined in terms of the types of E_1 and E_2 .

If f has types $s \rightarrow t$ and x has type s ,
then expression $f(x)$ has type t

Ex:- add(int a, float b)

{

}

fn¹ { float c, int d } $\rightarrow t$

{

add (2, 2.5)

}

[add (float, int)]

[\because int changes to float
float changes to int]

- Here f and π denote expression, $s \rightarrow t$ denote a function from s to t .
- This rule for functions with one argument carries over to functions with several arguments.
- ii) Type inference :-

- It generally determines the type of language construct from the way it is used.
- Ex:- $E_1 + E_2$ i.e., $2 + 5 =$ Datatype will be int
 $\text{int} \quad \text{int}$
 $\text{abc} + \text{abc} =$ Datatype will be string
(string) (string)

- There is no need to declare variables.
- Type inference are used in meta languages.

If $f(x)$ is an expression
then for some α and β , f has a type $\alpha \rightarrow \beta$
and x has type α

Type conversion or type casting:-

- A type cast is basically a conversion from one type to another.
- There are two types of conversions
 - 1) Implicit type conversion
 - 2) Explicit type conversion

1) Implicit type conversion:- (smaller to bigger)

- If a compiler converts one data^{type} into another type of data automatically.
- There is no data loss

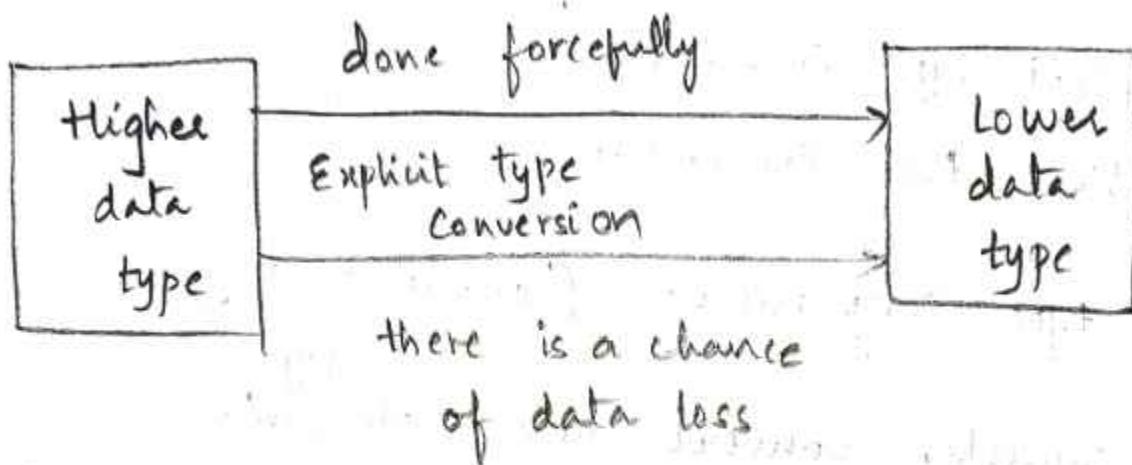
Ex:- `short a=20;`
~~int b=a;~~ // Implicit conversion

Assign:- `bool → char → short int → int → long → float`

2) Explicit type conversion:-

- When data of one type is converted explicitly to another type with the help of predefined functions.

- There is a data loss.
- Conversion done forcefully.
- Some conversions cannot be made implicitly
int . to short (" int range is more than short
so there is a chance of data loss)



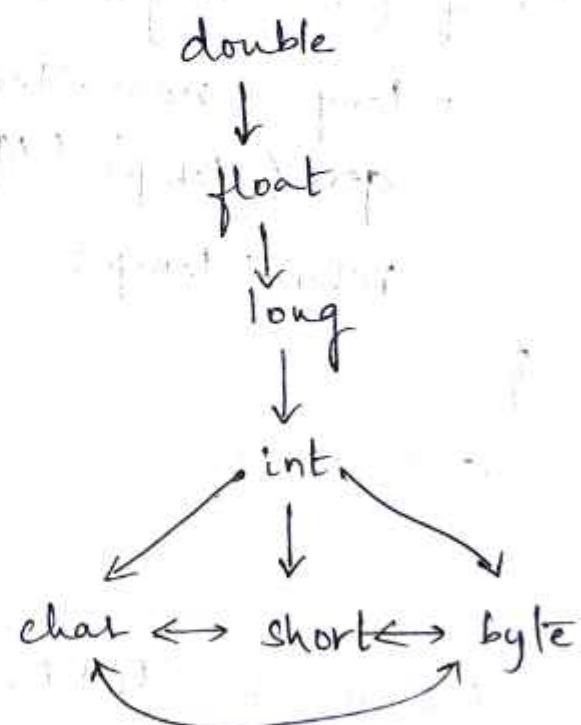
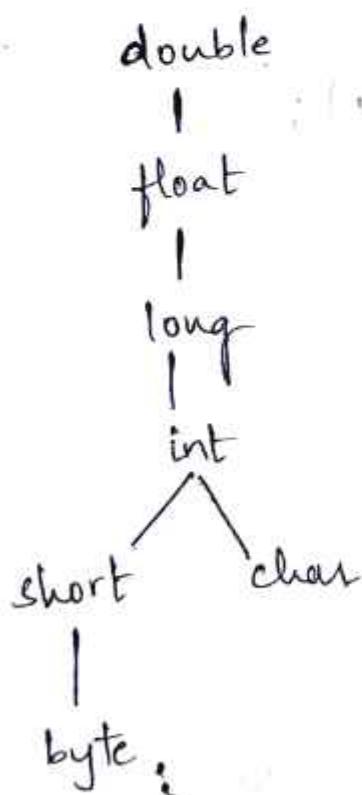
Ex:- $t_1 = (\text{float}) 2$

$t_2 = t_1 * 3.14$

Ex:- if ($E_1 \cdot \text{type} = \text{integer}$ and $E_2 \cdot \text{type} = \text{integer}$)
 $E \cdot \text{type} = \text{integer}$;
else if ($E_1 \cdot \text{type} = \text{float}$ and $E_2 \cdot \text{type} = \text{integer}$)
 $E \cdot \text{type} = \text{float}$;

- two Conversions
 - i) Widening Conversions
 - ii) Narrowing Conversions.

- widening conversions generally preserve information
- narrowing conversions generally lose information
- widening rules - any lower can be widened to higher type.
- a char can be widened to int or float
but char cannot be widened to short.
- narrowing rules - a type s can be narrowed to type t if there is a path from t.



1) $\text{max}(t_1, t_2)$ takes two types t_1 and t_2 and returns the maximum of two types in widening hierarchy.

2) $\text{widen}(a, t, w)$ generates type conversions if needed to widen the contents of an address a of type t , into a value of type w .

```
Addr { widen( Addr a, Type t, Type w )
    if (t=w) return a;
    else if (t = integer and w = float) {
        temp = new Temp();
        gen('temp' = '(float)' a);
        return temp;
    }
    else error;
}
```

→ Semantic Action $E \rightarrow E_1 + E_2$

```
 $E \rightarrow E_1 + E_2 \{ E \cdot \text{type} = \text{max}(E_1 \cdot \text{type}, E_2 \cdot \text{type}),$ 
 $a_1 = \text{widen}(E_1 \cdot \text{addr}, E_1 \cdot \text{type}, E \cdot \text{type});$ 
 $a_2 = \text{widen}(E_2 \cdot \text{addr}, E_2 \cdot \text{type}, E \cdot \text{type});$ 
 $E \cdot \text{addr} = \text{new Temp}();$ 
 $\text{gen}(E \cdot \text{addr}' = 'a_1 + 'a_2); \}$ 
```

Overloading of Functions and operators:-

An overloaded symbol has different meanings depending on its context. Overloading is resolved when a unique meaning is determined for each occurrence of a name.

Ex:- The + operator in Java denotes either string concatenation or addition, depending on the type of its operands.

void err() { -- }

void err(String s) { -- }

Intermediate code for switch statements (a) Three Address code
 Translation of switch statement:

Switch statement syntax:

switch (E)

{

case $v_1 : s_1$

case $v_2 : s_2$

case $v_{n-1} : s_{n-1}$

default : s_n

}

$t_1 := \text{switch}(x+y)$

2 case 1: $a = a + 2$;

break;

case 4: $b = b * 5$;

break;

case 6: $c = c / 2$;

break;

default: $d = d - 2$;

break;

}

Translation of switch statements:

Code to evaluate E into t
 goto test

L_1 : code for s_1
 goto next

L_2 : code for s_2
 goto next

L_{n-1} : code for s_{n-1}
 goto next

L_n : code for s_n
 goto next

test: if $t = v_1$, goto L_1
 if $t = v_2$, goto L_2

if $t = v_{n-1}$, goto L_{n-1} , goto L_n

Three Address code:

1. $t_1 = x + y$

2. goto L_1

next:

1. $t_1 = x + y$

15. $c = t_5$

2. If ($t_1 = 1$) goto L_2

16. goto Next

3. If ($t_1 = 4$) goto L_3

17. Next

4. If ($t_1 = 6$) goto L_4

5. $t_2 = d - 2$

18. $c = t_5$

6. $d = t_2$

19. goto Next

7. goto Next

8. $t_3 = a + 2$

9. $a = t_3$

10. goto Next

11. $t_4 = b * 5$

12. $b = t_4$

13. goto Next

14. $t_5 = c / 2$

H $t_i = v_{n-i}$ goto L_{n-i}
 goto L_n

next:

Intermediate code for procedures = (d) Three address code

D → define T id (F) { S } | \rightarrow S → adds stmt that returns the
 value of an expression.
 F → E | T id, F | \rightarrow E → adds function calls, with actual
 parameters A.
 Non-terminals D and T generates
 declarations and types.

E → id (A);

A → ε / E, A

float add()

(int a,

int b)

{

 return add();

}

\rightarrow Function definition generated by D consists
 of keyword define, a return type,
 the function name, formal parameters
 in parenthesis and function body
 consisting of statements.

\rightarrow Non-terminal F generates zero or more
 formal parameters.

where formal parameters consist of
 a type followed by identifiers

\rightarrow Non-terminal S → generates statements
 expressions.

\rightarrow In three-address code, a function call is unraveled into the
 evaluation of parameters in preparation for a call followed by call itself.
 and the parameters are passed by value.

Eg: If the given function is in the form of

P(A₁, A₂, A₃, ..., A_n) Eg₂: n = f(a[i]);

param A₁,

param A₂

param A_n

call P, n

Translated into three-address code as follows

1) $t_1 = i * 4$

2) $t_2 = a[t_1]$

3) param t₂

4) $t_3 = \text{call } f, 1$

5) $n = t_3$

P → ie function name.

n → no. of arguments.

- The first 2-lines compute the value of expression $a[i]$ into temporary t_2 ,
- line 3 makes t_2 ^{an} actual parameter for the call on line 4 of f with one parameters
- line 5 assign the value returned by the function call to t_3 .

→ functions types:-

- The type of function must encode the return type and types of the return type and the types of the formal parameters.
- Let "void" be a special type that represent no parameter or no return type.
- whenever the function is called the function name is ~~is~~ entered into the symbol table for use in the rest of the program.
- The formal parameters are stored in the Activation Record. For storing formal parameters the Activation Records are used.

Eg 2: void main()

```

{
    int x, y
    ...
    ...
    swap(&x, &y);
}

```

Three address code

1. call main
2. param &x
3. param &y
4. call swap, 2

void swap(int *a, int *b)

```

{
    int i;
    i = *b;
    *b = *a;
    *a = i;
}

```

Eg 3: float add() or float add(int a)

float add(int a, float b)

{

return add() or return add(x),

return add(x,y);

}

↑
formal
parameters

Actual parameters

UNIT-V

Machine-Independent optimizations:- The principal sources of optimization, Introduction of Data-flow analysis, foundations of Data-flow Analysis, constant propagation, partial Redundancy Elimination, Loops in flow Graphs.

=

Optimization is the process of transformation of code to an efficient code. Efficiency is in terms of space requirement and time for its execution without changing the meaning of the given code.

The following constraints are to be considered while applying the techniques for code improvement.

- ① The transformation must preserve the meaning of the program, that is, the target code should ensure semantic equivalence with source program.
- ② Program efficiency must be improved by a measurable amount without changing the algorithm used in the program.
- ③ When the technique is applied on a special format, then it is worth transforming.

Optimization can be classified as

(1) Local optimizations:-

Optimizations performed within a single basic block are termed as local optimizations. These techniques are simple to implement and does not require any analysis since we do not require any information relating to how data and control flows.

(2) Global optimizations:-

Optimizations performed across basic blocks is called global optimizations. These techniques are complex as it requires additional analysis to be performed across basic blocks. This analysis is called data-flow analysis.

optimization techniques can be applied to intermediate code or on the final target code. It is a complex and a time-consuming task that involves multiple sub phases, sometimes applied more than once.

* most compilers allow the optimization to be turned off to speed up compilation process.

* To apply optimization it is important to do control flow analysis and data flow analysis followed by transformations.

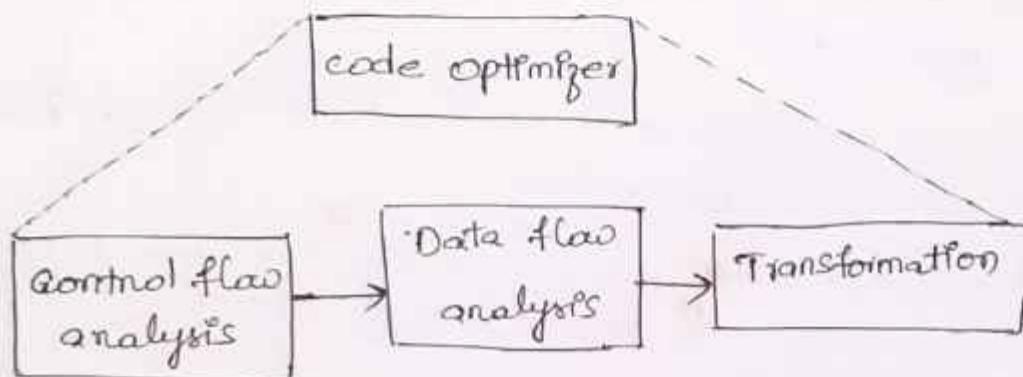


Fig: code optimization model.

control flow analysis:- It determines the control structure of a program and builds a control flow graph.

Data flow Analysis:- It determines the flow of scalar values and builds data flow graphs. The solution to flow analysis propagates data flow information along a flow graph.

Transformation:- Transformations help in improving the code without changing the meaning or functionality.

Flow Graph:-

A graphical representation of three address code is called flow graph. The node in the flow graph represent a single basic block and the edges represent the flow of control. These flow graphs are useful in performing the control flow and data flow analysis and to apply local or global optimization and.

2

code generation.

Basic Block:- A basic block is a set of consecutive statements that are executed sequentially. Once the control enters into the block then every statement in the basic block is executed one after the other before leaving the block.

Eg: for the statement $a = b + c * d / e$ the corresponding set of three address code is.

$$\begin{aligned}t_1 &= c * d \\t_2 &= t_1 / e \\t_3 &= b + t_2 \\a &= t_3.\end{aligned}$$

All these statements correspond to a single basic block.

Procedure to identify the Basic Blocks:

Given a three address code, first identify the leader statements and group the leader statements with the statements up to the next leader.

* To identify the leader use the following rules:

1. First statement in the program is a leader.
2. Any statement that is the target of a conditional or unconditional statement is a leader statement.
3. Any statement that immediately follows a conditional/unconditional statement is a leader statement.

Example:- Identify the basic blocks for the following code fragment

```
main()
{
    int p=0, n=10;
    int a[n];
    while ( i<=(n-1) )
    {
        a[p] = p * p;
        p = p+1;
    } return n;
}
```

The three address code for initialise function is as follows:

- (1) $i := 0$ \longrightarrow Leader 1 using rule 1
(2) $n := 10$
(3) $t_1 := n - 1$ \longrightarrow Leader 2 using rule 2
(4) if $i > t_1$ goto (12)
(5) $t_2 := i * i$ \longrightarrow Leader 3 using rule 3.
(6) $t_3 := 4 * i$
(7) $t_4 := a[t_3]$
(8) $t_4 := t_2$
(9) $t_5 := i + 1$
(10) $i := t_5$
(11) goto (3)
(12) return. \longrightarrow Leader 4 using rule 3.

1) $i = 0$
2) $n = 10$

B1

Basic Block 1 includes statements
(1) and (2)

3) $t_1 := n - 1$
4) if $i > t_1$ goto (12)

B2

Basic Block 2 includes statements
(3) and (4)

5) $t_2 := i * i$
6) $t_3 := 4 * i$
7) $t_4 := a[t_3]$
8) $t_4 := t_2$
9) $t_5 := i + 1$
10) $i := t_5$
11) goto (3)

B3.

Basic Block 3 includes statements
(5)-(11)

12) return.

B4.

Basic Block 4 includes statement
(12)

Fig: Basic Blocks.

Flow Graph:-

flow graph shows the relation between the basic block and its preceding and its successor blocks. the block with the first statement is B1. An edge is placed from block B1 to B2 , if block B2 could immediately follow B1 during execution or satisfies the following conditions.

- The last statement in B1 is either conditional or unconditional jump statement that is followed by the first statement in B2 or
- The first statement in B2 follows the last statement in B1 and is not an unconditional / conditional jump statement

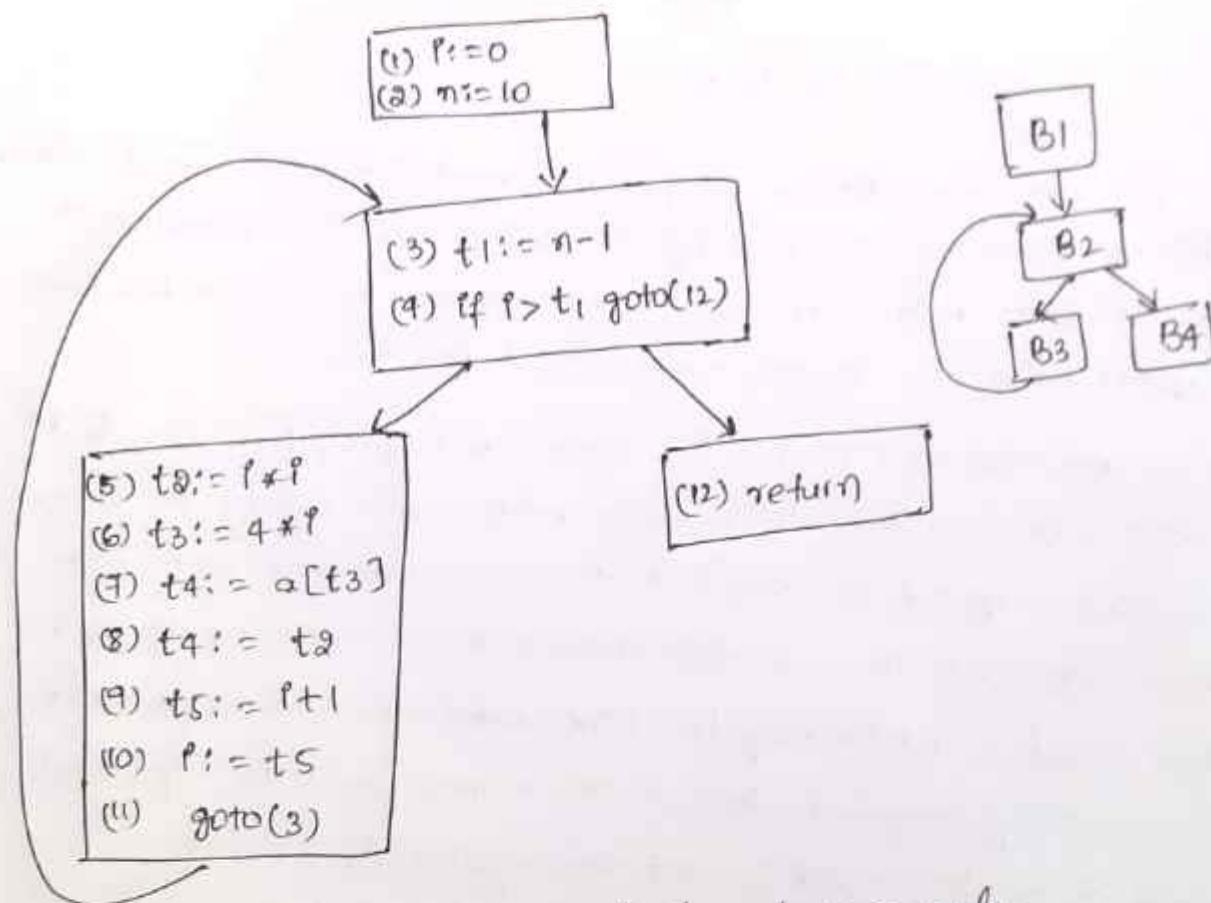


fig: flow graph for above example.

DAG representation of Basic Block.

A DAG is a useful data structure for implementing transformations within a basic block. It gives pictorial representation of how values computed at one statement are useful in computing the values of other variables.

- * It is useful in identifying common sub-expressions within a basic block.
- * A DAG has nodes, which are labeled as follows.
 - ① The leaf nodes are labeled by either identifiers or constants. If the operators are arithmetic then it always requires two r-values.
 - ② The labels of interior nodes correspond to the operator symbol.

Note: The DAG is not the same as a flow graph. Each node in a flow graph is a basic block, which has a set of statements that can be represented using DAG.

=

Construction of DAG.

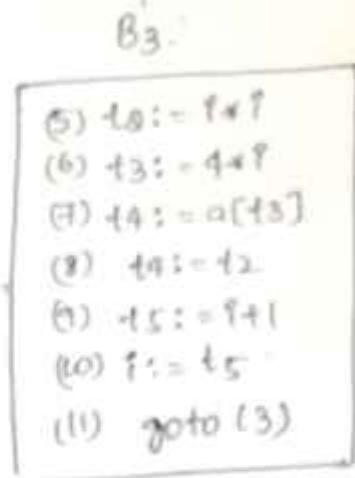
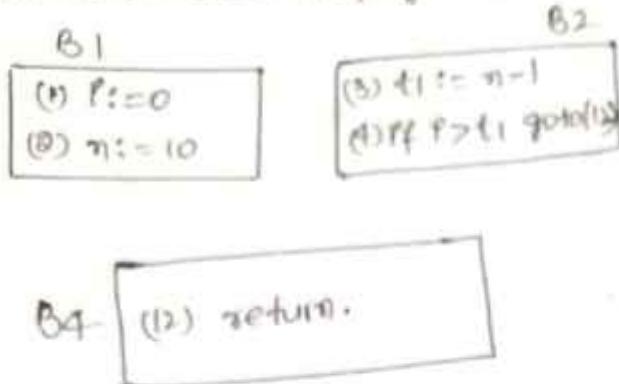
To construct DAG, we process each statement of the block. If the statement is a copy statement, that is, a statement of the form $a = b$, then we do not create a new node, but append label a to the node with label b .

* If the statement is of the form $a = b \text{ op } c$, then we first check whether there is a node with same values as $b \text{ op } c$. If so, we append the label a to that node. If such node does not exist then we first check whether there exists nodes for b and c which may be leaf nodes or an internal nodes if recently computed, then create a new node for 'op' and add it the left child b and the right child as ' c '.

* Label this new node with a . This would become the value of a to be used for next statements; hence, we mark the previously marked nodes with a as a_0 .

=

Example:- Let us consider Basic block B1 to B4 and construct DAG for each block step by step.



For Block B1

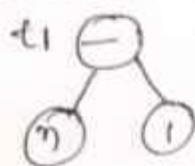
Fig(a)

Fig: DAG for Block B1

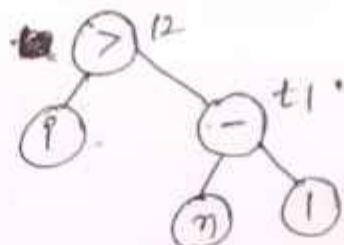
* for the statement ~~return~~, first we create a leaf labeled 4 and attach it to it.

For Block B2

For first statement, which we first create nodes for n and 1, then create node for operator (-) and label it as t1.



DAG for first statement
in Block B2

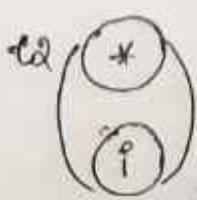


Fig(b): DAG for Block B2.

* Fig(b) shows the second statement, which is a conditional statement, we create a node for operator > and label it as t2 as when this condition is satisfied it should go to statement 12.

For Block B3

* for the first statement, we create nodes for i and use as right child, then create node for operator (*) and label it as t2.



Fig(a): DAG for Block B3.

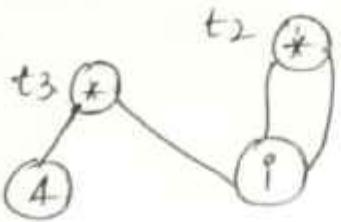
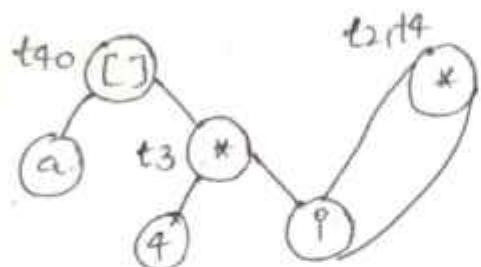
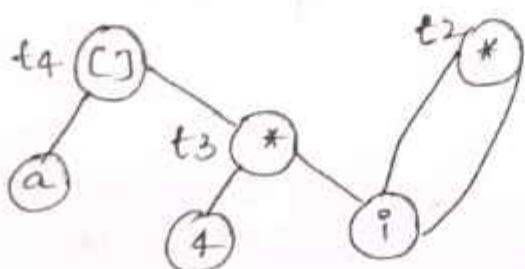


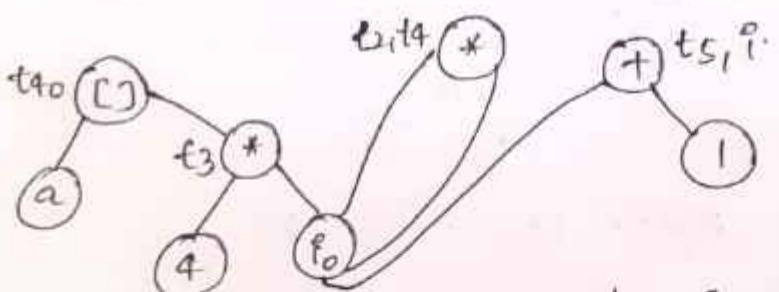
Fig: DAG for two statements in Block B3.



Fig(B): DAG for three statements in Block B3.



Fig(C): DAG for four statements in Block B3.



Fig(E): DAG for complete block B3.

Principle Sources of Optimization.

5

A transformation of a program is called local if it is applied within a basic block and global if applied across basic blocks.

+ there are different types of transformations to improve the code and these transformations depend on the kind of optimization required.

① Function - preserving transformations are those transformations that are performed without changing the function it computes.
+ these are primarily used when global optimizations are performed.

② Structure - preserving transformations are those that are performed without changing the set of expressions computed by the block.
+ many of these transformations are applied locally.

③ Algebraic transformations are used to simplify the computation of expression set using algebraic identities.
+ These can replace expensive operations by cheaper ones, for instance multiplication by 2 can be replaced by left shift.

Function preserving Transformations.

The following techniques are function-preserving transformations

1) common sub-expression Elimination

2) copy propagation.

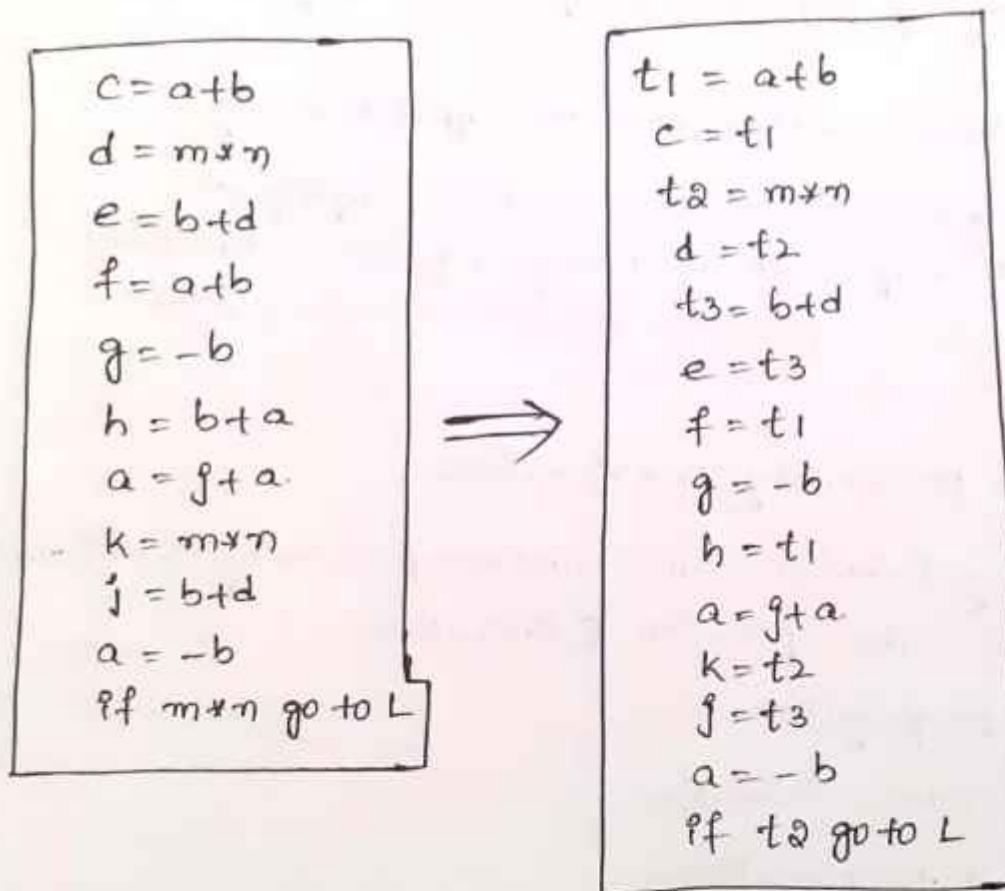
3) Dead code elimination

4) constant propagation.

(i) common sub-expression Elimination:-

An expression E is said to be common sub expression if E is computed before and the variables in the expression are not modified since its ~~computation~~ computation. If such expression is present then the re-computation can be avoided by using the result of the previous computation.

- * This technique can be applied both locally and globally
- * we need to maintain a table to store the details of expressions evaluated so far and use this information to identify and eliminate the re-computation of the same expression.
- * The common sub-expression elimination can be done within basic block by analyzing and storing the information of expression in the table until the operands in the expression are redefined
- * If any operand in the expression is redefined, then remove the expression from the table



Program code before and after common subexpression elimination

- * The above figure shows the optimized code on applying common subexpression elimination technique locally.

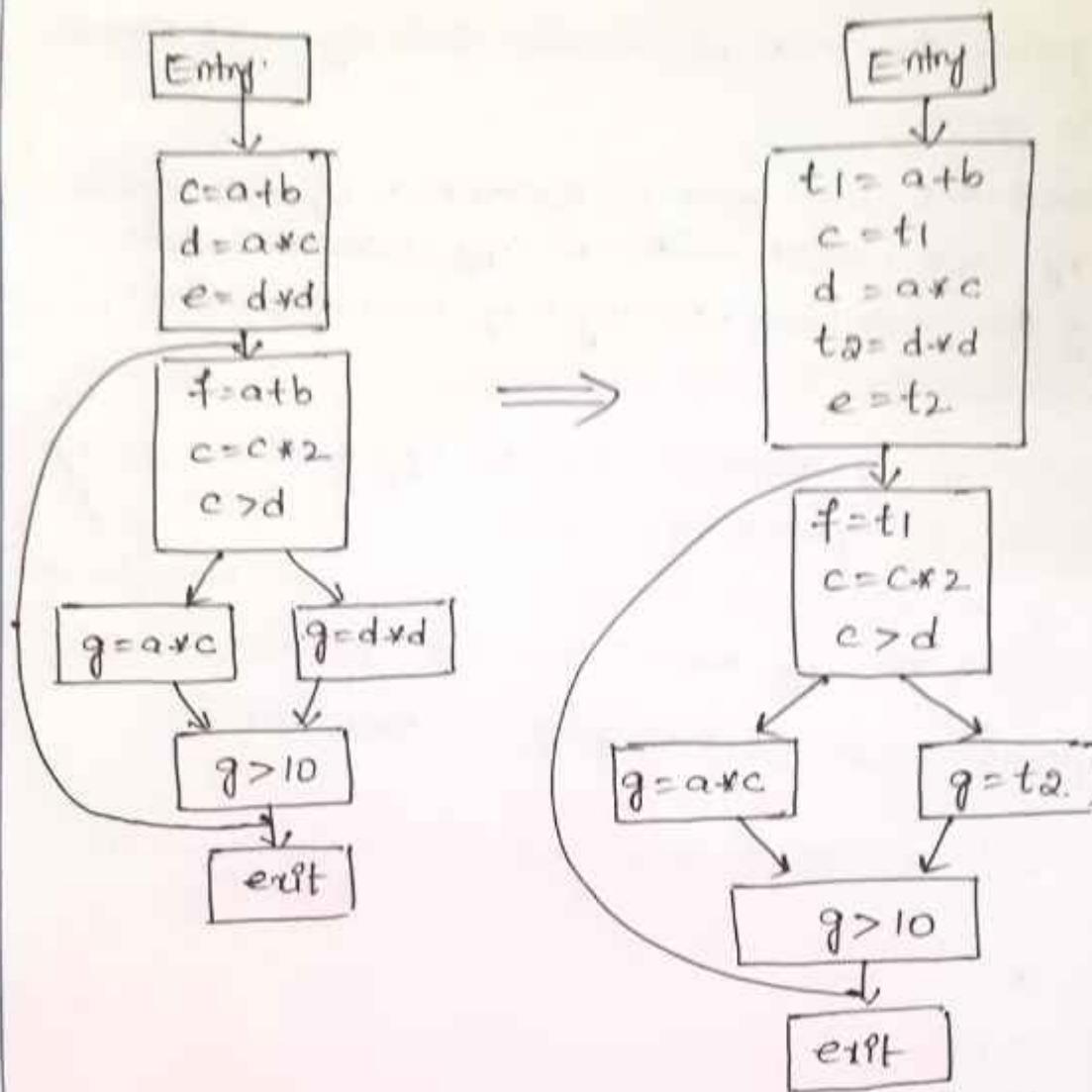


Fig: Example for global common sub expression Elimination.

* The above example shows the common-subexpression elimination globally by replacing the evaluated expression with a new temporary variable t_1 and t_2 and these variables are used wherever the same expression is used.

=

(ii) Copy propagation :-

A copy statement is a statement that is in the form $a = b$.

* copy propagation technique is applied to replace the later use of variable 'a' with the use of 'b' if original values of 'a' do not change after this assignment.

* This technique can be applied locally and globally.

- * This optimization reduces runtime stack size and improves execution speed.
- * To implement this, we need to maintain a separate table called copy table, which holds all copy statements. While traversing the basic block, if any copy statement is found, add this information into the copy table.
- * While processing any statement, check the copy table for variable a , which can be replaced variable b .
- * If there is an assignment statement that computes the value of a then remove the copy statement from the copy table.
- * Copy propagation helps in identifying the dead code.

Example:— Let the basic block contain the following set of statements.

$$\begin{aligned} Y &= X \\ Z &= Y+1 \\ W &= Y \\ Y &= W+Z \\ Y &= W \end{aligned}$$

Statement NO	Instruction	Updated Instruction	copytable content
1.	$Y = X$	$Y = X$	$\{(Y, X)\}$
2.	$Z = Y+1$	$Z = X+1$	$\{(Y, X)\}$
3.	$W = Y$	$W = X$	$\{(Y, X), (W, X)\}$
4.	$Y = W+Z$	$Y = X+Z$	$\{(W, X)\}$
5.	$Y = W$	$Y = W$	$\{(W, X), (Y, X)\}$

- * On the first statement 1, since it is a copy statement, the information is added in to the copy table.
- * Statement 5 uses the value X indirectly from Y , hence $GET(Y, \text{table})$ could return X and the statement is modified

7

Replacing Y with X.

- * The third statement is a copy statement and since Y is to replace X, we insert into the copy table W to be replaced by X
- * When statement 4 is processed, Y value is defined, and hence, details of Y are removed from the copy table.
- * The fifth statement is a copy statement and is inserted into the copy table.

(iii) Dead code Elimination:-

Code that is unreachable or does not affect the program is said to be dead code. Such code requires unnecessary CPU time, which can be identified and eliminated using this technique.

Example program:-

```
1) Put Var1;  
2) void Sample()  
3) {  
4)     Put 9;  
5)     p=1; /* dead store */  
6)     Var1=1; /* dead store */  
7)     Var2=2;  
8)     return;  
9)     Var1=3; /* unreachable */  
10) }
```

The code fragment after dead code elimination

```
Put Var1;  
void Sample()  
{  
    Var1=2;  
    return;  
}
```

→ Here the value of p is never used, and the value assigned to Var1 variable in statement 6 is dead store and the statement 9 is unreachable. These statements can be eliminated as they do not affect the program execution.

(iv) Constant propagation :-

Constant propagation is an approach that propagates the constant value assigned to a variable at the place of its use.

For example, given an assignment statement $t = c$, where c is a constant, replace later uses of t with uses of c , provided there are no intervening assignments to t .

This approach is similar to copy propagation and is applied at first stage of optimization. This method can analyze by propagating constant value in conditional statement, to determine whether a branch should be executed or not, that is, identifies the dead code.

Example: Let us consider the following example:

```
pi = 22/7
void area_per(int r)
{
    float area, perimeter;
    area = pi * r * r;
    perimeter = 2 * pi * r;
    print area, perimeter;
```

3.

In this example, we can notice some simple constant propagation results, which are as follows.

- In line 1 the variable π is constant and this value is the result of $22/7$, which can be computed at compile time and has the value 3.1413.
- In line 5 the variable π can be replaced with the constant value 3.1413.
- In line 6 since the value of π is constant, the partial result of the statement can be computed and the statement can be modified as $\text{perimeter} = 6.283 * r;$

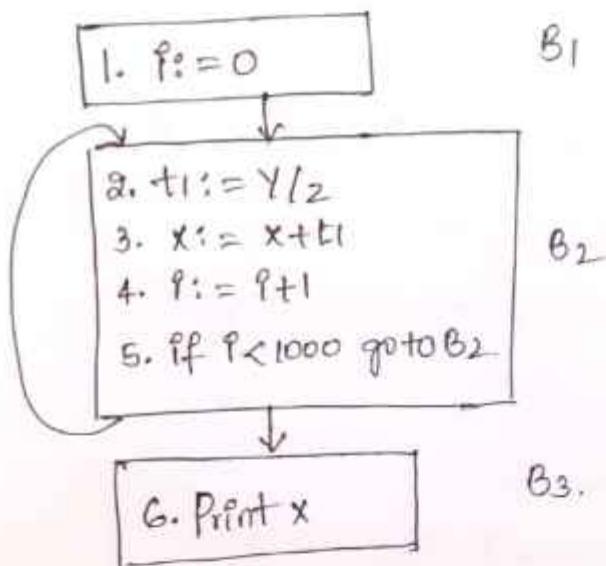
Loop optimization Techniques.

In most of the programs, 90% of execution time is inside the loops; hence loop optimization is the most valuable machine - independent optimization and are good candidates for code improvement.

Let us consider the example of a program that uses a loop statement:

```
for (int i=0; i<1000; i++)
    x = x+y/z;
    print x;
```

The optimized code of this program is represented in flow graph as



flow graph after code optimization.

The total number of statements that are executed are as follows

<u>Statement Number</u>	<u>Frequency of Execution</u>	<u>Total No. of Statements</u>
1	1	1
2-5	1000	4000
6	1	1

on the code, the number of statements that are executed are 4002. Instead if the code is written as follows then the total number of statements that are executed are only 1002.

1. $t_1 = y/2$
2. $z = z + t_1$
3. $z = z + t_1$
4. $z = z + t_1$
5. \dots
6. \dots
1001. $z = z + t_1$
1002. Print z .

The execution is three times better for later than the first form, but the space requirements is more. Inefficent code is generated in the loops for various reasons like

- Induction Variable usage to keep track of iteration
- Unnecessary computations made inside the iterative loops that are not effective.
- Use of high strength operators inside the loop.

The important loop optimizations are

- ① Elimination of loop invariant computations.
- ② Strength reduction
- ③ Code motion
- ④ Elimination of induction variables.

(1) A Loop Invariant computation:-

A loop invariant computation is one that evaluates the same value every time the statements in loop are executed.
* moving such loops out/side computational statements outside the loop leads to a reduction in the execution time.

Algorithm for elimination of loop invariant code.

9

Step 1: Identify loop-invariant code.

Step 2: An instruction is loop-invariant if, for each operand:

- (i) The operand is constant, OR
- (ii) All definitions for all the operands that reach the instruction inside the loop are made outside the loop, OR
- (iii) There is exactly one definition made using loop invariant variables inside the loop for an operand that reaches the instruction.

Step 3: move it outside the loop

- (i) move each instruction i to newly created preheader if and only if it satisfies these requirements of the loop.

Example:-

```
for (int i=0; i<1000; i++)  
    X = X + Y/z;  
    Print X;
```

In the above example, the value y and z remains unchanged as there is no definition for these variables in the loop. Hence, the computation of Y/z remains unchanged, which can be moved above the loop and the same code can be rewritten as follows.

```
t1 = Y/z ;  
for (int i=0; i<1000; i++)  
    X = X + t1 ;  
    Print X ;
```

(a) Induction Variables:-

Induction variables are those variables used in a loop, their values are in lock-step, and hence, it may be possible to eliminate all except one.

* There are two types of induction variables - basic and derived.

* Basic induction variables are those that are of the form:

$$I = I + c$$

where I is loop variable and c is some constant.

* Derived induction variables are those that are defined only once in the loop and their value is linear function of the basic induction variable.

For example:-

$$J = A * I + B$$

Here J is a variable that is depended on the basic induction variable I and the constants A and B . This is represented as a triplet (I, A, B) .

* Induction variable elimination involves three steps.

1) Detecting induction variable.

2) Reducing the strength of induction variable.

3) Eliminating induction variable.

Induction Variable Elimination:-

Some loops contain two or more induction variables that can be combined into one induction variable.

Example:- The code fragment below has three induction variables (i_1 , i_2 and i_3) that can be replaced with one induction variable, thus eliminating two induction variables.

```

int a[SIZE];
int b[SIZE];
void f(void)
{
    int i1, i2, i3;
    for (i1=0, i2=0, i3=0; i1<SIZE; i1++)
        a[i2++] = b[i3++];
}
return;

```

After induction variable elimination the code is rewritten as

```

int a[SIZE]
int b[SIZE]
void f(void)
{
    int i1;
    for (i1=0; i1<SIZE; i1++)
        a[i1] = b[i1];
}
return;

```

- * Induction variable elimination can reduce the number of additions (or subtractions) in a loop, and improve both runtime performance and code space.

Machine - Dependent Optimization

This optimization can be applied on target machine instructions. This includes register allocation, use of addressing modes, and peep hole optimization.

- * Instructions involving register operands are faster and shorter; hence, if we make use of more registers during target code generation, efficient code will be generated.

Peephole optimization:-

Generally code generation algorithms produce code, statement by statement. This may contain redundant instructions and suboptimal constructs. The efficiency of such code can be improved by applying peephole optimization, which is simple but effective optimization on target code.

- * The peephole is considered a small moving window on the target code. The code in peephole need not be contiguous. It improves the performance of the target program by examining and transforming a short sequence of target instructions.
- * The advantage of peephole optimization is that each improvement applied increases opportunities and shows additional improvements. It may need repeated passes to be applied over the target code to get the maximum benefit.

=

(i) Redundant loads and stores

The code generation algorithm produces the target code, which is either represented with single operand or two operands or three operands.

- * Let us assume the instructions are with two operands. The following is an example that gives the assembly code for the statement

$x = y + z$.

1. MOV Y, R0
2. ADD Z, R0
3. MOV R0, X

- * Instruction 1 moves the value of Y to Register R0, second instruction performs the addition of value in Z with the register content and the result of the operation is stored in the register.
- * The third instruction copies the register content to the location of x. At this point the value of x is available in both location of x and the register R0.

If the above algorithm is applied on the code $a = b + c, d = a + c$
then it generates the code given below.

1. MOV b, R₀
2. ADD C, R₀
3. MOV R₁, a
4. MOV a, R₀
5. ADD C, R₀
6. MOV R₀, d

Here we can say that 3 and 4 are redundant load and store instructions. These instructions will not affect the values before or after their execution. Such redundant statements can be eliminated and the resultant code is as follows.

1. MOV b, R₀
2. ADD C, R₀
3. ADD C, R₀
4. MOV R₀, d

(ii) Algebraic Simplification:-

There are few algebraic identities that occur frequently enough and are worth considering.

Look at the following statements.

$$x := t + 0$$

$$x := x + 1$$

They do not alter the value of x. If we keep them as it is, later when code generation algorithm is applied on it, it may produce statements that are of no use. Hence, such statements whether they are in three address code or target code can be removed.

(iii) Deadcode Elimination:-

Removal of unreachable code is an opportunity for peephole optimization. A statement immediately after an

unconditional jump or a statement that never get a chance to be executed can be identified and eliminated. Such code is called the dead code.

Ex: If define $x=0$

```
if(x)
{
    ... print value.
}
```

If this is translated to target code as

```
If x=1 goto L1
    goto L2.
```

L1: print value.

L2: -----

here, value will be never printed. So whatever code inside the body of "if(x)" is dead code; hence it can be removed.

(iv) Reduction in Strength:-

This optimization mainly deals with replacing expensive operations by cheaper ones for example.

$\rightarrow z^2 \Rightarrow z * z$

\rightarrow fixed-point multiplication and division by a power of 2.

\Rightarrow shift

\rightarrow floating point division by a constant \Rightarrow floating-point multiplication by a constant.

≡