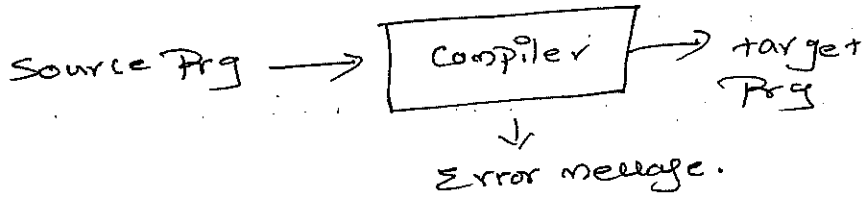


# Compiler Design

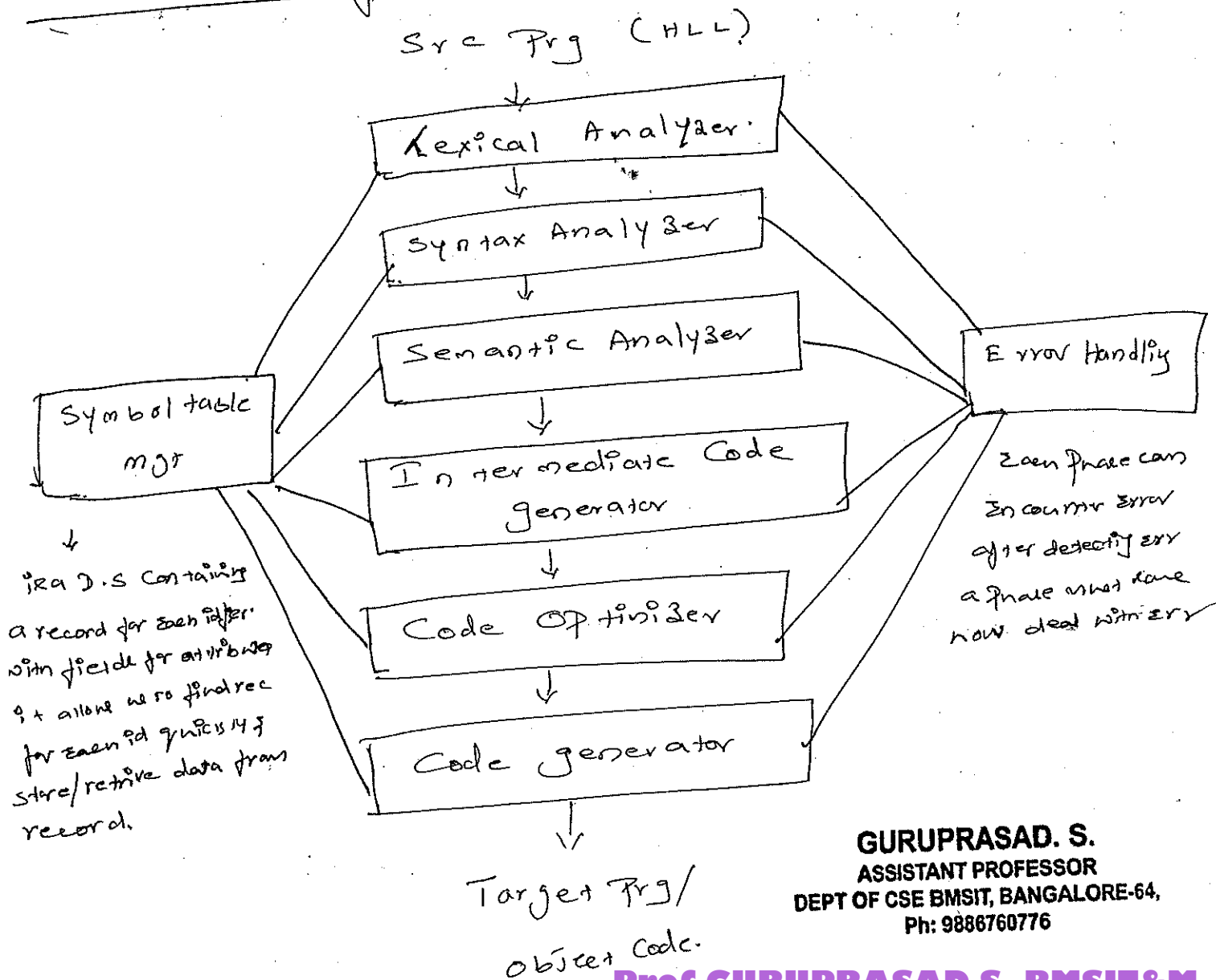
## Compiler:

Compiler is a program which accepts source program written in HLL and produces object code/target program.



Since compiler is a huge program, & it is difficult to understand the entire compilation process, so the process is divided into no. of modules called Phases.

## The Phases of Compiler / Structure of Compiler



**GURUPRASAD. S.**  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-64,  
Ph: 9886760776

**Prof GURUPRASAD S, BMSIT&M**

<https://hemanthrajhemu.github.io>

## Lexical Analyser

- It scans the Src Prg from left to right & breaks the Src Prg into meaningful tokens
- It removes extra white spaces like tab, space, new line
- removes all comments.
- Places / makes entry for every Variable, Constant, label into symbol table
- An entry for a Variable / const is an ordered pair  
< token name, attribute value >  
< id, 1 >

## Syntax Analyse

- It groups a set of tokens to identify a syntactical structure defined by language.
- It identifies the structure of the Prg from tokens obtained by Lexical Analyse.
- It detects syntax errors & produces parse tree / derivation tree

## Semantic Analyse

- It checks for the semantics of the identified syntactical structures.
- checks for type checking, scope of variables, pointers & usage, etc

## Intermediate Code generation

- Produces explicit intermediate representation of Src Prg for abstract m/c
- It should have two important properties  
Easy to produce  
Easy to understand & translate

eg:-  $D = A + B * C$

$T_1 = B * C$

$T_2 = A + T_1$

$D = T_2$

either three address rep<sup>n</sup> or two addr rep<sup>n</sup> is needed.

Code optimization:

→ It attempts to improve intermediate code so that the target code runs faster & consume less memory space.

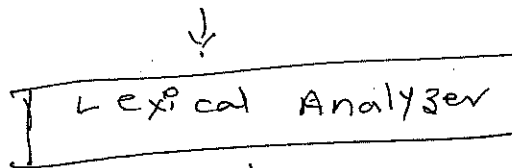
→ It uses techniques like local optimization, loop opt<sup>n</sup>, dead code elimination etc

Code generation:

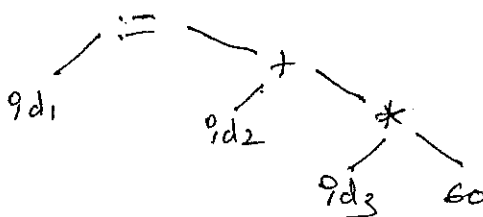
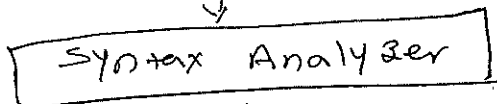
→ It converts the optimized intermediate code into sequence of m/c instr<sup>n</sup>.

→ This phase should know the m/c specific details & should utilize registers efficiently

eg:-  $Position := initial + rate * 60$



<id<sub>1</sub>> <=> <id<sub>2</sub>> <+> <id<sub>3</sub>> <\*> <60>

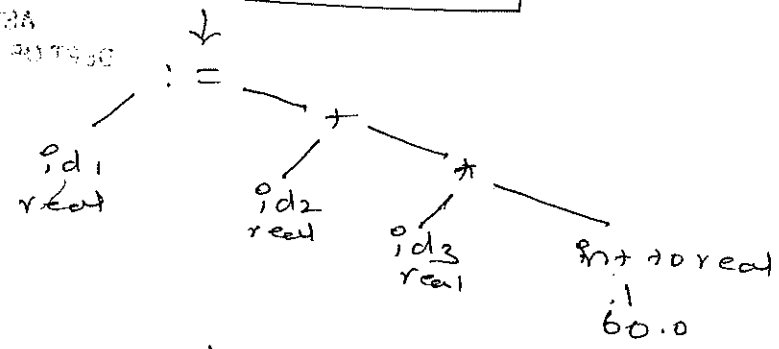


Symbol table

id	addr
Position	
initial	
rate	

Semantic Analyser

S. GURUPRASAD S.  
BMSIT & M  
BMSIT & M  
BMSIT & M



Intermediate Code gen

$T_1 = \text{int to real } (60)$   
 $T_2 = \text{id}_3 * T_1$   
 $T_3 = \text{id}_2 + T_2$   
 $\text{id}_1 = T_3$

Code of optimization

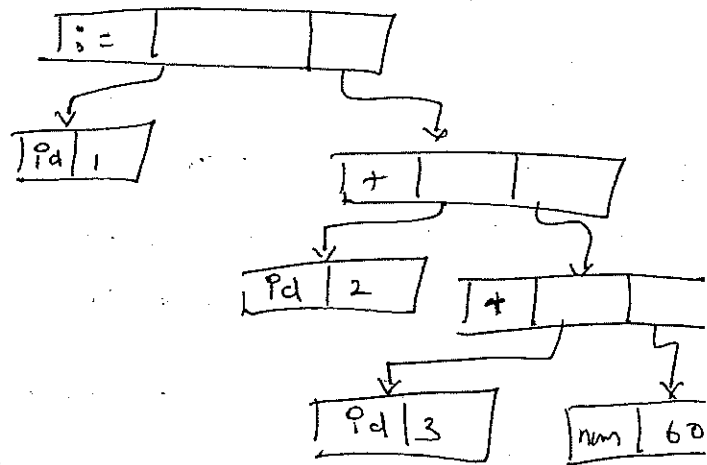
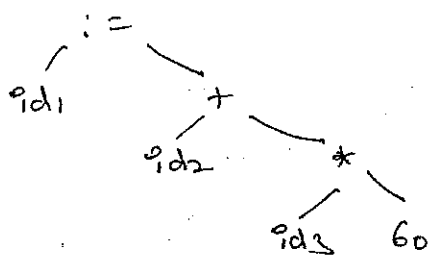
$T_1 = \text{id}_3 * 60.0$   
 $\text{id}_1 = \text{id}_2 + T_1$

Code generation

```

MOV F    id3, R2
MUL F    #60.0, R2
MOV F    id2, R1
ADD F    R2, R1
MOV F    R1, id1
  
```

# Data Structure for Syntax Tree



→ Each operator is a Node with two pointers to its left & right child.

→ A leaf is a record with two or more fields, one to idify token & other to keep track of info about token.

## The Analyse & Synthesise model of Compilation

The entire compilation process is split into two parts.

- ① Analyse +
- ② Synthesise.

→ The Analyse part breaks up the src prog into consistent pieces, analyse it & creates intermediate rep<sup>n</sup> of src prog.

→ The Synthesise part constructs the desired target prog from the intermediate rep<sup>n</sup>.

→ The Synthesise part requires to know abt m/c def<sup>n</sup> & the language specific kernel, so it uses a specialized technique.

**GURUPRASAD, S.**  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-64,  
Ph: 9886760776

**Prof GURUPRASAD S, BMSIT&M**

<https://hemanthrajhemu.github.io>

## Analysis of Src Prog

During analysis, the op<sup>n</sup> implied by src prog are determined & recorded in hierarchical structure called tree / syntax tree where each node rep<sup>n</sup> a op<sup>n</sup> & children rep<sup>n</sup> arguments of op<sup>n</sup>

In compilation, analysis consists of three phases:

### Linear / Lexical Analysis

Stream of characters of src prog is read from left to right & grouped into tokens having collective meaning.

eg:-  $Pos := initial + rate * 60$

LA will generate following tokens.

Pos - id

:= - assignment

initial - id

+

rate - id

\*

60 - num

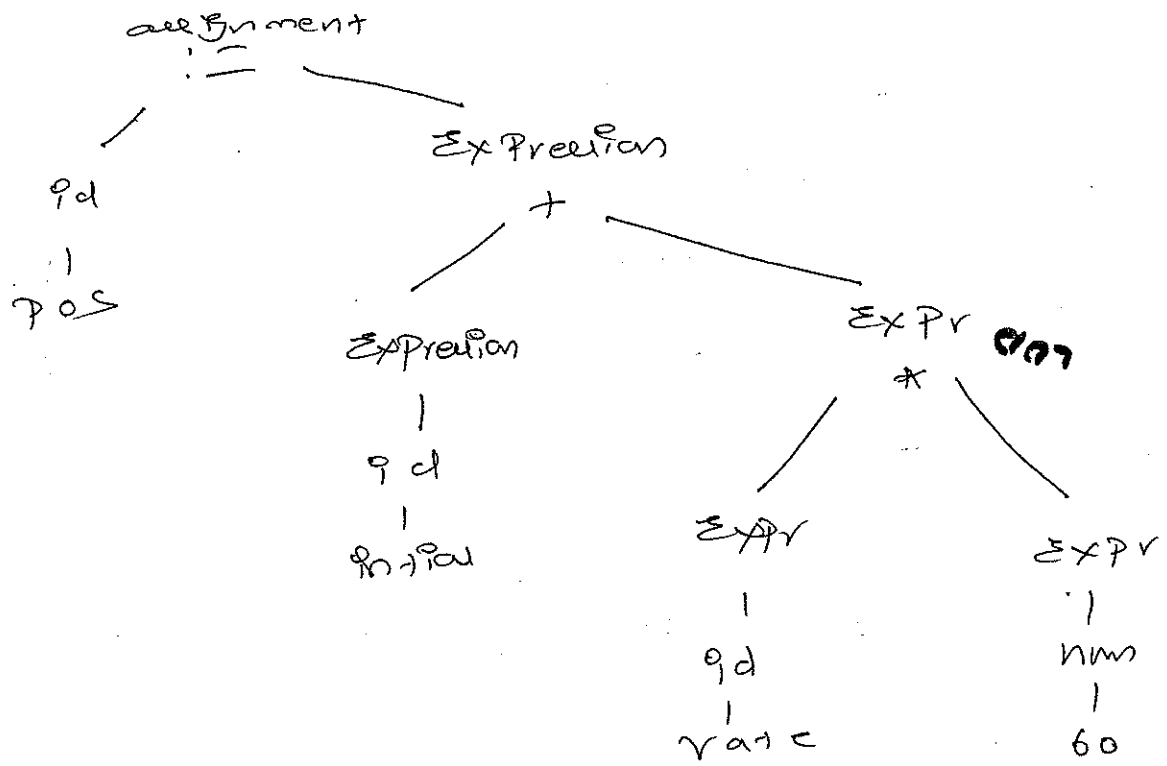
> tabs & space & comments are eliminated.

### Hierarchical Analysis / Syntax Analysis

> Src tokens are grouped together hierarchically into nested collection to identify a syntactical structure of prog

REPRODUCED FROM THE ORIGINAL SOURCE BY THE AUTHOR

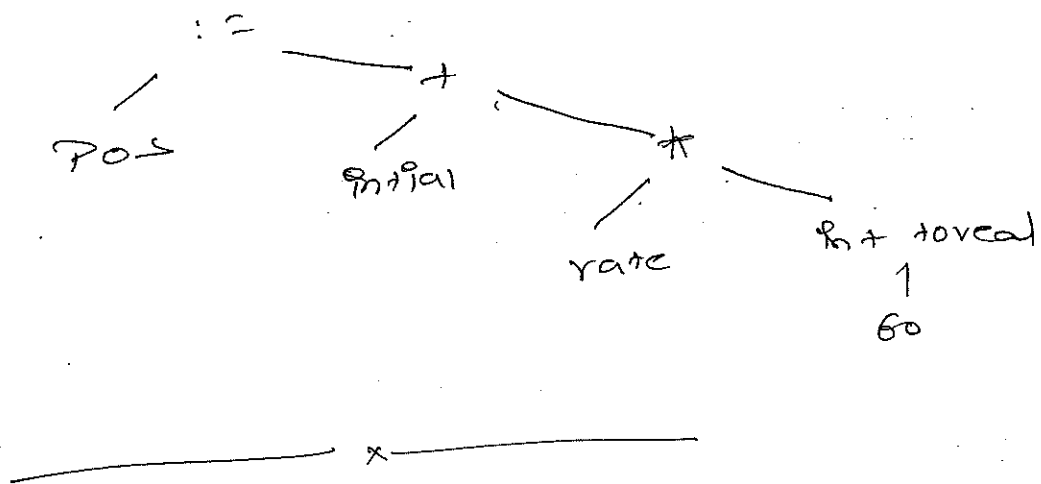
Ex :-



### ③ Semantic analysis

→ It checks the semantics of the identified syntactical structure by performing type checking, scope resolution etc.

Ex :-

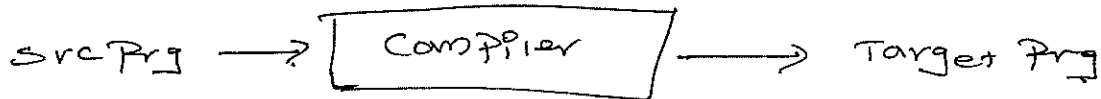


GURUPRASAD. S.  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64,  
 Ph: 9886760776

# The Language Processors

In addition to Compilers several other programs may be required to create executable target programs, like Pre-processor, assembler, loader & linker editors, Virtual machine etc.

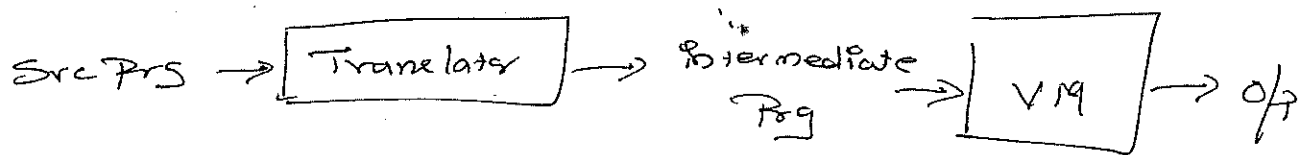
Compiler:



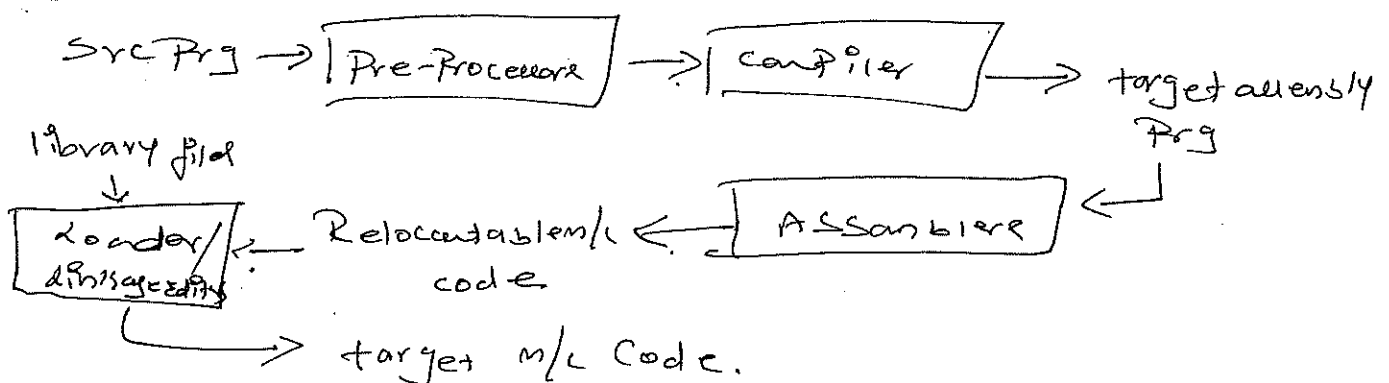
Interpreter:



VM:



Loader/Linkers:



GURUPRASAD S.  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT & M  
BANGALORE-56



## Tools used by Compiler

### ① Structure Editor :-

It takes a sequence of commands as I/P to build a Src Prog, it not only performs text creation & modif'n but also analyses Prog text, putting an appropriate hierarchical structure on the Src Prog.

It does additional tasks such as checking the I/P are correctly formed, can supply keyword automatically.

### ② Pretty Printer :-

It analyses the Prog & prints it in such a way that the structure of Prog becomes clearly visible. by indentation, changing font of comments etc.

### ③ Static Checker :-

It attempts to detect errors/bugs without actually running the Prog.

### ④ Interpreter :-

It converts each line of Src Prog to target Prog & executes it immediately.

## The Evolution of Programming Languages

The first electronic computers appeared in 1940s & were programmed in m/c language by sequence of 0's & 1's. The programming was slow, tedious & error prone.

## the move to Higher Level Languages

The first step towards user-friendly Prog lang was development of mnemonic & Assembly level languages in early 1950s, which converted mnemonic to repn diff op<sup>n</sup> & Assemblers converted mnemonic to target m/c code.

It was easier compared to m/c lang Prog but was m/c dependent mnemonic & instr<sup>n</sup>.

The major step towards HLL was made in latter half of 1950s with the development of languages like FORTRAN, COBOL.

In the following decades many more HLL were created with innovative features to help programming easier & robust.

To day there are thousands of Prog lang. they can be classified in variety of ways one class is by generation.

- I generation — m/c language
- II generation — Assembly language
- III generation — HLL like FORTRAN, COBOL, C, etc
- IV generation — lang designed for specific app<sup>n</sup> like NOMAD  
SQL
- V generation — language applied to logic & constraints like  
Prolog & OPS5

OO — is a lang that supports object oriented programming.

# Impact on Compilers

Since the design of Programming Language & Compiler are intimately related, the advances in Programming Language placed new demands on compiler writers. Compilers needed to promote wide HLL compatibility the extra overhead of Programming Language. A compiler must translate correctly the potentially infinite set of Programs.

## The Science of building Compilers

### 1 Modelling in compiler design & Implementation

The study of compilers is mainly a study of how we design the right mathematical models & choose the right algorithms, while balancing the need of generality & power against simplicity & efficiency.

Some of the fundamental models are finite state machine & Regular Expression which are useful for describing the lexical units of Programs.

CFA is used to describe the syntactic structure of Programming Language.

### 2 The Science of Code Optimization

Optimization refers to attribute that a compiler must be able to produce code that is more efficient than the obvious code.

Compiler optimization must meet the following design objectives.

- \* The optimization must be correct, i.e. it should preserve the meaning of compiled program.
- \* The optimization must improve the performance of many prog.
- \* The compilation time must be kept reasonable
- \* The effort required for opt<sup>n</sup> must be manageable or should be eq<sup>t</sup> to the improve in performance.

## Applications of Compiler Technology.

- \* Implementation of HLL Prog lang.
- \* Optimization for Computer architecture
- \* Design of new computer architectures & interfaces
- \* Program translation
- \* S/W Productivity tools.

## Compiler Construction Tools

- 1) Parser generator - automatically produce syntax analyzer from a grammatical description of Prog lang.
- 2) Scanner generator - produce Lexical Analyze from Reg Exp description of tokens of a language.
- 3) Syntax-directed translation engine - produce collection of subroutines & generate intermediate code.
- 4) Code generator - generate the target code from a collection of rules & m/c registers

5) Data flow analysis Engine - facilitates the gathering of instructions flow & how values are transmitted from one part to another.

## Programming Language: Basics

**GURUPRASAD. S.**  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-66.  
 Ph: 9888760776

### 1) Static/dynamic distinction

The most important issue that we face when designing a compiler for a language is what decisions can compiler make about a program.

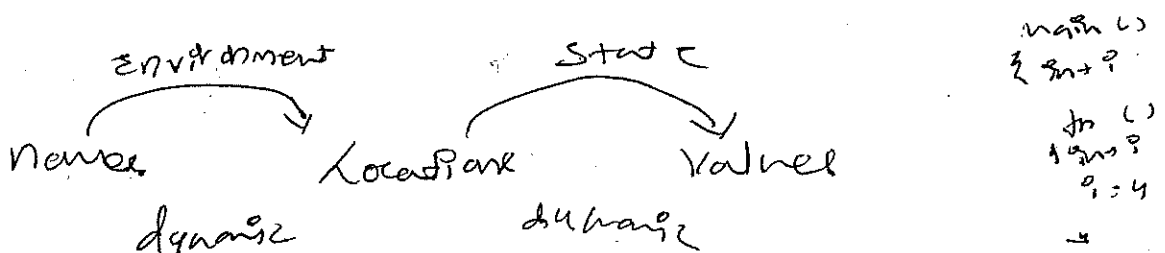
If a language uses a policy that allows compiler to decide on issue, then we say that they use Static Policy. Where the decisions are made during compile time.

If the compiler makes decisions during the execution then it is called Dynamic Policy.

### 2) Environment & State

Another important distinction we must make when discussing programming language is whether changes occurring while program will affect the value of data element or affect the interpretation of name for the data.

The association of names with locations in memory & then with values can be described by two mappings



## 1) Static Scope & Blocs Structure

The static scope are Public, Private & Protected.

**Blocs** are a group of declarations & statements in a file.

## 2) Explicit Access Control

Public, Private & Protected.

## 3) Dynamic Scope ⇒

```
#define a(x+y)
int x = 2;
void b() {int x = 1; printf("%d", a);}
void c() {printf("%d", a);}
void main() {b(); c();}
```

It refers to following policy: A use of name  $x$  refers to the declaration of  $x$  in the most recently called procedure with such a declaration.

Eg:- macro expansion in the 'C' Pre-processor.

## 4) Parameter Passing Mechanism

Call by Value - actual parameters are evaluated & copied into formal parameters. changes will not affect original value.

Call by reference - The address of actual parameters are passed to the called procedure the changes made will affect original value.

Call by name - used in early Algol like  
Algol 60 into call by value

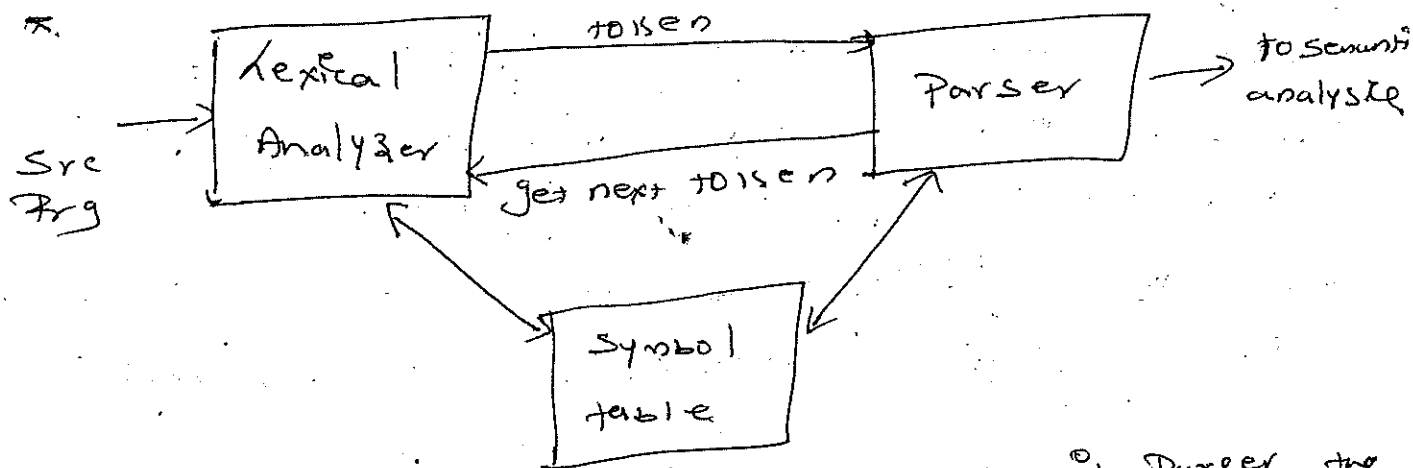
Copy Reference - into call by value but the values are copied back while returning

# Lexical Analyse

## The Role of Lexical Analyser

\* As the first phase of compiler the main task of LA is to read the i/p char from Src Prg, group them into lexeme, and produce o/p as a sequence of tokens for each lexeme.

\* It interacts with symbol table to enter the lexeme for identifiers constants etc. & read the info of certain lexeme.



The fig depicts the interaction of LA with Parser, the Parser sends get next token and the LA will read the char from Src Prg until next token is recognized & is given to Parser.

\* The LA performs stripping out of comments, whitespace

Some times LA is divided into cascade of two processes:

① Scanning consists of simple process that do not require tokenization of i/p such as del<sup>n</sup> of comments, whitespace etc

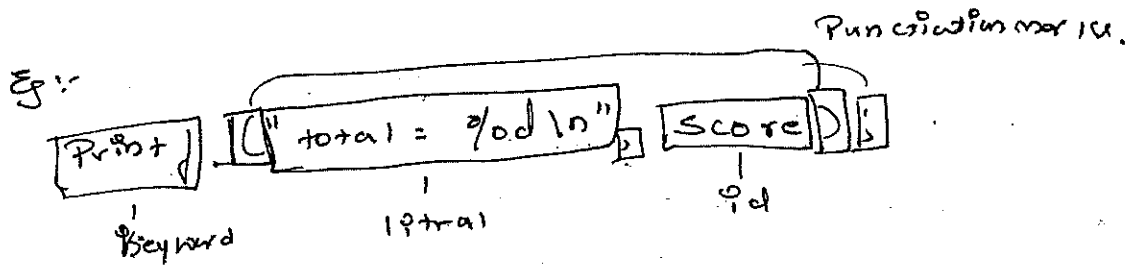
② LA proper is more complex portion where scanner produces the seq. of tokens as output.





## Lexeme:

Is a sequence of characters in the src prog that matches the pattern for a token & is identified by LA as an instance of that token



Each prog lang will have finite set of tokens types

- ① keyword
- ② operators
- ③ Identifiers
- ④ constants, numbers, literals
- ⑤ Punctuation symbols like , ; { } etc

**GURUPRASAD. S.**  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64,  
 Ph: 9886760776

## Attributes of Tokens

When more than one lexeme can match a pattern the LA must provide additional info about lexeme for subsequent phases.

Eg: token type id diff info is associated with it like, type locn first found, initial value etc. which is stored in ST & a ptr to ST will be provided with each id.

Eg:-  $E = M * C * + 2$  the tokens & attributes are

$\langle \text{id, ptr to ST entry for E} \rangle$	$\langle \text{id, ptr to ST entry of C} \rangle$
$\langle \text{assign\_op} \rangle$	$\langle \text{exp\_op} \rangle$
$\langle \text{id, ptr to ST entry for M} \rangle$	$\langle \text{num\_int\_val} \rangle$
$\langle \text{mul\_op} \rangle$	

# Lexical Errors

It is hard for LA to idfy errors <sup>in src code</sup> w/o aid of other components.

eg:  $f(x) = f(x) \dots$

the LA cant tell wheather  $f$  is mterpreting of 'id' or undeclared  $f$  name. in this case LA return  $f$  as id & next phrase detects it as error.

The simplest recovery strategy followed by LA is

"Panic mode" recovery, where compiler delete unrecog. character from i/p till it finds a well formed token.

Other Palliable error recovery actions are:

- ① Delete one character from remaining i/p
- ② Insert a missing character into the remaining i/p
- ③ Replace a character by another character.
- ④ Transpose two adjacent characters.

## Input Buffering

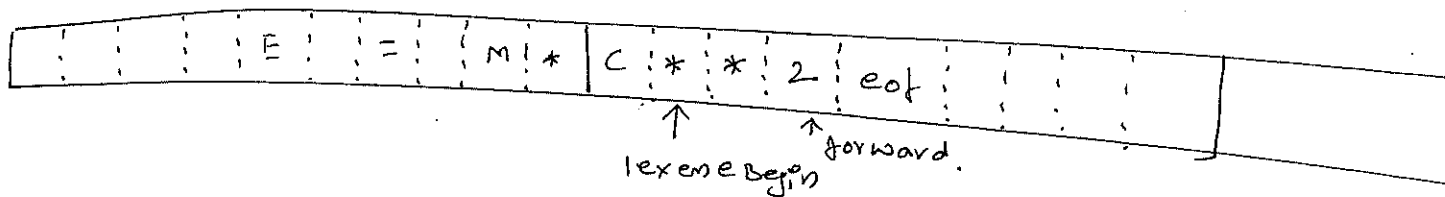
The task of reading the src prog has to be speeded up. but this task is made difficult because the LA must look one or more character beyond the next lexeme before it recognize the right lexeme.

eg:- to idfy an id it should scan till it gets a sym which is not a dig/letter/—

So to accomplish this task two-buffer scheme are used to handle large loops ahead safely.

## ① Buffer Pairs

To avoid overhead required to process single i/p char at a time this scheme uses two buffers, that are alternately reloaded.



Each buffer is of the same size  $N$ , usually (4096 bytes). With one system read command we can read  $N$  characters into a buffer, rather than using one system call per character.

If fewer than  $N$  characters remain in the i/p file, then sp/char. rep<sup>n</sup> by eof.

two pointers to the i/p are maintained.

- ① Pointer lexeme begin, marks the beginning of the current lexeme.
- ② Pointer forward scans ahead until a pattern match found. once a pattern is matched it is returned to parser & begin ptr is set to the next char after pattern & forward ptr is advanced till next

## Sentinels

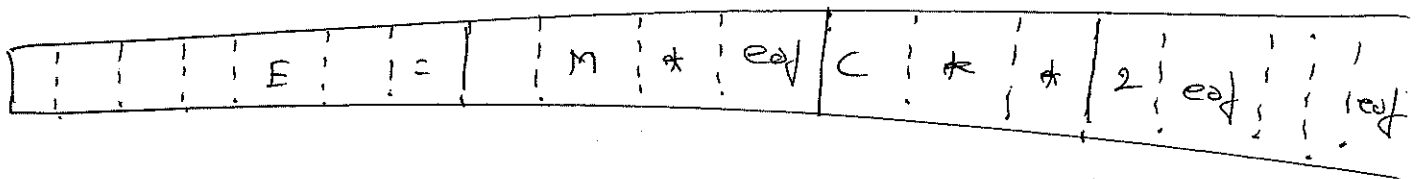
In i/p buffer pair, we must check each time we advanced forward, that we have not moved off one of the buffers, if we do then we must also reload the other buffer, thus for each character read, we make two calls.

\* one for the end of the buffer.

\* one to determine what char is read.

We can combine the buffer-end test with the test for the current char if we extend each buffer to hold a sentinel char at the end.

The sentinel is a special char that cannot be a part of the src Prog like eof.



i.e. one eof for each buffer. to repn end of buffer  
is one eof for end of Pattern

## Specification of Tokens

Regular Expressions are an important notation for specifying lexeme patterns.

## Strings & Languages

Alphabet: is a finite set of symbols  $\Sigma$

eg:-  $\{0, 1\}$  binary alphabet  $\{a-z A-Z\}$  - letter.

$\{0-9\}$  - digit

String:- is a finite sequence of symbols

drawn from alphabet. length of string is repn as  $|s|$

Empty string is denoted by  $\epsilon$  or  $\lambda$

## String-Related terms:

① Prefix of String  $S$ : is any string obtained by removing zero or more symbols from end of  $S$

Eg:  $S = \text{banana}$

Prefix:  $\text{ban}, \text{bana}, \text{ba}, \epsilon$

② Suffix of String  $S$ : is any string obtained by removing zero or more symbols from the beginning of  $S$

Eg:-  $S = \text{banana}$ .

Suffix:  $\text{nana}, \text{na}, \epsilon$

③ Substring of  $S$ : is obtained by deleting any prefix & any suffix of string  $S$ .

④ Proper Prefix, Suffixes & Substrings of  $S$ : are those of  $S$  that are not  $\epsilon$  or  $S$  itself.

⑤ Subsequence of  $S$ : is any string formed by deleting zero or more not necessarily consecutive positions of  $S$

Eg:-  $S = \text{banana}$

Subsequence:  $\text{bnna}, \text{aana}, \text{nn}$

⑥ Concatenation of string: if  $x$  &  $y$  are two strings then concatenation is repn by  $xy$

Eg:  $x = \text{dog}$      $y = \text{house}$      $xy = \text{doghouse}$

## Operations on languages

Language  $L$  - is set of finite strings formed by the fixed alphabet

$$\text{eg:- } \Sigma = \{0, 1\}$$

$$L = \{0, 0, 1, 11, 10, 1010, 0101, \dots\}$$

$$\Sigma = \{a, b\}$$

$$L = \{\text{alpha}, \text{beta}, \text{gamma} \dots a, b, aa, bb \dots\}$$

Def

The closure of a Language  $L$  is denoted by  $L^*$ . It is the set of strings got by concatenating  $L$  zero or more times.

$$L^* = \bigcup_{i=0}^{\infty} L^i \quad \text{or} \quad L = \Sigma^*$$

i.e. it includes all possible strings of language  $L$  and also null string  $\epsilon$

Positive Closure of  $L$  is denoted by  $L^+$   
includes all strings of alphabet except null string

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

Union of  $L$  &  $M$

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

Concatenation of  $L$  &  $M$

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

# Regular Expression

The RE are built recursively out of smaller RE using the following rules.

Each RE  $r$  denotes a language  $L(r)$

the rules that define RE over some alphabet  $\Sigma$  & languages that expr<sup>n</sup> denote.

Rule 1:  $r$   $L(r)$   
 $\epsilon$  is a RE and  $L(\epsilon) = \{\epsilon\}$

Rule 2:  
if  $a$  is a symbol of  $\Sigma$  then  $a$  is a RE where  
 $L(a) = \{a\}$

Rule 3:  
if  $r$  &  $s$  are two REs denoting  $L(r)$  &  $L(s)$  then  
 $(r) | (s)$  is a RE denoting  $L(r) \cup L(s)$

Rule 4:  
if  $r$  &  $s$  are two REs denoting  $L(r)$  &  $L(s)$  then  
 $(r)(s)$  is a RE denoting  $L(r)L(s)$

Rule 5:  
 $r^*$  is a RE denoting  $(L(r))^*$

Ex: -  $\Sigma = \{a, b\}$

$a b$ denotes $\{a, b\}$		$(a b)^*$
$(a b)(a b)$ denotes $\{aa, bb, ab, ba, bb\}$		$L = \{\epsilon, a, b, ab, \dots, aaaa\}$
$a^*$ denotes $\{\epsilon, a, aa, aaaa, \dots\}$		$a a^*b$

Prof GURUPRASAD S. BMSIT&M  
<https://hemanthrajhemu.github.io>

## Algebraic laws of RE

LAW	DESCRIPTION
$rs = sr$	$ $ is commutative.
$r(st) = (rs)t$	Concatenation is associative
$r (st) = (r s)t$	$ $ is associative.

- The unary operator  $*$  has highest precedence & is left associative
- concatenation has second highest precedence & left associative
- $|$  has lowest precedence & left associative.

## Regular Definition

To give names to certain RE & use these names in subsequent expressions.

If  $\Sigma$  is an alphabet of basic symbols, then regular definition is a sequence of definition of the form.

$$\begin{aligned}
 d_1 &\rightarrow r_1 && \text{where } d_i \text{ is new symbol \& } \\
 d_2 &\rightarrow r_2 && r_i \text{ is RE over alphabet } \Sigma \\
 &\dots && \\
 d_n &\rightarrow r_n
 \end{aligned}$$

eg:- identifiers are strings of letters, digits & underscores.

letter  $\rightarrow A|B|\dots|Z|a|b|\dots|z|_$

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

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



Eg: unsigned numbers. (int or float) string.

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

digits  $\rightarrow$  digit digit\*

optional fraction  $\rightarrow$  . digits | E

optional exponent (E (+|-) digits) | E

number  $\rightarrow$  digits optional fraction optional exponent

6.336 E 4      1.89 E -4

## Extension of Regular Expression

Stene introduced RE with basic operators of union & concatenation i.e. | and . but many extensions are added.

① one or more instance.

if  $r$  is an RE then  $(r)^+$  denotes  $(L(r))^+$

$r^+ = r^+ | E$  and  $r^+ = r r^* = r^* r$

② zero or one instance.

if  $r$  is an RE then  $r?$  denotes  $L(r?) = L(r) \cup \{\epsilon\}$

i.e.  $r? = r | E$

③ Character class

$a|b|c$  is denoted by  $[abc]$

and  $a|b|c \dots |z$  is denoted by  $[a-z]$

GURUPRASAD. S.

ASSISTANT PROFESSOR

DEPT OF CSE BMSIT, BANGALORE-64,

Ph: 9886760776

Prof GURUPRASAD S, BMSIT&M

<https://hemanthrajhemu.github.io>

# Summary of RE

RE is specified on alphabet  $\Sigma$ ,  
 following are the REs

$\epsilon$        $\forall \epsilon \in \Sigma^*$        $L(\epsilon) = \epsilon$

$a$        $\forall a \in \Sigma$        $L(a) = \{a\}$

$(r) | (s)$        $\forall r, s \in \Sigma^*$        $L(r) \cup L(s)$

$(r)(s)$        $\forall r, s \in \Sigma^*$        $L(r)L(s)$

$r^*$        $\forall r \in \Sigma^*$        $(L(r))^*$

$r^+$        $\forall r \in \Sigma^*$        $(L(r))^+$

$r?$        $\forall r \in \Sigma^*$        $L(r) \cup \{\epsilon\}$

$[a-z]$        $\forall a \in \Sigma$        $L = \{a | b | \dots | z\}$

Ex 1: letter  $\rightarrow [a-zA-Z_ ]$

digit  $\rightarrow [0-9]$

id  $\Rightarrow$  letter\_ (letter\_ | digit)^\*

Ex 2: digit =  $[0-9]$

digit =  $[0-9]^+$

num = ~~digit (digit)? (E [+ -]^\* digit)?~~

num = digit (digit)? (E [+ -]^\* digit)?

digit      fraction      exponent

# Recognition of Tokens

eg:  $stmt \rightarrow \text{if } expr \text{ then } stmt$   
|  $\text{if } expr \text{ then } stmt \text{ else } stmt$   
|  $\epsilon$

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

$term \rightarrow id \mid number$

id:

letter  $\rightarrow [a-z A-Z]$

dig  $\rightarrow [0-9]$

$id \rightarrow letter (letter|dig)^+$

numbers:

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

number  $= digit (-digit)^? (\epsilon [+ -]^? digit)^?$

$if \rightarrow \text{if } then \rightarrow \text{then } else \rightarrow \text{else}$

$relop \rightarrow < \mid > \mid < = \mid > = \mid = \mid < >$   
LT GT LE GE EQ NE

An intermediate step in construction of LA is to convert RE to Transition Diagram or DFA

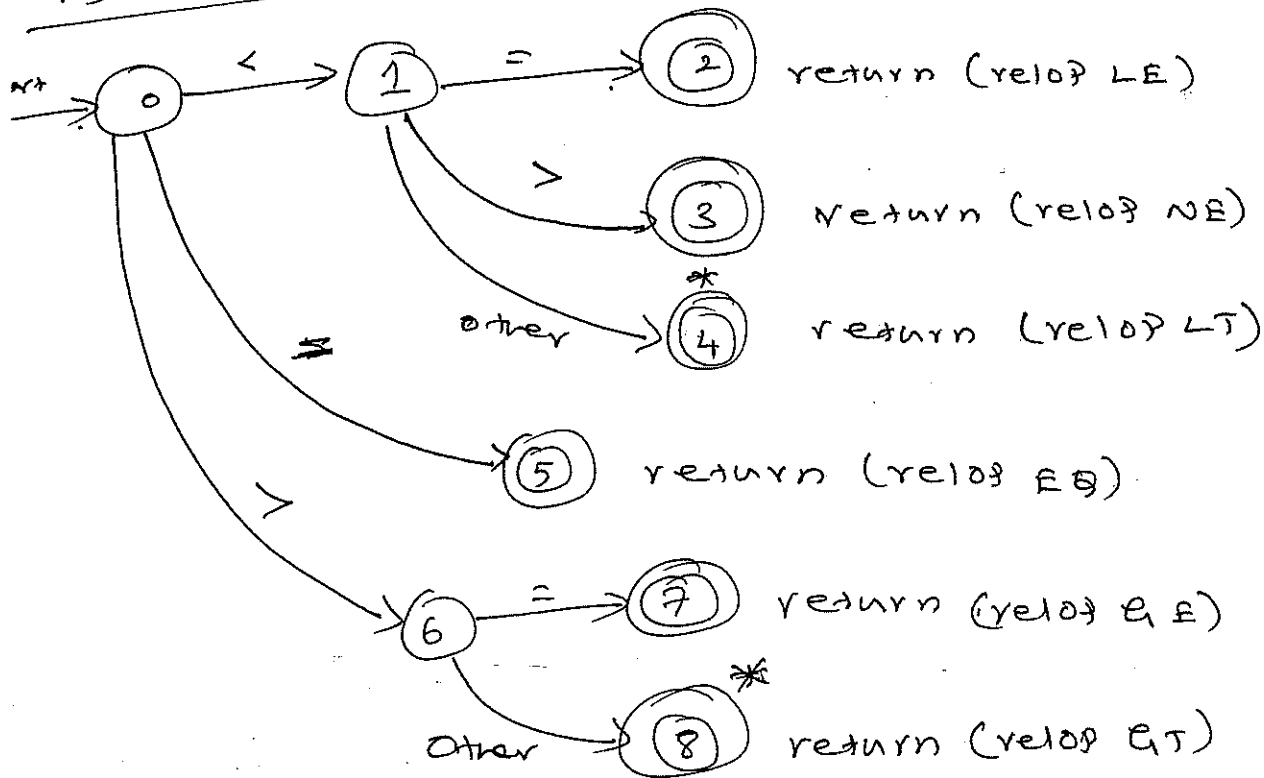
- \* TD have coll<sup>n</sup> of nodes & edges, called states.
- \* Edges are directed from one state to another.
- \* Each edge is labeled by a symbol.
- \* certain states are said to be accepting or final states, which says that lexeme has been found.

**GURUPRASAD. S.**  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-64,  
Ph: 9886760778

**Prof GURUPRASAD S, BMSIT&M**

<https://hemanthrajhemu.github.io>

# TD for Rel op



Indicates retract the forward pointer one position means  
 the basis.

the next step is to implementation of TD

TOKEN getRelop()

```

{
    state = 0;
    for (;;)
    {
        switch (state)
        {
            case 0 : c = getch();
                    if (c == '<') state = 1;
                    else if (c == '=') state = 5;
                    else if (c == '>') state = 6;
                    else fail();
                    break;
            case 1 : c = getch();
                    if (c == '=') state = 2;
                    else if (c == '>') state = 3;
                    else state = 4; break;
        }
    }
}
    
```

**GURUPRASAD. S.**  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64,  
 Ph: 9886760776

Case 2: return ~~AE~~; break;

Case 3: return NE; break;

Case 4: retract E);  
return LT; break;

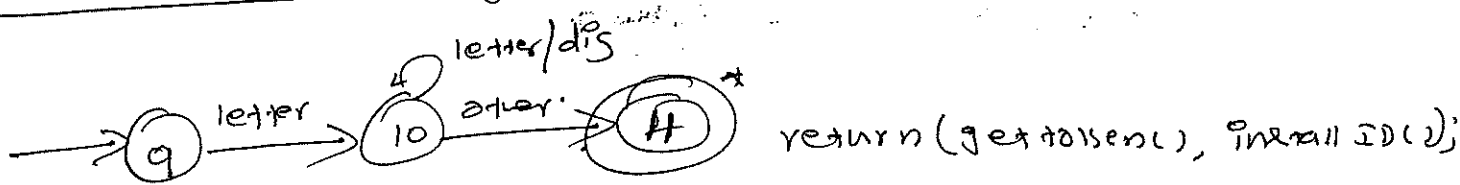
Case 5: return EQ; break;

Case 6: C = getch();  
if (C == '=' ) state = 7;  
else state = 8;  
break

Case 7: return QE; break;

Case 8: retract ();  
return ET;  
break;

## Recognition of Reserved words & Identifiers



In the TD recognize id's & keywords, there are two ways that we can handle reserved words that look like identifiers.

① Install the reserved words in the Symbol Table initially.

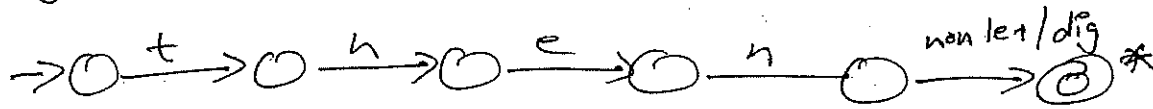
A field of the symbol Table entry indicates that these strings (keywords) are never ordinary identifiers & tells which tokens they represent.

When we find an identifier, a call to `InstallID` places it into ST & returns a pointer to symbol table entry for lexeme found.

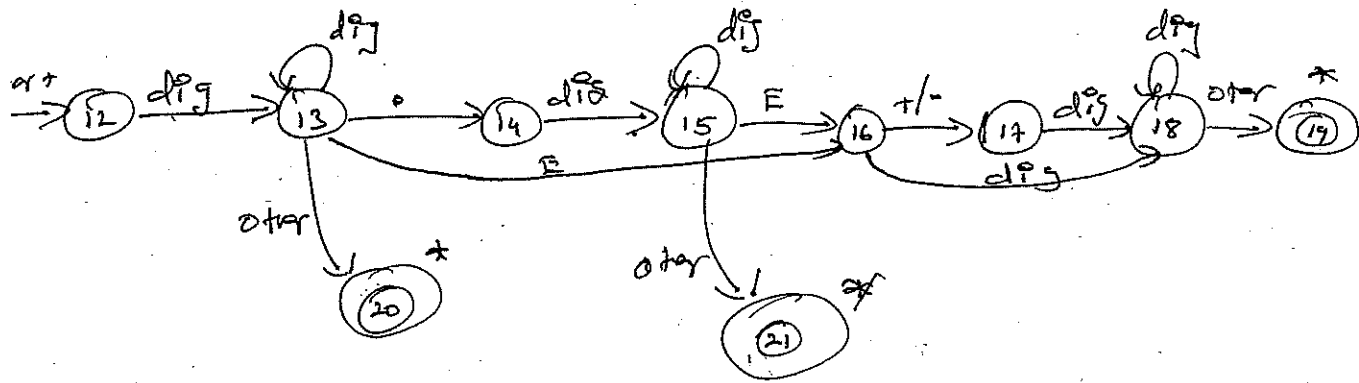
The function `getToken` examines the ST entry for lexeme found & returns whatever token - either id or keyword.

3) Create separate Transition Diagram for each keyword.

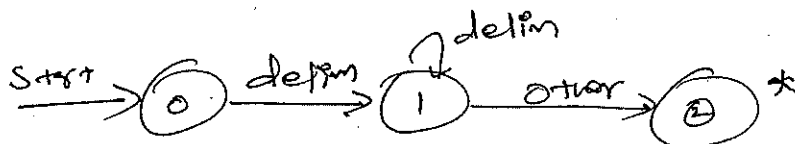
Eg :-



TD for unsigned numbers



TD for white space



delim = (blan|tab|newline)+

\_\_\_\_\_ X \_\_\_\_\_

# SYNTAX ANALYSIS

Objectives:

we learn Parsing methods.

- basic concepts
- techniques for hand impl<sup>n</sup>
- algorithms used by automated tools.
- Error recovery techniques.

Preamble:

Every Prog lang has precise rules that prescribe the syntactic structure of Prog.

Eg:- C Prog is made up of functions, functions out of decl<sup>n</sup> & stmts & stmts out of expressions & so on.

The syntax of Prog lang can be specified using Context Free Grammar CFG or BNF Backus Naur Form. Notations.

Grammar offer significant benefits for both language designers & Compiler Writers.

→ Grammar gives precise yet easy to understand syntactic SFA of Prog lang.

→ From certain clause of grammar we can construct automatically an efficient parser that determines the syntactic structure of Src Prog.

→ The structure of Prog lang designed by grammar will help in translating Src Prog to correct obj code & detect errors.

→ The grammar allows language to be evolved or developed iteratively by adding new constructs to perform new tasks.

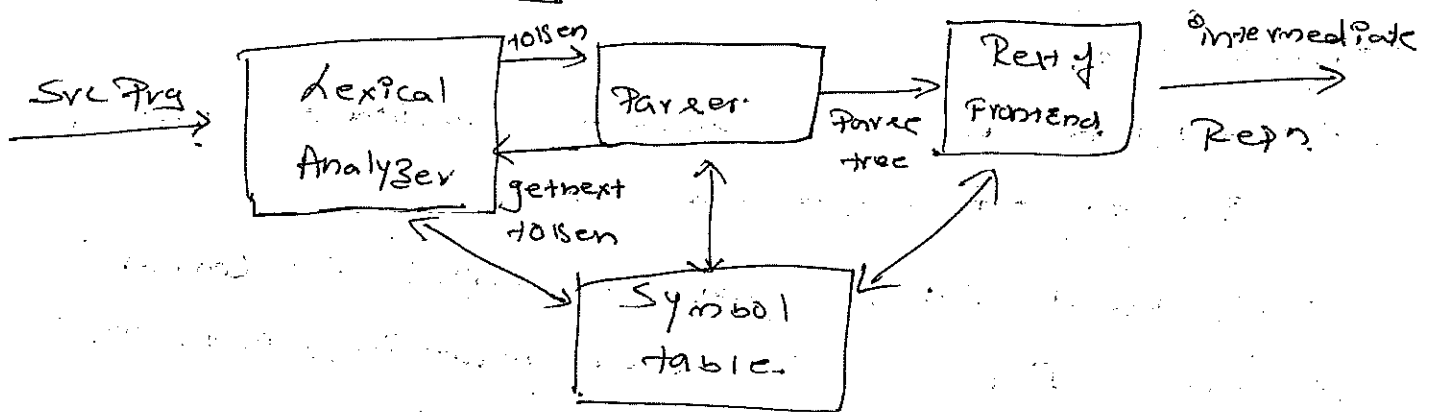
10

**GURUPRASAD. S.**  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-64,  
Ph: 9886760776

## Introduction:

- How Parser fit in Compiler.
- typical grammar for arithmetic expr?
- Parsing technique for expr?
- error handling.

## The Role of Parser:



The Parser obtains strings of tokens from LA, and verifies that string of tokens can be generated by the grammar of src. language.

The Parser should report any syntactical error in an intelligible fashion & to recover from commonly occurring errors to continue processing remainder of program.

The Parser should generate parse tree for well-formed programs & pass it to rest of compiler.

There are three general types of parsers:

- ① Universal Parsers.
- ② Top down Parsers
- ③ Bottom up Parsers



## ① Universal Parsing:

These methods such as Cocise-Vonger-Kalam's algorithm & Earley's algo can parse any grammar but are too in-efficient to use in production compilers.

## ② Top down Parsing

They build the parse tree from Top (root) to the Bottom (leaves).  
The i/p is scanned from left to right one symbol at a time.  
they work for a sub class of grammar,  
eg:- LL Parser.

## ③ Bottom up Parsing:-

They build parse tree from Bottom (leaves) towards Top (root).  
The i/p is scanned from left to right one symbol at a time, works for sub class of grammar.  
eg:- operator precedence, LR.  
uses of BU parsing is constructed w/ automated tools.

## Representative Grammar

The parsing of any lang construct like while int etc are easy as keyword guides the choice of grammar.  
So we concentrate on Expression because of associativity &

### Precedence of operators

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

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

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

These grammar belongs to the class of LR grammar that are suitable for BU parsing.

**GURUPRASAD. S.**  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-64,  
Ph: 9886760776

**Prof GURUPRASAD S, BMSIT&M**

<https://hemanthrajhemu.github.io>

The grammar is Left Recursive so cannot be used for Top Down Parsing, but after eliminating Left Recursion it can be used.

$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

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

GURUPRASAD. S.  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-64  
Ph: 9886760776

## Syntax Error Handling

If compiler had to process only correct programs the design & implementation would be very simple. but practically it is not possible.

Planning the error handling right from start can both simplify compiler structure & improve handling of errors.

Common programming errors that can occur at different levels:

- \* Lexical errors - include misspelling of id, keyword, operators
- \* Syntactic errors - misplaced semicolon, braces, etc.
- \* Semantic errors - type mismatch, improper use of pointer etc.
- \* Logical errors - like use of alignment op = in place of comparison op ==

Several parsing techniques like LL & LR methods detect an error soon as possible.

The goals of error handler in Parser:

- \* Report the presence of error clearly & accurately
- \* Recover from each error quickly enough to detect subsequent errors.
- \* Add minimal overhead to processing of correct programs.

## Error Recovery Strategies

Once the error is detected, the simplest approach for the parser is to quit with an informative error message. When it detects first error, additional errors are detected if it can recover from first error.

The various error recovery strategies are

### Panic mode Recovery:

On discovering an error, the parser discards  $i/j$  symbols one at a time until one of a designated set of tokens are synchronizing tokens are found. Like (i,  $\downarrow$ ) etc.

Panic mode correction often risks a considerable amount of  $i/j$  without checking it for additional errors.

### Phrase-level-recovery

On discovering an error the parser performs local correction on the remaining  $i/j$  by replacing prefix by some string that allows parser to continue.

Typical local correction is to replace comma by semicolon deleting extra semicolon etc

But it is difficult in the situation where error has occurred before the point of detection.

The goal of error handler in Parser:

- \* Report the presence of error clearly & accurately
- \* Recover from each error quickly enough to detect subsequent errors.
- \* Add minimal overhead to processing of correct programs.

## Error Recovery Strategies

Once the error is detected, the simplest approach for the parser is to quit with an informative error message. When it detects first error, additional errors are detected if it can recover from first error.

The various error recovery strategies are

### 1) Panic mode Recovery:-

On discovering an error, the parser discards  $n$  symbols, one at a time, until one of a designated set of tokens are synchronizing tokens are found. i.e. ( ; , ) etc.

Panic mode correction often risks a considerable amount of  $n$  without checking it for additional errors.

### 2) Phrase-level recovery

On discovering an error the parser performs local correction on the remaining  $n$  by replacing prefix by some string that allows parser to continue.

Typical local correction is to replace comma by semicolon deleting extra semicolon etc

But it is difficult in the situation where error has occurred before the point of detection.

## Error Production

By considering common errors, we can augment a grammar with productions that generate erroneous constructs.

When error occurs the parser can generate a parse tree of the erroneous construct that has been recognized.

## Global Correction

They are the algorithms for choosing minimal sequence of changes to obtain a globally least cost correction.

Given an incorrect string  $x$  and grammar  $G$ , the alg will find a parse tree for a related string  $y$  such that no of insertion, deletion & change of  $x$  to get  $y$  is as small as possible.

These methods are in general too costly to ~~execute~~ implement in terms of time & space.

## Context Free Grammars

CFG is used for formal description of structure of strings set of strings constructed using tokens.

Consists of terminals, nonterminals, start symbol & productions.

formally defined as a quadruple  $G = (N, T, P, S)$

$N$  - set of non terminals

$T$  - set of terminals

$P$  - set of productions

$S$  - start symbol

Productions are of the form  $A \rightarrow \alpha$ . Where  $A \in N$  &  
 $\alpha \in (N \cup T)^*$

Notational convention used:

$a, b, c \dots 1, 2, 3, \dots$  used for terminals  
 $+, -, *, = >$  used for operators  
 $A, B, C, \dots$  used for non terminals.

$S$  - is start state

$\Sigma, \beta, \gamma$  rep<sup>n</sup> grammar symbols like  $A \rightarrow \alpha_1 / \alpha_2 \dots$

~~Example~~ Eg:-

G:  
 $E \rightarrow E + T \mid E - T \mid T$   
 $T \rightarrow T * F \mid T / F \mid F$   
 $F \rightarrow (E) \mid id \mid num$

**GURUPRASAD. S.**  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-54.  
 Ph: 9886760776

Derivation:

Construction of parse tree is done by deriving the required string from the grammar rule, in every step the non terminal is replaced by its production.

Eg:- derive string  $ctr + (a-b)/z$ .

$E \rightarrow E + T$   
 $\rightarrow T + T / F$   
 $\rightarrow F + F / id$   
 $\rightarrow id + (E) / id$   
 $\rightarrow ctr + (E - T) / z$   
 $\rightarrow ctr + (T - F) / z$   
 $\rightarrow ctr + (F - b) / z$   
 $\rightarrow ctr + (a - b) / z$

Sem final form of G

sentence of G.

$A \Rightarrow \alpha \beta$  represents that  $\alpha \beta$  is derived from  $A$  in single step

$A \xrightarrow{*} \alpha \beta$  rep<sup>n</sup> that  $\alpha \beta$  is derived from  $A$  in zero or more steps.

$A \xrightarrow{+} \alpha \beta$  rep<sup>n</sup> that  $\alpha \beta$  is derived from  $A$  in one or more steps.

i)  $A \xrightarrow{*} \alpha$  and

$\alpha$  contains only terminals then it is sentence of grammar

$\alpha$  contains both terminals & non terminals then

it is sentinal form of  $G$ .

### Left most Derivation (LMD)

while deriving a string from a grammar  $G$  if we start replacing left most non terminal every time then it is called LMD

rep<sup>n</sup> by  $S \xrightarrow{L} \alpha$ .

Each step in the derivation is called Variable Prefix

& the portion (non-terminal) replaced is called HANDLE.

eg:-

$S \rightarrow S; S$

$S \rightarrow id = E$

$S \rightarrow \text{Print}(L)$

$E \rightarrow id | num$

$E \rightarrow E + E$

$L \rightarrow E$

$L \rightarrow L, E$

$E \rightarrow (S, E)$

Using LMD derive the string:

$X := (Y := 1, Y + 2); \text{Print}(X + 10, Y)$

$S \rightarrow S ; S$

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

$S \rightarrow \text{id} = (S, E) ; S$

$S \rightarrow \text{id} = (\text{id} = E, E) ; S$

$S \rightarrow \text{id} = (\text{id} = \text{num}, E) ; S$

$S \rightarrow \text{id} = (\text{id} = \text{num}, E + E) ; S$

$S \rightarrow \text{id} = (\text{id} = \text{num}, \text{id} + \text{num}) ; S$

$S \rightarrow \text{_____} \parallel \text{_____} ; \text{Print}(L)$

$S \rightarrow \text{_____} \parallel \text{_____} ; \text{Print}(L, E)$

$S \rightarrow \text{_____} \parallel \text{_____} ; \text{Print}(E, E)$

$S \rightarrow \text{_____} \parallel \text{_____} ; \text{Print}(E + E, E)$

$S \rightarrow \text{_____} \parallel \text{_____} ; \text{Print}(\text{id} + \text{num}, \text{id})$

$S \rightarrow \text{id} = (\text{id} = \text{num}, \text{id} + \text{num}); \text{Print}(\text{id} + \text{num}, \text{id})$

$S \xrightarrow{+L} X = (Y = 1, Y + 2); \text{Print}(X + 10, Y)$

### Right most Derivation (RMD)

while deriving string from grammar G if we replace the right most non-terminal every time then it is RMD

reph by  $S \xrightarrow{R} \alpha$

eg:-

**GURUPRASAD. S.**  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64,  
 Ph: 9886760776



$S \rightarrow S; S$

$S \rightarrow S; \text{Print}(L)$

$S \rightarrow S; \text{Print}(L, E)$

$S \rightarrow S; \text{Print}(L, id)$

$S \rightarrow S; \text{Print}(E, id)$

$S \rightarrow S; \text{Print}(E+E, id)$

$S \rightarrow S; \text{Print}(id+num, id)$

$S \rightarrow id = E ; \text{---||---}$

$S \rightarrow id = (S, E) ; \text{---||---}$

$S \rightarrow id = (S, E+E) ; \text{---||---}$

$S \rightarrow id = (S, id+num) ; \text{---||---}$

$S \rightarrow id = (id = E, id+num) ; \text{---||---}$

$S \rightarrow id = (id = num, id+num) ; \text{Print}(id+num, id)$

$S \xrightarrow{R} x = (y = 1, y + 2) ; \text{Print}(x + 10, y)$

Parse Tree:

> A parse tree is a graphical representation of derivation sequence showing how the string is derived from grammar in start state.

> Each interior node of parse tree represents the application

Productions

> No. of non-terminal & leaves represent terminal.

> The string derived from parse tree is called

Yield

Prof GURUPRASAD S, BMSIT&M

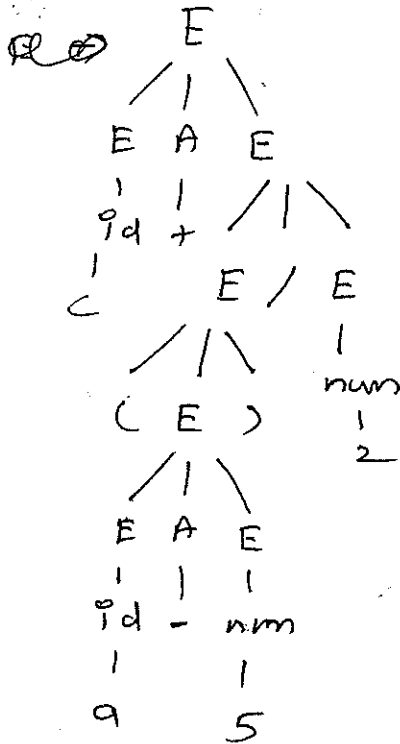
<https://hemanthrajhemu.github.io>

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

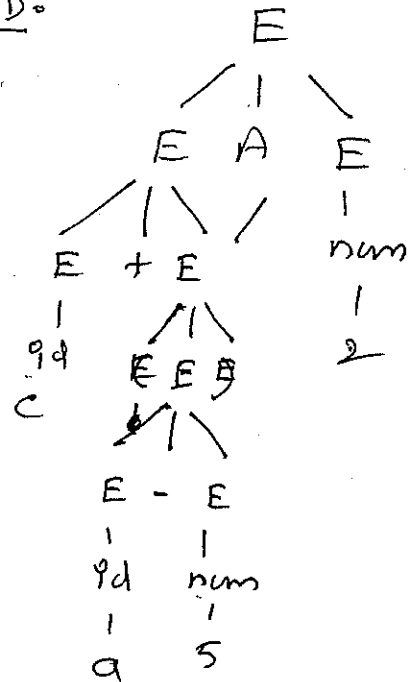
$A \rightarrow + \mid - \mid * \mid /$

Derive parse tree for  $c + (a - 5) / 2$  using RMD & LMD

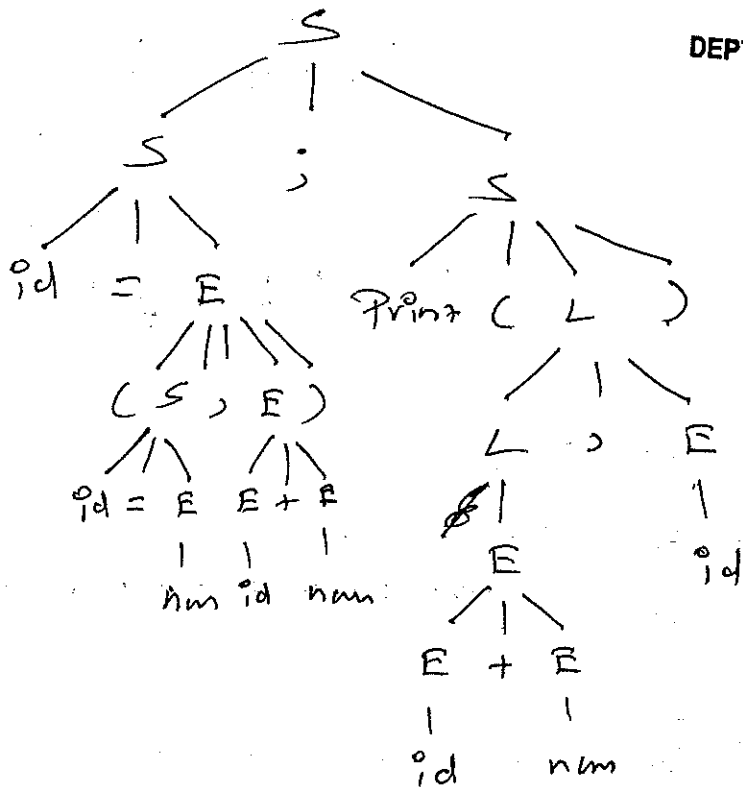
M D:



RMD:



— 2:



**GURUPRASAD. S.**  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64,  
 Ph: 9886760776

See another  
 posts in page  
 and

id / string is obtained by reading leaves from left to right

**Prof GURUPRASAD S, BMSIT&M**

<https://hemanthrajhemu.github.io>

# Ambiguity of Grammar

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

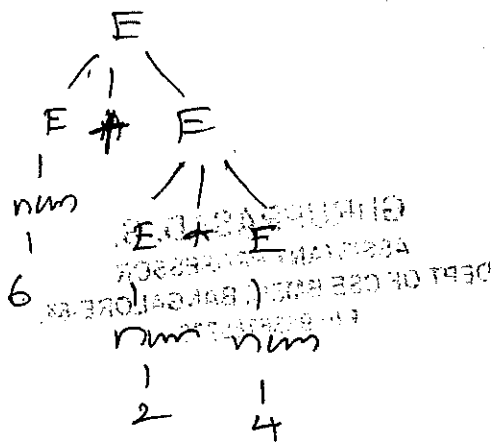
An ambiguous grammar produces more than one LMD or more than one RMD for same sentence.

Ex:-  $G: E \rightarrow EAE \mid (E) \mid id \mid num$

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

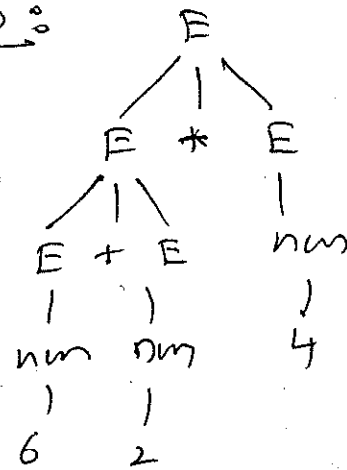
derive string  $6 + 2 * 4$  using LMD & RMD

LMD:



$$6 + (2 * 4) = \boxed{14}$$

RMD:



$$(6 + 2) * 4 = \boxed{32}$$

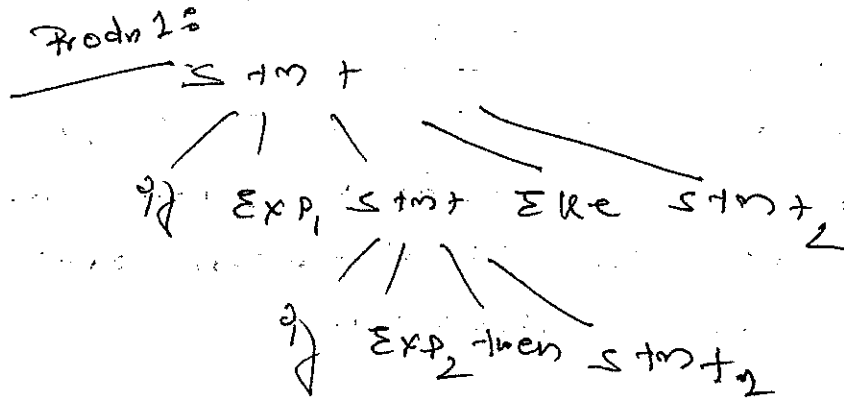
Ambiguous IF:

Context free

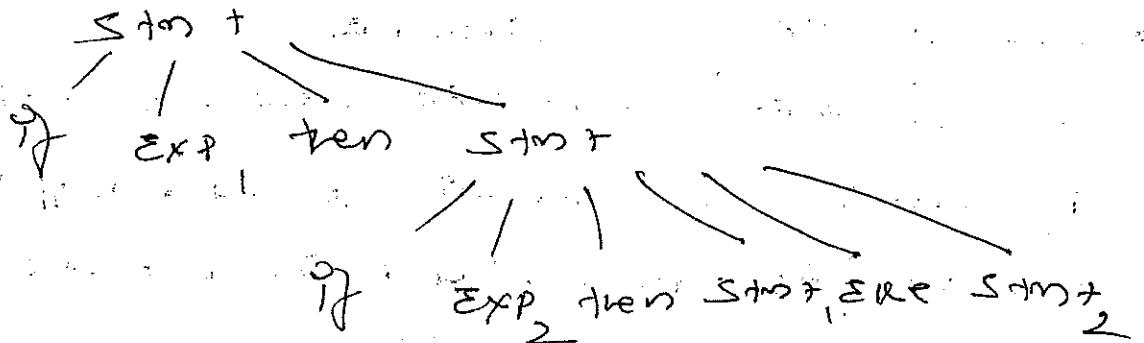
$G: S \rightarrow$  if  $\exists$  parse trees  $T_1, T_2$  such that  $T_1 \neq T_2$  and  $T_1 \Rightarrow w, T_2 \Rightarrow w$

Derive string: if  $\exists$  parse trees  $T_1, T_2$  such that  $T_1 \neq T_2$  and  $T_1 \Rightarrow w, T_2 \Rightarrow w$

2 derivations:



Prodn 2:



Eg:-  
 if  $(x = 0)$  then  
 if  $(y = 1/x)$  then ok := true  
 else  
 z = 1/x  
 end if

**GURUPRASAD. S.**  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64  
 Ph: 9886760776

if the else part is associated with first if stmt then we get divide by zero error

so the else part has to be attached to the second.

if stmt (there is ambiguity)

so Prodn 2 will yield correct string.

Ambiguity of grammar is because of 2 reasons.

① LMD & RMD EXP

② the way grammar is designed.

# CFL vs RE

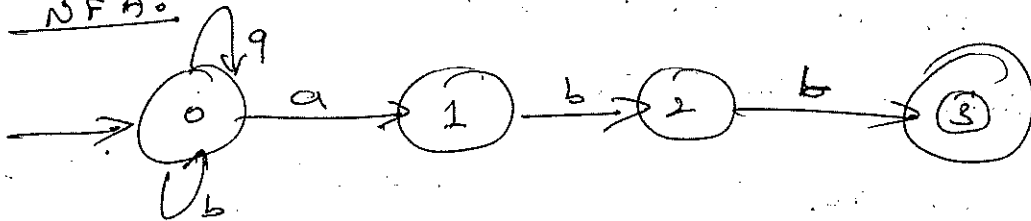
Grammars are more powerful notation than R.E, Every construct that can be described by R.E can be described by the grammar but not vice versa, Every Regular Language is a CFL. but not vice versa.

## Job Convert RE to CFL :-

- \* Construct NFA for the given RE.
- \* for each state  $i$  of the NFA create a non terminal  $A_i$
- \* if state  $i$  has a transition to state  $j$  on i/p  $a$ , add the production  $A_i \rightarrow a A_j$ . If state  $i$  goes to state  $j$  on i/p  $\epsilon$  add the production  $A_i \rightarrow A_j$
- \* if  $i$  is an accepting state add  $A_i \rightarrow \epsilon$ .
- \* if  $i$  is a start state make  $A_i$  a start symbol

Ex: ①  $(a|b)^* abb$

NFA:



CFLS

$$A_0 \rightarrow a A_0 \mid b A_0 \mid a A_1$$

$$A_1 \rightarrow b A_2$$

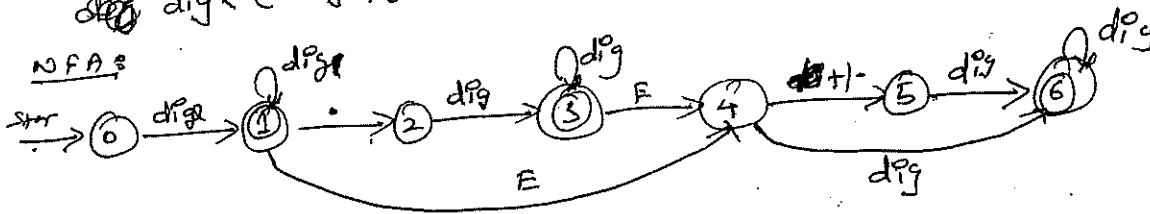
$$A_2 \rightarrow b A_3$$

$$A_3 \rightarrow \epsilon$$

Convert the RE to CFE

Ex 2: ~~dig~~ ~~(dig)~~

~~dig~~ ~~dig~~ (dig)\* (E [+]? dig)\*

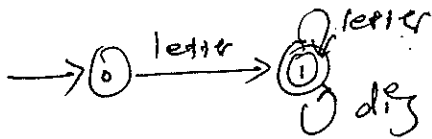


CFEs:

- $A_0 \rightarrow \text{dig } A_1$
- $A_1 \rightarrow \text{dig } A_1 \mid \cdot A_2 \mid \epsilon A_4$
- $A_2 \rightarrow \text{dig } A_3$
- $A_3 \rightarrow \text{dig } A_3 \mid \epsilon A_4$
- $A_4 \rightarrow \text{dig } A_6 \mid + A_5 \mid \cdot A_5$
- $A_5 \rightarrow \text{dig } A_6$
- $A_6 \rightarrow \text{dig } A_6 \mid \epsilon$

**GURUPRASAD. S.**  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64,  
 Ph: 9886760776

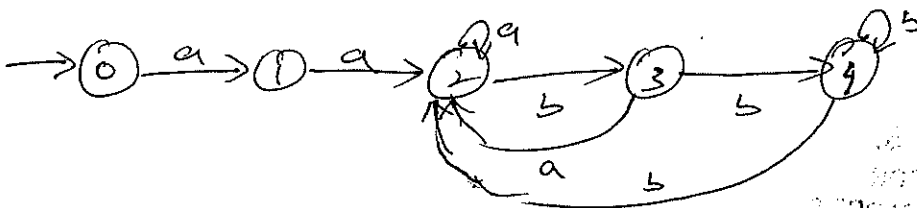
Ex 3: ~~id~~  $\text{id} \rightarrow \text{letter}(\text{letter} \mid \text{dig})^*$



- $A_0 \rightarrow \text{letter } A_1$
- $A_1 \rightarrow \text{letter } A_1 \mid \text{dig } A_1 \mid \epsilon$

Ex 4:

at least 2 a's      ending with 2 b's



GURUPRASAD S.  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64,  
 Ph: 9886760776



## Writing a Grammar:

→ Grammars are capable of describing most, but not all of the syntax of  
Prog lang

Eg:- id must be declared before its use can't be explicitly  
grammar

→ such things are checked by next phase of compiler.

## Lexical v/s Syntax Analyse

→ LA & SA are separated due to modularizing compiler

→ LA rules are quite simple so we use RE not CFG.

→ RE provides concise & easier to understand notation for tokens.

→ Efficient LA can be constructed automatically from RE than  
w/ CFG.

## Notes:

There are several transformations have to be applied to  
grammar to make it suitable for parsing like.

→ Eliminating Ambiguity

→ Eliminating left recursion

→ Left factoring.

## Eliminating Ambiguity

There are two solutions to Eliminate Ambiguity

① Re-write Grammar - for context sensitive Ambiguity

② ADD Precedence - for context free Ambiguity.

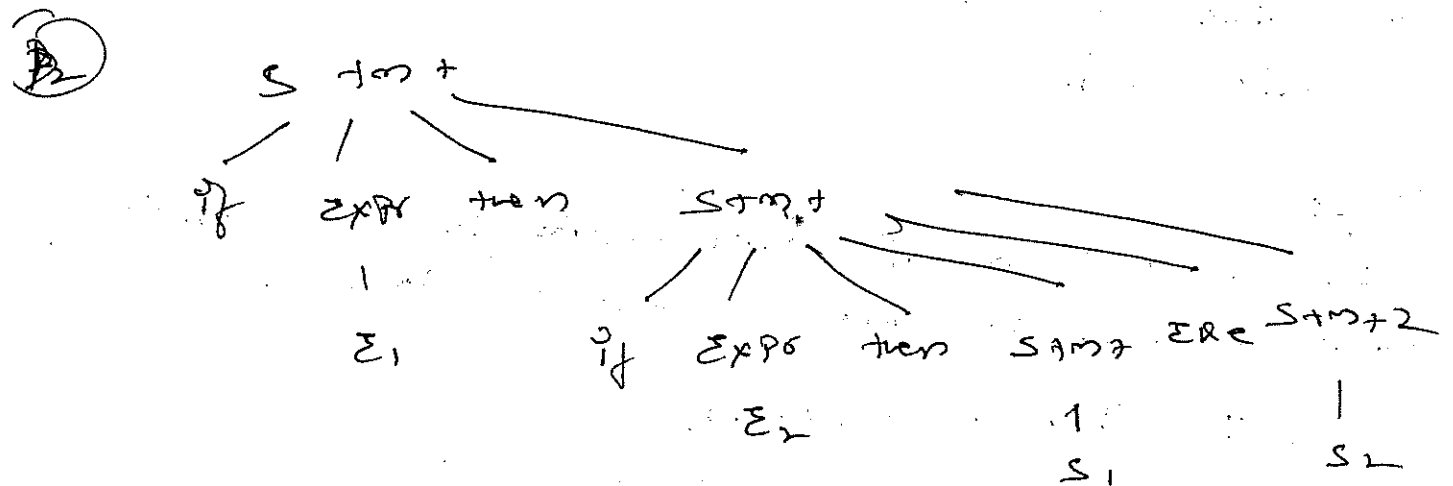
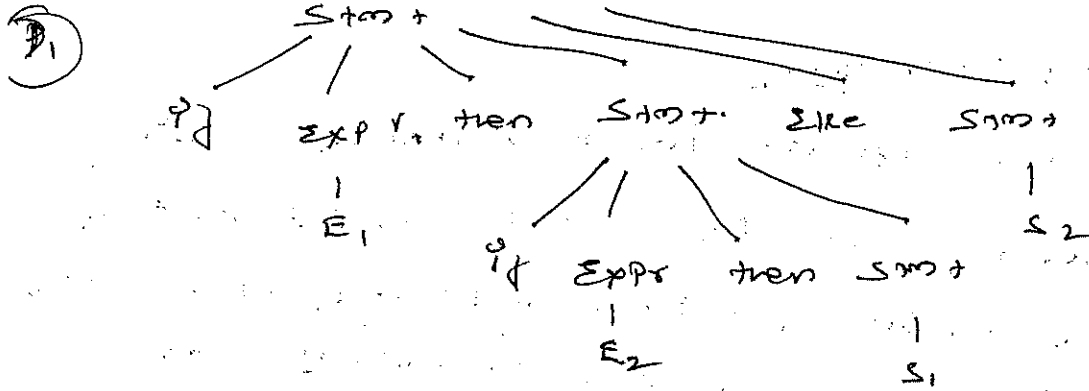


# Rewrite Grammar

eg:  $S \rightarrow \text{if } \text{EXPR} \text{ then } S_1 \text{ |}$

$\text{if } \text{EXPR} \text{ then } S_1 \text{ else } S_2 \text{ | other}$

So string:  $\text{if } E_1 \text{ then if } E_2 \text{ then } S_1 \text{ else } S_2$



Derivn 2 yields correct result the grammar is ambiguous as it has 2 diff parse trees.

So to eliminate ambiguity we re-write the Grammar

The idea is that a  $\text{then}$  appearing b/w  $\text{then}$  &  $\text{else}$  must be 'matched'

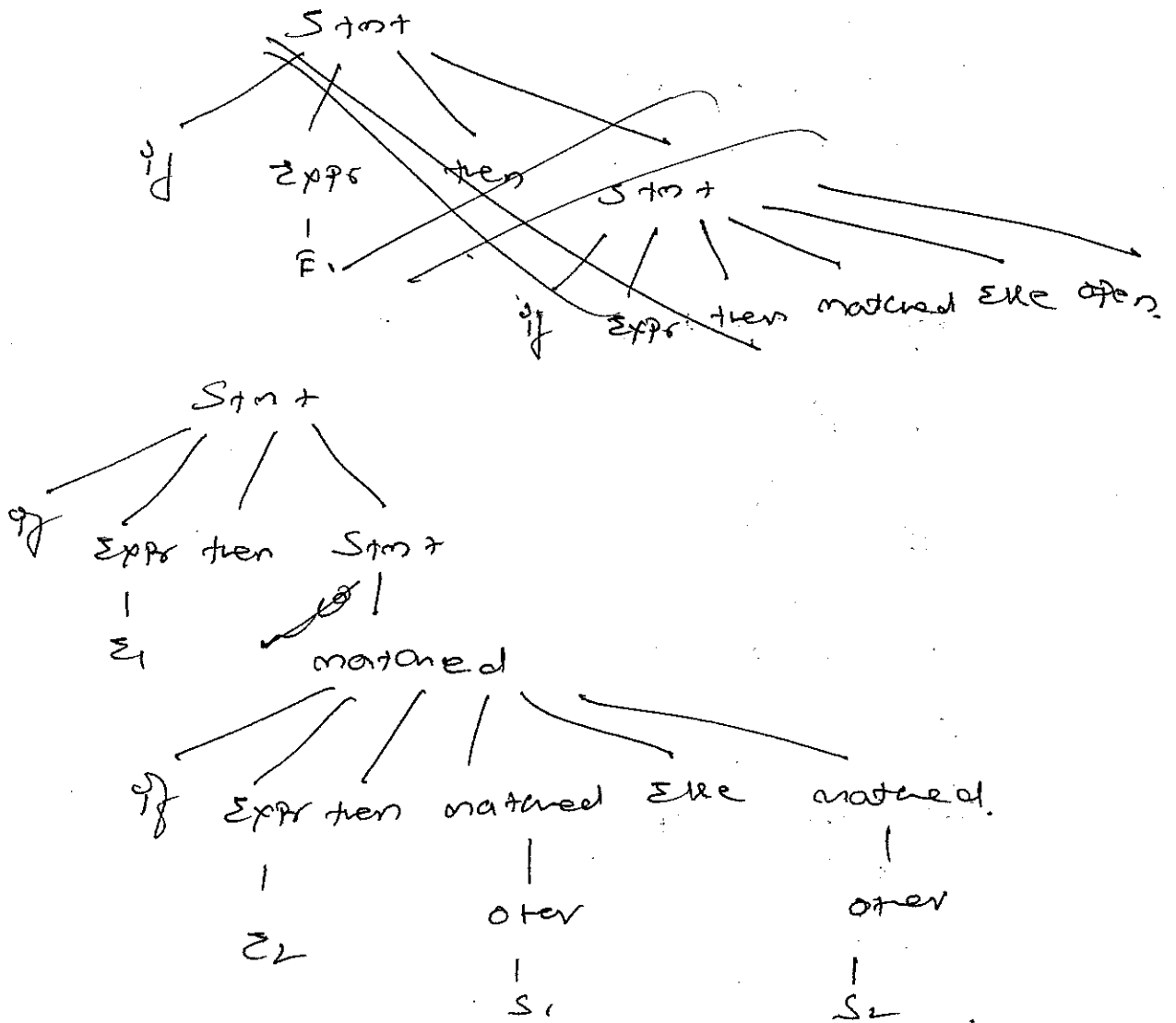
Start  $\rightarrow$  matcher | open

matched  $\rightarrow$  if Expr then matched else matched |  
other

Open Start  $\rightarrow$  if Expr then Start |

if Expr then matched else open

if E<sub>1</sub> then if E<sub>2</sub> then S<sub>1</sub> else S<sub>2</sub>



**GURUPRASAD. S.**  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-64,  
Ph: 9886760776

**Prof GURUPRASAD S, BMSIT&M**

<https://hemanthrajhemu.github.io>

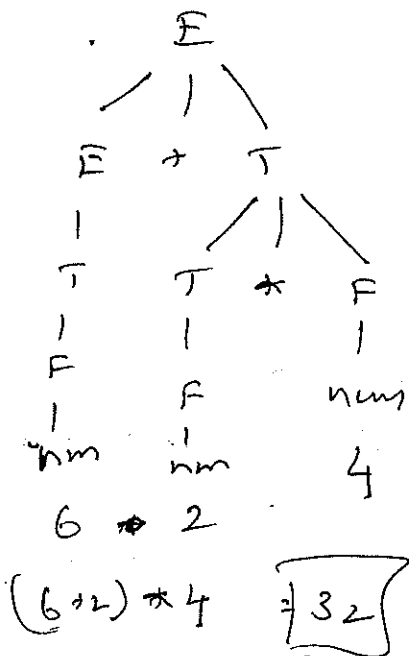
② Add Precedence

eg:-  $E \rightarrow EAE | (E) | id | num$   
 $A \rightarrow + | - | * | / | \uparrow$

if we rewrite the grammar by adding precedence such that the lower precedence operators appear towards root & higher precedence operators towards leaf

$E \rightarrow E + T | E - T | T$   
 $T \rightarrow T * F | T / F | F$   
 $F \rightarrow (E) | id | num.$

String  $6 + 2 * 4$



# ① Eliminating Left Recursion

A grammar is left recursive if it is of the form

$$A \Rightarrow A \alpha | \beta$$

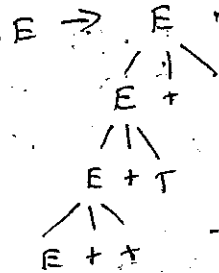
eg:-

$$E \Rightarrow E + T | T$$

$$T \Rightarrow T * F | F$$

$$F \Rightarrow (E) | id$$

because,



So to eliminate left recursion we rewrite the grammar as.

$$A \Rightarrow B A'$$

$$A' \Rightarrow \alpha A' | \epsilon$$

i.e.

$$A \Rightarrow A \alpha_1 | A \alpha_2 | \dots | A \alpha_n | B_1 | B_2 | \dots | B_n$$

↓

$$A \Rightarrow B_1 A' | B_2 A' | \dots | B_n A'$$

$$A' \Rightarrow \alpha_1 A' | \alpha_2 A' | \dots | \alpha_n A' | \epsilon$$

eg:-

$$E \Rightarrow T E'$$

$$E' \Rightarrow + T E' | \epsilon$$

$$T \Rightarrow F T'$$

$$T' \Rightarrow * F T' | \epsilon$$

$$F \Rightarrow (E) | id$$

**GURUPRASAD. S.**  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-54,  
 Ph: 9886760776

**Prof GURUPRASAD S, BMSIT&M**

<https://hemanthrajhemu.github.io>

Ex 2:

$$S \Rightarrow A a | b$$

$$A \Rightarrow A c | s d | \epsilon$$

Substitute  $S$  by  $A a | b$

~~$$S \Rightarrow A a | b$$

$$A \Rightarrow A c | A a d | b d | \epsilon$$

$$A \Rightarrow A c | A a d | b d | \epsilon$$

$$A \Rightarrow A c | A a d | b d | \epsilon$$~~

Eliminate left recursion.

~~$$S \Rightarrow A a | b$$

$$A \Rightarrow A c | A a d | b d | \epsilon$$~~

$$A \Rightarrow A c | b$$

$$\Downarrow$$

$$A \Rightarrow b A'$$

$$A' \Rightarrow c A' | \epsilon$$

$$S \Rightarrow A a | b$$

$$\left. \begin{array}{l} A \Rightarrow A c \\ A \Rightarrow A a d | b d \\ A \Rightarrow \epsilon \end{array} \right\} \begin{array}{l} \longrightarrow \\ \longrightarrow \\ \longrightarrow \end{array} \left. \begin{array}{l} A \Rightarrow A' \\ A' \Rightarrow c A' | \epsilon \\ A \Rightarrow b d A' \\ A' \Rightarrow a d A' | \epsilon \end{array} \right.$$

$$S \Rightarrow A a | b$$

$$A \Rightarrow A' | b d A' | \epsilon$$

$$A' \Rightarrow c A' | a d A' | \epsilon$$

CAJASRIUS LTD  
 ONE OF THE BEST  
 QUALITY PRINTING & PUBLISHING  
 AT BANGALORE

# Left factoring

Even though there is no left recursion in a grammar it may not be parsable for LR(0) LR(1) parser. If it is of the form

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

If  $\alpha$  is returned from LR it is difficult to tell whether to choose  $\alpha B$  or  $\alpha \beta$ .

So to eliminate this we do left factoring way

$$\begin{array}{l} A \rightarrow \alpha A' \\ A' \rightarrow B \mid \beta \end{array}$$

eg:-

$$\begin{array}{l} S \rightarrow aA \mid aB \\ \Downarrow \\ S \rightarrow aS' \\ S' \rightarrow A \mid B \end{array}$$

eg:-

$$\begin{array}{l} S \rightarrow \frac{ict_s}{\alpha} \mid \frac{ictses}{\alpha} \frac{es}{\beta} \frac{a}{\gamma} \\ C \rightarrow b \qquad \qquad \qquad \gamma \quad \beta \quad \gamma \\ \Downarrow \\ S \rightarrow ict_s S' \mid a \\ S' \rightarrow es \mid \epsilon \\ C \rightarrow b \end{array}$$

**GURUPRASAD. S.**  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64,  
 Ph: 9886760776

GURUPRASAD S.  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64,  
 Ph: 9886760776

**Prof GURUPRASAD S, BMSIT&M**

# Top down Parsing

## Left factoring Examples

$$\textcircled{1} S \rightarrow S_1 S_2 \mid S_3 S_4 \mid S_5 S_6$$

$$A \rightarrow \alpha \quad B \quad \gamma$$

$$A \rightarrow \alpha B \mid \alpha \gamma$$

$$\downarrow$$

$$A \rightarrow \alpha A'$$

$$A' \rightarrow B \mid \gamma$$

~~$$S \rightarrow S_1 S_2 \mid S_3 S_4 \mid S_5 S_6$$~~

$$S \rightarrow S_1 S_2 \mid S_3 S_4 \mid S_5 S_6$$

$$S \rightarrow S_1 S_2 \mid S_3 S_4 \mid S_5 S_6$$

$$\textcircled{2} E \rightarrow T + E \mid T$$

$$A \rightarrow \alpha \quad B \quad \gamma$$

$$E \rightarrow T E'$$

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

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

$$\textcircled{3} S \rightarrow id = E \mid id (E) \mid over$$

$$A \rightarrow \alpha \quad B \quad \gamma$$

$$S \rightarrow id, S' \mid over$$

$$S' \rightarrow = E \mid (E)$$

**GURUPRASAD. S.**  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64,  
 Ph: 9886760776

**Prof GURUPRASAD S, BMSIT&M**

<https://hemanthrajhemu.github.io>

# Top down Parsing

It can be viewed as Problem of Constructing a parse tree starting from root in Pre-order.

or  
Equivalent to finding the LMD for all i/p string

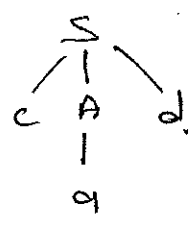
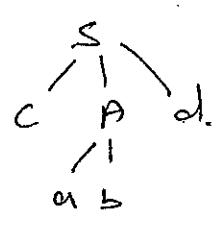
- Start at root
- Repeat until parse tree matches i/p string
- determine the correct production to be applied for a non terminal (key point)
- Proceed constructing left tree by expanding left-most NT. All terminal be reached.

## Recursive Decendent Parsing LL(0)

A RD Parsing Program consists of a set of Procedures one for each non-terminal. Ex<sup>n</sup> begins with Procedure of start symbol, which halts and announces success if its Procedure body reach the entire string.

RD Parser may require back tracking to find correct prod<sup>n</sup> to be applied.

Ex:-  
 $S \rightarrow cAd$  string cad  
 $A \rightarrow ab|a$



we need to backtrack



Eg: Recursive dependent Parenthesis

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) | id$

**GURUPRASAD. S.**  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-84,  
Ph: 9886760776

The Procedures for RD Parenthesis

Procedure E()

```
{
  T();
  EPRIME();
}
```

Procedure EPRIME()

```
{
  ch = getch();
  if (ch == '+')
  {
    T();
  }
  EPRIME();
}
```

Procedure T()

```
{
  F();
  TPRIME();
}
```

Procedure TPRIME()

```
{
  ch = getch();
  if (ch == '*')
  {
    F();
  }
  TPRIME();
}
```

Procedure F()

```
{
  ch = getch();
  if (ch == '(')
  {
    E();
  }
  ch = getch();
  if (ch == ')')
  {
    return TRUE;
  }
  else
  {
    return FALSE;
  }
}
```

else

```
ch = getch();
```

```
if (ch == 'id')
```

```
return TRUE;
```

```
else
```

```
return FALSE;
```

## Predictive Parsing LL (1)

The Parser predicts the next alternative to choose so the grammar need to be unambiguous. i.e. need to define one & only alternative at any point

- (S<sub>1</sub>) So the grammar need to be eliminated with
  - (i) Ambiguity
  - (ii) left factoring
  - (iii) left recursion.
- (S<sub>2</sub>) Two functions are used for finding parse table. so find
  - (i) first symbol
  - (ii) follow symbol.
- (S<sub>3</sub>) Construct parse table
- (S<sub>4</sub>) Parse the given string using stacks.

### First & Follow set

They are the set terminal symbol that help for filling parse table. with valid entry  $\epsilon$  used in error recovery.

→ First: if  $\alpha$  is the string of grammar symbols then  $First(\alpha)$  is the set of terminals that begins the derivation from  $\alpha$ . if  $\alpha \Rightarrow^* \epsilon$  then  $\epsilon \in First(\alpha)$

→ Follow(A) is the set of terminals 'a' that can appear immediately to the right of A such that there exist a derivation of the form  $S \xRightarrow{*} \alpha A \beta$

→ The set of terminals (including  $\epsilon$ ) that can appear at the left of any parse tree derived from a NT is the first of that NT.

→ The set of terminals (excluding  $\epsilon$ ) that follow a NT in any derivn is called FOLLOW set of that NT.

## Rules to create FIRST:

① If  $X$  is a terminal then  $\text{first}(X) = \{X\}$

② If  $X \rightarrow \epsilon$  then  $\epsilon \in \text{first}(X)$

③ If  $X \rightarrow Y_1 Y_2 \dots Y_n$  and  $Y_1 \dots Y_{i-1} \xrightarrow{*} \epsilon$

and  $a \in \text{first}(Y_i)$  then  $a \in \text{first}(X)$

i.e. if  $X \rightarrow Y$  &  $Y \xrightarrow{*} \epsilon$  and if  $a \in \text{first}(Y)$  then  $a \in \text{first}(X)$

~~If  $X \rightarrow Y_1 Y_2 \dots Y_n$  then  $\text{first}(X) = \text{first}(Y_1) \cup \text{first}(Y_2) \cup \dots \cup \text{first}(Y_n)$~~   
Rules to create Follow

① If  $S$  is a start symbol, then  $\$$  is in  $\text{Follow}(S)$

② If there is a prodn of the form  $A \rightarrow \alpha B \beta$  then everything in  $\text{first}(\beta)$  except  $\epsilon$  is placed in  $\text{Follow}(B)$   $\beta \neq \epsilon$

③ If there is a prodn of the form  $A \rightarrow \alpha B$  then everything in  $\text{Follow}(A)$  is in  $\text{Follow}(B)$   $\beta = \epsilon$

④ If there is a prodn of the form  $A \rightarrow \alpha B \beta$  where  $\beta \xrightarrow{*} \epsilon$  then everything in  $\text{Follow}(A) = \text{Follow}(B)$

## Itemsets:

① If  $X \rightarrow Y \alpha$  and  $Y \xrightarrow{*} \epsilon$  then

$$\text{first}(X) = \text{first}(Y) + \text{first}(\alpha)$$

and

if  $\alpha \xrightarrow{*} \epsilon$  then

$$\text{first}(X) = \text{first}(Y) + \text{first}(\alpha) + \text{first}(\alpha)$$

if all RHS produces  $\epsilon$  include  $\epsilon$  by default

in  $\text{first}(X)$

# Summary of first & follow rules

## first rules:

- ①  $X$  is term then  $first(X) = \{X\}$
- ② if  $X \rightarrow \epsilon$   $first(X) = \{\epsilon\}$
- ③  $X \rightarrow Y$   $first(X) = first(Y)$
- ④  $X \rightarrow YZ$   $Y \xrightarrow{*} \epsilon$   $Z \xrightarrow{*} \epsilon$  then  
 $first(X) = first(Y) + first(Z) + first(a)$

## follow rules:

Note: all rules are recursively applied on every prod<sup>n</sup>

- ①  $S \rightarrow \$$
- ②  $A \rightarrow \alpha B \beta$   $B \neq \epsilon$   
 $\uparrow$  first
- ③  $A \rightarrow \alpha B \beta$   $B \xrightarrow{*} \epsilon$   
 $\uparrow$  follow
- ④  $A \rightarrow \alpha B$   $B = \epsilon$   
 $\uparrow$  follow

**GURUPRASAD. S.**  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64,  
 Ph: 9886760776

Prob:  $E \rightarrow E + E \mid (E) \mid a$   $Q: E \rightarrow E * E \mid \epsilon * E \mid (E) \mid a$   
 $A \rightarrow + \mid * \mid ($   $\textcircled{S1}$  is ambiguous so rewrite with precedence

$E \rightarrow E + T \mid T$   
 $T \rightarrow T * F \mid F$   
 $F \rightarrow (E) \mid a$

$\left. \begin{array}{l} \textcircled{S2} \text{ eliminate} \\ \text{left recursion} \end{array} \right\}$

$E \rightarrow T E'$   
 $E' \rightarrow + T E'$   
 $T \rightarrow F T'$   
 $T' \rightarrow * F T'$   
 $F \rightarrow (E) \mid a$

$\textcircled{S3}$  first & follow

Compute the first & follow for the grammar

23)  $E \rightarrow TE' \quad E' \rightarrow +TE' \mid \epsilon \quad T \rightarrow FT' \quad T' \rightarrow *FT' \mid \epsilon \quad F \rightarrow (E) \mid id$

	E	E'	T	T'	F
first	( id	+ $\epsilon$	( id * $\epsilon$	( id	( id
follow	\$ )	\$ )	+ \$ )	+ \$ )	* + \$ )

to compute first:

Non Terminal	Production	First
<del>E</del>	$E \rightarrow TE'$ ③ $T \rightarrow \epsilon$ 4x	$first(E) = first(T)$ $\{ (, id \}$
E'	$E' \rightarrow +TE' \mid \epsilon$	$first(E') = \{ +, \epsilon \}$
T	$T \rightarrow FT'$ ③ $F \rightarrow \epsilon$ no 4x	$first(T) = first(F)$ $\{ (, id \}$
T'	$T' \rightarrow *FT' \mid \epsilon$	$first(T') = \{ *, \epsilon \}$
F	$F \rightarrow (E)$ $F \rightarrow id$	$first(F) = \{ (, id \}$

no rule 4 is applied.

to compute follows

NT	Prodn	Follow
E	$E \rightarrow TE'$ ① $S \rightarrow \$$ ② $A \rightarrow \lambda BB$ ③ $A \rightarrow \lambda B$ ④ $A \rightarrow \lambda$	$follow(E) = \$$ $follow(T) = follow(TE') = +$ $follow(T) = follow(E) = \$$ $follow(E') = follow(E) = \$$
E'	$E' \rightarrow +TE'   \epsilon$ ① $A \rightarrow \lambda BB$ $B \neq \epsilon$ $E' \rightarrow +TE'$ ② $A \rightarrow \lambda BB$ $B \rightarrow \epsilon$ ③ $E' \rightarrow +TE'$ ④ $A \rightarrow \lambda B$ $B = \epsilon$	$follow(T) = follow(TE') = +$ $follow(T) = follow(E') = \$$ $follow(E') = follow(E')$
T	$T \rightarrow FT'$ ① $A \rightarrow \lambda BB$ $B \neq \epsilon$ $T \rightarrow FT'$ ② $A \rightarrow \lambda B$ $B \rightarrow \epsilon$ ③ $T \rightarrow FT'$ ④ $A \rightarrow \lambda B$ $B = \epsilon$	$follow(F) = follow(FT') = *$ $follow(F) = follow(T) = + \$$ $follow(T') = follow(T) = + \$$
T'	$T' \rightarrow *FT'   \epsilon$ ① $A \rightarrow \lambda BB$ $B \neq \epsilon$ $T \rightarrow *FT'$ ② $A \rightarrow \lambda BB$ $B \rightarrow \epsilon$ ③ $T \rightarrow *FT'$ ④ $A \rightarrow \lambda B$ $B = \epsilon$	$follow(F) = follow(FT') = +$ $follow(F) = follow(T) = + \$$ $follow(T') = follow(T) = + \$$

F  
 $F \rightarrow (E)$   
 ①  $A \rightarrow \lambda B$   
 ②  $x$  ③  $x$

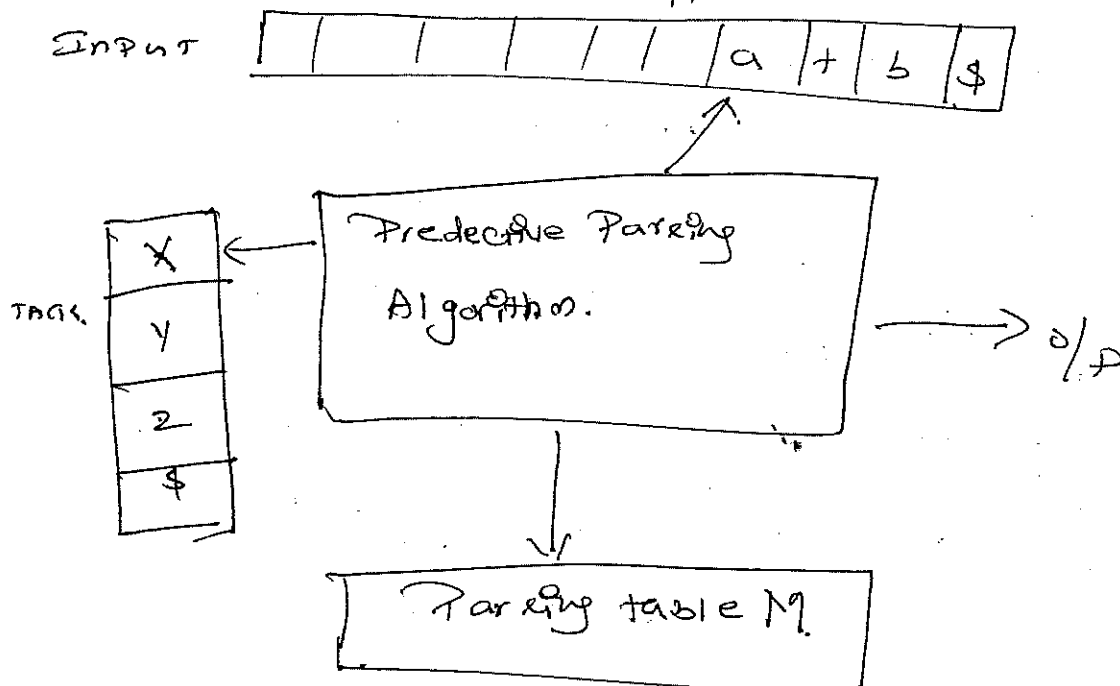
## Non-Recursive Predictive Parser LLD

It is built by maintaining a stack explicitly rather than implicitly via recursive calls.

The Parser produces LMD, if we know  $\theta/p$  that has been matched so far then there holds a sequence of grammar symbols  $\alpha$  such that

$$S \xrightarrow[\text{LMD}]{\alpha} \alpha$$

$\theta/p$



1) Input : The string to be parsed followed by \$

2) Stack : Initially \$ is placed on top of stack, later it may contain grammar symbols

3) Parsing table : It contains 2-D array  $M[A, a]$   
 $\begin{matrix} & \text{NT} & \text{ } \\ & \downarrow & \text{or } \$ \end{matrix}$

Every entry in table is either prod<sup>n</sup> or error.

**GURUPRASAD. S.**  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64.  
 Ph: 9886760776

**Prof GURUPRASAD S, BMSIT&M**

<https://hemanthrajhemu.github.io>

## Predictive Parsing Algorithm

The P.P Program determines  $X$  which is the symbol on top of the stack.  $\epsilon$ . the current i/p symbol. ~~as pop of stack~~ 'a'. these 2 symbols determine the action of parser.

There are 3 possibilities of action

- ① If  $X = a = \$$ , the Parser halts & announces successful completion of parsing.
- ② If  $X = a \neq \$$  then POP  $X$  from the stack & advance the i/p ptr to point next symbol.
- ③ If  $X$  is NT, the program consults the parsing table which may contain  $X$  prod<sup>n</sup> or error entry.  
If  $M[X, a] = X \rightarrow UVW$ , the parser replaces  $X$  on top of the stack by  $UVW$  in reverse order.

## Algorithm to Construct Predictive Parsing Table

i/p : Grammar  $G$

o/p : Parsing table  $M$

Procedure: For each prod<sup>n</sup>  $A \rightarrow \alpha$  of the grammar, do the following

(i) for each terminal  $a$  in  $\text{first}(\alpha)$  add  $A \rightarrow \alpha$  to  $M[A, a]$

(ii) If  $\epsilon \in \text{first}(\alpha)$  then for each terminal  $b$  in  $\text{follow}(A)$  add  $A \rightarrow \alpha$  to  $M[A, b]$  if  $\epsilon \in \text{first}(\alpha)$  &  $\$ \in \text{follow}(A)$  add  $A \rightarrow \alpha$  to  $M[A, \$]$  as well

(iii) Make undefined entries in table as error.



Construct Reductive Parsing table

$M[NST]$

$N$	$T$	$id$	$+$	$*$	$($	$)$	$\$$
$E$		$E \rightarrow TE'$			$E \rightarrow TE'$		
$E'$			$E' \rightarrow +TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
$T$		$T \rightarrow FT'$			$T \rightarrow FT'$		
$T'$			$T' \rightarrow \epsilon$ <del><math>T' \rightarrow \epsilon</math></del>	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
$F$		$F \rightarrow id$			$F \rightarrow (E)$		

(5) Parse the string  $id + id * id$

Parse stack	input	action
$\$ E$	$id + id * id \$$	$E \rightarrow TE'$
$\$ E' T$	$id + id * id \$$	$T \rightarrow FT'$
$\$ E' T' F$	$id + id * id \$$	$F \rightarrow id$
$\$ E' T' id$	$id + id * id \$$	match
$\$ E' T'$	$+ id * id \$$	$T' \rightarrow \epsilon$
$\$ E'$	$+ id * id \$$	$E' \rightarrow +TE'$
$\$ E' T +$	$+ id * id \$$	match
$\$ E' T$	$id * id \$$	$T \rightarrow FT'$
$\$ E' T' F$	$id * id \$$	$F \rightarrow id$

$\$ E' T' id$   
 $\$ E' T'$   
 $\$ E' T' F \star$   
 $\$ E' T' F$   
 $\$ E' T' id.$   
 $\$ E' T'$   
 $\$ E'$   
 $\$$

$id \star id \$$   
 $\star id \$$   
 $\star id \$$   
 $id \$$   
 $id \$$   
 $\$$   
 $\$$   
 $\$$

match  
 $T' \Rightarrow \star FT'$   
 match  
 $F \Rightarrow id$   
 match  
 $T' \rightarrow \epsilon$   
 $E' \rightarrow \epsilon$   
 accept

② Compute first & follow & obtain predictive parse table for the following grammar

$$S \rightarrow AB | PQ \cup$$

$$A \rightarrow \cup y | m$$

$$B \rightarrow bc$$

$$P \rightarrow \cup P | \epsilon$$

$$Q \rightarrow \cup Q | \epsilon$$

$$C \rightarrow e.$$

**GURUPRASAD. S.**  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64,  
 Ph: 9886760776

\* The grammar is not left recursive nor require left factoring

\* Then we need to compute first & follow

First sets

NT	Production	first
S	$S \rightarrow AB$	$first(S) = first(A)$ $\{x, m\}$
$\{x, p, q, m\}$	$S \Rightarrow PQx$	$first(S) = first(P) + first(Q) + first(x)$ $\{x, p, q\}$ none      none
A $\{x, m\}$	$A \Rightarrow xym$	$first(A) = \{x, m\}$
B $\{b\}$	$B \Rightarrow b$	$first(B) = \{b\}$
P $\{p, \epsilon\}$	$P \rightarrow P   \epsilon$	$first(P) = \{p, \epsilon\}$
Q $\{q, \epsilon\}$	$Q \rightarrow qQ   \epsilon$	$first(Q) = \{q, \epsilon\}$
C $\{e\}$	$C \rightarrow \epsilon$	$first(C) = \{e\}$

DEPT OF ELECTRONICS & COMMUNICATIONS ENGINEERING  
 BMSIT - BANGALORE  
 2020-21

	S	A	B	P	Q	C
first	x p q m	x m	b	p ε	q ε	e
follow	\$	\$ b	\$	\$ x	x	\$

Notes - if there are more than one symbol in  $\overline{A}$   
 don't cross  $B \xrightarrow{*} \epsilon$

Follow Set

NT	Production	Follow
S	$S \rightarrow AB$ $A \rightarrow \lambda B$ $B \neq \epsilon$ $S \rightarrow AB$ $B \neq \epsilon$ $A \rightarrow \lambda B$ $B = \epsilon$	$Follow(S) = \{ \epsilon \}$ $Follow(A) = Follow(B) - \epsilon$ $Follow(B) = Follow(S)$
	$S \rightarrow PQX$ $A \rightarrow \lambda B$ $B \neq \epsilon$ $S \rightarrow PQX$ $A \rightarrow \lambda B$ $B \neq \epsilon$	$Follow(Q) = Follow(X) = \epsilon$ $Follow(P) = Follow(Q) + Follow(X)$ $Follow(P) = \epsilon + \epsilon$ <del><math>Follow(P) = \epsilon</math></del>
A	$A \rightarrow xy \cdot m$	x
B	$B \rightarrow bC$ $A \rightarrow \lambda B$	$Follow(C) = Follow(B)$
P	$P \rightarrow pP$ $A \rightarrow \lambda B$	$Follow(P) = Follow(pP)$
Q	$Q \rightarrow qQ$ $A \rightarrow \lambda B$	$Follow(Q) = Follow(qQ)$
C	$C \rightarrow e$	x

GURUPRASAD. S.  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64,  
 Ph: 9886760776

Prof GURUPRASAD S, BMSIT

<https://hemanthrajhemu.github.io>

# Parse table :

NT	2/1 Symbols									
	x	y	m	p	q	e	b	\$	S	
S	<del>S → AB</del> S → PQx		S → AB	S → PQx	S → PQx					
A	A → xy		A → m							
B							B → bC			
P	P → ε			P → pP	P → ε					
Q	<del>Q → ε</del> Q → ε				Q → qQ					
C						C → e				

## Prob 3 :

$$S \rightarrow iC + SS' \mid a.$$

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

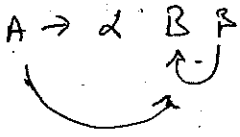
$$C \rightarrow b$$

	S	S'	C
first	{a}	{ε}	{b}
follow	{ε}	{ε}	{ε}

## first symbols

NT	Prodn	first set
S	$S \rightarrow iC + SS' \mid a.$	$\text{first}(S) = \{i, a\}$
S'	$S' \rightarrow es \mid \epsilon$	$\text{first}(S') = \{e, \epsilon\}$
C	$C \rightarrow b$	$\text{first}(C) = \{b\}$

follow set

NT	Prodn	follow set
S	$S \rightarrow i(c+SS')$ $A \rightarrow \alpha B \beta$  $S \rightarrow \frac{i(c+SS')}{\alpha \quad \beta}$ $S \rightarrow \frac{i(c+SS')}{i.B \quad B}$	$\text{follow}(S) = \{\epsilon\}$ $\text{follow}(S) = \text{first}(S')$ $\text{follow}(S) = \text{follow}(S)$ $\text{follow}(S') = \text{follow}(S)$ $\text{follow}(c) = \text{first}(+SS')$

\* Notes if there is more than one symbol in 'B' don't cross it  $\rightarrow \epsilon$

S'	$S' \rightarrow eS$ $A \rightarrow \alpha B \beta$	$\text{follow}(S) = \text{follow}(S')$
----	---	--

Notes: if there are multiple entries in Parse table (PP) then it is not LL(1) grammar

Parse table

NT	T/P					
	i	+	e	a	b	$\phi$
S	$S \rightarrow i(+SS')$			$i \rightarrow a$		
S'			$S' \rightarrow eS$			$S' \rightarrow \epsilon$
C					$C \rightarrow b$	

Grammar:

$S_{\text{stmt}} \rightarrow I \mid S_{\text{stmt}}$

$S_{\text{stmt}} \rightarrow \text{other}$

$I \mid S_{\text{stmt}} \rightarrow \text{if} (\text{Exp}) S_{\text{stmt}} \text{ Else part}$

$E \text{ re part} \rightarrow \text{ere } S_{\text{stmt}} \mid \epsilon$

$x P \rightarrow 0 \mid 1$

next symbols

NT	Production	first
$S_{\text{stmt}}$	$S_{\text{stmt}} \rightarrow I \mid S_{\text{stmt}} \mid \text{other}$	$\text{first}(S_{\text{stmt}}) = \{ \text{first}(I \mid S_{\text{stmt}}) + \{ \text{if other} \} \text{first}(\text{other}) \}$
$I \mid S_{\text{stmt}}$	$I \mid S_{\text{stmt}} \rightarrow \text{if} (\text{Exp}) S_{\text{stmt}} \text{ Else part}$	$\text{first}(I \mid S_{\text{stmt}}) = \{ \text{if} \}$
$E \text{ re part}$	$E \text{ re part} \rightarrow \text{ere } S_{\text{stmt}} \mid \epsilon$	$\text{first}(E \text{ re part}) = \{ \text{ere}, \epsilon \}$
$x P$	$\text{Exp} \rightarrow 0 \mid 1$	$\text{first}(\text{Exp}) = \{ 0, 1 \}$

Note again

If  $A \rightarrow \alpha B \beta$  & B is having more than one symbol then,

only consider rule (2) i.e

$A \rightarrow \alpha B \beta$   $\neq \epsilon$   $\text{follow}(B) = \text{first}(B)$

don't check rule (3) or rule (4)

If there are multiple entries

in parse table then it is not

LL(1) grammar.

	$S_{\text{stmt}}$	$I \mid S_{\text{stmt}}$	$E \text{ re part}$	$\text{Exp}$
$\text{first}$	$\text{if other}$	$\text{if}$	$\text{ere } \epsilon$	$0 \mid 1$
$\text{follow}$	$\$ \text{ ere}$	$\$ \text{ ere}$	$\$ \text{ ere}$	$)$

follow sets

NT	Production	follow
start	$start \rightarrow \underline{I} \overline{f} \overline{s} \overline{t} \overline{m} \overline{t}$ $A \rightarrow \alpha \overline{B}$	① follow(start) = {f} ② follow(I f s t m t) = follow(start)
I f s t m t	$I \overline{f} \overline{s} \overline{t} \overline{m} \overline{t} \rightarrow \overline{I} \overline{f} \overline{s} \overline{t} \overline{m} \overline{t} \overline{E} \overline{k} \overline{e} \overline{p} \overline{a} \overline{r} \overline{t}$ $A \rightarrow \alpha \overline{B}$	$follow(Exp) = \{first(I f s t m t) E k e p a r t\}$ $= \{ \}$
I f s t m t	$I \overline{f} \overline{s} \overline{t} \overline{m} \overline{t} \rightarrow \overline{I} \overline{f} \overline{s} \overline{t} \overline{m} \overline{t} \overline{E} \overline{k} \overline{e} \overline{p} \overline{a} \overline{r} \overline{t}$ $A \rightarrow \alpha \overline{B} \overline{R}$	$follow(start) = \{first(E k e p a r t)\}$
I f s t m t	$I \overline{f} \overline{s} \overline{t} \overline{m} \overline{t} \rightarrow \overline{I} \overline{f} \overline{s} \overline{t} \overline{m} \overline{t} \overline{E} \overline{k} \overline{e} \overline{p} \overline{a} \overline{r} \overline{t}$ $A \rightarrow \alpha \overline{B}$	$follow(start) = follow(I f s t m t)$
E k e p a r t	$E \overline{k} \overline{e} \overline{p} \overline{a} \overline{r} \overline{t} \rightarrow \overline{E} \overline{k} \overline{e} \overline{p} \overline{a} \overline{r} \overline{t} \overline{s} \overline{t} \overline{m} \overline{t}$ $A \rightarrow \alpha \overline{B}$	$follow(start) = follow(E k e p a r t)$

Parse table

	other	{f}	(	)	Ek e p a r t	0	1	\$
start	start → other	start → {f}						
I f s t m t		I f s t m t → {f(Exp)}						
E k e p a r t					E k e p a r t → Ek e p a r t Ek e p a r t → start Ek e p a r t → ε			E k e p a r t → ε
Exp						Exp → 0	Exp → 1	



## Error recovery in Predictive Parser

An error is detected during P.P when ~~terminal~~ <sup>terminal</sup> symbol on top of parse stack does not match the next <sup>terminal</sup> i/p symbol or when non terminal A is on top of the stack & a terminal i/p symbol  $\epsilon$  is in  $M[A, a]$  is error in parse table

## Panic mode Recovery

Panic mode recovery is based on idea of skipping symbols on the i/p until a token in selected set of synchronizing tokens appear. There are two ways to find the synchronizing tokens.

(i) As a starting point, place all symbols in  $follow(A)$  into the synchronizing set of NTA.

If we skip tokens until an element of  $follow(A)$  is seen at top A from stack it is likely that parse can continue.

(ii) If we add symbols in  $first(A)$  to the synchronizing set of NTA then it may be possible to resume parse

eg:-

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

	E	E'	T	T'	F
next	(id	+ E	(id	* E	(id
follow	\$ )	\$ )	+ \$ )	+ \$ )	* + \$ )

Parse string ) id \* + id

Parse stack	I/P	action
\$ E	) id * + id \$	error skip )
\$ E	id * + id \$	$E \rightarrow TE'$
\$ E' T	id * + id \$	$T \rightarrow FT'$
\$ E' T' F	id * + id \$	$F \rightarrow id$
\$ E' T' id	id * + id \$	match
\$ E' T'	id * + id \$	$T' \rightarrow * FT'$
\$ E' T' F *	* + id \$	match
\$ E' T' F	+ id \$	STOP F
\$ E' T'	+ id \$	$T' \rightarrow E$
\$ E' <del>⊙</del>	+ id \$	$E' \rightarrow + TE'$
\$ E' T +	+ id \$	match
\$ E' T	id \$	$T \rightarrow FT'$
\$ E' <del>FF</del>	id \$	$F \rightarrow id$
\$ E' T id	id \$	match
\$ E' T'		$T' \rightarrow E$
\$ E'		$E' \rightarrow E$

Prof GURUPRASAD S. RMSIT&M  
 Srinagar, J.K.

Prof GURUPRASAD S. RMSIT&M

	E	E'	T	T'	F
next	(id	+ E	(id	* E	(id
follow	)	)	+)	+)	* +)

Parse Stack ) id \* + id

Parse Stack	I/P	Action
\$ E	) id * + id \$	error skip)
\$ E	id * + id \$	$E \rightarrow TE'$
\$ E   T	id * + id \$	$T \rightarrow FT'$
\$ E   T' F	id * + id \$	$F \rightarrow id$
\$ E   T' F	id * + id \$	match
\$ E   T' F   *	* + id \$	$T' \rightarrow * FT'$
\$ E   T' F   * +	* + id \$	match
\$ E   T' F   * + id	+ id \$	ST TOP F
\$ E   T' F   * + id +	+ id \$	$T' \rightarrow E$
\$ E   T' F   * + id + )	+ id \$	$E' \rightarrow + TE'$
\$ E   T' F   * + id + ) (	+ id \$	match
\$ E   T' F   * + id + ) ( id	+ id \$	$T \rightarrow FT'$
\$ E   T' F   * + id + ) ( id +	id \$	$F \rightarrow id$
\$ E   T' F   * + id + ) ( id + )	id \$	match
\$ E   T' F   * + id + ) ( id + ) (	id \$	$T' \rightarrow E$
\$ E   T' F   * + id + ) ( id + ) ( +	id \$	$E' \rightarrow E$
\$ E   T' F   * + id + ) ( id + ) ( + )	id \$	match

© GANESH MURUGU  
 2023/11/14 10:34  
 ALL RIGHTS RESERVED TO HIM

Prof GURUPRASAD S, BMSIT&M

# Phrase level Recovery

Just recovery techniques are implemented by filling in the blank entries in the Predictive parse table with pointers to error routines. These routines may change or insert or delete symbols on the  $\$$  or  $\epsilon$  case appropriate error messages they may also pop from stacks.



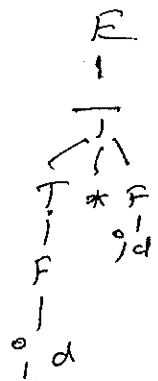
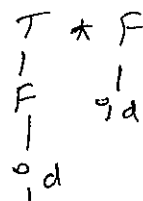
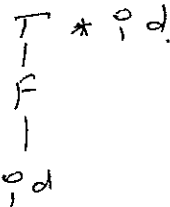
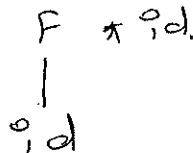
## Unit (3) Bottom up Parsing

A bottom-up parse corresponds to the construction of a parse tree for the  $\$$  string beginning at the leaves (bottom) and working up towards the root (top).

Ex: - G:  $E \rightarrow E + T \mid T$   
 $T \rightarrow T * F \mid F$   
 $F \rightarrow (E) \mid id$

Show the BU parse for  $id * id$

$id * id$



$id + id$

**GURUPRASAD. S.**  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64,  
 Ph: 9886760776

**Prof GURUPRASAD S, BMSIT&M**

<https://hemanthrajhemu.github.io>

## Reduction:

Bottom up Parsing is a process of reducing a string  $w$  to the start symbol of the grammar. At each reduction step a specific sub-string matching the body of the production is replaced by the Non Terminal so the key decision in the BU Parsing is about when to reduce & about what prodn to apply at the Parsing Procedure.

## Handle Pruning

Bottom up Parsing scans symbols from left to right & constructs a RMD in reverse order:

A "HANDLE" is defined as sub string that matches the body of the production and whose reduction represents one step along the reverse of RMD.

eg:-  $id_1 * id_2$

String	Handle	reducing Prodn
$id_1 * id_2$	$id_1$	$F \rightarrow id_1$
$F * id_2$	$F$	$T \rightarrow F$
$T * id_2$	$id_2$	$F \rightarrow id_2$
$T * F$	$T * F$	$T \rightarrow T * F$
$T$	$T$	$E \rightarrow T$
$E$	Start Symbol	so string accepted.

formally  $\text{if } S \xrightarrow{*} \alpha A \omega \Rightarrow \alpha \beta \omega$  then the Prod<sup>n</sup>  
 $A \Rightarrow \beta$  in the Derivation follow  $\alpha$  be a handle of  $\alpha \beta \omega$

## Shift Reduce Parser

There are one of the techniques used in BU Parser.  
 In S-R-Parser we use stacks, i/p buffer, where  
 $\$$  is used to mark the bottom of stack & right  
 end of the i/p.

### Initial Configuration

<u>Stacks</u>	<u>Input</u>
$\$$	$w \ \$$ string

Parser operates by shifting zero or more i/p symbols  
 on to the stacks until an handle is on top of the stacks.

The Parser then reduce  $\beta$  to the LHS of appropriate Prod<sup>n</sup>

The Parser repeated this cycle until it has  
 detected an error or stacks contains start symbol  
 & i/p is empty. then Parser stops & announces successful  
 completion of parsing / string accepted.

### Final Configuration

<u>Stacks</u>	<u>i/p</u>
$\$ S$	$\$$

# The action made by Shift reduce Parser are

- ① Shift : The next i/p symbol is shifted onto the top of stack.
- ② Reduce - The Parser knows that handle is on top of stack & should be replaced with appropriate LHS of Prod<sup>n</sup>.
- ③ Accept - Start symbol on top of stack & i/p buffer is empty.
- ④ Error - the Parser discovers the error & call error recovery routine.

eg:-

$E \Rightarrow E + E \mid E * E \mid (E) \mid id$

Stacks	i/p	action
\$	id + id \$	shift + id
\$ id	* id \$	$E \Rightarrow id$
\$ E	* id \$	shift + *
\$ E *	id \$	shift + id
\$ E * id	\$	$E \Rightarrow id$
\$ E * E	\$	$E \Rightarrow E * E$
\$ E	\$	accept

Ex 2.  $id + id * id$

States	OP	actions
\$	$id + id * id \$$	Shift + id
$\$ id$	$+ id * id \$$	$E \rightarrow id$ reduce
$\$ E$	$+ id * id \$$	Shift + +
$\$ E +$	$id * id \$$	Shift + id
$\$ E + id$	$* id \$$	$E \rightarrow id$ reduce
$\$ E + E$	$* id \$$	$E \rightarrow E + E$ reduce
$\$ E$	$* id \$$	Shift + *
$\$ E *$	$id \$$	Shift + id
$\$ E * id$	$\$$	$E \rightarrow id$ reduce
$\$ E * E$	$\$$	$E \rightarrow E * E$ reduce
$\$ E$	$\$$	accept

Ex 3  $id * id + id$

GURUPRASAD S  
 BMSIT & M  
 BANGALORE  
 INDIA



# Conflict during Shift reduce parser

There are few class of CFG for which SR Parser can't be used.

Every SR Parser for such a grammar can reach a configuration in which the parser cannot decide whether to Shift or to Reduce (Shift-Reduce Conflict) in spite of knowing the entire stack contents & next  $i/p$  symbol.

The other conflict is the parser cannot decide which of the several reductions to make. is known as Reduce-Reduce Conflict

Technically these grammars are not in the LR(1)s, where is the look ahead operator & any ambiguous grammar cannot be LR grammar

## LR Parser

Example:  $S \rightarrow aABe$   $A \rightarrow Abc$   $A \rightarrow b$   $B \rightarrow d$  string  $abbcde$

Parser action	$i/p$	action
\$	abbcde \$	Shift
\$a	bbcde \$	Shift
\$ab	bcd e \$	$A \rightarrow b$
\$a <b>Ab</b>	bcd e \$	Shift
\$a <b>Ab</b>	cde \$	Shift
\$a <b>Abc</b>	d e \$	$A \rightarrow Abc$
\$aA	de \$	

## LR-Parser

These parsers are more general than SP parser

L - scan input from left to right

R - obtain RMD in reverse.

### Advantages:

- LR parser can be constructed to synchronize virtually all parsing language constructs for which CFE exists.
- it is more general non-backtracking SP parser.
- The class of grammar that can be used using LR methods is proper subset of grammar that can be parsed with Predictive parser.
- An LR parser can detect error at soon of possible.

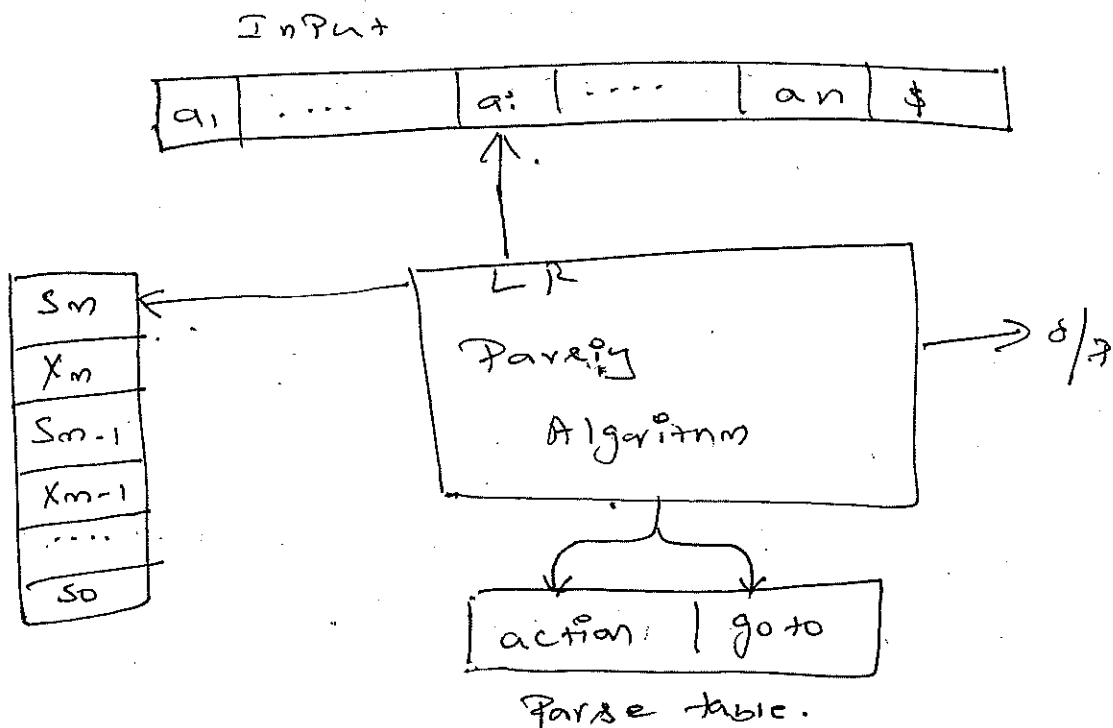
### Disadvantages:

- It needs lot of effort to construct LR parsers by hand for typical programming lang grammar.
- So we need specialized tool for LR parser generation.
- If the grammar contains ambiguous or other constructs it will fail.

There are 3 techniques for constructing LR Parser task

- ① Simple LR (SLR) - easy to implement & least powerful
- ② Canonical LR - most powerful & expensive
- ③ Look Ahead LR (LALR) - is intermediate in power & cost & will work on most practical grammars & can be implemented at moderate effort.

### LR Parsing Algo



\* It contains of an i/p, o/p, stack, driver program & Parse table that has two parts (Action & goto)

\* The driver prog is same for all LR parsers only Parse Table will change.

\* I/P - I/P is read from left to right one symbol at a time. till it ends with  $\$$

\* Stack - contains string of the form  $S_0 X_1 S_1 X_2 S_2 \dots X_n S_n$   
 where  $X_i$  is a grammar symbol &  $S_i$  is a state

Parser Tables: Consist of two parts, parsing action functions & a goto function.

Driver Prg: The driver Prg behaves as follows.

It determines  $S_m$  the state currently on top of the stack &  $a_i$  the current  $i$ th symbol, then consults the action  $[S_m, a_i]$  from parser table which can have one of the four values.

- (i) Shift  $S$  where  $S$  is a state
- (ii) Reduce by a grammar  $A \Rightarrow \alpha$ .
- (iii) Accept
- (iv) Error.

The function goto takes a state & grammar symbol as arguments & produces a state.

Configuration of LR parser is a pair where first component is the stack & whose second component is  $i/p$

$(S_0 X_1 S_1 X_2 S_2 \dots X_m S_m, a_i a_{i+1} \dots a_n \$)$

The next move of the parser is determined by reading  $a_i$ , the current  $i$ th symbol &  $S_m$  the state on top of the stack & then consulting parser action table. entry action  $[S_m, a_i]$

4 types of moves are as follows

- (i) If action  $[S_m, a_i] = \text{Shift } S$  the parser executes a shift move.
- (ii) If action  $[S_m, a_i] = \text{reduce } A \Rightarrow \alpha$  then top 'r' grammar symbols & 'r' state from the stack where 'r' is length of  $R$ .

(iii) If action  $[S_m a_i] = \text{accept}$  Parser is complete.

(iv) If action  $[S_m a_i] = \text{Error}$  the Parser has discovered an error & call an error recovery routine.

~~Construction of LR(0) items~~

### Construction of SLR Parser Table

Item: An LR(0) item or item of a grammar  $G$  is a Prod<sup>n</sup> of  $G$  with a dot at some pos<sup>n</sup> on the RHS of Prod<sup>n</sup>.

eg:-  
 $A \rightarrow XYZ$   
 $A \rightarrow \cdot XYZ$   
 $A \rightarrow X \cdot YZ$   
 $A \rightarrow XY \cdot Z$   
 $A \rightarrow XYZ \cdot$

if  $A \rightarrow \epsilon$   
then  
 $A \rightarrow \cdot$

Item indicates how much of the Prod<sup>n</sup> we have seen at a given pt. in parsing process.

eg:-  $A \rightarrow X \cdot YZ$  means we have seen  $Y$  &  $Z$  derivable from  $X$  & we have to see the string derivable from  $YZ$ .

The coll<sup>n</sup> of LR(0) items is called canonical LR(0) coll<sup>n</sup> for grammar.

The grammar  $G$  is augmented before LR(0) items by giving new start state  $S'$  where  $S' \rightarrow S$  which indicates Parser when it should stop parsing and announce all after it.

GURUPRASAD. S.  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-64.  
Ph: 9886760776

Prof GURUPRASAD S, BMSIT&M

<https://hemanthrajhemu.github.io>

## Algorithm to construct LR(0) items

Procedure  $items(G')$

begin

$C = \text{closure}(S' \rightarrow \cdot S)$

repeat

for each item  $I$  in  $C$  & each grammar symbol  $X$   
such that  $GOTO(I, X)$  is not empty. ~~if not~~

{

add  $GOTO(I, X)$  to  $C$

}

until no more sets of items can be added to it.

End.

## The Closure of $I$

If  $I$  is a set of items for grammar  $G$  then  $\text{Closure}(I)$  is a set of items constructed from  $I$  by the following:

(i) Initially every item in  $I$  is added to  $\text{Closure}$

(ii) If  $A \rightarrow \alpha \cdot B\beta$  is in  $\text{Closure}(I)$  &  $B \rightarrow \gamma$  is a production add the item  $B \rightarrow \cdot \gamma$  to  $I$  if it is not already there. Apply this rule until no more new items can be added to  $\text{Closure}(I)$ .

## The Junctions

Begin

$J = I$

repeat

for each item  $A \rightarrow \alpha \cdot B\beta$  in  $J$  & each prodn  $B \rightarrow \gamma$

of  $G$  such that  $B$  is not in  $J$  do

{  
add  $B \rightarrow \cdot \gamma$  to  $J$

}

Prof GURUPRASAD S, BMSIT&M

<https://hemanthrajhemu.github.io>

until no more items can be added to  $I$

return  $J$

end

Goto operation

$\text{goto}(I, X)$  where  $I$  is set of items  $X$  is grammar symbol

$\text{goto}(I, X)$  is defined to be the closure of set of all items

$[A \rightarrow \alpha X \beta]$  such that  $[A \rightarrow \alpha \cdot X \beta] \in I$ .

$\text{goto}(I, X) = \text{closure}(A \rightarrow \alpha X \beta) \mid A \rightarrow \alpha \cdot X \beta \in I$

Algorithm to construct SLR Parity table

I/P: Canonical colln. of set of items  $\text{follow}(I)$

O/P: SLR Parity table consisting of ACTION & GOTO

Method:

Let  $\{I_0, I_1, \dots, I_n\}$  be the set of items the states of parser are  $0, 1, 2, \dots, n$  where state  $i$  is constructed from  $I_i$

The parity actions for state  $i$  are determined as follows.

Step 1: If  $A \rightarrow \alpha \cdot a \beta$  in  $I_i$  &  $\text{goto}(I_i, a) = I_j$  then set action to Shift  $j$  where  $a$  is terminal or non terminal

Step 2: If  $A \rightarrow \alpha \cdot$  is in  $I_i$ ; then set action  $[i, a]$  to

Reduce  $A \rightarrow \alpha$  for each  $a$  in  $\text{follow}(A)$

Step 3: If  $S' \rightarrow s$  is in  $I_i$ ; then set action  $[i, \$]$  to accept

where  $S'$  is start symbol

Prob:-  $S' \Rightarrow S$

$S \Rightarrow (L)$

$S \Rightarrow \alpha$

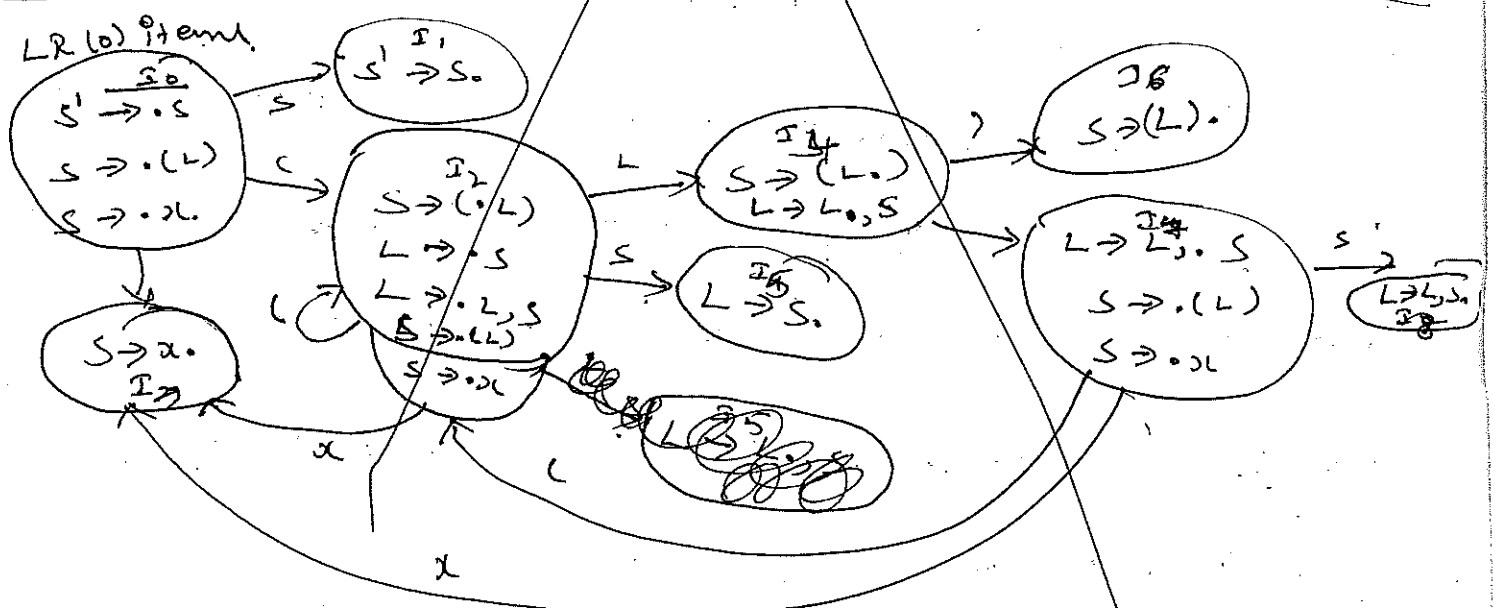
$L \Rightarrow S$

$L \Rightarrow L, S$

	$S$	$L$
first	$( \alpha$	$( \alpha$
follow	$\$ ) ,$	$) ,$

compute first & follow set

NT	first	follow
$S$	$S \rightarrow (L)$ $A \rightarrow \alpha B$	$follow(S) = first( )$
$L$	$L \rightarrow S$ $A \rightarrow \alpha B$	$follow(L) = follow(S)$
$L$	$L \rightarrow L, S$ $A B A$	$follow(L) = first(S)$
	$L \rightarrow L, S$ $A \alpha B$	$follow(S) = follow(L)$





Construct Parse table

	action				ε <sub>0</sub> → ε <sub>0</sub>		
	(	)	,	x	\$	S	L
0	S <sub>2</sub>			S <sub>3</sub>		1	
1					accept		
2	S <sub>2</sub>			S <sub>3</sub>		5	4
3		r <sub>2</sub>	r <sub>2</sub>		r <sub>2</sub>		
4		S <sub>6</sub>	S <sub>7</sub>				
5		r <sub>3</sub>	r <sub>3</sub>				
6		r <sub>1</sub>	r <sub>1</sub>		r <sub>1</sub>		
7	S <sub>2</sub>			S <sub>3</sub>		8	
8		r <sub>4</sub>	r <sub>4</sub>				

Parse string (x, x)

Stacks	I/P	action
\$	(x, x) \$	S shift
\$ (2	x, x) \$	S shift
\$ (2 x 3	, x) \$	reduce S → x
\$ (2 x 3 x x	, x) \$	reduce L → S
\$ (2 S 5	, x) \$	S shift
\$ (2 L 4	x) \$	S shift
\$ (2 L 4, 7	) \$	reduce S → x
\$ (2 L 4, 7 x 4	\$	reduce S → ( )
\$ (2 L 4, 7) x 4	\$	
\$		

GURUPRASAD S  
 ASSISTANT PROFESSOR  
 DEPT OF CSE & I.T. ENGINEERING  
 BMSIT

Prob 1%

(S1) augmented grammar.

Grammar :

- $S' \rightarrow S$
- 1)  $S \rightarrow (L)$
- 2)  $S \rightarrow \alpha$
- 3)  $L \rightarrow S$
- 4)  $L \rightarrow L, S$

$$first(S) = \{ ( \alpha \}$$

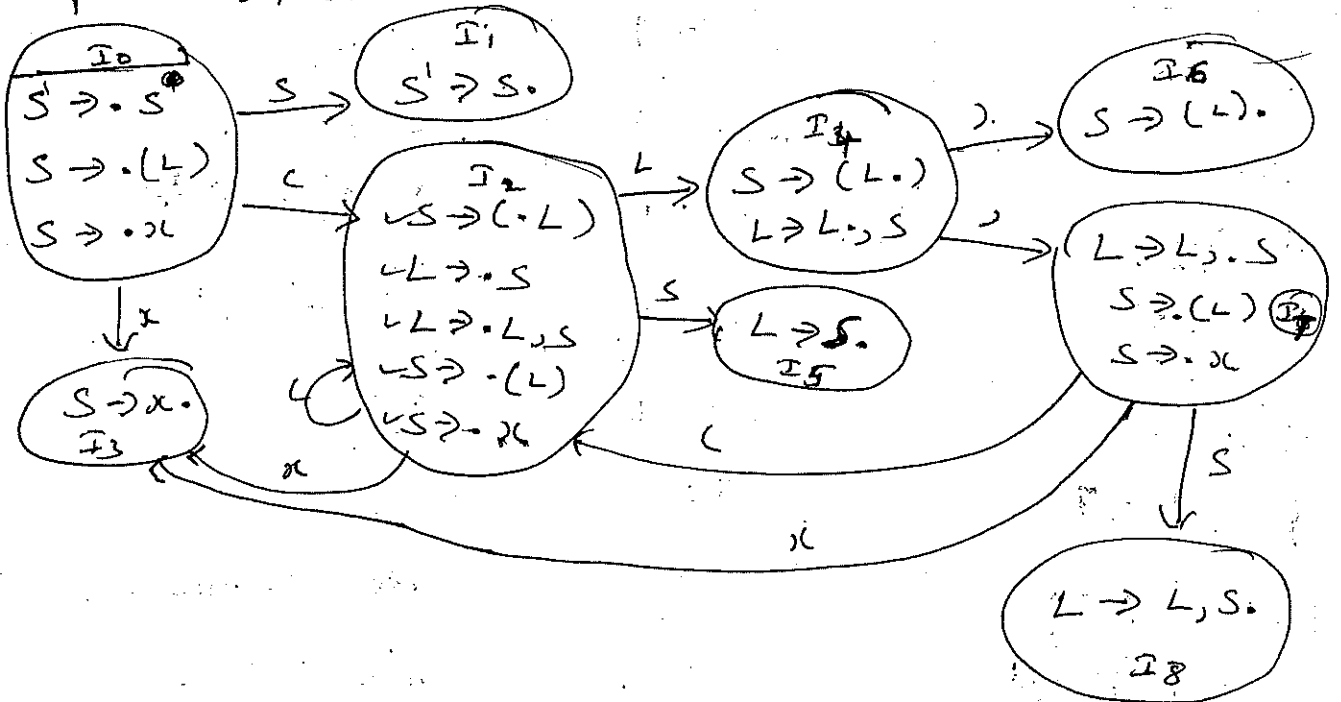
$$first(L) = \{ ( \alpha \}$$

	S	L
first	( $\alpha$	( $\alpha$
follow	\$ ) ,	) ,

(S2) find first & follow

NT	Productions	follow
S	$S \rightarrow (L)$ $A \rightarrow \alpha B B$	$S \rightarrow \$$ $follow(L) = first(\alpha)$
L	$L \rightarrow L, S$ $A \rightarrow \alpha B B$  $L \rightarrow L, S$ $A \rightarrow \alpha B$	$follow(L) = first(\alpha, S)$  $follow(S) = follow(L)$

(S3) find LR(0) items.



GURUPRASAD. S.  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64  
 Ph: 9886760776

S4 construct Parse table

	action				proto		
	(	)	,	x	\$	S	L
0	S <sub>2</sub>			S <sub>3</sub>		1	
1					accept		
2	S <sub>2</sub>			S <sub>3</sub>		5	4
3		r <sub>2</sub>	r <sub>2</sub>		r <sub>2</sub>		
4		S <sub>6</sub>	S <sub>7</sub>				
5		r <sub>3</sub>	r <sub>3</sub>				
6		r <sub>1</sub>	r <sub>1</sub>		r <sub>1</sub>		
7	S <sub>2</sub>			S <sub>3</sub>		8	
8		r <sub>4</sub>	r <sub>4</sub>				

S5 Parse string

(x, x)

2/p

action

Stacks

0

0(2

0(2x3

0(2S5

0(2L4

0(2L4,?

0(2L4,7x3

0(2L4,7x8

0(2L4

0(2L4)8

0S1

(x, x) \$

x, x) \$

, x) \$

, x) \$

, x) \$

x) \$

) \$

) \$

) \$

\$

\$

Shift

Shift

reduce S → x

reduce L → S

Shift

Shift

reduce S → x.

reduce L → L, S

Shift

reduce S → (L)

accept

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

$E \rightarrow E$   
 1)  $E \rightarrow E + T$   
 2)  $E \rightarrow T$   
 3)  $T \rightarrow T * F$   
 4)  $T \rightarrow F$   
 5)  $F \rightarrow (E)$   
 6)  $F \rightarrow id$

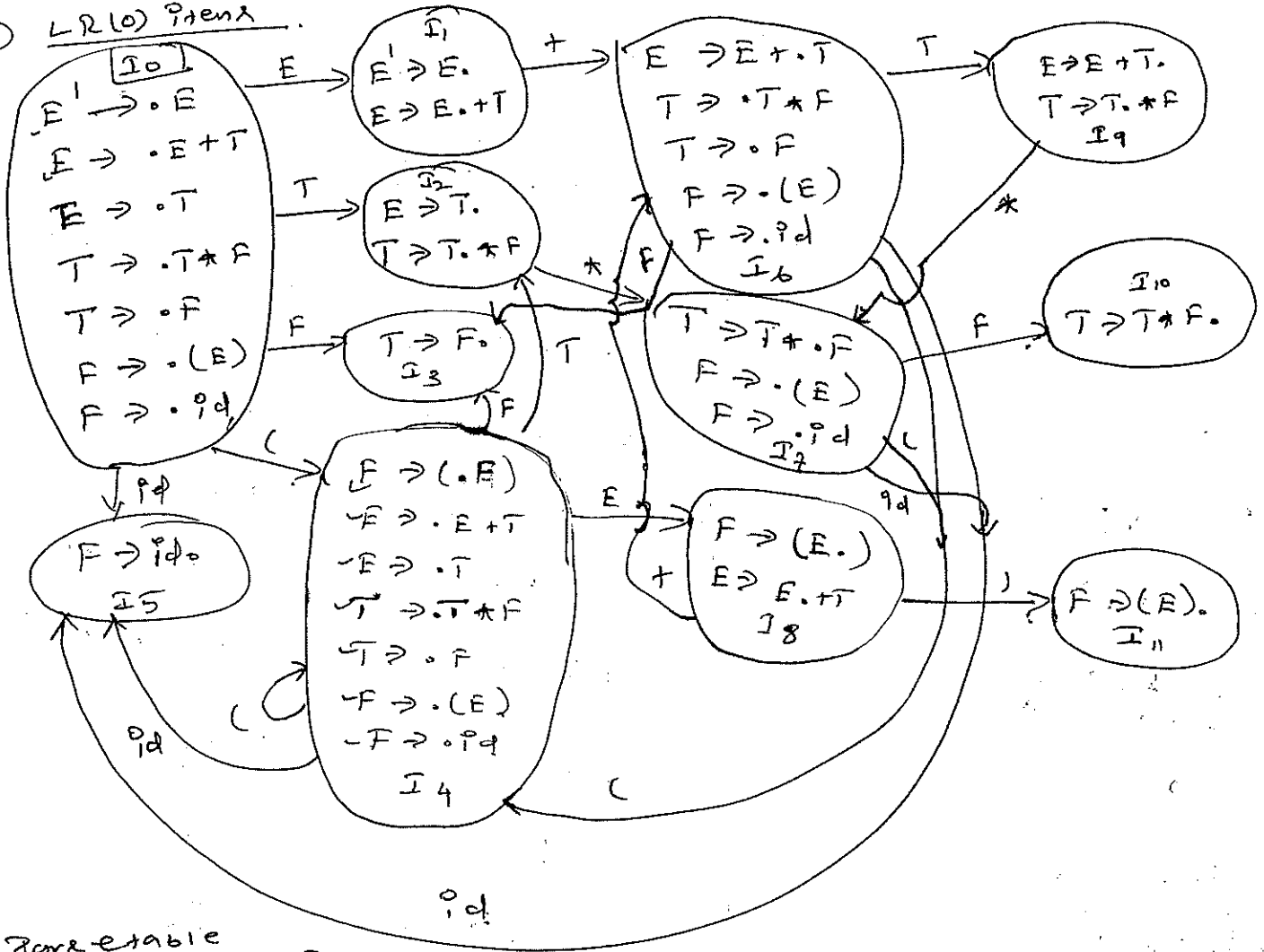
	E	T	F
first	{id}	{id}	{id}
follow	{+}	{*}	{*)}

$first(F) = \{id\}$   
 $first(T) = \{id\}$   
 $first(E) = \{id\}$

NT	Prodn	follow
E	$E \rightarrow E + T$ $A \rightarrow \alpha \overline{B} \beta$ $E \rightarrow E + T$ $A \rightarrow \alpha \overline{B} \beta$ $E \rightarrow T$ $A \rightarrow \alpha \overline{B}$	$follow(E) = first(+T) = +$ $follow(T) = follow(E)$ $follow(E) = follow(E)$
T	$T \rightarrow T * F$ $A \rightarrow \alpha \overline{B} \beta$ $T \rightarrow T * F$ $A \rightarrow \alpha \overline{B}$ $T \rightarrow F$ $A \rightarrow \alpha \overline{B}$	$follow(T) = first(*F) = *$ $follow(F) = follow(T)$ $follow(F) = follow(T)$
F	$F \rightarrow (E)$ $A \rightarrow \alpha \overline{B} \beta$	$follow(E) = first({}) = \epsilon$

S2

LR(0) items



S3 parse table

	←	accept		→	goto				
	+	*	(	id	⊘	⊘	E	T	F
0			S4	S5			1	2	3
1	S6				accept				
2	r2	S7			r2	r2			
3	r4	r4			r4	r4	8	2	3
4			S4	S5					
5	r6	r6			r6	r6		9	3
6			S4	S5				⊘	10
7			S4	S5					
8	S6					S11			
9	r1	S7			r1	r1			
10	r3	r3			r3	r3			
11	r5	r5			r5	r5			

Parse tree  $pd + pd + id$

Stacks	P   T	Action
0	$pd + pd + id \$$	Shift
0 $pd5$ x x	$+ pd + id \$$	reduce $F \rightarrow id$
0 $F3$ x x	$+ pd + id \$$	reduce $T \rightarrow F$
0 T 2	$+ pd + id \$$	Shift
0 T 2 * 7	$pd + id \$$	Shift
0 T 2 * 7 $pd5$ x x	$+ pd \$$	reduce $F \rightarrow id$
0 T 2 * 7 F 10 x x x x x x	$+ pd \$$	reduce $T \rightarrow T * F$
0 T 2 x x	$+ pd \$$	reduce $E \rightarrow T$
0 E 1	$+ pd \$$	Shift
0 E 1 + 6	$pd \$$	Shift
0 E 1 + 6 $pd5$ x x	$\$$	reduce $F \rightarrow id$
0 E 1 + 6 $F3$ x x	$\$$	reduce $T \rightarrow F$
0 E 1 + 6 T 7 x x x x x x	$\$$	reduce $E \rightarrow E + T$
0 E 1	$\$$	accept

GURUPRASAD. S.  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-6  
Ph: 9886760776

Prof GURUPRASAD S, BMSIT&M

<https://hemanthrajhemu.github.io>

Prob 3:  $S' \rightarrow \cdot S$

- 1)  $S \rightarrow eA$
- 2)  $S \rightarrow d$
- 3)  $A \rightarrow ad$
- 4)  $A \rightarrow a$

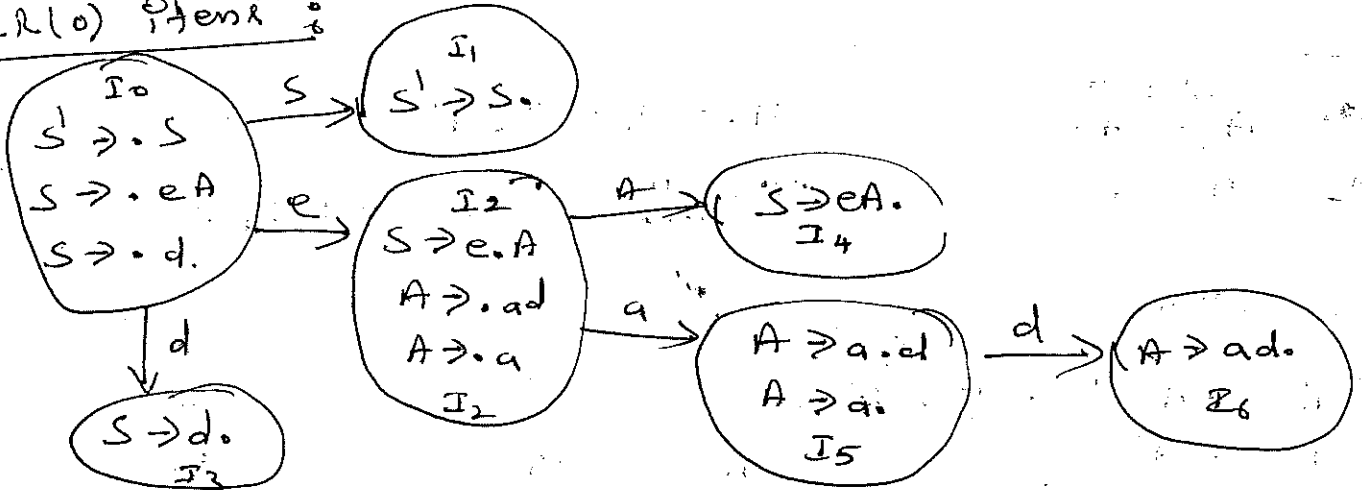
$$\text{FIRST}(S) = \{e, d\}$$

$$\text{FIRST}(A) = \{a\}$$

	S		A
first	e	d	a
follow	\$		\$

NT	Produ	follow
S	$S \rightarrow eA$ $A \rightarrow \frac{a}{x} B$	follow(A) = follow(S)

LR(0) items:



Parse table

	e	a	d	\$	S	A
0	S <sub>2</sub>		S <sub>3</sub>		1	
1				accept		
2		S <sub>5</sub>				4
3				r <sub>2</sub>		
4				r <sub>4</sub>		
5			S <sub>6</sub>	r <sub>4</sub>		
6				r <sub>3</sub>		





Parse steps	(a)	P/P	actions
		(a)\$	shift
0		(a)\$	shift
0(2		a)\$	shift
0(2(2		)\$	reduce $A \rightarrow a$
0(2(2a3 xx		)\$	shift
0(2(2A4		)\$	reduce $A \rightarrow (A)$
0(2(2A4)5 xx x xxx		)\$	shift
0(2A4		\$	reduce $A \rightarrow (A)$
0(2A4)5 xv x xxx		\$	accept
0A1			

Prob 5:  
 $S \rightarrow CC$   
 $C \rightarrow cC$   
 $C \rightarrow d$   
 $w = cdcd$

Prob 6:  
 $S \rightarrow L^*R$   
 $S \rightarrow R$   
 $L \rightarrow *R$   
 $L \rightarrow id$   
 $R \rightarrow L$

$w = id = id$

**GURUPRASAD. S.**  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64,  
 Ph: 9888760776

Prob 5:

$$S' \rightarrow S$$

$$1) S \rightarrow CC$$

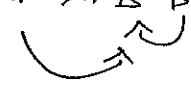


$$2) C \rightarrow cC$$

$$3) C \rightarrow d$$

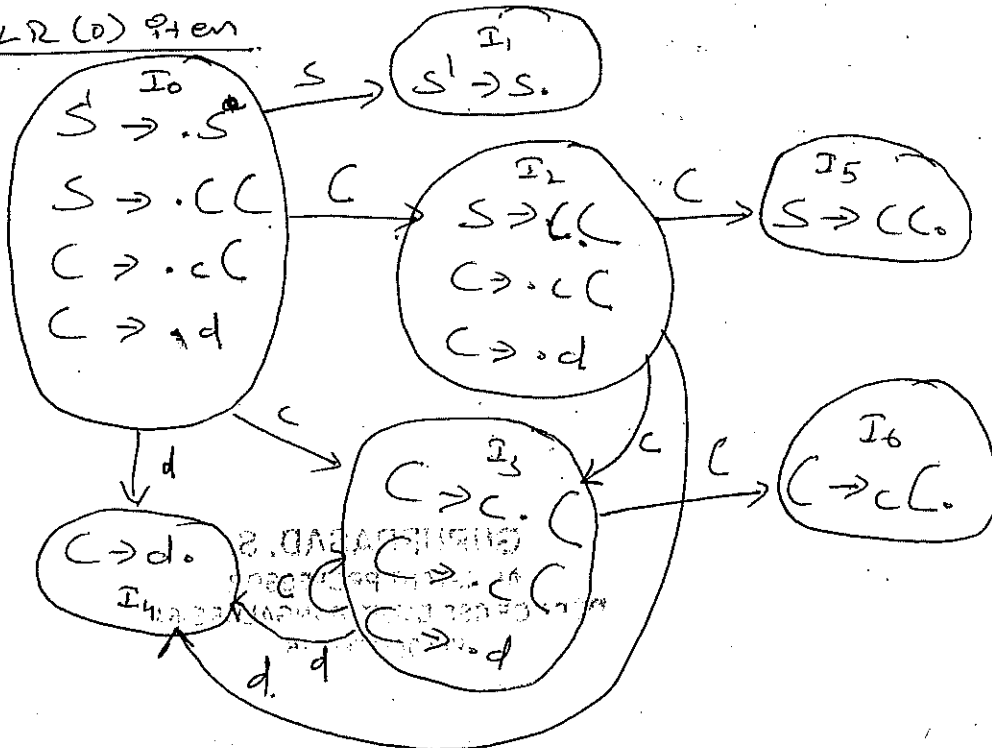
$$\text{direct}(S) = \text{direct}(C)$$

$$\text{direct}(C) = \{c, d\}$$

	S	C
direct	cd	cd
follow	\$	\$ cd

NT	Prodn	follow
S	$S \rightarrow CC$ $A \rightarrow \overline{B} \overline{B}$ 	$S \rightarrow \phi$ $\text{follow}(C) = \text{follow}(S)$ $\text{follow}(C) = \text{follow}(S)$
C	$S \rightarrow \overline{C} C$ $\overline{A} \overline{B}$ 	$\text{follow}(C) = \text{follow}(S)$
C	$C \rightarrow cC$ $A \rightarrow \overline{A} B$ 	$\text{follow}(C) = \text{follow}(C)$

LR(0) Item



	action			goto		
	c	d	\$	S	C	
0	S <sub>3</sub>	S <sub>4</sub>		1	2	
1			accept			
2	S <sub>3</sub>	S <sub>4</sub>			5	
3	S <sub>3</sub>	S <sub>4</sub>			6	
4	r <sub>3</sub>	r <sub>3</sub>	r <sub>3</sub>			
5			r <sub>1</sub>			
6	r <sub>2</sub>	r <sub>2</sub>	r <sub>2</sub>			

stack	i/p	action
0	cdcd\$	shift
0c3	dcd\$	shift
0c3d4	cd\$	reduce C → d
0c3d4 x^	cd\$	reduce C → cL
0c3c6	cd\$	shift
0c3c6 xxx^	d\$	<del>shift</del>
0c2	d\$	reduce C → d
0c2c3	\$	reduce C → cL
0c2c3d4	\$	reduce S → c'c'
0c2c3d4 xx	\$	
0c2c3c6	\$	reduce S → c'c'
0c2c3c6 *xxx^	\$	
0c2c5	\$	
0s1	\$	accept

GURUPRASAD. S.  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64.  
 Ph: 9886760776

Prof GURUPRASAD S, BMSIT&M

<https://hemanthrajhemu.github.io>

Prob 6/2

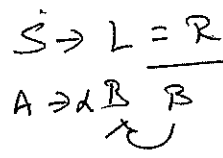
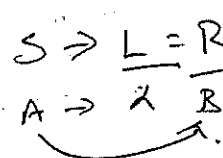
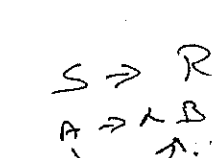
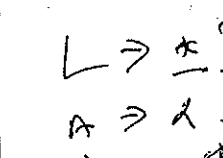
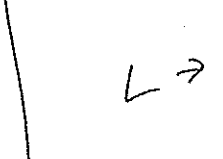
- $S \rightarrow S$   
 1)  $S \rightarrow L = R$   
 2)  $S \rightarrow R$   
 3)  $L \rightarrow \star R$   
 4)  $L \rightarrow Pd$   
 5)  $R \rightarrow L$

$$f_{LR+}(S) = f_{LR+}(L) = \{+, Pd\}$$

$$f_{LR+}(L) = \{+, Pd\}$$

$$f_{LR+}(R) = f_{LR+}(L) = \{+, Pd\}$$

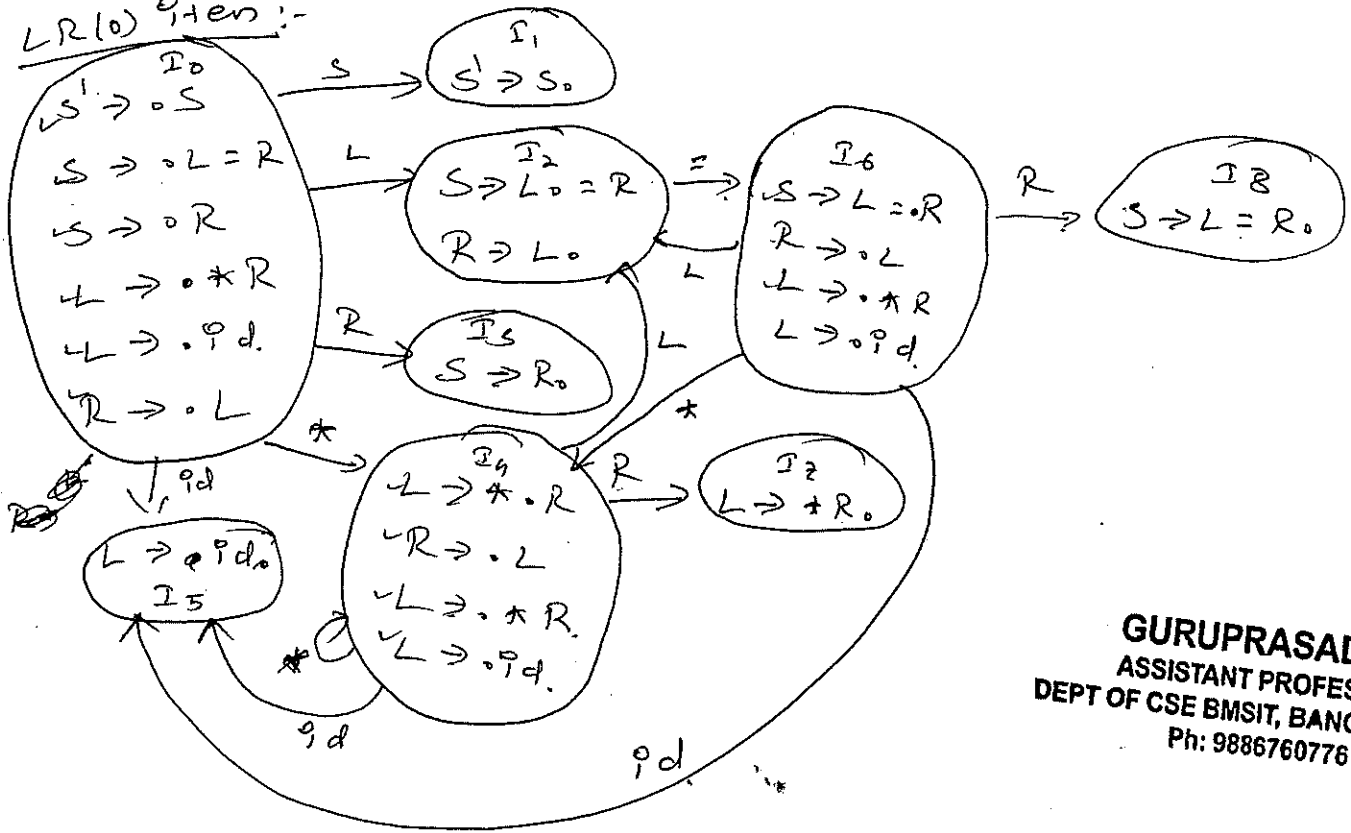
	S	L	R
$f_{LR+}$	$\{+, Pd\}$	$\{+, Pd\}$	$\{+, Pd\}$
$f_{011}$	\$	= \$	\$ =

NT	Prodn	$f_{011}$
S	$S \rightarrow L = R$ $A \rightarrow \star B$  $S \rightarrow \underline{L} = \underline{R}$ $A \rightarrow \star B$  $S \rightarrow R$ $A \rightarrow \star B$ 	$f_{011}(L) = f_{LR+}(L = R) = \{+, Pd\}$ $f_{011}(R) = f_{011}(S)$ $f_{011}(R) = f_{011}(S)$
L	$L \rightarrow \star R$ $A \rightarrow \star B$  $L \rightarrow Pd$	$f_{011}(R) = f_{011}(L)$ $f_{011}(L) = f_{011}(R)$
R	$R \rightarrow L$ $A \rightarrow \star B$ 	

Parse table

0

LR(0) items :-



**GURUPRASAD. S.**  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64.  
 Ph: 9886760776

	*	id	=	\$	S	L	R
0	S <sub>4</sub>	S <sub>5</sub>			1	2	3
1				accept			
2			<del>r<sub>5</sub> s<sub>6</sub></del>	r <sub>5</sub>			
3				r <sub>2</sub>			
4	S <sub>4</sub>	S <sub>5</sub>				2	7
5			r <sub>4</sub>	r <sub>4</sub>			
6	S <sub>4</sub>	S <sub>5</sub>				2	8
7			r <sub>3</sub>	r <sub>3</sub>			
8				r <sub>1</sub>			

Stacks	Q/P	action
0	Qd = Qd \$	Shift
Qids x x	= Qd \$	reduce L → Qd
0 L 2	= Qd \$	Shift
0 L 2 = 6	Qd \$	Shift
0 L 2 = 6 Qd 5 x x	\$	reduce L → Qd
0 L 2 = 6 L 2 x x	\$	reduce R → L
0 L 2 = 6 R 8 x x x x x x	\$	reduce S → L = R
0 S 1	\$	accept

End of ans 3

# Syntax Directed Translation

## Semantic Analysis

Why:

- Syntactic correctness was checked during parsing
- Another level of correctness is not captured i.e.
  - \* Var is declared/not
  - \* Types are consistent
  - \* if  $x = y$  is  $y$  assignable to  $x$
  - \* fn calls have right no. & type of parameters
  - \* if  $P \rightarrow Q$  is  $Q$  is member of OBP
  - \* is var  $x$  is initialized before use etc

What:

→ Semantic actions include

(i) Symbol table handling

- \* maintain info abt declared symbols.
- \* info abt types
- \* info abt scopes.

(ii) Checking context cond<sup>ns</sup>

- \* scope rule
- \* type checking

(iii) Invocation of code generation routines.

It is done during reduction during parsing.

How:

It is done by using a technique called Syntax Directed Definition (SDD)

The semantic analysis involves both the description of analysis & the implementation algorithm

In parsing CF is used as description & BOPD generated as algorithms

bht

**Prof GURUPRASAD S, BMSIT&M**

<https://hemanthrajhemu.github.io>

**GURUPRASAD. S.**  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-54  
Ph: 9886760776

There are no standard method to specify semantics of lang as it varies from lang to lang

One method used is to identify attributes of  $\forall$  lang entities that must be computed and to write attribute equations or semantic rules

This method is called Syntax Directed Definition

In SDD attributes are directly associated with grammar symbols

If  $X$  is a grammar symbol and  $a$  is an attribute then  $X.a$  represents the value of  $a$  associated to  $X$

In SDD grammar symbols (T & NT) are attached with attributes based on info of  $\forall$  lang constructs

Values of these attributes are computed by semantic rules associated with grammar productions.

Each node in parse tree now acts as record holding attribute values

Attributes in SDD could be of any kind like.

- numbers
- table reference
- types
- name etc

Each occurrence of grammar symbol will have separate instance of attribute.

### Types of Attributes

There are two types of attribute.

- ① Synthesized
- ② Inherited



## ① Synthesized Attribute

\* The value of these attributes is computed by the value of attributes of its children in parse tree.

eg:-  $E \rightarrow E + num | num$

$$E_2.val = E_1.val + num.val$$

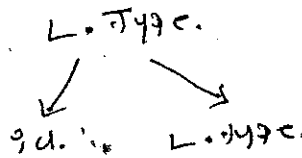
\* It is used with LR Parser.

\* The SDD is S-attribute of every attr is synthesized.

## ② Inherited Attribute

\* The value of attribute is computed from the attribute value of its parent / siblings in parse tree.

eg:-  $D.type$   
 $\swarrow \searrow$   
 $T.type \rightarrow L.type$



$D \rightarrow TL^*$

$T \rightarrow int | id$

$L \rightarrow L, id | id$

\* The SDD is L-attribute of val of attr is synthesized / inherited

## Semantic Rules in SDD

They are the rules define the value of attribute of a grammar symbol.

Production.	Semantic Rules
$L \rightarrow E_n$	$L.val = E_n.val$
$E \rightarrow E_1 + T$	$E.val = E_1.val + T.val$
$E \rightarrow T$	$E.val = T.val$
$T \rightarrow T_1 * F$	$T.val = T_1.val * F.val$
$T \rightarrow F$	$T.val = F.val$
$F \rightarrow (E)$	$F.val = E.val$
$F \rightarrow digit$	$F.val = digit.lexVal$

# Evaluating an SDD at Nodes of a Parse Tree

The parse tree showing values of its attributes is called an "Annotated parse tree".

To Construct Annotated Parse Tree:-

⊕ Before we evaluate an attribute at a node of parse tree we must evaluate all the attribute upon which its value depends.

So with synthesized we can evaluate all children before parent in bottom up manner.

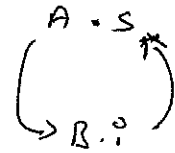
but with inherited & synthesized attribute there will be no order

eg:  $A \rightarrow B$

Semantic rule

$$A.s = B.i$$

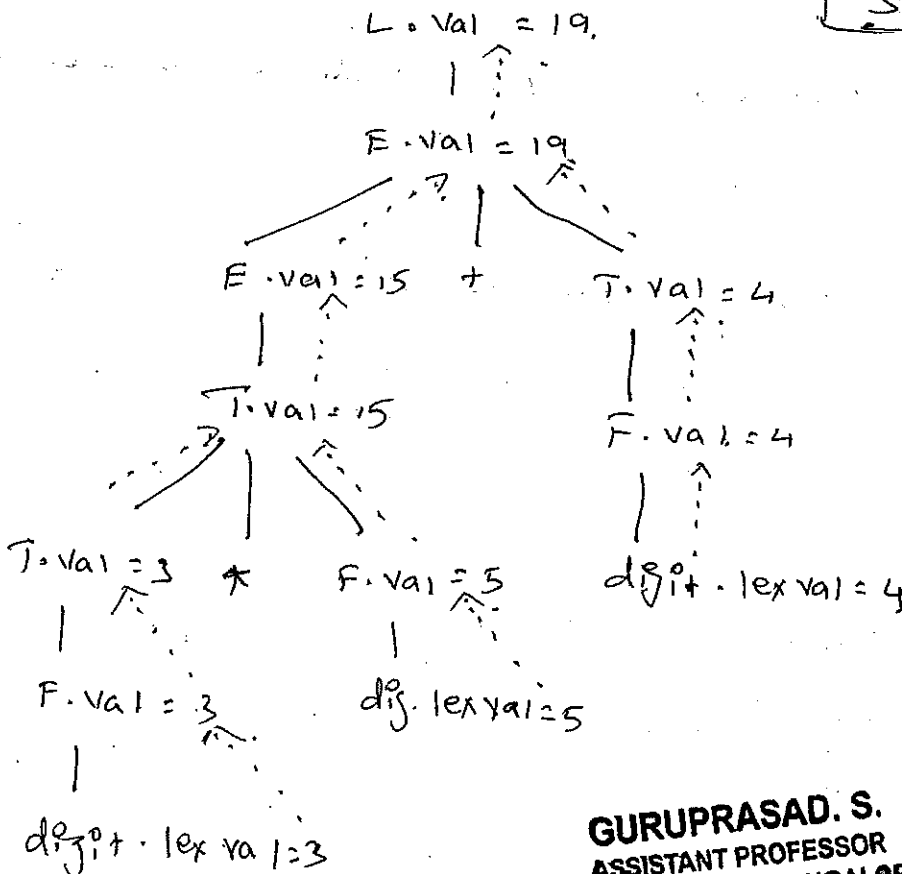
$$B.i = A.s + 1$$



this is not possible to evaluate due to circularity.

Annotated parse tree. eg continued...

$3 * 5 + 4 = 19$



**GURUPRASAD. S.**  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64,  
 Ph: 9886760776

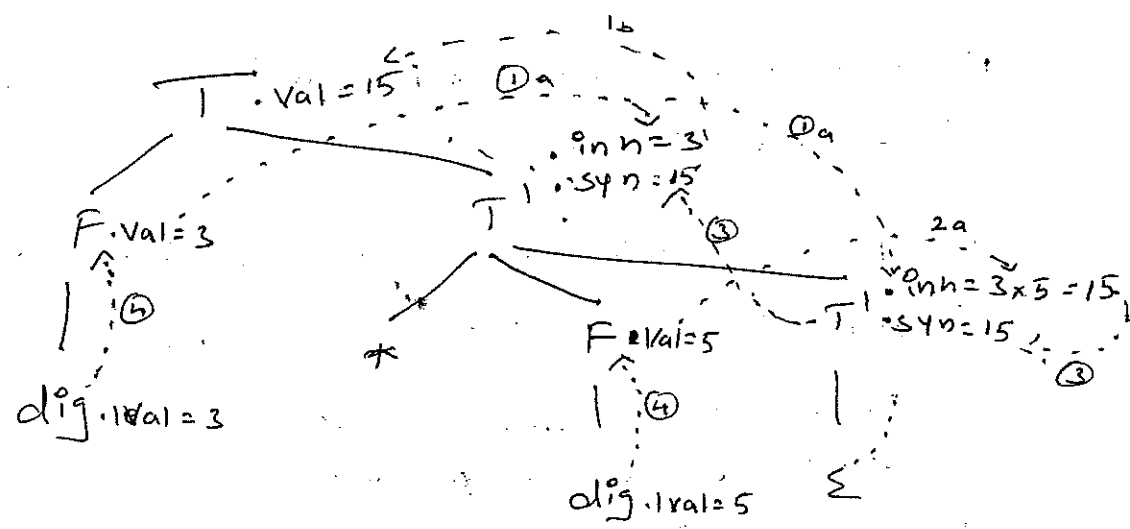
**Prof GURUPRASAD S, BMSIT&M**

<https://hemanthrajhemu.github.io>

Ex 2

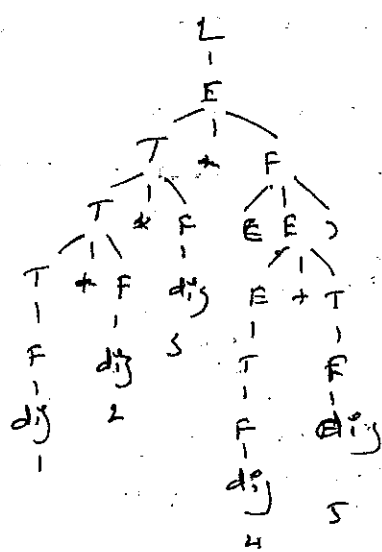
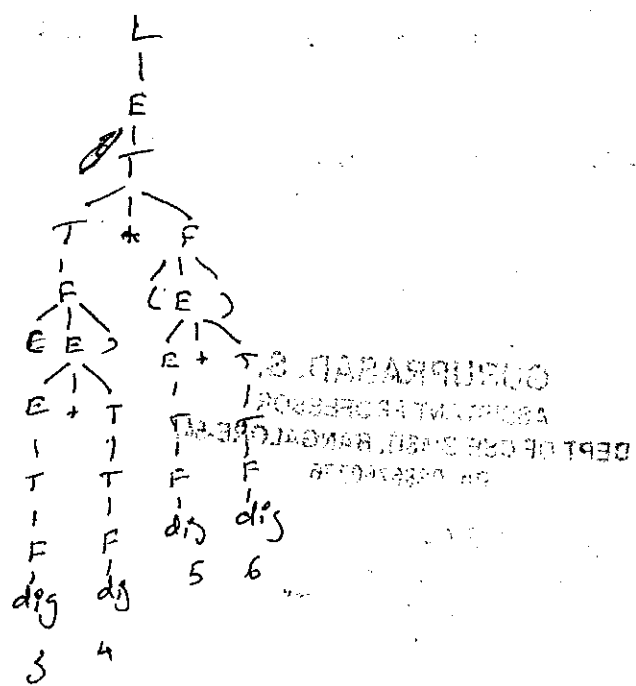
Production	Semantic Rules
$T \rightarrow FT'$	$T'.inh = F.val$ $T.val = T'.syn$
$T' \rightarrow * FT'$	$T'.inh = T.inh * F.val$ $T'.syn = T'.syn$
$T' \rightarrow \epsilon$	$T'.syn = T'.inh$
$F \rightarrow dig$	$F.val = dig.lexval$

3 \* 5



① (3+4) \* (5+6)

② 1 \* 2 + 3 \* (4+5)



## Evaluation orders of SDDs.

### Dependency Graph:

Dependency Graphs are useful tool for <sup>determining</sup> an evaluation order for attrib instances in given parse tree.

show how values are calculated in annotated parse tree.

It depicts the flow of info among the attrib instance in a particular parse tree.

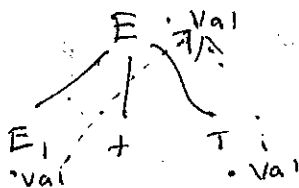
An edge from one attrib to another means that the value of first need to be computed before the second.

\* for each parse node with symbol  $X$  the DG has node associated with  $X$

\* If the semantic rule of Prod<sup>n</sup>  $P$  define value of synthesized attrib  $A.b$  in terms of val of  $X.c$  then DG will have edge from  $X.c$  to  $A.b$  (where  $X$  is always child)

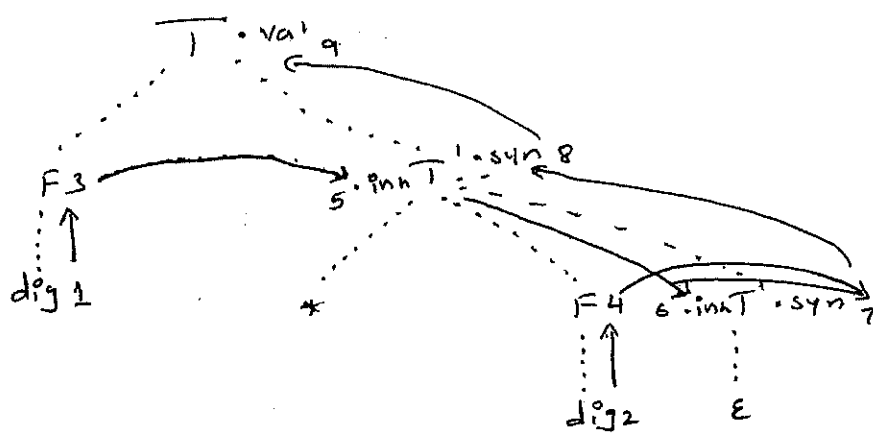
\* If semantic rule of Prod<sup>n</sup>  $P$  define value of inherited attrib  $B.c$  in terms of val of  $X.a$  then DG has an edge from  $X.a$  to  $B.c$  (where  $X$  is parent or sibling)

Eg: Synthesized.



**GURUPRASAD. S.**  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-64.  
Ph: 9886760776

Ex: Inheriteds  
order root + down  
1-9



DG specify the order in which we can evaluate the attrib at various nodes of parse tree.

If DG has edge from node  $M$  to  $N$  then  $M$  has to be evaluated before  $N$ , each ordering makes a directed graph to linear order and is called as Topological Sort

If there is a cycle in the graph then there are no topological sorts, then there is no way to evaluate the SDD.

### S- attributed Definition

An SDD is S- attributed if every attrib is Synthesized.

When an attrib is S- attributed, we can evaluate it's attrib in any BU order i.e we apply Post order traversal

It is imple with LR parser. with parsing stacks.

### L- attributed Definition

SDD is L- attributed if each attrib must be either

(i) Synthesized

(ii) Inherited but with rule that

if from  $A \rightarrow X_1 X_2 \dots X_n$  then  $X_i$  can be computed only by

(a) inherited attrib associated with  $A$

(b) either inherited or synthesized attrb associated with  $X_1 \dots X_{p-1}$  located to left of  $X_p$

(c) either inherited or synthesized attrb associated with  $X_p$  itself but no way cycles in DG.

## Semantic Rules with Controlled Side Effects

In practice translation involve side effects.

eg: calculator print result

enter st for type of id etc.

We shall control the side effects in SDD to allow constant evaln with DG. with

\* Permit incidental side effects that do not constrain attribute evaluation.

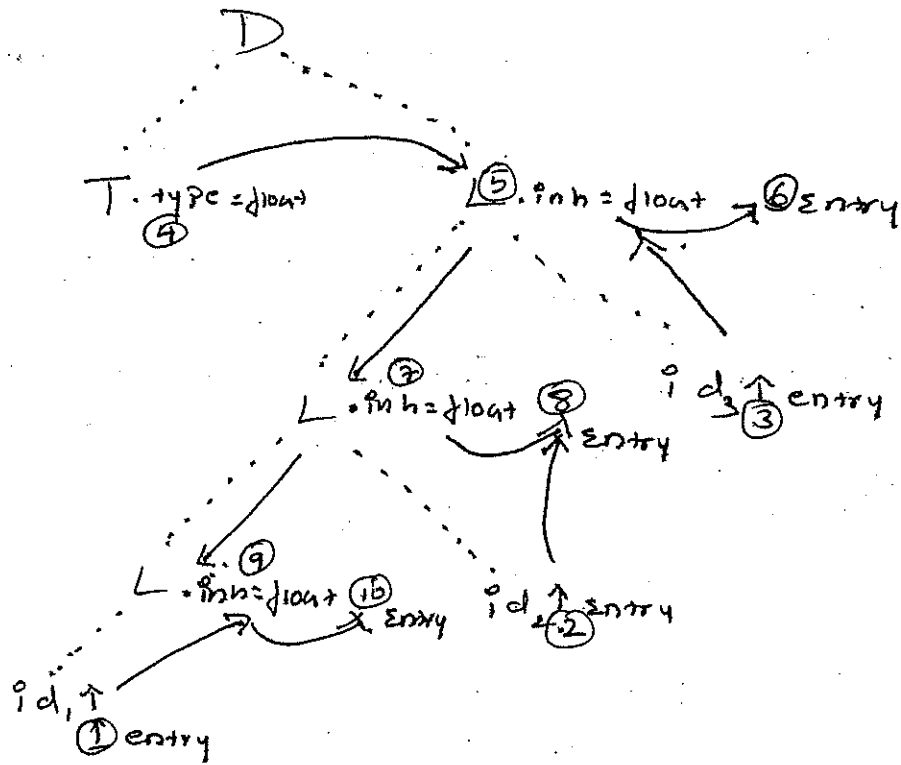
\* Constrain the allowable evaluation orders, so that the same translation be produced for any allowable order.

Eg:-

<u>Production.</u>	<u>Semantic Rules</u>
$D \rightarrow TL$	$L.inh = T.type$
$T \rightarrow int$	$T.type = integer$
$T \rightarrow float$	$T.type = float$
$L \rightarrow L_1, id$	$L_1.inh = L.inh$ add type (id.entry, L.inh)
$L \rightarrow id$	add type (id.entry, L.inh)

D - Decl<sup>n</sup> T - type L - list of int

w: float i<sub>d1</sub>, i<sub>d2</sub>, i<sub>d3</sub>



eg in + a b c. float w x y z

## Application of Syntax Directed Translation

Compiler use Syntax tree as intermediate rep<sup>n</sup> i.e. SDD converts  
 if string to syntax tree - the compiler walks through these tree using  
 another set of rules to translate it into intermediate code.

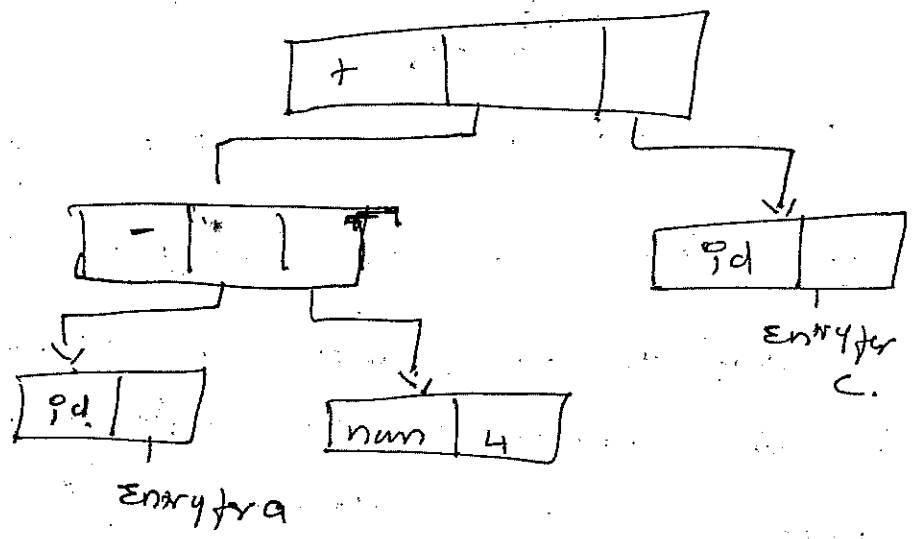
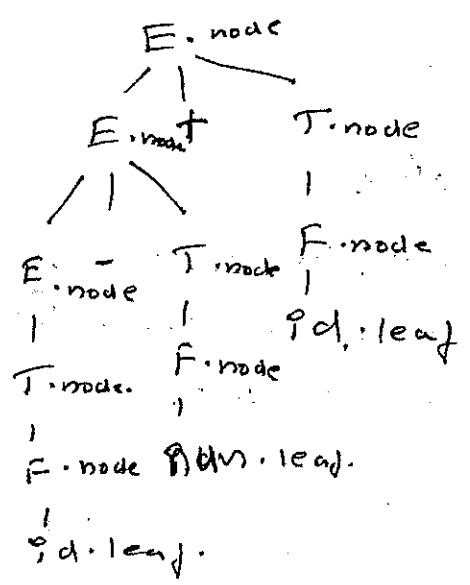
### Construction of Syntax tree.

- \* ~~Each object will have an 'op' field. that is the label of nodes~~
- \* Each node in syntax tree rep<sup>n</sup> one object, it will be record/construct holding info abt obj & ptr to child nodes.
- \* If node is leaf it holds lexical value for the leaf.
- \* If node is an interior node. it holds ptr to its children

Node (op, c<sub>1</sub>, c<sub>2</sub> ... c<sub>n</sub>)

Ex:- Rodn	Semantic Rules
$E \rightarrow E_1 + T$	$E.node = \text{new Node}('+', E_1.node, T.node)$
$E \rightarrow E_1 - T$	$E.node = \text{new Node}('-', E_1.node, T.node)$
$E \rightarrow T$	$E.node = T.node$
$T \rightarrow (E)$	$T.node = E.node$
$T \rightarrow id$	$T.node = \text{new Leaf}(id, id.entry)$
$T \rightarrow num$	$T.node = \text{new Leaf}(num, num.val)$

a - 4 + I



- $P_1 = \text{new leaf}(id, \text{entry a})$
- $P_2 = \text{new leaf}(num, 4)$
- $P_3 = \text{new node}('-', P_1, P_2)$
- $P_4 = \text{new leaf}(id, \text{entry c})$
- $P_5 = \text{new Node}('+', P_3, P_4)$

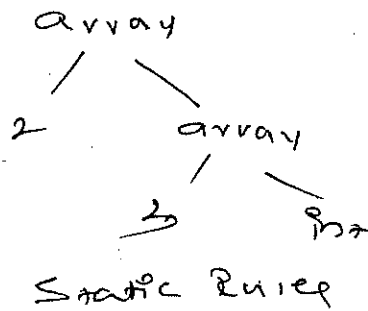




ip2: Structure of type.

In C  $\text{int}[2][3]$  array of 2 array of 3 integers.

Reqn are



Prodm

$T \rightarrow BC$

$T \cdot t = C \cdot t$

$C \cdot b = B \cdot t$

$B \rightarrow \text{int}$

$B \cdot t = \text{int}$

$B \rightarrow \text{float}$

$B \cdot t = \text{float}$

$C \rightarrow [\text{num}] C_1$

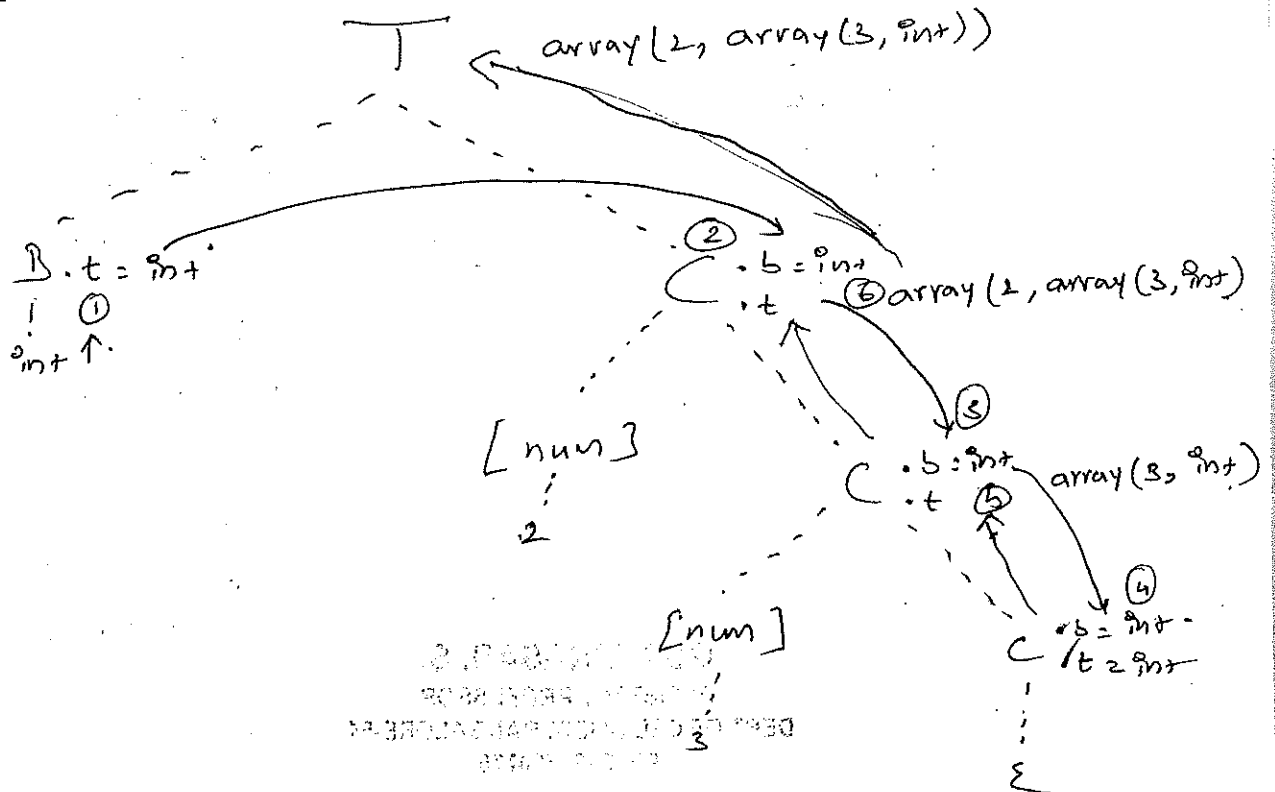
$C \cdot t = \text{array}(\text{num.val}, C_1 \cdot t)$

$C_1 \cdot b = C \cdot b$

$C \rightarrow \epsilon$

$C \cdot t = C \cdot b$

$\text{int}[2][3]$



# Syntax Directed Translation Scheme

SDT schemes are complementary notation to SDD, it is the implementation of SDD

SDT contains CFG embedded with semantic actions appear at any position in production body

Any SDT can be imp<sup>d</sup> by first building a parse tree & then performing the actions from left to right in depth first order.

SDTs are implemented on two important classes of SDDs

- 1) underlying grammar is LR Parsable & SDD is S-attributed.
- 2) underlying grammar is LL Parsable & SDD is L-attributed.

The objective is to convert semantic rules in SDD into SDT with actions executed at right time during parse.

SDT imp<sup>d</sup> during parsing use a marker nonterminal  $M$  where  $M \rightarrow \epsilon$

## Fortran translation Scheme

SDT for S-attributed SDD will contain actions at the end of prod<sup>n</sup> and are executed along with reduction of the body to head of prod<sup>n</sup> they are called Fortran SDT's

Ex:

- 1)  $L \rightarrow E_n$  { Print (E.val) }
- 2)  $E \rightarrow E_1 + T$  { E.val = E<sub>1</sub>.val + T.val }
- 3)  $E \rightarrow T$  { E.val = T.val }
- 4)  $T \rightarrow T_1 * F$  { T.val = T<sub>1</sub>.val \* F.val }
- 5)  $T \rightarrow F$  { T.val = F.val }
- 6)  $P \rightarrow (E)$  { P.val = E.val }
- 7)  $P \rightarrow \text{dig}$  { P.val = dig.lexval }

**GURUPRASAD. S.**  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-64.  
Ph: 9886760776

**Prof GURUPRASAD S, BMSIT&M**



## SDT's with action before Production

An action may be placed at any place within the body of productions & is performed immediately after all symbols to its left are processed.

thru if we have prodn  $B \rightarrow X \{a\} Y$  the action  $a$  is done after we have recognized  $X$  (if  $X$  is terminal) or all terminals derived by  $X$  (if  $X$  is non terminal)

In BU parse perform action 'a' as soon as the occurrence of  $X$  on parse stack.

In TD parse perform action 'a' before we attempt to expand the occurrence of  $Y$ .

Not all SDT's can be imple during parsing

g:  $L \rightarrow E_n$

$E \rightarrow \{ \text{Print}(' + ') \} E_1 + T$

$E \rightarrow T$

$T \rightarrow \{ \text{Print}(' * ') \} T * F$

$T \rightarrow F$

$F \rightarrow (E)$

$F \rightarrow \text{digit} \{ \text{Print}(\text{dig. lex. val}) \}$

g) is imple to imple thru SDT either TD/BU parse as parser has to print instance of  $* +$  before knowing whether it appears or not

So we rather non terminal  $M_2, M_4$  for 2 2 4 productions

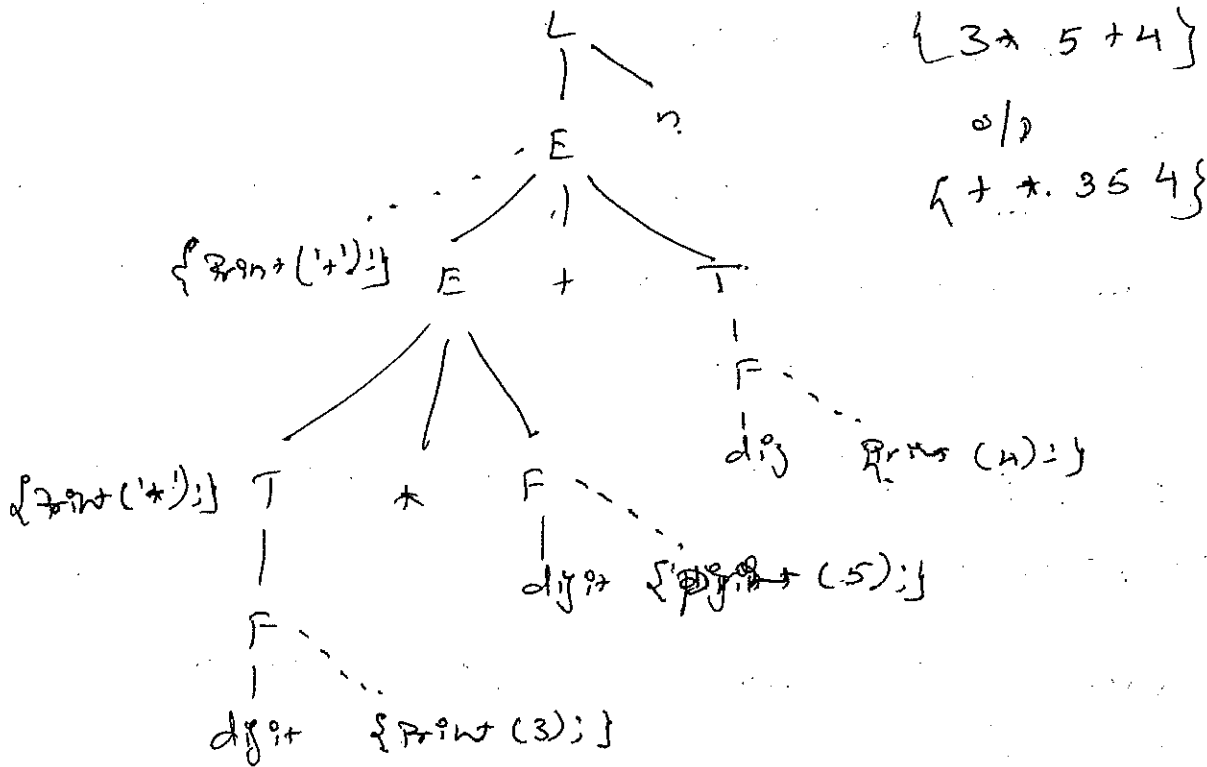
where  $M_L \rightarrow \epsilon$  &  $M_4 \rightarrow \epsilon$  & shift digit

**Prof GURUPRASAD S, BMSIT&M**

<https://hemanthrajhemu.github.io>

**GURUPRASAD. S.**  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-64.  
Ph: 9886760776

then performs the order traversal of the tree & at each node labeled by an action the visited perform action



### Eliminating Left Recursion from SDT

Since we can't parse a grammar with left recursion in TD  
 it had to be eliminated.

In case of SDT also we should eliminate left recursion to know  
 the order in which the action in an SDT are performed.

so

$$A \Rightarrow A\alpha | B$$

$$A \Rightarrow \beta R$$

$$R \Rightarrow \alpha A | \epsilon$$

eg:-

$$E \Rightarrow E_1 + T \{ \text{Print}('+') \}$$

$$E \Rightarrow T$$

$$E \Rightarrow TR$$

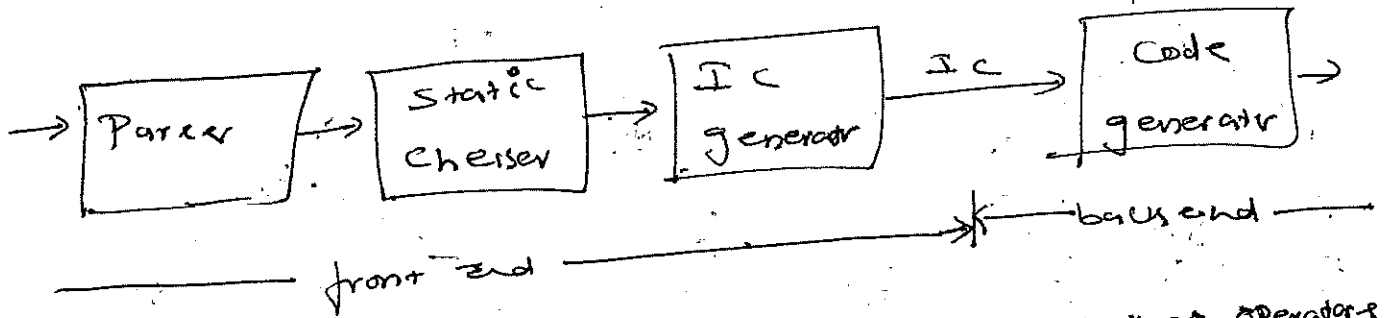
$$R \Rightarrow + T \{ \text{Print}('+') \} R | \epsilon$$

End of unit 5

# Intermediate Code Generation

In analysis & synthesis model of compiler front end analyze the src prog & create intermediate repr for which back end generates target code.

So details of src prog is confined to front end while target m/c details are confined to back end.



Static checking includes type checking to ensure that operators are applied to compatible operands. also include additional syntactic checks.

Before src prog is converted to target prog a sequence of intermediate repr are constructed. Some high level close to program & low level close to target m/c repr like syntax tree, 3-addr code, etc.

## Variants of Syntax Tree

The nodes in syntax tree repr constructs of src prog; the children of a node repr the meaning of all const. of construct

A Directed Acyclic Graph (DAG) is an expr; identifies common sub expressions (sub expr that occur more than once) of the expr.

DAG can be constructed with same technique used

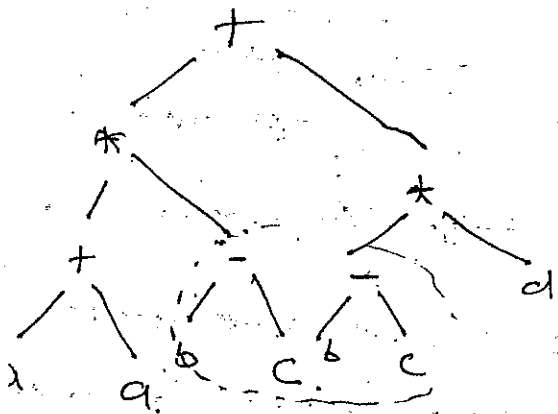
Construct syntax tree / Abstract syntax tree

# DAE for Expression

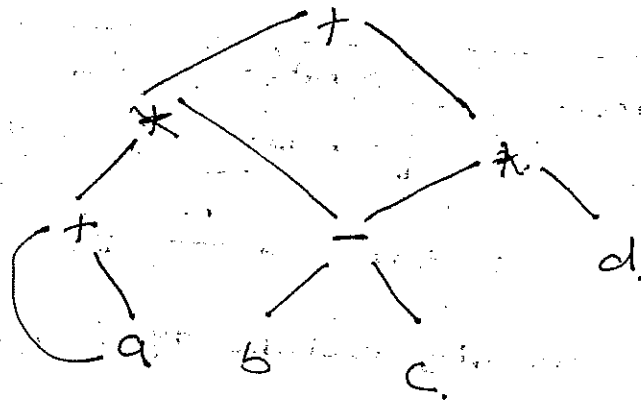
It rep<sup>n</sup> expr<sup>n</sup> & give compiler important info regarding the generation of efficient code to evaluate the expr<sup>n</sup>

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

AST:

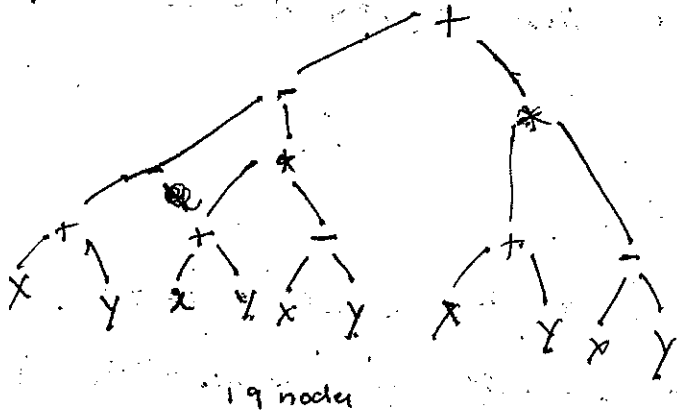


DAGs

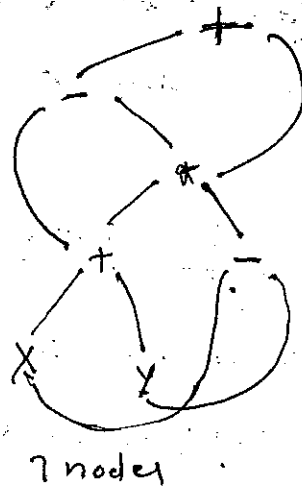


$$((x+y) - ((x+y) * (x-y))) + ((x+y) * (x-y))$$

AST

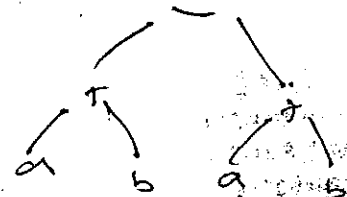


DAE:



$$(a+b) + (a+b)$$

AST



DAE





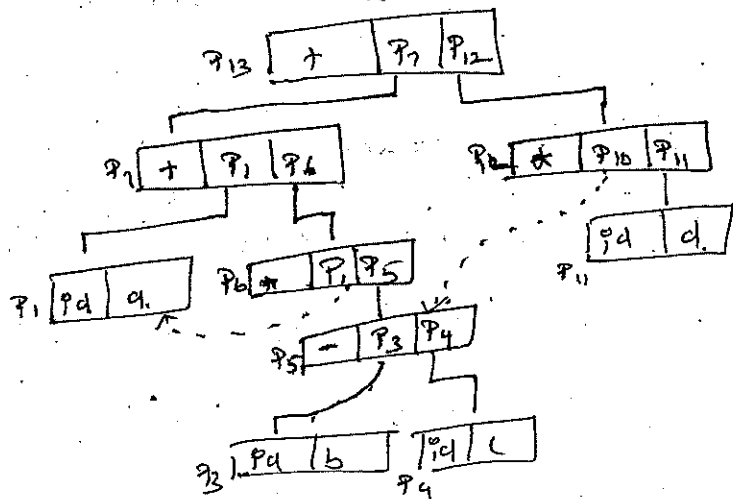
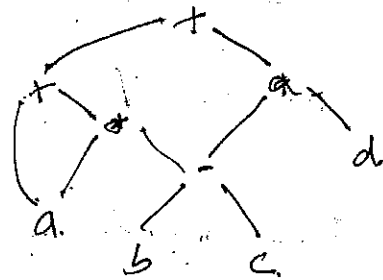
# Construction of Syntax tree for simple Expr<sup>n</sup>

The SDD given below can construct either Syntax tree/DAG

Production	Semantic Rules
$E \rightarrow E_1 + T$	$E.node = \text{newNode}('+', E_1.node, T.node)$
$E \rightarrow E_1 - T$	$E.node = \text{newNode}('-', E_1.node, T.node)$
$E \rightarrow T$	$E.node = T.node$
$T \rightarrow (E)$	$T.node = E.node$
$T \rightarrow id$	$T.node = \text{newLeaf}(id, id.entry)$
$T \rightarrow num$	$T.node = \text{newLeaf}(num, num.val)$

Construct Syntax tree for  $a + a * (b - c) + (b - c) * d$ .

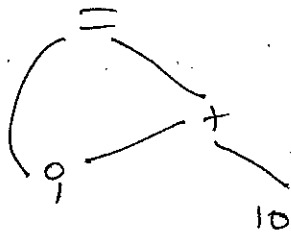
- $P_1 = \text{newLeaf}(id, \text{entry } a)$
- $P_2 = \text{newLeaf}(id, \text{entry } a) = P_1$
- $P_3 = \text{newLeaf}(id, \text{entry } b)$
- $P_4 = \text{newLeaf}(id, \text{entry } c)$
- $P_5 = \text{newNode}('-', P_3, P_4)$
- $P_6 = \text{newNode}('*', P_2, P_5)$
- $P_7 = \text{newNode}('*', P_1, P_6)$
- $P_8 = \text{newLeaf}(id, \text{entry } b) = P_3$
- $P_9 = \text{newLeaf}(id, \text{entry } c) = P_4$
- $P_{10} = \text{newNode}('-', P_8, P_9) = P_5$
- $P_{11} = \text{newLeaf}(id, \text{entry } d)$
- $P_{12} = \text{newNode}('*', P_{10}, P_{11})$
- $P_{13} = \text{newNode}('+', P_7, P_{12})$



**GURUPRASAD. S.**  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-81,  
 Ph: 9886760776

~~as above~~  
The value number method to construct DAG.

Often the nodes of a syntax tree or DAG are stored in array of records as shown below.



val no	op	l	r
1	op		
2	num	10	
3	+	1	2
4	=	1	3

array for i

Nodes of DAG for  $i = i + 10$  allocated in an array

Algorithm for value number method

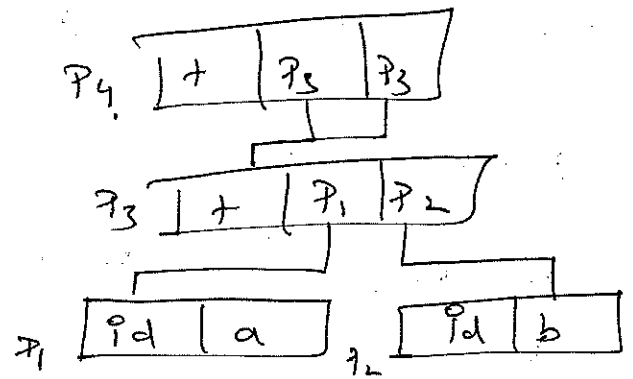
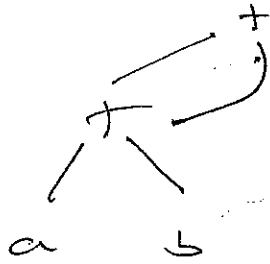
Input: label op, node l and node r

Output: The value number of a node in the array with signature  $\langle op, l, r \rangle$

Method: Search the array for a node M with label op, left child l, and right child r. If there is such a node return the value number of M. If not create in the array a new node N with label op, left child l and right child r and return the value number.

A more efficient approach is to use a hash table in which the nodes are put into "bucket".

eg:-  $a + b + (a + b)$



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

$P_2 = \text{newleaf}(\text{id}, \text{entry } b)$

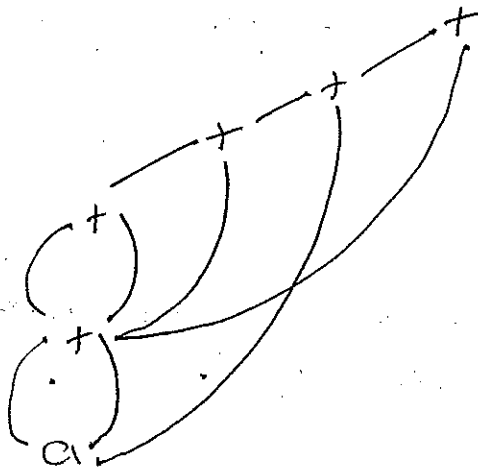
$P_3 = \text{newnode}(+, P_1, P_2)$

$P_4 = \text{newnode}(+, P_3, P_3)$

Value map

1	id		→ entry a
2	id		→ entry b
3	+	1   2	
4	+	3   3	

eg:-  $a + a + ((a + a + a + (a + a + a + a)))$



1	id		→ entry a
2	+	1   1	
3	+	2   2	
4	+	3   2	
5	+	4   1	
6	+	5   2	

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

$P_2 = \text{newnode}(+, P_1, P_1) \text{ } a + a$

$P_3 = \text{newnode}(+, P_2, P_2) \text{ } a + a + a$

$P_4 = \text{newnode}(+, P_3, P_2)$

$P_5 = \text{newnode}(+, P_4, P_1)$

$P_6 = \text{newnode}(+, P_5, P_2)$

**GURUPRASAD. S.**  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-8  
 Ph: 9886760776

**Prof GURUPRASAD S, BMSIT&M**

<https://hemanthrajhemu.github.io>

# Three address code

It is of the form  $A = B \text{ op } C$ .

$x + y + z$  might be translated into sequence of 3 address instructions

$$t_1 = y + z \quad t_1, t_2 \text{ are temp names}$$

$$t_2 = x + t_1$$

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

3-address code

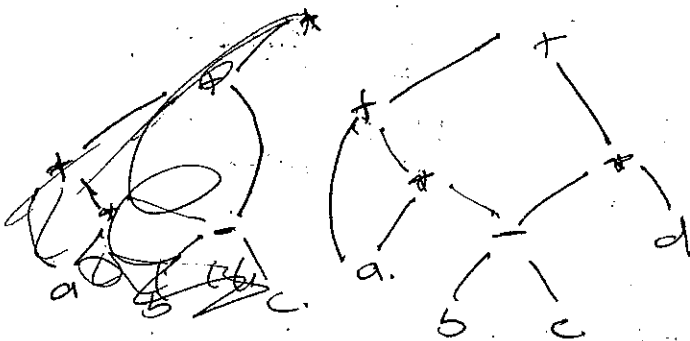
$$t_1 = b - c$$

$$t_2 = t_1 * a$$

$$t_3 = t_2 + a$$

$$t_4 = t_1 * d$$

$$t_5 = t_3 + t_4$$



Three address code is built from two concepts  
(i) address (ii) instructions.

Address can be one of the following

- 1) a name - Src Prog name appears as 3 address code, in imp/d name is replaced by a ptr to ST entry
- 2) a constant -
- 3) compiler-generated temporary names

## Instructions of 3-address code

1) Assignment - of the form  $x = y \text{ op } z$

or of the form  $x = \text{op } y$   
binary                      unary

2) Copy instr - of the form  $x = y$

3) unconditional Jump - of the form if x goto L

4) Conditional Jump - of the form if x.relop.y goto L

5) Procedure call & return

6) Indexed Copy instr<sup>n</sup> -  $x = y[i]$ ,  $x[i] = y$

7) address & pointer assignment of the form  $x = &y$   $x = *y$

Eg:-  
do  $i = i + 1$   
while ( $a[i] < v$ );

L:  $t_1 = i + 1$   
 $i = t_1$   
 $t_2 = i * 4$   
 $t_3 = a[t_2]$   
if  $t_3 < v$  goto L

100:  $t_1 = i + 1$   
101:  $i = t_1$   
102:  $t_2 = i * 4$   
103:  $t_3 = a[t_2]$   
104: if  $t_3 < v$  goto 100

a) Symbolic labels

b) Position numbers

## Quadruples

The intermediate rep<sup>n</sup> of 3-addr instrs in Data Structures can be done in 3 ways (i) Quadruples (ii) Triples (iii) indirect triples

Quadruple has 4 fields OP, arg1, arg2, result

Operator - contains internal code for operators.

$A = -B * (C + D)$

$t_1 = -B$

$t_2 = C + D$

$t_3 = t_1 * t_2$

$A = t_3$

OP	Arg1	Arg2	Result
uminus	B	-	$t_1$
+	C	D	$t_2$
*	$t_1$	$t_2$	$t_3$
=	$t_3$	-	A

Prof GURUPRASAD S, BMSIT&M

<https://hemanthrajhemu.github.io>

GURUPRASAD, S.  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-64,  
Ph: 9886760776

Notes unary operators & alignment op do not use arg2

Param operators do not use either arg1 or result

uncondl & condn jmps put the target label in result

The contents of arg1, arg2 & result are normally point to ST

Entries for the names, so temporary names must be entered into ST

when they are created.

Triples

One 3-address stmt are repr by a structure with only 3 fields

OP, arg1, arg2, where arg1 & arg2 are arguments of OP.

they are either ptr to ST or ptr to structure.

We use parenthesized name to repr ptrs into triple structure

while ST ptr are repr by the name then self.

1)  $A = -B * (C + D)$

$t_1 = -B$

$t_2 = C + D$

$t_3 = t_1 * t_2$

$A = t_3$

OP	arg1	arg2
(0) unimod	B	-
(1) +	C	D
(2) *	(0)	(1)
(3) =	A	

2)  $A[L] = B$

OP	arg1	arg2
(0) [L]	A	1
(1) =	(0)	B

3)  $A = B[L]$

OP	arg1	arg2
(0) [L]	B	1
(1) =	(0)	A

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

$t_1 = b - c$

$t_2 = a * t_1$

$t_3 = a + t_2$

$t_4 = t_1 * d$

$t_5 = t_3 + t_4$

OP	arg1	arg2
(0) unimod	b	c
(1) *	a	(0)
(2) +	a	(1)
(3) *	(0)	d
(4) +	(2)	(3)

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

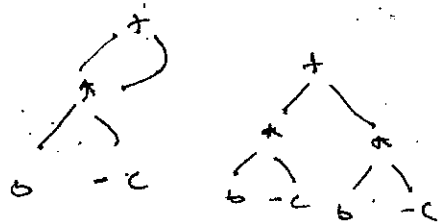
Dag:

$$t_1 = -c$$

$$t_2 = b * t_1$$

$$t_3 = t_2 + t_2$$

$$a = t_3$$



3 addr:

$$t_1 = -c$$

$$t_2 = b * t_1$$

$$t_3 = -c$$

$$t_4 = b * t_3$$

$$t_5 = t_2 + t_4$$

$$a = t_5$$

OP

arg 1

arg 2

OP	arg 1	arg 2
(0) unary	c	-
(1) *	b	(0)
(2) unary	c	-
(3) *	b	(2)
(4) +	(1)	(3)
(5) =	a	(4)

### Indirect Triplee

Consists of 1 key pointers to triplee rather than key the triplee themselves

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

35	(0)
36	(1)
37	(2)
38	(3)
39	(4)
40	(5)
	.....



OP	arg 1	arg 2
(0) unary	c	-
(1) *	b	(0)
(2) unary	c	-
(3) *	b	(2)
(4) +	(1)	(3)
(5) =	a	(4)

### Static Single Assignment Form

SSA is an intermediate rep<sup>n</sup> that facilitates certain code optimization. The diff. b/w SSA & 3-addr code is

(a) all assignments in SSA are to var with distinct names

(b) SSA uses a notational convention called  $\phi$  functions to combine the two def<sup>n</sup> of x

Q) Ex: ~~(a+b-c)~~ (a+b-c) + (e - a+b-c\*d)

P = a+b  
 q = P-c  
 r = q\*d  
 s = e-r  
 v = P+q

P<sub>1</sub> = a+b  
 q<sub>1</sub> = P<sub>1</sub>-c  
 r<sub>2</sub> = q<sub>1</sub>\*d  
 r<sub>3</sub> = e-r<sub>2</sub>  
 q<sub>2</sub> = P<sub>3</sub>+q<sub>1</sub>  
 SSA

⊥  
 t<sub>1</sub> = a+b  
 t<sub>2</sub> = t<sub>1</sub>-c  
 t<sub>3</sub> = t<sub>2</sub>\*d  
 t<sub>4</sub> = e-t<sub>3</sub>  
 t<sub>5</sub> = t<sub>4</sub>+t<sub>2</sub>

Q) Ex: if (flag) x = -1; else x = 1;

y = x \* a;

SSA

if (flag) x<sub>1</sub> = -1; else x<sub>2</sub> = 1;  
 x<sub>3</sub> =  $\begin{cases} x_1 & \text{true} \\ x_2 & \text{false} \end{cases}$

Ex: Ex: ① a + (b+c) AST, tax tree, DAG, QT, T, IT vj

② a[i] = b\*c - b\*d

③ x = \*p + \*y

Type Design  
 type checks - runtime  
 type appn - static  
 type expn - no expn  
 type equivalence -  
 type checks - syntn  
 type convn - narrow

GURUPRASAD. S.  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64,  
 Ph: 9886760776

Prof GURUPRASAD S, BMSIT&M

<https://hemanthrajhemu.github.io>



# Type Declaration

The appl<sup>n</sup> of types can be grouped under (i) type checking & (ii) Translation

→ Type checking: we use logical rules to reason about the behaviour of a program at run time

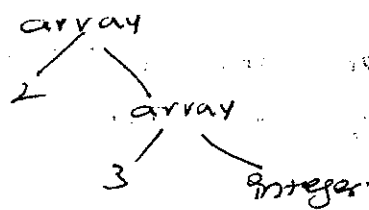
Ex:  $x \& y$  needs both opnds to be logical/boolean

→ Translation appl<sup>n</sup>: from the type of name, a compiler can determine the storage that will be needed for the name at runtime.

## Type Expression

Types have structure, which we rep<sup>n</sup> using "type expression" - a type expr<sup>n</sup> is either formed by a basic type or formed by applying an operator called a type constructor to a type expr<sup>n</sup>

Ex:-  $\text{int}[2][3]$  can be read as array of 2 array of 3 integers.  
Each is written as  $\text{array}(2, \text{array}(3, \text{integer}))$



operator array takes two parameters  
number & type.

The following are def<sup>n</sup> of type expr<sup>n</sup>:

- \* A basic type is a type expr<sup>n</sup> like int, float, char, boolean, void.
- \* A type name is type expr<sup>n</sup> -  $\text{int float}$
- \* A type expr<sup>n</sup> can be formed by applying the array type constructor to number & a type expr<sup>n</sup>.  $\text{int}[2][3]$
- \* Structural record is a data structure with named fields. A type expr<sup>n</sup> can be formed by applying record type constructor to field names & their types.  $\text{int a}$ ,  $\text{string b}$
- \* A type expr<sup>n</sup> can be formed by using union type constructor → for function types.  $\text{int fun}(a, b)$
- \* A convenient way to rep<sup>n</sup> a type expr<sup>n</sup> is to use graph or value number method.

# Type Equivalence

Many type checking rules have the form "if two type Expr<sup>n</sup> are equal then return a certain type else error".  
 This rule is used to say whether two type Expr<sup>n</sup> are equivalent.

When the Expr<sup>n</sup> are rep<sup>d</sup> by graphs two types are structurally equivalent iff the following cond<sup>n</sup> are true.

- \* They are the same basic type.
- \* They are formed by applying the same constructor to structurally eq<sup>t</sup> types.
- \* one is a type name that denotes other.

## Storage layout for local names.

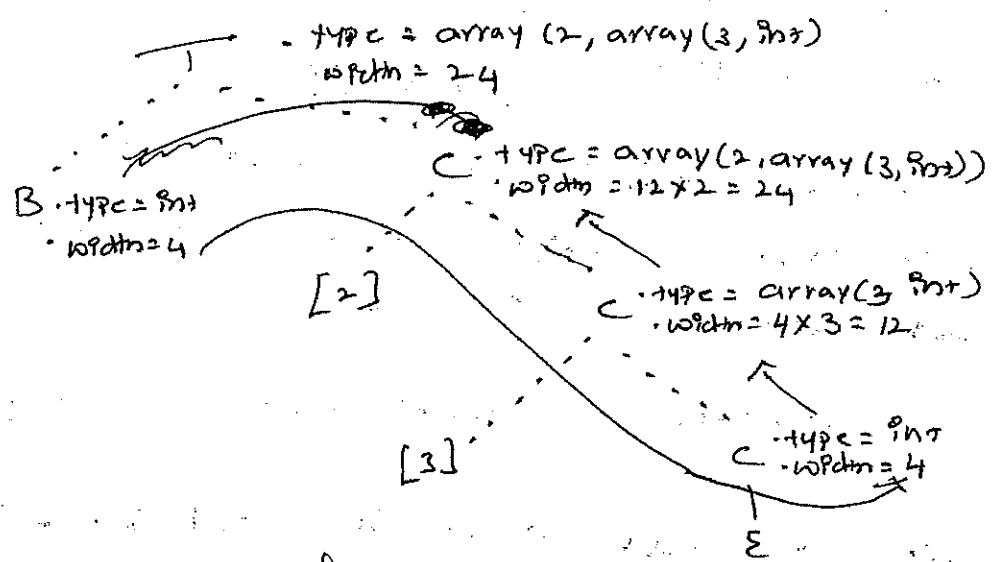
From the type of a name, we can determine the amount of storage that will be needed for the name at runtime.

At compile time, we can use these amounts to assign each name a relative address. The type & relative addr. are saved in symbol table entry for name.

The width of a type is the no. of storage units needed for objects of that type.

Form	Semantic Rule
$T \Rightarrow B C$	$\{ t = B.type; w = B.width; \}$
$B \Rightarrow int$	$B.type = integer \quad B.width = 4;$
$B \Rightarrow float$	$B.type = float \quad B.width = 8;$
$C \Rightarrow E$	$C.type = t \quad C.width = w;$
$C \Rightarrow [num] C_1$	$C.type =$ $\{ array (num.val, C_1.type);$ $C.width = num.val * C_1.width;$

int [2] [3]



### Sequence of Declaration

We can calculate the relative address of a variable in SymTab. This is achieved by using the method top.put which creates a SymTab entry for any variable. top indicates the current symbol table.

Initially offset is 0. As each new name is entered into SymTab with its relative address is set to current value of offset which is incremented by width of it.

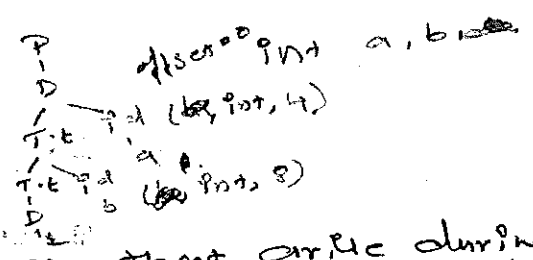
Ex: - P → D { offset = 0 }

D → T id { top.put (id.lexeme, T.type, offset) }

offset = offset + T.width

D → P<sub>1</sub>

D → E



### Translation of Expression

Here we explore kernel that arrive during the

translation of expression and statements.

An Expr<sup>n</sup> with more than one operator will

translate into Expr<sup>n</sup> with at most one operator per sub-expression

# Translation of Expressions

## Operator within Expression ⊕

The SDP builds up 3-addr code for an assignment stmt  $S$  using attrib code for  $S$  and 3-addr code for  $E$ .

$S$ .code =  $E$ .code - 3-addr code for  $S \rightarrow E$

$E$ .addr = addr that will hold value of  $E$ .

Prod $n$	Semantic Rules
$S \rightarrow id = E$	$S$ .code = $E$ .code    $gen(top.get(id.lexeme) = E.addr)$
$E \rightarrow E_1 + E_2$	$E$ .addr = new Temp() $E$ .code = $E_1$ .code    $E_2$ .code    $gen(E.addr = E_1.addr + E_2.addr)$
$E \rightarrow - E_1$	$E$ .addr = new Temp() $E$ .code = $E_1$ .code    $gen(E.addr = - E_1.addr)$
$E \rightarrow (E)$	$E$ .addr = $E_1$ .addr $E$ .code = $E_1$ .code
$E \rightarrow id$	$E$ .addr = $top.get(id.lexeme)$ $E$ .code = ...

## Add using array elements

Array elements can be accessed quickly if they are stored in blocks of sequential contiguous locations, i.e. 0, ..., n-1  
 if width of each array element is  $w$  then  $i$ th element of array  $A$  begins at location

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

where  $\text{base}$  is the relative addr of  $A[0]$  i.e.  $A$

is for 2-dimensional array

$$A[i_1][i_2] \text{ the relative addr is } \text{base} + i_1 \times w_1 + i_2 \times w_2$$

In general for  $k$  dimensional array

$$\text{base} + i_1 \times w_1 + i_2 \times w_2 + \dots + i_k \times w_k$$

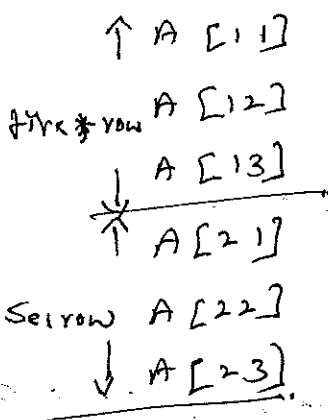
A two dimensional array is normally stored in one of two forms

(1) row major (row by row)

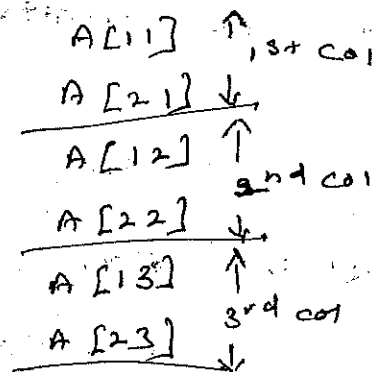
(2) column major (column by column)

11	12	13
21	22	23

array  $A[2][3]$



Row major



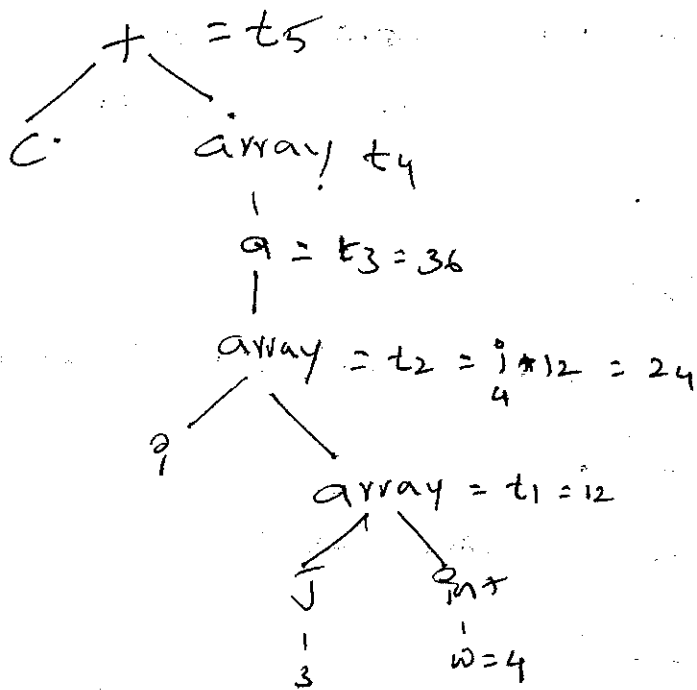
Col major

GURUPRASAD. S.  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64,  
 Ph: 9886760776

Prof GURUPRASAD S, BMSIT&M

<https://hemanthrajhemu.github.io>

Consider the exprn  $C + a[i][j]$  with  $i=2$   $j=3$   
 the parse tree is



$$t_1 = j * 4 = 12$$

$$t_2 = i * t_1$$

$$i = 2 \quad \frac{2 * 12}{4} = 24$$

$$t_3 = t_1 + t_2$$

$$t_4 = a[t_3]$$

$$t_5 = C + t_4$$

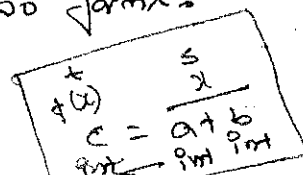
## Type checking

To do type checking, a compiler needs to assign a type expr to each component of the src prog, the compiler must then determine that these type expr conform to a collection of logical rules that called the type system for src prog. Type checking has potential for catching errors in prog.

## Rules for type checking

Type checking can arise in two forms

- (i) Synthesis
- (ii) Inference



(i) Type synthesis builds up the type of an expr from the types of its sub exprs. It requires names to be declared before they are used.  
 The type  $E_1 + E_2$  is defined in terms of types of  $E_1$  &  $E_2$

A typical rule for type synthesis,

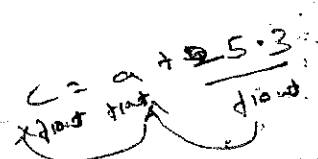
If  $f$  has type  $S \rightarrow t$  and  $x$  has type  $S$ ,

then exprn  $f(x)$  has type  $t$

where  $f$  &  $x$  denote exprn and  $S \rightarrow t$  denotes a function  $f$  from  $S$  to  $t$

(ii) Type inference determines the type of a language construct from the way it is used.

Eg:-  $\text{null}(x)$   $\text{let } x \text{ let } |x| + |x| \text{ end}$  then  $x$  must be of some type.



A typical rule for type inference is,

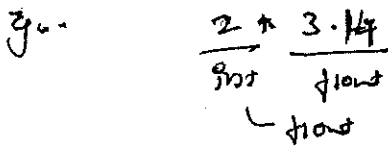
If  $f(x)$  is an exprn,

then for some  $A$  and  $B$ ,  $f$  has type  $A \rightarrow B$  and  $x$  has type  $A$

is used in languages like C where decl'n of name is not required.

### Type Conversion:

Type conversion is needed when types of the operands in an exprn are different because the operations / members will be diff for diff type.



so  $t_1 = \text{float}(2)$

$t_2 = t_1 * 3.14$

The rule is if  $E \rightarrow E_1 + E_2$

if  $(E_1 \text{ type} = \text{integer and } E_2 \text{ type} = \text{integer})$   $E \text{ type} = \text{integer}$

else if  $(E_1 \text{ type} = \text{float and } E_2 \text{ type} = \text{integer})$   $E \text{ type} = \text{float}$

There are 100 major schemes for type conversion

**GURUPRASAD, S.**  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-64,  
Ph: 9886760776



(i) widening

(ii) narrowing

**Prof GURUPRASAD S, BMSIT&M**

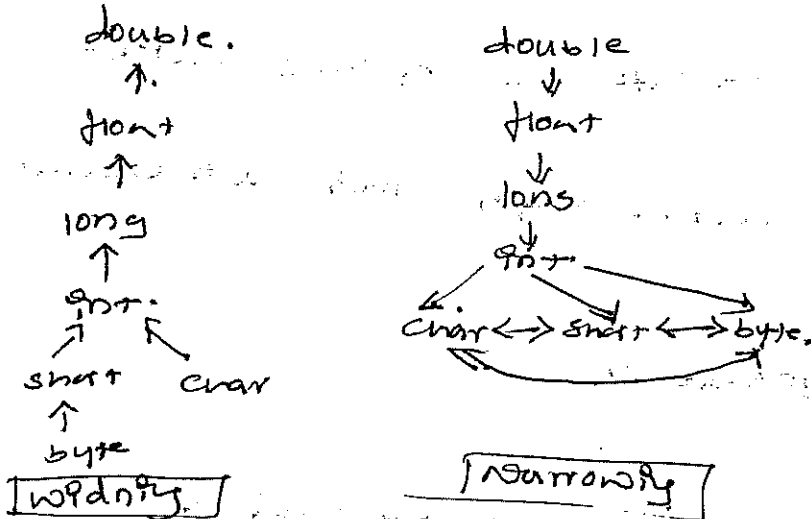
<https://hemanthrajhemu.github.io>

implicit / coercion conversion - done by compiler.

explicit / cast conversion - done by Programmer.

Widening - information re to preserve info so convert from lower to higher.

Narrowing - may lose info so convert from higher to lower.



We can use two functions  $\max(t_1, t_2)$  - returns max value of type.  
 $\text{widen}(a, t, w)$  a - addr t - type w - widened type.

### Unification:

Unification is the problem of determining whether two expressions  $S$  &  $t$  can be made identical by substituting expr for the variables in  $S$  &  $t$ . If  $S$  and  $t$  have constants but no variables, then  $S$  &  $t$  unify iff they are identical.

Eg:- consider two expressions.

$$((d_1 \rightarrow x_2) \times \text{let}(d_3)) \rightarrow \text{let}(x_2) \quad \text{--- (1)}$$

$$((d_3 \rightarrow d_4) \times \text{let}(d_3)) \rightarrow d_5 \quad \text{--- (2)}$$

The following substitution  $S$  is the most general unifier for these expressions

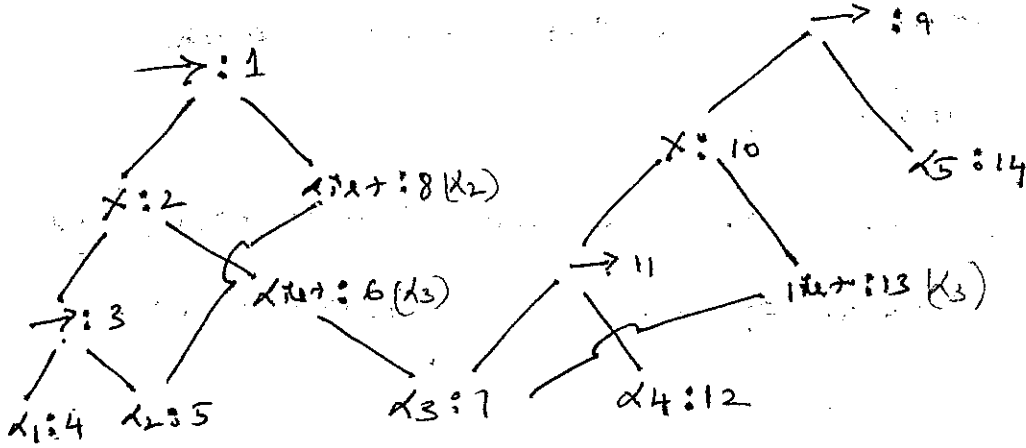
$x$	$S(x)$
$d_1$	$x_1$
$d_2$	$d_2$
$d_3$	$d_1$
$d_4$	$d_2$
$d_5$	$\text{let}(x_2)$



The substitution made in two type expr is.

$$(K_1 \rightarrow \alpha_2) \times \text{let}(\alpha_1) \rightarrow \text{let}(\alpha_2) \quad \text{--- (3)}$$

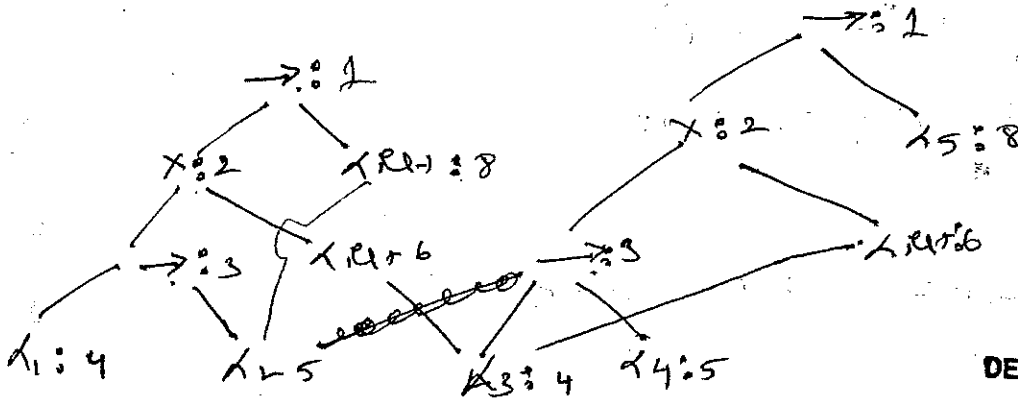
Suppose that two expr are rep<sup>d</sup> by initial graph where each node is in its own equivalence class.



$x = a + b$   
 $y = c + d$   
 $(1, 9)$      $(4, 7)$   
 $(2, 10)$      $(5, 4)$   
 $(8, 14)$      $x = a + b$   
 $(3, 11)$      $z = c + d$   
 $(6, 13)$

Now if alg of unifications is applied to compute (1,9), it notes that nodes 1 & 9 both rep<sup>n</sup> same operators i.e. let merges into same equivalence class then it calls unify(2,5,10) and unify(8,14)

The equivalence class after unifications is the resulting graph.



**GURUPRASAD. S.**  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64,  
 Ph: 9886760776

### Algorithm

**Input:** a graph rep<sup>d</sup> by a dyic and a pair of nodes m & n to be unified.

**Output:** boolean value true if the expr<sup>n</sup> rep<sup>n</sup> by the nodes m & n

unify else false.

**Prof GURUPRASAD S, BMSIT&M**

<https://hemanthrajhemu.github.io>

isectnode: boolean unify (Node m, Node n)

```
{  
    s = find(m); t = find(n);  
    if (s == t) return true;  
    else if (nodes s and t represent same base type)  
        return true;  
    else if (s is an op-node with children s1, t1 and  
            t is an op-node with children t1, t2)
```

```
        {  
            union(s, t);  
            return unify(s1, t1) and unify(s2, t2);  
        }
```

```
    else if s or t represent a variable  
        {  
            union(s, t);  
            return true;  
        }
```

```
    else return false;
```

find(n): returns the representative node of equivalence class containing n

union(m, n): merges the equivalence class containing node m with n

2016-20 Batch, Dept. of EEE

SL. USN	NAME OF THE CANDIDATE	QUOTA	I SGPA	II SGPA	III SGPA	IV (800)	V	VI	VII	VIII	Total	Percentage	Class	Only 1st yr per	Only 2nd yr per	Only 3rd yr per	REMARKS
1	1BY16EE001	AICTE-IKSS	6.08	4.25	422	452								5.17	54.63		
2	1BY16EE002	Comed-K	8.67	9.08	671	689								8.88	85.00		
3	1BY16EE003	Comed-K	6.67	6.42										6.55			Left
4	1BY16EE004	CET	7.5	8.58	584	644								8.04	76.75		
5	1BY16EE005	AKANGSHA P	1.75	1.68										1.72			Lost Eligibility
6	1BY16EE006	ANAMIKA BASU	1.25	3.1	257	299								2.18	34.75		Lost Eligibility
7	1BY16EE007	ANIRUDH ARAVIND	7.5	8.25	529	572								7.88	68.81		
8	1BY16EE008	ANIRUDHA K S	5.17	6	461	527								5.59	61.75		
9	1BY16EE009	AWATI SPANDANA	6.42	8.08	522	552								7.25	67.13		
10	1BY16EE010	B ASHWINI	5.58	7.14	553	588								6.36	70.06		
11	1BY16EE011	B MANJUNATH	8.75	8.75	651	673								8.75	82.75		
12	1BY16EE012	BALAJI SUBRAMANYAM M K	8.17	7.75	554	567								7.96	70.06		
13	1BY16EE013	BHARGAVREDDY B	7.67	7.25	500	530								7.46	64.38		Lost Eligibility
14	1BY16EE014	BHASKARA RAJU A V	5.33	4.07	408	292								4.70	43.75		
15	1BY16EE015	CHAITRHA B C	6.92	7	489	554								6.96	65.19		
16	1BY16EE016	DEEPA M S	8.33	8.58	595	645								8.46	77.50		
17	1BY16EE017	ENTHA ROHITH	8.25	8.17	601	645								8.21	77.88		COB
18	1BY16EE018	G SANKESH BOTHRA															
19	1BY16EE019	SAUTAMA BHARADWAJ	8.42	8.33	593	626								8.38	76.19		
20	1BY16EE020	GOVARDHAN	6.33	6.25	482	553								6.29	63.44		
21	1BY16EE021	H VISHAL	3.75	3.2	348	357								3.48	44.06		Lost Eligibility
22	1BY16EE022	HARRISHAN U	8.58	8.58	626	660								8.58	80.38		
23	1BY16EE023	HARSHIT VAISH	8.25	8.25	606	644								8.25	78.13		
24	1BY16EE024	KACHARLA BALASAI NIKHIL	7.17	6.58	502	506								6.88	63.00		
25	1BY16EE025	KALPIT CHATURVEDI	5.08	4.17	397	427								4.63	51.50		Left
26	1BY16EE026	KALSHIK GOWDA H N	7.42	8.42	555	438								7.92	62.06		
27	1BY16EE027	KAVYA P	8	8.67	559	605								8.34	72.75		
28	1BY16EE028	KEERTHI GURURAJ G S	7.33	8.58	551	580								2.82	30.19		Lost Eligibility
29	1BY16EE029	KISHUJ SINHA	3.25	2.38	249	234								8.00	70.44		
30	1BY16EE030	KUSHAGRA DHAWAN	7.58	8.42	554	573								5.21	46.81		
31	1BY16EE031	MADHU CHANDRA M	6.33	4.08	355	394								7.25	64.13		
32	1BY16EE032	MADHUSUDHAN H K	7.17	7.33	485	541								2.34	47.06		Lost Eligibility
33	1BY16EE033	MALLARI SURESH	1.17	3.5	333	420								6.79	71.31		
34	1BY16EE034	MOHAMMAD ADIL ANSARI	6.5	7.08	566	575								8.00	76.63		
35	1BY16EE035	NAVEDYA DIJA	7.67	8.33	569	657								7.67	65.69		
36	1BY16EE036	NIKITA CHAUHAN	7.92	7.42	473	578								5.71	53.44		Lost Eligibility
37	1BY16EE037	NISHA CHAURASIA	6.17	5.25	387	468											Lost Eligibility
38	1BY16EE039	PULKIT KUMAR DAGUR	0	0													Lost Eligibility
39	1BY16EE040	RAKESH M	0.5	0													
40	1BY16EE041	ROHAN CHINNI C L	7.33	8	582	645								7.67	76.69		
41	1BY16EE042	RUCHITHA D	5.33	7	443	531								6.17	60.88		
42	1BY16EE043	RUPANISHA KHARE	4.33	4.29	367	417								4.31	49.00		
43	1BY16EE044	SADIQ HASAN ABBASI	7.42	7.67	566	575								7.55	71.31		
44	1BY16EE045	SARTHAK GUPTA	6.58	7.67	516	564								7.13	67.50		
45	1BY16EE046	SAUMYA MISRA	6.33	6.92	454	526								6.63	61.25		
46	1BY16EE047	SAVALEE SANJAY KHANDAL	7.5	9.08	549	611								8.29	72.50		
47	1BY16EE048	SESHA GOPAL S	7.08	6.58	536	603								6.83	71.19		



## Control Flow

The translation of statements such as if-else-stmt and while-stmt related to the translation of boolean expr.

In programming, boolean exprn are used to

(i) Alter the flow of control - boolean expr are used as condnl expr in stmts that alter the flow of control. For eg: if (E) then S, the expr E must be ~~evaluated~~ true if stmt S is reached.

(ii) Compute logical values - a boolean expr can return true or false values.

## Boolean Expressions

Boolean exprn are composed of boolean operators. The (E<sub>1</sub>, ||, !, & ) applied to elements that are boolean variables or relational expressions, relational expressions are of the form E<sub>1</sub> rel op E<sub>2</sub>, where E<sub>1</sub> & E<sub>2</sub> are arithmetic exprn.

Consider the boolean exprn gen by the grammar.

$B \rightarrow B || B | B \& \& B | ! B | (B) | E_1 \text{ rel } E_2 | \text{true} | \text{false}$ .

The rel op may be <, <=, =, !=, >, >= rel op rel op

In exprn B<sub>1</sub> || B<sub>2</sub> if either B<sub>1</sub> or B<sub>2</sub> is true entire exprn is true but in B<sub>1</sub> & B<sub>2</sub> if either B<sub>1</sub> or B<sub>2</sub> is false entire exprn is false.

The semantic defn of the prog lang determines whether all parts of boolean expr must be evaluated.

## Short Circuit Code

In short-circuit (jumping) code, the boolean operators  $\&\&$ ,  $\|\|$ ,  $!$  translate into jumps, - the operators themselves do not appear in the code. Instead the value of a boolean expr is repn by a fallthru in the code sequence.

Ex:-

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

translated into

$\text{if } (x < 100) \text{ goto } L_2$

$\text{if } \text{false } x > 200 \text{ goto } L_1$

$\text{if } \text{false } x \neq y \text{ goto } L_1$

$L_2 : x = 0$  - true

$L_1 :$  - false

## Flow Control Statements

Now consider the translation of boolean exprs into 3-addr code in the context of the grammar.

$S \rightarrow \text{if } (B) S_1$

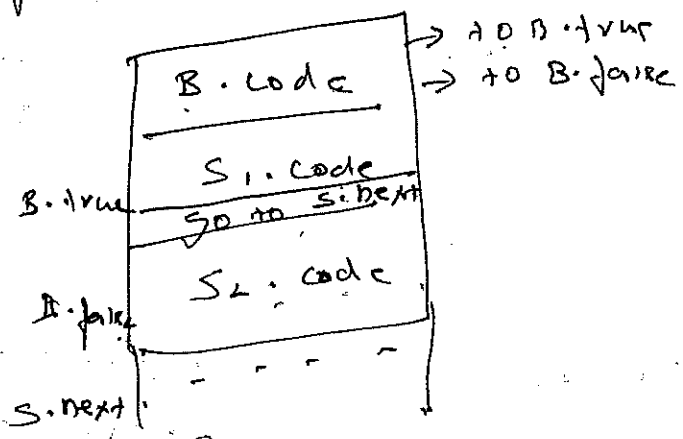
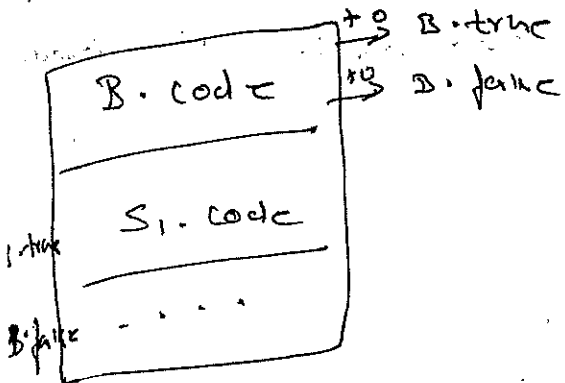
$S \rightarrow \text{if } (B) S_1 \text{ else } S_2$

$S \rightarrow \text{while } (B) S_1$

In this grammar 'B' repn boolean exprs & 'S' repn a statement both B & S have synthesized attrib code

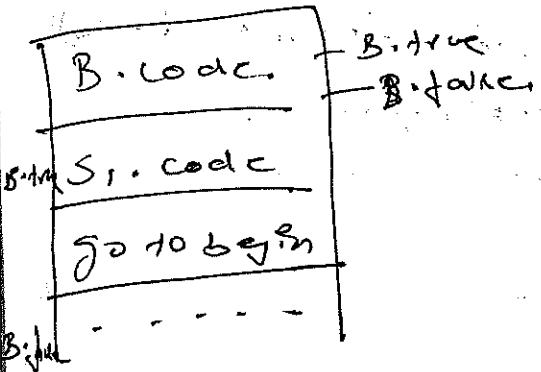
which gives the transn into 3-addr code

The traversal of  $\{j(B)S_1$  consisting of B-code followed by  $S_1$  code within B-code are jumps based on the value of B. If B returns control flow to the first item of  $S_1$  code & if B returns control flow to the item immediately followed  $S_1$  code.



(b) if - else

(a) if.



GURUPRASAD. S.  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-64,  
Ph: 9886760776

(c) while.

The SDD for flow of control statements produce 3- code for each exp.

Prodn	Semantic Rule
$P \rightarrow S$	$S.next = newLabel()$ $P.code = S.code    label(S.next)$
$S \rightarrow assign$	$S.code = assign.code$
$S \rightarrow \{j(B)S_1$	$B.true = newLabel()$ $B.false = S1.next$ $S.code = B.code    label(B.true)    S1.code$

Prof GURUPRASAD S, BMSIT&M



$S \rightarrow \text{if}(B) S_1 \text{ else } S_2$

$B.\text{true} = \text{newlabel}()$

$B.\text{false} = \text{newlabel}()$

$S_1.\text{next} = S_2.\text{next} = S.\text{next}$

$S.\text{code} = B.\text{code} \parallel \text{label}(B.\text{true}) \parallel S_1.\text{code}$

$\parallel \text{gen}(\text{'goto' } S.\text{next}) \parallel \text{label}(B.\text{false}) \parallel S_2.\text{code}$

$S \rightarrow \text{while}(B) S_1$

$\text{begin} = \text{newlabel}()$

$B.\text{true} = \text{newlabel}()$

$B.\text{false} = S_1.\text{next}$

$S_1.\text{next} = \text{begin}$

$S.\text{code} = \text{label}(\text{begin}) \parallel B.\text{code} \parallel \text{label}(B.\text{true})$

$\parallel S_1.\text{code} \parallel \text{gen}(\text{'goto' } \text{begin})$

$S \rightarrow S_1 S_2$

$S_1.\text{next} = \text{newlabel}()$

$S_2.\text{next} = S.\text{next}$

$S.\text{code} = S_1.\text{code} \parallel \text{label}(S_1.\text{next}) \parallel S_2.\text{code}$

$\text{newlabel}()$  creates a new label each time it is called.  
 $\text{label}(L)$  attaches label  $L$  to next 3-addr stmt to be generated.

**GURUPRASAD. S.**  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-64,  
Ph: 9886760776

**Prof GURUPRASAD S, BMSIT&M**

<https://hemanthrajhemu.github.io>



The semantic rules for bool expr translation

Semantic Rule

Rules

$B \rightarrow B_1 \parallel B_2$

$B_1.true = B.true$   
 $B_1.false = newlabel()$   
 $B_2.true = B.true$   
 $B_2.false = B.false$   
 $B.code = B_1.code \parallel label(B_1.false)$   
 $\parallel B_2.code$

GURUPRASAD S.  
 ASSISTANT PROFESSOR  
 DEPT OF COMPUTER SCIENCE  
 BMSIT, BANGALORE

$B \rightarrow B_1 \&\& B_2$

$B_1.true = newlabel()$   
 $B_1.false = B.false$   
 $B_2.true = B.true$   
 $B_2.false = B.false$   
 $B.code = B_1.code \parallel label(B_1.true) \parallel B_2.code$

$B \rightarrow ! B_1$

$B_1.true = B.false$   
 $B_1.false = B.true$   
 $B.code = B_1.code$

$B \rightarrow E_1 \text{ relop } E_2$

$B.code = E_1.code \parallel E_2.code \parallel$   
 $gen('if' E_1.addr \text{ relop } E_2.addr \text{ goto}$   
 $B.true) \parallel$   
 $gen('goto' B.false)$

$B \rightarrow true$

$B.code = gen('goto' B.true)$

$B \rightarrow false$

$B.code = gen('goto' B.false)$

if  $(x < 100 \vee x > 200 \wedge x \neq y)$   $x = 0$

\*

if  $x < 100$  goto L2

goto L3

L3: if  $x > 200$  goto L4

goto L1

L4: if  $x \neq y$  goto L2

goto L4

L2:  $x = 0$

L1:

**GURUPRASAD. S.**  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-64,  
Ph: 9886760776

### Boolean Value & Jumping Code.

The use of Boolean expressions to alter the flow of control is similar. A boolean expression may also be evaluated for its value as in assignment statements even as  $x = \text{true}$  or  $x = a \& b$ .

A clear way of handling both roles of bool expressions is first build a syntax tree for expressions using either of following approaches

(i) Use two passes: Construct a complete syntax tree first. Then walk the tree in depth first order, applying the translations specified by semantic rules.

(ii) Use one pass for stmts! but two passes for expressions with this approach, we would translate E in while (E)s, before S<sub>i</sub> is examined.

## Basis Patching :

The key problem when generating code for bool expr is flow of control stmt is that of matching a jump stmt with the target of jump. So the approach to solve is basis patching, in which labels of jumps are treated as synthesized attributes. Specifically when a jump is generated, the targets of jump are temporarily left unspecified. Each such jump is put on a list of jumps of all jumps on a list have the same target label. The labels are filled when proper label can be determined.

## One Pass Code Generation using Back Patching

Back Patching can be used to generate code for bool expr & flow of control stmt in one pass.

The ~~needed~~ synthesized attributes true list & false list of NT 'B' are used to manage labels in jumping code for boolean expr.

B.true list - list of ~~targets~~ <sup>jump</sup> to reach if B is true

B.false list - list of ~~targets~~ <sup>jump</sup> to reach if B is false

S.next list - list of jumps to instrn immediately follows

Instrns are generated into an array & labels will be indices into this array & jns are used to manipulate list of jumps.

(i) make list(i) :- Create the list containing only i, an index into the array of instrns, & it returns ptr to newly created list

(ii) merge(p1, p2) :- Concatenate the list pointed to by p1 & p2 & return a ptr to concatenated list

(iii) back patch(p, i) :- Invert i as the target label for solving the list

# Basis Function for Boolean Exprn

The following are the translation scheme suitable for generating code for bool exprn during BU pass.

M - marker NT :- called semantic action to pick up at appropriate times & the index of next lrt to be generated.

B - Boolean E - Exprn

$B \rightarrow B_1 \parallel M B_2 \mid B_1 \& M B_2 \mid \neg B_1 \mid (B) \mid E_1 \vee E_2 \mid \text{true} \mid \text{false}$

$M \rightarrow \epsilon$

## The translation schemes

$B \rightarrow B_1 \parallel M B_2$

$\{$  basis Patch (B, false lrt, M, next);

$B$ . true lrt = merge ( $B_1$ . true lrt,  $B_2$ . true lrt);

$\rightarrow B$ . false lrt =  $B_2$ . false lrt;  $\}$

$B \rightarrow B_1 \& M B_2$

$\{$  basis Patch ( $B_1$ . true lrt, M, next);

$B$ . true lrt =  $B_2$ . true lrt;

$\rightarrow B$ . false lrt = merge ( $B_1$ . false lrt,  $B_2$ . false lrt);  $\}$

$B \rightarrow \neg B_1$

$\{ B$ . true lrt =  $B_1$ . false lrt;

$\rightarrow B$ . false lrt =  $B_1$ . true lrt;  $\}$

$B \rightarrow (B_1)$

$\{ B$ . true lrt =  $B_1$ . true lrt;

$\rightarrow B$ . false lrt =  $B_1$ . false lrt;  $\}$

$B \rightarrow E_1 \text{ rel } E_2$

```

{ B.true | ret = make | ret (next | next);
  B.false | ret = make | ret (next | next);
  emit ('if 'E1.addr' rel op E2.addr' goto -');
  emit ('goto -');

```

$B \rightarrow \text{true}$

```

{ B.true | ret = make | ret (next | next);
  emit ('goto -');

```

$B \rightarrow \text{false}$

```

{ B.false | ret = make | ret (next | next);
  emit ('goto -');

```

$M \rightarrow \epsilon$

```

{ M.instr = next | next;

```

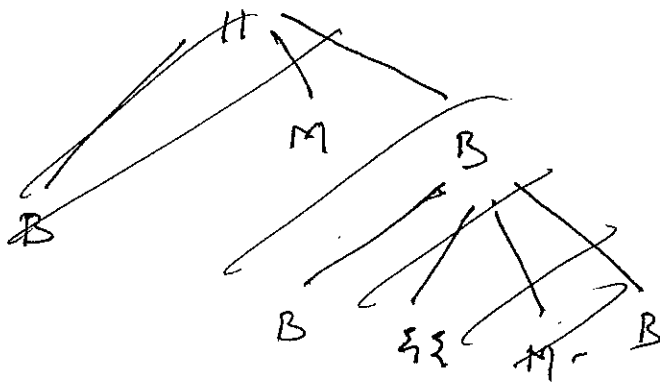
Basis Patching:

in if (B) S ... Jump to S if B is true else jump to next stmt of S  
 so B must be translated first then label has to be found where to  
 jump if true & if false if these labels are forward referenced  
 then we need two passes to compute & align labels so to avoid  
 this problem Basis Patching is used wherein jumps are left  
 unspecified temporarily — see back

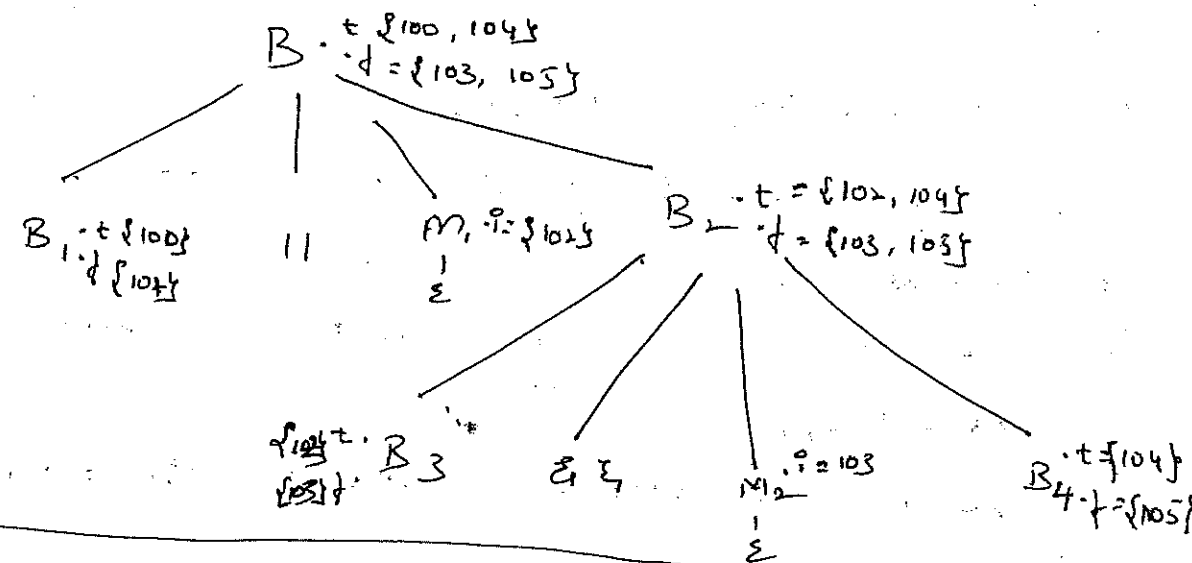
**GURUPRASAD. S.**  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64,  
 Ph: 9886760776

eg:-

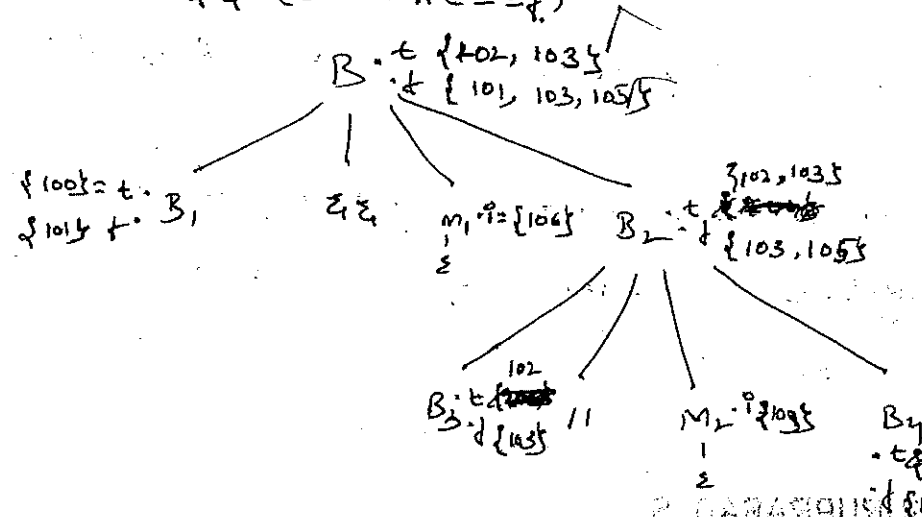
$x < 100 \parallel x > 200$  & &  $x = y$



100: if  $x < 100$  goto - {100} <sup>fall</sup>  
 101: goto - {102}  
 102: if  $x > 200$  goto - {104}  
 103: goto - {106}  
 104: if  $x = y$  goto - {100}  
 105: goto 106  
 fall 106:  
 ...  
 10200:



eg:-  $a == b$  & &  $(c == d \parallel e == f)$



100: if  $(a == b)$  goto 101  
 101: goto 106  
 102: if  $(c == d)$  goto 200  
 103: goto 104  
 104: if  $(e == f)$  goto 200  
 105: goto 106

Case 3

- ③  $(a == b \parallel c == d) \parallel e == f$
- ④  $(a == b \& \& c == d) \& \& e == f$

DEPT OF SOFTWARE ENGINEERING  
 ASSISTANT PROFESSOR  
 GURUPRASAD S.

## Switch Statement

The switch or case statements are available in a variety of languages as shown below

Switch (C)

{

Case  $V_1$  :  $S_1$

Case  $V_2$  :  $S_2$

⋮

Case  $V_{n-1}$  :  $S_{n-1}$

default :  $S_n$

}

## Translation of Switch Statement

The basic translation of switch code is:

- ① Evaluate the expr E
- ② Find the value  $V_j$  in the list of case that is the value of expr
- ③ Execute the stmt  $S_j$  associated with value found.

## SDT of Switch Stmt

The switch stmt is translated to the intermediate code as shown.

The test appears at the end so that the simple code generator can recognize the multi way branch and can generate efficient code for it.

**GURUPRASAD. S.**  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-64,  
Ph: 9886760776

Code to Evaluate E into t

go to test

L<sub>1</sub>: code for S<sub>1</sub>

go to next

L<sub>2</sub>: code for S<sub>2</sub>

go to next

:

L<sub>n-1</sub>: code for S<sub>n-1</sub>

go to next

L<sub>n</sub>: code for S<sub>n</sub>

go to next

test: if t = v<sub>1</sub> goto L<sub>1</sub>

if t = v<sub>2</sub> goto L<sub>2</sub>

if t = v<sub>n-1</sub> goto L<sub>n-1</sub>

goto L<sub>n</sub>

next:

## Intermediate Code for Procedures

In 3-address code a function call is translated into the evaluation of parameters in preparation for a call, followed by call itself, and the parameters are passed by value.

Eg:- if a re array of int & f is a fn from int to int

then the call

$m = f(a, b)$

translate into following 3-address code:

$t_1 = r_4$

$m = t_3$

$t_2 = a[t_1]$

Param t<sub>2</sub>

$\& = \text{call } f$



Function definitions & calls can be translated using concepts:

- ① Function types: the type of a fn must encode the return type & type of formal parameters. let void be a special type that req<sup>n</sup> no parameters or no return type.
- ② Symbol Table: let 'S' be the top<sup>†</sup> symbol table when fn def<sup>n</sup> is reached the fn name is inserted into 'S' for use in the rest of prog.
- ③ Type checking: within expr<sup>n</sup>, a fn is treated like any other operator. for eg if a fn with parameter of type real then the integer 2 is coerced to a real in the call f(2).
- ④ Function call: when generating 3-address code for fn call  $f(e_1, e_2, \dots, e_n)$  it is sufficient to generate 3-addresses for evaluating or reducing parameters  $e_i$  to address, followed by param instr<sup>n</sup> for each parameter.

---

End of Part (b) unit (6)

**GURUPRASAD. S.**  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-64,  
Ph: 9886760776

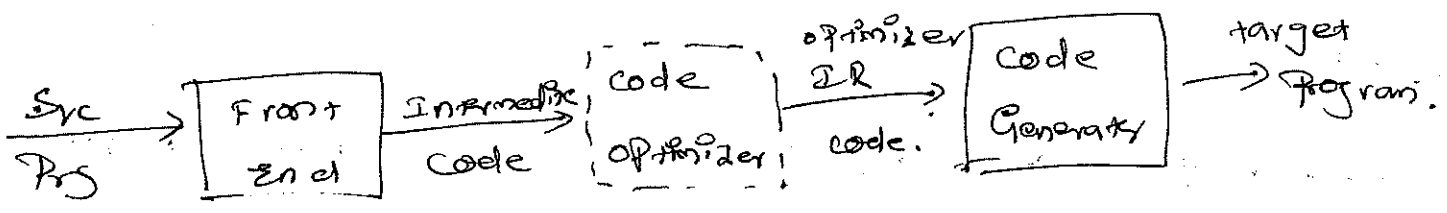
**Prof GURUPRASAD S, BMSIT&M**

<https://hemanthrajhemu.github.io>



# Code Generation

A Code Generator is the final phase of Compiler model. It takes IR code. Produced by front end & relevant ST info to produce semantically equivalent target program.



A code generator has three primary tasks: instruction selection, register allocation and instruction ordering.

## Issues in the Design of a Code Generator

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

- (i) Input to the code generator.
- (ii) IR code
- (iii) Instruction set
- (iv) Register allocation
- (v) Instruction order.

The input to the code generator is IR code of src prog, produced by front end. There are many choices for IR include 3-addr, register, quad, triple, indirect triple, virtual m/c reg such as byte codes, stack-m/c code, graphical reg such as syntax tree, DAG.

We assume that the front end has scanned & parsed & translated the src prog into a relatively low level IR. We also assume that syntactic, semantic errors have been detected already.

**GURUPRASAD. S.**  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64,  
 Ph: 9886760776

**Prof GURUPRASAD S, BMSIT&M**

2) Programs to be compiled separately.  
3) Programs are re-allocatable m/c lay for an e/p allow

4) Programs are placed in a fixed location in memory, can be non-dynamically executed initially.  
5) Programs are placed in a fixed location in memory, can be relative

6) Object code produced can be either (i) absolute  
7) relative

8) Stack based m/c operations are done by pushing operations  
9) This is Java Virtual Machine (JVM)

10) Address mode is complex from rest.

11) RISC m/c has few regs, 2-odd address, variety of address modes & simple from rest.

12) RISC m/c has many registers, 3-odd from, simple & stack based.

13) The main common target m/c architectures are RISC, CISC

14) The main set architecture of the target m/c has a. The main set architecture of the target m/c has a. The main set architecture of the target m/c has a. The main set architecture of the target m/c has a. The main set architecture of the target m/c has a.

## Instruction Selection

The Code Generator must map the IR to a code sequence that can be executed by the target m/c. The complexity of performing this mapping is determined by factors such as

- \* The level of IR
- \* The nature of instr set architecture.
- \* The desired quality of generated code.

The nature of the instr set of target m/c has strong impact. An instr selection eg:- uniformity & completeness of instr set are important factors.

Instr speed & m/c idiom are another factors.

eg:- Every 3-addr code of the form  $X = Y + Z$  are translated to code sequence like

LD R0, Y ; R0 = Y

ADD R0, R0, 2 ; R0 = R0 + 2

ST X, R0 ; X = R0

→ these instructions produce redundant LD & store. eg:

$a = b + c$   
 $d = a + c$

LD R0, b

ADD R0, R0, c

ST a, R0

LD R0, a

ADD R0, R0, c

ST d, R0

eg:-  $a = a + 1$  can be written as

LD R0, a

ADD R0, R0, 1

ST a, R0

this can be mapped using single instr INC

**GURUPRASAD. S.**  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-64,  
Ph: 9886760776

**Prof GURUPRASAD S, BMSIT&M**

<https://hemanthrajhemu.github.io>

A 16-bit problem in C, is deciding what value to hold in not registers. Registers are fastest compared with on memory. Also not held in reg. need to reside in memory. Program of reg are important. The use of reg is divided into 2 problems. Register allocation: Select the set of variables that will reside in registers at each point in program. Register assignments: Pick the specific register that a variable will reside.

3-operand code eg.

$$\begin{aligned}
 t &= a + b \\
 t &= t + c \\
 t &= t + d
 \end{aligned}$$

Only shift.  $\text{R}/\text{L} \text{ } \text{R} \text{ } \text{R}$  or error in second shift

Shorter + more easily code in.

L R, a  
 A R, b  
 M R, c  
 D R, d  
 ST R, t  
 L R, a  
 A R, b  
 A R, c  
 SRA R, 32  
 D R, d  
 ST R, t

Notes SRDA means shift right double arithmetic, SRDA R0, R2, Rn, R1, R0 & clear R0

(V) Evaluation order -

The order in which computations are performed can affect the efficiency of the target code. Same computation order require fewer registers to hold intermediate results than others.

The target Language

Familiarity with target m/c and its with set of prerequisites for designing a good code generator.

Our target computer models a 3 adar m/c with load and store opns, computation opns, jump opns & condnl jump opns.

The foll are the kind of instructions available:

(i) load opns - LD dest, addr LD R, n addr → dest

(ii) store opns - ST n, r  
r → n  
r3 to 1000

(iii) computation opns. OP dest, sr1, sr2.  
op - ADD, SUB

(iv) un condnl jump - BR L

(v) condnl jump - Bcond r, L eg. BLT2 r L  
rg  
Jump if r val is less than zero

Ex:-  $X = Y - 2$   
LD R1, Y  
LD R2, 2.  
SUB R1, R1, R2  
ST X, R1

$b = a[i]$  Each i is 8 bytes.  
LD R1, i  
MUL R1, R1, 8  
LD R2, a(R1)  
ST b, R2

**GURUPRASAD. S.**  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-64.  
Ph: 9886760776

Some common cost measure are the begin of compilation time  
 State, run time of code execution of program files  
 For something we have cost of space to be one file or  
 cost associated with the other mode of opndk.

Program & Cost

$R_1, R_2, R_3, R_4, R_5, R_6, R_7, R_8, R_9, R_{10}$

508  $R_1, R_2$

LD  $R_2, Y$

LD  $R_1, X$

LD  $R_2, Y$

ST  $(R_1), R_2$

LD  $R_2, Y$

LD  $R_1, R$

$X = Y$

ST  $R_2, R_2$

LD  $R_2, (R_2)$

LD  $R_1, R$

$X = R_1$

ST  $(R_1), R_1$

MUL  $R_2, R_2, R_2, R_2, R_2, R_2, R_2, R_2$

LD  $R_2, Y$

GURUPRASAD, S.  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64,  
 Ph: 9886760776

//  $cost_{total} = C_1 + cost_{code}(R_2) = R_1$



Eg:-

① LD R0, #1 - calc + 1 no mem opn

② LD R0, #1 calc + 2 one mem opn

③ LD R1, #100 (R2) Load R1 by (content(R2) + 100 + content(R2))

Calc R2 3 at 100 R2 + 100 R1 = word following R2

Eg:- ①  $x = b * c$

$y = a + x$

LD R0, b

LD R1, c

MUL R0, R0, R1

ST x, R0

LD R1, a

LD R2, x

ADD R1, R1, R2

ST y, R1

LD R1, x

LD R2, #5

MUL R2, R2, R1

ST b(R2), R1

$b[5] = x$

②  $x = a[i]$

$y = b[j]$

$a[i] = y$

$b[j] = x$

LD R0, i

MUL R0, R0, #8

ST x, a(R0)

$x = a[i]$

LD R1, j

MUL R1, R1, #8

ST y, b(R1)

$y = b[j]$

LD R0, y

LD R1, i

MUL R1, R1, #8

ST a(R1), R0

$a[i] = y$

③  $y = *q$

$q = q + 4$

$*p = y$

$p = p + 4$

LD R0, p

~~LD R0, p~~

ST y, 0(R0)

$y = *p$

LD R1, p

LD R2, y

~~ST R2, 0(R1)~~

ST 0(R1), R2

$*p = y$

LD R0, q

ADD R0, R0, #4

ST q, R0

$q = q + 4$

LD R0, y

ADD R0, R0, #4

ST p, R0

$p = p + 4$

can it return.

Records are created & destroyed during procedure  
A Dynamically managed area: Stack for holding activation

A Dynamically managed area: Heap for holding data objects that  
are created allocated & freed during program

Global constants & other data generated by compiler.

Statically determined data areas: Data area for holding

Executable target code.

Statically determined area: Code that holds the

code that holds the code area and data area.

The executable program runs in its own logical address space that

is address in target code.

Here we know how names in the IR can be converted

addresses to the target code

```
LD R0, 0
LD R1, 0
SUB R1, R1, R0
BLT2 R1, L2
```

```
!P? > n goto L2
```

```
LD R0, 0
LD R1, 0
SUB R0, R0, R1
BLT2 R0, L1
ST 2, 0
BR L2
L1 ST 2, 1
```

GURUPRASAD, S.  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-84,  
Ph: 9886760776

## Static Allocation

Code gen for simplified Proc call & return focus on the following 2-addr Stmt.

(i) call callee (ii) return (iii) halt (iv) action.

Ex:-

code for C	100: Action 1
action 1	120: ST 364, #140 (store ret addr i.e. 140 in locn 364)
call P	132: BR 200 (addr of P)
action 2	140: ACTION 2
halt	160: HALT
	⋮
code for P	200: ACTION 3
action 3	220: BR *364
return	

## Stacks Allocation

Static allocation can become stack allocation by using relative addresses for storage in activation records. ~~However stack allocation~~

However in stacks allocn the pos<sup>n</sup> of an act<sup>n</sup> record for a procedure is not known until run time.

When a procedure is called the calling procedure increments the stack pointer & transfers the control to the called procedure. After control returns to the caller, we decrement stack pointer, thereby deallocating the activation record of called procedure.

Ex:- LD SP, # Stack Start

Code for the first procedure

HALT

A procedure call sequence increments stack pointer, saves the return addr & transfers control to called procedure

ST    4 SP, # here + 16  
 save return addr

BP    callee . code area.

the called procedure transfer control to ret addr using

BR \* 0 (SP)

and decrement SP by

506 SP, SP, # caller . retaddr.

Model form

action 1

Call

action 2

halt

Code ptr #

action 3

return



100 : LD SP, # 600    #initial stack

108 : Action 1

128 : ADD SP, SP, # 8    # increase

136 : ST    # SP, # 152

144 : BR 200

152 : SUB SP, SP, # 8    # increase

160 : Action 2

180 : HALT

200 : Action 3

220 : BR \* 0 (SP)

## Basic Blocks & Flow Graphs

Basic blocks are the graphical representation of IR code that is helpful for code generation. which is constructed as follows:

(!) Start after the IR code. into basic blocks which are

regions of consecutive 3-addr IR with the properties that

① The flow of control can only enter the B.B through the

that exists in the blocks.

⑤ Control will leave the blocks without halting or branching, except in the last instruction in the blocks.

⑥ The basic blocks become the nodes of a flow graph whose edges indicate with blocks can follow with other blocks.

### Basic Blocks

The first job is to partition a sequence 3-addr instr into basic blocks

Algorithm: Partitioning 3-addr instr into B.B

Input: A sequence of 3-addr instr

Output: A list of B.B for the sequence in which each instr is assigned to exactly one B.B

Method:

Step ① - determine the set of Leaders, the first instr of B.B with the following rules.

① The first instr is a Leader

② The target of a conditional / un-conditional goto is a leader.

③ The instr which immediately follows a condn / uncondn goto is a Leader.

Step ② - For each Leader construct the B.B which consists of a Leader & all the instr up to next Leader or end of the program.

Note: Any instr not placed in a B.B can never be executed & may be removed.

- 1)  $prod = 0$   $L_1$
- 2)  $i = 1$
- 3)  $t_1 = i * 4$   $L_2$
- 4)  $t_2 = 9 [t_1]$
- 5)  $t_3 = ~~prod~~ * 4$
- 6)  $t_4 = 6 [t_3]$
- 7)  $t_5 = t_2 * t_4$
- 8)  $t_6 = prod + t_5$
- 9)  $prod = t_6$

```

10)  $t_7 = i + 1$ 
11)  $i = t_7$ 
12)  $i < 20$  goto 3
13) end

```



```

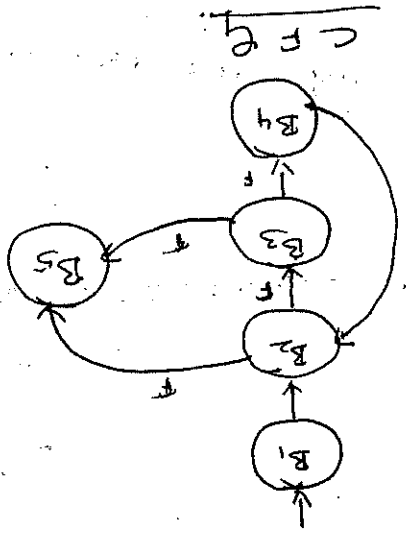
1)  $prod = 0$   $L_1$ 
2)  $i = 1$ 

```

while ( $i < 20$ );  
 $i = i + 1$

do  
 $prod = prod + a [i] * b [i]$   
 $i = i + 1$

GURUPRASAD S.  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64,  
 Ph: 9886760776



$i = 1$   
 $i = 10$   
 $a = b + c + d$   
 $i = i + 1$   
 $i = i - 1$

```

1)  $i = 1$ 
2)  $i = 10$ 
3)  $i > 5$  goto 13
4)  $i < 6$  goto 15
5)  $t_1 = b + c$ 
6)  $t_2 = t_1 + d$ 
7)  $a = t_2$ 
8)  $t_3 = i + 1$ 
9)  $i = t_3$ 
10)  $t_4 = i - 1$ 
11)  $i = t_4$ 
12)  $i + 3$ 
13) end

```

3 order code:

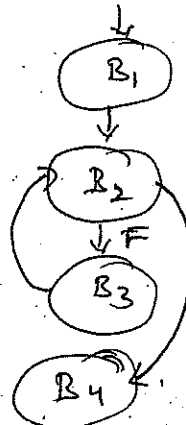
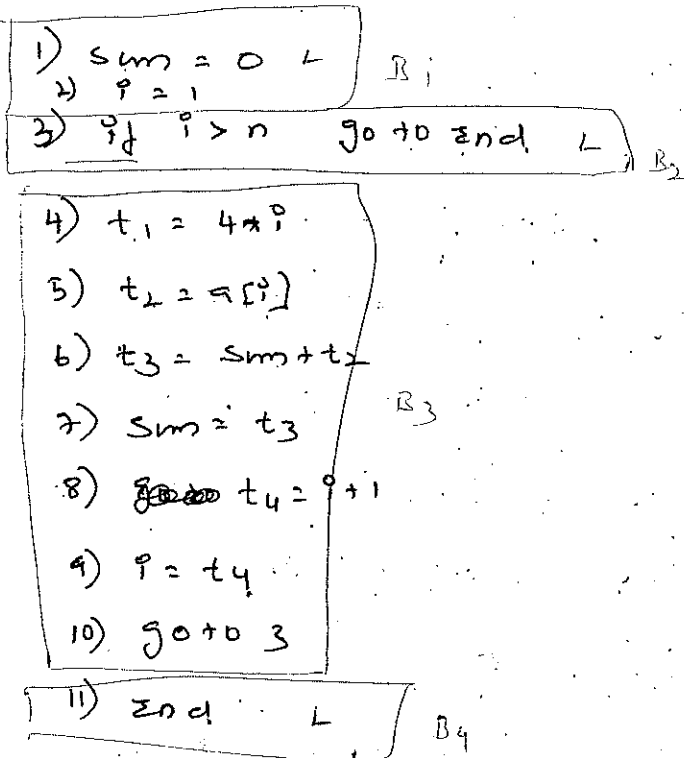
Prof GURUPRASAD S, BMSIT & M

eg:- 1) Generate the intermediate code for the following construct  
<https://hemanthrajhemu.github.io>

Sum = 0

for  $i = 1$  to  $n$  do

Sum = Sum +  $a[i]$



**GURUPRASAD. S.**  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-64,  
Ph: 9886760776

### Next Use Information

Knowing when the value of a variable will be used next is essential for generating good code.

Algorithm to determine the liveness and next use info for each statement in B.B

**Input:** A basic Block of 3-addr stmt with S.T initially live for all non temp variables in B.

**Output:** At each stmt  $i: x = y + z$  in B, we attach the liveness and next use info of  $x, y, z$ .

**Methods:** Start at the last stmt in B and scan back wards to the beginning of B. At each stmt  $i: x = y + z$  in B do the following

- (i) attach the stmt  $i$  the info currently found in the SymbolTable regarding the next use & liveness of  $x, y, z$ .

(ii) In the Symbol Table, set  $Y_3$  to "love" & next set

of  $Y_3$  to 9.

### Flow Graphs

Once an intermediate code is transformed to B.B. we can then control flow by a flow graph or control flow graph.

The nodes of CFG are the B.B.

There is an edge from block B to block C iff it is possible for the flow graph to the blocks to form a path from the last block to block C.

There are 2 ways each edge could be inserted

DEPT OF CSE BMSIT BANGALORE

\* Store related condition/in condition jump from the end of B to C

Labels of C

\* C immediately follows B in the original order of B-order flow

and B doesn't end in an uncondition jump

often we add too nodes, called entry & exit but don't correspond to executable intermediate form

Ex:

for  $i = 1$  to 5 do

for  $j = 1$  to 5 do

if  $i < j$  then

for  $k = 1$  to 5 do

if  $k < i$  then

for  $i = 1$  to 5 do  
for  $j = 1$  to 5 do  
if  $i < j$  then



matrix is stored in row major order. so  $a[i][j] =$

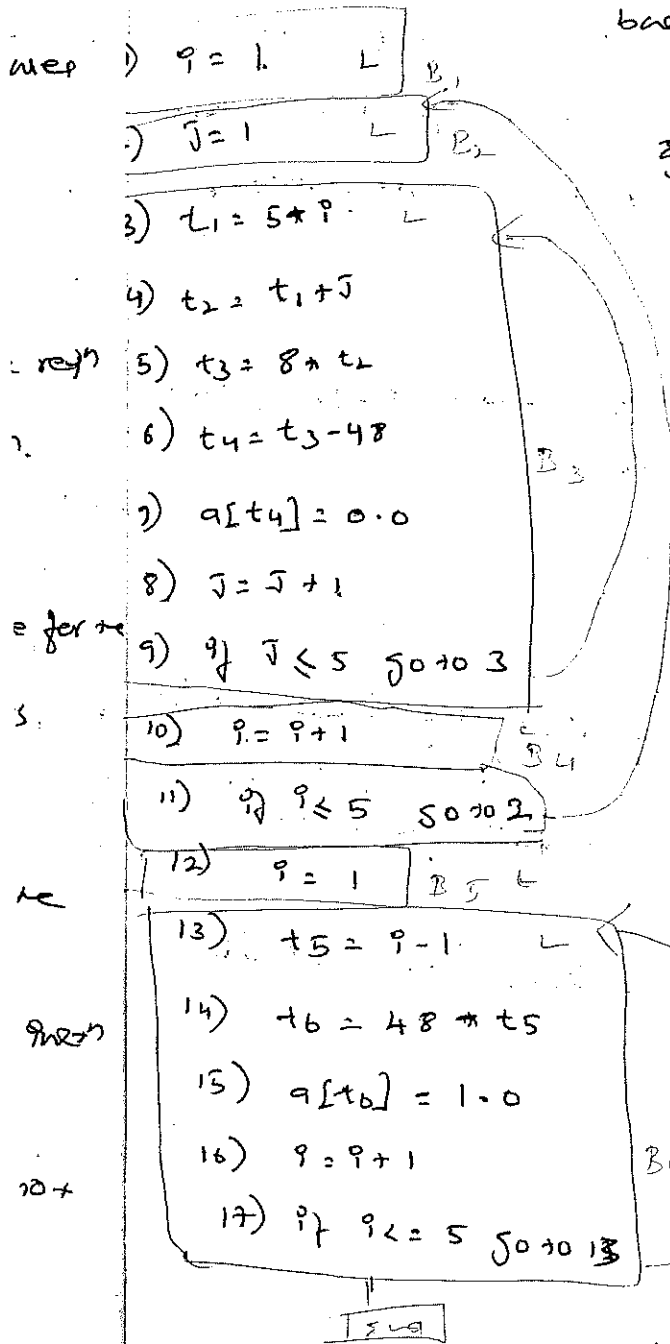
$$\text{base} + \left( \left( \frac{i}{5} * 9 \right) + j \right) * 8 - 48$$

no. of rows

Eg:-  $a[2,2] = 200 + \left( \left( \frac{5+2}{5} + 2 \right) * 8 \right) - 48$   
 $= 248$

$$\text{base} + (i_1 * n_2 + i_2) * w$$

we'll



**GURUPRASAD. S.**  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64,  
 Ph: 9886760776

Loops:

For lang constructs like while stmt do-while stmt & for stmt naturally give rise to loops many code transformations depend upon the pattern of loops in a flow graph.

A set of nodes  $L$  in a flow graph is a loop if -

DEPT OF CSE BMSIT, BANGALORE-64  
 ASSISTANT PROFESSOR  
 GURUPRASAD S.

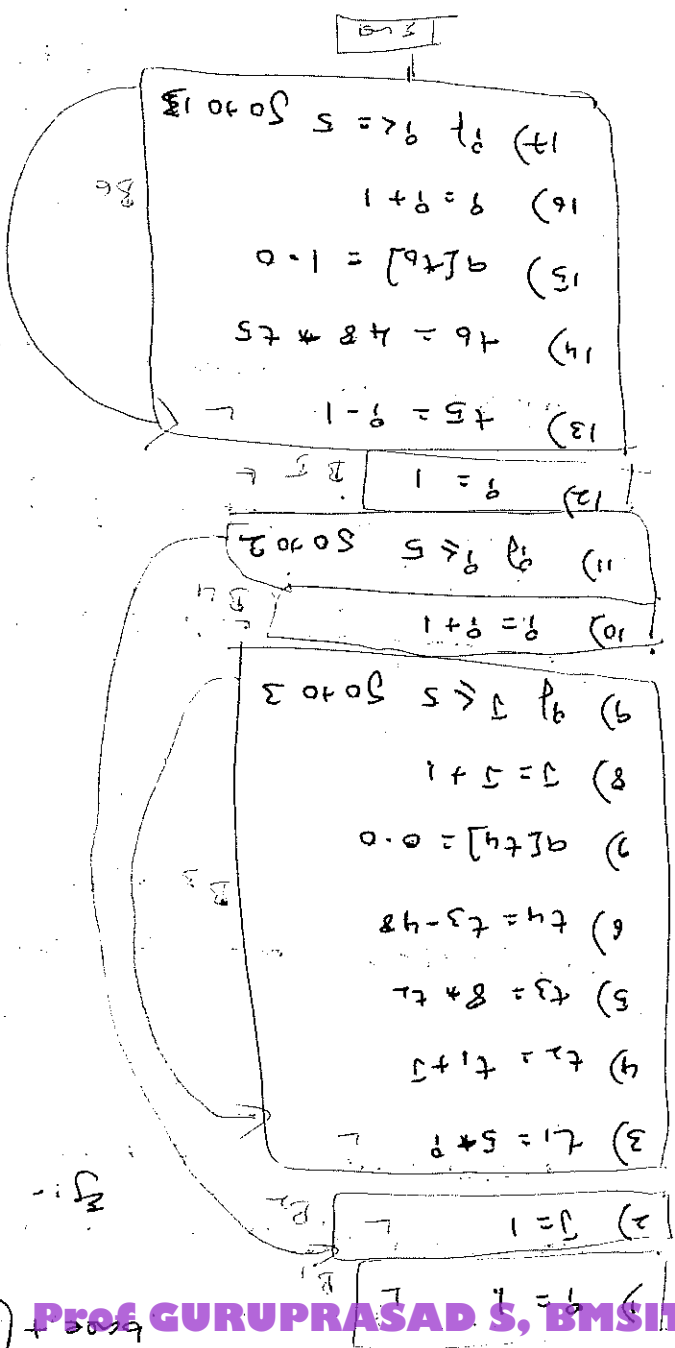
DEPT OF CSE BMSIT, BANGALORE-64  
 PH: 9886760776

A set of nodes L in a flow graph G is

the path of loops in a flow graph.

For loops construct flow while stmt do-while stmt if-else stmt normally give rise to loops many code transformation depend upon

Loops:



GURUPRASAD S.  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64  
 Ph: 9886760776

$$\text{base} + (i * n_2 + j * n_3)$$

$$= 248 = 200 + ((5 + 2) * 2) * 8 - 48$$

(9) There is a node in  $L$  called the loop entry with the property that no other node in  $L$  has a greater outdegree than  $L$  i.e. Every path from the entry of the loop goes to any node in  $L$  goes to the loop entry.

(10) Every node in  $L$  has a non-empty path completely within  $L$  to the entry of  $L$  i.e. self loop.

### Matrix multiplication

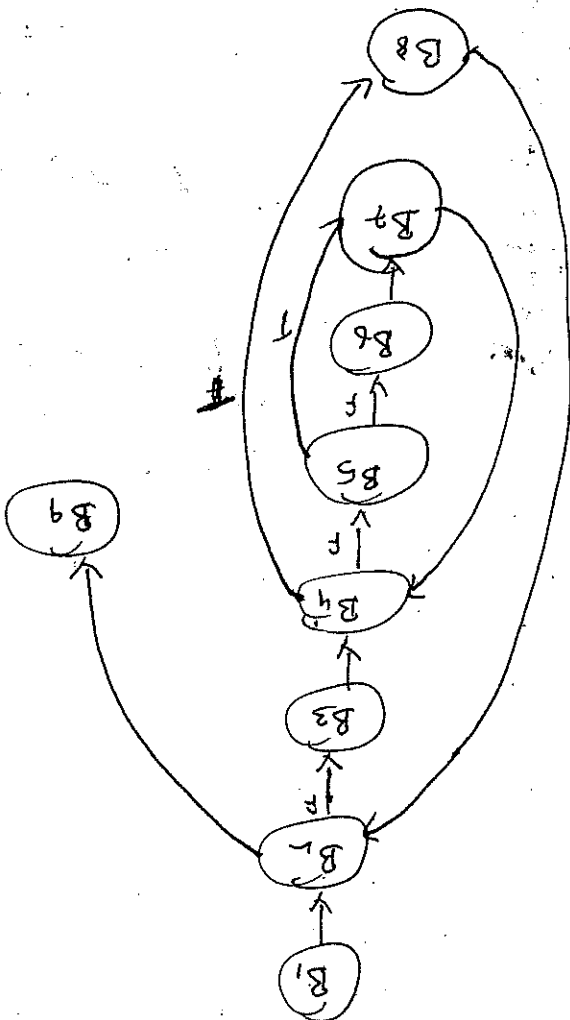
Bubble Sort for ascending order

```

for (i=1; i<n; i++) // To keep track of pass no's
{
    for (j=0; j<n-i; j++) // To compare the elements
    {
        if (a[j] > a[j+1]) // checking the condition
        {
            temp = a[j];
            a[j] = a[j+1];
            a[j+1] = temp;
        }
    }
}

```

GURUPRASAD. S.  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64.  
 Ph: 9886760776



CFG for sort

Param 0  
 Param 1  
 Param n  
 3) call sort  
 4) i = 1  
 5) j = 0  
 L6) ~~goto B1~~ (i > n) goto B4  
 L7) ~~return i~~ i = n - 1  
 L8) if (j > i) goto B2  
 L9) i2 = i + j  
 10) t3 = a[t2]  
 11) t4 = j + 1  
 12) t5 = i + t4  
 13) t6 = a[t5]  
 14) if (t3 < t6) goto B8  
 L15) temp = t3  
 16) t3 = t6  
 17) t6 = temp  
 L18) t7 = j + 1  
 19) j = t7  
 20) goto B8  
 L20) t8 = i + 1  
 21) i = t8  
 22) goto B6  
 L23) return



A substantial improvement in the running time of code can be achieved by local optimization in each B.B. and more can be achieved in region Global optimization which looks for info thru b/p blocks.

### The DAG repn of B.B.

Many local optimizations begin by transforming B.B. into DAG. We can construct DAG for expr in each B.B. as follows:

(1) There is a node in DAG for each of basic value of variable in the blocks (leaf node)

(2) There is node  $N$  (non-leaf) associated with each stmt  $S$  in the children of  $N$  corresponds to statement. That are leaf defn prior to  $S$  of the opnd used by  $S$ .

(3) Node  $N$  is labeled by operator labeled by  $S$  and applied to also the leaf var for constant & leaf defn within blocks

(4) Certain nodes are designated as  $\phi$  nodes whose variables are live on exit from blocks

The DAG representation of B.B. lets us perform several code improvements transformation on code repn by the blocks

(i) Elimination of local common sub expression i.e. the terms that compute a value that has already been computed.

(ii) Eliminate dead code i.e. terms that compute a value never needed.

(iii) ~~Reduce~~ Re-order same that do not depend on one another

(iv) Apply algebraic laws to re-order opnd & simplify computation.

# Finding local Common Sub-Expressions

Common Sub-Exprn is noticing a new node  $M$  is about to be added, whether there is an existing node  $N$  with same children in the same order and with same operator.

If so  $N$  may be used in place of  $M$ .

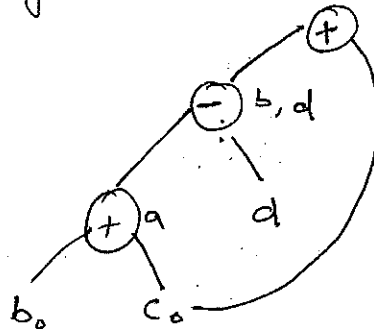
Ex:-

$$a = b + c$$

$$b = a - d$$

$$c = b + c$$

$$d = a - d$$



Since there are only 3 non leaf nodes we can represent all same of B.B as

$$a = b + c$$

$$b = a - d$$

$$\frac{c = d + c}{\text{Correct}} \quad \text{or} \quad \frac{c = b + c}{X}$$

as  $b$  is not lve on exit so we use node  $d$  rather than  $b$

Ex 2:-

$$a = b + c \quad \text{--- (1)}$$

$$b = b - d$$

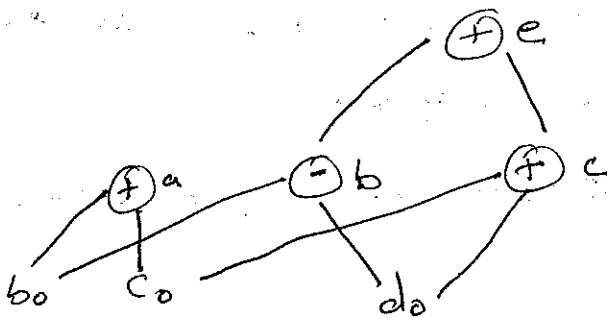
$$c = c + d$$

$$e = b + c \quad \text{--- (4)}$$

The values of  $b$  &  $c$  change

before it is used in 4<sup>th</sup>

Start to new node



GURUPRASAD S.  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64  
 Ph: 9886760776

$t_1 = 4 + 9$   
 $t_2 = 9 [t_1]$   
 $t_3 = t_2 + t_5$   
 $t_4 = 6 [t_1]$   
 $t_5 = t_2 + t_4$   
 $t_6 = t_3$   
 $Prod = t_6$

$t_1 = 4 + 9$   
 $t_2 = 9 [t_1]$   
 $t_3 = t_2 + t_5$   
 $t_4 = 6 [t_1]$   
 $t_5 = t_2 + t_4$   
 $t_6 = t_3$   
 $Prod = t_6$

$t_4 = 6 [t_1]$

$t_2 = 9 [t_1]$

$t_1 = 4 + 9$



$t_1 = 4 + 9$   
 $t_2 = 9 [t_1]$   
 $t_3 = t_2 + t_5$   
 $t_4 = 6 [t_1]$   
 $t_5 = t_2 + t_4$   
 $t_6 = t_3$   
 $Prod = t_6$

$t_1 = 4 + 9$

$Prod = t_6$

$t_6 = t_3$

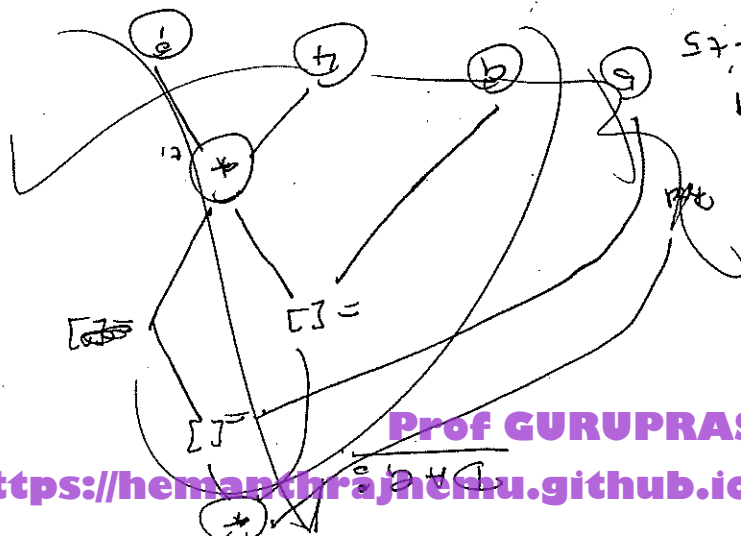
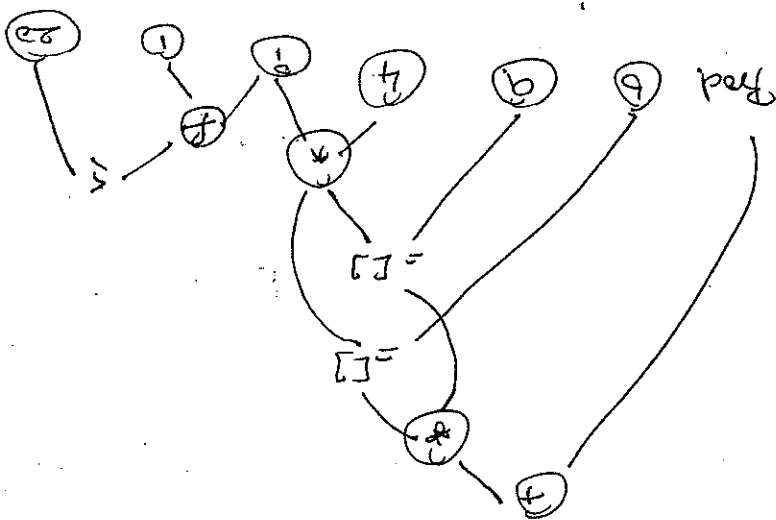
$t_5 = t_2 + t_4$

$t_4 = 6 [t_1]$

$t_3 = t_2 + t_5$

$t_2 = 9 [t_1]$

$t_1 = 4 + 9$





# Dead Code Elimination

The op<sup>n</sup> of DAG, corresponding to dead code elim<sup>n</sup> can be imp<sup>d</sup> as follows:

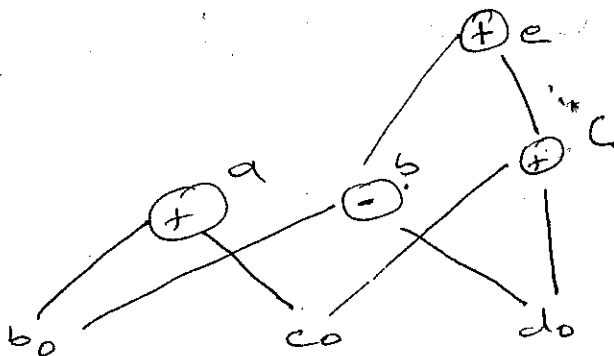
- ① Delete from DAG any root that has no live variables attached, a root node is one with no ancestor.
- ② Repeated application of this transformation will remove all nodes from the DAG that correspond to dead code.

eg:  $a = b + c$

$b = b - d$

$c = c + d$

$e = b + c$



I

GURUPRASAD. S.  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-56  
Ph: 9886760776

$a, b$  are live but  $c$  and  $e$  are not,

we can remove  $e$  and  $c$  becomes root as  $a$  is also not live we can eliminate it.

to solve:

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

**GURUPRASAD. S.**  
 ASSISTANT PROFESSOR,  
 DEPT OF CSE BMSIT, BANGALORE-56,  
 Ph: 9886760776

① Algebraic Identities are very important class of math on basic blocks we may apply algebraic identities such as.

$$x + 0 = 0 + x = x$$

$$x - 0 = x$$

$$x * 1 = 1 * x = x$$

$$x / 1 = x$$

② Another class of algebraic identities include local reduction in strength i.e replacing a more expensive operator by cheaper one as.

$\frac{x^2}{x^2}$ $\frac{x * x}{x * x}$ $\frac{x * 0.5}{x * 0.5}$	$\frac{x^2}{x^2}$ $\frac{x * x}{x * x}$ $\frac{x / 2}{x / 2}$
Cheaper	Expensive

③ Third class of related operators. Constant folding, here we evaluate constant exprn at compile time. I replace the exprn by its value.

$$x_i - 2 * 3.14 \quad \text{re replaced by value } 6.28$$

④ Data constant can help to apply another algebraic transformation such as commutativity and associativity.

⑤ Commutative  

$$x * y = y * x$$
  
 If exprn  $y * x$  has already present then can be read forward of creating node for  $x * y$ .

⑥  $x \geq y$  can be reversed by  $x - y$

also associative

$$x_i - a = b + c$$

$$x_i + d + e = c + d + e$$

$$\Rightarrow t = c + d$$

$$a = b + c$$

$$e = t + e$$

$$q = b + c$$

$$e = \overline{q + d}$$

etc

$Prod = 0$   
 $i = 0$   
 Eg:-  $T_1 = 4 * i$   
 $T_2 = a[T_1]$   
 $T_4 = b[T_1]$   
 $T_5 = T_2 * T_4$   
 $Prod = Prod + T_5$   
 $i = i + 1$   
 $i \leq 20$  goto 3



$Prod = 0$   
 $T_1 = 0$   
 $T_1 = T_1 + 4$   
 $T_2 = a[T_1]$   
 $T_4 = b[T_1]$   
 $T_5 = T_2 * T_4$   
 $Prod = Prod + T_5$   
 $T_1 = T_1 + 1$   
 $i \} T_1 \leq 80$  goto B2

$T_1 = 4 * i$  Increment  $i$  value by 4 Every time.  
 So we can replace by cheaper expr'n  $T_1 + 4$   
 so  $i \leq 20$  is written as  $T_1 \leq 80$   $\therefore (20 * 4)$

## Representation of Array Reference

to solve: $x = a + b + c + d + e + f$ $y = a + c + e$
---

The proper way to repr array access in DAG is as follows.

- (i) An assignment of form an array  $x = a[i]$ , is repr by creating a node with operator  $= [ ]$  & two children representing the initial value of array  $a$  & index  $i$ . Variable  $x$  becomes label of this new node.
- (ii) An assignment of form  $a[i] = y$  is repr by new node  $[ ] =$  with children repr  $a$  &  $i$  &  $y$ . ~~there is no label~~ there is no label. this node will kill all nodes currently created.

Eg:-  $x = a[i]$   
 $a[5] = y$   
 $z = a[i]$

After we perform whatever operation possible while constructing DAG or by manipulating DAG. We must re-construct the 3-addr code for B.B from which we built DAG.

Re-assembling B.B from DAG

The =\* will fill all other nodes so for connected to DAG. associated with 'addr' or arguments.

the operator =\* must reuse all nodes that are currently we don't know what 'p' or 'q' point to - so

When we assign indirectly through a 'p' or 'q'  $x = *p$  or  $*q = x$

Pointer assignment & Procedure Call

GURUPRASAD, S.  
ASSISTANT PROFESSOR  
DEPT OF CSE BMSIT, BANGALORE-64.  
Ph: 9886760776

to solve:  
 $a[p] = 6$   
 $*q = c$   
 $d = a[p]$   
 $e = *q$   
 $*p = a[p]$

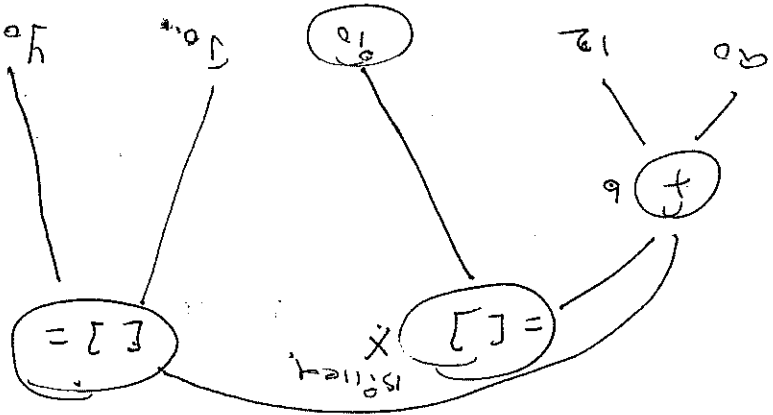
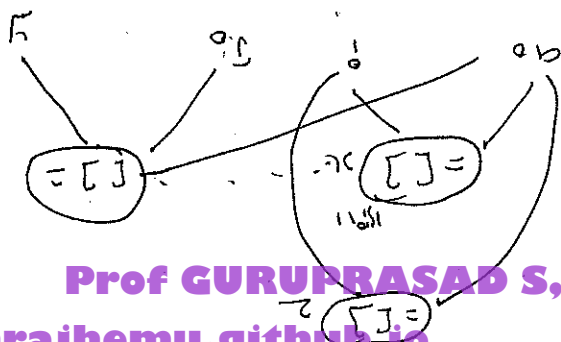


Fig 2:  
 $b = 12 + q$   
 $x = b[p]$   
 $b[p] = y$



Node x is visited by  $a[p] = y$



One of the primary reason during code gen is deciding how to use registers to best advantage.

There are four registers for used registers

(!) In most arch, some or all of ops of an op must be in registers in order to perform ops.

(00) Reg reuse Good temporary

(01) Reg are used to hold (global) values that are compiled in one basic block & used in other block.

(02) Reg are often used to help the run-time average reg #.

The o/c form are of the form

LD reg, mem

ST mem, reg

OP reg, reg, reg

Register & address decoders

The two decoder structures used are used in gen of code arcs

(!) Address Decoders (AD):- keeps track of the location

the current value of the variable can be found. the locn

may be reg, mem or stack or combination. the info can be stored

in ST entry for that variable name.

**GURUPRASAD S.**  
 ASSISTANT PROFESSOR  
 DEPT OF CSE BMSIT, BANGALORE-64,  
 Ph: 9886760776

(10) Register Descriptor - Keeps track of the variable name, whose current value is in the register, we assume that initially, all register descriptors are empty. As the codegen program runs each reg will hold the value of zero or more names.

## The Code-Generation Algorithm

An essential part of alg is function  $getReg(I)$ , which selects reg for each mem loc<sup>n</sup> associated with 3-addr code of instr<sup>n</sup>  $I$ .  
 $getReg(I)$  has access to all the Regs and Address Descriptors of all variables in B.B.

M/c instr<sup>n</sup> for op<sup>n</sup>:

For a 3-addr instr<sup>n</sup> such as  $x = y + z$ , do the following

(i) use  $getReg(x = y + z)$  to select reg for  $x, y$  and  $z$  call them

$R_x, R_y, R_z$

(ii) If  $y$  is not in  $R_y$  then issue an instr<sup>n</sup> LD  $R_y, y$  where  $y$  is mem loc<sup>n</sup> for  $y$

(iii) If  $z$  is not in  $R_z$  then issue LD  $R_z, z$  where  $z$  is mem loc<sup>n</sup> for  $z$

(iv) Issue instr<sup>n</sup> ADD  $R_x, R_y, R_z$

M/c instr<sup>n</sup> for Copy stmt

For stmt of the form  $x = y$   $getReg(x)$  will use same reg for both  $x$  &  $y$

if  $y$  is not in reg generate LD  $R_y, y$  if already in  $R_y$  do nothing.

If the var is used temporarily only within blocks then at the end of block value is forgot & reg is empty.

If var is live on exit or we don't know about live here, then we assume variable is needed and generate ST x, r.

Manages Register & Address Descriptor.

As the code gen his never load store & other instructions it needs to update the register & address descriptor. The rules are as follows:

1) For the first LD R, x.

a) Change RD for reg R to hold only x.

b) Change AD for x by adding reg R as address.

2) For the first ST x, r.

Change AD of r to indicate its own mem locn

3) For an op even as ADD R<sub>1</sub>, R<sub>2</sub>, R<sub>3</sub> of x = y + z

a) Change RD for R<sub>2</sub> to hold z only

b) Change AD for r so that it only locn R<sub>2</sub>

c) Remove R<sub>2</sub> from AD if any var other than z

4) When we process copy stmt x = y with load y and reg y

a) Add y to RD for R<sub>2</sub>

b) Change AD for x so that it only locn R<sub>2</sub>



Ex:-  
 $t = a - b$   
 $u = a - c$   
 $v = t + u$   
 $a = d$   
 $d = v + u$

$t_1 = a - b$   
 $t_2 = a - c$   
 $t_3 = t_1 + t_2$   
 $a = d$   
 $d = t_3 + t_2$

$R_1, R_2, R_3$

define:  $t, u, v$  are temp and local to block

~~local to block~~ ~~are~~ ~~a, b, c, d~~

$a, b, c, d$  are live on exit from block

stmt	Codegen	R D	AD																																		
$t = a - b$	LD $R_1, a$ LD $R_2, b$ SUB $R_2, R_1, R_2$	$R_1 \quad R_2 \quad R_3$ <table border="1"> <tr><td>a</td><td>b</td><td></td></tr> <tr><td>a</td><td>t</td><td></td></tr> </table>	a	b		a	t		<table border="1"> <tr><td>a</td><td>b</td><td>c</td><td>d</td><td>t</td><td>u</td><td>v</td></tr> <tr><td>a</td><td>b</td><td>c</td><td>d</td><td></td><td></td><td></td></tr> <tr><td>a</td><td>b</td><td>c</td><td>d</td><td>t</td><td>u</td><td>v</td></tr> <tr><td>a</td><td>b</td><td>c</td><td>d</td><td><math>R_2</math></td><td></td><td></td></tr> </table>	a	b	c	d	t	u	v	a	b	c	d				a	b	c	d	t	u	v	a	b	c	d	$R_2$		
a	b																																				
a	t																																				
a	b	c	d	t	u	v																															
a	b	c	d																																		
a	b	c	d	t	u	v																															
a	b	c	d	$R_2$																																	
$u = a - c$	LD $R_3, c$ SUB $R_1, R_1, R_3$	$R_1 \quad R_2 \quad R_3$ <table border="1"> <tr><td>u</td><td>t</td><td>c</td></tr> </table>	u	t	c	<table border="1"> <tr><td>a</td><td>b</td><td>c</td><td>d</td><td>t</td><td>u</td><td></td></tr> <tr><td>a</td><td>b</td><td><math>R_3</math></td><td>d</td><td><math>R_2</math></td><td><math>R_1</math></td><td></td></tr> </table>	a	b	c	d	t	u		a	b	$R_3$	d	$R_2$	$R_1$																		
u	t	c																																			
a	b	c	d	t	u																																
a	b	$R_3$	d	$R_2$	$R_1$																																
$v = t + u$	ADD $R_3, R_2, R_1$	$R_1 \quad R_2 \quad R_3$ <table border="1"> <tr><td>u</td><td>t</td><td>v</td></tr> </table>	u	t	v	<table border="1"> <tr><td></td><td></td><td></td><td></td><td>t</td><td>u</td><td>v</td></tr> <tr><td>a</td><td>b</td><td>c</td><td>d</td><td><math>R_2</math></td><td><math>R_1</math></td><td><math>R_3</math></td></tr> </table>					t	u	v	a	b	c	d	$R_2$	$R_1$	$R_3$																	
u	t	v																																			
				t	u	v																															
a	b	c	d	$R_2$	$R_1$	$R_3$																															
$a = d$	LD $R_2, d$	$R_1 \quad R_2 \quad R_3$ <table border="1"> <tr><td>u</td><td>a, d</td><td>v</td></tr> </table>	u	a, d	v	<table border="1"> <tr><td></td><td></td><td></td><td></td><td>t</td><td>u</td><td>v</td></tr> <tr><td>a</td><td>b</td><td>c</td><td>d <math>R_2</math></td><td></td><td><math>R_1</math></td><td><math>R_3</math></td></tr> </table>					t	u	v	a	b	c	d $R_2$		$R_1$	$R_3$																	
u	a, d	v																																			
				t	u	v																															
a	b	c	d $R_2$		$R_1$	$R_3$																															
$d = v + u$	ADD $R_1, R_1, R_3$	$R_1 \quad R_2 \quad R_3$ <table border="1"> <tr><td>d</td><td>a</td><td>v</td></tr> </table>	d	a	v	<table border="1"> <tr><td>a</td><td>b</td><td>c</td><td>d</td><td>t</td><td>u</td><td>v</td></tr> <tr><td>a</td><td>b</td><td>c</td><td><math>R_2</math></td><td></td><td></td><td><math>R_3</math></td></tr> </table>	a	b	c	d	t	u	v	a	b	c	$R_2$			$R_3$																	
d	a	v																																			
a	b	c	d	t	u	v																															
a	b	c	$R_2$			$R_3$																															
exit	ST $a, R_2$ ST $d, R_1$	<table border="1"> <tr><td>d</td><td>a</td><td>v</td></tr> </table>	d	a	v	<table border="1"> <tr><td>a</td><td>b</td><td>c</td><td>d</td><td>t</td><td>u</td><td>v</td></tr> <tr><td><math>R_2</math></td><td>b</td><td>c</td><td>d <math>R_1</math></td><td></td><td></td><td><math>R_3</math></td></tr> </table>	a	b	c	d	t	u	v	$R_2$	b	c	d $R_1$			$R_3$																	
d	a	v																																			
a	b	c	d	t	u	v																															
$R_2$	b	c	d $R_1$			$R_3$																															

For prefix Rn Rn by by the value are

(1) If y is currently in reg, then it is a reg already containing y as Rn, don't have to load to Rn

(2) If y is not in reg, but there is reg currently empty

then one such reg as Rn

(3) The different case if y is not in Reg if there are no empty reg then we should pick one of the available reg and make it reg to reuse.

Let R be a candidate reg & V is another var. that

the R is R say Rn then possibilities are:

(a) If the AD for V says that V is same where beside R,

then use one R.

(b) If V is R the value being computed by Rn is Rn

not use one of the other operands of Rn then we are OK

(c) otherwise if V is not used later i.e. after Rn

and if V is like an arg from the basic then we are OK

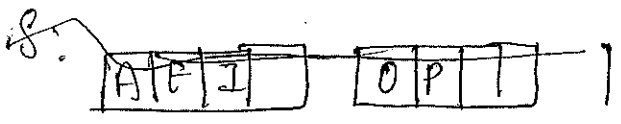
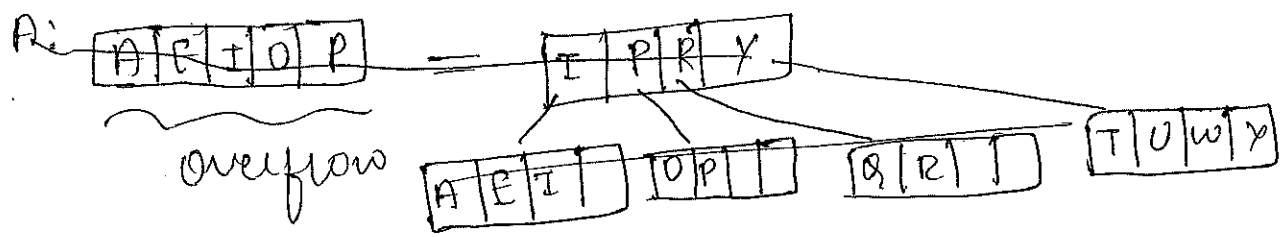
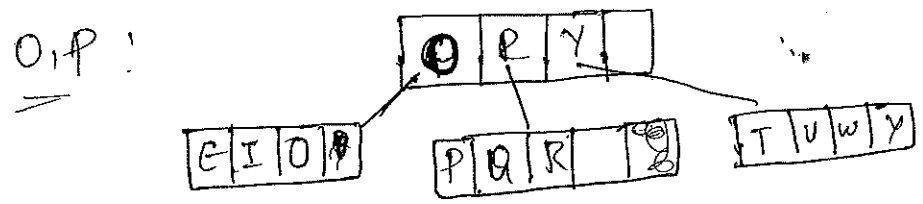
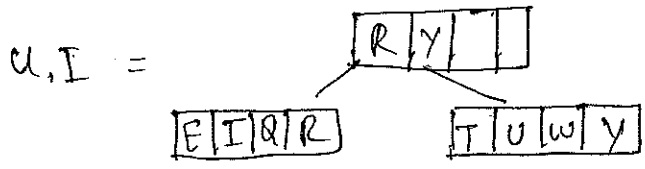
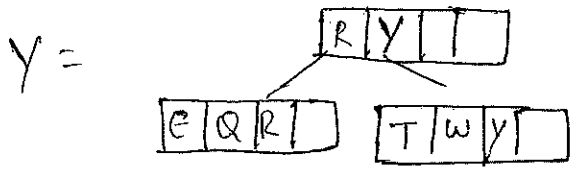
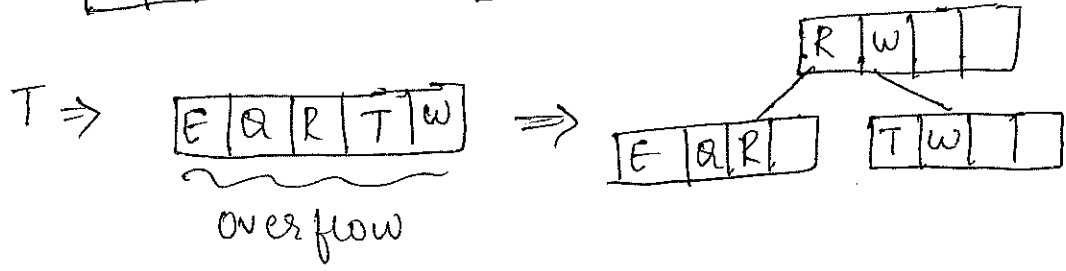
(d) If we are not OK by one of the last two cases then

we need to generate the store Rn ST V, R place a copy of V in the main memory that is called spill.

End of unit 8

Q W E R T Y U I O P A S D F G H J K L  
 Z X C V B N M

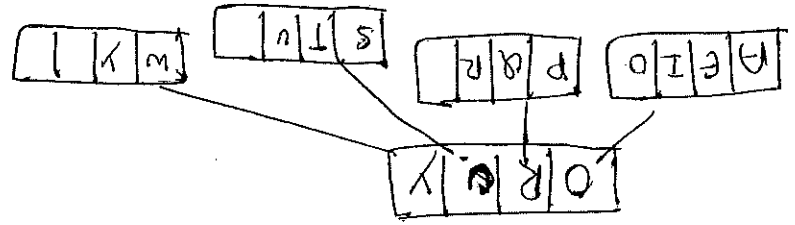
Q W E R ⇒ E Q R W



S: O R Y

A F I O P A R S T U W Y

question



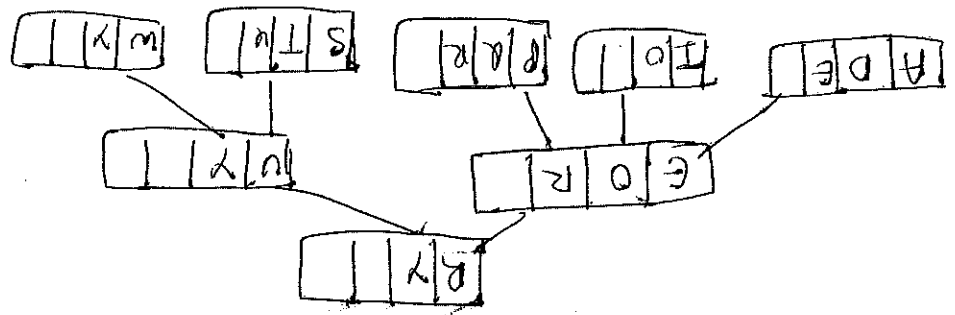
D:

A F I O P A R S T U W Y

question

O R Y

A F I O P A R S T U W Y



P.g

A F I O

P A R

S T U

W Y

#

A F I O

P A R

S T U

W Y

J: =>

O R Y

A F I O

P A R

S T U

W Y

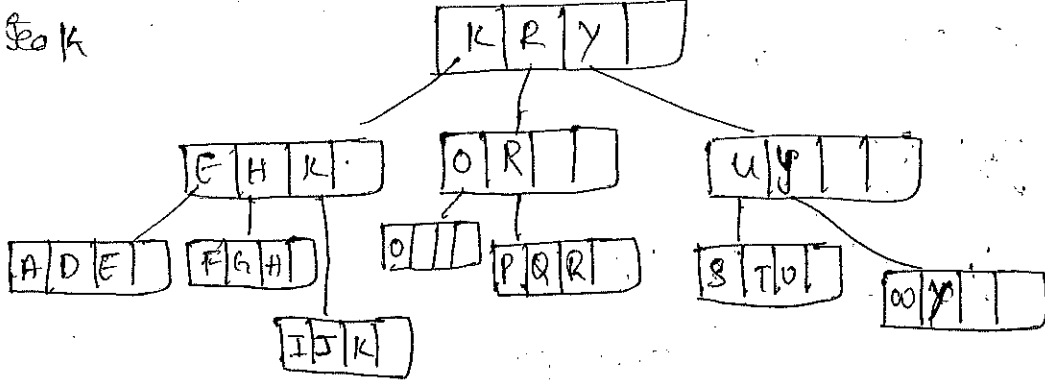
I O

P A R

S T U

W Y

Seq K



L :

\* All opt'n depend on data flow analysis.

\* DFA is a technique to derive info about flow data.

along  $\text{prog exn path}$ .

\* Eg:-  $X = a + b$ ,  $X = a + b$  are both  $X, a, b$  are having same value or different - determined by DFA.

### DFA as summary/Abstraction

\*  $\text{prog exn}$  is series of transformations, each start transform

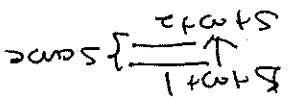
o/p state to new o/p state.

\* o/p state -  $\text{prog pt}$  before - state and o/p state -  $\text{prog pt}$  after

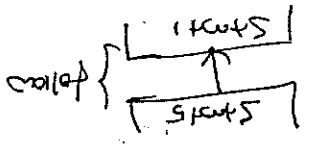
\* To analyze  $\text{prog}$  behavior. we must consider all possible  $\text{prog}$

paths

\* When in a basic block - the  $\text{prog pt}$  after a state is same as the  $\text{prog pt}$  before next state



\* Between  $B_1$  and  $B_2$  the  $\text{prog pt}$  of last state of  $B_1$  is immediately followed by first state of  $B_2$



the execution path (path)  $P_1$  to  $P_n$  is a sequence of

paths  $P_1, P_2, \dots, P_n$  such that

$P_i$  precedes  $P_{i+1}$

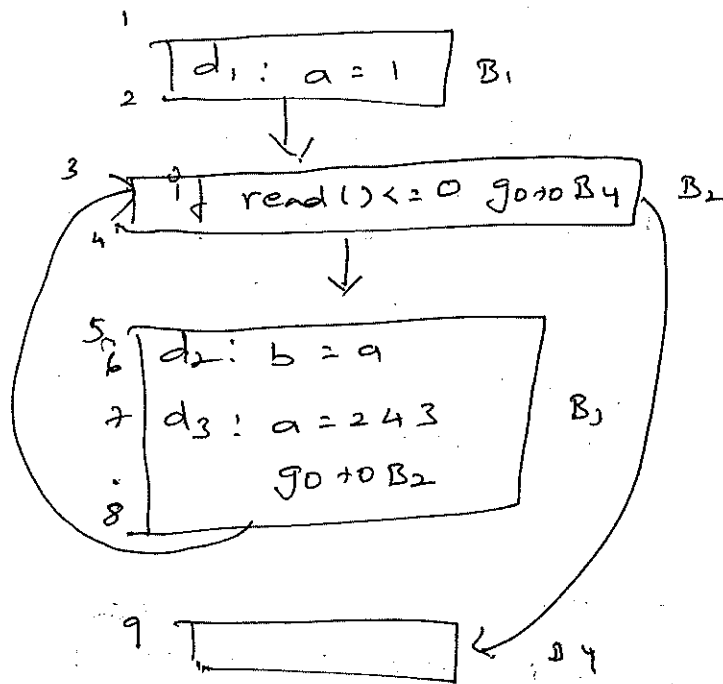
$P_i$  is the end of some block and  $P_{i+1}$  is start of successor blocks

Different (in phase) no. of paths may be possible in a  $\text{prog}$

An analysis choose different state & diff info & no analysis

is possible repr of the state.

eg:-



Different possible paths  $\{1, 2, 3, 4, 9\}$ ,  $\{1, 2, 3, 4, 5, 6, 7, 8, 3, 4, 9\}$   
initially  $a=1$  reaches till 5  $d_2$  in first run in next  
run  $d_3$  reaches 3 so  $\{d_1, d_3\}$  are "reaching definition"

## The Dataflow Analysis Schema

For each app of DFA

\* at every prog pt, a dataflow value is associated, that  
reprs all possible prog states observed from that pt.

\* data flow for stmt 'S' is reprd as  $IN[S]$  and  $OUT[S]$   
reprs states before S and after S

two factors affect  $IN[S]$  and  $OUT[S]$

① semantics of stmt 'S' [transf & function]

② control flow.

It depends on all dependencies to reach the leader  
 but the decision is different between blocks  
 same new blocks and out [5] of out [5] did this

$$IN[5] = OUT[5]$$

of [5] can have 5...5n then

data block value are derived from block of out

Can't do this

$$IN[5] = OUT[5]$$

data block value before 5n

of 5 convert data block value of the 5n to new

2) backwards

$$a = a + 1$$

$$OUT[5] = IN[5]$$

new data block value after 5

of 5 raises data block value before 5n+5 and produce

1) forward

transfer function case in 2 blocks

be equal to some value.

$$a = a$$

of the 5n+ both b & a will



If block B contains  $s_1 \dots s_n$  then,  
 $IN[B] = IN[s_1]$        $OUT[B] = OUT[s_n]$

$$f_B = f_{s_n} \circ \dots \circ f_{s_2} \circ f_{s_1}$$

$$OUT[B] = f_B(IN[B])$$

$$IN[B] = \bigcup_{P \text{ (predecessor of } B)} OUT[P]$$

forward flow

$$IN[B] = f_B(OUT[B])$$

$$OUT[B] = \bigcup_{S \text{ (successor of } B)} IN[S]$$

backward flow

## Reaching Definitions

- \* Where in a program each variable 'x' is defined.
- \* A definition d reaches point P if there is a path from d to P such that d is not killed along the path.
- \* Kill of definition 'd' happens if 'x' is having any other definition along the path.
- \* A definition 'd' is any stmt that assigns value to 'x'

$$d: u = v + w$$

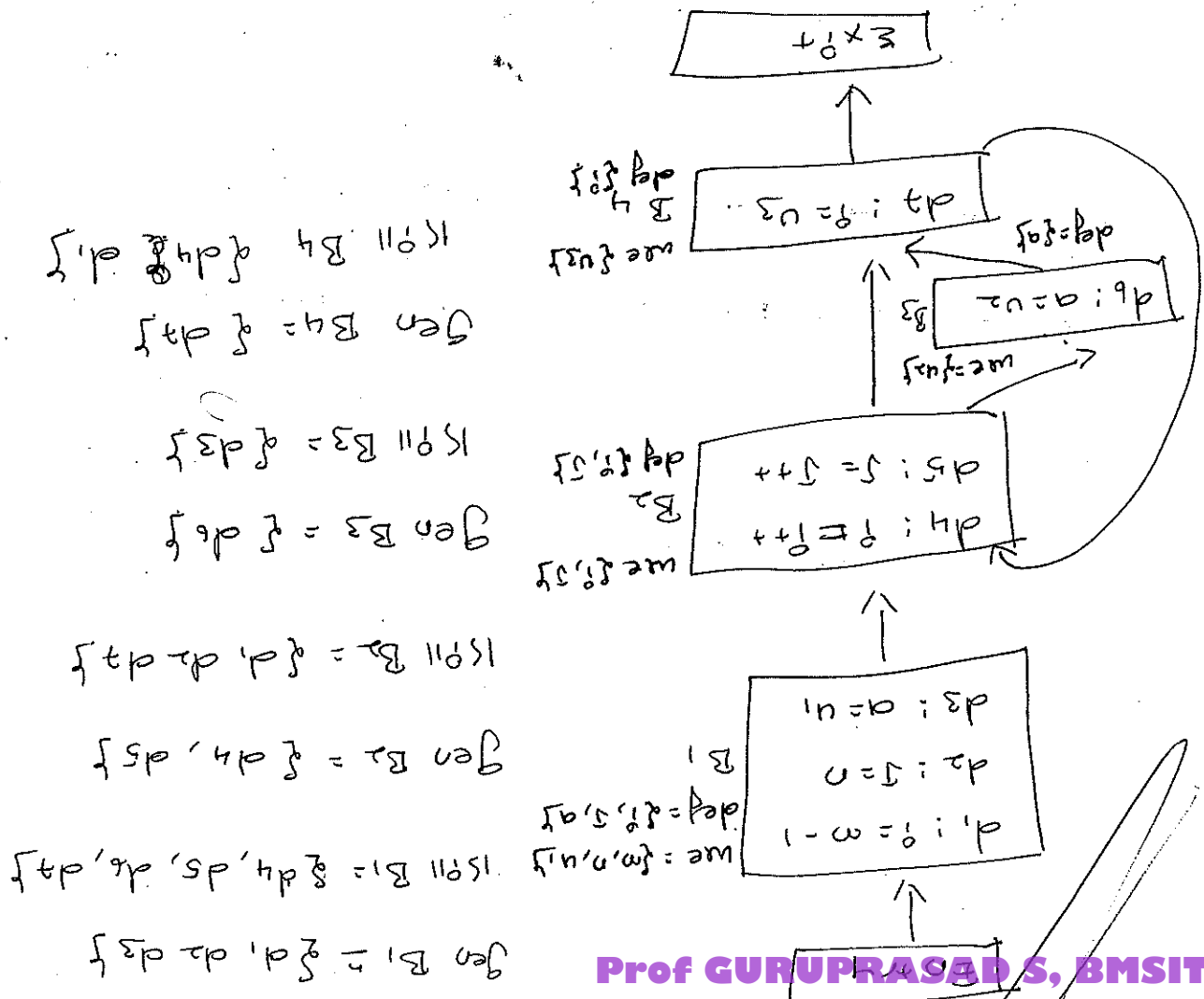
d generates a defn d of var u and kills all other defns of u

$$f_d(x) = \text{gen } u \cdot (x - \text{kill } d)$$

A: B.B generated set defined and its all set defn  
 gen set - downwards expand - are set of defn  
 generated in B.B (gen have precedence over B.B)

$d_1: a = 3$   
 $d_2: a = 4$

gen B<sub>1</sub> =  $d_2$   
 $d_1: a = 3$   
 $d_2: a = 4$



## Live Variable Analysis

Variable  $x$  is Live <sup>at pt P</sup> if  $x$  at pt  $P$  could be used along some path in flow graph starting at  $P$  else  $x$  is dead

Live var info is useful in register allocation for basic blocks.

Data flow equations can be defined w.r.t

(i)  $IN[B]$  and  $OUT[B]$

(ii) transfer function of  $B$

\*  $def B$  - set of var defined in  $B$  prior to any use in  $B$

\*  $use B$  - set of values used in  $B$  prior to  $def$  of var

So

Var in  $use B$  is live on entrance to  $B$

Var in  $def B$  is dead in  $B$  as path starts from  $B$

No variables are live on exit

Var is live coming into block  $B$  if it is used before redefinition

Var is live coming out of the block if it is not redefined in the block **OR** (it is used in the successive block.)

Out[n] = U<sub>s</sub> success[n] in [s]

In[n] = we[n] U (out[n] - def[n])  
Prof GURUPRASAD S, BMSIT&M  
<https://hemanthrajhemu.github.io>

# Available Expression

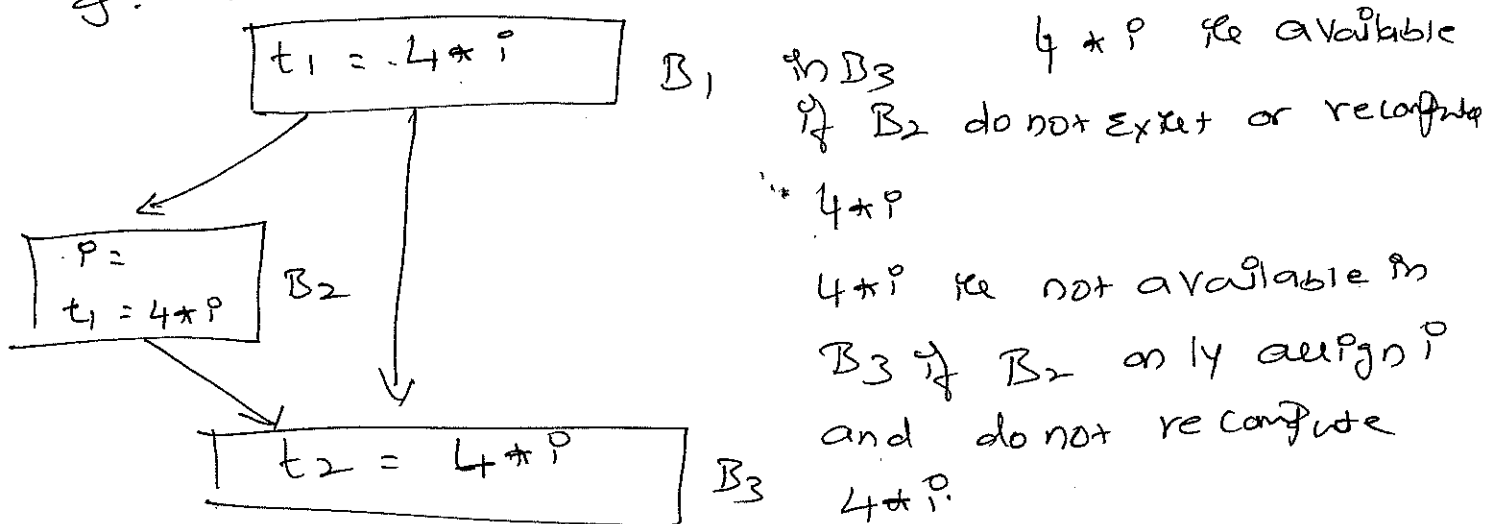
An expression  $x+y$  is available at point  $P$  if every path from entry to  $P$  evaluates  $x+y$  and after last evaluation  $x$  and  $y$  are not changed

A block generates  $x+y$  if it evaluates it

A block kills  $x+y$  if it assigns  $x$  or  $y$  and does not recompute  $x+y$

We! to detect global common sub Expressions

Eg:-



Eg:-

