

# *Compiler Design*

*This download is Sponsored by:*



*Contributed By:*  
**Sibananda Achari**  
Trident Academy of Technology

## **Disclaimer**

This document may not contain any originality and should not be used as a substitute for prescribed textbook. The information present here is uploaded by contributors to help other users. Various sources may have been used/referred while preparing this material. Further, this document is not intended to be used for commercial purpose and neither the contributor nor LectureNotes.in is accountable for any issues, legal or otherwise, arising out of use of this document. The contributor makes no representation or warranties with respect to the accuracy or completeness of the contents of this document and specifically disclaim any implied warranties of merchantability or fitness for a particular purpose. By proceeding further, you agree to LectureNotes.in Terms of Use.

This document was downloaded by: Vishvanath Surwshe on 04th May, 2017.

At LectureNotes.in you can also find

- 1. Previous Year Questions for BPUT**
- 2. Notes from best faculties for all subjects**
- 3. Solution to Previous year Questions**



# *Compiler Design*

Topic:

## *Introduction To Compiler Design*

Contributed By:

*Sibanya Achary*

Trident Academy of Technology

## → Introduction to Compilers

9/12/13

### Introduction to Compilers

- Need of Compiler
- Structure of Compiler
- Tools for compiler writing

#### Types of Compiler —

##### → one-pass Compiler:

- It passes through the source code of each compilation unit (phase) only once.

##### → Multi-pass Compiler:

- Processes the source code of a program several times.
- Multipass compilers are sometimes called whole compiler.

##### → Load and go Compiler:

- An compiling technique in which there is no stop between loading and execution.

##### → Optimizing Compiler:

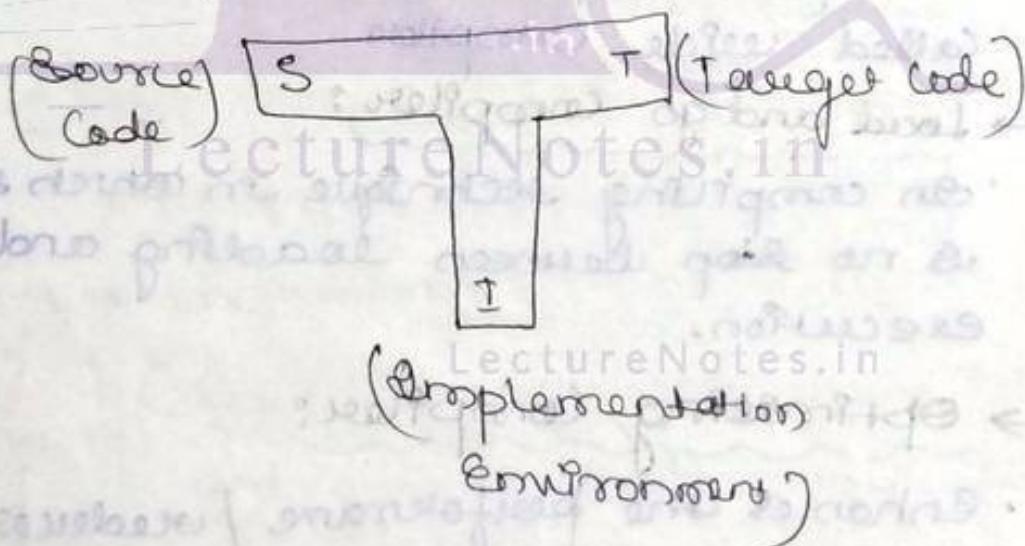
- Enhances the performance / reduces the size of the resulting machine program. (Optimizes the program code).
- They require multiple passes in order to analyse the entire program and maximise the use of code throughout.

- Can store & execute by writing simple instructions that can be executed fast.

### → Native Compiler -

- Compiles programs for the same architecture of OS that it is running on.
- Platform is a combination of CPU architecture and OS.
- Compiles programs that can be executed on same platform.

### J-Program -



### → Cross Compiler -

- Executes in an environment and generates code for another.

Ex: Microsoft visual studio. Previously a native compiler.

→ Embedded devices such as mobile phones, washing machines uses cross compilers. (11)

### → Semantics Analysis

→ This phase checks the program for semantics error and gathers type information for the successive phases.

Eg - Using a variable that is not declared is checked during semantics analysis.

#### • Type Checking -

→ Check types of operands (possibly imposing type conversions); No real no. for index or of bounds array etc.

→ CFG cannot be used to specify semantics rules. At the end of semantic analysis, syntax directed output is obtained.

→ Normally, semantic information cannot be represented by CFL used in syntax analysis. CFG used in Syntax Analysis are integrated with attributes (semantic rules). The result is, syntax directed translation and attribute grammar.

#### • Synthesis Phase -

→ Starts with Intermediate code generation.

### → Intermediate Code Generation

→ Intermediate code is machine independent.

→ An intermediate code is generated as a program for an abstract machine.

- The Intermediate code should be easy to translate into the target program.
- As Intermediate code we consider the three-address code, similar to assembly sequence of instructions with atmost three operands such that:
- There is atmost one operator, in addition to the assignment. Thus, we make explicit the operators precedence.
  - temporary names must be generated to compute intermediate operations.

Eg -

```

if (a <= b)
    a = a - c;
}
C = b * c;

```

### 3-Address Code

```

-t1 = a > b
if -t1 goto L0
-t2 = a - c;
a = -t2
L0 : -t3 = b * c;
C = -t3;

```

### Code Optimisation

- optimising the space consumed by variables and program.
- Different compilers have different optimisation techniques.

### Code Generation

- Generates the target code consisting of assembly code.

- Memory locations are selected for each variable.

(12)

- Instructions are translated into a sequence of assembly instructions.
- Variables and intermediate results are assigned to memory registers.

### Error Handling

- Lexical analyser reports invalid characters.
- Syntax analyser reports invalid token sequences.
- Semantics analyser reports invalid scope errors.

### Compiler Construction Tools

These all -

- Scanner generators
- Parser generators
- Syntax directed translation engines
- Code generator generator
- Data flow analysis engines

LectureNotes.in



---

Mindfire is a software service provider,  
with unrelenting focus on optimal software development  
and delivery.

---



# *Compiler Design*

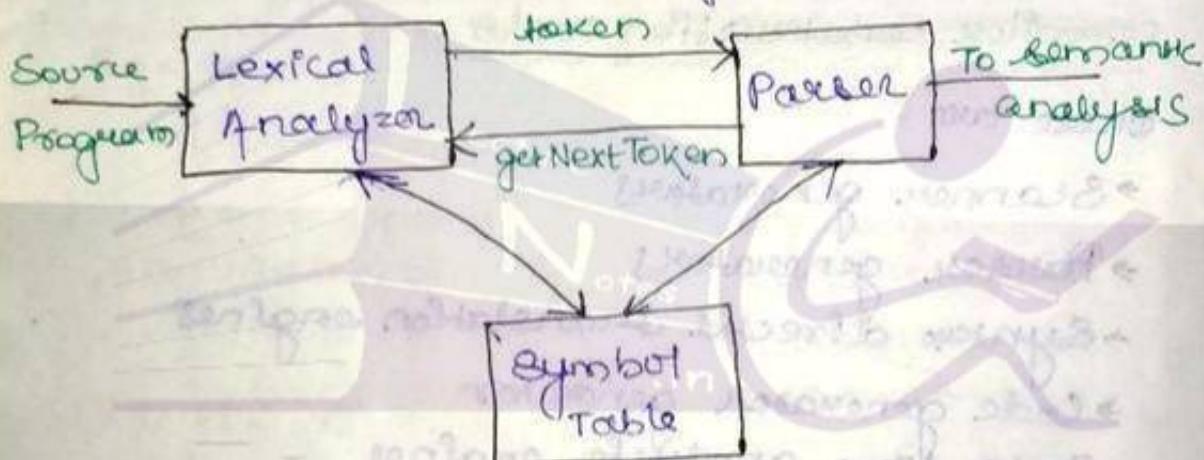
Topic:  
*Lexical Analysis, Finite Automata*

Contributed By:  
***Sibananda Achari***  
Trident Academy of Technology

## Lexical Analysis

The role of lexical analyser

- The first phase of the compiler
- Reads the input characters and produces as output a sequence of tokens that the parser uses for syntax analysis
- Strips out source program comments & white space information
- Keeps track of blank, tab & los.
- Generates error messages from the compiler with the source programs



Tokens, Patterns and Lexemes

- A token is a pair a token name and an optional token value.
  - A pattern is a description of the form that the lexemes of a token may take.
  - A lexeme is a sequence of characters in the source program that matches the pattern for a token.
- lexemes
- |            |  |                       |
|------------|--|-----------------------|
| expression | $a = b * c + 2$                              | → number              |
|            | $\downarrow$<br>$\downarrow$<br>$\downarrow$ | b add operator        |
|            |  | m ult-operator        |
|            |  | a s signment operator |

### Symbol Table -

- Matrix like structure having rows and columns.
- Scanning the entire strings from left to right.
- Here,

a

LectureNotes.in

- i denotes pointer to the symbol table entry for a.

→

1	a	-	-
2	b	-	-
3	c	-	-

- For operators we are not maintaining a pointer to the symbol table.
- for identifiers we are maintaining a pointer to the symbol table.

Let if (expr)

else

So, the regular expression for above is -

$S \rightarrow \text{if } (\text{expr}) \text{ statement} | \text{else } S \text{ statement}$

If  $S \rightarrow \text{if } \text{expr } \text{stmt} | E$ , then compiler will generate error as rule is not specified in the regular expression.

(13)

(15)

→ Literal means a sequence of characters enclosed in double quotes.

### Example

Tokens	Pattern	Lexemes
if	characters i,f	if
else	characters e,l,s,e	else
Comparison	<ox> or <= ox> = ox == ox !=	<=, !=
id	letter followed by letter and digits	pr, score, D2
number	Any numeric Constant	3.14159, 0, 6.02e23
literal	Anything but surrounded by "	"case dumped"

### Attributes for tokens

- Apart from token itself lexical analyzer passes other information regarding token, which is called attributes.
- It helps in compilation when more than one pattern matches a lexeme
- Eg : E = M \* C \*\* 2  
 <rd, pointer to st entry for E>  
 <assign-op>

<Id, pointer to ST for M>

<mult-op>

<Id, pointer to ST entry for C>

<exp-op>

<number, Integer Value 2>

(Ch-3 Pg-110-129)

### Input Buffering

→ It is desirable for the lexical analyzer to read its input from an input buffer because many characters beyond the next token may have to be examined before the next token itself can be determined.

→ 2 pointers are used in the I/P buffer -

\* lexeme Beginning is a pointer marks the beginning of the token being discovered.

\* look ahead pointer scans ahead of the beginning point, until the token is discovered.

lexeme Beginning

↓  
look ahead pointer

$$\text{Begin}^1 = l + 1 ; j = j + 1 \dots$$

L.B.      L.A.

↓  
 $\text{Begin}^1 = l + 1 ; j = j + 1 \dots$

$$\text{Begin}^2 = l + 1 ; j = j + 1 \dots$$

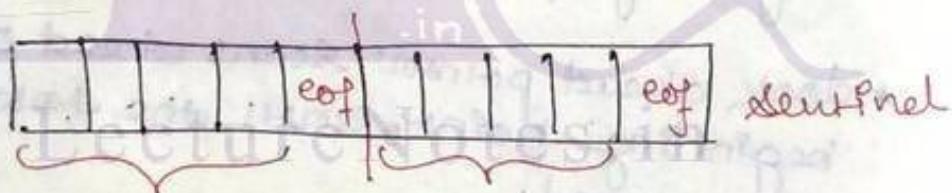
There are different buffering scheme which are used commonly by compilers.

→ One buffer scheme -

There are problems if a lexeme crosses the buffer boundary. To scan the rest of lexeme, the buffer has to be refilled thereby overwriting the first part of the lexeme.

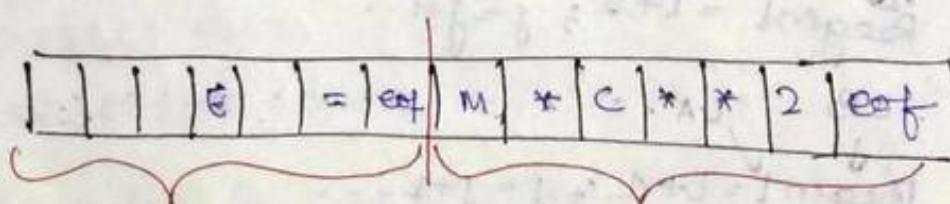
→ Two buffer scheme -

Here, buffers 1 and 2 are scanned alternately. When the end of the current buffer is encountered the other buffer is filled. Hence, the problem encountered in previous scheme is solved.

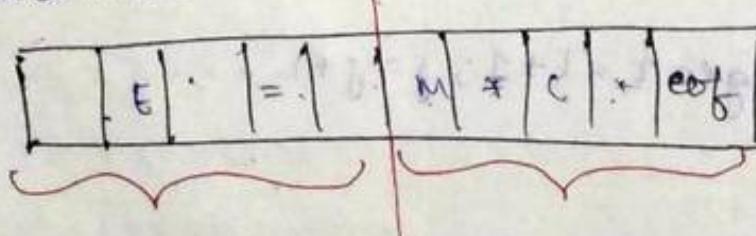


Let

$$E = M * C * * 2$$



Sometimes -



## Lexical Error

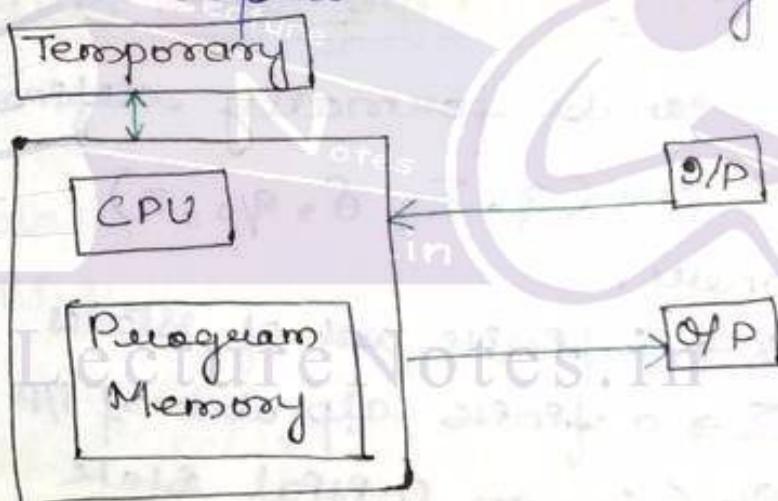
13/12/13 (n)

Some errors are out of power of lexical analyzer to recognize.

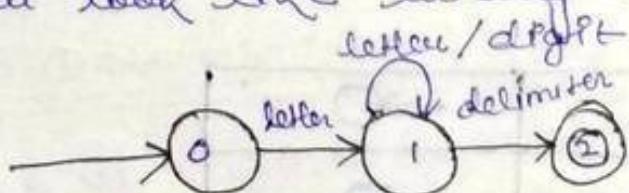
## Lexical Analysis

- Finite State Machine
- Regular Expression
- Lexical Analyzer
- Automaton - It is a machine responsible for processing the input and generating output.

It represents CPU + Program memory



- In finite state automaton, temporary memory is absent. Then the structure will look like Turing machine & PDA.



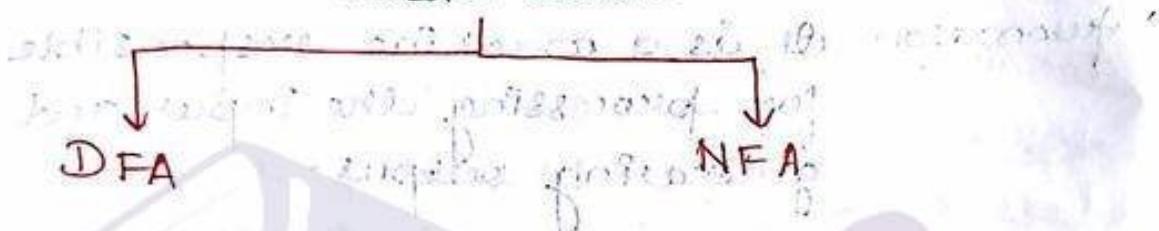
- An identifier can be defined as a string of letters and digits that begins with a letter.

→ State Automaton: An automatic transition from one state to another state or set of state on an input.

→ Finite State: Due to finite no. of state automaton

LectureNotes.in

### Finite automata



### Deterministic Finite State Automata

→ A DFA can be formally defined as

$$A = (Q, \Sigma, \delta, q_0, F)$$

where,

↳ Q → a finite set of states

↳ Σ → a finite alphabet of I/P symbols

↳  $q_0 \in Q$ , an initial state

↳  $F \subseteq Q$  → a set of final states

↳  $\delta: (Q \times \Sigma) \rightarrow Q$  → a transition function.

→

Initial State	$\xrightarrow{q_0}$
Final State	○
Other State	○
Transition/move	$\xrightarrow{\quad}$



---

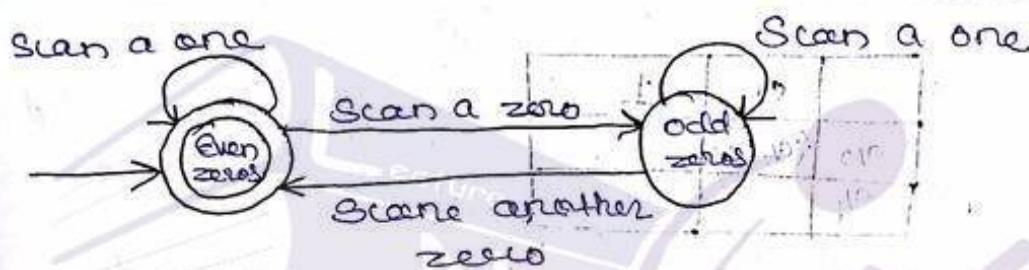
Mindfire is a software service provider,  
with unrelenting focus on optimal software development  
and delivery.

---

- The movement will be from left to right. ②
- No temporary m/m but program memory having I/P step.
- Eg:



- Eg: Check whether a binary no. has even no. of zeros or not.



- Generally DFA this is called a machine which accept a language.
- The language is the set of string the machine accept. which is denoted as L(M).
- $L(M) = \{ \text{Binary string having even no. of zeros} \}$

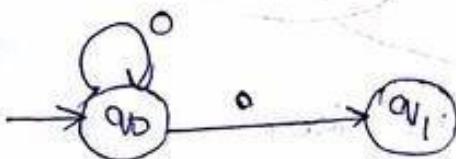
### Non-Deterministic Finite State Automata

- NFA is a finite state machine where from each state and a given input symbol the automaton may jump into several possible next states.

(21)

→ This distinguishes it from the DFA, where the next possible state is uniquely determined.

→ Eg -



→ Here,  $\delta = Q \times \Sigma \rightarrow P(Q)$

→ State Table

	0	1
q0	{q0, q1}	
q1		

→ A NFA is represented formally by a 5-tuple  $(Q, \Sigma, \Delta, q_0, F)$  consisting of

where,

• a finite set of states  $Q$

• a finite set of I/P symbols  $\Sigma$

• a transition relation

$$\Delta : Q \times \Sigma \rightarrow P(Q)$$

• an initial (or start) state  $q_0 \in Q$

• a set of states distinguished as accepting or final states  $F \subseteq Q$ .

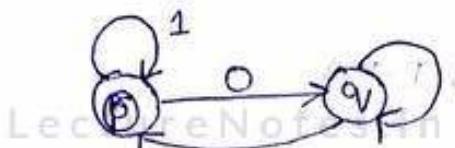
→ In a Deterministic DFA, reading

one symbol from address increases

→ Transition function takes 2 input.

- an input tape
- current state

→ Eg -  $\delta(P, 0) = \{q_1\}$ ,  $\Sigma = \{0, 1\}$



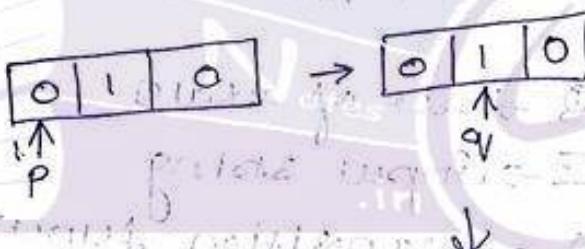
$$\delta(P, 1) = \{P\}$$

$$\delta(q_1, 0) = \{P\}$$

$$\delta(q_1, 1) = \{q_1\}$$

to generate a language for the automata

def



0 | 1 | 0

A language is a set of strings accepted by the finite automata.

### E-NFA (ε-Transition of a NFA)

→ allows us to make a transition with the empty string.

→ spontaneously make a transition without receiving an I/P symbol

→ allows a transformation to a new state without consuming any input symbols. (23)

DFA	G-NFA
$\delta(q_i, \epsilon) = \{q_j\}$	$\delta(q_i, \epsilon) = \{q_j\}$ $q_i, q_j \in Q$ $q_i \neq q_j$

$G\text{-NFA} \rightarrow NFA \rightarrow DFA$

LectureNotes.in

→ A NFA - G is represented formally by a 5-tuple  $(Q, \Sigma, \Delta, q_0, F)$  consisting of where,

$Q \rightarrow$  set of states

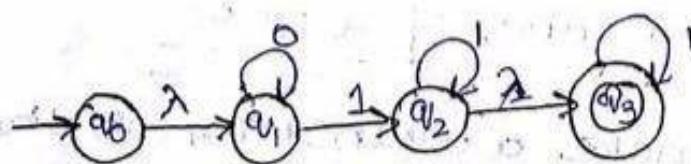
$\Sigma \rightarrow$  Input string

$\Delta \rightarrow$  transition function

$q_0 \rightarrow$  initial state

$F \rightarrow$  final state

Conversion of G-NFA to NFA



Here,  $M = \{ (q_0, q_1, q_2, q_3), \{0, 1, 2\}, \{q_0, q_3\} \}$

→ This is an equivalent NFA to the given G-NFA.

$M' = \text{NFA without } \delta \text{ moves or } \epsilon\text{-transitions}$  (28)

Step I : To find meet the  $\epsilon$  closure of the given automata -

$$\epsilon\text{-closure } (\bar{q}_0) = \{ \bar{q}_0 \} \cup \delta(\bar{q}_0) = \{ \bar{q}_0, q_1 \}$$

$$\epsilon\text{-closure } (\bar{q}_1) = \{ \bar{q}_1 \}$$

$$\epsilon\text{-closure } (\bar{q}_2) = \{ \bar{q}_2 \} \cup \delta(\bar{q}_2) = \{ \bar{q}_2, q_3 \}$$

$$\epsilon\text{-closure } (\bar{q}_3) = \{ \bar{q}_3 \}$$

Here,  $M' = \{ \bar{q}_0, q_1, q_2, q_3, \bar{q}_2, q_3, \bar{q}_3, q_0, q_1, q_2, q_3, \bar{q}_0, \bar{q}_1, \bar{q}_2, \bar{q}_3 \}$

E-NFA or NFA do DFA conversion  
 For each NFA we can find an equivalent DFA which accepts the same language.

### Subset construction algorithm

Initially,  $\epsilon$ -closure ( $\bar{q}_0$ ) is the only state in Dstates, and it is unmarked, while (there is an unmarked state in Dstates)

§ mark  $T_0$

for each i/p symbol  $a$ )

§  $U = \epsilon\text{-closure } (\text{move}(T, a))$ ,

if ( $U$  is not in Dstates)

add  $U$  as an unmarked state to Dstates;

Dstate  $[T, a] = U$ ,

to be deal in L.P. + i/p

(25)

With all states to be ~~stack~~

19/12/13 (a)

Conversion of NFA to DFA

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

D states -

$$T = \text{marked}$$

for each input 0 and 1

(move ( $\emptyset, q_0, q_1, 0$ ))

$$= E\text{-clo}(\text{move}(\emptyset, q_0, q_1, 0) \cup \text{move}(\emptyset, q_1, 0))$$

$$= E\text{-clo}(\emptyset \cup q_1)$$

$$= E\text{-clo}(q_1)$$

$$= \{q_1\}$$

$$\text{det } \{q_1\} = B$$

$$\text{Now, } T = \{q_0, q_1\}$$

$$B = \{q_1\}$$

The next step is -

$$E\text{-clo}(\delta(\{q_0, q_1\}, 1))$$

$$= E\text{-clo}(\delta(\{q_0\}, 1) \cup \delta(\{q_1\}, 1))$$

$$= E\text{-clo}(\emptyset \cup q_2)$$

$$\text{det } \{q_2\} = E\text{-clo}(q_2)$$

$$= \{q_2, q_3\}$$

$$\text{det }$$

$\{q_2, q_3\}$  is labelled as C.

(26)

Now,  $\text{① } \{q_0, q_1\}$

$$B = \{q_1\}$$

$$C = \{q_2, q_3\}$$

The next step is check for B -

$$\epsilon\text{-clo}(\delta(\{q_1\}, 0))$$

$$= \epsilon\text{-clo}(q_1)$$

$$= \{q_1\} = B \text{ (already exists)}$$

Next,

$$\epsilon\text{-clo}(\delta(\{q_1\}, 1))$$

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

$$= \{q_2, q_3\} = C \text{ (already exists)}$$

Now,  $\text{① } \{q_0, q_1\}$

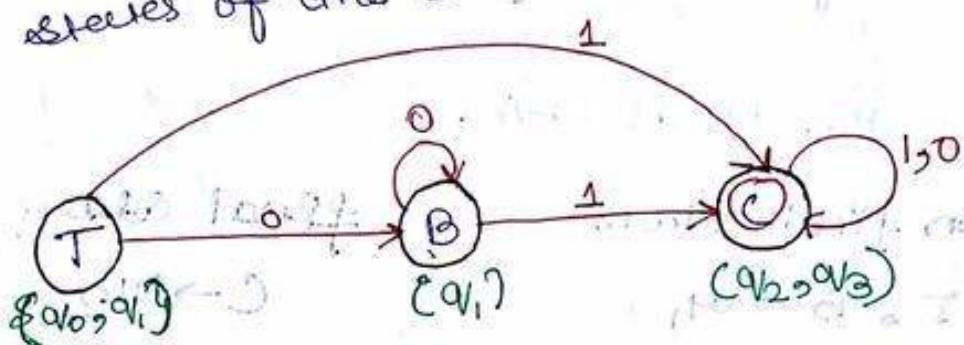
$$\textcircled{B} = \{q_1\}$$

$$\textcircled{C} = \{q_2, q_3\}$$

$$\epsilon\text{-clo}(\delta(C, 0)) = \{\phi\} \text{ (states)}$$

$$\epsilon\text{-clo}(\delta(C, 1)) = \{q_2, q_3\} = C$$

$\therefore 3$  states are discovered which are the states of the DFA:



$$D\text{Trans}[T, 0] = \cup$$

$$D\text{Trans}[T, 1] = B \cup C$$

$$D\text{Trans}[B, 1] = C$$

$$D\text{Trans}[C, 1] = C$$

## Minimization of DFA

### Algorithm:

- Partition the set of states into 2 groups:
    - $G_1$ : set of accepting states
    - $G_2$ : set of non accepting states
  
  - For each new group  $G_i$ 
    - partition  $G_i$  into subgroups such that
      - \* states  $s_1$  &  $s_2$  are in the same group
      - \* iff for all i/p symbols  $a$ ,
      - \* states  $s_1$  &  $s_2$  have transitions to states in the same group
  
  - Start state of the minimized DFA is the start state of the group containing the start state of the original DFA.
  
  - Accepting states of the minimized DFA are the groups containing the accepting states of the original DFA.
- from the final DFA,
- |                        |                     |
|------------------------|---------------------|
| non final states       | final states        |
| $T, B \rightarrow G_1$ | $C \rightarrow G_2$ |



---

Mindfire is a software service provider,  
with unrelenting focus on optimal software development  
and delivery.

---

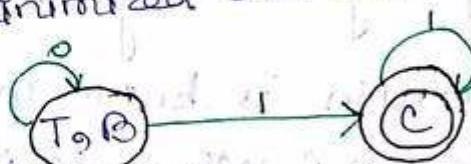
$$\delta(T_0, 0) = B \quad S_1$$

$$\delta(B, 0) = B \quad S_2$$

$$\delta(T_0, 1) = C \quad S_1 \quad (29)$$

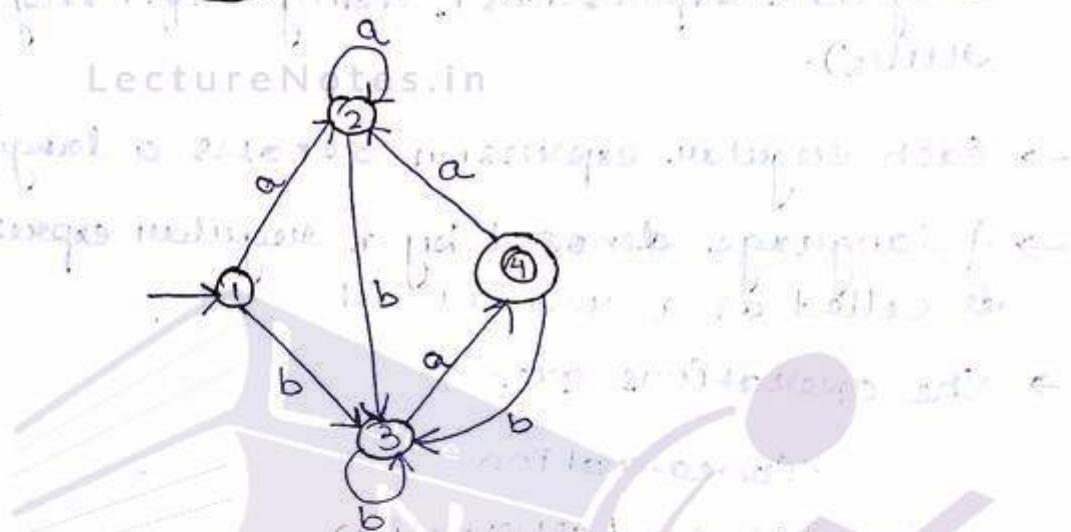
$$\delta(C, 1) = C \quad S_2$$

so minimized DFA is,



Q

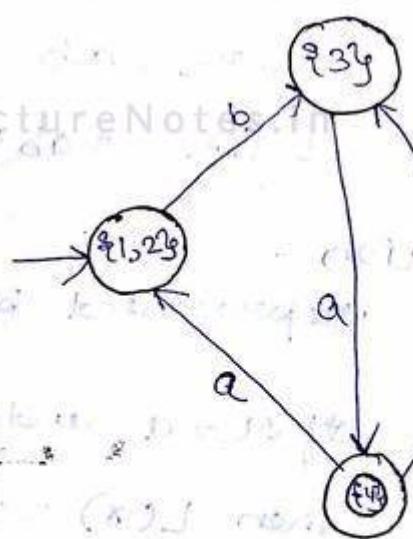
LectureNotes.in



Groups :  $\{1, 2, 3\}$  &  $\{4\}$

$\{1, 2, 3\}$  &  $\{3\}$   
 (no more partitioning)

a	b
$1 \rightarrow 2$	$1 \rightarrow 3$
$2 \rightarrow 2$	$2 \rightarrow 3$
$3 \rightarrow 4$	$3 \rightarrow 3$



(29)



---

Mindfire is a software service provider,  
with unrelenting focus on optimal software development  
and delivery.

---



# *Compiler Design*

Topic:

***Lexical Analysis Regular Expressions***

Contributed By:

***Sibanya Achary***

Trident Academy of Technology

# Regular Expressions

- We use regular expression to describe tokens of a programming language.
- A regular expression is built up of simple regular expressions (using defining rules).
- Each regular expression denotes a language.
- A language denoted by a regular expression is called as a regular set.
- The operations are -
  - Concatenation
  - Union / alternation
  - Kleen Closure

## → Concatenation -

If  $\alpha_1 = a$ ,  $\alpha_2 = b$

$$\alpha_1 \cdot \alpha_2 = ab$$

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

## → Union -

represented by + or | → pipeline

If  $\alpha_1 = a$  and  $\alpha_2 = b$ .

$$\text{then } L(\alpha) = \{a, b\}$$



A regular language is a formal language that is accepted by a DFA, NFA and read-only TM.

A regular language can be described by a regular expression built on simple regular grammar.

### Operations on languages

→ Concatenation -

$$L_1 L_2 = \{s_1 s_2\} : s_1 \in L_1 \text{ & } s_2 \in L_2 \}$$

→ Union -

$$L_1 \cup L_2 = \{s_1 | s_1 \in L_1 \text{ or } s_1 \in L_2\}$$

→ Exponentiation -

$$L^0 = \{\epsilon\} \quad L^1 = L \quad L^2 = LL$$

$$w = ab, w^0 = \epsilon, w^2 = w \cdot w = abab$$

→ Kleene Closure : Concatenate L with itself any no. of time

$$L^* = \bigcup_{i=0}^{\infty} L^i \quad (\text{It includes empty string with other strings generated})$$

$$L = \{a, b\} \cup \{\epsilon\}$$

→ Positive closure : Concatenate L with itself at least one or more

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

(It excludes empty string)

$$(a, b)^* = \{\epsilon, ab, a, a, ab, b, b, \dots\}$$

$$(a, b)^+ = \{a, b, aa, ab, ba, bb, \dots\}$$

# Converting a regular expression into a NFA (Thomson's Construction)

Conversion of expression  $\rightarrow$  NFA.

sup:

Regular Expression to DFA directly

23/12/13

(17S-180)

Augment the regular expression or with a special end symbol # to make accepting states important the new expression is #

$\delta \# \rightarrow$  augmenting  
functions -

$$\delta = (a|b)^* a$$

$$\delta \# \rightarrow (a|b)^* a \# \rightarrow \text{leave node}$$



Syntax tree

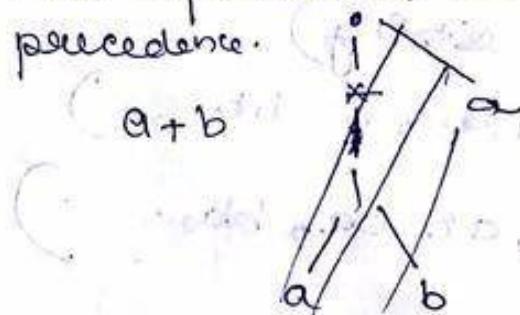
To construct a syntax tree we need precedence ( $* \cdot | \circ$ )

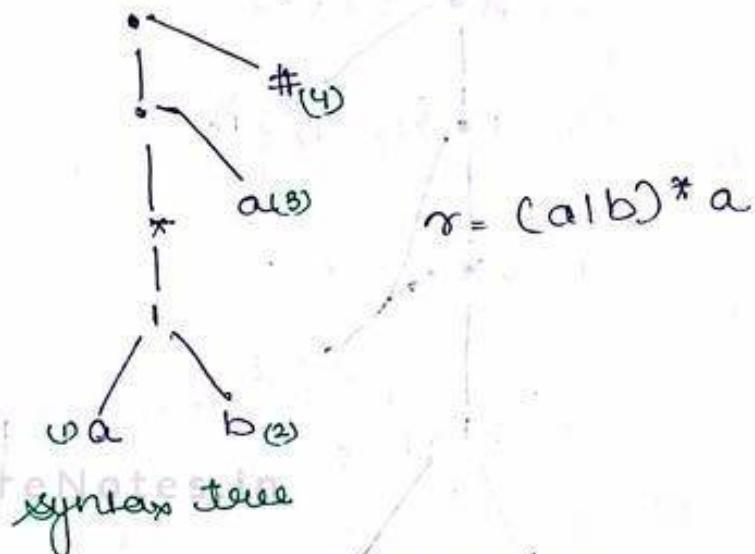
$\star \rightarrow 1$  (asterisk)

$\circ \rightarrow 2$  (concatenation)

$1 \rightarrow 3$  (concatenation)

For expression, we have () has the highest precedence.





we numbered the leaf nodes.

Now, the final statement is to calculate the nodes.

**nullable:** The subtree at node  $n$  generates language including the empty string.

**Firstposition:** The set of positions that can match the first symbol of a string generated by the subtree at node  $n$ .

**Lastposition:** The set of position that can match that last symbol of a string generated by the subtree at node  $n$ .

**Followposition:** The set of position that can follow position  $i$  in the tree.

(1) ~~soql~~ (2) ~~soqlsp~~ ~~soqlsp~~

I calculate 1<sup>st</sup> & last pos<sup>2</sup> of each node. (34)

{1,2,3} o {4}

{1,2,3} o {3} {4} {4} {4}

{1,2,3} \* {2,2,3}

{1,2,3} {1,2}

{1,3} {1,3} {2,2,2,2}

(no itself)  
 last pos<sup>2</sup>  
 first pos<sup>2</sup>  
 (no itself)

How to calculate firstpos & lastpos?

P	nullable(m)	firstpos(n)	lastpos(m)
leaf labelled E	true	o	o
leaf labelled with pos <sup>2</sup>	false	{1}	{1}
c <sub>1</sub> , c <sub>2</sub>	nullable(c <sub>1</sub> ) or nullable(c <sub>2</sub> )	firstpos(c <sub>1</sub> ) or firstpos(c <sub>2</sub> )	lastpos(c <sub>1</sub> ) or lastpos(c <sub>2</sub> )
c <sub>1</sub> , c <sub>2</sub>	nullable(c <sub>1</sub> ) 3rd nullable(c <sub>2</sub> ) nullified	if nullable(c <sub>1</sub> ) 1 <sup>st</sup> pos(c <sub>1</sub> ) 1 <sup>st</sup> pos(c <sub>2</sub> ) else 1 <sup>st</sup> pos(c <sub>2</sub> )	if nullable(c <sub>2</sub> ) 1 <sup>st</sup> pos(c <sub>1</sub> ) lastpos(c <sub>2</sub> ) else lastpos(c <sub>2</sub> )
c <sub>1</sub> *	true	firstpos(c <sub>1</sub> )	lastpos(c <sub>1</sub> )

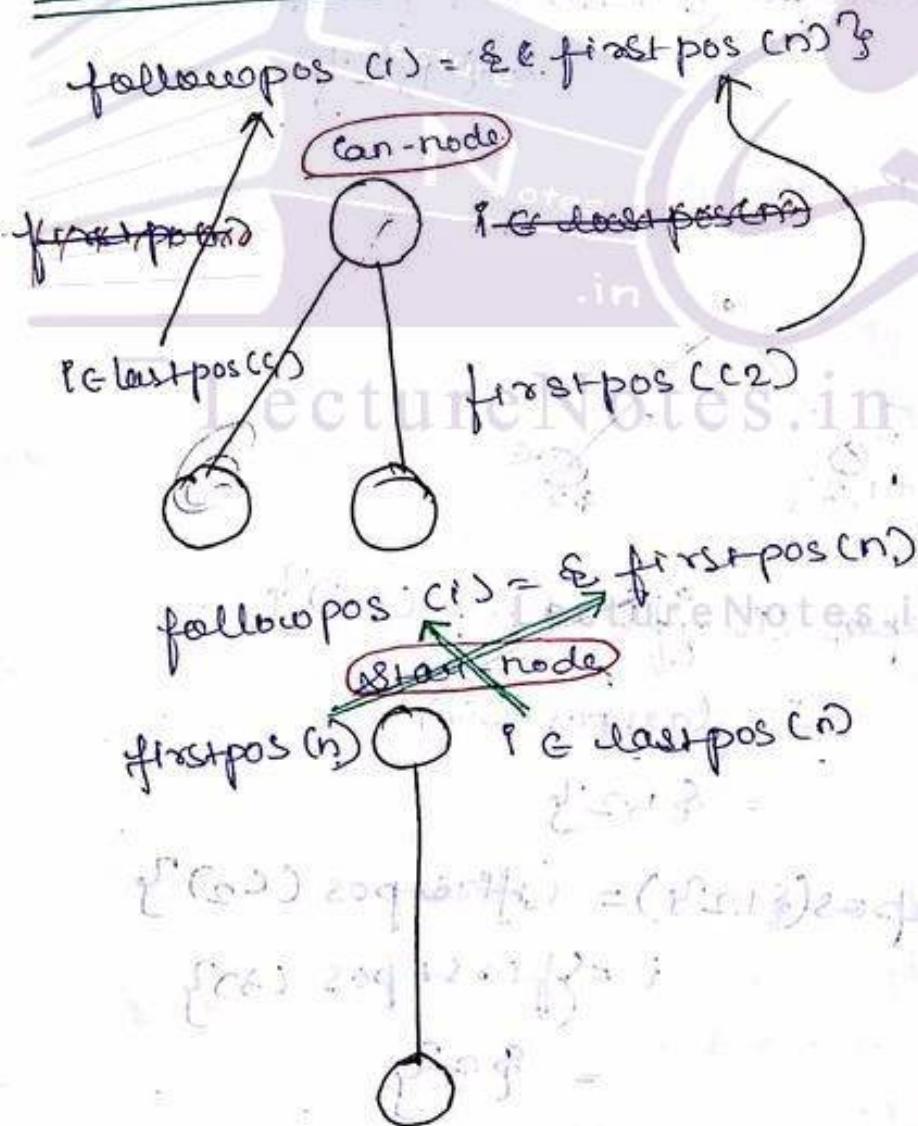
## How to evaluate followpos?

(35)

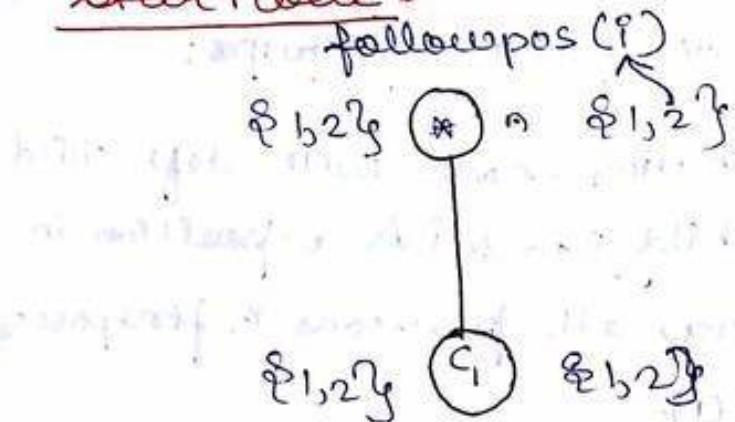
2 rules define the function followpos:

- If  $n$  is concatenation-node with left child  $c_1$  and right child  $c_2$ , &  $i$  is a position in  $\text{lastpos}(c_1)$ , then all positions in  $\text{firstpos}(c_2)$  are in  $\text{followpos}(r)$ .
- If  $n$  is a star-node, and  $i$  is a position in  $\text{lastpos}(n)$ , then all positions in  $\text{firstpos}(n)$  are in  $\text{followpos}(n)$ .

### VISUALISATION.



(35)

Start node :

LectureNotes.in

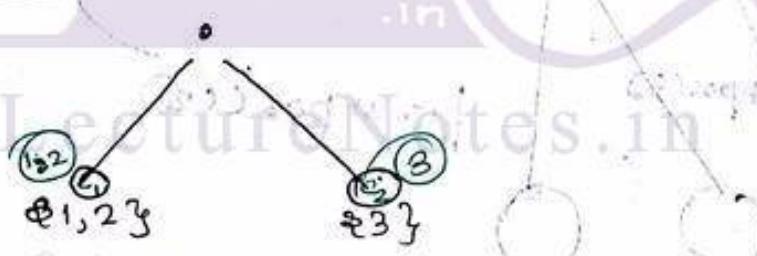
$$\text{if } i \in \text{lastpos}(C_1) \rightarrow * \text{ (value = lastpos } \{1, 2\})$$

$i = \{1, 2\}$

$$\text{followpos}(r) = \{ \text{firstpos}(r) \}$$

$$\text{followpos}(\{1, 2\}) = \{ \text{firstpos}(*) \}$$

$$= \{1, 2\}$$

Concatenation node :

$$\text{followpos}(r) = \{ \text{firstpos}(C_2) \}$$

$$i = \text{lastpos}(C_1)$$

$$= \{1, 2\}$$

$$\text{followpos}(\{1, 2\}) = \{ \text{firstpos}(C_2) \}$$

$$i = \{ \text{firstpos}(C_3) \}$$

$$= \{3\}$$

## "followpos"

Then we define the function followpos for the positions leafnode and not for the nonleaf node.

06/12/13

$S\# = (a \mid \epsilon) b c^*$

Find out - nullable

- firstpos

- lastpos

- followpos

- create DFA using the algorithm

## 2) Lexical errors

→ Some errors are out of power of lexical analyzer to recognize.

$f_i(a = f(x))$

→ Here  $f_i$  is taken as an identified

\* The advantage of two buffer is it overcomes overflow.

## Error Recovery

→ Panic mode : When lexical analyzer is not able to continue with the process of compilation, it resorts to panic mode of recovery.

→ In this mode the following actions are performed to identify a token.

- successive characters are ignored until we reach to a well formed token.

- delete one character from the remaining input.
- insert a missing character into the remaining input
- replace a character by another character.
- Transpose two adjacent characters.

## Lex

→ Lex is a scanner generator

- Input is description of patterns & actions.
- Output is a C program which contains a function yylex() which, when called, matches patterns and performs actions per input.
- Typically, the generated scanner performs lexical analysis and produces tokens for the (YACC-generated) parser.

(YACC - Yet Another C Compiler)

## The Lex and Flex Scanner Generators

→ Lex & its newer cousin flex are scanner generators.

→ Systematically translate regular definitions into C source code for efficient scanning



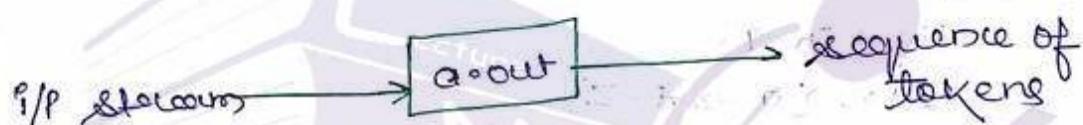
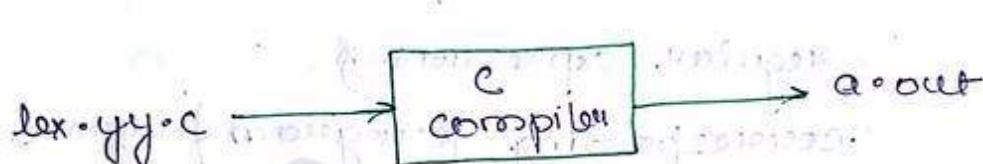
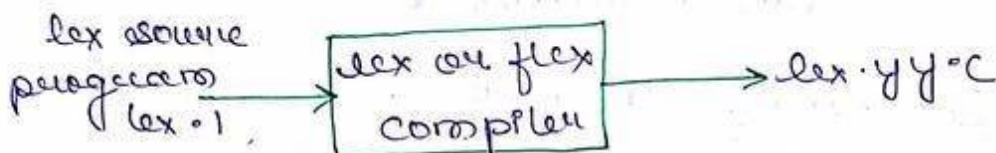
---

Mindfire is a software service provider,  
with unrelenting focus on optimal software development  
and delivery.

---

→ Generated code is easy to integrate in C applications. (39)

### Creating a lexical Analyzer with Lex & FLEX



### - Lex Specification (Pg-143)

- A lex specification consists of 3 parts:
  1. regular definitions, C declarations
  2. translation rules
  3. user-defined auxiliary procedures

- The translation rules are of the form -

```
P1 { action1; }  
P2 { action2; }  
⋮  
Pn { actionn; }
```

(39)

$$\delta(q_i, a) = q_j$$

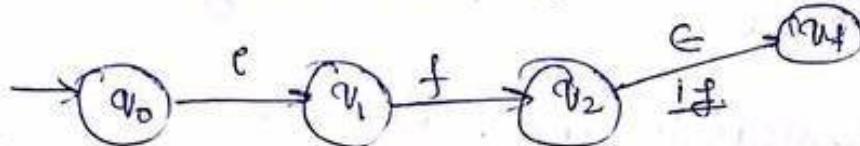


Fig. transition diagram to generate an identifier if

### Lex Source Program

→ Lex source is a table of

- regular expressions &
- corresponding program fragments

e.g.

digit [0-9]

letter [a-zA-Z]

%{

%} %-

{ letter } { letter } | ( digit ) +

10

%%

main()

{

yylex(); ← Invokes the lexical analyzer

}

cout << the matching lexeme

printf ("id : %s\n", yytext);

printf ("new line\n")

The table is translated to a C pgm (lex.yy.c)

### Lex fragment C program

→ Consider the Lex fragment:

- `%{ printf ("read 'a'\n"); }`

- `b%{ printf ("read 'b'\n"); }`

→ After compiling we obtain a binary executable which when executed on the off.

- The code produces

`abghijkl read 'a'`

`ghijkl read 'b'`

`fghijkl read 'b'`

### Example - A lex program

```
%{  
    int abc_count, xyz_count;  
}  
%
```

```
ab [cC] %{ abc_count++; }  
xyz %{ xyz_count++; }  
in %;
```

```
% %.  
main()  
{  
    abc_count = xyz_count = 0;  
    yylex();  
    printf ("%d occurrences of abc or abc\n",  
           abc_count);  
    printf ("%d occurrences of xyz\n", xyz_count);  
}
```

abc\_count = xyz\_count = 0;

yylex();

printf ("%d occurrences of abc or abc\n", abc\_count);

printf ("%d occurrences of xyz\n", xyz\_count);

for NOT having oil bearing remaining  
padding required to print analogously



---

Mindfire is a software service provider,  
with unrelenting focus on optimal software development  
and delivery.

---



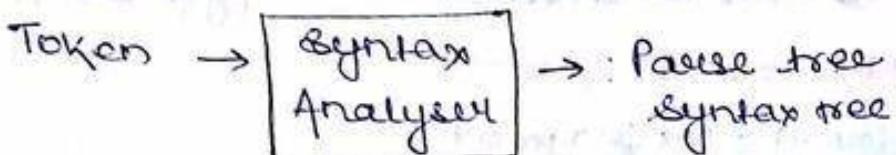
# *Compiler Design*

Topic:  
*Syntax Analysis*

Contributed By:  
***Sibananda Achari***  
Trident Academy of Technology

## SYNTAX ANALYSIS

(12)



→ Tool used for syntax and (syntactic constraint)  
Yacc is YACC.

→ SA<sup>2</sup> = parser ← interacts with → symbol table

→ CPU is type-2 grammar, recognized by PDA.

→ DPDA → has odd palindrome recognition.

→ NDPDA → has even palindrome recognition

→ CPU consists of → production rules

- set of NT

- set of T

- starting variable,

→  $A \rightarrow \alpha$   
 $\downarrow$   
 $NT$       ↳ string of grammar symbols (T & NT)

→  $\alpha\beta\gamma$  → denotes string of grammar symbols.

\* The symbols that are used as terminals are specified in the following form -

→ lower case letters early in the alphabet as

a, b, c, ...

→ operators as +, \*, /, etc.

→ parentheses, comma

→ digits like 0 to 9 & bold face identifiers as id, if

\* The symbols for NT are -

→ uppercase letters early in the alphabet

→ start symbol

→ lower case italics as expr, stmt.

→ uppercase letters as x, y, z represent grammar symbol i.e. either T or NT.

② →  $\alpha, \beta, \gamma$  represent string of grammar symbols

## Difference between SA & LA

(43)

- Separating the syntactic structure of a language into lexical and non lexical parts provides a convenient way of modularising front end of compiler.
- RE are concise & easy to understand notation for tokens.
- By using RE, constructs like identifiers, keywords, constants and whitespaces can be described in a proper manner. However, CFG are used for describing nested string as if else else stat ; balanced parenthesis etc.

## Ambiguity

- Derivation of parse tree -

it follows 2 ways -

- leftmost derivation,  $\xrightarrow{\text{LMD}}$

- rightmost derivation,  $\xrightarrow{\text{RMD}}$

eg -

$$\begin{array}{c}
 E \\
 | \\
 E + E \\
 | \quad | \\
 id \quad id
 \end{array} \Rightarrow \begin{array}{l}
 E = E + E \xrightarrow{\text{RMD}} E + \text{id} \xrightarrow{\text{RMD}} \text{id} + \text{id} \\
 \uparrow \qquad \uparrow \\
 \text{id} + E \\
 \uparrow \\
 \text{id} + \text{id}
 \end{array}$$

If there exists 2 or more parse tree for a given CFG, then grammar is said to be ambiguous, needs to be resolved.

~~eg -~~  $E \rightarrow E + E \mid E * E \mid E/E \mid (E) \mid \text{id}$

- This is an ambiguous grammar. To resolve :

- precedence

- associativity need to be considered

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

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

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

T = term

F = factor

(43)

## Elimination of left recursion:

A grammar is said left recursive if it has a NT A such that there is a derivation  $A \xrightarrow{*} A\alpha$ .

Rule to overcome:

For  $A \rightarrow A\alpha | \beta$  including  $\epsilon$  then reduce left recursion by  $A \rightarrow \beta A'$ ,

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

LectureNotes.in provides notes, assignments, question papers, and much more for all engineering students.



LectureNotes.in

Top rank notes, solved examples and previous year question papers for all engineering students. Visit us at [www.lecturenotes.in](http://www.lecturenotes.in)

Engineering notes, solved examples and previous year question papers for all engineering students.

Engineering notes, solved examples and previous year question papers for all engineering students.

Engineering notes, solved examples and previous year question papers for all engineering students.

Engineering notes, solved examples and previous year question papers for all engineering students.

Engineering notes, solved examples and previous year question papers for all engineering students.

Engineering notes, solved examples and previous year question papers for all engineering students.

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

$\Downarrow$  eliminate left recursion

$$A \rightarrow BA'$$

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

To resolve ambiguity, we need to rewrite it using LR(0) -

So, from the last example,

i.e.

$$E \rightarrow E+E | E * E | (E) | id$$

$$\cancel{E \rightarrow E + T | T}$$

$$T \rightarrow T * F | F$$

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

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

$$\text{Ans} \rightarrow E \rightarrow TE' \quad \text{LectureNotes.in}$$

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

$$T \rightarrow FT'$$

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

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

$$1 \quad S \rightarrow Aa | b$$

$$A \rightarrow AC | Sd | \epsilon$$

$$\rightarrow S \rightarrow Aa | b$$

$$A \rightarrow AC | Sd | \epsilon$$

$$A \rightarrow AC$$

$$A \Rightarrow Sd$$

$$\rightarrow Aa$$

(Pg - 213)

The algorithm says that-

→ Arrange the non terminals in some sequence.

$S, A \text{ or } A_1 S$

→ two loops are used -

for ( $i = 1$  to  $n$ )

{

for ( $j = 1$  to  $i-1$ )

{

if ( $A_i \rightarrow A_j \gamma$ )

{

if ( $A_j \rightarrow \delta_1 \delta_2 \delta_3$ )

replace  $A_i \rightarrow \delta_1 \gamma \delta_2 \delta_3$

}

$A \rightarrow \text{sd}$  replaced by  $A \rightarrow \text{ad}$

→ Eliminate left recursion.

$A_i \rightarrow \text{non-terminal}$

$i = 2 \quad A_2 = A_2 = A$

$j = 1 \quad A_2 \rightarrow S \gamma$

$S \rightarrow A \alpha | b$

Now we need to replace it -  $| S \alpha \leftarrow A$

$A_2 \rightarrow S \delta$

$A \rightarrow A \alpha d$

$A \rightarrow b d$

14b

Now, we are getting the production rules as -

$$S \rightarrow A\alpha | b$$

$$A \rightarrow Ac$$

$$A \rightarrow Aad$$

$$A \rightarrow bd$$

now, we need to eliminate left recursion  
for the above grammar.

$$A \rightarrow \overset{\alpha}{A} \underset{\beta}{c} | \overset{\alpha}{A} \underset{\beta}{ad} | \overset{\beta}{bd} | \epsilon$$

$$S \rightarrow A\alpha | b$$

$$A \rightarrow \overset{\alpha}{bd} \underset{\beta}{A'} | A'$$

$$A' \rightarrow \overset{\alpha}{CA'} | \underset{\beta}{ad} A' | \epsilon$$

### Left factoring:

- It is a grammar transformation that is useful for producing a grammar suitable for predictive parsing or top-down parsing.
- When the choice between two alternative predictions of  $A$  ( $A \rightarrow \alpha\beta_1 | \alpha\beta_2$ ) is not clear then, the decision of using the production rules to obtain eight choice.
- In such a situation, we have to back track & check for other choice.
- Backtracking introduces overhead & this overhead needs to be eliminated on left factored.

Let a CFG be there -

$$S \rightarrow CAD$$

$$A \rightarrow ab/a$$

$$\omega = CAD$$

$S \xrightarrow{LHS} CAD$  } <sup>\*</sup> Sentential (string of terminal & non-terminal symbols)

$\xrightarrow{LHS} CAD$  Sentence (string of terminal symbols)

\*  $S \xrightarrow{LHS} CAD$

$$\xrightarrow{LHS} CAD$$

In left factoring,

$$A \rightarrow \alpha \beta_1 | \alpha \beta_2$$

e.g.:  $A \rightarrow ab/a$   $\downarrow$

left factored

$$\boxed{\begin{array}{l} A \rightarrow \alpha A' \\ A' \rightarrow \beta_1 | \beta_2 \end{array}}$$

$$A \rightarrow \alpha A' \quad (\because \alpha \rightarrow \alpha)$$

$$A' \rightarrow \beta_1 | \beta_2$$

A parse tree always starts from S symbol.

## Dangling if - else problem

(4)

Let the CFG be -

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

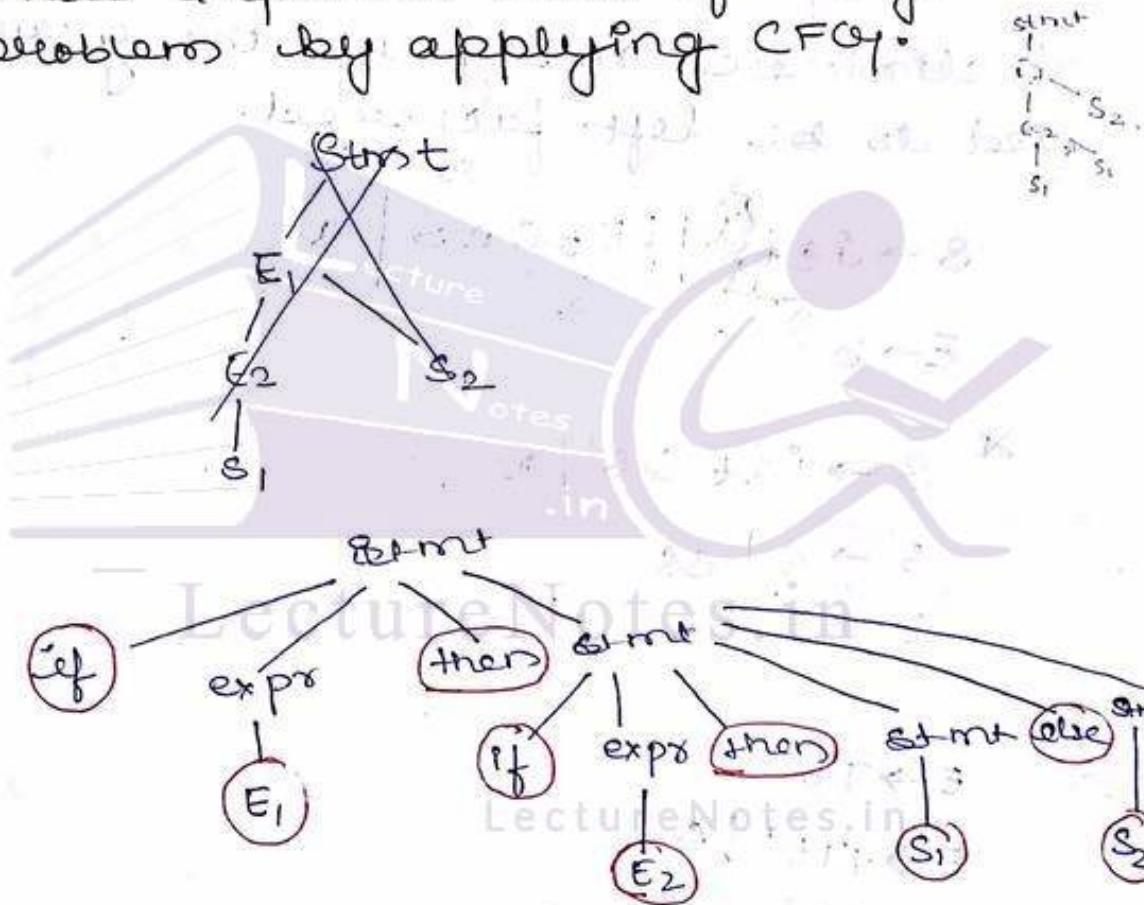
$\text{if expr then stmt else stmt}$

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

Eg if  $E_1$  then if  $E_2$  then  $S_1$  else  $S_2$

Draw a parse tree of the given problem by applying CFG.

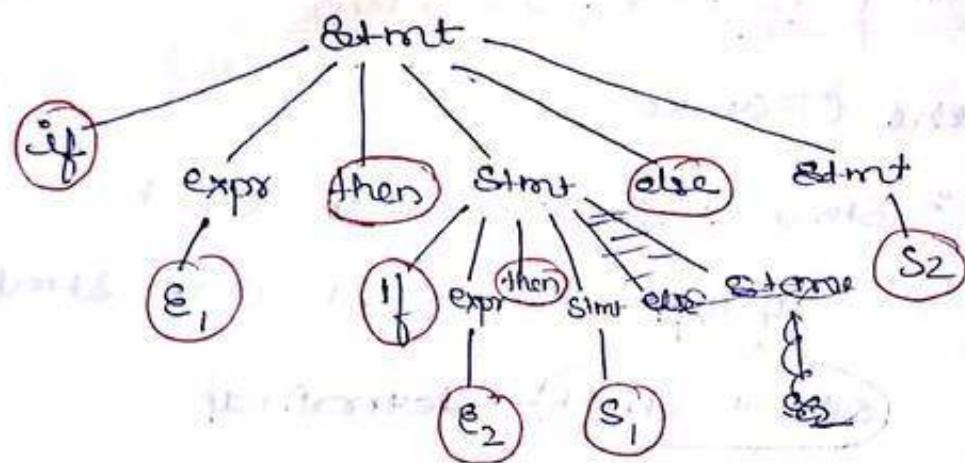
$\Rightarrow$



Parse tree contains every token generated by the grammar

terminal forms of PEGs are also called as leaves

(4)



∴ The above given grammar is ambiguous.  
To eliminate ambiguity the grammar need to be left factored.

$$S \rightarrow i e t s^* | i e t s^* s^* | a$$

$\alpha$

$$E \rightarrow b$$

$$S \rightarrow i e t s^* s^* | a$$

$$s^* \rightarrow g | es$$

$$E \rightarrow b$$

$$E \rightarrow T E'$$

$$E \rightarrow + T E' | G$$

$$T \rightarrow F T'$$

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

$$F \rightarrow C ( E ) | id$$

For the above CFG, [ $w = id + id * id$ ], generate a parse tree using leftmost derivation.



---

Mindfire is a software service provider,  
with unrelenting focus on optimal software development  
and delivery.

---



# *Compiler Design*

Topic:  
***Top Down Parsing***

Contributed By:  
***Sibananda Achari***  
Trident Academy of Technology

## Types of Parsing Techniques

02/01/14

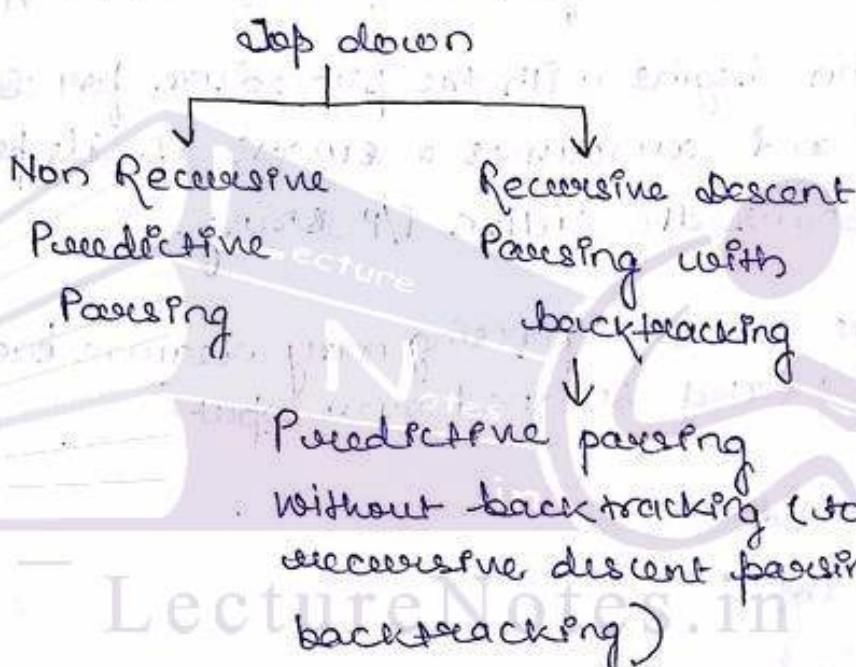
(51)

→ Top-down (leftmost derivation)

→ Bottom-up (rightmost derivation)

### Top down parsing:

- It generates parse tree starting from the root and creating nodes of the tree in pre-order.
- It can be viewed as finding leftmost derivation of an input string.



### Recursive Descent Parsing

For every non-terminal, a procedure is present and it needs to be called.

#### Algorithm

void A()

{

    Choose a A-production  $A \rightarrow x_1 x_2 \dots x_k$

    for ( $i=1$  to  $k$ )

        if ( $x_i$  is a non terminal)

            call procedure  $x_i()$ ;

(51)

(52)

else if ( $x_i$  = current I/P symbol)  
 advance to the next I/P symbol  
 else  
 /\* repeat error \*/

{

}

→ A recursive descent parsing consists of a set of procedures (one for each non terminal).

- Execution begins with the procedure for start symbol and announces a success if its procedure body scans the entire I/P string.
- Recursive descent parsing may require backtracking i.e. repeated scanning over input.

$$\text{Q } S \rightarrow cAd$$

$$A \rightarrow ab|a$$

$$\omega = Cad$$

- Algorithm matches the production  $S \rightarrow cAd$  with  $\omega = cad$  (given string).

$$x_1 \quad x_2 \quad x_3$$

When  $x_1 = A$  is encountered, it uses another algorithm to make use of the production  $A \rightarrow ab|a$ .

When  $a$  is encountered, the loop is entered again to encounter  $b$ .

(52)

But 'ab' doesn't match with 'd' in  $w = cad$ ,  
therefore, backtracking occurs.

Then the next prediction that  $P_S, A \rightarrow a$  is  
advanced so. It matches, then i's value is  
incremented and  $a_1$  is also matched.

→ Based on the concept of DFS (Depth First  
Search)

## Recursive Predictive Parsing

FIRST (calculated for  $\alpha$ )

FOLLOW (calculated for NT (non-terminal))

Eg - If  $A \rightarrow \alpha$ , then we need to calculate  
FOLLOW(A) and FIRST( $\alpha$ ).

→ FIRST( $\alpha$ )

In this set of terminals that are derived  
from  $\alpha$ . Gives the first terminal -

Eg -  $S \rightarrow cAd$

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

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

For finding out the FOLLOW of a non-terminal,  
RHS of production rule needs to be checked.

In the above example,  $\text{FOLLOW}(A) = \{d\}$

Gives the terminal that immediately follows the  
given non-terminal.

NOTE :

(Pg - 217 - 221)

FIRST( $\alpha$ ), where  $\alpha$  is any string of grammar symbols is defined as the set of terminals that begin strings derived from  $\alpha$ .

$\rightarrow$  If  $\alpha^* \Rightarrow G$  (i.e.  $\alpha$  produces  $G$  in one or more steps) then,

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

$\rightarrow$  If  $A^+ \Rightarrow C\gamma$

then

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

$$\text{FIRST}(A) = \text{FIRST}(C\gamma)$$

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

$\rightarrow$  Three steps to be used to compute FIRST( $x$ ) (where,  $x$  is a grammar symbol), applying the following rules until no more terminals or  $\epsilon$  can be added to any first set are:

Step-I

If  $x$  is a terminal,  $\text{FIRST}(x) = \{ x \}$

Step-II

If  $x$  is a NT and for  $x \rightarrow y_1 \dots y_n$

If  $y_1 \rightarrow G$ , then  $\text{FIRST}(x) = \text{FIRST}(y_2, y_3 \dots y_n)$

Step-III

If  $y_i = \epsilon$  for all  $i$ , then  $\text{FIRST}(x) = \{ \epsilon \}$

FIRST( ) and FOLLOW( )

08/01/14  
 (55)

$$A \rightarrow \alpha \cap \epsilon$$

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

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

&  $F \rightarrow CE | Id$  • find  $\text{FIRST}(CE)$  &

$$\text{FIRST}(Id).$$

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

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

$$\text{FIRST}(F) = \{\epsilon, Id\}$$

Q  $A \rightarrow C Ca$

$$\text{FIRST}(CcCa) = \{\epsilon, C\}$$

$$\text{FIRST}(A) = \text{FIRST}(\alpha) = \{\epsilon, C\}$$

Rule: if a long string is there -

$$A \rightarrow Y_1 Y_2 Y_3 \dots Y_n$$

then,  $\text{FIRST}(Y_1 Y_2 Y_3 \dots Y_n) = \{Y_1\}$

- If  $Y_1 \in \epsilon$ , then in R.E

$$Y_1 \in Y_2 Y_3 \dots Y_n = Y_2 Y_3 \dots Y_n$$

$$\text{FIRST}(Y_1 Y_2 Y_3 \dots Y_n) = \{\epsilon\}$$

- If  $Y_1 \dots Y_n \notin \epsilon$

then

$$\text{FIRST}(Y_1 \dots Y_n) = \{\epsilon\}$$

(55)

## 56 FOLLOW

- FOLLOW is calculated for a non-terminal only.
- It is the set of terminals,  $a$ , that appear immediately to the right of
- ⇒  $A$  is some sentential form.
- Let us assume, there exists a derivation,
- ⇒ If  $S \xrightarrow{*} \alpha A \beta \gamma$ , then  $\text{FOLLOW}(A) = \alpha$
- ⇒ If  $S \rightarrow \alpha A$ , then  $\text{FOLLOW}(A) = \{\$\}$
- ⇒ If  $\text{Exp} \rightarrow \text{Exp} * \text{Exp} | \text{Exp} + \text{Exp} | \text{id}$ ,  
let Exp be the start symbol  $S$   
then,  $\text{FOLLOW}(S) = \{\$\}$
- ⇒ If  $A \rightarrow \alpha B \beta$ , then everything in  
 $\text{FIRST}(\beta) - \epsilon$  is in  $\text{FOLLOW}(B)$ .  
i.e.  $\text{FOLLOW}(B) = \text{FIRST}(\beta) - \epsilon$

$\Rightarrow$  If  $A \rightarrow \alpha B$  or  $A \rightarrow \alpha B \beta$  and (57)

FIRST( $B$ ) contains  $\epsilon$  then everything  
in FOLLOW( $A$ ) is in FOLLOW( $B$ ).

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

$\rightarrow$  To find the FIRST and FOLLOW of the production rule, the grammar should be non ambiguous. Then left recursion should be eliminated if exist and the grammar should be left factored.

Ex:  $E \rightarrow E * E \mid E + E \mid (E) \mid id$

It is ambiguous as we cannot define whether we should go for \* or + for derivation.

So we need to resolve it.

$\therefore +$  has lowest precedence  
so we will go with -

$$E \rightarrow E^A \underset{\alpha}{+} T \mid P$$

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

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

Now, in the above new production we will check for left recursion.

(57)

(53)

∴ loop is created as the first symbol is repeated again and again as  $E \rightarrow E + T | T$ .

∴ The left recursion rule is applied -

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

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

~~from the above production -~~

$$T \rightarrow F T'$$

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

$$F \rightarrow (E) | \epsilon$$

Now, the above new production should be left factored.

The rule is -

$$A \rightarrow \alpha \beta, | \alpha \beta_2$$

In the above grammar, the no production is as above mentioned rule so no left factor is required.

\$ should be always at the end.

(54)

	<u>FIRST( )</u>	<u>FOLLOW( )</u>	(5)
E	(, id	), +, \$	\$
E'	+ , ε	), +	\$
T	(, id	+ , ), *	\$
T'	* , ε	+ , ), *\$	\$
F	(, id	* , +, )	\$

F C<sub>5</sub>Rd

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

`FIRST(id) = {id}`

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

(v) We will check for (E)  
as (d) do not satisfy  
the condition -  $A \rightarrow \neg A \wedge B$

$$E \xrightarrow{\alpha B} E'$$

If FIRST(B) is E,

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

$$E \xrightarrow{\alpha} T^B \epsilon^B$$

If FIRST( $E'$ ) contains  $\epsilon$   
then everything except  $\epsilon$   
is in FOLLOW( $E'$ ).

$$\textcircled{I} \quad e' \rightarrow +\overset{\times}{T} \overset{B}{E}' | e$$

$$\cancel{G' \rightarrow + T E' | C}$$

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

15

E C91d

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

$$T \rightarrow T^{\alpha} \beta^{\beta} \cdot \alpha B \beta$$

$$\text{FOLLOW}(T') = \text{FOLLOW}(T)$$

$$FOLLOW(F) = FIRST(T') - t$$

$$\text{IV} \quad T' \rightarrow *F T' \in$$



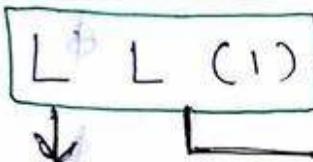
---

Mindfire is a software service provider,  
with unrelenting focus on optimal software development  
and delivery.

---

⑥①

## LL(1) Grammar



Scanning the input from left to right stands for 1 input symbol at each step to make parsing action description.

→ 2 types of parsing are three-left & right.

The grammar should be non ambiguous, left factored and do not have left recursion is called LL(1) grammar.

Parsing table is a 2D array where rows are defined by non terminal and columns are defined by terminal.

3-173729 1:60:03:03

317729 17

⑥②

# Construction of Recursive Predictive PT (6)

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

## Steps to build the Parsing Table

→ For each production  $A \rightarrow \alpha$  of the grammar, do the following -

- If  $\epsilon$  is in  $\text{FIRST}(\alpha)$ , then for each terminal in  $\text{FOLLOW}(A)$  add production rule  $A \rightarrow \alpha$  to the parser table (M $[A, B]$  for the non-terminal A & terminal B).

M [A, b].

- If  $\epsilon$  is in  $\text{FIRST}(\alpha)$  and  $\$$  is in  $\text{FOLLOW}(A)$ , then add the production rule  $A \rightarrow \alpha$  to the parser table.

M [A, \$].

62

- For each terminal  $a \in \Sigma$ , in  $\text{FIRST}(\alpha)$  add  $A \rightarrow \alpha$  to  $M[A, a]$ .

Q

$$S \xrightarrow{\alpha, \beta, \gamma} i E S S' | a$$

$$S' \xrightarrow{\alpha, \beta, \gamma} e S | \epsilon$$

$$E \rightarrow b$$

Find the  $\text{FIRST}()$  &  $\text{FOLLOW}()$  and then construct the predictive parsing table.

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

	$\text{FIRST}()$	$\text{FOLLOW}()$
S	i, a	e, \$
S'	e, ε	e, \$
E	b	

	$\epsilon$	e	t	a	b	\$
S	$S \xrightarrow{\epsilon} i E S S'$	$S \xrightarrow{\epsilon} e S$		$S \xrightarrow{\epsilon} a$		
S'		$S' \xrightarrow{\epsilon} e S$				$S \xrightarrow{\epsilon} a$
E			$E \xrightarrow{\epsilon} b$		$E \xrightarrow{\epsilon} b$	

63

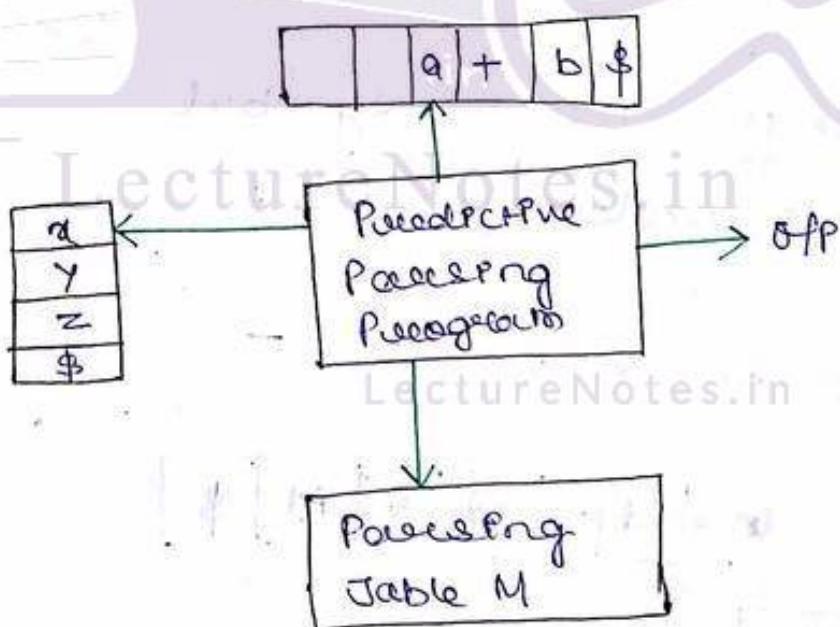
## Non-Recursive Predictive Parsing

13/01/14  
63

It can be built by maintaining a stack explicitly. The parser may mix / induces a left most derivation if  $w$  is the input that has been matched so far after the stack holds a sequence of grammar symbols  $\alpha$  such that

$$S \xrightarrow{*} w\alpha$$

The parser ~~scans~~ is table driven. It uses predictive parsing table  $M$ . The table driven predictive parser has an I/P buffer, a stack containing parser has an I/P buffer, a stack containing a sequence of grammar symbols, a parsing table and an output string.

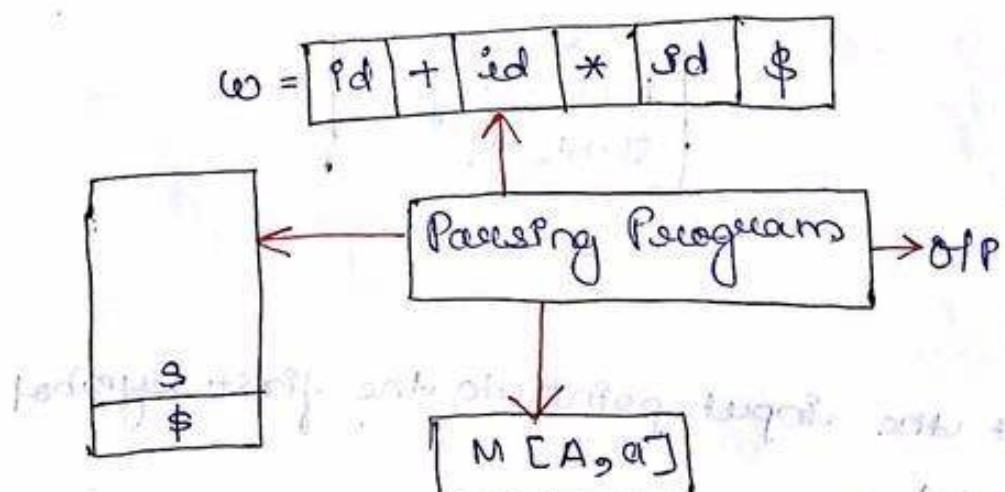


- I - Set the input point to the first symbol of  $w$ , where  $w \rightarrow I/P$  string provided

(64)

- II - Set  $X$  to the TOS
- III - while ( $X \neq \$$ )
- IV - If ( $X$  is a) pop the stack and advance the input pp.
- V - else if ( $X$  is a decimal) eval( $X$ )
- VI - else if ( $M[A, a] = X \rightarrow Y_1, Y_2, \dots, Y_k$ )  
off the specification  $X \rightarrow Y_1, Y_2, \dots, Y_k$
- VII - Push  $Y_k, Y_{k-1}, \dots, Y_1$  onto the stack such that  $Y_1$  is on the top
- VIII - }  
IX - ends while //end while

Q  $w = id + id * id$



(65)

Matched	Stack	G/P	O/P / Act Pos
	E \$	id + id * id \$	
	T E' \$	id + id * id \$	E → TE'
id	FT' E' \$	id + id * id \$	T → FT'
	id T' E' \$	id + id * id \$	f → id
	T' E' \$	+ id * id \$	matched
	E' \$	+ id * id \$	T' → E
id +	+ T E' \$	+ id * id \$	E' → + TE'
	TE' \$	id * id \$	matched
	FT' E' \$	id * id \$	T → FT'

The CFC for E is -

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

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

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

$$M[E, id] =$$

$$E \rightarrow TE'$$

$$\begin{array}{|c|} \hline E = X \\ \hline \$ \\ \hline \end{array}$$

id + id	id T' E' \$	id * id \$	F → id
	T' E' \$	* id \$	matched
id + id *	* FT' E' \$	* id \$	T' → * FT'
	FT' E' \$	* id \$	E' → + TE' matched
id + id * id	id T' E' \$	id \$	F → id
	T' E' \$	\$	matched
	E' \$	\$	T → E
	\$	\$	E' → E

## ⑥ Types Of Errors

Common programming errors can occur at different levels of compilation are.

### → Lexical Errors :

Eg - misspelling of identifiers, keywords / operators and missing codes around the text string.

### → Syntactic Errors :

Eg - missing braces / extra braces, misplaced semicolon.

### → Semantic Errors :

Eg - It includes mismatch between operators and operands.

### → Logical Errors :

Eg - It can be anything from incorrect reasoning on the part of the programmer.

→ The error handler in a parser has to report the presence of errors clearly and accurately, recover from each error quickly to detect subsequent errors and at minimal overhead to process the program.

### → The different error recovery strategies

Ques -

- Panic Mode Recovery
- Phrase Level Recovery



---

Mindfire is a software service provider,  
with unrelenting focus on optimal software development  
and delivery.

---



# *Compiler Design*

Topic:  
***Bottom Up Parsing***

Contributed By:  
***Sibananda Achari***  
Trident Academy of Technology

## → Bottom - Up Parsing

- Operator Precedence Parsing ex.
- LR Parsers

• 15-01-14

- Bottom up Parsing corresponds to construction of a parse tree for an input string.
- Beginning at the leaves and working up towards the root.
- Bottom - up parser are also known as shift-reduce parsers.
- Shift - reduce parsers are of 2 types -
- Operator precedence parser : It is applied to small class of grammar.
  - LR Parsers : It scans the I/P string from left to right and generates a rightmost derivation tree.

(69)

(68)

## Q Why Shift-reduce?

→ Bottom up parsing can be thought as a process of reducing a string  $w$  to the start symbol  $S$  of the grammar.

At each reduction step, a specific substring matching the body of the production is replaced by the non terminal at the head of the production.

There are 4 types of actions performed by Shift-reduce parser.

→ Shift: It shifts the next T/P symbol on the top of the stack.

→ Reduce: The right end of the string to be reduced must be present at the TOS. Locate the left end of the string within the stack and replace it by non terminal.

→ Accept: After successful completion of parsing, the parser announces accept.

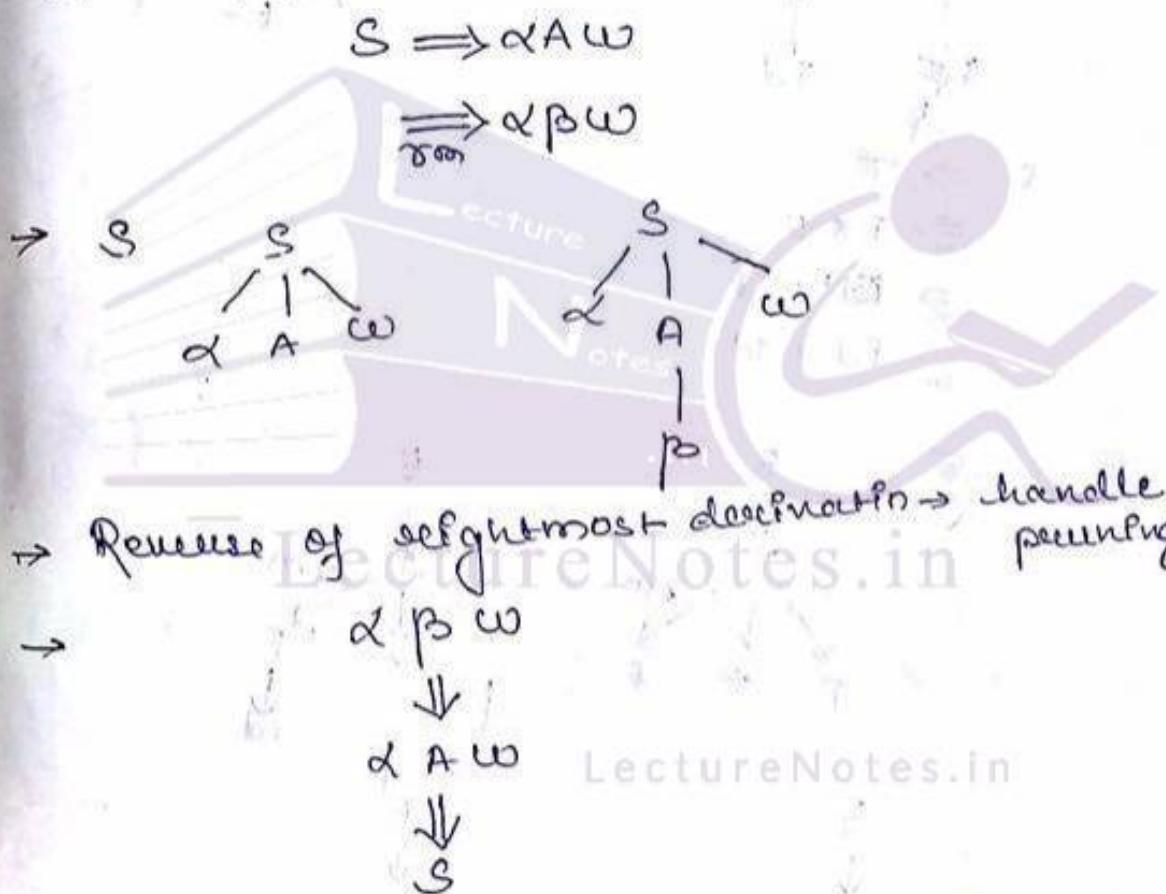
→ Error: In case of syntax error, calls for an error recovery routine.

(69)

## Handle Punning

(69)

- Handle Punning is related to the replacement of a handle within (it is a substring which matches the body of the production) by a non terminal present at the head of the production.
- It can be obtained from leftmost derivation in reverse. Formally, if  $S \rightarrow^* \alpha A \omega$  i.e.



→ Reverse of leftmost derivation  $\rightarrow$  handle punning

$$\begin{array}{c} \alpha \beta \omega \\ \downarrow \\ \alpha A \omega \\ \Downarrow \\ S \end{array}$$

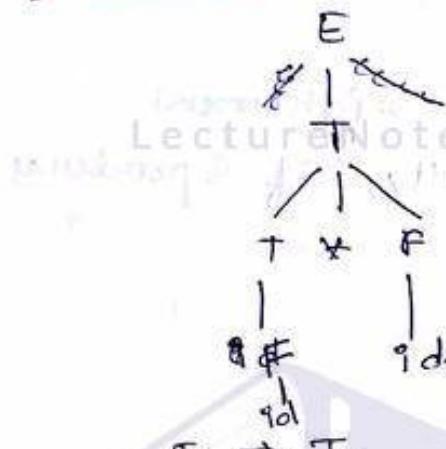
→ It identifies a position following  $\alpha$  where the handle of the string  $\alpha \beta \omega$  can be identified and  $\beta$  is replaced by the non terminal  $A$ .

(67)

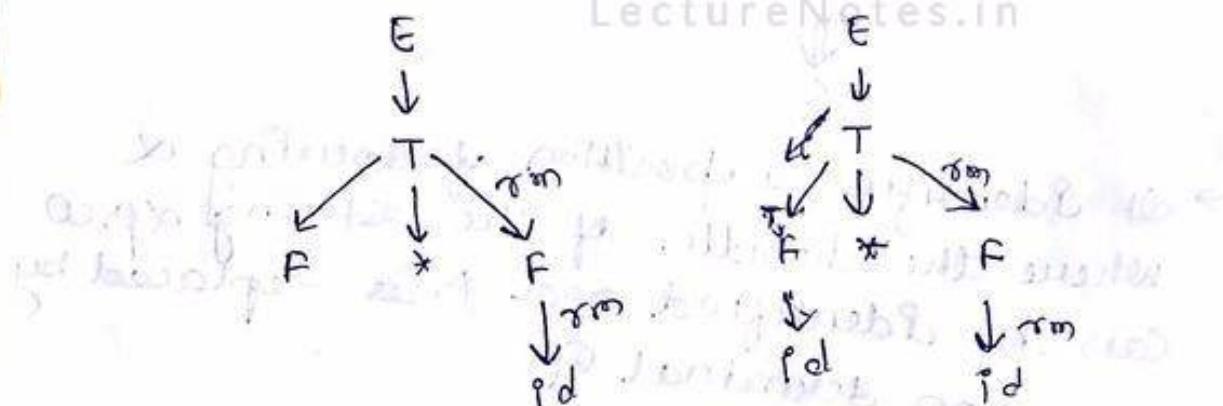
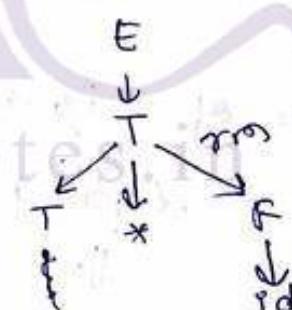
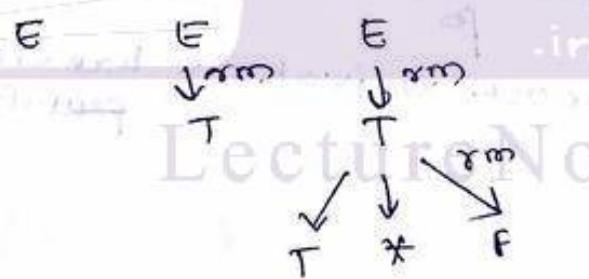
(70)

Input

id \* id

Handle $E \rightarrow E + T \mid T$  $T \rightarrow T * F \mid F$  $F \rightarrow (E) \mid id$ Right-most derivation of id \* id

$E \Rightarrow T$   
 $\Rightarrow T * F$   
 $\Rightarrow T * id$   
 $\Rightarrow id * id$



(70)

<u>Output</u>	<u>Handle</u>	<u>Replacing production</u> (71)
$id_1 * id_2$	$id_1$	$F \rightarrow id$
$F * id_2$	$F$	$T \rightarrow F$
$T * id_2$	$id_2$	$F \rightarrow id$
$T * F$	$T * F$	$T \rightarrow T * F$
$T$	$T$	$E \rightarrow T$
$E$		
<u>Using the Input string <math>id * id \\$</math>, define the steps of action performed by shift-reduce parser. Using a stack, initially has endmarker \$.</u>		
<u>Stack</u>	<u>Output</u>	<u>Action</u>
\$	$id_1 \$$	$F \rightarrow id$
$id_1 \$$	$id_1 * id_2 \$$	$T \rightarrow F$
$id_1 * id_2 \$$	$id_1 * id_2 \$$	$F \rightarrow id$
$F * id_2 \$$	$T * id_2 \$$	$T \rightarrow T * F$
$T * id_2 \$$	$T * id_2 \$$	$E \rightarrow T$
$T * F \$$		

(72) Intermediate Code Generation 20/01/14  
→ In compiler, we generate Intermediate code of target code.

→ Intermediate code generates the target code.

### Benefits of Intermediate Code generation

- It is machine independent. As it is independent, a compiler for a different machine can be used by attaching it for the new machine.
- A machine independent code optimizer can be applied to the Intermediate code representation.

3 address code :

- There are different forms of Intermediate code representation. One of them is three address code.
- It is represented as per the following format.  
$$x := y \text{ op } z$$
- There will be always one operator & two operands on right hand side. And it permits one operand on LHS.
- Here,  $x, y, z$  can be constant, temporary variables, memory addresses or variable names.
- Here, temporary variables are generated by the compiler.

(72)

Statements	These address code format <span style="border: 1px solid black; border-radius: 50%; padding: 2px;">73</span>
→ Assignment statements having binary operators.	$x := y \text{ op } z$
→ Assignment statements having unary operator	$x = \text{op } z$
→ Copy Statement / Assignment statement	$x = y$
→ unconditional jump statements	goto $\pm$ (concrete L is the label)
→ conditional jump statements	if $x \text{ relOp } y$ goto L ↳ relational operator
→ procedure call / function call statement	param $x_1$ , param $x_2$ , : param $x_n$ call P, n → no. of arguments Procedure name
→ indexed statement Eg: $x = arr[z]$ $arr[z] = x$	$x = y[i]$ $x[i] := y$
→ pointer / Address statement	$x := \& y$ $*x := y$ $p := *y$

(74)

Eg:  $M = x * y * z$

Find its three address code.

$\Rightarrow$  The three address code is

$$t_1 := x * y$$

$$t_2 := t_1 * z$$

$$M := t_2$$

### Implementation of Three Address Code

$\rightarrow$  Three address statement is an abstract form of intermediate code.

$\rightarrow$  In compiler, this statements can be implemented as records with fields for the operator and the operands.

$\rightarrow$  3 such representations are there -

- Quaduples
- Triples
- Indirect Triples

Eg

Find the three address code (TAC) of the following statement:

$$a = \underbrace{b * -c}_{(1)} + \underbrace{b * -c}_{(2)}$$

⇒ The TAC for the above is -

- 0  $t_1 = -c$
- 1  $t_2 = b * t_1$
- $a = t_1 + t_2$
- 2  $t_3 = -c$
- 3  $t_4 = b * t_3$
- 4  $t_5 = t_2 + t_4$
- 5  $a = t_5$

### Quadruple representation:

Index	operator	arg 1	arg 2	result
0	u-	$-c$		$t_1$
1	*	$b$	$t_1$	$t_2$
2	u-	$-c$		$t_3$
3	*	$b$	$t_3$	$t_4$
4	+	$t_2$	$t_4$	$t_5$
5	assign	$t_5$		$a$

### Triple Representation:

Index	operator	arg 1	arg 2
0	u-	$c$	
1	*	$b$	(0)
2	u-	$c$	
3	*	$b$	(2)
4	+	(1)	(3)
5	assign	$a$	(4)

### Indirect Triple Representation

Table-I.

	op	arg 1	arg 2
21	u-	$c$	
22	*	$b$	(21)
23	u-	$c$	
24	*	$b$	(23)
25	+	(22)	(24)
26	assign	$a$	(25)



---

Mindfire is a software service provider,  
with unrelenting focus on optimal software development  
and delivery.

---

76

Table - II

	Statement
0	21
1	22
2	23
3	24
4	25
5	26

Q Find the triple representation of the following statements:

$$\triangleright x[i] = y$$

$$\triangleright x = y[i]$$

$\Rightarrow$  The TAC for the above is -

index	op	arg 1	arg 2
0	[ ] =	x	i
1	assign	(0)	(0)y
2			

index	op	arg 1	arg 2
0	[ ] =	y	i
1	assign	x	(0)

# Operator Precedence Parsing

23/01/14 (77)

## → Precedence Relations

- bottom-up parsers for a large class of CFG can be easily developed using operator grammars.
- operator grammars have the property that no production right side is empty or has two adjacent nonterminals.
- It should not produce  $\epsilon$ -production.
- This property enables the implementation of efficient operator-precedence parsers. These parser rely on the following 3 precedence relation:

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

- These operator precedence relations allow to delimit the handles in the right sentential forms :  $<$  marks the left end,  $=$  appears in the interior of the handle, and  $>$  marks the right end.

Q What are the disadvantages of Operator Precedence Parsing?

- It is hard to handle tokens like - signs which has two different precedences
- It is applied only to small class of

(77)

(75)

## Operator grammar.

→ Relationship between an grammar for the language to be parsed and the operator precedence parser is very tedious.

Thus

### Ways to determine Precedence Relations

→ There are 2 common ways -

- Using operator precedence relations
- Construct unambiguous grammar for the language which reflects both precedence and associativity

### Using operator precedence relation

→ The idea behind it is to delimit the handle of a right sentential form with < marking left end = appearing inserted to the handle if any and > marking the right end.

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

It is not in the sentential form as no non-terminal is there.

so we need to delimit the

handle.

→ id has higher precedence over +.

(76)

points more on the bottom of the page

$$\text{Let } \stackrel{a_i}{12} \stackrel{a_{i+1}}{>} + \stackrel{a_{i+2}}{<} \stackrel{E}{\text{ }} + \stackrel{E}{\text{ }}$$

(79)

→ Operator precedence parser belongs to the shift reduce parser.

→ \$ has less precedence as compared to any terminal.

$$\therefore \rightarrow \$ \stackrel{\text{id}_1}{\geq} + \stackrel{\text{id}_2}{\geq} * \stackrel{\text{id}_3}{\geq} \stackrel{\text{id}_4}{\geq} \$$$

$$\text{relation} \Rightarrow \$ \leq \text{id}_1 \geq + \leq \text{id}_2 \geq * \leq \text{id}_3 \geq \$$$

Step-I: Scan the string from the left end until the leftmost > is encountered.

Step-II: Then scan backward over any = until < is encountered.

Step-III: The handle contains everything to the left of the first > and to the right of <. Including any surrounding non-terminal.

meaning  
To check whether, two adjacent non terminals appear in a right sentential form.

$$\rightarrow \$ \leq \text{id}_1 \geq + \leq \text{id}_2 \geq * \leq \text{id}_3 \geq \$$$

$$\rightarrow \$ E - + < \text{id}_2 \geq * < \text{id}_3 \geq \$$$

$$\rightarrow \$ E + E * < \text{id}_3 \geq \$$$

$$\rightarrow \$ E + E * E \$ \leftarrow \text{Handle}$$

(79)

(80)

 $\rightarrow \$ + * \$$ 

Now, we need to establish a relationship for the above -

 $\rightarrow \$ < \cdot + < \cdot * > \$$ 

$\uparrow$   $< \cdot >$  is scanned

$\rightarrow \$ < \cdot + \cdot > \$$

$\uparrow$   $\therefore E * E$  should be precedence first

 $\rightarrow \$ E + E \$$ 
 $\rightarrow \$ < \cdot + \cdot > \$$ 

$\uparrow$   $< \cdot >$  is scanned so move back

 $\rightarrow \$ E \$$ 

The terminals identified are

+ \* \$

Now, we will construct a table where no. of rows = no. of columns.

	Pd	+	*	\$
Pd	$\Rightarrow$	$\Rightarrow$	$\Rightarrow$	$\Rightarrow$
Pd +	$< \cdot$	$= \Rightarrow$	$< \cdot$	$\cdot >$
Pd *	$< \cdot$	$\cdot >$	$= \Rightarrow$	$\cdot >$
Pd \$	$< \cdot$	$< \cdot$	$< \cdot$	<del>\$</del>

(80)

Q pd +-\* / ^ \$

(81)

→ Operator Precedence Parsing Algorithm

Initialize. Set  $ip$  to point to the first symbol of w\$.

Repeat : Let  $a$  be the top stack symbol, and  $b$  be the symbol pointed to by  $ip$ .

If  $s$  is on the top of the stack (accept) and  $ip$  points to  $s$  then  
else

Let  $a$  be the top terminal on the stack, and  $b$  the symbol pointed to by  $ip$ .

If  $a < b$  or  $a = b$  then  
push  $b$  onto the stack  
advance  $ip$  to the next input symbol.

else if  $a > b$  then

repeat  
pop the stack  
until the top stack terminal is replaced by  $<$   
by the terminal most recently popped

else error ()

end

(81)

Steps of Derivation	Stack	Input String	Action
	\$	id <sub>1</sub> + id <sub>2</sub> * id <sub>3</sub> \$	input to the algorithm is string w\$ and precedence table
E → Pd <sub>1</sub>	Pd <sub>1</sub> \$	4 id <sub>2</sub> * id <sub>3</sub> \$	accept shift
E → Pd	Pd <sub>2</sub> + \$	+ id <sub>2</sub> * id <sub>3</sub> \$	reduce
		* id <sub>2</sub> * id <sub>3</sub> \$	a = id → b = +
		* id <sub>3</sub> \$	shift
		* Pd <sub>3</sub> \$	reduce
			shift

The production rule is -

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

- popped  
- action

Sp → pointer pointing to an input string

$$a = \$ \quad b = id$$

$$a < b$$

larger prec - pop  
less prec - push

E → id	id <sub>3</sub> * + \$	Pd <sub>3</sub> \$	shift
E → E*E	* + \$	\$	reduce
E → E+E	+ \$	\$	reduce
	\$	\$	reduce

NOTE:

(83)

Let  $\Theta_1$  and  $\Theta_2$  be two operators.

If the relationship between the two operators has same precedence then we can write it as

$$\Theta_1 \circ \Theta_2 \text{ or } \Theta_2 \circ \Theta_1$$

→ If the operators have same precedence i.e.  $\Theta_1 = \Theta_2$  then we can write it as  $\Theta_1 \circ \Theta_2 \text{ or } \Theta_2 \circ \Theta_1$  if they are left associative.

Eg - + - \* /

→ If the operators have same precedence but right associative then

$$\Theta_1 < \Theta_2 \text{ and } \Theta_2 < \Theta_1$$

Eg - exponentiation (<sup>1</sup>)

→ If  $\Theta_1 \circ \Theta_2$  and they are left associative then the precedence relation can be written as -

$$\Theta_1 \circ \Theta_2 \text{ and } \Theta_2 < \Theta_1$$

→ If  $\Theta_1 < \Theta_2$  and the relationship is in combination of left and right associativity:

$$\Theta_1 < \Theta_2 \text{ and } \Theta_2 \circ \Theta_1$$

(83)

(8A)

Eg - \* and +

→ Make  $\circ$ , is having less precedence than  $\text{id}$

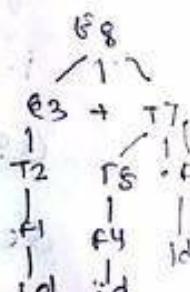
i.e.  $\circ < \text{id}$  and  $\text{id} > \circ$ ,

$\$ < \text{id}$  and  $\text{id} > \$$

$( < \text{id} \text{ and } \text{id} > )$

$( < ( \text{ and } ) > )$

$( = )$



Q  $(\text{id}_1 * (\text{id}_2 \uparrow \text{id}_3) - \text{id}_4 / \text{id}_5)$

(Pg-158-173)

→ Operator Precedence Relation from associativity and Precedence.

Algorithm for Constructing Precedence function

→ Create function  $f_a$  for each grammar terminal  $a$  and for the end of string symbol  $\$$ .

→ Partition the symbols into groups so that

$f_a$  and  $f_b$  are in the same group

if  $a = b$  (there can be symbols in the same group even if they are not connected by this relation).

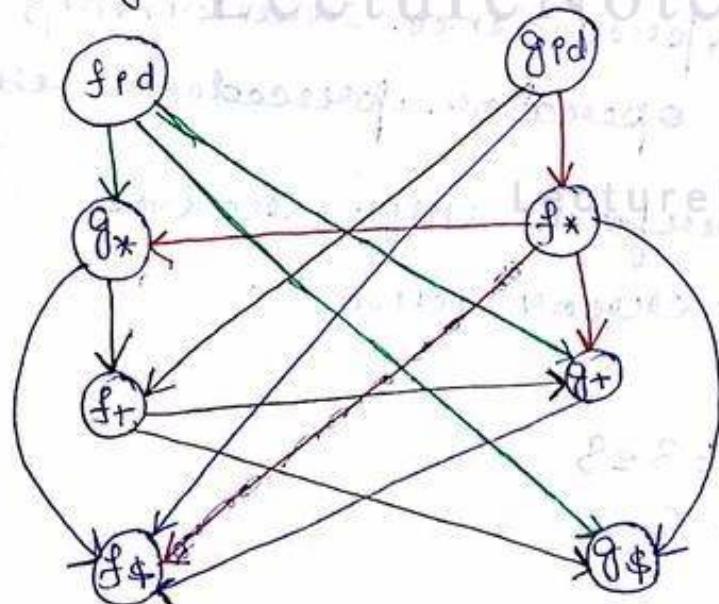
(8B)

- Create a directed graph whose nodes are in the groups next to each symbols a & b.
- do : place an edge from the group of  $g_b$  to the group of  $f_a$ , if  $a < b$ , otherwise if  $a > b$  place an edge from the group of  $f_a$  to that of  $g_b$
- If the constructed graph has a cycle then no precedence functions exist. When there are no cycles collect the length of the longest paths from the groups of  $f_a$  and  $g_b$  respectively.

30/1/14

$$\text{Eq} - E \rightarrow E+E \mid E * E \mid \text{id}$$

O/P string - id + id \* id &  
Every a should be mapped to b -



After & no outgoing way is there, therefore we calculate the path from the starting point to \$.

85



---

Mindfire is a software service provider,  
with unrelenting focus on optimal software development  
and delivery.

---

(86)

The precedence table -

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

Computation Of First Terminal

Last Terminal, Leading and

Trailing of a non-terminal

for a given non-production rule

Step-I: Compute the first terminal & last terminal.

Step-II: Compute leading and trailing.

Step-III: From the leading and the trailing, establish precedence relationship and build operator precedence table.

Step-IV: finally perform shift, reduce accept or reject action.

Q

$$S \rightarrow rCtSeS$$

$$S \rightarrow rCtS$$

$$S \rightarrow a$$

$$C \rightarrow b$$

for more information visit [www.lecturenotes.in](http://www.lecturenotes.in)

→ The grammar is ambiguous and this grammar defines precedence relationship between a pair of terminals that can have two precedence relationship.

The First Terminal can be calculated for a given non terminal over the terminals that can be first in a string derived from that non terminal.

Similarly, for each non terminal those terminals that are the last terminal in a string derived from that non terminal.

There are 2 rules to calculate first terminals and the last terminals for a given non terminal.

First Terminal	Last Terminal
<p>Rule 1:</p> $A \xrightarrow{*} \gamma a \delta$ where, $\gamma$ is either $\epsilon$ or a single non terminal then add $a$ to the first last	

(33)

Rule 1: First terminal A is equal to a such that

$$a) A \xrightarrow{*} \gamma a \delta$$

where,  $\gamma$  is either  $\epsilon$ /non terminal.

Add a to the first terminal of  $A=a$  such that  $A \xrightarrow{*} \gamma a \delta$  and  $\gamma$  is  $\epsilon$ /single non terminal.

Rule 2:

Last terminal of  $A=a$  such that

$$A \xrightarrow{*} \gamma a \delta$$

and,  $\delta$  is  $\epsilon$ /single non terminal

First Terminal	Last Terminal
s i a	e, t, a
c b	b

$$\bullet C \xrightarrow{*} \underbrace{b}_{\gamma a \delta}$$

$$\bullet S \xrightarrow{*} \underbrace{a}_{\gamma a \delta}$$

$$\bullet S \xrightarrow{*} \underbrace{i \gamma a t}_{\gamma a \delta} \underbrace{s}_{\delta}$$

$$\bullet S \xrightarrow{*} \underbrace{i c t}_{\gamma a \delta} \underbrace{s}_{\gamma a \delta}$$

$$\bullet S \xrightarrow{*} \underbrace{\epsilon}_{\gamma a \delta} \underbrace{c t s e}_{\gamma a \delta} \underbrace{s}_{\delta}$$

$$\bullet S \xrightarrow{*} \underbrace{i c t s e}_{\gamma a \delta} \underbrace{s}_{\gamma a \delta}$$

$$\bullet S \xrightarrow{*} \underbrace{i c t s e}_{\gamma a \delta} \underbrace{s}_{\gamma a \delta}$$

(34)

## Algorithm to Construct Operator Precedence Relation

31/1/14 (89)

Rule 1: leading (A) =  $\{ \alpha \mid A \xrightarrow{*} \alpha \delta \}$

where,  $\delta \rightarrow$  either  $\epsilon / \alpha$  & single non-terminal

Rule 2: trailing (A) =  $\{ \alpha \mid A \xrightarrow{*} \gamma \alpha \delta \}$

where,  $\delta \rightarrow$  either  $\epsilon / \alpha$  & single non-terminal.

Rule 3: If  $\alpha$  is in the leading of A and there is a production

$$B \rightarrow A\alpha$$

then  $\alpha$  is in the leading of B.

Rule 4: Compute the leading & trailing for the first & last terminal.

Rule 5: The leading is use to establish  $\leftarrow$  relationship and the trailing helps to calculate  $\rightarrow$  relationship.

$a \leftarrow b$  - leading

$a \rightarrow b$  - trailing

Rule 6: Set  $\$$  is less than  $a$  i.e.  $\$ < a$  in leading (S).

• Set  $b \rightarrow \$ + b$  in trailing (S).

where  $S \rightarrow$  start symbol in production  
•  $\rightarrow$  terminal symbol and non-terminal

(89)

Rule 7: Choose a production.

Let  $A \rightarrow x_1 x_2 \dots x_n$

For  $i = 1, 2, \dots, n-1$

$\rightarrow$  if ( $x_i$  and  $x_{i+1}$  == terminal)

then [Set  $x_i = x_{i+1}$ ]

$\rightarrow$  if ( $i < n-2$  &  $x_i$  and  $x_{i+2}$  ==

terminal and  $x_{i+1}$  == non terminal)

then Set  $x_i = x_{i+2}$

$\rightarrow$  if ( $x_i$  == terminal and  $x_{i+1}$  == non terminal)

then Set  $x_i < a$  &  $a$  is leading ( $x_{i+1}$ )

$\rightarrow$  if ( $x_i$  == non terminal and  $x_{i+1}$  == terminal)

then Set  $a > x_{i+1}$  &  $a$  is trailing ( $x_{i+1}$ ).

g  $S \rightarrow i C + S e S$

$S \rightarrow i C + S$

$S \rightarrow a$

$C \rightarrow b$   
Leading

Trailing

First Terminal	Last Terminal
$i, a$	$t, e, a$
$b$	$b$

Once table satisfies for a particular relationship then use another value of  $i$ .

(10)

Now, we need to establish relationship between \$ and S. So for that - 91

(i) Leading (\$) -

$$\{ \{ p, a \} \} \quad \$ <_o i \quad \$ <_o a$$

	i	b	t	a	e	\$
i			≡	✗		
b			•>			
t	<		✗	✗	≡	•>
a	✗		✗		•>	•>
e			✗	<	•>	•>
\$	<			✗		

(ii) Trailing (\$) = {t, e, a}

$$t > \$ \quad e > \$ \quad a > \$$$

(iii) Leading (C) = {b, i}

$$\$ <_o b$$

Let the production be -

$$S \rightarrow i C t S e S$$

$$A \rightarrow x_1 x_2 x_3 x_4 x_5 x_6$$

$$if i = 1 \text{ to } 5 \quad (\because n = 6 \quad \therefore D-1=5)$$

$$\rightarrow i=1 \quad x_1=i \quad x_2=C$$

$$1 \leq i \leq 4 \quad x_1=i \quad x_3=t \quad x_2=C$$

$$\rightarrow i=2 \quad x_2=C \quad x_3=b \quad x_4=e$$

$$\rightarrow i=3 \quad x_3=t \quad x_4=s \quad x_5=e$$

$$\rightarrow i=4 \quad x_4=s \quad x_5=e \quad x_6=s$$

$$\rightarrow i=5 \quad x_5=e \quad x_6=s$$

(92)

$$S \rightarrow i C t S$$

A  $x_1 x_2 x_3 x_4$

for  $i = 1 \text{ to } 3$

$$i=1 \quad x_1=i \quad x_2=C \quad x_3=t$$

$$i \div t$$

$$i=2 \quad x_2=C \quad x_3=t$$

$$t \geq b > t$$

$$i=3 \quad x_3=t \quad x_4=S$$

$$* < i, a$$

$$\begin{array}{c} S \rightarrow | a \\ A \rightarrow x_1 \end{array}$$

Stack	Relation	Input	Action
\$	<	fbfbfaea\$	shift
i \$	<	bfbfbfaea\$	shift
b i \$	>	tfbfaea\$	reduce
i \$	=	tfbfaea\$	shift
t i \$	<	fbfaea\$	shift
t i f	<	baea\$	shift
b t i f	>	aea\$	reduce
t i f	=	aea\$	shift
t i f	<	ea\$	shift
a t i f	>	a\$	reduce
t i f	=	a\$	shift
e t i f	<	\$	shift
et i f	>	\$	reduce
e t i f	>	\$	reduce
t i f	>	\$	reduce
i f	>	\$	reduce

(92)

$\cdot \text{t} \$$	$\Rightarrow$	\$	reduce	(93)
$\cdot i \$$	$\Rightarrow$	\$	reduce	
$\cdot \$$		\$	accept	

### Assignment

Resolve the ambiguity in the given grammar.  
 Build operator precedence parsing table  
 using leading and trailing and then perform  
 shift-reduce parsing for input - 1bt1bt a &  
 using the non ambiguous precedence rule.

Shift reducing Parsing = Bottom up Parsing

### Conflicts During Shift-Reduce Parsing

- There are CFGs for which shift-reduce parser cannot be used.
- Stack contents and the next input symbol may not decide action:
  - Shift/reduce : Whether make a shift operation or a reduction.
  - Reduce/reduce : The parser cannot decide which of several reductions do make.
- If a shift-reduce parser cannot be used for a grammar, that grammar is called a non-LR(0) grammar.
  - ↓
  - left-right scanning
  - leftmost derivation
  - K look ahead.
- An ambiguous grammar cannot be LR grammar.



---

Mindfire is a software service provider,  
with unrelenting focus on optimal software development  
and delivery.

---



# *Compiler Design*

Topic:  
***LR Parsing Algorithm***

Contributed By:  
***Sibananda Achari***  
Trident Academy of Technology

## 2. LR Parsers

→ There are three types of LR parsing are there -

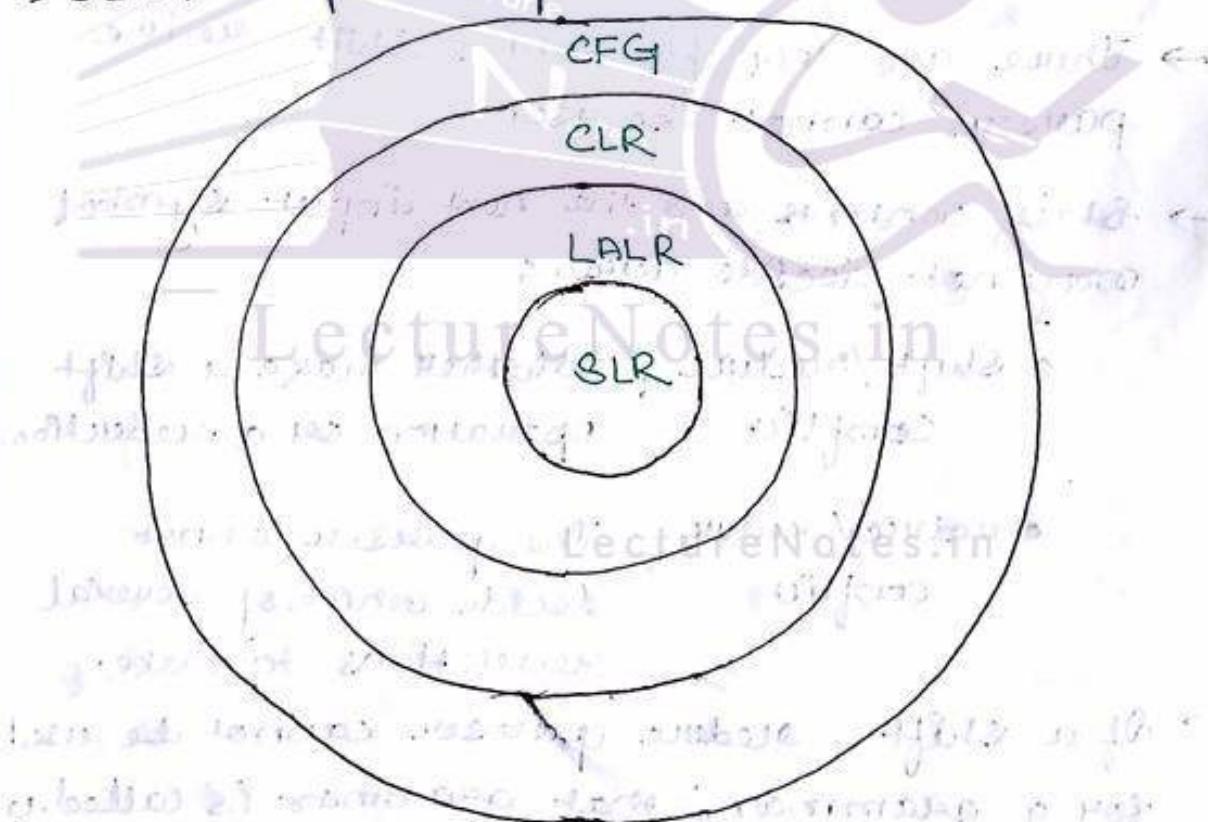
- SLR
- CLR
- LALR

→ CLR : Canonical LR parser

most commonly used LR parser

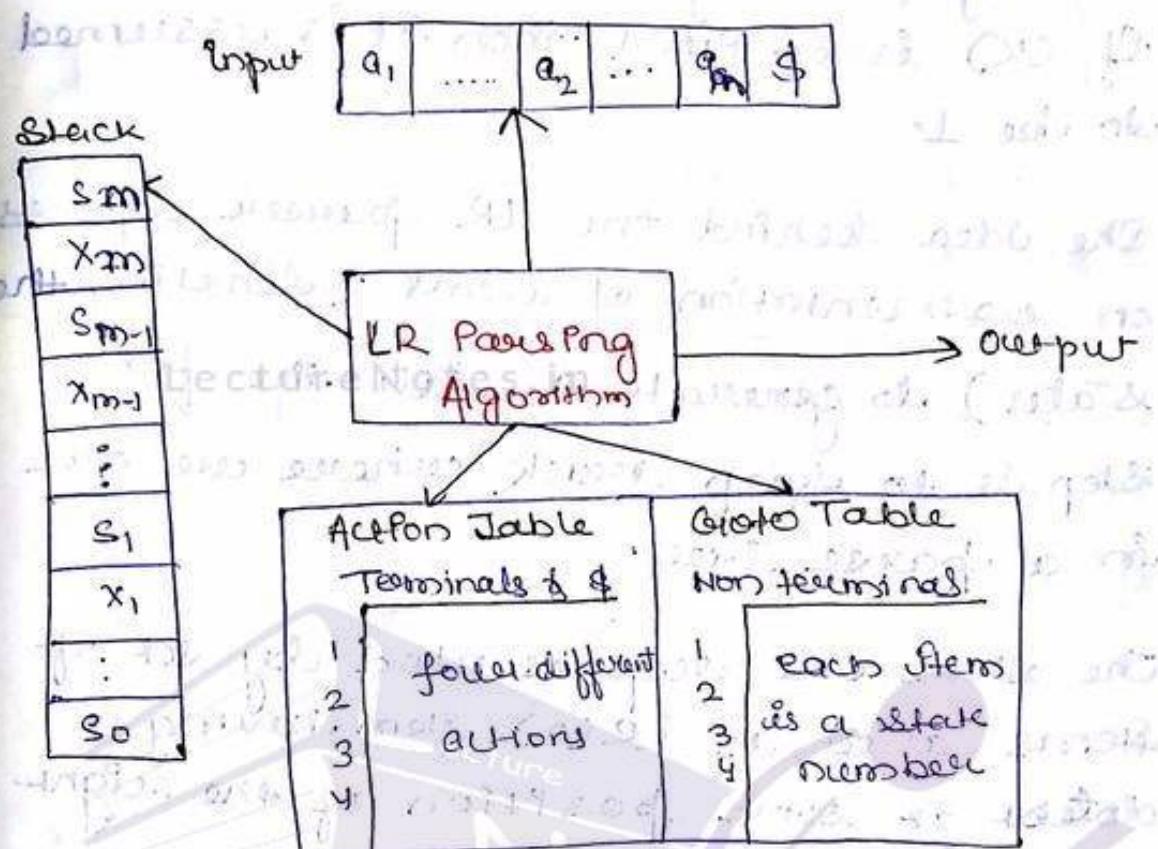
→ LALR : look ahead LR parser

→ SLR : simple LR parser



# LR Parsing Algorithm

03/02/14 (95)



- The LR parsers are table driven parsers very much like non recursive LL parsers.
- The concept of LR( $k$ ) parsing defines scanning the input from left to right and constructing a leftmost derivation in reverse.
- And  $k$  denotes the no. of input symbols of look ahead that are used in making parsing decisions.
- The general case of  $k$  is taken as  $k=0$  or  $k=1$ .

⑥

- Generally the case considered for making parsing decision is  $K \leq 1$ .  
If  $(K)$  is omitted then it is assumed to be 1.
- The idea behind the LR parser depends on representation of items (denoting the states) to generate a state the first step is to keep track where we are in a parse tree.
- The states are represented by set of items called as LR(0). Item having a dot(.) at some position of the right side.

→ Eg :

$$A \rightarrow aBb$$

Possible LR(0) items

$$A \rightarrow \cdot aBb$$

$$A \rightarrow a \cdot B b$$

$$A \rightarrow aB \cdot b$$

$$A \rightarrow aBb \cdot$$

→ They are also known as items in short.

If there is a production  $A \xrightarrow{\text{LR}(0)} C$  then it generates only one item, set of items / items i.e.  $A \rightarrow \cdot$

⑦

- This representation indicates how much of a production we have seen at a given point in the parsing process. (97)
- To construct collection of sets of LR(0) items also known as canonical LR(0) collection is the basis for constructing SLR parser.

We define an augmented grammar.

$G'$  is  $G$  with a new production rule

$$S' \rightarrow S$$

where:  $S' \rightarrow$  new start symbol.

NOTE: The new starting symbol indicates the parser when it should stop parsing and generates a message accept of the input.

Acceptance occurs only when the parser is about to reduce

$$S' \rightarrow S$$

→ We take the help of 2 types of functions.

• CLOSURE

• GOTO

(97) 22

(93)  $\rightarrow$  CLOSURE

If  $I$  is a set of items of a grammar  
G then closure of  $I$  is the set of  
items constructed from  $I$  ~~are~~ by  
using two rules:

Rule 1:

Initially add every item in  
 $I$  to closure ( $I$ ).

Rule 2:

If  $A \rightarrow \alpha \cdot B \beta$  and there is a  
production rule  $B \rightarrow \gamma$  then  
add item  $B \rightarrow \gamma$  to the  
closure ( $I$ ) if not already  
there.

$\rightarrow$  GOTO ( $I, x$ )

GOTO ( $I, x$ ) is define to be the closure  
of a set of all items  $[A \rightarrow \alpha \cdot x \beta]$   
such that  $[A \rightarrow \alpha \cdot x \beta]$  is in  $I$ .

which says that  $x$  can be non terminal  
or terminal but  $B$  is always a  
non terminal.

(94)

## Actions of a LR parser

(99)

→ shift : Shifts the next input symbol and the state  $s$  onto the stack.

$$(S_2 \ x_1 \ s_1 \dots x_m \ s_m, a_i, a_{i-1} \dots a_n \ S) \rightarrow$$

$$(S_2 \ x_1 \ s_1 \dots x_m \ s_m \ a_i \ s, a_{i-1} \dots a_n \ S)$$

→ reduce :  $A \rightarrow \beta$  (or  $\alpha$  where  $\alpha$  is a production number)

→ pop  $\beta$  ( $=\tau$ ) items from the stack.

→ others push  $A$  and  $s$  where  $s = \text{goto}[S, A]$

$$(S_2 \ x_1 \ s_1 \dots x_m \ s_m, a_i, a_{i-1} \dots a_n \ S) \rightarrow$$

$$(S_0 \ x_1 \ s_1 \dots x_{m-\tau} \ s_{m-\tau}, A \ s, a_i \dots a_n \ S)$$

→ output is the reducing production

→ reduce  $A \rightarrow \beta$ .

→ accept : Parsing successfully completed

→ error : Parser detected an error (an empty entry in the action table)

(Pg - 240)

## (SLR) Parsing Tables for Expression Grammar

$$E \rightarrow E-T \quad G_1: E \rightarrow E+E \mid E * E \mid (E) \mid id$$

$E \rightarrow T$       ambiguity need to be removed

$T \mid$

(99)

(100)

$$(1) E \rightarrow E + T$$

$$(2) E \rightarrow T$$

$$(3) T \rightarrow T * F$$

$$(4) T \rightarrow F$$

$$(5) F \rightarrow (E)$$

$$(6) F \rightarrow id$$

Now we need to have augmented grammar.

G1.

$$E' \rightarrow E = I$$

$$E \rightarrow E + T = I$$

Next step is to build collection of a set of LR(0) items.

$$\text{closure } (I) = \text{closure } (\{E' \rightarrow E\})$$

$$I_0 : \{ E' \rightarrow \cdot E \}$$

$$E \rightarrow \cdot E - T$$

$$E \rightarrow \cdot T$$

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

$$T \rightarrow \cdot F$$

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

$$F \rightarrow \cdot id$$

Initial state of DFA  
(I0)

} closure

(100)



---

Mindfire is a software service provider,  
with unrelenting focus on optimal software development  
and delivery.

---

I<sub>1</sub>:  $\text{goto } (I_0, E) = \text{closure } (\{ E' \rightarrow \cdot E \})$  (10)

$$= \{ E' \rightarrow E \cdot, E \rightarrow E \cdot + T \}$$

prefixes Viable prefixes

Viable prefixes are prefix sentential forms that can appear on the stack of a shift reduce parser. It is always possible to add terminals to the end of viable prefixes to obtain a right sentential form.

I<sub>2</sub>:  $\text{goto } (I_0, T) = \text{closure } (\{ E \rightarrow \cdot T \})$

$$= \{ E \rightarrow T \cdot, T \rightarrow T \cdot * F \}$$

I<sub>3</sub>:  $\text{goto } (I_0, F) = \text{closure } (\{ T \rightarrow \cdot F \})$

$$= \{ T \rightarrow F \cdot \}$$

I<sub>4</sub>:  $\text{goto } (I_0, C) = \{ F \rightarrow C \cdot E \},$

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

$$E \rightarrow \cdot T,$$

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

$$T \rightarrow \cdot F,$$

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

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

⋮

I<sub>5</sub>:  $\text{goto } (I_0, \text{id}) = \{ F \rightarrow \text{id} \cdot \}$

$$\text{(b2) } I_0 \rightarrow E' \rightarrow E \xrightarrow{\text{GOTO}(I_0, E)} \left( \begin{array}{l} E' \rightarrow E^* \\ E \rightarrow E \cdot + T \end{array} \right) I_1 \xrightarrow{\text{GOTO}(I_1, T)} \left( \begin{array}{l} E \rightarrow E \cdot + T \xrightarrow{\text{GOTO}(I_1, T)} \\ E \rightarrow E \cdot + T \end{array} \right) I_2$$

$$\begin{array}{c}
 E \rightarrow \cdot T \\
 \xrightarrow{\text{GOTO}(I_2, T)} \left( \begin{array}{l} E \rightarrow T^* \\ T \rightarrow T \cdot * F \end{array} \right) I_2 \\
 \xrightarrow{\text{GOTO}(I_2, T \cdot * F)} \left( \begin{array}{l} T \rightarrow T \cdot * F \\ T \rightarrow T \cdot * F \end{array} \right) I_3
 \end{array}$$

$$\begin{array}{c}
 T \rightarrow \cdot F \\
 \xrightarrow{\text{GOTO}(I_3, F)} \left( \begin{array}{l} T \rightarrow F^* \\ T \rightarrow T \cdot F \end{array} \right) I_4
 \end{array}$$

$$\begin{array}{c}
 F \rightarrow \cdot (E) \\
 \xrightarrow{\text{GOTO}(I_4, E)} \left( \begin{array}{l} F \rightarrow (\cdot E) \\ E \rightarrow E + T \end{array} \right) I_5
 \end{array}$$

$$\begin{array}{c}
 E \rightarrow E + T \\
 \xrightarrow{\text{GOTO}(I_5, E + T)} \left( \begin{array}{l} E \rightarrow E + T \\ E \rightarrow E + T \end{array} \right) I_5
 \end{array}$$



---

Mindfire is a software service provider,  
with unrelenting focus on optimal software development  
and delivery.

---



# *Compiler Design*

Topic:  
***SLR Parser***

Contributed By:  
***Sibananda Achari***  
Trident Academy of Technology

Construction of the Canonical LR(0) collection.

$\rightarrow C \text{ is } \frac{1}{2} \text{ closure } (S^* \rightarrow \cdot S))^\beta$

repeat the followings until no more set of LR(0) items can be added to C

for each  $l$  in C and each grammar symbol  $x$ .

if  $\text{goto}(l, x)$  is not empty and not in C

add  $\text{goto}(l, x)$  to C

$\rightarrow$  goto function is a DFA on the sets in C.

Augmented

$$\begin{aligned} & C(E' \rightarrow \cdot E) \\ & = E' \rightarrow \cdot E \end{aligned} \quad \left. \begin{array}{l} \text{kernel items} \\ \text{non-kernel items} \end{array} \right\}$$

I<sub>0</sub>:

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

$$E \rightarrow \cdot T$$

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

$$T \rightarrow \cdot F$$

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

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

LectureNotes.in

non-kernel items

goto (I<sub>0</sub>, E)

I<sub>1</sub>:

$$E' \rightarrow E \cdot$$

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

Kernel Items

(10A)  $\text{goto } (I_0, T)$

$I_2 :$   $E \rightarrow T$   
 $T \rightarrow T * F$

} Kernel Items

$\text{goto } (I_0, F)$

$I_3 :$   $T \rightarrow F$  - KI

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

$I_4 :$   $F \rightarrow \text{id}$  - KI

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

$I_5 :$

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

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

$I_6 :$

$E \rightarrow E + \cdot T$  - KI

$T \rightarrow \cdot T * F$

$T \rightarrow \cdot F$

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

$F \rightarrow \cdot (E)$

} Non Kernel Item

$\text{goto } (I_2, *)$

$I_7 :$

$T \rightarrow T * \cdot F$  - KI

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

$F \rightarrow \cdot (E)$

} Non KI

(10)

105  
goto ( $I_5$ ,  $\epsilon$ )

$I_8$ :

$F \rightarrow (E^*) - KI$

goto ( $I_6$ ,  $T$ )

$I_9$ :

$E \rightarrow E + T^* \rightarrow KI$

$T \rightarrow T^* \star F \rightarrow \text{no } KI$

goto ( $I_7$ ,  $F$ )

$I_{10}$ :

$T \rightarrow T * F \rightarrow KI$

goto ( $I_8$ ,  $)$ )

$I_{11}$ :

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

### Constructing SLR Parsing Table

- Construct the canonical collection of sets of LR(0) items for  $G^*$ .  $C \leftarrow \{ I_0, \dots, I_n \}$ .
- Create the parsing action table as follows.

- If  $a$  is a terminal  $A \rightarrow \alpha \cdot a \beta$  is in  $I_i$  and  $\text{goto } (I_i, a) = I_j$ , then  $\text{action}[r, a]$  is shift  $s$ .

- If  $A \rightarrow \alpha \cdot$  is in  $I_i$ , then  $\text{action}[r, a]$  is reduce  $A \rightarrow \alpha$  for all  $a$  in  $\text{FOLLOW}(A)$ , where  $A \neq S'$ .

- If  $S' \rightarrow s \cdot$  is in  $I_i$  then  $\text{action}[r, \$]$  is accept.

105

(106)

- If any conflicting actions are generated by the rules, the grammar is not SLR(1).

### Parsing Table

State	action										goto
	id	+	*	(	)	\$	E	T	f		
0	S4	S6		S5			1	2	3		
1		S6					Accept				
2		S2	S7		S2	S2					
3		S4	S4		S4	S4					
4	S\$			S5			8	2	3		
5		S6	S6		S6	S6					
6	S\$			S5			9	3			
7	S4			S5							10
8		S6			S11						
9		S11	S7		S1	S1					
10		S13	S3		S3	S3					
11		S18	S5		S5	S5					

I Action [1, \$] = accept

$$A \rightarrow \alpha$$

$$A \rightarrow \alpha \cdot \beta$$

↑

II Action [0, id] = S4

I0 and goto(I0, id) = Iy

Ans.

III  $I_0$  and goto  $CIO$ ,  $C = I_5$   
 Action  $[0, C] = S5$

Reduction

$$A \rightarrow \alpha.$$

$$A \neq E'$$

I  $E \rightarrow T.$

$I_2$  and Action  $[2, T] =$

follow  $(E) = (, \$)$

Action  $[2, , ] = \gamma_2$

Action  $[2, \$] = \gamma_2$

II  $T \rightarrow F.$

follow  $(T) = (, ;, +, \$)$

Action  $[3, +] = \gamma_4$

Action  $[3, , ] = \gamma_4$

Action  $[3, \$] = \gamma_4$

### SLR(1) Grammar

→ An LR parser using SLR(1) parsing tables for a grammar  $G_1$  is called as the SLR(1) parser for  $G_1$ .

→ If a grammar  $G_1$  has an SLR(1) parsing table, it is called SLR(1) grammar (or SLR grammar is short).

(108)

→ Every SLR grammar is unambiguous, but every unambiguous grammar is not a SLR grammar because it may generate shift-reduce conflict or reduce-reduce conflict.

Pg - 255

Ex - 4.4.2

### Shift/reduce and reduce/reduce conflicts

- If a state doesn't know whether it will make a shift operation or reduction for a terminal, we say that there is a shift/reduce conflict.
- If a state doesn't know whether it will make a reduction operation using the production rule  $i$  or  $j$  for a terminal, we say that there is a reduce/reduce conflict.
- If the SLR parsing table of a grammar  $G$  has a conflict we say that grammar is not SLR grammar.

### Canonical LR Parser

- A LR(1) item is

$$A \rightarrow \alpha \cdot \beta \gamma$$

where  $\alpha$  is the look-ahead of the LR(1) stem.

- ( $\alpha$  is a terminal or end marker)

- (109)
- To avoid some of invalid deductions, the states need to carry more information.
  - Extra information is put into a state by including a terminal symbol as a second component in an item.
  - Such an object is called LR(1) Item.
    - It refers to the length of the second component.
    - The look-ahead has no effect in an item of the form  $[A \rightarrow \alpha \cdot \beta, a]$ , where  $\beta$  is not  $\epsilon$ .
    - But an item of the form  $[A \rightarrow \alpha \cdot a]$  calls for a deduction in  $A \rightarrow \alpha$  only if the next input symbol is  $a$ .
    - The set of such  $a$ 's will be a subset of FOLLOW( $A$ ), but it could be a proper subset.
  - closure ( $I$ ) is :  $C$  where  $I$  is a set of LR(1) items
    - every LR(1) item in  $I$  is in closure ( $I$ ).
    - If  $A \rightarrow \alpha \cdot B \beta \cdot a$  in closure ( $I$ ) and  $B \rightarrow \gamma$  is a production rule of  $G$ . Then  $B \rightarrow \gamma \cdot b$  will be in the closure ( $I$ ) for each terminal  $b$  in FIRST( $Ba$ ). (109)



---

Mindfire is a software service provider,  
with unrelenting focus on optimal software development  
and delivery.

---

(11D)

### Goto operation

If  $I$  is a set of LR(1) items, and  $x$  is a grammar symbol (terminal or non-terminal), then  $\text{goto}(I, x)$  is defined as follows:

If  $A \rightarrow \alpha \cdot X \beta \cdot a \in I$

then every item is closure

### Construction of the Canonical LR(1) Collection

#### Algorithm

$C$  is  $\{$  closure  $\{S^1 \rightarrow S \cdot S\}\}$ .

repeat the following until no more set of LR(1) items can be added to  $C$ .

for each  $I$  in  $C$  and each grammar symbol  $x$

if goto( $I, x$ ) is not empty.

then add all items in goto( $I, x$ ) to  $C$ .

end if

(II) search in  $C$  for a first occurrence of  $a$  in  $S^1$ .

if it is found then add  $a$  to  $S^1$  and remove  $a$  from  $S^1$ .

if it is not found then add  $a$  to  $S^1$  and remove  $a$  from  $S^1$ .

if it is not found then add  $a$  to  $S^1$  and remove  $a$  from  $S^1$ .

if it is not found then add  $a$  to  $S^1$  and remove  $a$  from  $S^1$ .

if it is not found then add  $a$  to  $S^1$  and remove  $a$  from  $S^1$ .

calculate the FIRST( $C$ ), FOLLOW( $S$ ), FIRST( $S$ ), FOLLOW( $C$ ).

$$\text{I}_0 : \quad S' \rightarrow S, \$ \quad \Rightarrow \quad S' \rightarrow S, \$ \text{ is of the form } \\ A \rightarrow a \cdot X\beta, a \text{ where } \beta \text{ is } e \\ S \rightarrow C \cdot C, \$ \\ C \rightarrow c \cdot C, \$ \\ C \rightarrow d$$

$$\text{So, FIRST}(C \cdot) = \text{FIRST}(\$) \\ = \$$$

### LR(1) item (contd.)

13/02/14

→ When  $\beta$  (in the LR(1) item  $A \rightarrow \alpha \cdot \beta, a$ ) is not empty, the look-ahead does not have any effect.

→ When  $\beta$  is empty ( $A \rightarrow \alpha \cdot, a$ ), we do the operations reduction by  $A \rightarrow \alpha$  only if the next input symbol is  $a$  (not for any terminal in FOLLOW( $A$ )).

### Construction of the canonical LR(1) collection

#### Algorithm

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

repeat the followings until no more set of LR(1) items can be added to  $C$ .

for each  $I$  in  $C$  and each grammar symbol  $x$

if  $\text{goto}(I, x)$  is not empty and not in  $C$ .

add  $\text{goto}(I, x)$  to  $C$ .

goto function is a DFA on the sets in  $C$ .

(112)

I  $\rightarrow$   $S' \rightarrow S$

$S' \rightarrow CC$

$C \rightarrow CC$

$C \rightarrow d$

$\Rightarrow I_0: \text{closure} \{ \{ S' \rightarrow S, \$ \} \} =$

$(S' \rightarrow S, \$)$

$(S \rightarrow C \cdot C, \$)$

$(C \rightarrow C \cdot C, C/d)$

$(C \rightarrow \cdot d, C/d)$

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

$(S' \rightarrow S \cdot , \$)$

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

$(S \rightarrow C \cdot C, \$)$

$(C \rightarrow C \cdot C, \$)$

$(C \rightarrow \cdot d, \$)$

$I_3: \text{goto}(I_0, C)$

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

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

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

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

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

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

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

(112)

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

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

$C \rightarrow \cdot cc, \$$

$C \rightarrow \cdot d, \$$

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

$C \rightarrow d \cdot, \$$

$I_8 : \text{goto}(I_3, c)$

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

$I_9 : \text{goto}(I_6, c)$

$C \rightarrow cc \cdot, \$$

14/02/14

### Look ahead LR Parser

LALR parser can be created from CLR parsing techniques based on the idea of combining the core states.

$I_3 \rightarrow$

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

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

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

$I_6 \rightarrow$

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

$C \rightarrow \cdot cc, \$$

$C \rightarrow \cdot d, \$$

$I_{36} \rightarrow c \cdot c, c/d / \$$

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

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

(114)

### Core States

→ For each core (the 1<sup>st</sup> component) present among the set of LR(1) items, find all the sets having same core and replace the sets by their union.

I<sub>4</sub>-

$$C \rightarrow \underline{d} \cdot, c/d$$

I<sub>7</sub>-

$$C \rightarrow \underline{d} \cdot, \$$$

⇒ I<sub>47</sub>-

$$C \rightarrow \underline{d} \cdot, c/d, \$$$

I<sub>8</sub>-

$$C \rightarrow \underline{cc} \cdot, c/d$$

I<sub>9</sub>-

$$C \rightarrow \underline{cc} \cdot, \$$$

⇒ I<sub>89</sub>-

$$C \rightarrow \underline{cc} \cdot, c/d, \$$$

In total we have 7 sets of states after LALR parsing which was before 10 states due to CLR parsing.

\$

(M)

~~Good Recovery~~

	Action			Carrot	
	c	d	\$	s	c
1	s3	s4		1	2
2	s6	s7			5
3	s3	s4			8
4	r3	r3		r1	
5					
6	s6	s7			9
7			r3		
8	r2	r2		r2	
9					



---

Mindfire is a software service provider,  
with unrelenting focus on optimal software development  
and delivery.

---



# *Compiler Design*

Topic:

***Error Recovery In Predictive Parsing***

Contributed By:

***Sibanya Achary***

Trident Academy of Technology

## (116) Error Recovery in Predictive Parsing

Error Recovery first do the stack of a table driven predictive parser. Since, it makes explicit the terminals and non-terminal that the parser hopes to match with the remainder of the input. An error is detected during predictive parsing when the terminal on the top of the stack does not match the next input symbol or when the non-terminal A is on the top of stack, a terminal 'a' is the next input symbol and in the parsing table  $M[A, a]$  is empty.

## LL Parsing Technique

→ To handle error detected, the symbols in the first and follow are used as the synchronizing tokens and for the empty cells the tokens are represented as synch (synchronizing tokens).

If the parser looks up entry  $M[A, a]$  and finds that it is blank then the input symbol 'a' is skipped.

(116)

- (17)
- If the entry is synch then the non-terminal on the TOS is popped to resume parsing.
  - If a token on the TOS does not match the input symbol then we pop the token from the stack.

	FIRST( $\cdot$ )	FOLLOW( $\cdot$ )
E	(, pd	), synch \$
E'	+ , E	), +
T	(, pd	+ , ),
T'	* , E	+ , )
F	(, pd	* , + , )

	pd	+	*	(	)	*	\$
E	$E \rightarrow TE'$				$E \rightarrow TE'$	synch	synch
E'		$E' \rightarrow +TE'$				$E' \rightarrow E$	$E' \rightarrow E$
T	$T \rightarrow FT'$	synch		$T \rightarrow FT'$	synch	synch	
T'		$T' \rightarrow E$	$T' \rightarrow *FT'$		$T' \rightarrow E$	$T' \rightarrow E$	
F	$F \rightarrow (pd)$	synch	synch	$F \rightarrow (E)$	synch	synch	

Stack	Input	Remark
E \$	) id * + Pd \$	M [E, ]) = empty, skip
E \$	id * + id \$	
TE' \$	Pd * + id \$	
FT'E' \$	id * + Pd \$	F → Pd
T'E' \$	* + Pd \$	
*FT'E' \$	* + id \$	
FT'E' \$	+ id \$	M [f, +] = synch F popped
T'E' \$	+ Pd \$	
E' \$	+ id \$	
+TE' \$	+ Pd \$	
TE' \$	id \$	
FT'E' \$	Pd \$	F → Pd
@T'E' \$	Root \$	
T'E' \$	\$	
E' \$	\$	
\$	\$	

## LR Error Recovery Techniques

17/02/14

- An LR Parser will detect errors when it consults the parsing action table and finds an error entry.

- All empty entries in the action table (119)  
are error entries.
- Errors are never detected by consulting the goto table.
- An LR parser will announce error as soon as there is no valid continuation for the scanned input position.
- A CLR parser will recede after announcing an error.
- However, SLR and LALR will make several reductions before announcing an error.
- All the LR parsers will never shift an erroneous input symbol ~~however~~ onto the stack.
- There are 2 modes of recovery -
- Panic mode recovery
  - Phased level recovery
- Panic mode recovery:
- It scans down the stack until a stage is left with goto on a particular non-terminal is found.
- $\text{goto } [S, A] = \text{action}$
- Discard 0 or more input symbols until a symbol 'b' is found in the FOLLOW(A).

(120)

- The parser stack non-terminal A and the state goto [S, A] and set resumes normal parsing.

### Please Level Recovery:

- Each entry in the action table is marked with specific error routine.
- An error routine inserts or deletes symbols onto the stack or input for the following cases -
- \* Missing Parentheses - missing left/right parenthesis
  - \* Missing Operators
  - \* Missing Operands.

### Error Recovery in Operator Precedence Parsing:

- There are 2 points in parsing process at which an operator precedence parser can discover syntactic errors.
- If no precedence relation holds between the terminal on the top of the stack and between the current input.
  - If a handle is been found but there is no production with the handle as a right side.

(120)

## Handling Error Reduction:

(121)

→ The error checker for reductions need to check the following:

- If  $+, -, *, /$  is reduced, it checks that non terminals appear on both the sides. If not, it issues a message missing operand.
- If id is reduced, it checks that there is no <sup>non</sup> terminal to the right or left. If there is, it generates missing operation.
- If parenthesis is reduced, it checks that there is a non terminal between parenthesis. If not, it can cause no expression inside parenthesis.

## YACC (Yet Another Compiler Compiler)

→ It is a LALR parser generator.

→ It takes ambiguous grammar and resolves conflict.

### Reduce - Reduce Conflict

→ By choosing conflicting production, listed <sup>top</sup>, in the YACC specification.

### Shift/Reduce Conflict

→ It is resolved in favour of shift.

(122)

(122)

- YACC resolves shift/reduce conflict by attaching precedence and associativity to each production involved in a conflict as well as to each terminal involved in conflict.
- It must choose between shifting each input symbol and reducing by production  $A \rightarrow \alpha$ .
- YACC reduces in 2 cases -

#### Case-I

If the precedence of the production is greater than that of the input symbol  $a$ .

#### Case-II

If the precedences are same and the associativity is left. Otherwise, there is a shift.

- Normally, the precedence of the production is taken to be same as that of its right most terminal.

Eg

$$E \rightarrow E + E \mid E * E \mid (E) \mid \text{Id}$$

To construct the parsing table we need to augment the grammar.

(12)

So the augmented grammar is -



---

Mindfire is a software service provider,  
with unrelenting focus on optimal software development  
and delivery.

---

$E' \rightarrow E$

$E \rightarrow E+E \mid E * E \mid CE) \mid Pd$

(123)

I<sub>0</sub>:

$E' \rightarrow E$

$E \rightarrow \cdot E+E$

$E \rightarrow \cdot E * E$

$E \rightarrow \cdot (E)$

$E \rightarrow \cdot Pd$

LectureNotes.in

→ We can create LR-parsing tables for ambiguous grammar -

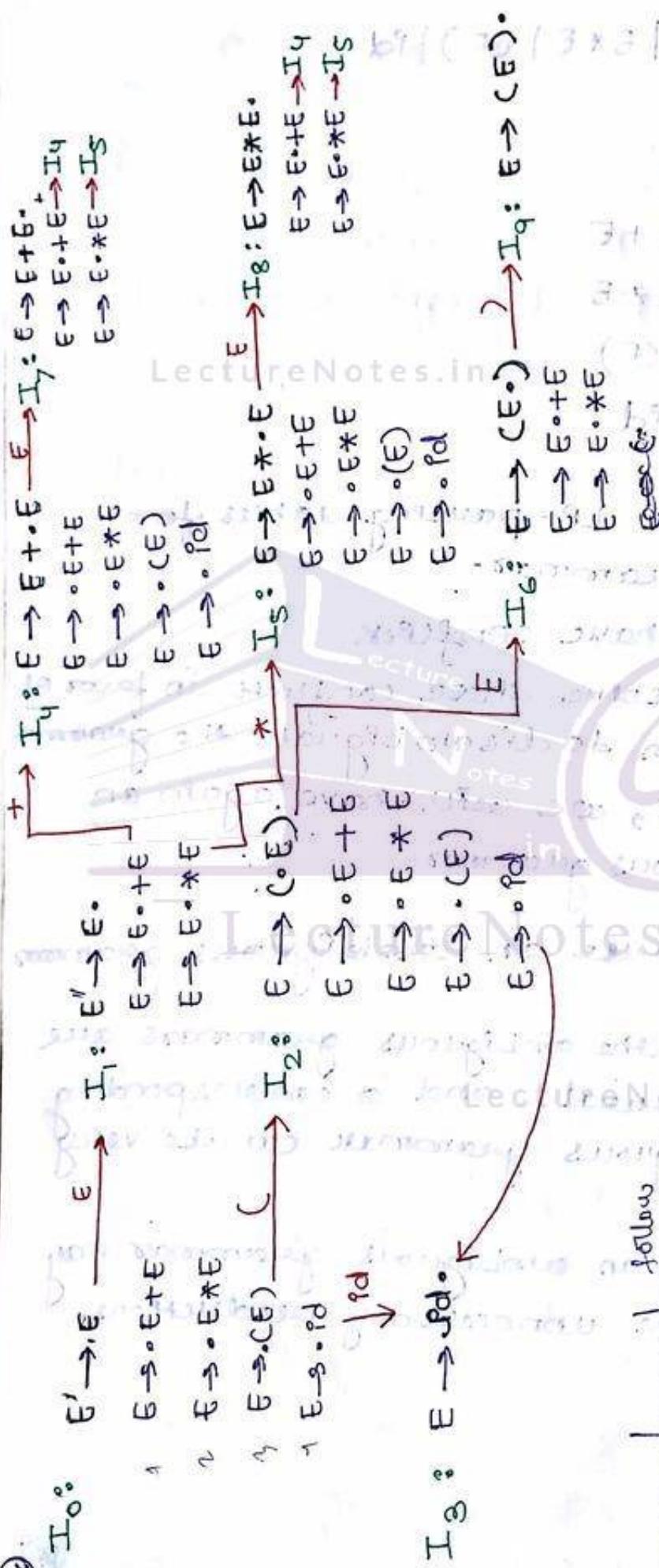
- They will have conflicts
- We can resolve these conflicts in favor of one of them to disambiguate the grammar.
- At the end, we will have again an unambiguous grammar.

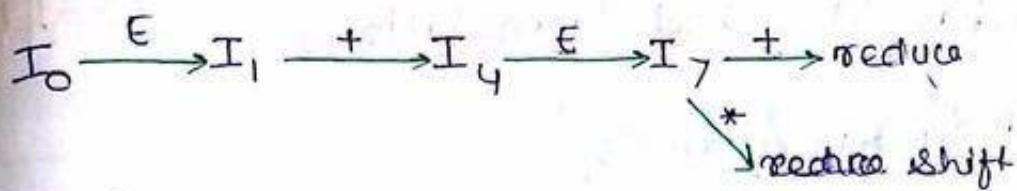
→ We want to use an ambiguous grammar because -

- Some of the ambiguous grammars are much natural, and a corresponding unambiguous grammar can be very complex.
- Usage of an ambiguous grammar may eliminate unnecessary reductions.

(123)

124



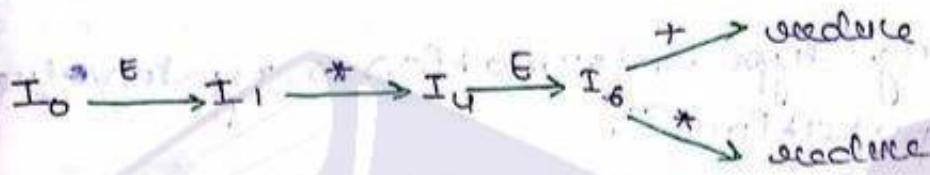


$\rightarrow$  when current token is +

- shift  $\rightarrow$  + is right-associative
- reduce  $\rightarrow$  + is left-associative

$\rightarrow$  when current token is \*

- shift  $\rightarrow$  \* has higher precedence than +
- reduce  $\rightarrow$  + has higher precedence than \*



### SLR Parsing Table for Ambiguous Grammar

		Action						
		id	+	*	(	)	\$	E
0		S3			S2			1
1			S4	S5			acc	
2	S3				S2			6
3			S4	S4		S4	S14	
4	S3				S2			7
5	S3				S2			8
6			S4	S5		S9		
7			S1	S5		S1	S11	
8			S2	S2		S2	S12	
9			S3	S3		S3	S13	



---

Mindfire is a software service provider,  
with unrelenting focus on optimal software development  
and delivery.

---



# *Compiler Design*

Topic:

***Syntax Directed Translation And Symbol Table***

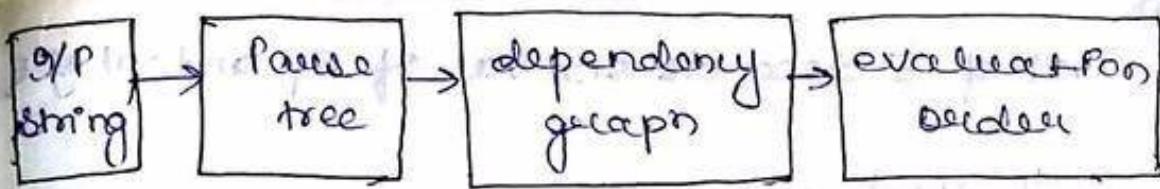
Contributed By:

***Sibnaranda Achari***

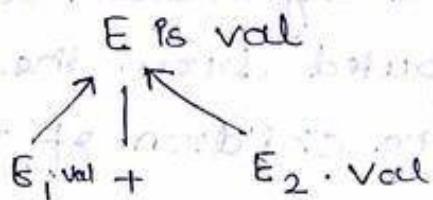
Trident Academy of Technology

Beyond Syntax Analysis

- Semantic analysis is the phase where we collect information about the type of expression and check for type related errors.
- The more information we collect at run time, the less overhead we had at run time.
- Collecting type information may involve "computation".
  - what is the type of  $x+y$  given the types of  $x$  and  $y$ ?
- Tool : attribute grammars
  - CFG
  - Each grammar symbol has associated attributes.
  - The grammar is augmented by rules (semantic actions) that specify how the values of attributes are computed from other attributes, which is called semantic rule.
  - The process of using semantic actions to evaluate attributes is called syntax directed translation.



Eg:  $E = E_1 + E_2$



$E$  value is dependent on  $E_1$  and  $E_2$ .

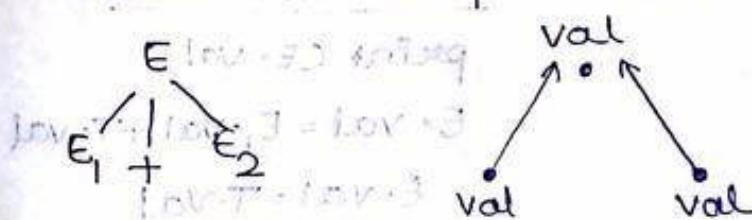
### Syntax directed definitions-

- It is a generalization of a CFD in which each grammar symbol has an associated set of attributes.
- **Synthesized attribute**
- **Inherited attribute**
- An attribute can represent a string, a no; a type, a memory location or relation.
- The value of an attribute at a parse tree node is defined by a semantic rule.

Eg  $E \rightarrow E + E$

$E \rightarrow E_1 + E_2$

Parse tree - Dependency graph -



$\text{val } 2 \times \text{val }, T = \text{INV } T$

$\Rightarrow \text{val } T \leq T$  (122)

128

Val of  $E$  depends on Val of  $E_1$  and Val of  $E_2$

### Synthesized attributes

- The value of a synthesized attribute at a node is computed from the values of attributes at the children of that node in the parse tree.
- A syntax directed definition that uses synthesized attributes is said to be an S-attribute definition.
- A parse tree for an S-attribute definition can always be annotated by evaluating the semantic rules for the attributes at each node bottom up, from the leaves to the root.
- A parse tree showing the value (S) of S attribute (S) is called an annotated parse tree.

### Syntax directed definition of a simple deck

#### Calculator

##### Production

$$L \rightarrow E_n$$

$$E \rightarrow E_1 + T$$

$$E \rightarrow T$$

$$T \rightarrow T_1 * F$$

##### Semantic rule

$$\text{print}(E \cdot \text{val})$$

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

$$E \cdot \text{val} = T \cdot \text{val}$$

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

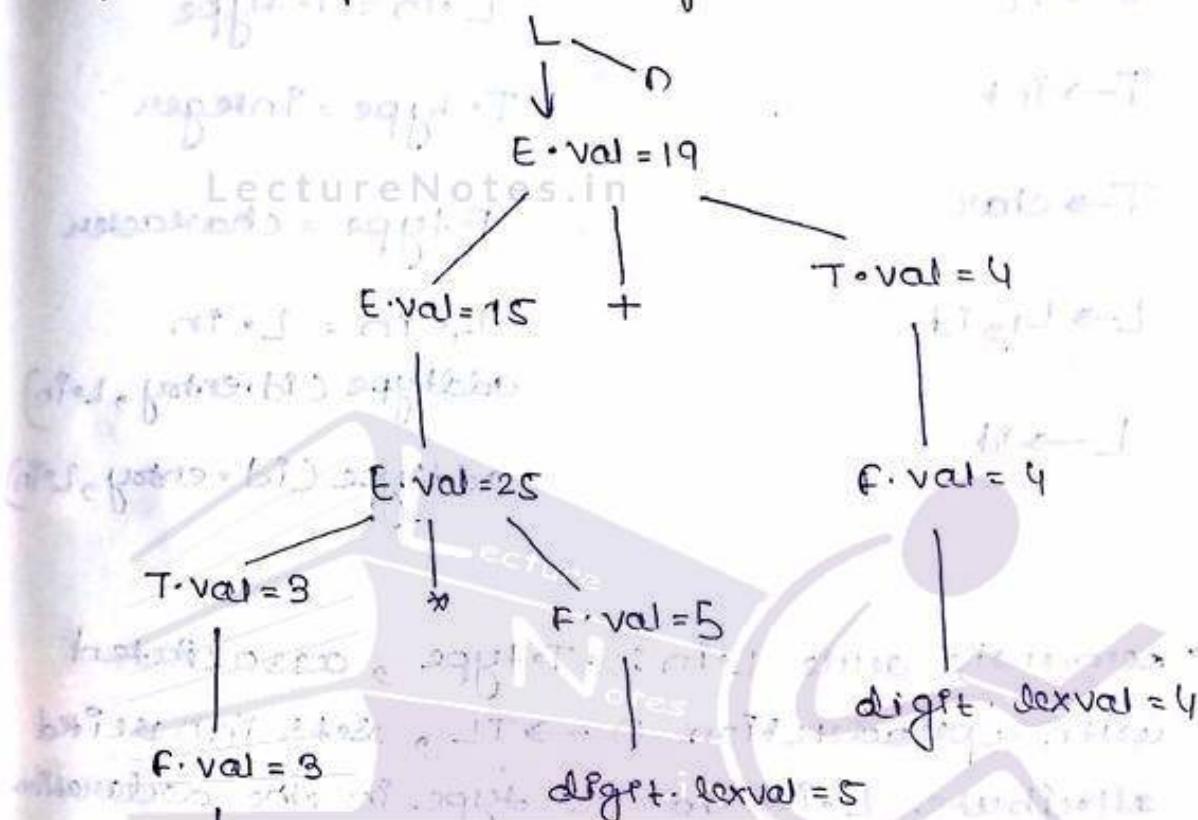


$T \rightarrow f$  $T \cdot \text{val} = f \cdot \text{val}$ 

(124)

 $f \rightarrow (\epsilon)$  $f \cdot \text{val} = \epsilon \cdot \text{val}$  $f \rightarrow \text{digit}$  $f \cdot \text{val} = \text{digit} \cdot \text{lexval}$ 

(125)

Annotated parse tree for  $3 * 5 + 4 \cdot 7$ 

### Inherited attributes

- The value of an inherited attribute is computed from the values of attributes at the siblings and parent of that node.
- An inherited value at a node in a parse tree is defined in terms of attributes at the parent and/or siblings of that node.
- It is used to express the dependence of a programming language construct.

(130)

## Syntax-directed definition of a grammar of declarations

## Production

⊕ → TL

T → Int

$T \rightarrow \infty$

$L \rightarrow L_{\text{1D}}$

1 → 1

### Semantic scale

L<sup>•</sup>T<sub>D</sub> = T<sup>•</sup> type

T-type = integer

T-type = character

$$L_1 \circ \varphi = L \circ \varphi$$

addtype C id·entry , L·in)

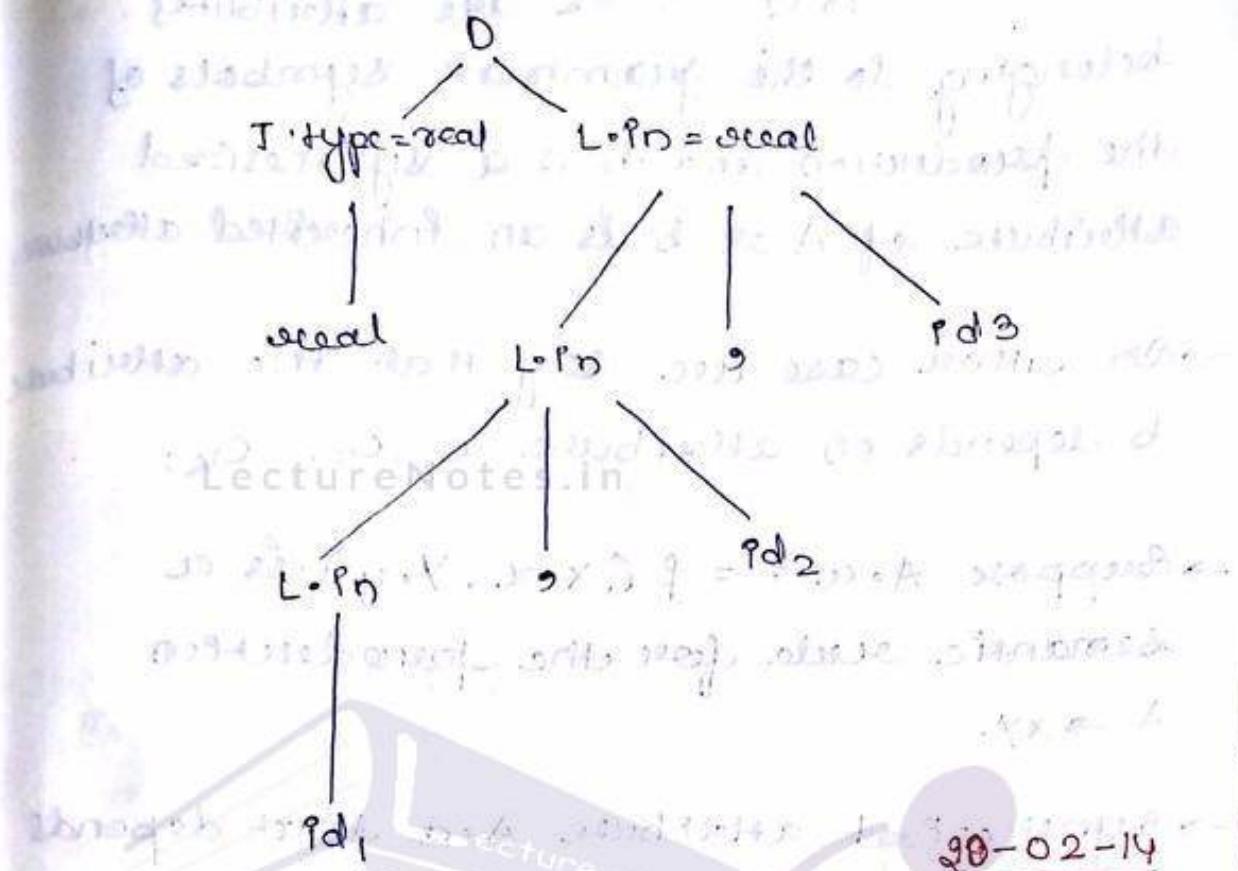
addtype Cid-enjoy, Lfin

→ Semantic rule  $L \cdot \text{in} := T \cdot \text{type}$ , associated with precondition  $D \rightarrow TL$ , sets predeclared attribute  $L \cdot \text{in}$  to the type in the declaration

→ Rules associated with the prediction for cell procedure addtype to add the type of each identifier to its entry in the symbol table.

→ The compiler then passes this type down the parent tree using the inherited attribute with `Link` which shows it is now being used.

## Anotated parse tree for $\text{seal Pd1, Pd2, Pd3}$ (13)



## Dependency Graphs

- The interdependencies among the inherited and synthesized attributes at the nodes in a parse tree can be depicted by a directed graph called a dependency graph.
- If an attribute  $b$  at a node  $i$  in a parse tree depends on a attribute  $c$ , then the semantic rule for  $b$  at that node must be evaluated after the semantic rule that defines  $c$ .
- If a production  $A \rightarrow \alpha$  has associated with it a set of semantic rules of the form -

$$b := f(c_1, c_2, \dots, c_k)$$

where  $c_1, c_2 \dots c_k$  are attributes belonging to the grammar symbols of the production and  $b$  is a synthesized attribute of  $A$  or  $b$  is an inherited attribute.

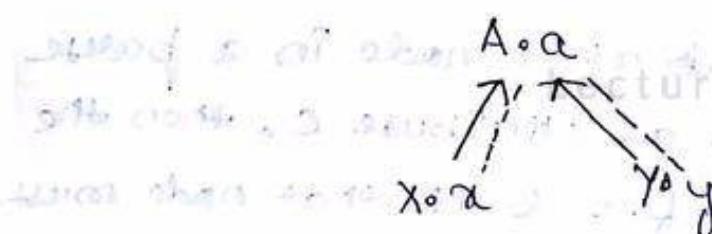
→ In either case, we say that the attribute  $b$  depends on attributes  $c_1, c_2 \dots c_k$ .

→ Suppose  $A \cdot a \stackrel{?}{=} f(x \cdot x, y \cdot y)$  is a semantic rule for the production  $A \rightarrow xy$ .

→ Synthesized attribute  $A \cdot a$  of  $A$  depends on the attributes  $x \cdot x$  and  $y \cdot y$ .

→ The dependency among node is the solid lines.

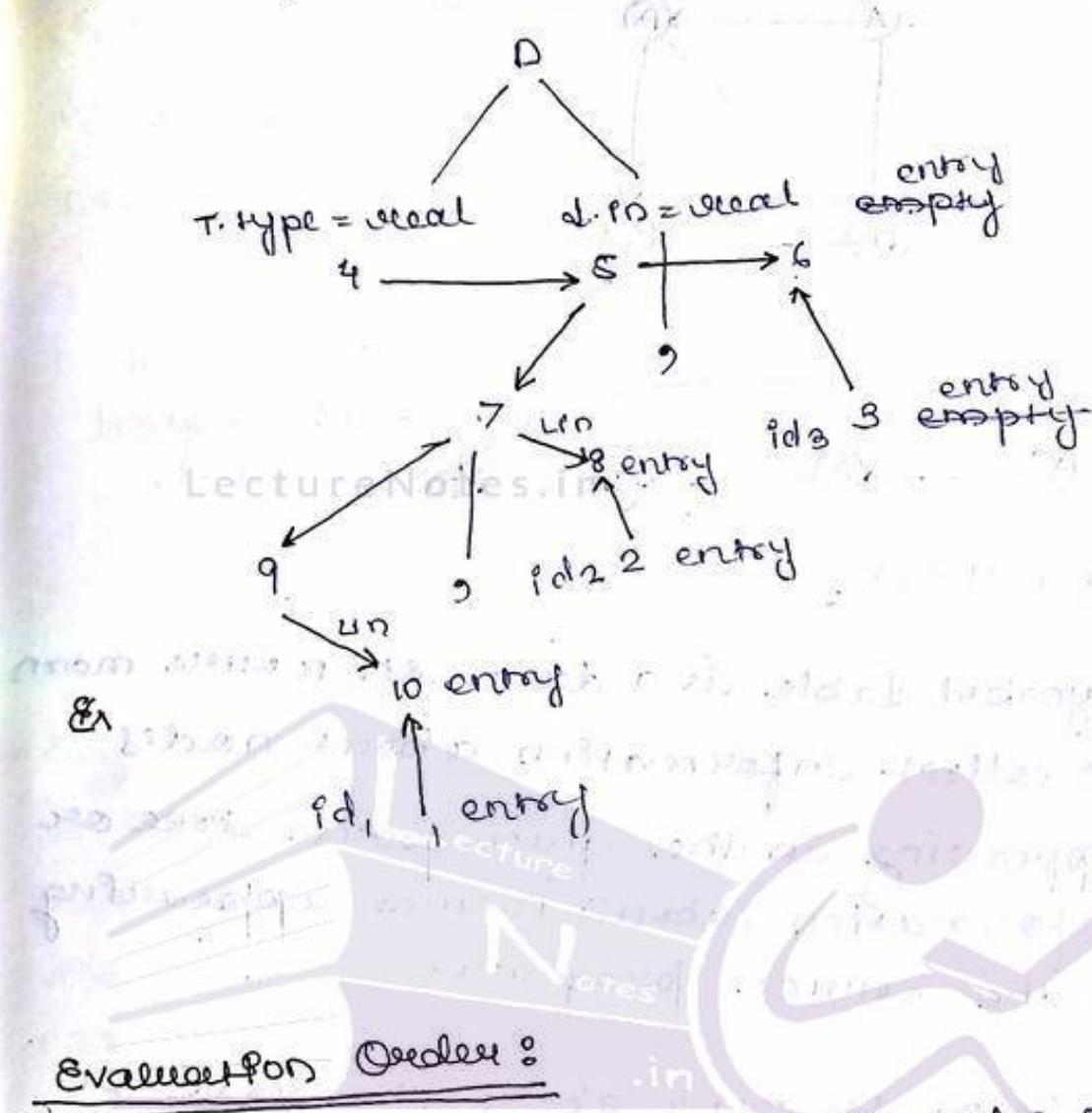
→ The parent node is represented by dotted lines.



Applications of SOT -

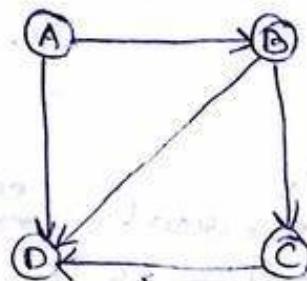
- Disambiguate overloaded operators
- Type checking
- Uniqueness checking
- Name checking
- Type coercion

## Dependency Graphs for real id1, id2, id3 (133)



- Evaluation order of semantic rules is obtained by a topological sort of a directed acyclic graph.
- In the topological sort, the dependent attributes  $c_1, c_2, \dots, c_k$  in a semantic rule  $b = f(c_1, c_2, \dots, c_k)$  are available at a node before  $f$  is evaluated.
- Evaluation of the semantic rules in this order yields the translation of the input string.

(133)

Eg -

### Symbol Table

- Symbol Table is a data structure meant to collect information about nodes appearing in the given parse tree or information about names appearing in the source program.
- It keeps the track about the scope or binding information about names.
- Each entry in the symbol table has a pair of information.

name	description type
id	integer

- Information consists of attributes such as datatype, location etc.





---

Mindfire is a software service provider,  
with unrelenting focus on optimal software development  
and delivery.

---

→ Whenever a name is encountered it is checked in the symbol table to see if edges already present. If information is not there, a new entry is created. (135)

int p11  
char c

LectureNotes.in

name	description type
if	keyword if
id1	integer
c	character

→ In some cases, the symbol table record is created by the lexical analyzer as soon as the name is encountered in the input and the attributes of the name are entered.

When the declarations are processed.

Operations on a symbol table.

- Determines whether a given name is in the table.

- Add a new name to the table.

- Delete a name from the table.

- Access information associated to a given name.

- Add new information for a given name.

(135)

(36)

→ There are various approaches to symbol table organisation.

- Linear List
- Search Tree
- Hash table

LectureNotes.in

Symbol	Type	Value
pi	float	3.141592653589793
phi	float	1.618033988749895
infinity	float	inf

LectureNotes.in



Linear List

- It is the simplest approach in the symbol table organisation.
- The new names are added to the table in order they arrive.
- A name is searched for its existence linearly.
- The average no. of comparisons required are proportional to the no. of entries in the symbol table.
- It takes less space but more access time.

Search Tree

- More efficient than linear trees.
- It provides two links - left and right which point to the record in the search tree.
- A new name is added at a proper position such that it can be accessed alphabetically.
- The time of searching is proportional to  $\log_2 n$ .

## Hash Table

(125)

- It is a table of  $K$  pointers from 0 to  $K-1$  that points to the symbol table to search a value where we find the hash value of the name by applying suitable hash function.
- If a record is non-existent then the record for the name is created and added to the list.

## The Scope of Information in the symbol table

- Each name possess a region of validity within the source program called the scope of the name.
- The rules governing the scope of a name in a block structure language are as follows:
  - A name declared within a block  $B$  is valid only within the block  $B$ .
  - If block  $B_1$  is nested under block  $B_2$ , then any name that is valid for  $B_2$  is also valid for  $B_1$ .

(140)

- (37)
- The rules become more complicated when a complex program is analysed.
  - To handle such a situation multiple symbol table for each active block is maintained with the following information :
    - \* Each table is list of names and their, all associated attributes and the tables are organised on the stack.
    - \* Whenever a new block is entered, a new table is pushed onto the stack.
    - \* When a declaration is completed the table on the stack is searched for the name.
    - \* If name is not found, it is inserted.
    - \* When a reference is generated, it is searched in all tables starting from the top

41



---

Mindfire is a software service provider,  
with unrelenting focus on optimal software development  
and delivery.

---



# *Compiler Design*

Topic:  
***Intermediate Code Generation***

Contributed By:  
***Sibananda Achari***  
Trident Academy of Technology

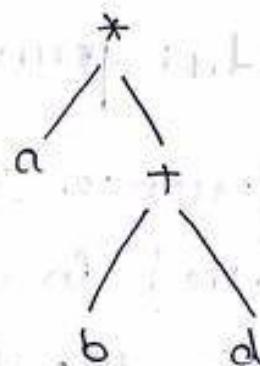
## Intermediate Code Generation

- ICCP are machine independent codes but are close to machine instructions.
- The given source program is converted to an equivalent program in an Intermediate language by the [ICG].  
The intermediate language can be syntax tree, postfix notation, 3-address code.

### Syntax Tree

Syntax tree is a variant of parse tree where each leaf node represents an operand and each internal node represents an operator.

Consider a sentence.  $a * (b + d)$



(141)

Semantic rules for the syntax tree -

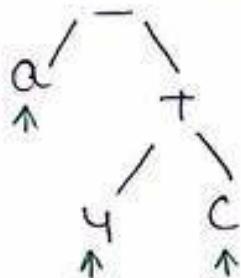
$$\begin{array}{l} E \rightarrow E_1 \text{ op } E_2 \\ E \rightarrow (E_1) \\ E \rightarrow \text{id} \end{array} \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} \text{production rule}$$

~~LectureNotes.in~~  
node (op, left, right) } 2 common  
leaf (id, entry) } functions

$$\begin{array}{l} E \cdot \text{val} := E_1 \cdot \text{val} \text{ op } E_2 \cdot \text{val} \\ E \cdot \text{val} := E_1 \cdot \text{val} \\ E \cdot \text{val} := \text{LEAF } (\text{id}, \text{entry}) \end{array} \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} \text{semantic rules}$$

Q       $a - 4 + c$

Define all the semantic rules for the given expression.  
 $\Rightarrow$  Translation:



- (P1) LEAF (id, entry - a)
- (P2) LEAF (id, entry - c)
- (P3) LEAF (number, val 4)

(142)

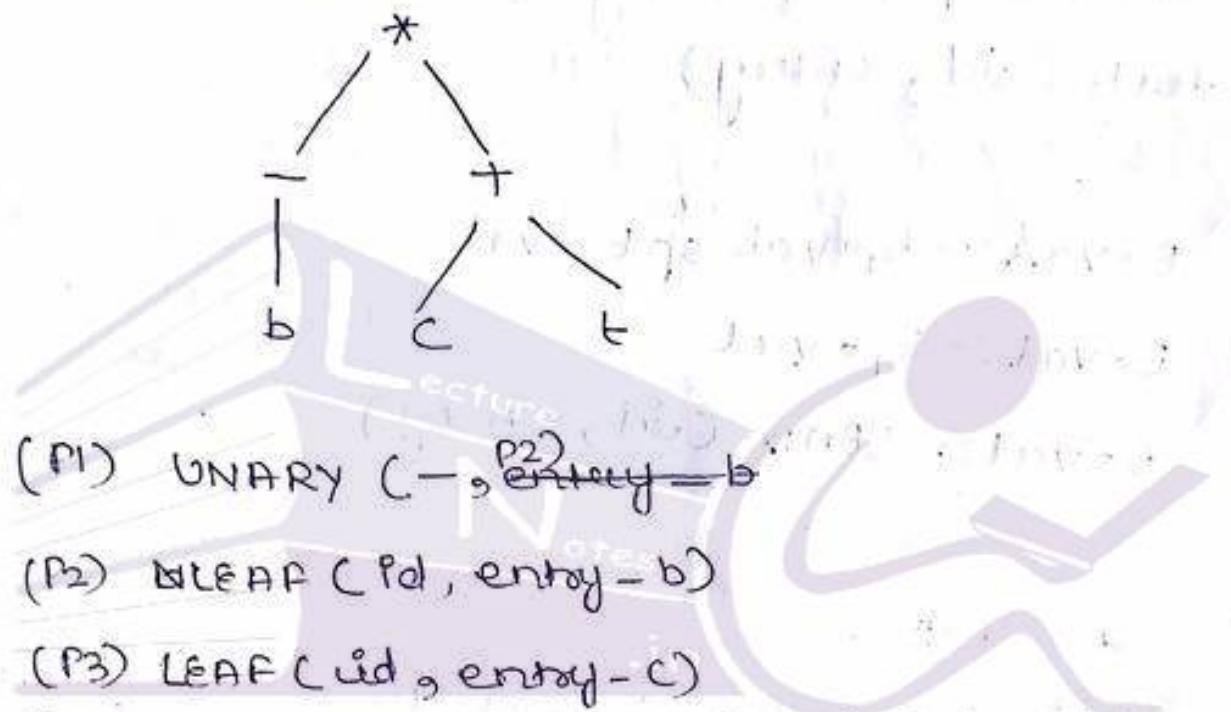
(14<sup>2</sup>)

(P4) NODE (+, P3, P2)

(P5) NODE (-, P1, P4)

$$Q \quad - b * (c + d)$$

Define translations.



(P1) UNARY (-, entry = b)

(P2) LEAF (id, entry = b)

(P3) LEAF (id, entry = c)

(P4) LEAF (id, entry = d)

(P5) NODE (+, P3, P4)

(P6) NODE (\*, P1, P5)

### Postfix Notation

Eg: 4 \* 5 + 3

Convert the expression from infix to  
postfix.

$$\Rightarrow ((4 * 5) + 3)$$

$$= ((4 * 5) + 3)$$

$$= (4 * 5) + 3$$

$$= 4 * 5 + 3$$

(14)

$$\mathcal{E} \rightarrow E_1 \text{ op } E_2$$

Postfix Notation -  $E ::= E_1 \text{ } E_2 \cdot \text{op}$

$E \cdot \text{val} ::= E_1 \cdot \text{val} \mid E_2 \cdot \text{val} \mid \text{op}$

→ It is another useful form of Intermediate code if the language is mostly expression.

$$\rightarrow E \rightarrow E_1 \text{ op } E_2$$

$$E \rightarrow (E_1)$$

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

$$\rightarrow E \cdot \text{val} ::= E_1 \cdot \text{val} \mid E_2 \cdot \text{val} \mid \text{op}$$

$$\rightarrow E \cdot \text{val} ::= E_1 \cdot \text{val}$$

$$\rightarrow E \cdot \text{val} ::= \text{id}$$

### 3- Address Code

Usually contains three addresses, two for operands, one for result.

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

x, y, z are names, constants or compiler-generated temporaries,

(145)

145

op stands for an operator such as fixed or floating point arithmetic operator.

### Representation of ICG

There are 4 forms of representation:

- Quadruple
- Triple
- Indirect Triple
- Static Single Assignment (SSA)

$$\Sigma \text{ (i) } x + y * z$$

$$\text{(ii) } A = -B * (C + D)$$

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

$$T_1 := C + D$$

$$T_2 := -B$$

$$T_3 := T_1 * T_2$$

$$A := T_3$$

$$\Rightarrow \underline{x + y * z}$$

$$T_1 := y * z$$

$$T_2 := x + T_1$$

1

146

## Types of Address - Address Statements

(115)

- Assignment statements of the form :  $x := y \text{ op } z$   
where op is a binary arithmetic or logical operation;
- Assignment statements of the form :  $x := \text{op } y$   
where op is a unary operation.  
Essential unary operations include unary minus, logical negation, and shift operators;
- Copy statements of the form :  $x := y$   
where the value of y is assigned to x;
- The unconditional jump goto L. The 3-address statement with label L is the next to be executed.
- Conditional jumps such as :  
 $\boxed{\text{if } x \text{ condition } y \text{ goto L}}$   
This instruction applies a relational operator ( $<, =, >, \leq, \text{ etc.}$ ) to x and y, and executes the statements with label L next if x stands in relation w.r.t. y to y;
- Procedure calls : call up, n and return values from functions ;

(147)

(vi) secondary. Their typical use is the following:

param  $x_1$

param  $x_2$

..

param  $x_n$

call p, n

generated as part of a call of the function:

$p(x_1, x_2, \dots, x_n)$ . The integer  $n$  indicates the number of actual parameters.

24/02/14

→ Indexed assignments of the form:

$x[i] = y$

$x[ij] := y$

The first one sets  $x$  to the value in the location  $i$  memory units beyond  $y$ . The statement  $x[ij] := y$  sets the content of the location  $i$  units beyond  $x$  to the value of  $y$ ;

→ Address and pointer assignments:

$x := &y$

$x := *y$



---

Mindfire is a software service provider,  
with unrelenting focus on optimal software development  
and delivery.

---

## Implementation of TAC Expression

(147)

$$\text{eg } a = -b * (c+d)$$

$$0 \leftarrow t_1 := c+d$$

$$1 \leftarrow t_2 := -b$$

$$2 \leftarrow t_3 := t_2 * t_1$$

$$3 \leftarrow a := t_3$$

### Final representation:

Index	OP	Arg 1	Arg 2	Res
(0)	+	c	d	$t_1$
(1)	uminus	b	"	$t_2$
(2)	*	$t_1$	$t_2$	$t_3$
(3)	=	$a+b$	$t_3$	a

### Simple representation:

Index	OP	Arg 1	Arg 2
(0)	+	c	d
(1)	uminus	b	"
(2)	*	(0)	(1)
(3)	=	(2)	a

Let avoid putting temporary names in symbol table.

(148)

## 118' Indirect triple representation:

array of  
It consists of pointers to array of triples

(0)	16
(1)	17
(2)	18
(3)	19

	op	Arg1	Arg2
(16)	+	c	d
(17)	Umplus	b	
(18)	*	(0)	
(19)	=	(2)	a

Generate TAC for the given expression

$\Rightarrow$   $x = id \rightarrow -a * b + -a * b \rightarrow E$

 $t_1 := -a$ 
 $t_2 := t_1 * b$

$t_3 := Umplus a$

$t_4 := t_3 * b$

$t_5 := t_2 + t_4$

$x := t_5$

Converting the above into postfix form

$s \rightarrow id = E \rightarrow \text{starting root}$

$\textcircled{S} \quad x = E$

id · place =  $E \cdot \underbrace{\text{place}}_{\text{value}}$

(14)

$E \cdot \text{code}$

TA statements

$S \rightarrow Id = E \rightarrow \text{Production}$

Semantic rules -

$S \cdot \text{code} := E \cdot \text{code} \amalg \text{gen}(Id \cdot \text{place} \> ::= \> E \cdot \text{place})$

→  $E \cdot \text{place}$  → the name that will hold the value of  $E$

→  $E \cdot \text{code}$  → the sequence of three-address statements evaluating  $E$ .

→ gen → generates the three address code.

Expression appearing instead of variables when passed to gen and quoted operand like "+".

→  $\amalg$  → symbol is used for concatenation, only are added to show the relationship bet<sup>n</sup> LHS & RHS. Such relationship is called semantic rules

Production

Lectures

$E \rightarrow E_1 + E_2$

$E \cdot \text{place} ::= \text{newtemp};$

$E \cdot \text{code} ::= E_1 \cdot \text{code} \amalg E_2 \cdot \text{code}$

$\amalg \text{gen}(E \cdot \text{place} ::=$

$E_1 \cdot \text{place} + E_2 \cdot \text{place})$

→ newtemp → generates temporary names.

(15)

$E \rightarrow E_1 * E_2$ $E \rightarrow -E_1$ $E \rightarrow (E_1)$ $E \rightarrow id$	$E \cdot place := \text{newtemp};$ $E \cdot code = E_1 \cdot code \parallel E_2 \cdot code \parallel$ $\text{gen}(E \cdot place) \cdot := ? E_1 \cdot place \cdot *$ $E_2 \cdot place)$  $E \cdot place := \text{newtemp};$ $E \cdot code = E_1 \cdot code \mid \text{gen}(E \cdot place)$ $\cdot := ? 'minus' E_1 \cdot place)$  $E \cdot place := E_1 \cdot place;$ $E \cdot code = E_1 \cdot code$  $E \cdot place := \text{fd} \cdot place;$ $E \cdot code := ?$
--	---

→ To distinguish between expression nodes on RHS & LHS  
 the expression nodes on RHS is represented with  
 $E_1$ .

### TAC for Boolean Expressions

- Boolean expressions are composed of the Boolean operators (and, or & not)
- Boolean expressions have 2 primary purposes
  - to compute logical values
  - to be used as conditional expressions

→ Relational expressions are of the form: (151)

A operator B where A & B are arithmetic expressions. operator can be of the following types:  $<$ ,  $>$ ,  $\leq$ ,  $\geq$ ,  $=$ ,  $\neq$  etc.

→ Two methods exist there for translating boolean functions -

- encode true & false numerically & evaluate a boolean expression.
- implement boolean expression by a flow of control.

→ Consider a boolean expression generated by the following grammar:

$$E \rightarrow E \text{ or } E \mid E \text{ and } E \mid \text{not } E \mid (CE) \mid \\ \text{Id} \text{ or } \text{id} \mid \text{true} \mid \text{false}$$

→ The precedence of operator is as follows

• not is highest, then and, then or as lowest.

→ Boolean expression have different interpretations depending on the context.

(152)

\* Compute logical values

Eg - a OR b AND NOT c

a OR b AND NOT c

JA Stmt :-

t1 := NOT c

t2 := b AND t1

t3 := a OR t2

### Semantic Rules / Translation Rules

production rule	Semantic rule
$E \rightarrow E_1 \text{ or } E_2$	$\{ E \cdot \text{place} := \text{newtemp};$ $\text{emit } (E \cdot \text{place} \leftarrow E_1 \cdot \text{place}$ $\text{or } E_2 \cdot \text{place}) \}$
$E \rightarrow E_1 \text{ and } E_2$	$\{ E \cdot \text{place} := \text{newtemp};$ $\text{emit } (E \cdot \text{place} \leftarrow E_1 \cdot \text{place}$ $\text{and } E_2 \cdot \text{place}) \}$
$E \rightarrow \text{not } E_1$	$\{ E \cdot \text{place} = \text{newtemp};$ $\text{emit } (E \cdot \text{place} \leftarrow \text{not } E_1 \cdot \text{place}) \}$
$E \rightarrow (E_1)$	$\{ E \cdot \text{place} := E_1 \cdot \text{place} \}$

(153)

$E \rightarrow id_2 \text{ vceop } id_2$

18/25

$\{ E \cdot place = \text{newtemp};$   
 emit ( "if"  $id_2 \cdot place$  vceop  
 $id_2 \cdot place$  "goto" nextstat + 3 );  
 emit (  $E \cdot place$  " $:=$ " "0" );  
 emit ( "goto" nextstat + 2 );  
 emit (  $E \cdot place$  " $:=$ " "1" ) }

$E \rightarrow \text{decne}$

$\{ E \cdot place := \text{newtemp};$   
 emit (  $E \cdot place$  " $:=$ " "1" ) }

$E \rightarrow \text{false}$

$\{ E \cdot place := \text{newtemp};$   
 emit (  $E \cdot place$  " $:=$ " "0" ) }

Eg:  $a \leftarrow b$

For the above given eg: - a OR b AND NOT c

$t1 = \text{NOT } c$

if  $b$  AND  $t1$

goto L true

L1: if  $a$  OR  $t1$

goto L true

L false : - -

L true : - -

(15/25)

Ex :  $a < b$  or  $c < d$  and  $e < f$

if  $a < b$  goto Ltrue

goto L1

L1 : if  $c < d$  goto Ltrue

goto Lfalse

L2 : if  $e < f$  goto Ltrue

goto Lfalse

Lfalse :

Ltrue :

short  
circuit  
code

### TAC for Boolean Expression

100: if  $a_1 < b$  goto 103

101:  $t_1 := 0$

102: goto 104

103:  $t_1 := 1$

104: if  $c < d$  goto 107

105:  $t_2 := 0$

106: goto 108

107:  $t_2 := 1$

108: if  $e < f$  goto 111

109:  $t_3 := 0$

110: goto 112

111:  $t_3 := 1$

112:  $t_4 := t_2 \text{ and } t_3$



---

Mindfire is a software service provider,  
with unrelenting focus on optimal software development  
and delivery.

---

### Short - Circuit :

(155)

In short-circuit or jumping code, the boolean operators AND, OR, NOT decompiles into jumps. The operators themselves do not appear in the code. Instead the value of the expression is represented by the position.

Eg - if ( $x < 100 \text{ || } x > 200 \text{ ?? } x != y$ )  
 $x = 0 ;$

### short - circuit code :

iffalse  $x < 100$  goto L2

iffalse  $x > 200$  goto L1

L1 :

L2 :

Find the lines replaced and missed in the above short circuit code.

⇒ if  $x < 100$  goto Ltrue  
    goto L1

L1 : if  $x > 200$  goto L2  
    goto Lfalse

L2 : if  $x != y$  goto Ltrue  
    goto Lfalse

Lfalse :  
Ltrue :

(157)

(156) SSA (Static Single Assignment)

25/02/14

$$P = a + b$$

$$q_1 = P - c$$

$$P = q_1 * d$$

$$P = e - P$$

$$q_1 = P + q_1$$

all the assignments have distinct names

SSA -

$$P_1 = a + b$$

$$q_{11} = P_1 - c$$

$$P_2 = q_{11} * d$$

$$P_3 = e - P_2$$

$$q_{12} = P_3 + q_{11}$$

→ Static Single Assignment is an intermediate representation that facilitates certain code optimization to distinct aspects. Distinguish SSA from 3-address code.

→ The first is that all assignments in SSA are two variables with distinct names. Hence the term static single assignment.

→ The same variable may be defined in two different control flow path in a program.

Eg: if (flag)

$x = -1;$

else

$x = 1;$

$y = x * a;$

$\Rightarrow \text{if } (\text{flag})$

$x_1 = -1;$

else

$x_2 = 1;$

$y = \phi(x_1, x_2) * a;$

→ The  $\phi$  function returns the value of the argument that corresponds to the control flow path that was taken to get to the assignment.

### Type Expressions

→ A basic type is a type expression.

→ A typename is a type expression.

→ A type expression can be formed by applying the array type constructor to a number and a type expression.

Eg array (number, type expression)

$\text{int } a[2];$

(157)

(159)

(55)  $\Rightarrow a(2, \text{integer})$

- The type expression can be formed by using the type constructor for function types.
- Type expression may contain variables whose values are type expression.

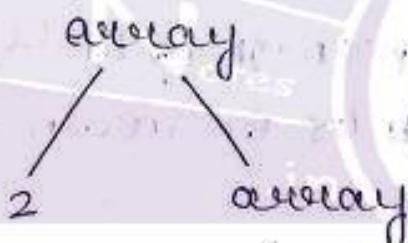
Eg:  $\text{int } [2][3];$

$\Rightarrow \text{array}(2, \text{type expression});$

$\text{array}(2, \text{array}(3, \text{integer}));$



Syntax tree



Production rule:

$$D \rightarrow T \cdot id \cdot ; \mid D \mid e$$

*Name of the record*

$$B \rightarrow \text{int} \mid \text{float}$$

$$T \rightarrow BC$$

$$C \rightarrow \epsilon \mid [\text{num}] \mid C$$

eg: `int [2][3];`

$D \Rightarrow T \text{ Id}$

$\Rightarrow B C \text{ Id}$

$\Rightarrow \text{int } C \text{ Id}$

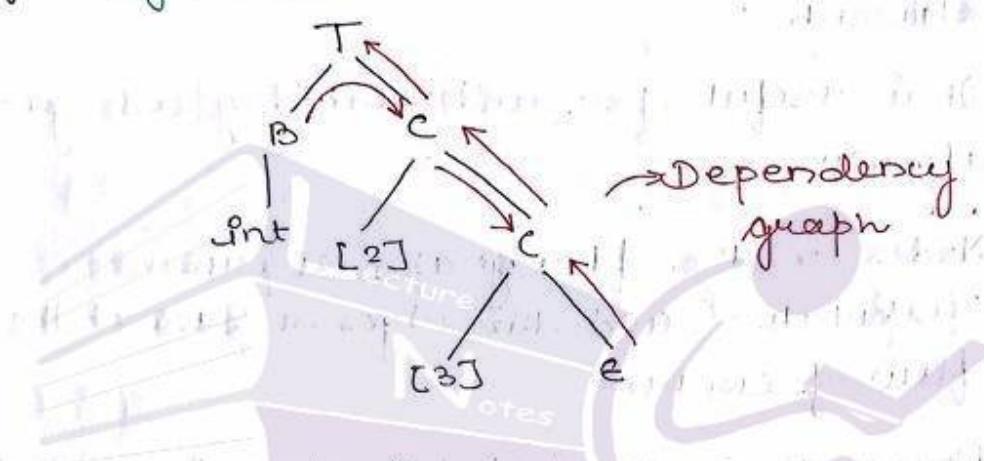
$\Rightarrow \text{int } [\text{num}] C \text{ Id}$

$\Rightarrow \text{int } [\text{num}] [\text{num}] \text{ Id}$

$\Rightarrow \text{int } [2][3] \text{ Id};$

eg:

Syntax tree



Type expression

$$\begin{cases} t = B \cdot \text{type}, \\ w = B \cdot \text{width} \end{cases} \xrightarrow{T \rightarrow BC}$$

$$\begin{cases} B \cdot \text{type} = \text{integer}, \\ B \cdot \text{width} = 4 \end{cases} \xrightarrow{B \rightarrow \text{int}}$$

$$\begin{cases} B \cdot \text{type} = \text{float}, \\ B \cdot \text{width} = 8 \end{cases} \xrightarrow{B \rightarrow \text{float}}$$

$$\begin{cases} C \cdot \text{type} = t, \\ C \cdot \text{width} = w \end{cases} \xrightarrow{C \rightarrow e}$$

$$\begin{cases} \text{array}(\text{number} \cdot \text{value}, C1 \cdot \text{type}), \\ C \cdot \text{width} = \text{number} \cdot \text{value} * C1 \cdot \text{width} \end{cases} \xrightarrow{C \rightarrow e} [\text{num}] C$$

eg: `C2, integer`  $w = 2 * 4 = 8$



---

Mindfire is a software service provider,  
with unrelenting focus on optimal software development  
and delivery.

---



# *Compiler Design*

Topic:

***Code Generation And Code Optimization***

Contributed By:

***Sibanya Achary***

Trident Academy of Technology

## Basic Blocks And Flow Graphs

- Flow Graph:
- A graph representation of a 3-Address statements.
- It is useful for understanding code generation algorithm.
- Nodes in the flow graph represent computation and the edges represent the flow of control.
- Flow graphs are useful to find inner loops where program is expected to spend most of its time.
- The basic blocks can be generated from the flow graph in the following manner:

Step - I: Partition the Intermediate code into basic blocks.

Basic blocks are maximal sequence

3-address instructions with the following properties.

- The flow of control can enter the basic block through the first instruction in the block i.e. there are no jumps into the middle of the block.
- Control will leave the block at the end, (last instruction) without halt or possibility of branching.

Step-II: The basic blocks become the nodes of the flow graph whose edges indicates which block can follow which other blocks.

Algorithm 1 : Partition into basic blocks.

Input : A sequence of three-address statements

Output : A list of basic blocks with each three-address statement in exactly one block.

Method :

- We first determine the set of leaders, the first statement of basic blocks.

The rules we use are as follows:

- The 1<sup>st</sup> statement is a leader.
- Any statement that is the target of a conditional / unconditional goto is a leader.

begin

pseud = 0;

l = 1;

do begin

pseud = pseud + a[i] \* b[l];

l = l + 1;

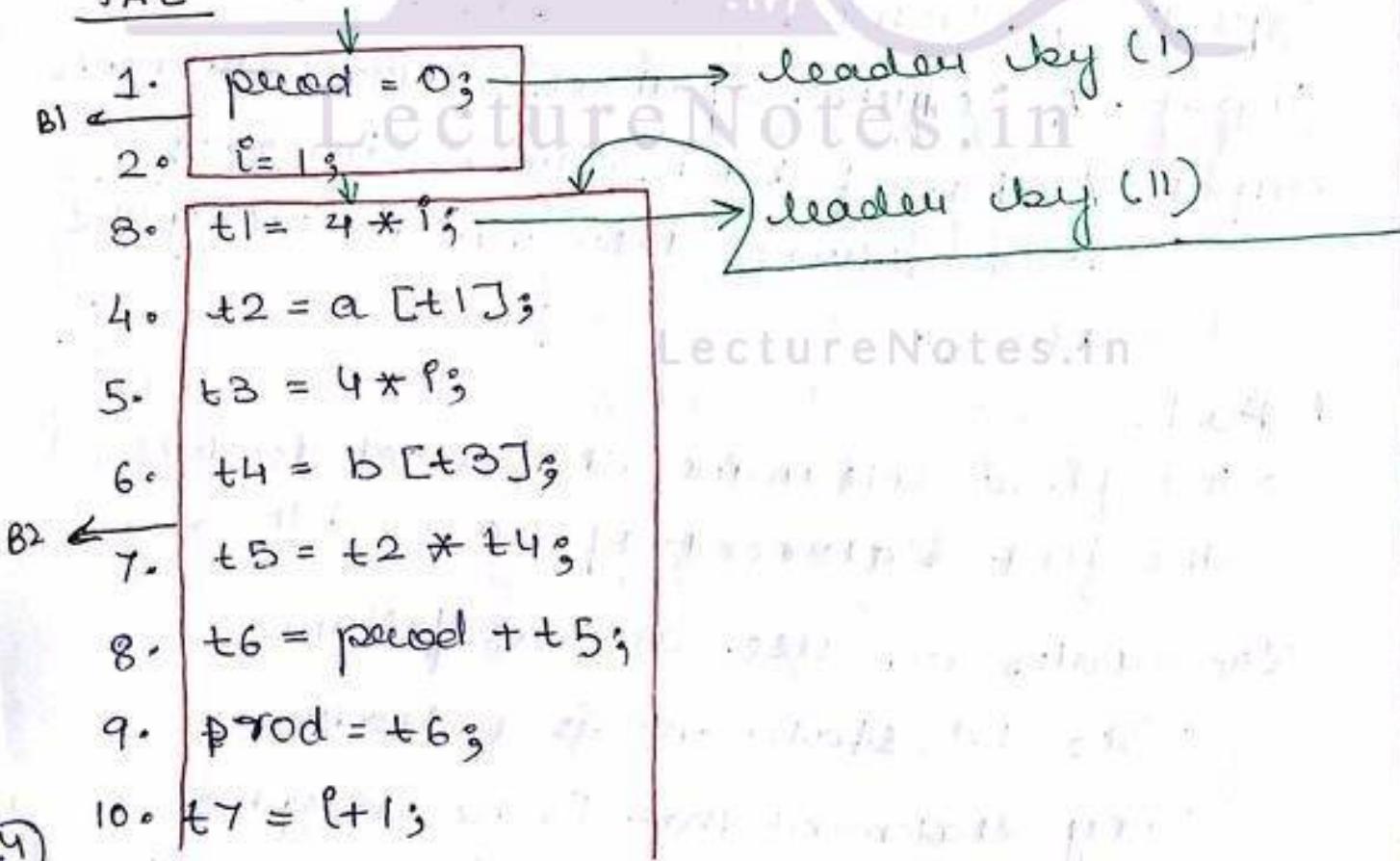
end

while (i &lt;= 20)

end.

If there is a statement involving do-while, while, for then the statement need to be change into if statement with conditional jump.

The 3-address code for the above given code is -

JAC

11.  $p = t \gamma;$   
 12. If ( $p \leq 20$ ) goto 3

(163)

Now we need to identify the leaders -

preced = 0; leader by (1)

$t1 = 4 * p$ ; leader by (11)

LectureNotes.in

$\exp(x)$

```

  { int p = 1;
    for (i = 2; p <= x; i++)
    {
      p = p * i;
    }
    r = p + 1;
  }
  }
```

convert it into TAC.

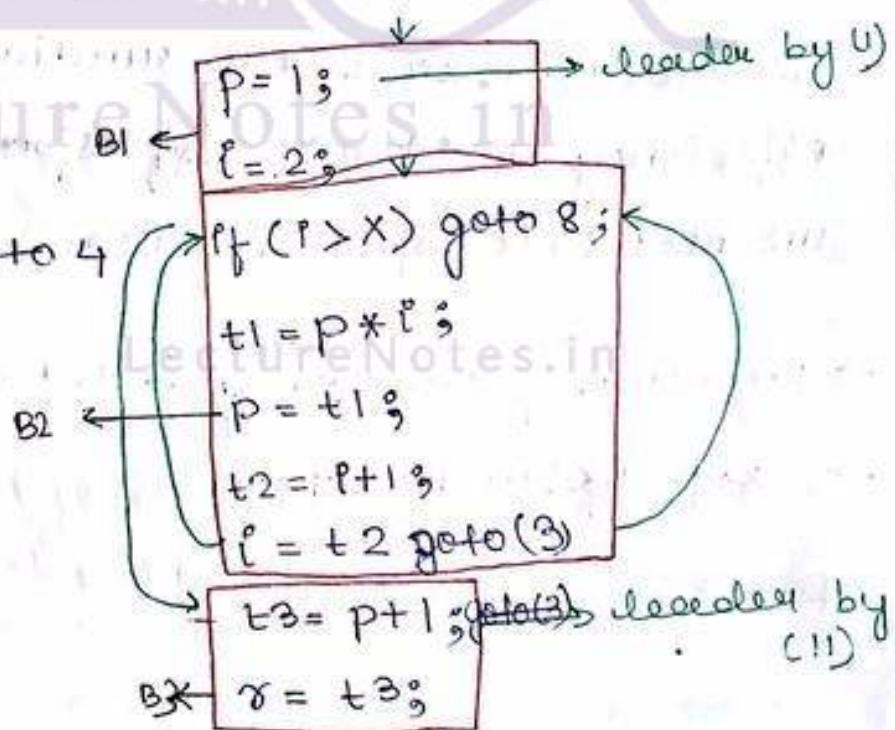
1.  $\exp(x)$

2.  $p = 1$

3. if ( $i \leq x$ ) goto 4

4.

5.



loop = { B2 }

(165)

(16y) TAC -  
 $\exp(x)$   
{  
1.  $p = 1;$   
2.  $l = 2$  goto 4  
3.  $t1 = p * l$   
4.  $p = t1;$   
5.  $t2 = l + 1;$   
6.  $l = t2;$   
7. if ( $i \leq n$ ) goto 3  
8.  $t3 = p + 1;$   
9.  $y = t3;$

### Principle Sources of Optimisation

- We want to reduce the runtime and to improve efficiency in terms of space and time so we need code optimization.
- Function restructuring, transformations
- Loop optimization

### Function restructuring

- Common subexpression elimination
- Copy propagation
- (16b) → Dead code elimination
- Constant folding

## Dead Code Elimination :

(185)

- If a piece of code define at a particular program point is not used anywhere in the program then that code is called as **dead code** and it should be removed from the program.

eg

$$c = a * b$$
$$\boxed{x = a} \rightarrow \text{dead code}$$

$$\begin{array}{c} \circ \\ \circ \\ \circ \\ d = a * b \end{array}$$

$$\Rightarrow \boxed{c = a * b}$$
$$\boxed{d = a * b} \rightarrow \text{common subexpression}$$

$$\Rightarrow t1 = a * b$$

$$c = t1$$

$$d = t1$$

## Common Subexpression Elimination :

An expression need not be evaluated if it was previously computed and the values of the variables in this expression have not changed since earlier computation.

eg -  $a = \boxed{b * c};$  common subexpression  
 $d = \boxed{b * c} + x - y;$

(167)

(16)  $\Rightarrow t1 = \text{do} * c$

$a = t1$

$d = a + x - y \rightarrow$  if  $a$  not written, then  $a = t1$   
or,  $d = t1 + x - y$  would become a dead code.

If in the question, dead code correction  
hasn't been mentioned then we could have  
written  $a$  not statement.

28/02/14

## Compile time evaluation

69  $x = 12.4;$

:

$y = x / 2.3; \rightarrow$  execution time action

We can improve the efficiency of the program by  
shifting execution time actions to compile time  
action.

$y = 12.4 / 2.3 \rightarrow$  compile time action

We can evaluate an expression by a single  
value known as constant folding.

If a variable is assigned a constant value and  
used in an expression without changing its assigned  
value then some portion of the expression can be  
evaluated using constant value.

## Variable Propagation

Eg.  $c = a * b;$  → used  
 defined  $x = a;$   
 $\vdots$   
 $d = x * b;$   
 $\Rightarrow d = a * b;$  Notes.in  
 (replace.  $x$  by  $a$ )

If a variable is assigned to another variable, we use one in place of another. This will be useful to carry out other optimizations that otherwise isn't possible.

If we replace.  $x$  by  $a$ , then  $a * b$  and  $x * b$  will be identified as common subexpression and the common subexpression need to be eliminated as follows:

$$\begin{aligned} t1 &= a * b \\ c &= t1 \\ \vdots \\ d &= t1 \end{aligned}$$

## Code Motion

We can improve the execution time of the program by reducing the evaluation of frequency of expressions.

Evaluation of expression is moved from one part of the program to another in such a way that it is evaluated less frequently. (169)

(168) Loops are executed several times till an exiting loop invariant statement gets off the loop to optimize the code.

$a = 200;$

while ( $a > 0$ )

{

$b = a + y;$  → loop invariant statement

if ( $a \% b == 0$ )

printf ("%d", a)

}

$a = 200;$

$b = a + y;$

while ( $a > 0$ )

{

if ( $a \% b == 0$ )

printf ("%d", a);

}

### Introduction Variable & Strength Reduction

An induction variable may be defined as an integer scalar variable which is used in loop for some assignment statements.



---

Mindfire is a software service provider,  
with unrelenting focus on optimal software development  
and delivery.

---

(169)  
Strength reduction means replacing high strength operator by a low strength operator.

Ex:

$$l = 1$$

while ( $i < 10$ )

{

$y = l * 4;$

}

$$\Rightarrow l = 1;$$

$$t1 = 4;$$

while ( $t1 < 40$ )

{

$$y = t1;$$

$$t1 = t1 + 4;$$

}

Use: Algebraic Identifiers - certain computation that look different to the compiler and were not identified as common subexpression, they are actually same.

eg

$b \oplus c$ . will be detected different from  
 $c \oplus b$ .

Certain operations like multiplication & addition will produce same result and they can be treated as common subexpression.

(171)

(17)

## Peephole Optimization

A statement by statement code generation strategy offers produce target code that contains redundant instructions and sub-optimal constructs.

A simple but effective technique for improving the target code is peephole optimization.

It is a method for trying to improve the performance of the target program by examining a short sequence of target instructions called the peephole and replacing these instructions by a short or faster sequence of instruction whenever possible.

Peephole is a small moving window on the target program that need to be optimization.

### Techniques of Peephole Optimization.

- Redundant instruction optimization
- Flow of control optimization
- Redundant Instruction Elimination

Eg

- 1)  $\rightarrow \text{Mov } a, R_0$
- 2)  $\rightarrow \text{Mov } R_0, a$

We can delete instruction 2 because whenever instruction 2 is executed, instruction 1 will ensure that the value of register R0 is already there in it. (171)

## → Flow of Control Optimization

- Unreachable Code

# define debug 0

if (debug)

{

    print debugging information

}

=

lecture

N  
otes

# define debug 0

if (debug == 1) goto L1

goto L2

L1 : Print debugging information

L2 :

↓ optimization

If (debug != 1) goto L2

Print debugging information

L2 : ..

The unnecessary jumps can be eliminated in the intermediate code by replacing the jump sequences. (172)

## Undecidability

**Decidability:** A problem with two answers yes or no is said to be decidable if the corresponding language is recursive. That language is also said to be decidable.

A problem or a language is said to be undecidable if it is not decidable.

- ⇒ The intersection of two regular language is decidable.
- ⇒ The intersection of two CFL is undecidable.
- ⇒ Let  $G$  be an unrestricted grammar. Then the problem of whether or not  $L(G) = \emptyset$  is undecidable.
- ⇒ Let  $M$  be any TM, then the problem of whether or not  $L(M)$  is finite is undecidable.
- ⇒ Let  $M$  be any TM, then the problem of  $L(M)$  is regular is undecidable.
- ⇒ Let  $M_1$  and  $M_2$  be two TMs, then is  $L(M_1) \subseteq L(M_2)$  is undecidable.

(173)

$\Rightarrow$  Let  $G_1$  be any unrestricted grammar and  $G_2$  be any regular grammar then the problem of whether or not  $L(G_1) \cap L(G_2) = \emptyset$  is undecidable.

$\Rightarrow$  Let  $G_1$  &  $G_2$  be two CFGs, then

$L(G_1) \cap L(G_2) = \emptyset$  and  $L(G_1) \subseteq L(G_2)$  are undecidable.

$\Rightarrow$  The post correspondence problem PC is undecidable.

$\Rightarrow$  The modified PCP is also undecidable.

$\Rightarrow$  The elimination of ambiguity from a CFG is undecidable.

Q Which of the following is/are undecidable?

1.  $G_1$  is a CFG. Is  $L(G_1) = \emptyset$ ?

2.  $G_1$  is a CFG. Is  $L(G_1) = \Sigma^*$ ?

3.  $M$  is a TM. Is  $L(M)$  regular?

4.  $A$  is a DFA and  $N$  is an NFA.

$$L(A) = L(N) ?$$

① 3 only

② 1, 2 & 3 only

③ 1 & 3 only

④ 2 & 3 only

## P & NP Class

(174)

- P-Class: A language  $L$  is said to be in P-class if there exist a deterministic TM  $M$  such that  $M$  is of the time complexity  $P(n)$  for some polynomial  $P$  and  $M$  accepts  $L$ . Here,  $n$  is the no. of moves taken by TM to accept the string.
- Eg: Accept the string.

2-SAT

PATH

Relatively Prime

→ P-class problems can be solved efficiently.

- NP-Class: A language  $L$  is said to be in NP-class if there exist a non-deterministic TM  $M$  such that  $M$  is of time complexity  $P(n)$  for some polynomial  $P$  and  $M$  accepts  $L$ .

Eg:

SAT

PATH

composite

CLIQUE

(175)

## → (35) Polynomial Time Reduction:

If we can transform instances of one problem in deterministic polynomial time into instances of a second problem that has the same answer yes/no. Then we will say the first problem is polynomial time reducible to second problem.

## → NP-Complete:

A language  $L$  is said to be in NP-complete if it satisfies the following -

- (i)  $L$  must be in NP-class.
- (ii) For every language  $L'$  in NP, if there is a polynomial time-reduction from  $L'$  to  $L$ .

→ Eg: 3-SAT

SAT

PATH

CLIQUE

## → NP-Hard :

(17c)

A language  $L$  is said to be in NP-hard,  
if

(i)  $\forall L' \in \text{NP}$ . if  $L' \leq_p L$

(ii)  $L$  may or may not be in NP class.

Eg : - TSP

⇒ Every NP-complete is also NP-hard but not  
Vice-versa.

### NOTE:

- If  $L_1 \leq_p L_2$  and if  $L_1$  is NP-complete then  
 $L_2$  is in NP-hard if we don't know  
anything about  $L_2$  that it is in NP or not.
- If  $L_1 \leq_p L_2$  and if  $L_1$  is NP-complete and  
given  $L_2$  in NP-class then  $L_2$  is also in  
NP-complete.
- If  $L_1 \leq_p L_2$  and iff  $L_1$  is undecidable  
 $\Rightarrow L_2$  is also undecidable.
- If  $L_1 \leq_p L_2$  and iff  $L_2$  is decidable  
 $\Rightarrow L_1$  is also decidable.

(17d)

(77)

Q Let 'S' be an NP-complete, Q and R be two other problems not known to be in NP. Q  $\leq_p$  S and S  $\leq_p$  R. Then which of the following is true?

- (A) Q is NP-complete
- (B) Q is NP-Hard
- (C) R is NP-complete
- (D) R is NP-Hard

Q Ram and Shyam have been asked that a certain problem  $\Pi$  is NP-complete.

Ram shows a polynomial time-reduction from 3-SAT problem to  $\Pi$  and Shyam shows a polynomial time reduction from  $\Pi$  to 3-SAT. Which of the following can be inferred from these reductions?

$$\begin{aligned} 3\text{-SAT} &\leq_p \Pi_{\text{NP-complete}} \\ \Pi &\leq_p 3\text{-SAT} \end{aligned}$$

NP-complete

- (A)  $\Pi$  is NP-Hard
- (B)  $\Pi$  is in NP but is not NP-complete
- (C)  $\Pi$  is NP complete
- (D)  $\Pi$  is neither NP-hard nor in NP.

I consider 3 decision problems  $P_1, P_2, P_3$ .

Given that -  $P_1$  is decidable and  $P_2$  is undecidable  
which one of the following is true?

- (a)  $P_3$  is decidable if  $P_1 \leq_p P_3$ .
  - (b)  $P_3$  is undecidable if  $P_3 \leq_p P_2$ .
  - (c)  $\cancel{P_3}$  is undecidable if  $P_2 \leq_p P_3$ .
  - (d)  $P_3$  is decidable if  $P_3 \leq_p \overline{P_2}$ .



---

Mindfire is a software service provider,  
with unrelenting focus on optimal software development  
and delivery.

---