

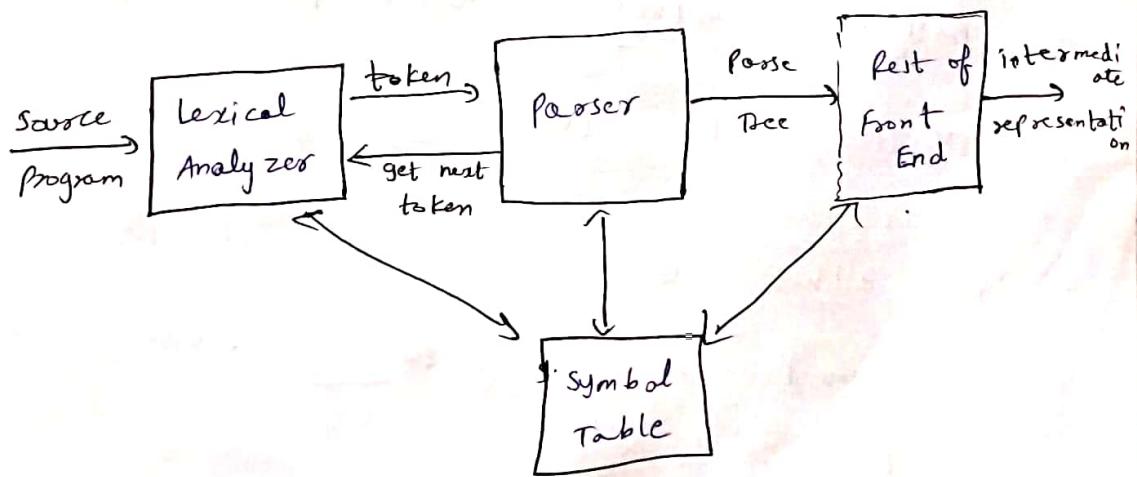
Unit - II

Syntax Analysis

1) Introduction

1.1 Role of Parser

- The parser obtains a string of tokens from the lexical analyzer and verifies that the string of tokens can be generated by the grammar for the source language.
- Parser is expected to report any syntax errors in an intelligible fashion and to recover from commonly occurring errors to continue processing the remainder of the program.



→ There are three types of parsers for grammars: universal

Top-down

Bottom-up

→ universal parsing methods such as Cocke-Younger-Kasami algorithm and Earley's algorithm can parse any grammar.

→ These are inefficient to use in production compilers.

→ The methods used in Compilers can be classified as either Top-down (or) Bottom-up

→ Top-down methods build parse trees from the top (root) to the bottom (Leaves)

→ Bottom-up methods build parse trees from Leaves (bottom) to the root (top)

→ In either case, input to the parser is scanned from left to right, one symbol at a time.

1.2 Representative Grammars

→ Constructs that begin with keywords

like while or int, are, easy to parse,

because the keyword guides the choice

of the grammar production that must be

applied to match the input.

Because of associativity and precedence the

expressions are difficult to parse.

Eg :-

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

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

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

The above expression grammar belongs to class of LR grammars that are suitable for Bottom up parsing

This grammar can be adapted to handle additional operators and additional

levels of precedence -

The above grammar is Left recursive

So it cannot be used for Top-down parsing.

Non left Recursive grammar

$$E \rightarrow T E'$$

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

$$T \rightarrow F T'$$

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

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

→ Non left Recursive grammar can be used for Top-down parsing.

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

→ The above grammar treats * and + alike, so it is useful for handling ambiguities during parsing.

1.3 Syntax Error Handling

→ A compiler is expected to assist the programmer in locating and tracking

down errors.

→ Planning the Error Handling right from the start can both simplify the structure of compiler and improve its handling of errors

→ Common programming errors

(i) Lexical errors include misspelling of identifiers, keywords or operators.

(ii) Syntactic errors include misplaced semicolon (or) extra (or) missing braces ("{}" or "{}")

(iii) Semantic errors include type mismatches between operators and operands - (^{undeclared variable})

(iv) Logical errors can be anything from incorrect reasoning on the part of programmer to the use ~~of~~ in a program " \neq " instead of ~~com~~ " $=$ ".

→ The precision of parsing methods allows syntactic errors to be detected very efficiently -

→ Parsing Methods such as LL and LR methods detect an error as soon as

possible, that is, when the stream of tokens from the lexical analyzer cannot be passed further according to grammar for the language.

→ The error handling

- (i) report the presence of errors clearly and accurately.
- (ii) Recover from each error quickly enough to detect subsequent errors.
- (iii) Add minimal processing of tokens to the overhead to the remaining part of correct programs.

1-4 Error Recovery

(i) Panic - Mode Recovery

→ On Discovering an error, the parser discards input symbols one at a time until one of designated set of tokens is found.

→ Synchronizing Tokens are usually delimiters such as ~~several~~ semicolon (;) or $\{$, whose

role in source program is clear and unambiguous.

→ while panic mode correction often skips a considerable amount of input without

Checking for additional errors (not go into infinite loop)

→ Advantage (not go into infinite loop)

(ii) Phrase-Level Recovery

→ on Discovering an error, a parser may perform local correction on remaining input, that is, it may replace a prefix of remaining input by some string that allows the parser to continue.

→ A typical local correction is to replace a comma, by a semicolon, delete an extra semicolon, or insert a missing semicolon.

→ It can be used in several error-repairing compilers, as it can correct any input string.

→ Drawback Difficulty it has in coping with

situations in which the actual error has occurred before the point of detection.

(iii) Error Productions

- By anticipating common errors that might be encountered, we can augment the grammar for the language at hand with ~~the~~ productions that generate erroneous constructs.
- A parser constructed from a grammar augmented by these error productions detects the anticipated errors when an error production is used during parsing.

(iv) Global Correction

- The aim is to make as few changes possible in processing an incorrect input string.
- There are algorithms for choosing a minimal sequence of changes to obtain a globally

least-cost correction
Eg: Given an incorrect input string x and a grammar G , these algorithms will find a parse tree for a related string y , such

that the no of insertions, deletions and changes of tokens required to transform x into y is as small as possible.

→ These methods are costly in terms of time and space

1.2 Context Free Grammars

→ Grammars are used to describe the syntax of programming language constructs like expressions and statements.

Eg:

stmt → if (expr) stmt else stmt

The formal Definition of a Context-Free Grammar

→ CFG consists of terminals, non-terminals, a start symbol and productions.

(i) Terminals are the basic symbols from which strings are formed. The term "Token" is the synonym for "Terminal".

Eg: if, else, "(", ")" are terminals in the above grammar.

(ii) NonTerminals are syntactic variables that denote set of strings.

In the above eg $\text{stmt} \Rightarrow \text{expr}$ by ~~variables~~ are
non Terminals

→ The set of strings denoted by NonTerminals help define the language generated by grammar.

(iii) In a grammar, one Non-Terminal is distinguished as start symbol, and set of strings it denotes is the language generated by grammar.

(iv) The productions of a grammar specify the manner in which the terminals and Non-Terminals can be combined to form strings.

Each production consists of called the head (or)

(a) A NonTerminal left side of the production; this production defines some of strings denoted by head.

(b) The symbol →

(c) A body on right side consisting of zero or more terminals and nonterminals

The components of the body describe one way in which strings of the nonterminal at the head can be constructed,

Grammar for simple arithmetic Expressions

Eg:-

$\text{expr} \rightarrow \text{expr} + \text{term}$

$\text{expr} \rightarrow \text{expr} - \text{term}$

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

$\text{term} \rightarrow \text{term} * \text{factor}$

$\text{term} \rightarrow \text{term} / \text{factor}$

$\text{term} \rightarrow \text{factor}$

$\text{factor} \rightarrow (\text{expression})$

$\text{factor} \rightarrow \text{id}$

Terminals: id, +, -, *, /, (,)

Notational Conventions

(1) These symbols are Terminals:

(a) Lower case letters (a, b, c)

(b) operator symbols (+, * ---)

(c) punctuation symbols (", ", " ---)

(d) digits 0, 1, ..., 9

(e) Bold face strings (if id)

(2) These symbols are NonTerminals:

(a) Upper case letters (A, B, C)

(b) The letter "S" (start symbol)

(c) Lower case italic names

(d) Superscript symbols (F, i, F)

(expr or stnt)

can be used to
represent non-terminals
for contexts

- (3) uppercase letters (x, y, z) represent Grammars
- (4) symbols (either NT or T)
- (5) lowercase letters in alphabets ($u, v, -z$) represents
- (6) strings of terminals (α, β, γ) represent strings
- (7) lowercase greek letters ($\alpha \rightarrow d$) α is head & d is body of grammar symbols. ($A \rightarrow d_1, A \rightarrow d_2, \dots, A \rightarrow d_K$)
- (8) A set of productions with common head may be written as $A \rightarrow d_1 | d_2 | \dots | d_K$
- (9) unless stated otherwise, the head of first production is start symbol,

Eg:

$E \rightarrow E + T E - T T$	E - start symbol
$T \rightarrow T * F T / F F$	$F \& T$ - Non Terminals
$F \rightarrow (E) id$	$id, +, -, *, /, (,),$ id - Terminals

Derivations

- The construction of a parse tree can be made precise by taking a derivational view, in which productions are treated as rewriting rules.
- Beginning with start symbol, each rewriting step

replaces a Non-Terminal by body of one its products.

→ The Derivational view corresponds to the top-down construction of the Parse-Tree, but precision afforded by derivations will be especially helpful when bottom-up parsing is discussed.

→ Bottom-up Parsing is related to class of derivations known as "Rightmost" derivations, in which rightmost NT is rewritten at each step.

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

$$E \Rightarrow -E \quad (E \text{ derives } -E)$$

$$E \Rightarrow -E \Rightarrow -(E) \Rightarrow \text{id}$$

→ we call ~~this~~ such a sequence of replacements a derivation of (id) from E.

$\alpha \xrightarrow{*} \gamma$ for any string α , and

If $\alpha \xrightarrow{*} \beta$ and $\beta \Rightarrow \gamma$, then $\alpha \xrightarrow{*} \gamma$

$\xrightarrow{*}$ (derives in zero or more derivations)

$\xrightarrow{+}$ (derives in one or more steps)

If $s \xrightarrow{*} \alpha$, where s is start symbol of grammar G ,

- \mathcal{L} is Sentential form of G ;
- Sentential form may contain both terminals and nonterminals and may be empty.
 - A sentence of G is a sentence form with no nonterminals.
 - The language generated by a grammar is its set of sentences.
 - A string of terminals w is in $L(G)$, the language generated by G , if and only if w is a sentence of G ($s \Rightarrow^* w$).
 - A language generated by a C-F-G is called Context free language.
 - If two grammars (C-F-G) generate same language, the grammars are said to be equal.
- Eg: string $-(id + id)$ is sentence of G because there is a derivation:

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

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

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

→ Each N.T is replaced by the same body in the two derivations, but the order of replacements is different.

(1) In left most derivations, the leftmost Non-Terminal in each sentential is always chosen.

If $\alpha \Rightarrow \beta$ is a step in which the leftmost non terminal in α is replaced by β .

(2) In right most derivations, the rightmost Non-terminal is always chosen, $\alpha \Rightarrow \beta_{rm}$.

→ Right most Derivations are also called

Canonical derivations.

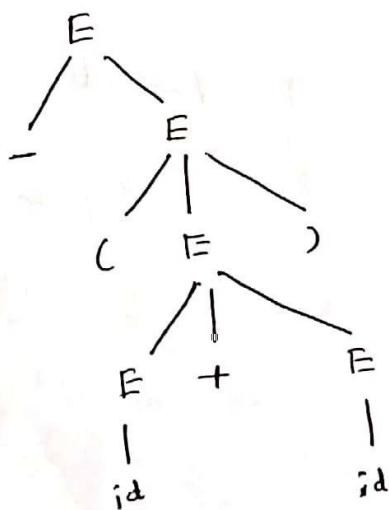
Parse Trees and Derivations

→ A parse Tree is a graphical representation of a derivation that filters out the order in which productions are applied to replace Non-Terminals.

→ Each interior node of a parse tree represents the application of a production.

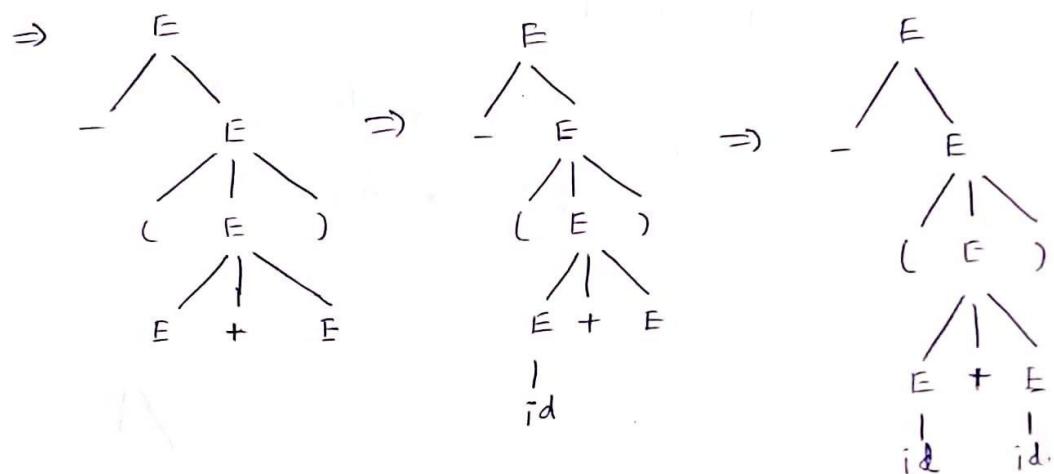
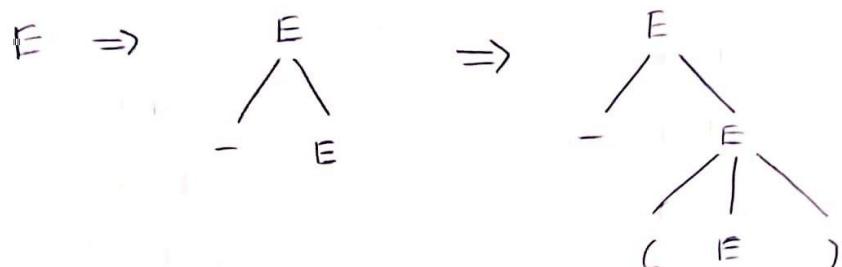
- The interior node is labeled with the nonterminal 'A' in the head of production.
- The children of the node are labeled, from left-to-right, by symbols in the body of production by which this 'A' was replaced during the derivation.
- The leaves of a parse tree are labeled by nonterminals (or) Terminals and read from left to right, constitute a sentential form, called yield (or) Frontier of the Tree.

Eg:



Parse tree for $-(id + id)$.

Ex: sequence of parse trees for derivations.



Ambiguity

- A grammar produces more than one parse tree for some sentence is said to Ambiguous.
- An Ambiguous Grammar is one that produces more than one leftmost derivation (or) more than one rightmost derivation for the same sentence.

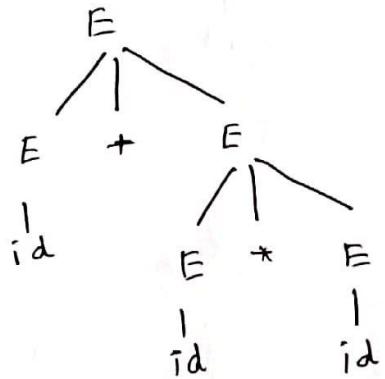
Eg^r sentence $id + id * id$ has two left most derivations

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

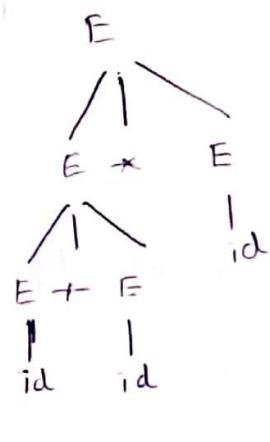
$$\begin{aligned} E &\Rightarrow E * E \\ &\Rightarrow E + E * E \\ &\Rightarrow id + E * E \\ &\Rightarrow id + id * E \\ &\Rightarrow id + id * id \end{aligned}$$

→ For most parsers, it is desirable that the grammar be made unambiguous.

Eg^r parse trees



(a)



(b)

Two parse trees for $\underline{id + id * id}$

Derivations

Verifying Language generated by Grammar

→ A Grammar (G) generates a Language L

has two parts:

i) show that every string generated by G is in L

ii) That every string in L can indeed be generated by G .

Eg: $s \rightarrow (s) s | \epsilon$

The above grammar generates all strings of balanced parenthesis, and only such strings.

$$S \xrightarrow{\text{em}} (s) s \xrightarrow[\text{em}]{\text{em}} (\pi) S \xrightarrow[\text{em}]{\text{em}} (\pi) Y$$

$\rightarrow (\pi) Y$ must be balanced.

Context-Free Grammars vs Regular Expressions

Context-Free Grammars are more powerful notation

→ Grammars are more powerful than regular Expressions.

→ Every construct that can be described by

a regular Expression can be described by a grammar but not vice versa.

Eg:- $R.E = (a|b)^* abb$ and Grammar

$$A_0 \rightarrow a A_0 | b A_0 | a A_1$$

$$A_1 \rightarrow b A_2$$

$$A_2 \rightarrow b A_3$$

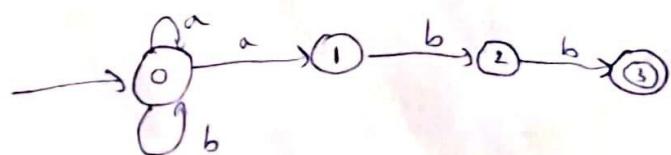
$$A_3 \rightarrow \epsilon$$

→ The grammar was constructed from

N.F.A

- 1) For each state i of NFA, create a non-terminal A_i
- 2) If state i has a transition to state j on input a , add the production $A_i \rightarrow a A_j$
If state i goes to state j on input b ,
add the production $A_i \rightarrow b A_j$
- 3) If i is an Accepting state, add $A_i \rightarrow \epsilon$
- 4) If i is the start state, make A_i be start symbol of grammar.

N.F.A for $(a|b)^* abb$



writing a grammar

→ several Transformations that could be applied to get grammar more suitable for parsing

→ Techniques are (i) Elimination of Ambiguity
(ii) Left-Recursion Elimination
(iii) Left-factoring -

→ The above techniques are useful for rewriting the grammars so they become suitable for top-down Parsing

lexical vs syntactic Analysis

(1) separating the syntactic structure of a language into lexical and Non-Lexical parts provide a convenient way of modularizing the front end of compiler into two manageable sized components.

(2) The lexical Rules of a language are quite simple and to describe them we don't need a notation as powerful as grammar

(3) Regular Expressions generally provide a more concise and easier-to-understand notation for tokens than grammars.

(4) More efficient lexical Analyzers can be constructed automatically from the regular Expressions than from Arbitrary Grammars.

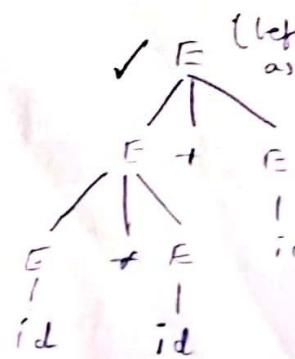
→ Regular Expressions are most useful for describing the structure of constructs such as identifiers, constants, keywords and whitespaces
→ Grammars, are most useful for describing Nested structures such as balanced parenthesis, matching begin-end's, corresponding if then else and soon

Eliminating Ambiguity

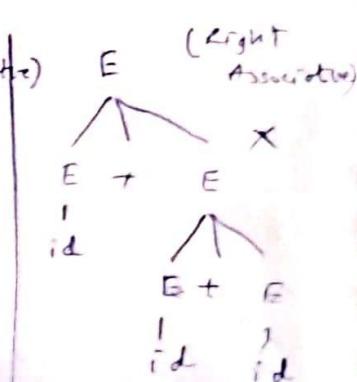
Eg: $E \rightarrow E + E$

/E * E

/ id

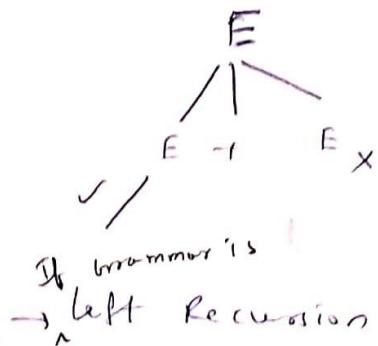


Not following Associativity
(Left to Right)



precedence solution
Solutions

solution restrict ^{to} left associativity using left recursion



so grammar will be

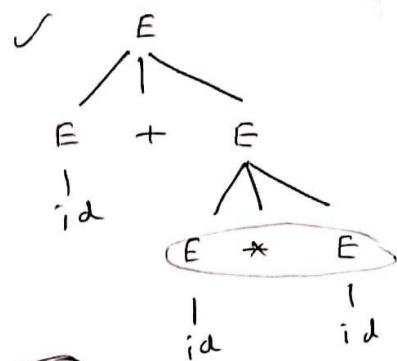
$$E \rightarrow E + id \mid id$$

operator is

then \rightarrow Left associativity.

not following (precedence)

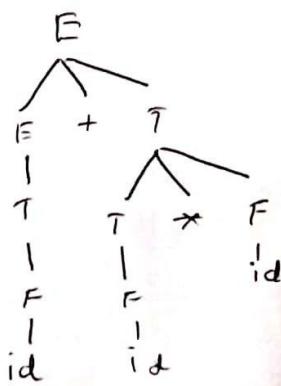
$id + id * id$



precedence
solution
Solution

use levels. (Higher levels '+' and '*' in lower levels)

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow id \end{aligned}$$



→ → Low level
+ → high level

→ \rightarrow Highest precedence should be at low level

Eg:

$bExp \rightarrow bExp \text{ or } bExp$
| $bExp$ and $bExp$
| not $bExp$
| True
| True

unambiguous grammar

$E \rightarrow E \text{ or } F/F$
 $F \rightarrow F \text{ and } G/G$
 $G \rightarrow \text{Not } G \mid \text{True} \mid \text{False}$

Eg:

Ambiguous

$R \rightarrow R + R$
| $R R$
| a^*
| a
| b
| c

unambiguous
Respected

$E \rightarrow E + T \mid T$
 $T \rightarrow TF \mid F$
 $F \rightarrow F^* \mid a \mid b \mid c$

Eg:

$A \rightarrow A \$ B \mid B$
 $B \rightarrow B \# C \mid C$
 $C \rightarrow C @ D \mid D$
 $D \rightarrow d$

\$, #, @ → left associative

@ > # > \$ (precedence)

Eg:
 $E \rightarrow E * F \mid F + E \mid F$
 $F \rightarrow F - F \mid id$

* → left associative
+ → right associative

- → (* ≡ +)

Left Recursion

→ A grammar is Left Recursion if it has a Non-Terminal A such that there is a derivation $A \xrightarrow{+} A\alpha$ for some string α .

Direct Left Recursion

$$A \rightarrow A\alpha$$

Indirect Left Recursion

$$S \rightarrow A\alpha$$

$$A \rightarrow S\beta$$

→ Top-down parsing methods cannot handle left-recursive grammars, so transformation is needed to eliminate left recursion.

(Why?)

→ $A \rightarrow A\alpha | \beta$ → Left Recursive Grammar

$A \rightarrow \beta A'$
 $A' \rightarrow \alpha A' | \epsilon$ → Non-Left Recursive productions

→ If there are multiple A productions

$A \rightarrow A\alpha_1 | A\alpha_2 | A\alpha_3 | \dots | A\alpha_m | \beta_1 | \beta_2 | \dots | \beta_n$

where no β_i begins with A .

$$\left. \begin{array}{l} A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A' \\ A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_m A' \mid \epsilon \end{array} \right\}$$

Non Recursive Grammar

Advantage

→ we are able to generate the same language even after removing left recursion.

Disadvantage

→ The above procedure only eliminates Direct Left Recursion but not Indirect ^{left} Recursion (Not Immediate)

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

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

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

↓
No L.R.

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

L.R

$$E \rightarrow T E'$$

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

L.R
α β

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

$$T \rightarrow F T'$$

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

Final Grammar

$$\begin{array}{l} E \rightarrow T E' \\ E' \rightarrow + T E' \mid \epsilon \\ T \rightarrow F T' \\ T' \rightarrow * F T' \mid \epsilon \\ F \rightarrow (E) \mid id \end{array}$$

Eg:- $S \rightarrow S_0 \underbrace{S_1 S_2}_{\alpha} | \underbrace{\epsilon}_{\beta}$ (L-R)
 Not L-R

$$S \rightarrow \alpha S'$$

$$S' \rightarrow \alpha S_1 S_2 S' | \epsilon$$

Eg:- $L \rightarrow L, \underbrace{S}_{\alpha} | \underbrace{S}_{\beta}$ (L-R)
 Not L-R

$$L \rightarrow S L'$$

$$L' \rightarrow , S L' | \epsilon$$

Eg:- $S \rightarrow S X$ (L-R)
 $S \rightarrow S S b$
 $S \rightarrow X S$
 $S \rightarrow a$

$$S \rightarrow S \underbrace{X}_{\alpha_1} | S \underbrace{S b}_{\alpha_2} | \underbrace{X S}_{\beta_1} | \underbrace{a}_{\beta_2}$$

Not L-R

$$S \rightarrow X S S' | a S'$$

$$S' \rightarrow X S' | S b S' | \epsilon$$

Eg:- $A \rightarrow A \underbrace{A}_{\alpha_1} | \underbrace{A b}_{\alpha_2}$ (L-R)

Not L-R

$$A' \rightarrow A A' | b A' | \epsilon$$

Left factoring (Eliminating Non-Determinism)

→ Left factoring is a grammar transformation that is useful for producing a grammar suitable for predictive, or Top-down parsing.

→ when the choice between two alternative A-productions is not clear, we may be able to rewrite the productions to defer the decision until enough of input has been seen that we can make right choice.

Grammar

$$A \rightarrow \alpha \beta_1 | \alpha \beta_2 | \dots | \alpha \beta_n | \gamma, \text{ where } \gamma$$

represents all alternatives that don't begin with α

α

$$A \rightarrow \alpha A' | \gamma$$

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

Eg:-

$$\begin{array}{l} S \rightarrow i E T S \\ \cancel{S} \rightarrow \cancel{i E T S} c S \\ | a \\ E \rightarrow b \end{array}$$

After left Factoring

$$\begin{array}{l} S \rightarrow i E T S' | a \\ S' \rightarrow \epsilon | c S \\ E \rightarrow b \end{array}$$

Egr-

→ Eliminating Ambiguity will not eliminate left factoring

Eg:

$$S \rightarrow \underline{a} S S b S$$

$$\quad | a \quad S a S b$$

$$\quad | \underline{a} \underline{b} b$$

$$\quad | b$$

$$S \rightarrow a S' | b$$

$$S' \rightarrow \underline{s} s b s | \underline{s} a s b | b b$$

After left factoring,

$$S \rightarrow a S' | b$$

$$S' \rightarrow s s'' | b b$$

$$s'' \rightarrow s b s | a s b$$

Eg:-

$$S \rightarrow \underline{b} S S a a S$$

$$\quad | \underline{b} S S a S b$$

$$\quad | \underline{b} S b$$

$$\quad | a$$

$$S \rightarrow b S S' | a$$

$$S' \rightarrow \underline{s} a S b | b$$

$$\quad | \underline{s} a a S$$

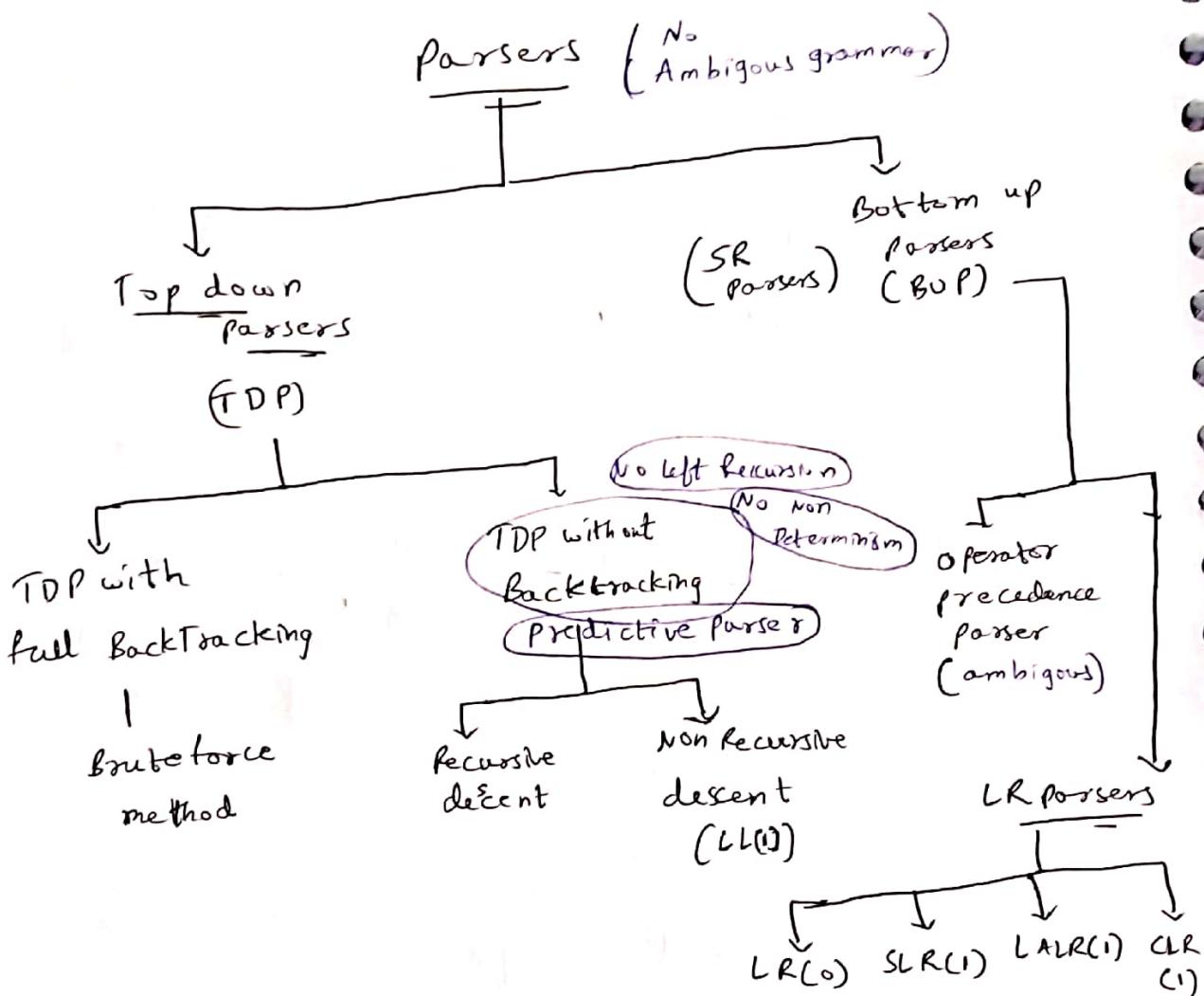
After left factoring

$$S \rightarrow b S S' | a$$

$$S' \rightarrow s a S'' | b$$

$$S'' \rightarrow a S | \$ b$$

Types of parsers



Eg: $S \rightarrow a A B e$

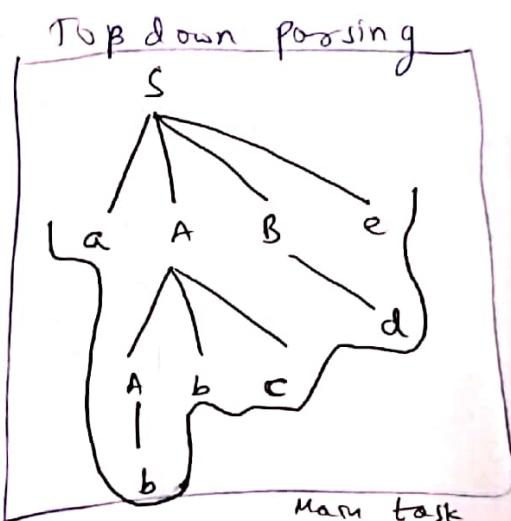
$A \rightarrow A b c \mid b$

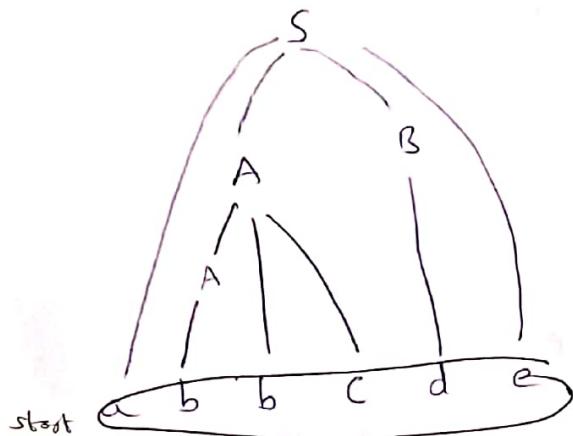
$B \rightarrow d$

w abbCde

leftmost derivations

$S \Rightarrow a A B e$
 $\Rightarrow a A b c B e$
 $\Rightarrow a b b c B e$
 $\Rightarrow a b b c d e$





Reverse of Right most derivation

Bottom up parsers

Main Task

(when to reduce)

$$\begin{aligned}
 S &\Rightarrow [a A B] e \\
 &\Rightarrow a [A B] e \\
 &\Rightarrow a [A b c] d e \\
 &\Rightarrow a [b] b c d e
 \end{aligned}$$

Top Down Parsing

→ Top-Down parsing can be viewed as the problem of constructing a parse tree for the input string, starting from the root and creating the nodes of parse tree in preorder.

→ Top-Down Parsing can be viewed as finding the leftmost derivation for the input string.

→ At each step of top-down parse, the key

problem is that of determining the production to be applied for a nonterminal, say A.

→ Once an A-production is chosen, the rest

of the parsing process consists of "matching"

the terminal symbols in the production body with input string.

→ Top Down parser is constructed for the grammar if it is free from Ambiguity

↳ Recursion

Recursive Descent parser

→ Parse is constructed from Top

→ Input is read from left to right.

→ The input is recursively passed for

preparing a parse tree with or without Backtracking.

Backtracking

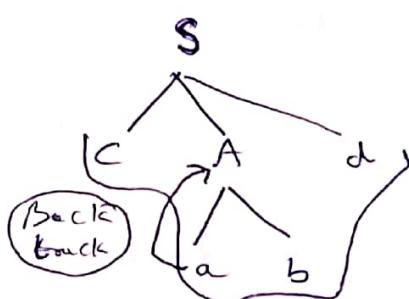
→ The parse tree is started from root node &

input string is matched against the production rules for replacing them.

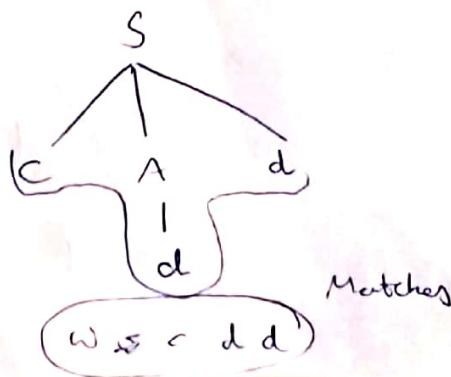
Ex: $S \rightarrow CAd$

$A \rightarrow ab|d$

$$w = cd\boxed{d}$$



$$w = cabd$$



Limitations

→ If the given grammar has more no. of alternatives then the cost of Back-Tracking is high.

LL Parser

→ It accepts LL grammar

→ It is denoted by $LL(k)$

Scans input from left to right

left most derivation

no. of look aheads $k = 1$

$$\therefore LL(k) = LL(1)$$

→ A grammar "G" is $LL(1)$ if there are two distinct productions

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

(i) For no terminal α and β derive string beginning with α .

(ii) Atmost one of α and β can derive Empty string.

(iii) If $\beta \Rightarrow^* \epsilon$ then α does not derive any string beginning with a terminal in $Follow(A)$

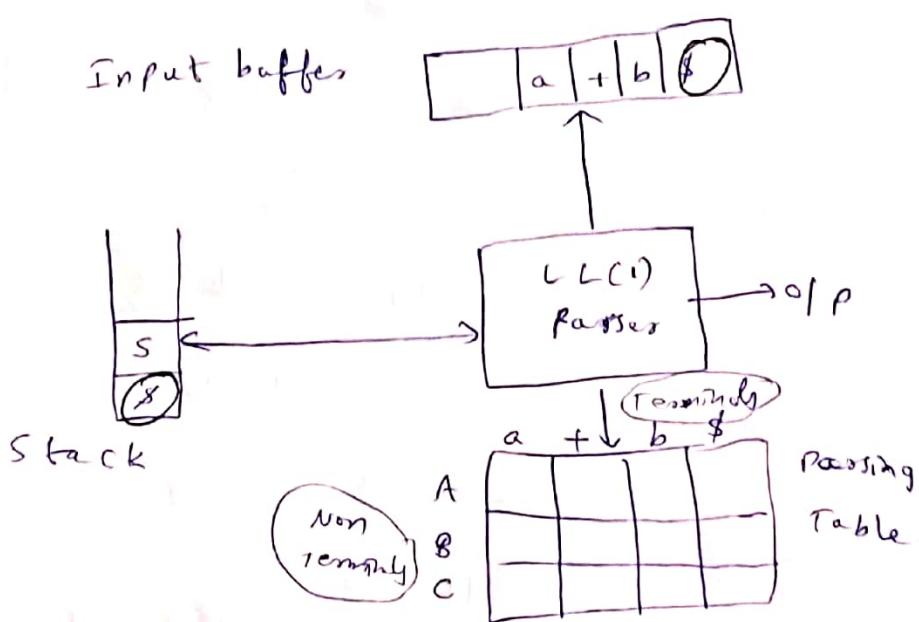
→ LL(1) uses data structure

(i) Input buffer

(ii) stack

(iii) parsing Table

Structure of LL(1)



Construction of predictive LL(1)

① FIRST() / Leading()

FOLLOW() / Trailing()

② Predictive parsing table by using FIRST & FOLLOW functions

③ Parse the input string with the help of the table.

Eg:- $S \rightarrow a A B C D$

$$A \rightarrow b$$

$$B \rightarrow c$$

$$C \rightarrow d$$

$$D \rightarrow e$$

$$\text{First}(S) = a$$

$$\text{First}(A) = b$$

$$\text{First}(B) = c$$

$$\text{First}(C) = d$$

$$\text{First}(D) = e$$

To compute FIRST(x), apply following rules until no more terminals (ϵ) can be added to any FIRST set.

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

ii) If x is nonterminal, and $x \rightarrow y_1 y_2 \dots y_k$ is a production for some $k \geq 1$, then place $\underline{\underline{a}}$ in $\text{FIRST}(x)$ if for some i , $\underline{\underline{a}}$ is in $\text{FIRST}(y_i)$; and $\underline{\underline{\epsilon}}$ is in all of $\text{FIRST}(y_1), \dots, \text{FIRST}(y_{i-1})$

$y_i \dots y_{i-1} \xrightarrow{*} \epsilon$. If $\underline{\underline{\epsilon}}$ is in ~~FIRST~~ $\text{FIRST}(y_i)$

for $j = 1, 2, \dots, k$, then add $\underline{\underline{\epsilon}}$ to $\text{FIRST}(x)$.

iii) If $x \rightarrow \epsilon$ is a production, then add ϵ to $\text{FIRST}(x)$.

FIRST(x) -

→ To compute FOLLOW(A) for all nonterminals A, apply following rules until nothing can be added to any FOLLOW set.

i) place \$ in FOLLOW(S), where S is the start symbol. (\$ is the right end marker in input)

ii) If there is a production $A \rightarrow \alpha B \beta$, then everything in FIRST(β) except ϵ is in FOLLOW(B).

(iii) If there is a production $A \rightarrow \alpha B$ or $A \rightarrow \alpha B \beta$ where FIRST(β) contains ϵ , then everything in FOLLOW(A) is in FOLLOW(B).

Eg:

Grammar	
$S \rightarrow ABCDE$	
$A \rightarrow a \epsilon$	
$B \rightarrow b \epsilon$	
$C \rightarrow c$	
$D \rightarrow d \epsilon$	
$E \rightarrow e \epsilon$	

	First	Follow
S	{a, b, c}	{\$}
A	{a, ϵ }	{b, c}
B	{b, ϵ }	{c}
C	{c}	{d, e, \$}
D	{d, ϵ }	{e, \$}
E	{e, ϵ }	{\$}

~~Suppose~~ A → BC ~~FOLLOW(C) = FOLLOW(A)~~

Eg:

Grammar

Eg:

Eg:

Eg:

Eg:

marks A,
added

Egr

Grammar

$$S \rightarrow Bb \mid Cdd$$

$$B \rightarrow ab \mid \epsilon$$

$$C \rightarrow cc \mid \epsilon$$

	First	Follow
S	{a, b, c, d}	{\\$}
B	{a, \epsilon}	{b}
C	{c, \epsilon}	{d}

Egr

$$E \rightarrow T E'$$

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

$$T \rightarrow FT'$$

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

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

	First	Follow
E	{id, ()}	{+, *, (,)}
E'	{+, \epsilon}	{+, \\$}
T	{id, ()}	{+, \\$, *)}
T'	{*, \epsilon}	{+, \\$, *)}
F	{id, ()}	{*, +, \\$, *)}

Egr

$$S \rightarrow A \overset{a}{\underset{\epsilon}{\mid}} B \overset{b}{\underset{\epsilon}{\mid}} C \overset{c}{\underset{\epsilon}{\mid}} d \overset{d}{\underset{\epsilon}{\mid}} B \overset{a}{\underset{\epsilon}{\mid}} S$$

$$A \rightarrow d \overset{a}{\underset{\epsilon}{\mid}} a \overset{b}{\underset{\epsilon}{\mid}} C \overset{c}{\underset{\epsilon}{\mid}}$$

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

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

	First	Follow
S	{d, g, h, b, a}	{\\$}
A	{d, g, h, \epsilon}	{h, g, \\$}
B	{g, \epsilon}	{a, s, h, g}
C	{h, \epsilon}	{h, g, s, b}

Egr

$$S \rightarrow a A \overset{b}{\underset{\epsilon}{\mid}} B \overset{b}{\underset{\epsilon}{\mid}} b$$

$$A \rightarrow c \mid \epsilon$$

$$B \rightarrow d \mid \epsilon$$

	First	Follow
S	{a}	{\\$}
A	{c, \epsilon}	{d, b}
B	{d, \epsilon}	{b}

Egr

$$S \rightarrow a B \overset{D}{\underset{\epsilon}{\mid}} h$$

$$B \rightarrow c C$$

$$C \rightarrow b \overset{E}{\underset{\epsilon}{\mid}} f \mid \epsilon$$

$$D \rightarrow E \overset{F}{\underset{\epsilon}{\mid}} f$$

$$E \rightarrow g \mid \epsilon$$

$$F \rightarrow f \mid \epsilon$$

	First	Follow
S	{a}	{\\$}
B	{c}	{g, f, h}
C	{b, \epsilon}	{g, f, h}
D	{g, f, \epsilon}	{h}
E	{g, \epsilon}	{f, h}
F	{f, \epsilon}	{h}

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

(1) The ~~green~~ productions should be placed in first of Right hand side in the Table.

First of Right hand side is $A \rightarrow E$ like then it should be placed in Follow of (A) in the Table.

Eg: $S \rightarrow (S) | E$

Parsing Table		()	\$
S	$S \rightarrow (S)$	$S \rightarrow E$	$S \rightarrow E$	

Eg: $S \rightarrow AaAb | BbBa$
 $A \rightarrow E$
 $B \rightarrow E$

	First	Follow
S	{a}	{\$}
A	{e}	{a, b}
B	{e}	{b, a}

	a	b	\$
S	$\frac{S \rightarrow AaAb}{A \rightarrow E}$	$\frac{S \rightarrow BbBa}{B \rightarrow E}$	
A	$A \rightarrow E$	$A \rightarrow E$	
B	$B \rightarrow E$	$B \rightarrow E$	

Parse Table

stack implementation by using

Input string = id + id * id *

Matched	STACK	Input	Action
	E \$	id + id * id *	output $E \rightarrow TE'$
	T E' \$	id + id * id *	output $T \rightarrow FT'$
	FT' E' \$	id + id * id *	output $F \rightarrow id$
	id T' E' \$	id + id * id *	match id
id	T' E' \$	+ id * id *	output $T' \rightarrow E$
id	E' \$	+ id * id *	output $E' \rightarrow TE'$
id	+ TE' \$	+ id * id *	match +
id +	T E' \$	id * id *	output $T \rightarrow FT'$
id +	FT' E' \$	id * id *	output $F \rightarrow id$
id +	id T' E' \$	* id *	match id
id + id	T' E' \$	* id *	output $T' \rightarrow FT'$
id + id	* FT' E' \$	* id *	match *
id + id *	FT' E' \$	id *	output $F \rightarrow id$
id + id *	id T' E' \$	id *	match id
id + id + id	T' E' \$	\$	

$id \rightarrow id * id$	$E' \neq$	\neq	output $T' \rightarrow c$
$id + id * id$	\neq	\neq	output $E' \rightarrow c$

Error Recovery in LL(1) [Predictive Parsing]

→ An error is detected during predictive parsing

When the terminal on top of the stack does not match the next input symbol (or) when

NonTerminal A is on top of the stack,

a is the next input symbol, and $M[A, a]$

is error (i.e. the parsing Table entry is empty).

Panic Mode

→ It is based on the idea of skipping symbols on the input until a token in a selected set of synchronization tokens appears.

→ The sets should be chosen so that the parser recovers quickly from errors that are likely to occur in practice.

→ Some

(i) A:

in F

Non te

(ii) It

as sy

(iii) If

the sy

it m

to A

in th

(iv) If

string,

used a

This

that

sec

(v) If

matche

Issue

insert

→ Some heuristics are as follow:

(i) As a starting point, place all symbols in FOLLOW(A) into the synchronizing set ~~for~~ for Nonterminal A.

(ii) It is not enough to use FOLLOW(A) as synchronizing set ~~for~~ for A.

(iii) If we add symbols in FIRST(A) to the synchronizing set for nonterminal A, then it may be possible to resume parsing according to A if a symbol in FIRST(A) appears in the input.

(iv) If a nonterminal can generate that empty string, then the production deriving ϵ can be used as a default.

This approach reduces the number of nonterminals. This approach reduces the number of nonterminals. That have to be considered during errors, that have to be considered during errors, recovery.

(v) If a terminal on top of stack cannot be matched, a simple idea is to pop the terminal, issue a message saying that the terminal was inserted, and continue parsing.

→ The Parsing Table is to be checked

i) If parser looks up entry M[A]

that it is blank, then the input symbol is skipped.

ii) If the entry is "synch", then the

token on top of the stack is popped

attempt to resume parsing.

iii) If token on top of the stack is

the input symbol, then we pop the token

stack.

Parsing Table updated

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

Synchronization tokens added to the

$$\text{Input string} = \boxed{\underline{id} * + id}$$

STACK	INPUT	REMARK
$E \$$	$\underline{id} * + id \$$	error, skip)
$\underline{E} \$$	$\underline{id} * + id \$$	$E \rightarrow T E'$
$\underline{T} E' \$$	$\underline{id} * + id \$$	
$\underline{F} T' E' \$$	$\underline{id} * + id \$$	by
$\underline{id} T' E' \$$	$\underline{*} + id \$$	
$\underline{T'} E' \$$	$\underline{*} + id \$$	
$\underline{* F} T' E' \$$	$\underline{*} + id \$$	
$\underline{F} T' E' \$$	$\underline{+} id \$$	error, M[F, +] = synch
$\underline{T'} E' \$$	$\underline{+} id \$$	F has been popped
$\underline{E'} \$$	$\underline{+} id \$$	
$\underline{+ T} E' \$$	$\underline{+} id \$$	
$\underline{T} E' \$$	$\underline{id} \$$	
$\underline{F} T' E' \$$	$\underline{id} \$$	
$\underline{id} T' E' \$$	$\underline{id} \$$	
$T' E' \$$	$\$$	
$E' \$$	$\$$	
$\$$	$\$$	

Reduction:

Phrase - Level - Recovery

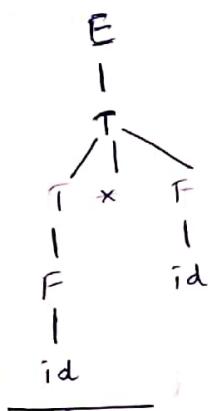
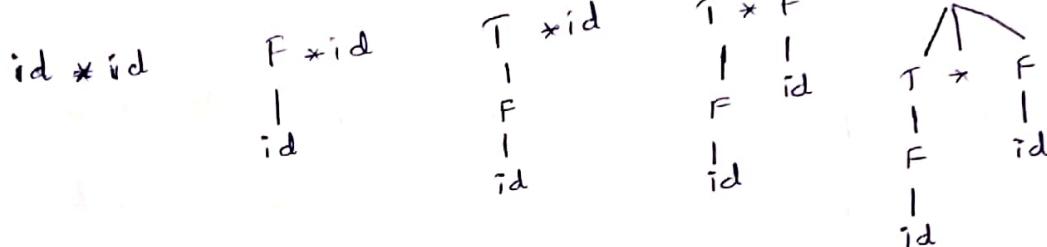
- It is implemented by filling in the blank entries in the predictive parsing table with pointers to error routines.
- These routines may change, insert (or) delete symbols on input and issue appropriate error messages.
- The steps carried out by the parser might then not correspond to the derivation of any word in Language at all.
- we must ensure that there is no possibility of infinite loops.
- checking that any recovery action eventually results in an input symbol being consumed is a good way to protect against ~~such~~ loops;

Bottom up Parsing

- In bottom up parsing method, the input string is taken first and we try to reduce this string with the help of grammar and try to obtain the start symbol.
- The process of parsing stops(halts) successfully as soon as we reach to start symbol.
- The parse tree is constructed from bottom to up that is from leaves to root.
- In this process, the input symbols are placed at the leaf nodes after successful parsing.
- The bottom up parse tree is created starting from ~~to~~ leaves, the leaf nodes together are reduced further to internal nodes, these internal nodes are further reduced and eventually a root node is obtained.
- In this process, the parser tries to identify RHS of production rule and replace it by corresponding LHS. This is called Reduction.

→ The task of bottom up parsing is to find the productions that can be used for reduction.

Eg:



Reductions

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

→ At each reduction step, a specific substring matching the body of a production is replaced by

the nonterminal at the head of that production.

→ The key decisions during bottom-up parsing are about when to reduce and about what production to apply, as parse proceeds.

Handle pruning

→ Bottom-up parsing during a left-to-right scan of the input constructs a rightmost derivation in reverse.

→ "Handle" is a substring that matches the body of a production, and whose reduction represents one step along the reverse of a rightmost derivation

<u>Right sentential form</u>	<u>Handle</u>	<u>Reducing Production</u>
$id_1 * id_2$	id_1	$F \rightarrow id$
$F * id_2$	F	$T \rightarrow F$
$T * id_2$	id_2	$F \rightarrow id$
$T * F$	$T * F$	$E \rightarrow T * F$

→ A handle of a right sentential form γ is a production $A \rightarrow \beta$ and a position of γ where the string β may be found and replaced by A to produce the previous right

Sentential form in a rightmost derivation of r .

Eg:

$$S \xrightarrow{r_m} aA_\underline{B}e \xrightarrow{r_m} a_\underline{A}de \xrightarrow{r_m} a_\underline{\Lambda}bcde \rightarrow a_\underline{b}bcde$$

Right sentential form

$$a \begin{matrix} b \\ | \\ \Lambda \end{matrix} bcde : r = abcd, \Lambda \rightarrow^{\beta} b, \text{Handle} = b$$

$$a \begin{matrix} b \\ | \\ \Lambda \end{matrix} bcde : r = a\underline{Ab}cde, \Lambda \rightarrow Abc, \text{Handle} = Abc$$

$$a \begin{matrix} b \\ | \\ \Lambda \end{matrix} de : r = aA_\underline{de}, B \rightarrow d, \text{Handle} = d$$

$$a A Be : r = a A Be, S \rightarrow a A Be, \text{Handle} = a A Be$$

grammar

$$S \rightarrow aABe$$

$$A \rightarrow Abclb$$

$$B \rightarrow d$$

P RUNNING THE HANDLE

→ removing the children of left Hand side non Terminal from the parse tree is called

Handle pruning

→ A Rightmost derivation in Reverse can be obtained ~~using~~ by Handle Pruning

Steps to Follow:

→ start with a string of terminals 'w'

that is to be parsed.

→ let $w = Y_n$, where Y_n is the n^{th} right sentential form of an unknown RMD.

→ To reconstruct the RMD in reverse,

locate handle β_n in Y_n ; Replace β_n by (Y_{n-1}) to get $(n \rightarrow)^{\text{th}}$ RSD

Eg: grammar

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

Right sentential form	Handle	Reducing production
$id_1 + id_2 * id_3$	id_1	$E \rightarrow id$
$E + id_2 * id_3$	id_2	$E \rightarrow id$
$E + E * id_3$	id_3	$E \rightarrow id$
$E + E * E$	$E + E$	$E \rightarrow E + E$
$E * E$	$E * E$	$E \rightarrow E * E$
(E) start symbol		

Shift-Reducing Parser (SR)

- It is a form of Bottom-up parsing in which a stack holds grammar symbols and an input buffer holds the rest of the string to be parsed.
- The Handle always appears at the top of stack just before it is identified as Handle.
- we use '\$' to mark the bottom of the stack and also the right end of the input.
- Initially stack is empty, and string w is on input, as follows:

<u>STACK</u>	<u>INPUT</u>
\$	$w \$$

- During left-to-right scan of the input string, the parser shifts zero (or) more input symbols onto top of the stack, until it is ready to reduce a string β of grammar symbols on top of the stack.
- It then reduces β to the head of the appropriate production.
- The parser repeats this cycle until it has

detected an error or until the stack contains the start symbol and input is empty;

STACK
\\$ S

INPUT
A

→ upon Entering this configuration, the parser halts and announces successful completion of parsing.

→ Primary operations are shift and reduce

→ There are four possible Actions

- (i) shift
- (ii) reduce
- (iii) Accept
- (iv) error

shift - shift the next input symbol onto the top of the stack.

Reduce - The right end of the string to be reduced must be at the top of the stack. locate the left end of string within the stack and decide with what nonterminal to replace the string.

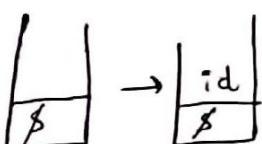
Accept - Announce successful completion of parsing.

Error → Discover a syntax error and call
an error recovery routine.

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

Eg⁵

STACK	CONTENT	INPUT	ACTION
\$		id ₁ + id ₂ * id ₃ \$	shift
\$ id ₁		+ id ₂ * id ₃ \$	Reduce by $E \rightarrow id$
\$ E		+ id ₂ * id ₃ \$	shift
\$ E +		id ₂ * id ₃ \$	shift
\$ E + id ₂		* id ₃ \$	Reduce by $E \rightarrow id$
\$ E + E		* id ₃ \$	shift
\$ E + E *		id ₃ \$	shift
\$ E + E * id ₃		\$	Reduce by shift $E \rightarrow id$
\$ E + E * E		\$	Reduce $E \rightarrow E * E$
\$ E + E		\$	Reduce $E \rightarrow E + E$
\$ (E)	start symbol	\$	Accept.



Bottom →

Top Down Parsing

Recursive Descent parser

→ A parser that uses collection of recursive procedures for parsing the given input string.

is called Recursive Descent (RD) Parser.

→ In this grammar CFG is used to build the recursive routines.

→ The RHS of the production rule is directly converted to a program.

→ For each non terminal a separate procedure is written and body of procedure (code) is RHS of the corresponding non terminal.

Basic steps for Construction of RD Parser

The RHS of the rule is directly converted into program code symbol by symbol.

(i) If the input symbol is non terminal then a call to procedure corresponding to the non terminal is made.

- (ii) If the input symbol is terminal then it is matched with the lookahead from input. The lookahead pointer has to be advanced on matching of the input symbol.
- (iii) If the production rule has many alternatives then all these alternatives has to be combined into a single body of procedure.
- (iv) The parser should be activated by a procedure corresponding to the start symbol.

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

$$T \rightarrow * \text{num } T \mid \epsilon$$

```

procedure E
{
  if lookahead = num then
  {
    match(num);
    T;           // call to procedure T
  }
}

```

```

else
  error;
  if lookahead = $
  {
    declare success; // return on success
  }
}

```

is
ling
ties
e.
l.

```
    else
        error;
```

}] / End of procedure E /

procedure T

{

if lookahead = '*' {

match('*');

if lookahead = 'num' {

match(num);

T;

}

else

error;

}

else NULL

}

procedure match(token t)

{

if lookahead = t

lookahead = next_token; *// lookahead pointer is advanced //*

else

error;

}

procedure error

{

print("Error");

3	*	4	8
---	---	---	---

$E \rightarrow \text{num } T$

(a)

3	*	4	8
---	---	---	---

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

(b)

3	*	4	8
---	---	---	---

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

(c)

3	*	4	8
---	---	---	---

Declare success

(d)

→ The parser will be activated by calling procedure E . since the first input character ③ is matching with num the procedure $\text{match}(\text{num})$ will be invoked and then the lookahead will point to next token. And a call to procedure T is given.

→ A match with '*' is found hence
 $\text{lookahead} = \text{next_token}$

→ Now '4' is matching with num hence again

the procedure for match(num) is fulfilled.

→ Then procedure for T is invoked and T is replaced by E.

→ As lookahead pointer points to \$ we quit by reporting success.

(Introduction to LR Parsing & Simple LR₀) *

Why use stack for LR-parsers

→ Any handle will always appear on the top of the stack and the parser need not search within the stack at any time.

Conflicts in LR-parsing

→ Two decisions decide a successful LR-parsing
(i) locate the substring to reduce
(ii) which production to choose when multiple productions with the selected substring on RHS exist.

→ LR parser may reach a configuration in which knowing the contents of stack and input buffer, still the parser cannot decide.

→ (a) whether to perform a shift or) a Reduce operations (shift Reduce conflict)

(b) ~~what~~ which out of the several possible reductions to make (Reduce-Reduce conflict)

Eg:- $\frac{S}{S \rightarrow \text{if expr then stmt}}$
 $\quad \quad | \quad \quad \frac{E}{\text{if expr then stmt}} \quad \frac{S}{\text{else stmt}}$
 $\quad \quad | \quad \quad \frac{E}{\text{other}}$

STACK

INPUT

$\$ \text{ if } E \text{ then } S$
valid handle

else - \$
→ should not be reduced

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

Eg:- STACK

Input

Action

\emptyset

$x * x \$$

shift

$\$ x$

$* x \$$

reduce by $F \rightarrow x$

$\$ F$

$* x \$$

reduce by $T \rightarrow F$

$\$ T$

$* x \$$

shift

$\$ T x$

$x \$$

shift

$\$ T * x$

$\$$

reduce by $F \rightarrow x$

$\$ T * E$

$\$$

reduce by $T \rightarrow T * F$

$\$ T$

$\$$

reduce by $E \rightarrow T$

$\$ E$

$\$$

Accept

Eg:-

STACK

\$

\$ x

\$ F

\$ T

\$ E

\$ E *

\$ E * x

\$ E * F

\$ E * T

\$ E * E

INPUT

x x x \$

x x \$

x x \$

x x \$

x \$

\$

\$

\$

\$

ACTION

shift

Reduce

F \rightarrow xL

Reduce F \rightarrow T

Reduce E \rightarrow T

shift

shift

reduce F \rightarrow xL

reduce T \rightarrow F ~~Rede~~

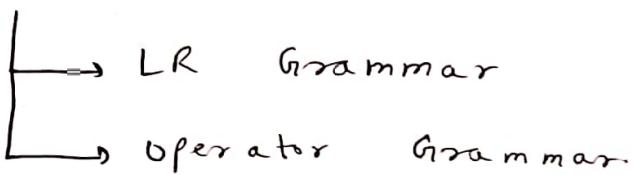
reduce E \rightarrow T

Error



Operator Grammar

→ shift reduce parsers can be built successfully using 1 for 2 main classes of grammar.



Properties of Operator Grammar

- 1) No production in the grammar has ϵ on (Right Hand Side)
- 2) No 2 Non Terminals appear together on RHS of any production.

Eg:-

$$\begin{aligned} E &\rightarrow EA E \mid (E) \mid -E \mid id \\ A &\rightarrow + \mid - \mid * \mid / \mid ^\wedge \end{aligned}$$

→ violates Rule 2 → not O.P.G

$$\boxed{\begin{aligned} E &\rightarrow E+E \mid E-E \mid E\times E \mid E/E \mid E^\wedge E \\ E &\rightarrow (E) \mid -E \mid id \end{aligned}} \quad \text{O.P.G}$$

Eg:-

$$\begin{aligned} S &\rightarrow SAS \mid a \quad \xrightarrow{\text{not O.P.G}} \\ A &\rightarrow bSb \mid b \\ S &\rightarrow SbSbS \mid Sbs \mid a \quad \xrightarrow{\text{O.P.G}} \end{aligned}$$

Operator Precedence Parsing

→ we define 3 precedence relation operations between pairs of terminals.

$p <- q$ p gives more precedence than q

$p \doteq q$ p has same precedence as q

$p \rightarrow q$ p takes precedence over q.

Use → They help us in selection of handles

How to Assign Relations

→ using associativity and precedence

→ For all terminals (including \$) we design an operator precedence table.

(operator)

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

I/P

TOS

	id	+	*	\$
id	↓	→	→	→
+	<.	→	<.	→
*	<.	→	→	→
\$	<.	<.	<.	ACC

*

$$\begin{array}{l} id \rightarrow + \\ id \rightarrow * \\ id \rightarrow \$ \end{array}$$

*

$$\begin{array}{l} \$ \leftarrow id \\ \$ \leftarrow + \\ \$ \leftarrow * \end{array}$$

Symbol on TOS

Symbol I/P → str

If I/P string symbol → TOS : PUSH I/P symbol
else If P string symbol < TOS : POP and Reduce TOS
else error

Eg:

F/P string

id + id + id \$

+ id + id \$

+ id + id \$

- id + id \$

+ id \$

+ id \$

+ id \$

id \$

\$

\$

\$

stack

\$

~~id~~ \$

Reduce
①

~~id~~ \$

\$ +

\$ + id

Reduce
②

\$ +

Reduce
③

\$

\$ +

\$ + id

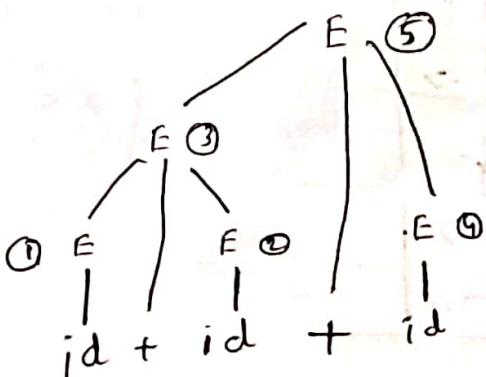
Reduce
④

\$ +

Reduce
⑤

\$

Accept,



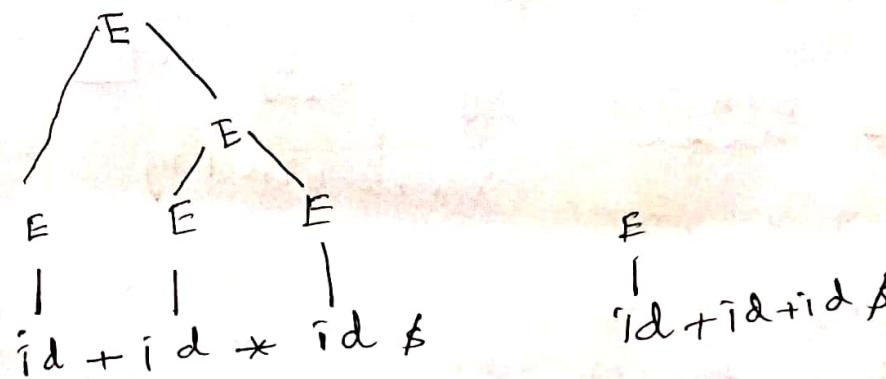
operator precedence parser is the only parser
that will handle ambiguous grammar

Advantage of operator precedence parsing

- Simplicity of the method
- Error detecting capability of the parser
- Ability to parse Ambiguous grammar (single parse tree)

Disadvantage of operator Precedence parsing.

- only a small class of grammars can be parsed using operator precedence technique.
- Difficult to handle operators which may have different roles (unary & binary plus/minus) in different situations
- As the no. of terminals increase the size of precedence table also increases (solution operator precedence function)



Introduction to LR Parsing

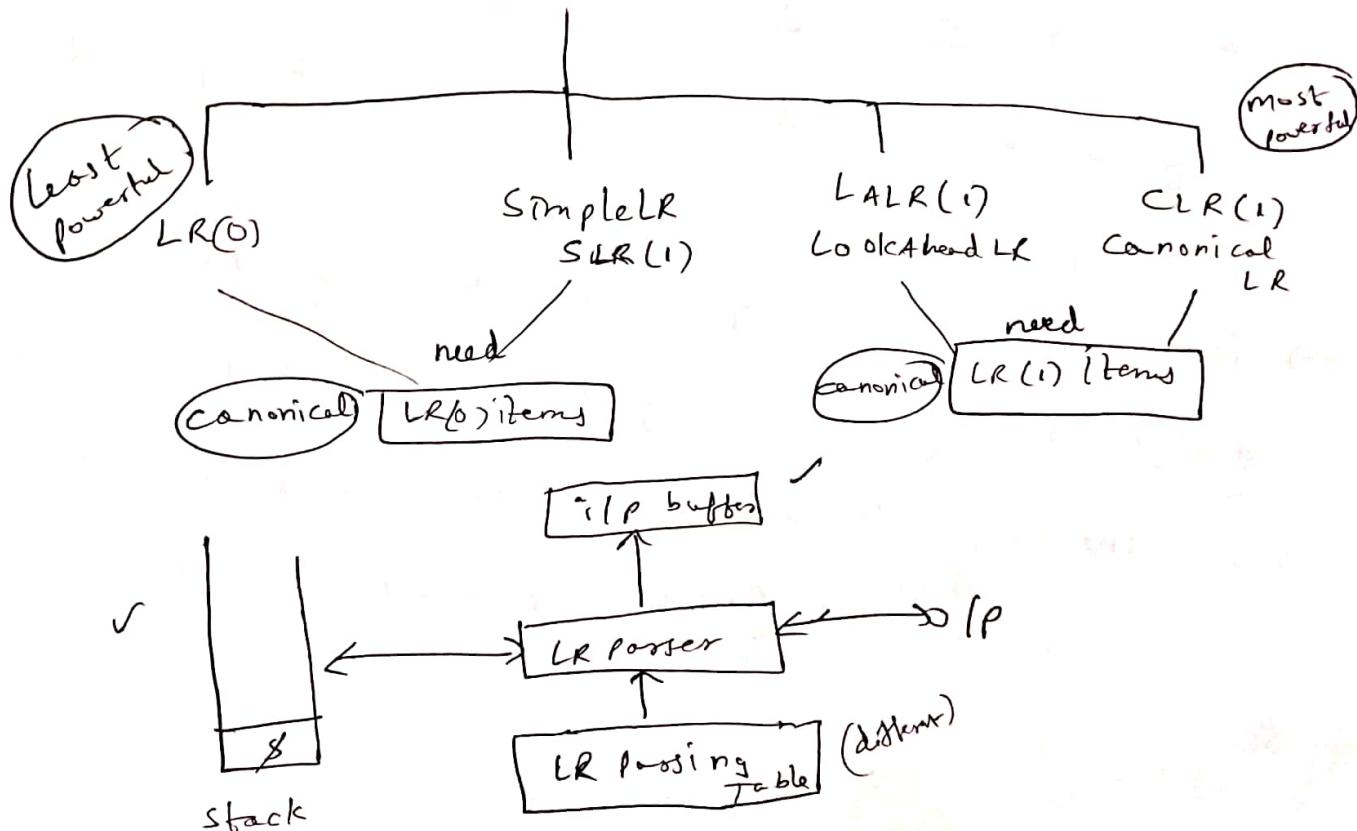
→ Most prevalent type of bottom-up parser today is based on concept called LR(k) parsing.

K - no. of IP symbols or lookahead used in

L - left to right scanning of input making decisions.

R - for construction rightmost derivation in reverse.

LR parser



Benefits of using LR(k) parsing

→ Most generic Non-Backtracking shift Reduce parser

Technique:

→ These parsers can recognize all programming languages for which context-free grammars can be written.

→ They are capable of detecting syntactic error as soon as possible in scanning of input.

Types of LR parsers

(i) Simple LR parser (SLR)

① easy to implement

② least powerful

③ may fail to work on some grammars accepted by CLR, LALR.

(ii) Canonical LR parser (CLR)

④ most powerful LR parser

⑤ most expensive

(iii) Lookahead LR parser (LALK)

⑥ intermediate to SLR and CLR in terms of power and cost.

Factories Eg:-

$$S \rightarrow AA$$

$$A \rightarrow aA/b$$

→ Add one more production

$$\boxed{\begin{array}{l} S' \rightarrow S \\ S \rightarrow AA \text{ ①} \\ A \rightarrow aA/b \text{ ②} \end{array}}$$

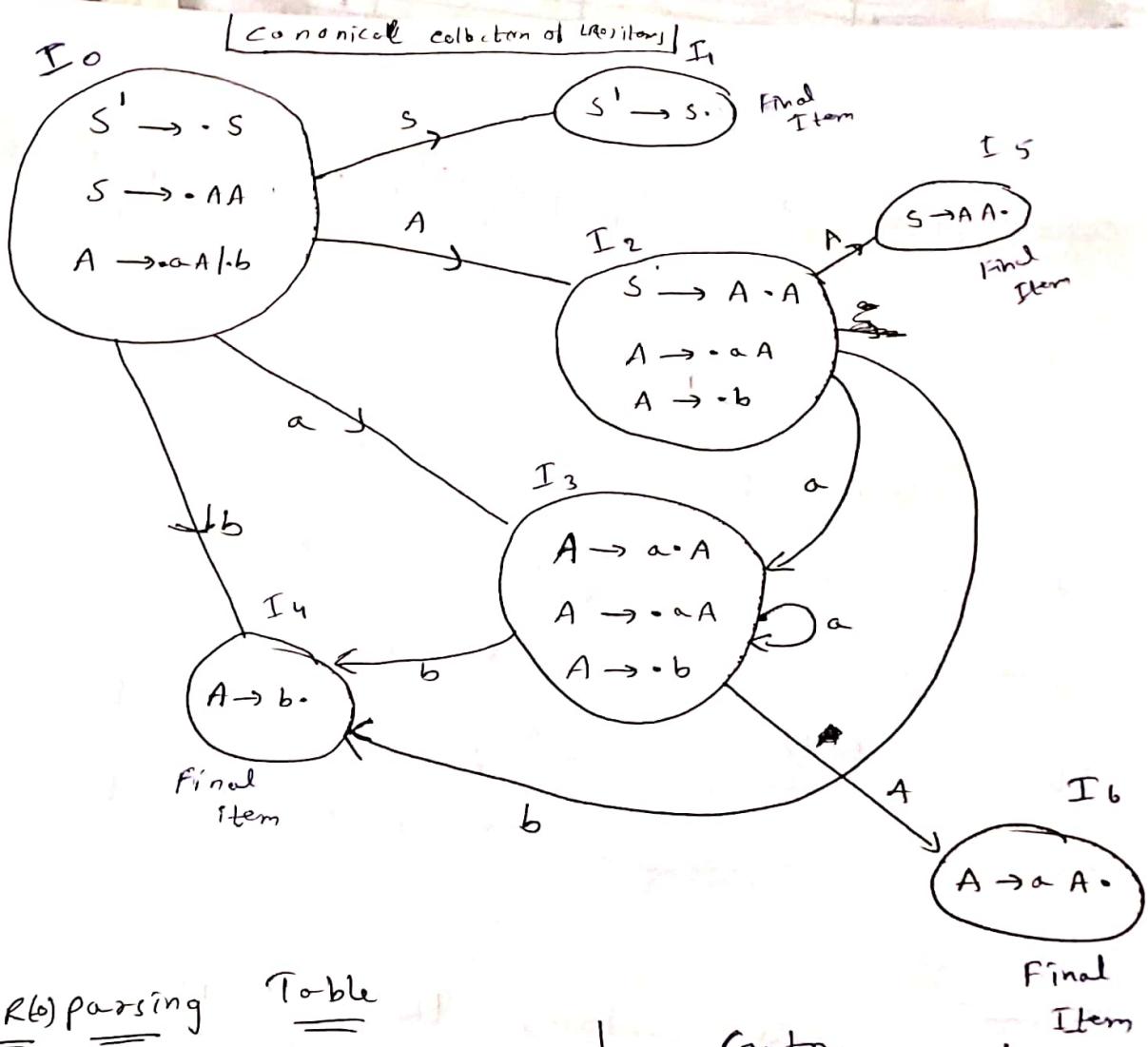
Augmented grammar

Various steps involved in Parsing

- For the given I/P string write context free grammar
- Check the ambiguity of the grammar
- Add Augment production in the given grammar
- Create canonical collection of LR(0) items
- Draw a Data flow diagram (DFA)
- Construct LR(0) Parsing Table.

Create Canonical collection of LR(0) items

- An LR(0) item is a production G_i with dot at some position on the right hand side of production
- LR(0) items is useful to indicate that how much of the ~~I/P~~ i/p has been scanned up to a given point in the process of parsing
- In LR(0), we place the reduce node in entire row.



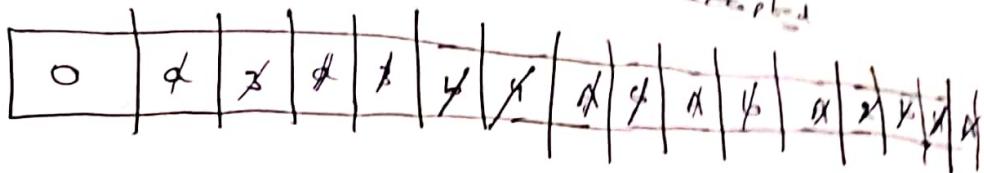
		Action			Goto	
		a	b	\$	A	S
0		S_3	S_4		2	
1				Accept		
2		S_3	S_4		5	
3		S_3	S_4		6	
4		r_3	r_3	r_3		
5		r_1	r_1	r_1		
6		r_2	r_2	r_2		

Eg :-

$a \atop a \atop a \atop b \atop b \atop b \atop \$$

String

is accepted



8

$\boxed{o} \boxed{s} \boxed{i}$

accept

LR(0) Table

- If a state is going to another state on a terminal. it is correspond to a shift.
- If a state is going to some other state on a variable it is correspond to a reduce.

go to move

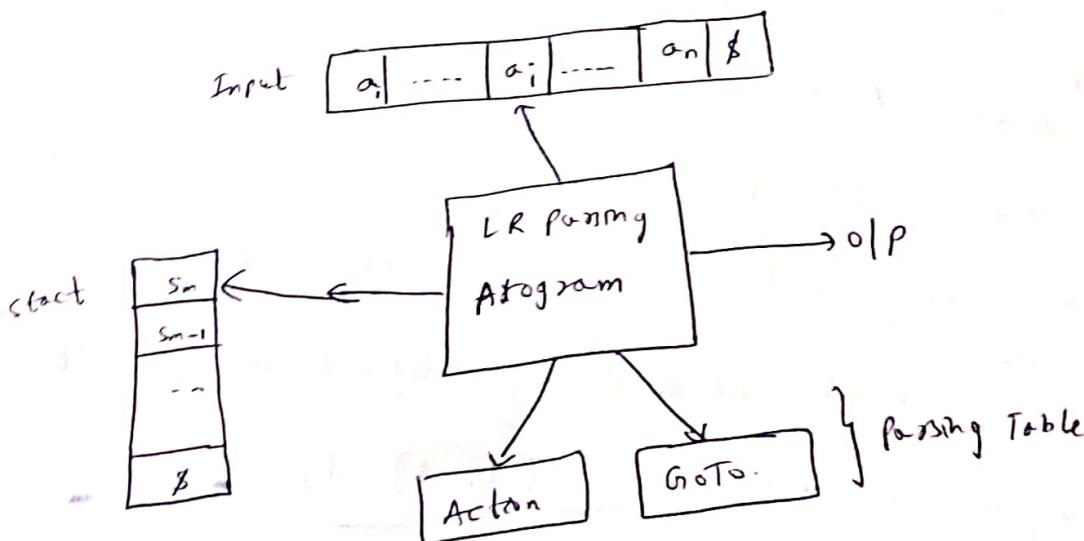
- If a state contains the final item in a particular row then write reduce node completed.

Parsing I/P string

<u>steps</u>	<u>Line</u>	<u>state</u>	<u>posting stack</u>	<u>Input</u>	<u>Action</u>
(1)		\emptyset	$\$ \emptyset o$	<u>aabb\$</u>	shift a2
(2)			$\$ o a_3$	<u>abb\$</u>	shift a3
(3)			$\$ o a_3 a_3$	<u>bb\$</u>	shift b4
(4)			$\$ o a_3 a_3 b_4$	<u>b\$</u>	reduce r ₃ (A → a) ₄
(5)			$\$ o a_3 a_3 A_6$	<u>b\$</u>	reduce r ₂ (A → a) ₄
(6)			$\$ o a_3 A_6$	<u>b\$</u>	reduce r ₁ (A → a) ₄

(7)	$\$ \rightarrow A_2$	$b b$	shift by
(8)	$\$ \rightarrow A_2 b_4$	\underline{b}	reduce $\tau_3 (A \rightarrow b)$
(9)	$\$ \rightarrow A_2 A_5$	$\underline{\underline{b}}$	reduce $\tau_1 (S \rightarrow AA)$
(10)	$\$ \rightarrow S$	$\underline{\underline{\underline{b}}}$	Accept

Components of LR Parser



→ Parsing Table changes from one LR parser to another.

Structure of LR Parsing Table

→ The Parsing Table consists of two parts

(a) Parsing Action Function (ACTION)

(b) goto Function (GOTO)

→ The Action function takes as arguments a state i and a terminal a (or $\$$ if the input end marker).

Action [i, a] can have four forms

Action [i, a] can have four forms

(i) Shift j , where j is a state - The action taken by the parser effectively shifts input a to the stack, but uses state j to represent a .

(ii) Reduce $A \rightarrow \beta$. The action of the parser effectively reduces β on the top of stack to read A .

(iii) Accept. The parser accepts the input and finishes parsing.

(iv) Error. The parser discovers an error in its input and takes some corrective action -

→ we extend GOTO Function, defined on sets of items, to states: if $\underline{\text{GOTO}}[I_i, A] = I_j$, then $\underline{\text{GOTO}}$ also maps A to state i and nonterminal A to state j .

Behavior of LR Parser

→ The next move of the parser from the configuration is determined by reading a_i , the current ip symbol and s_m the state on top of stack and then consulting the entry $\underline{\text{ACTION}}[s_m, a_i]$ in parsing action Table.

The configurations resulting after each of the four types of move are as follows:

(i) If ACTION [s_m, a_i] = shift δ , the parser executes shift move, it shifts the next state s onto the stack, entering the configuration.

$$(S_0 S_1 \dots s_m, a_{i+1} \dots a_n)$$

(ii) If ACTION [s_m, a_i] = reduce $A \rightarrow B$, then parser executes reduce move, entering the configuration

$$(S_0 S_1 \dots s_{m-r}, s, a_i a_{i+1} \dots a_n)$$

$r \rightarrow$ length of B

$$s = \text{GOTO}[s_{m-r}, A]$$

(iii) If ACTION [s_m, a_i] = accept, parsing is completed.

(iv) If ACTION [s_m, a_i] = error, parser has discovered an error and calls an error recovery routine.

Constructing SLR Parsing Table

Input : An augmented grammar G' .

Output : SLR parsing Table functions ACTION & GOTO
for G' .

Method :

(1) Construct $C = \{I_0, I_1, \dots, I_n\}$, collection of sets
of LR(0) items for G' .

(2) state i is constructed from I_i - The parsing
actions for state i are determined as follows.

(i) If $[A \rightarrow d \cdot a\beta]$ is in I_i and $\text{GOTO}(I_i, a) = I_j$,

then set $\text{ACTION}[i, a]$ to shift j. (a must be terminal)

(ii) If $[A \rightarrow d \cdot]$ is in I_i , then set $\text{ACTION}[i, a]$
to "use $A \rightarrow d$ " for all a in FOLLOW(A)

(A may not be S')

(iii) If $[S' \rightarrow s \cdot]$ is in I_i , then set $\text{ACTION}[i, \$]$

to Accept

SLR(1) Parsing Table

→ only change in Reduce moves

→ Reduce moves are placed in ~~Follow(A)~~

If production is " $A \rightarrow \alpha$ " not in entire row.

	Action			Goto	
	a	b	\$	A	S
0	s_3	s_4		2	1
1			Accept		
2	s_3	s_4		5	
3	s_3	s_4			6
4	r_3	r_3	r_3		
5			r_1		
6	r_2	r_2	r_2		

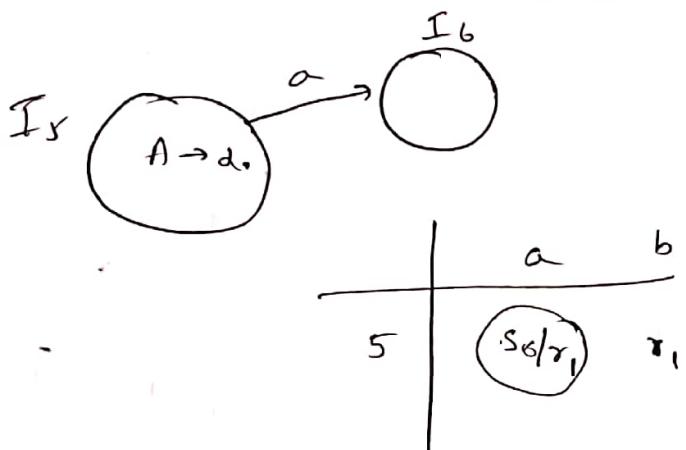
$I_4 \rightarrow A \xrightarrow{④} b \cdot$ → reduce move should be placed
in $\text{Follow}(A) = \{a, b, \$\}$

$I_5 \rightarrow S \xrightarrow{①} \cdot A A \cdot$ → reduce move should be placed
in $\text{Follow}(S) = \{\$\}$

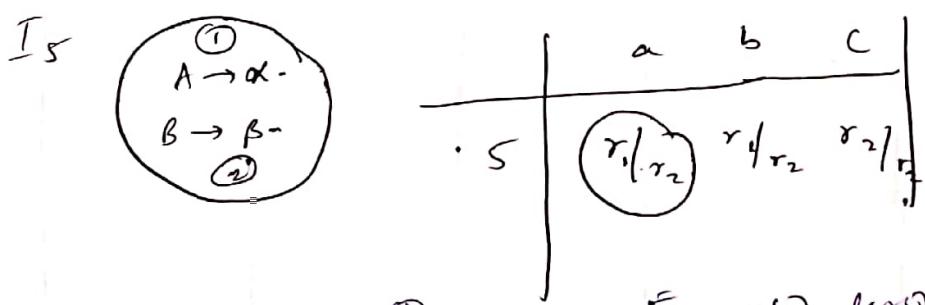
$I_6 \rightarrow A \xrightarrow{②} a A \cdot$ → reduce move should be
placed in $\text{Follow}(A) = \{a, b, \$\}$

In $LR(0)$ Parser

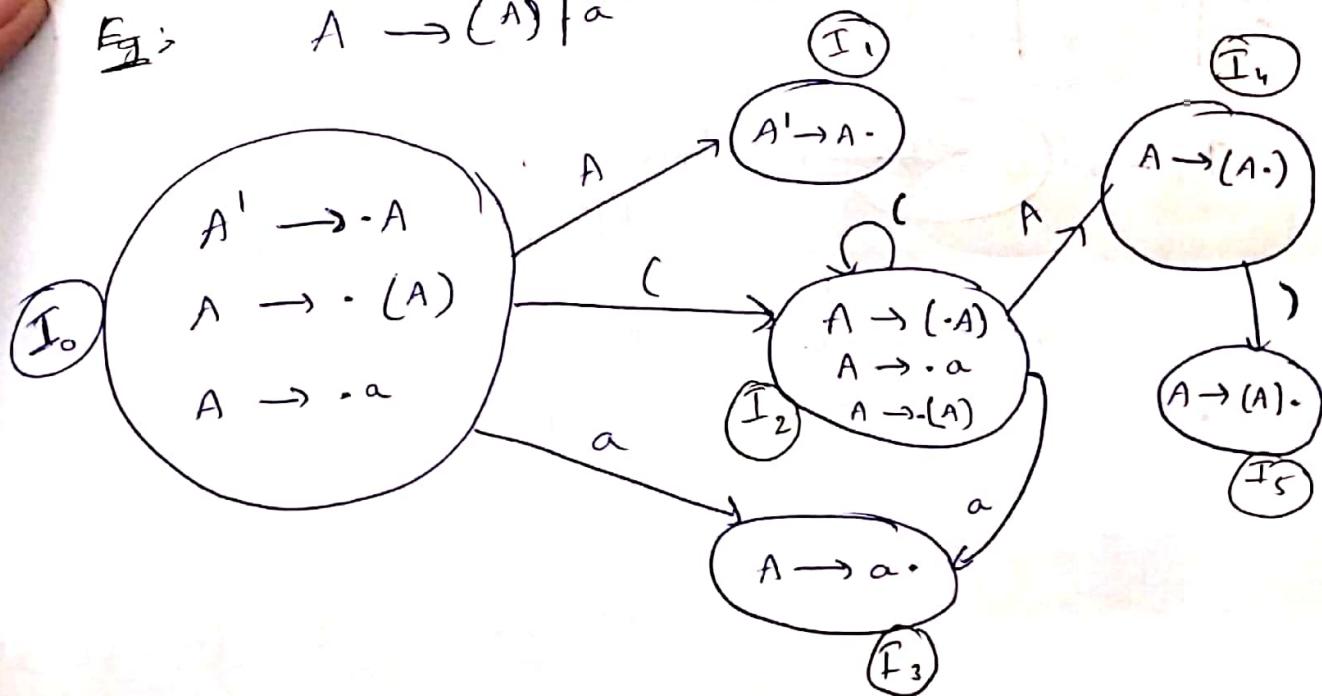
→ we may get s/r conflict when



→ we may also get r/r conflict when



$E_1:$ $A \xrightarrow{\text{①}} (A) + a^{\text{②}}$ [SLR(1) Parsing Table]



	Action				Goal
0	()	a	β	A
1				Accept	I
2	s_2		s_3		4
3		r_2		r_2	
4		s_5			
5		r_1		r_1	

I_3 and I_5 final items

$$I_3 \xrightarrow{\text{place}} r_2 \text{ in } A \rightarrow a \quad \text{Follow}(A) = \{ \), \$ \}$$

$$I_5 \xrightarrow{\text{place}} r_1 \text{ in } A \rightarrow (A) \cdot \quad \text{Follow}(A) = \{ \beta, \$ \}$$

step	passing stack	Input	Action
(1)	\$ 0	(a) β	shift 2
(2)	\$ 0 (2	a) \$	shift 3
(3)	\$ 0 (2 a 3) \$	reduce 2 $A \rightarrow a$
(4)	\$ 0 (2 A 4) \$	shift 5
(5)	\$ 0 (2 A 4) 5	\$	reduce 1 $A \rightarrow (A)$
(6)	\$ 0 A !	\$	Accept

Grammar

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

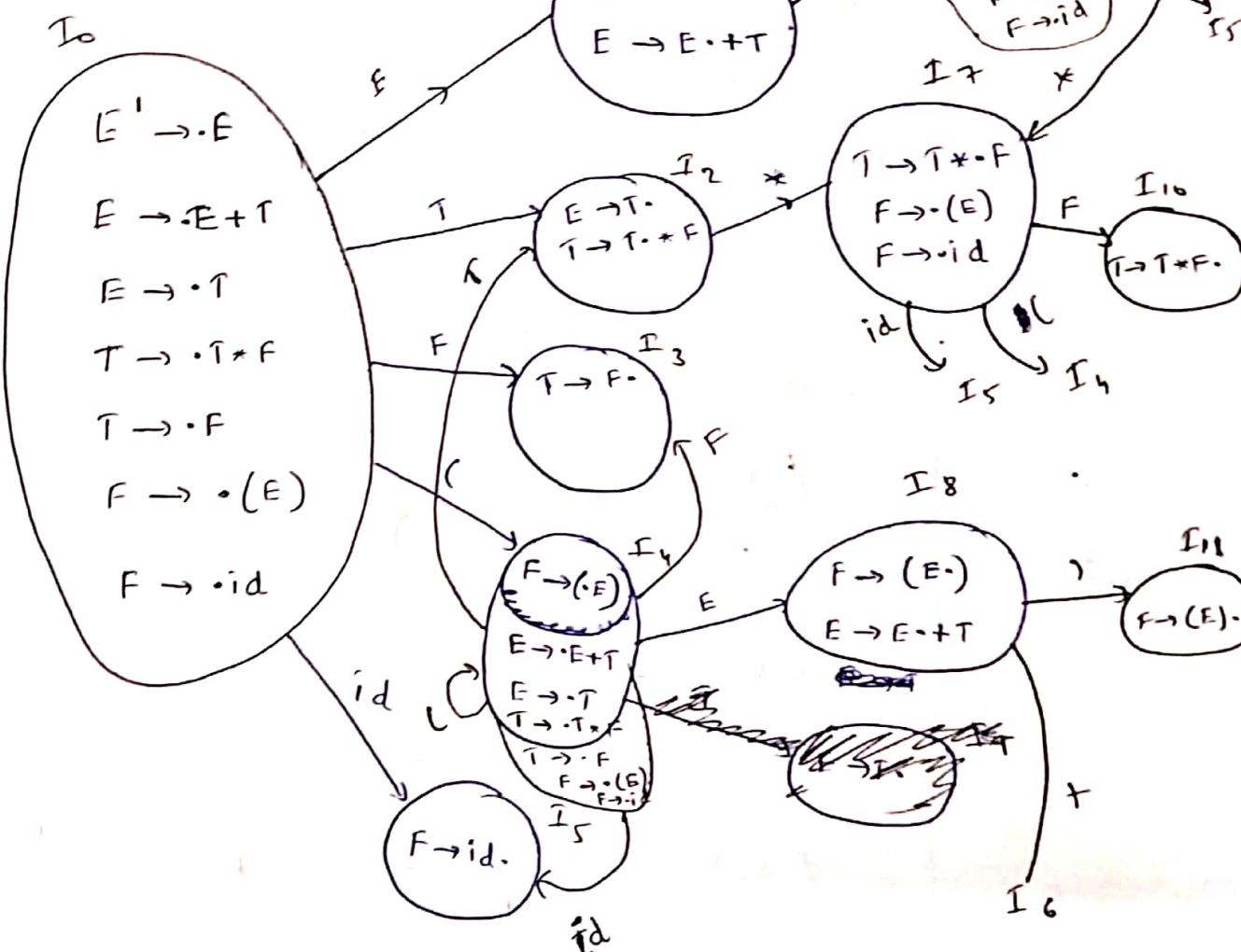
$$E \rightarrow T^2,$$

$$T \rightarrow T * F^3$$

$$T \rightarrow F^4$$

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

$$F \rightarrow id^6$$



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

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

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

Parsing Table

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

Egr	S Stack	Parsing Stack	IP string	Action
1	\$ 0		id * id + id *	shift 5
2	\$ 0 id \$		* id + id *	Reduce 6 ($E \rightarrow id$)
3	\$ 0 F 3		* id + id *	Reduce 4 ($T \rightarrow F$)
4	\$ 0 T 2		* id + id *	shift 7
5	\$ 0 T 2 * 7		id + id *	shift 5
6	\$ 0 T 2 * 7 id 5		+ id *	Reduce 6 ($F \rightarrow id$)

7)

$$\$0 \boxed{T2 \times 7 F10}$$

tidy

8)

$$\$0 \boxed{T2}$$

tidy

9)

$$\$0 \boxed{EI}$$

tidy

10)

$$\$0 EI + 6$$

tidy

11)

$$\$0 EI + 6 \boxed{ids}$$

\$

12)

$$\$0 EI + 6 \boxed{F3}$$

\$

13)

$$\$0 \boxed{EI + 6 T9}$$

\$

14)

$$\$0 \boxed{EI}$$

\$

CLR(1) and LALR(1) [LR(0) items]

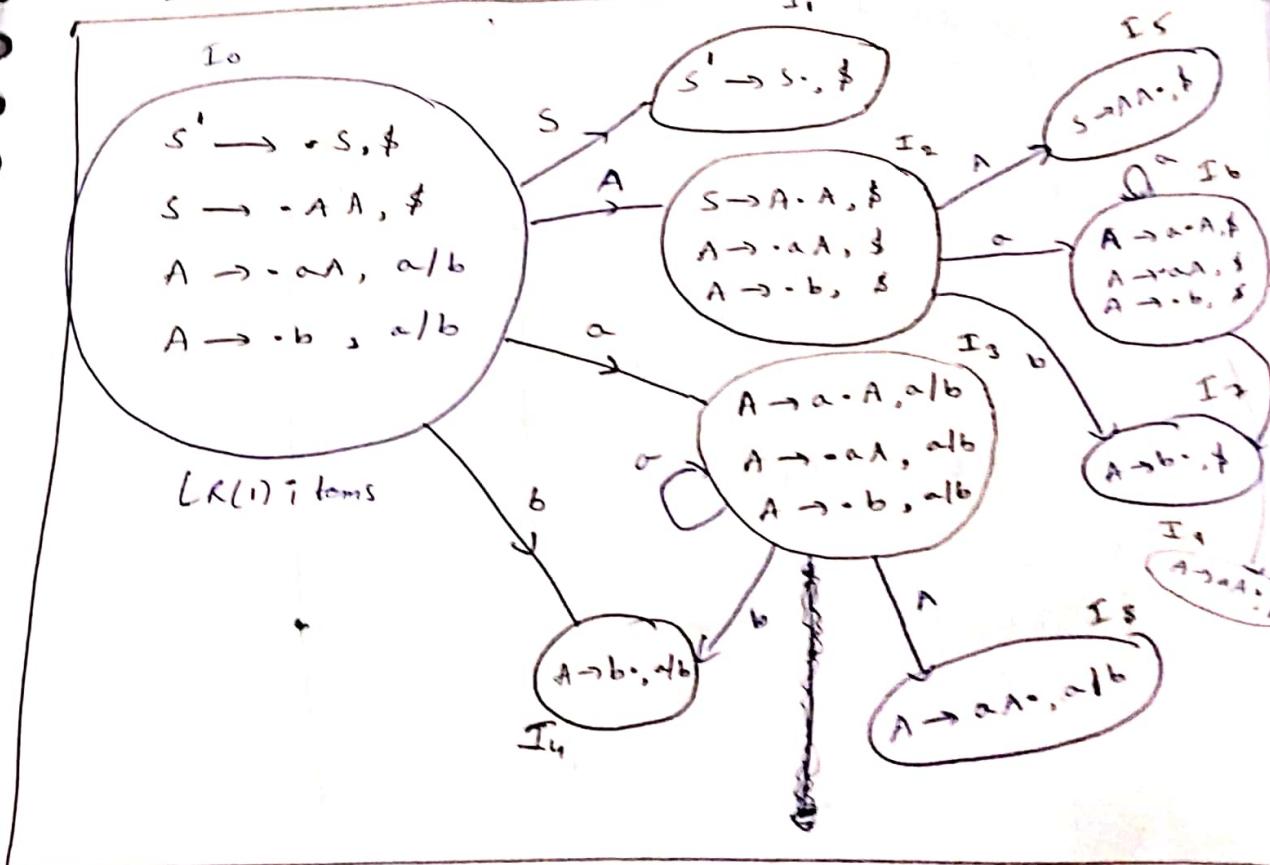
→ LR(1) item = LR(0) item + look ahead

Grammar

$$\begin{array}{l} S \rightarrow A \cdot A \\ A \rightarrow aA/b \end{array}$$

Rule

$$\begin{array}{l} A \rightarrow a \cdot B \beta, \alpha/b \\ B \rightarrow \cdot \gamma, \text{First}(B \cdot \alpha/b) \end{array}$$



Canonical collection of LR(1) items

Parsing Table		<u>C LR(1)</u>			Goto	
		a	b	*	A	S
0		s_3	s_4			
1				Accept		
2		s_6	s_7		5	
3		s_3	s_4		8	
4		r_3	r_3			
5						
6		s_6	s_7		9	
7						
8		r_2	r_2			
9						

No. of states =

$$C LR(1) \geq LR(0) = \frac{LA LR(1)}{SLR(1)}$$

LALR Parsing Table

	a	b	β	s	A
0	S_{36}	S_{47}			ϵ
1			Accept		
2	S_{36}	S_{47}			5
36	S_{36}	S_{47}			89
47	r_3	r_3	r_3		
5			r_1		
89	r_2	r_2	r_2		

LR(1) items

SL-conflict

$$\begin{aligned} A &\rightarrow \alpha - \alpha\beta, c/d \\ B &\rightarrow r., \alpha/\$ \end{aligned}$$

RR conflict

$$\begin{aligned} A &\rightarrow d., a \\ B &\rightarrow d., a \end{aligned}$$

Eg:-

$$E \rightarrow BB$$

$$B \rightarrow \underset{\text{C}}{B} d$$

I₀

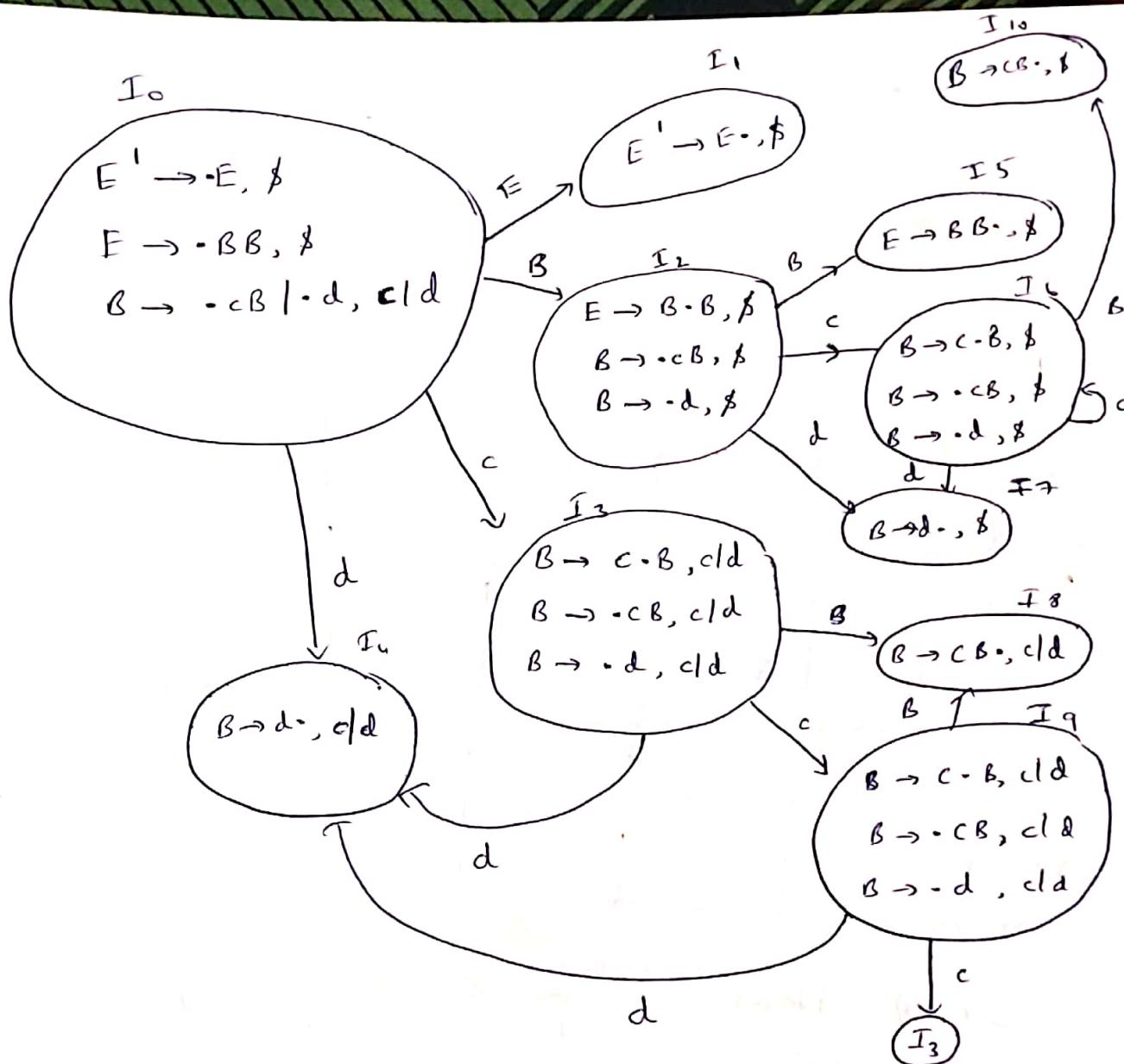
Canonical LR(1) items

$$E^1 \rightarrow \cdot E, \$$$

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

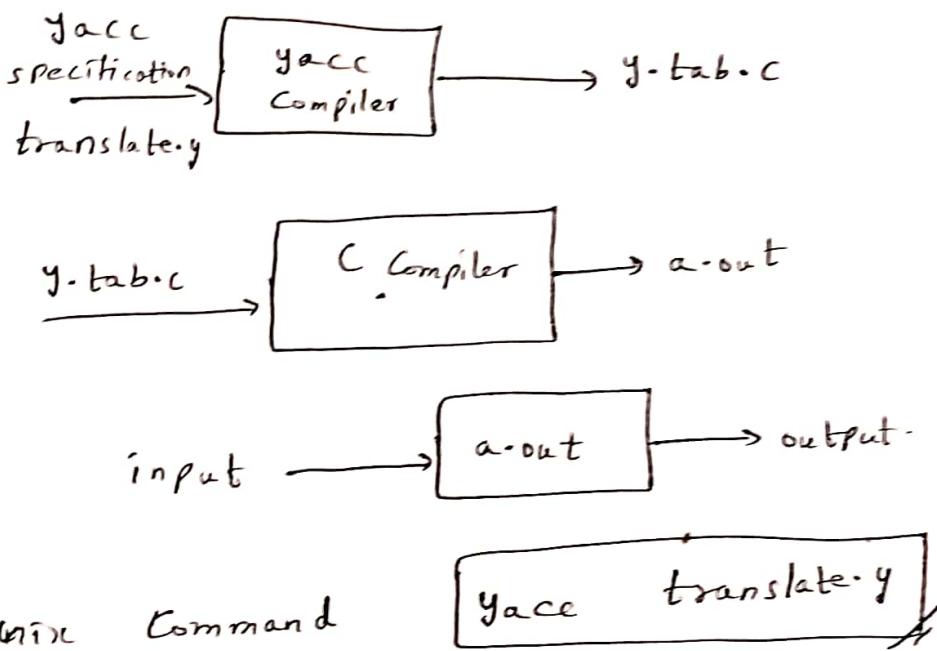
$$B \rightarrow \cdot CB, c/d$$

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



Parser Generators

→ Yet Another Compiler - Compiler (YACC)



→ Unix Command

yacc translate.y

→ The above command translates the file

translate.y into a C-program called y.tab.c using

LALR method.

→ y.tab.c is a representation of an LALR parser

written in C, along

→ By compiling y.tab.c along with ly library

that contains the LR parsing program using command

CC y.tab.c -ly

→ The yacc has three parts

declarations

%%

Translation rules

%%

supporting C routines

→ The Declarations Part

→ These are two sections in the declarations

Part of a YACC program (Both are optional).

In the first section, we put ordinary

In the first section, we put ordinary

C declarations, delimited by %% and %%.

→ The Translation Rules Part

→ In the part of the YACC specification

In the part of the YACC specification

after the first %% pair, we put the translation rules.

Each rule consists of a grammar production

and associated semantic action.

<head> → <body₁> | <body₂> | ... | <body_n>

would be written in YACC as

<head> : <body₁> { <semantic action₁> }
| <body₂> { <semantic action₂> }
| :

| <body_n> { <semantic action_n> }

→ A YACC semantic action is a sequence of 'c' statements. In a semantic action, the symbol $\$\$$ refers to the attribute value associated with the Nonterminal of the head, while $\$_i$ refers to value associated with i^{th} grammar symbol (Terminal or Nonterminal) of the body.

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

↓
expr : expr '+' term $\{ \$\$ = \$1 + \$3; \}$

| term
;

→ The 'term' is the third grammar symbol of the body, '+' is second symbol.

→ The semantic Action associated with the first production adds value of expr and term of body and assign result to expr of head.

$$\boxed{\$\$ = \$1 + \$3;}$$

$$\boxed{\$\$ = \$1;}$$

} → default action.

The Supporting C-Routines part

- The third part of a YACC specification consists of supporting C-routines.
- A lexical analyzer by ~~name~~ yylex() must be provided.
- The lexical analyzer yylex() produces tokens consisting of a token name and its associated attribute value.
- The attribute value associated with a token is communicated to the parser through a YACC-defined variable yyval.

YACC specification of a simple Desk calculator

% \$

#include <ctype.h>

% }

% token DIGIT

% ~%

line : expr `n' { printf("%d\n", \$1); }

;

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

| term

;

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

1 factor

:

factor : 'c' expr ')' { \$y = \$2; }

1 DIGIT

:

'.'

Yylex()

{

int c;

c = getchar();

if (isdigit(c))

{

yyval = c - '0';

return DIGIT;

}

return c;

}

→ YACC will ~~not~~ resolve all parsing
action conflicts using two following rules:

(i) A reduce/reduce conflict is resolved
~~minimum~~

by choosing the conflicting production
listed first in the YACC specification.

(ii) A shift/reduce conflict is resolved in ~~less~~ favour of shift.

Eg:-

% left '+' [left Associative]

+ & - same precedence

→ YACC resolves shift/reduce conflict by attaching a precedence & associativity to each production involved in a conflict, as well as to each terminal involved in a conflict.

→ If it must choose between shifting input symbol α and reducing by production $A \rightarrow \alpha$, YACC reduces if the precedence of the production is greater than α , or if the production precedences are the same and the associativity of production is left, otherwise shift action is performed.