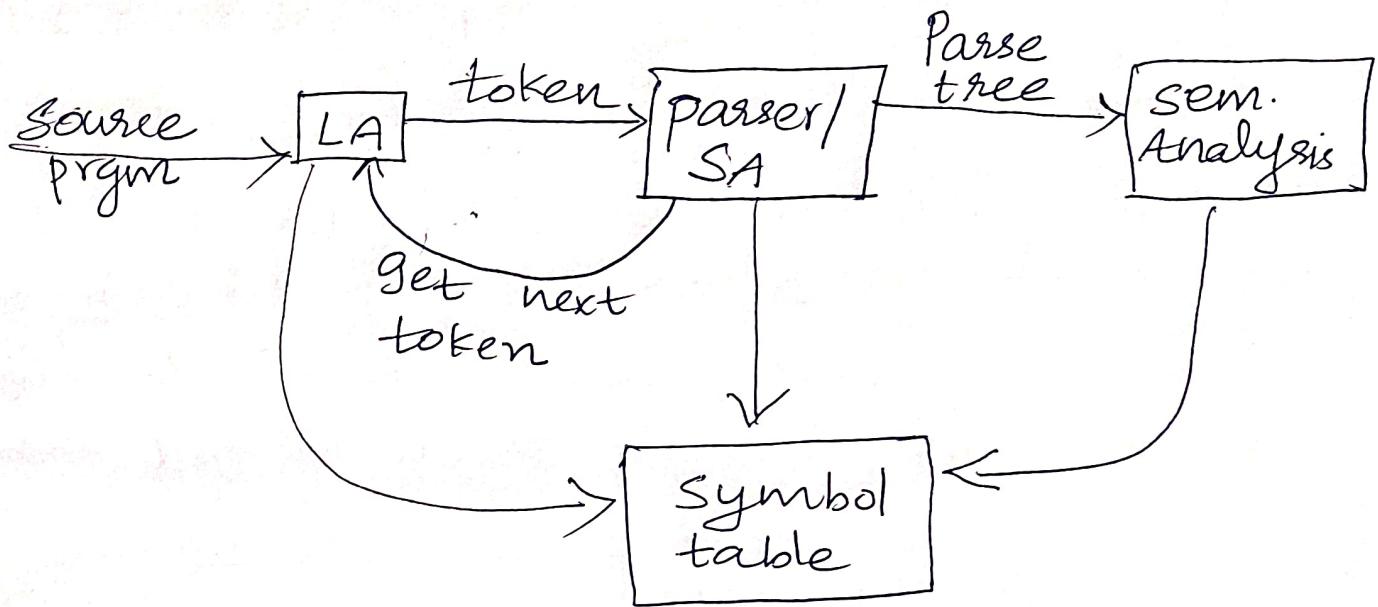
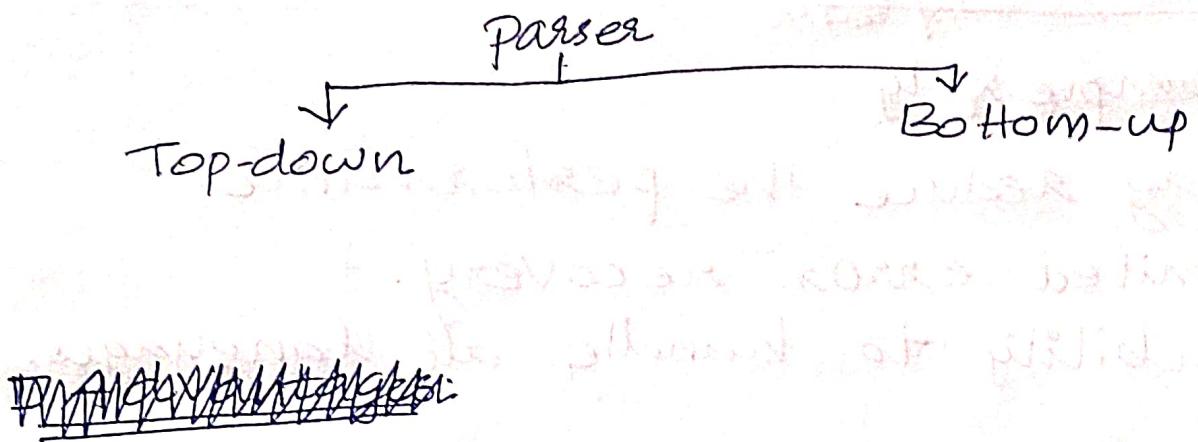


② Syntax Analysis / Parsing



- It is the 2nd phase of compilation process,
- primary goal is to verify the syntactical correctness of the source code .
- i/p : tokens
o/p : parsetree / Abstract syntaxtree.
- During this phase, the parser will check whether the input string is belong to the language generated by the Context Free Grammar (CFG).
If the syntax is correct, it will move forward; otherwise error.

- Two types:



- Features of syntax Analysis

- Syntax trees
- CFG
- Top-down and Bottom-up
- Error-detection.
- Basic optimization.

- The PDA is used to design the parser.

- Advantages of SA

- Structural Validation
(to check if the source code follows the grammatical rules of the programming lang, which helps to detect and report errors. in the source code)
- Improved code generation
- Easier semantic analysis.

Disadvantages

- Complexity
- May reduce the performance.
- Limited error recovery.
- Inability to handle all languages.

Eliminating Ambiguity

Consider the grammar:

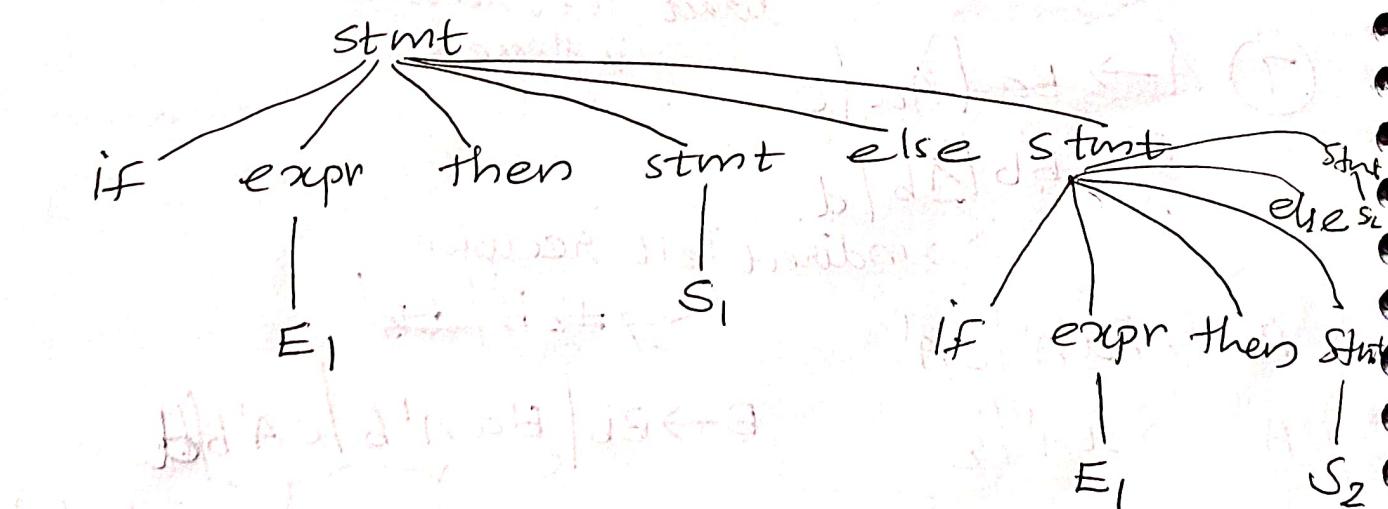
stmt \rightarrow if expr then stmt /
if expr then stmt else stmt /
other.

(other \rightarrow any other statement.)

Create parse tree for

if E_1 then S_1 , else if E_2 then S_2 else S_3

[for expr \rightarrow add E_1 and E_2]
[for stmt \rightarrow add S_1 and S_2]



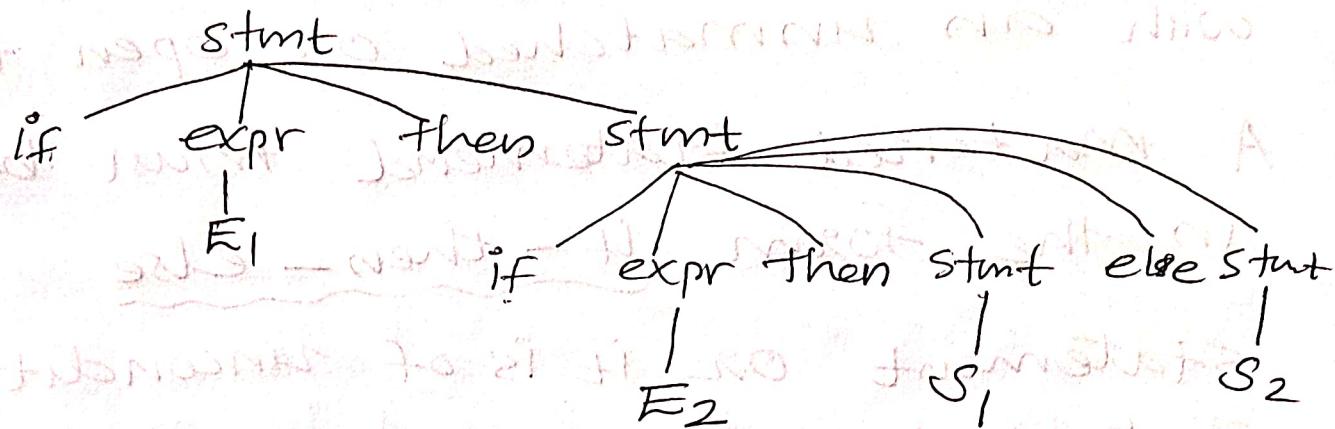
~~Prakash~~

It is ambiguous, because same string can be derive using two parse tree.

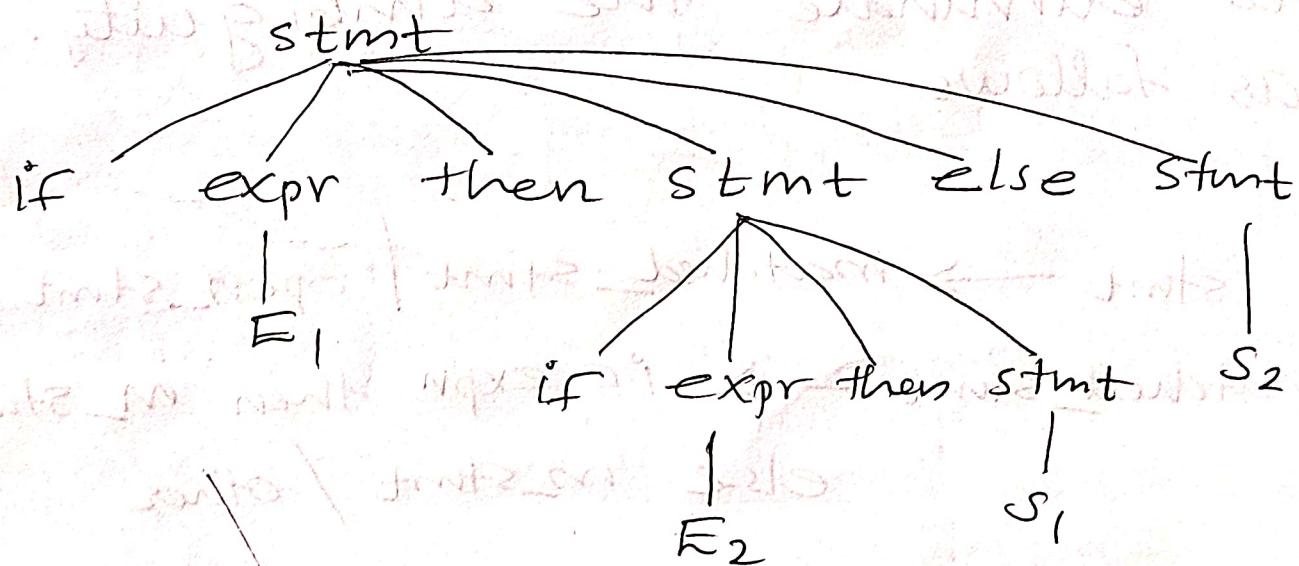
Eg: For the string

if E_1 then if E_2 then S_1 else S_2 .

(i)



(ii)



So the above grammar is ambiguous.

So the grammar can be call as

'dangling - else grammar'.

must be in the form:
if-then-else
ie, a statement appearing b/w a
then and an else must be 'matched';
ie, interior statement must not end
with an unmatched or open then.
A matched statement must be
in the form if-then-else
statement or it is of unconditional
statement.

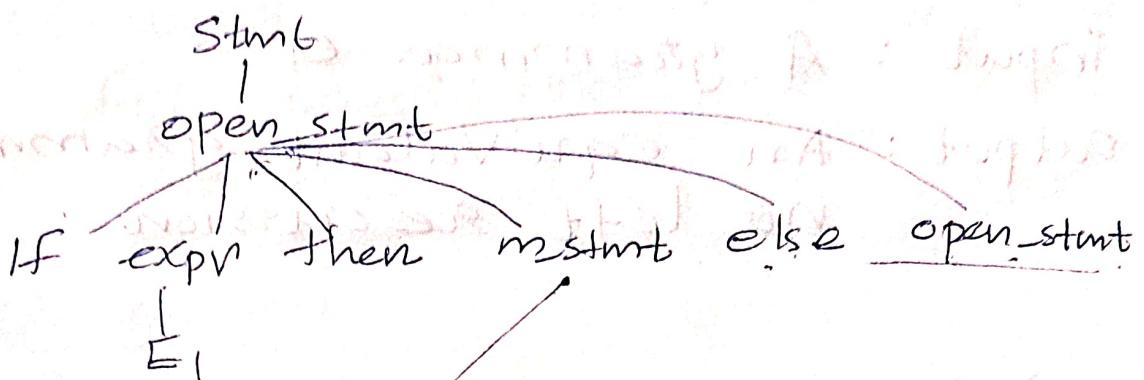
So we can rewrite the grammar
to eliminate the ambiguity.
as follows.

stmt → matched_stmt / open_stmt

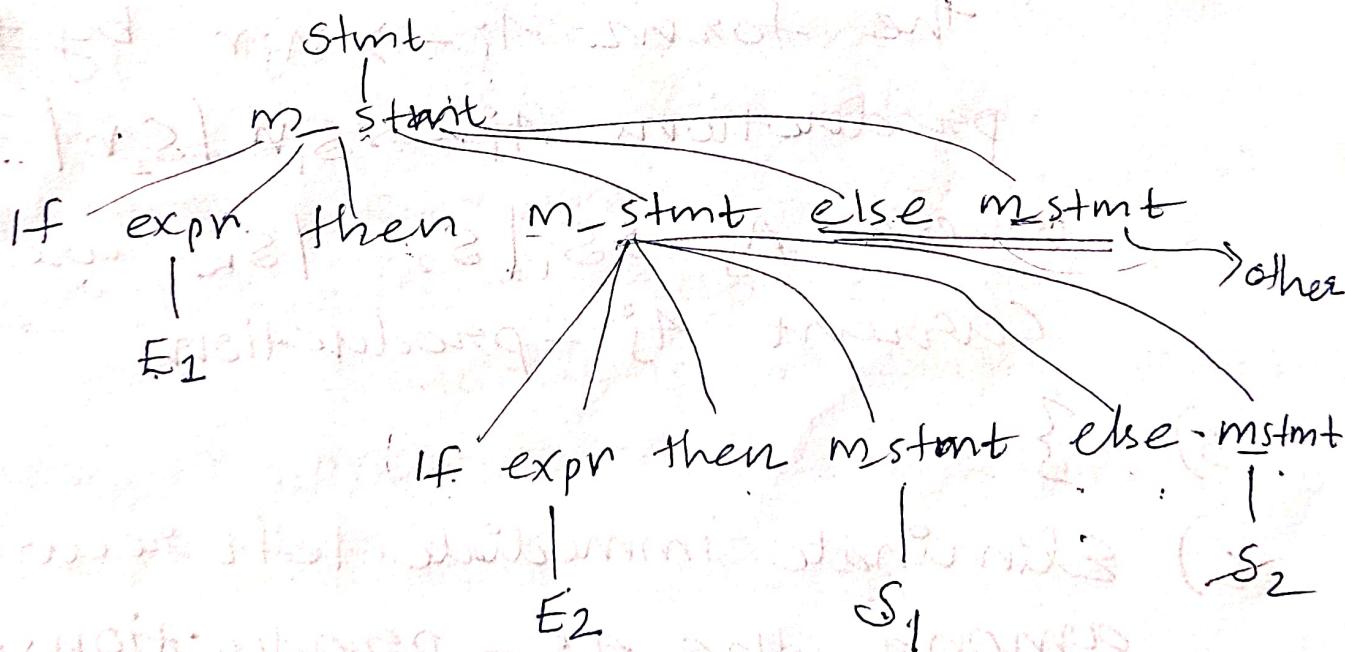
matched_stmt → if expr then m_stmt
else m_stmt / other

open_stmt → if expr then stmt /
if expr then m_stmt else
open_stmt.

The above is an unambiguous grammar for if-then-else statements.



If E₁. then if E₂ then S₁ else S₂



Algorithm for eliminating left recursion:

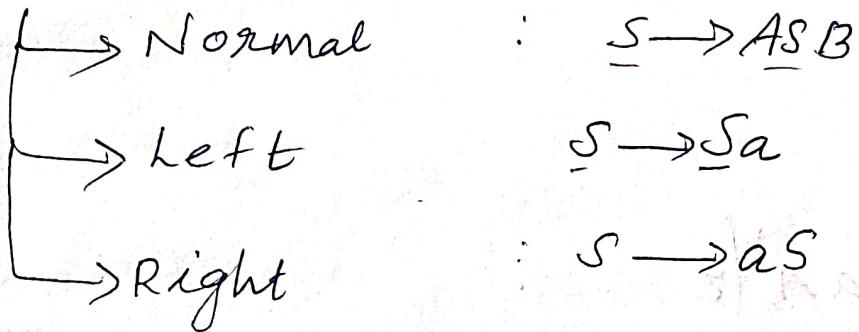
input : A grammar G

output : An equivalent grammar with no left recursion

- 1) Arrange the nonterminals in some order $A_1, A_2 \dots, A_n$.
- 2) for (each i from 1 to n) {
- 3) for (each j from 1 to $i-1$) {
- 4) replace each production of the form $A_i \rightarrow A_j r$ by the productions $A_i \rightarrow S_1 r | S_2 r | \dots | S_k r$, where $A_j \rightarrow S_1 | S_2 | \dots | S_k$ are all current A_j -productions.
- 5) }
- 6) Eliminate immediate left recursion among the A_i -productions.
- 7) }

Removal of Left Recursion

Recursion



Removal :-

If it is in the form $A \rightarrow A\alpha / \beta$

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

Examples

⇒ Direct left recursion

① $S \rightarrow a/b$

$S \rightarrow b S'$
 $S' \rightarrow a S' / \epsilon$

③ $A \rightarrow A b c d e$

$A \rightarrow e A'$
 $A' \rightarrow b c d A' / \epsilon$

② $S \rightarrow S a / d$

$S \rightarrow d S'$
 $S' \rightarrow a S' / \epsilon$

④ $A \rightarrow A_B d / A_a a$

$B \rightarrow B_e b$

$A \rightarrow aA'$

$A' \rightarrow BdA' / aA' / \epsilon$

$B \rightarrow bB'$

$B' \rightarrow eB' / \epsilon$

⑤ $E \rightarrow E + T / T$

$T \rightarrow T * F / F$

$F \rightarrow id$

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow id$

Indirect left recursion

(6)

$$S \rightarrow A\alpha / \beta$$

$$A \rightarrow Sd$$

Replace S with $A\alpha / \beta$

$$S \rightarrow A\alpha / \beta$$

$$A \rightarrow A\alpha d / \beta d$$

Now, remove left recursion of A as usual method.

$$S \rightarrow A\alpha / \beta$$

$$A \rightarrow \beta d A'$$

$$A' \rightarrow \alpha d A' / \epsilon$$

(7) $A \rightarrow Ba / Aa / c$

$$B \rightarrow Bb / Ab / d$$

direct left recursion
is there.

indirect left recursion

$$\{ A \rightarrow Ba A' / CA'$$

$$B \rightarrow AB B' / dB'$$

$$A' \rightarrow aa' / \epsilon$$

$$B \rightarrow Bb / Ba A' b / CA' b / d$$

prod. of A

$$B \rightarrow CA' b B' / dB'$$

$$B' \rightarrow b B' / a A' b B' / \epsilon$$

II. Left Factoring

It is a grammar transformation that's useful for producing a grammar suitable for predictive, or top-down parsing.

Algorithm

Input: A grammar G_1 .

Output: An equivalent left-factored grammar.

Method:

For each non-terminal A , find the longest prefix α common to two or more of its alternatives. If $\alpha \neq \epsilon$, there's a non-trivial common prefix — replace all of the A -productions $A \rightarrow \alpha\beta_1 | \alpha\beta_2 | \dots | \alpha\beta_n | r$ where r represents all alternatives that do not begin with α , by

$$A \rightarrow \alpha A' | r$$

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

Here, A' is a new non-terminal.

Repeat.

Normally, left factoring is applying for grammar with common prefixes.

Eg: $A \rightarrow \alpha\beta_1 | \alpha\beta_2 | \alpha\beta_3$

Here, α is common.

- These kind of grammar creates a problem to topdown parsers.
- The parser cannot decide which production must be chosen to parse the string in hand.

$$A \rightarrow \alpha\beta_1 | \alpha\beta_2 | \alpha\beta_3$$

Removal of left factoring can be done as;

$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta_1 | \beta_2 | \beta_3$$

Examples:

1) $x \rightarrow abc x / abx$

$$\begin{array}{l} x \rightarrow ab x' \\ x' \rightarrow cx/x \end{array}$$

2) $s \rightarrow \underline{iEts} / \underline{iEtses} / a$
 $E \rightarrow b$

$$\begin{array}{l} s \rightarrow \underline{iEts} s' / a \\ s' \rightarrow es / \epsilon \\ E \rightarrow b \end{array}$$

3) $A \rightarrow \underline{a}AB / \underline{a}BC / \underline{a}AC$

Ans: $\overline{A} \rightarrow a A'$ Step 1

$$A' \rightarrow \underline{AB} / BC / AC$$

$$A \rightarrow a A' \quad \longrightarrow \quad \overline{A} \rightarrow a A' \quad \text{Step 2}$$

$$A' \rightarrow \underline{AA''} / BC$$

$$A'' \rightarrow B / C$$

4) $S \rightarrow a \text{ ssbs} / a \text{sasb} / abb / b$

Step 1

$S \rightarrow as' / b$

$s' \rightarrow \cancel{\text{ssbs}} / \cancel{\text{sasb}} / \cancel{bb} / \cancel{b}$

Step 2

$S \rightarrow as' / b$

~~$S \rightarrow ssbs /$~~

$s' \rightarrow ss'' / bb$

$s'' \rightarrow Sbs / a Sb$

5) $S \rightarrow a \text{ lab} / a \text{bc} / a \text{bcd}$

Step 1

$S \rightarrow as'$

$s' \rightarrow b / bc / bcd / \epsilon$

Step 3

$S \rightarrow as'$

~~$s' \rightarrow b s'' / \epsilon$~~

~~$s'' \rightarrow a b a d / \epsilon$~~

$s'' \rightarrow c s''' / \epsilon$

$s''' \rightarrow d / \epsilon$

Step 2

$S \rightarrow as'$

$s' \rightarrow b s'' / \epsilon$

$s'' \rightarrow c / cd / \epsilon$

Now, can stop!

Parse tree and Abstract Syntax tree

Consider the grammar,

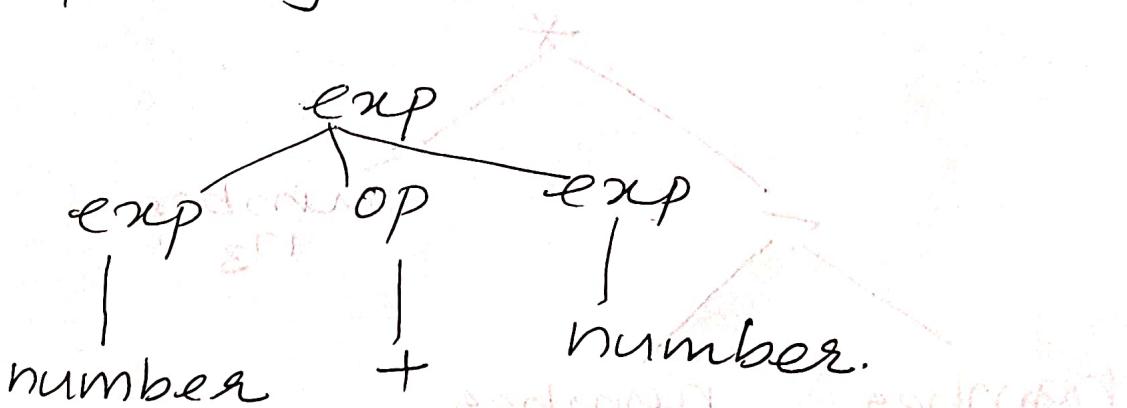
$$\begin{aligned} \text{exp} &\rightarrow \text{exp op exp} \mid (\text{exp}) \mid \text{number} \\ \text{op} &\rightarrow + \mid - \mid * \end{aligned}$$

The derivation for the string

'number + number' become

$$\begin{aligned} \text{exp} &\Rightarrow \text{exp op exp} \\ &\Rightarrow \text{number op } \cancel{\text{exp}} \\ &\Rightarrow \text{number } + \cancel{\text{exp}} \\ &\Rightarrow \text{number } + \text{number}. \end{aligned}$$

corresponding parse tree,



We can apply pre-order numbering (LMD) or post order numbering (RMD). for deriving the same string.

Abstract syntax tree

For number + number,

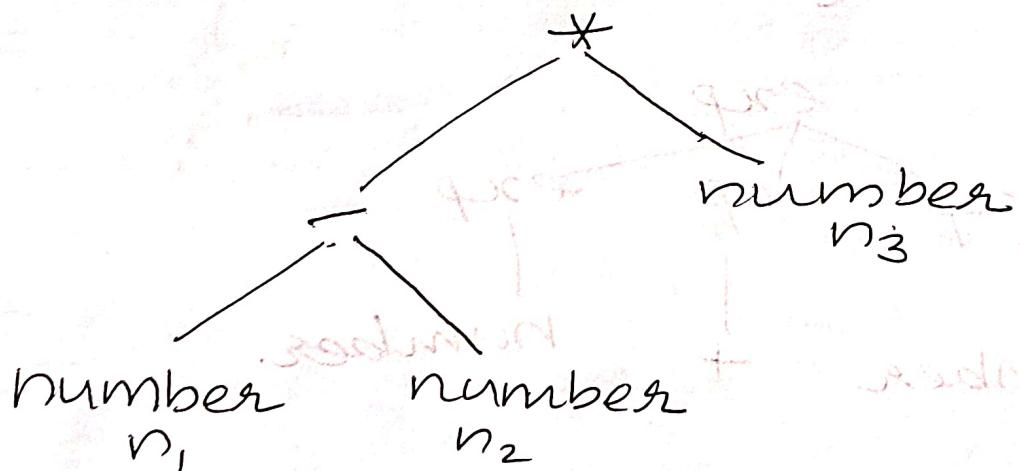
Abstract ST become;



For $3 + 4$,



for $(\text{number} - \text{number}) * \text{number}$,



- A parser will go through all the steps represented by a parse tree, but will usually only construct an abstract syntax tree.
- They contain all the information needed for translation, in a more efficient form than parse tree.

Syntax Error Handling

Common programming errors can occur at many different levels.

I) Lexical errors

- Misspelling of identifiers, keywords or operators.
- Missing quotes around text intended as a string.

2) Syntactic errors

- misplaced semicolons or extra/missing braces.

3) Semantic errors

- Type mismatches b/w operators & operands,

4) Logical errors

→ Anything from incorrect reasoning on the part of the programmer, to the use in a C program of the assignment operator '=' instead of the comparison operator '=='.

→ $\&$, $\&$
→ $\|$

We've passing methods, which allows syntactic errors to be detected very efficiently.

(Methods such as LL and LR can detect errors as soon as possible)

The error handler in a parser has goals such as:

- Report the presence of errors clearly (line, column) and accurately.
- Recover from each error quickly enough to detect subsequent errors.
- Add minimal overhead to the processing of correct programs.

Error-Recovery Strategies

Once an error is detected, how should the parser recover?

The simplest method is for the parser to quit with an error msg when it detects the first error.

There're different error recovery strategies:

1) Panic-mode recovery

On discovering an error, the parser discards symbols one at a time until one of a designated set of 'synchronizing tokens' is found.

→ semicolon, ; etc.

(to guarantee not to go into infinite loop).

2) Phrase-level recovery

A parser may perform local correction on the 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 extraneous semicolon
- insert a missing semicolon.

- The choice of the local correction is left to the compiler designer.
- phrase-level replacement has been used in several error-repairing compilers, as it can correct any input string. But it's difficult to correct.

3) Error Productions

Find from where these errors coming, rewrite or replace those production rules of that particular grammar.

4) Global correction.

Too costly.

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 the parse tree in preorder.
- It uses left most derivation to construct a parse-tree.
- Construction of parse tree starts at the root and proceeds towards the leaves.
- Efficient parsers can be constructed more easily by hand using top-down methods.
- The main problem in TDP is to determine the production that has to be applied for a non-terminal (say A). Once we choose A-production, the parsing process consists of matching terminal symbols in the production body with the input string.
- Applicable to small class of languages.

(1) Recursive - Descent Parsing

- It is a general parsing technique, but not widely used as it is not efficient (slow).
- It involves backtracking.
- So more recursive, i.e; it may require ~~several~~ repeated scans over the input.

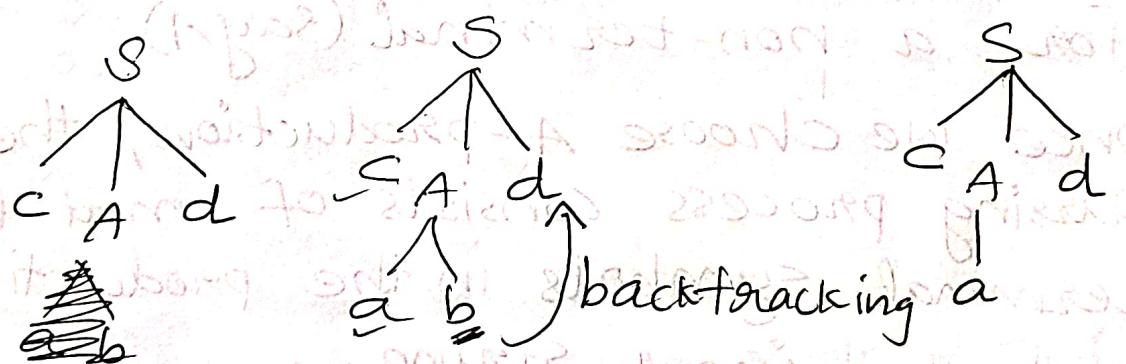
Example:-

Consider the grammar

$$S \rightarrow CAD$$

$$A \rightarrow ab/a$$

Construct a parse tree for $w = cad$



$w = cad$

$c a b d$

not matching.

$w = cad$



(2) Predictive - Parsing / LL(1)

Predictive - parsing or non - recursive descent parsers needing no backtracking, can be constructed for a class of grammars called LL(1).

LL(1) → lookahead symbol.
↓
for left most derivation.

Scanning the input from left to right

'1' - for using one input symbol for lookahead at each step to make parsing action decisions.

→ The class of LL(1) grammars is rich enough to cover most programming constructs.

i.e; No left - recursive or ambiguous grammar can be LL(1).

→ Most popular and powerful grammar in top - down parsing.

Lookahead symbol:-

The next terminal that a compiler will try to match in the input.

I. Find FIRST() and FOLLOW().

II. Create Parsing Table / M Table.

In top-down parsing, FIRST() and FOLLOW() allows us to choose which production to apply, based on the next input symbol.

FIRST()

Rule 1: If x is a terminal, then

$$\text{FIRST}(x) = \{x\}.$$

Rule 2:

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

Rule 3:

If x is a non-terminal and $x \rightarrow y_1 y_2 \dots y_k$ is a production for some $k \geq 1$, then place a in FIRST(x) if for some i , a is in FIRST(y_i), and ϵ is in all of FIRST(y_1), ..., FIRST(y_{i-1}). i.e; $y_1, y_2, \dots, y_{i-1} \Rightarrow \epsilon$

If ϵ is in FIRST(y_j) for all $j = 1, 2, \dots, k$, then add ϵ to FIRST(x).

FOLLOW()

Rule 1:

Place \$ in FOLLOW(\$), where S is the start symbol, and \$ is the input right end-marker.

Rule 2:

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

Rule 3:

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

Consider the grammar.

$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

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

Step 1: Eliminate left recursion and s factoring

Step 2: Find FIRST() and FOLLOW().

	FIRST()	FOLLOW()
E	id, (\$,)
E'	+ , ε	\$,)
T	id, (+ , \$,)
T'	* , ε	+ , \$,)
F	id, (* , + , \$,)

Step 3: LL(1) Table / Predictive parse table / M table.

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

Construction of a predictive parsing table:

Algorithm

Input : Grammar G.
Output : Parsing table M.

Method

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

- 1) For each terminal a in $\text{FIRST}(\alpha)$,
Add $A \rightarrow \alpha$ to $M[A, a]$.
- 2) If ϵ is in $\text{FIRST}(\alpha)$, then for each terminal,
 b in $\text{FOLLOW}(A)$, add $A \rightarrow a$ to $M[A, b]$.

If ϵ is in $\text{FIRST}(\alpha)$ and $\$$ is in $\text{Follow}(A)$,
add $A \rightarrow \alpha$ to $M[A, \$]$ as well.

If after performing the above, there's
no production at all in $M[A, a]$, then
Set $M[A, a]$ to error (an empty entry).

Input

.				a+b	\$
---	--	--	--	-----	----

Predictive
Parsing
Program

Output

x
y
z
\$

Stack

Parsing
Table
 M

Fig: Model of a table driven Predictive Parser.

Example 2

$$S \rightarrow (L) / a$$

$$L \rightarrow L, S / S$$

Step 1: Elimination of left recursion

$$S \rightarrow (L) / a$$

$$L \rightarrow S L'$$

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

Step 2: Elimination of left factoring.
Here, no common prefixes.

Step 3: Find FIRST() and FOLLOW().

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

$$\text{FIRST}(L) = \{ (, a \} \}$$

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

$$\text{FOLLOW}(S) = \{ ,) \} \}$$

$$\text{FOLLOW}(L) = \{) \} \}$$

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

(S is the starting symbol.)

LL(1) (parsing table.

	()	a	,	\$
S	$S \rightarrow (L)$		$S \rightarrow a$		
L	$L \rightarrow SL'$		$L \rightarrow SL'$		
L'		$L' \rightarrow \epsilon$		$L' \rightarrow SL'$	

(a) \$

Stack

Up starting

Action

\$S

(a) \$

$S \rightarrow (L)$

\$) L C

(a) \$

POP

\$) L

a) \$

$L \rightarrow SL'$

\$) L' S

a) \$

$S \rightarrow a$

\$) L' a

a) \$

POP

\$) L'

(\$)

$L' \rightarrow \epsilon$

\$)

)

POP

\$

\$

Accepted.

Example - 3

$$S \rightarrow aAB \mid bA \mid \epsilon$$

$$A \rightarrow aAb \mid \epsilon$$

$$B \rightarrow bB \mid \epsilon$$

Construct parsing table for the given grammar. Also check whether 'aabbb' is accepted or not.

	FIRST()			FOLLOW()
S	a, b, ϵ			{ \$ }
A	a, ϵ			{ b, \$ }
B	b, ϵ			{ \$ }
	a	b		\$
S	$S \rightarrow aAB$	$S \rightarrow bA$	$S \rightarrow \epsilon$	
A	$A \rightarrow aAb$	$A \rightarrow \epsilon$		$A \rightarrow \epsilon$
B	$B \rightarrow bB$	$B \rightarrow bB$		$B \rightarrow \epsilon$

stackInput stringAction $\$ S$

aa bb \$

 $S \rightarrow aAB$ $\$ BA_a$

aabb \$

POP

 $\$ BA$

abb \$

 $A \rightarrow aAb$ ~~$\$ BbAba$~~ ~~a b b \$~~

POP

 $\$ BbA$

bb \$

 $A \rightarrow \epsilon$ $\$ Bb$

bb \$

POP

 $\$ B$

b \$

 $B \rightarrow bB$ $\$ Bb$

b \$

POP

 $\$ B$

\$

 $B \rightarrow \epsilon$ $\$$

\$

Accepted