## 2.1 Role of Syntax Analyzer or Parser

The analysis phase of a compiler breaks up a source pgm into constituent pieces and produces an intermediate representation for it, called intermediate code. The synthesis phase translates the intermediate code into the target pgm.

Analysis is organized around the "syntax" of the language to be compiled. The syntax of a pgming language describes the proper form of its pgms, while the semantics of the language defines what its pgms mean, that is, what each pgm does when it executes

- Context-free grammars or BNF (Backus Naur Form) are widely used for specifying syntax
- Informal descriptions and suggestive examples are used for specifying semantics.

- The parser obtains a string of tokens from the lexical analyzer and verifies that the string of token names can be generated by the grammar for the source language.

- The parser constructs a parse tree and passes it to the rest of the compiler for further processing.
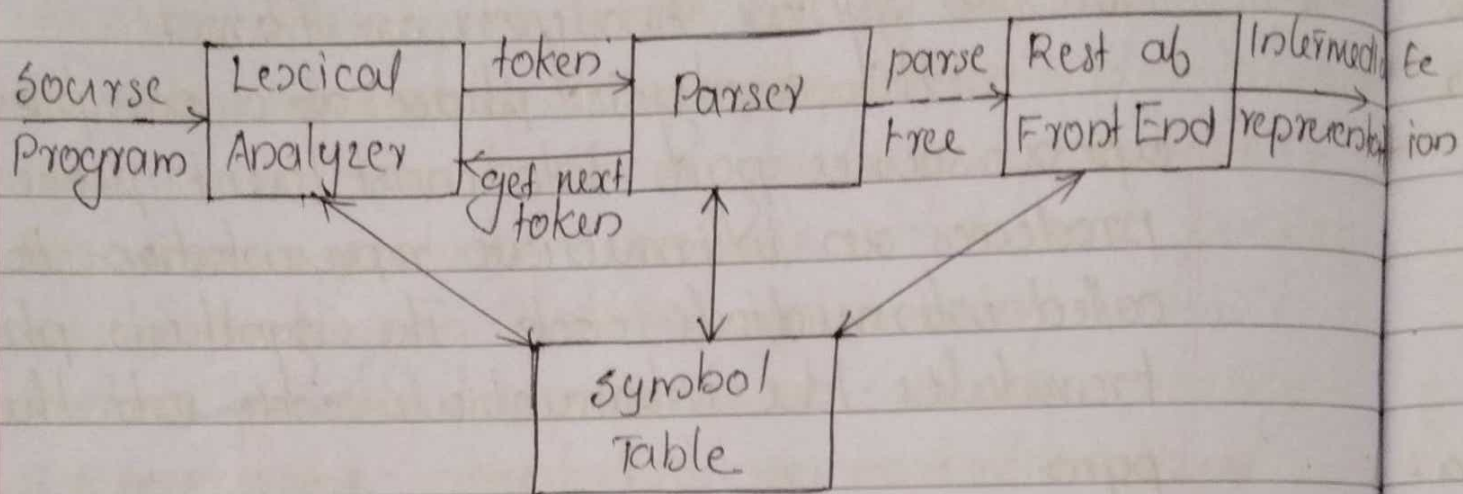
Fig : position of parser in compiler model.

– The parser should report any syntax error and recover from commonly occuring errors to continue processing the remainder of the pgm.

A grammar generally describes the hierarchical structure of most pgming language constructs.

–There are three general types of parsers for grammars.

1. Universal – can parse any grammar. Too inefficient to use in production compilers

2. Top-down – build parse tree from top (Root) to the bottom (leaves)    moves from top to leaves

3. Bottom-up – start from the leaves and work their way up to the root.    moves from leaves to root

The I/p to the parser is scanned from left ts right one symbol at a time.

## Definition ab Grammars / Context-free grammar.

A context-free grammar has both components (V, T, P, S).

1) A set ab non-terminal symbols, sometimes referred to as tokens. Tokens are elementary symbols ab the language debined by the grammar

2. A set ab nonterminals - sometimes called syntactic variables. Each nonterminal represent or set ab strings ab terminals (eg; stmt expr).

3. Productions, where each pdn consists ab a non-terminal, called the head or leftside ab the pdn; an arrow (→ or ::=), and a sequence ab terminals or nonterminals, called body or right side ab the pdn.

4. Start symbol. one non-terminal is designated as the start symbol

- grammars are specified by listing their pdns, with the pdns bor the start symbol listed birst

eg.
list ⟶ list + digit
list ⟶ list - digit
list ⟶ digit
digit ⟶ 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

can be also written as
list → list + digit / list - digit
digit

0, 1, 2, 3, 4, 5, 6, 7, 8, 9, +, - are terminals.

here, list, digit are nonterminals

- A grammar derives strings by beginning with the start symbol and repeatedly replacing a non terminals by the body of a pdn for that nonterminal.

eg: Consider the context-free grammar

$$S \rightarrow SS+ | SS* | a$$

show how the string $aa+a*$ can be generated by this grammar

Ans.
$$S \rightarrow SS*$$
$$\rightarrow SS+S*$$
$$\rightarrow aS+S*$$
$$\rightarrow aa+S*$$
$$\rightarrow aa+a*$$

Parsing is the pdm of taking a string of terminals and figuring out how to derive it from the start symbol of the grammar, and if it cannot be derived from the start symbol of the grammar then reporting syntax errors within the string

## Parse trees

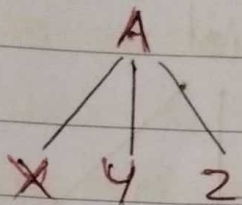- A parse tree pictorially shows how the start symbol of a grammar derives a string in the language

eg: If nonterminal A has a pdn

$$A \rightarrow XYZ,$$ then the parse tree may have

an interior node labeled A with three children labeled x, y and z from left to right



Thus, given a context free grammar, a parse tree according to grammar is a tree with the following properties:

1. The root is labeled by the start symbol
2. Each leaf is labeled by a terminal or $\epsilon$
3. Each interior node is labeled by a nonterminal
4. If A is the nonterminal labeling some interior nodes and $x_1, x_2 \cdots x_n$ are the labels of the children of that node from left to right, then there must be a pdn $A \rightarrow x_1 x_2 \cdots x_n$. Here, $x_1, x_2 \cdots, x_n$, each stand for a symbol that is either a terminal or a nonterminal.

   If $A \rightarrow \epsilon$ is a pdn, then a node labeled A may have a single child labeled $\epsilon$.

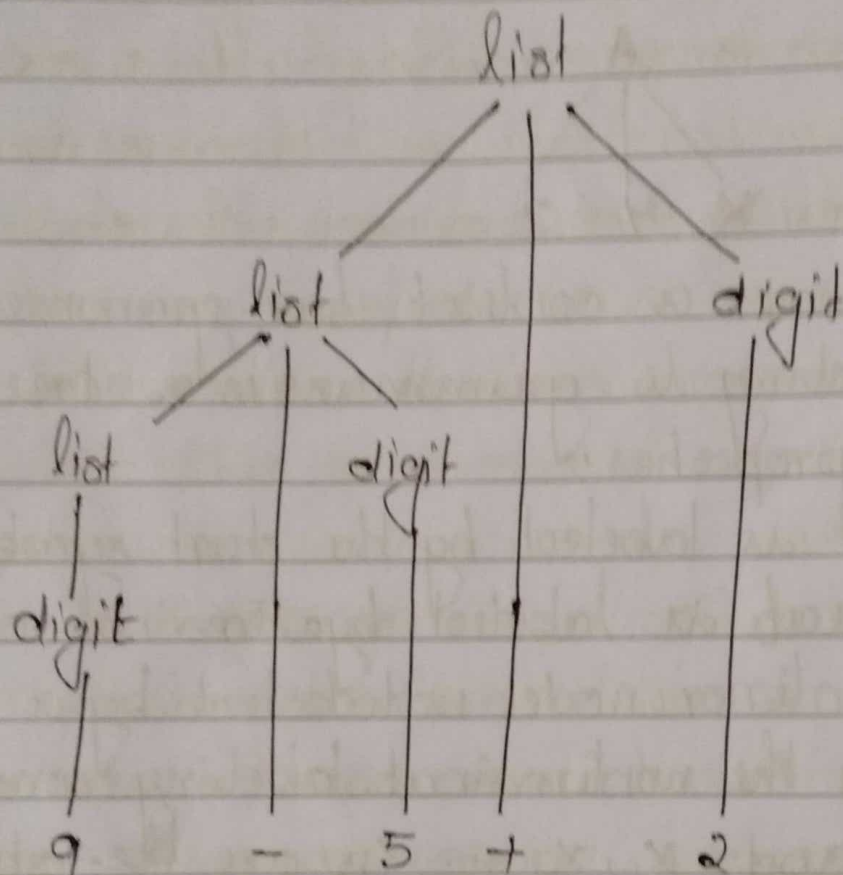   Consider the following pdn,

   $$list \rightarrow list + digit$$
   $$list \rightarrow list - digit$$
   $$list \rightarrow digit$$
   $$digit \rightarrow 0|1|2|3|4|5|6|7|8|9$$

Parse tree for 9-5+2 according to the grammar



- The following grammar is used for simple arithmetic expressions.

expression → expression + term
expression → expression - term
expression → term
term → term * factor
term → term / factor
term → factor
factor → (expression)
factor → id

expression, factor and term are non terminals
expression is the start symbol.
id, + - * ( ) - terminals.

## Notational Conventions

The following symbols are terminals

a) Lowercase letters such as a, b, c etc
b) Operator symbols such as +, -, * etc
c) Punctuation symbols such as parentheses, comma -
d) digits such as 0, 1, ..., 9
e) Boldface strings such as id, or if

The following symbols are nonterminals

a) Uppercase letters such as A, B, C ...
b) letter s is usually the start symbol
c) lowercase italic names such as expr or stmt
~~d)~~ when discussing parsing construct, uppercase letters may be used to represent nonterminals eg: expr, term
eg: expr, term → E, T resp.

So the above grammar can be replaced by

$$E \rightarrow E+T \mid E-T \mid T$$
$$T \rightarrow T*F \mid T/F \mid F$$
$$F \rightarrow (E) \mid \text{id}.$$

here, E, T and F are nonterminals, with E as the start symbol. The remaining symbols are terminals.

## Derivations:

Productions are treated as rewriting rules. Beginning with the start symbol, each rewriting

step replaces a nonterminal by body of one of its polns.

Eg. Consider, the following grammar, with a single nonterminal $E$

$$E \rightarrow E+E \mid E*E \mid -E \mid (E) \mid id$$

The following sequence represents the derivation of $-(id)$ from $E$

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

Such a sequence of replacement is called derivation of $-(id)$ from $E$. This derivation provides a proof that the string $-(id)$ is one particular instance of an expression.

<span style="color:red">Definition of derivation</span>

Consider a nonterminal $A$ in the middle of a sequence of grammar symbols, as in $\alpha A \beta$ where $\alpha$ and $\beta$ are arbitrary strings of grammar symbols. Suppose $A \rightarrow \gamma$ is a poln, then we can write $\alpha A \beta \Rightarrow \alpha \gamma \beta$.

The symbol $\Rightarrow$ means "derives in one step"

eg! $E \Rightarrow -E$, $E \Rightarrow -(E)$

$\overset{*}{\Rightarrow}$ means derivation in zero or more steps

eg! $E \overset{*}{\Rightarrow} -(id)$, $E \overset{*}{\Rightarrow} E$, $E \overset{*}{\Rightarrow} (E)$

$\xRightarrow{+}$ means, derivation in one or more steps.

eg: $E \xRightarrow{+} -(id)$

Sentential form:

A sentential form may contain both terminals & nonterminals, and may be empty.

eg: $S \xRightarrow{*} \alpha$,    s - start symbol

α - sentential form of G

A sentence of G is a sentential form with no non-terminals. The language generated by a grammar is its set of sentences.

Thus, ω is in L (G), iff ω is a sentence of G (or $S \xRightarrow{*} \omega$)

A language that can be generated by a grammar is said to be a context-free language. If two grammars generate the same language, the grammars are said to be equivalent.

The string $-(id + id)$ is a sentence of grammar because there is a derivation

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

The strings $E, -E, -(E), \ldots, -(id, id)$ are all sentential forms of this grammar

we write $E \xRightarrow{*} -(id + id)$ to indicate that

$-(id + id)$ can be derived from $E$.

$$LMD \quad E \Rightarrow E*E$$
$$\Rightarrow E+E*E$$
$$\Rightarrow id+E*E$$
$$\Rightarrow id+id*E$$
$$id+id*i^+$$

$$E*E$$
$$\Rightarrow E*E$$
$$\Rightarrow E*id$$
$$\Rightarrow E+E*id$$
$$\Rightarrow E+id*id$$
$$\Rightarrow id+id*i^-$$

- There are two types of derivations.

1. **Leftmost derivation**: For each derivation the leftmost nonterminal in each sentential form is replaced. If $\alpha \Rightarrow \beta$ is a step in which the leftmost nonterminal in $\alpha$ is replaced, we write $\alpha \underset{lm}{\Rightarrow} \beta$

eg: $E \underset{lm}{\Rightarrow} -E \underset{lm}{\Rightarrow} -(E) \underset{lm}{\Rightarrow} -(E+E) \underset{lm}{\Rightarrow} -(id+E)$

$$\underset{lm}{\Rightarrow} -(id+id)$$

2. **Rightmost derivation**: For each derivation, the rightmost nonterminal in each sentential form is replaced, we write $\alpha \Rightarrow \beta$ in this case.

eg: $E \underset{rm}{\Rightarrow} -E \underset{rm}{\Rightarrow} -(E+E) \underset{rm}{\Rightarrow} -(E+id) \underset{rm}{\Rightarrow} (id+id$
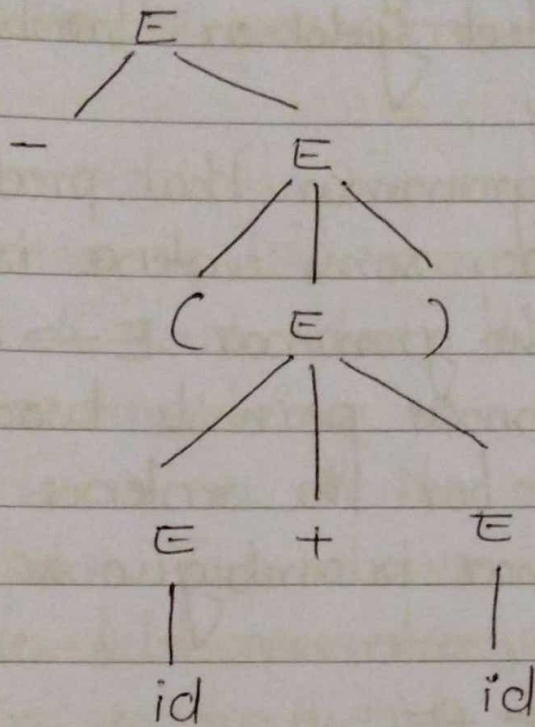
Rightmost derivations are some times called canonical derivations.

**Parse Trees and Derivations:**
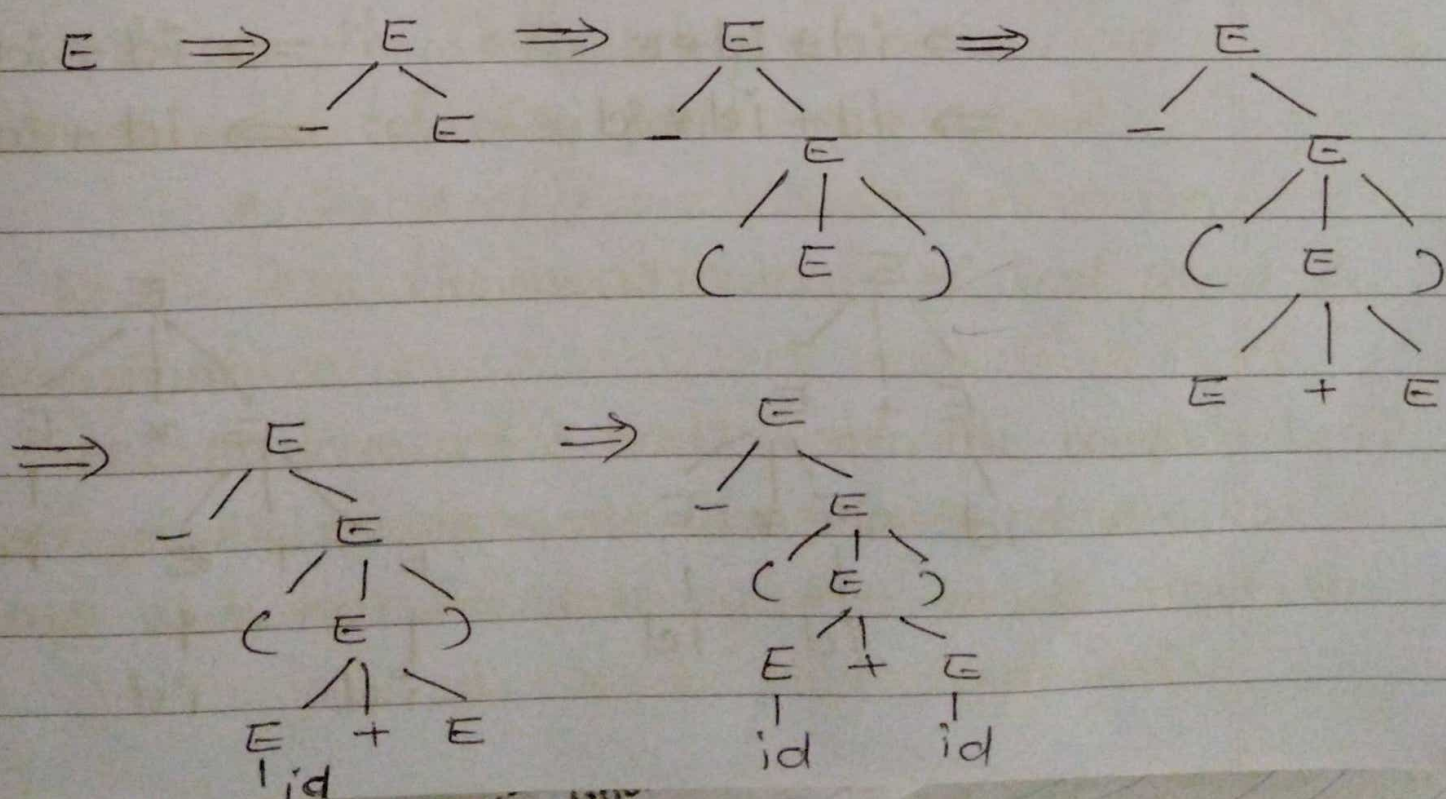
- A parse tree is a graphical representation of a derivation that filters out the order in which productions are applied to replace nonterm...

- Each interior node of parse tree represents the application of a pdn.

- The children of the node are labeled from left to right, by the symbols in the body of the pdn.

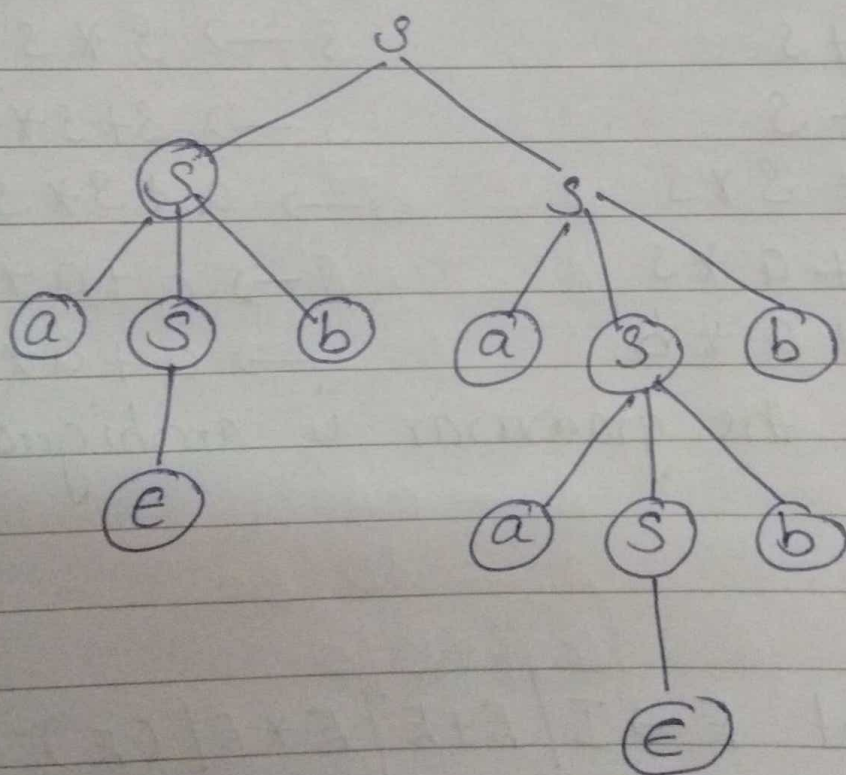eg. parse tree for $-(id+id)$



Sequence of parse trees for $-(id+id)$

# Example for Parse Tree

$CFG = (\{S\}, \{a,b\}, \{P = S \to SS \,/\, aSb \,/\, \epsilon\}, \{S\})$

find the derivation and parse tree for the string abaabb.

$S \to SS \to aSbS \to abS \to abaSb \to abaaSb$
$\to abaabb$

# Ambiguous Grammar

A grammar is said to be ambiguous if there exist two or more derivation tree for a string w.

<u>Example</u>  $G = (\{S\}, \{a+b, +, *\}, P, S)$

where 'p consists of $S \rightarrow S+S / S*S / a/b$ the string $a+a*b$ can be generated as

| | |
|---|---|
| $S \rightarrow S+S$ | $S \rightarrow S*S$ |
| $\rightarrow a+S$ | $S \rightarrow S+S*S$ |
| $\rightarrow a+S*S$ | $S \rightarrow a+S*S$ |
| $\rightarrow a+a*S$ | $S \rightarrow a+a*S$ |
| $\rightarrow a+a*b$ | $\rightarrow a+a*b$ |

Thus the grammar is ambiguous.

# Eliminating left recursion

Let G be a context $A \to \alpha \cdot B\beta, a/b$
free grammar. A production $B \to \gamma \cdot c/d$
of G is said left recursive if it has the form
$\text{First}(\beta) = \{c, d\}$
$$A \to A\alpha$$
where A is a non-terminal and $B \to \gamma \cdot a/b$
$\alpha$ is a string of grammar symbols.

$\to$  $A \to A\alpha / \beta$    by    $A \to \beta A'$
$$A' \to \alpha A' / e$$

eg.
$$E \to E+T / T$$
$$T \to T*F / F$$
$$F \to (E) / id$$

$$E \to TE'$$
$$E' \to +TE' / e$$
$$T \to FT'$$
$$T' \to *FT' / e$$
$$F \to (E) / id$$

# Eliminating left factoring

$A \to \alpha\beta_1 / \alpha\beta_2$ are two A production

$$A \to \alpha A'$$
$$A' \to \beta_1 / \beta_2.$$

① $S \rightarrow iEtS \mid iEtSeS \mid a$

$E \rightarrow b$

$\downarrow$

$S \rightarrow iEtSS' \mid a$

$S' \rightarrow eS \mid \epsilon$

$E \rightarrow b$

② $A \rightarrow aA'$      $A \rightarrow aAB \mid aB\epsilon \mid aAc$

$A' \rightarrow$

$\downarrow$

$A \rightarrow aA'$

$A' \rightarrow AB \mid Bc \mid Ac$

$\downarrow$

$A \rightarrow aA'$

$A' \rightarrow AD \mid Bc$

$D \rightarrow B \mid c$

③ $S \rightarrow bSSaaS \;/\; bSSaSb \;/\; bSb \;/\; a$

$\Downarrow$

$S \rightarrow bSS' \;/\; a$
$S' \rightarrow SaaS \;/\; SaSb \;/\; b$

$\Downarrow$

$S \rightarrow bSS' \;/\; a$
$S' \rightarrow SaA \;/\; b$
$A \rightarrow aS \;/\; Sb$

- For most parsers, it is desirable that the grammar be made unambiguous, for if it is not, we cannot uniquely determine which parse tree to select for a sentence.

- In other cases, it is convenient to use carefully chosen ambiguous grammars, together with disambiguating rules that "throw away" undesirable parse trees, leaving only one tree for each sentence.
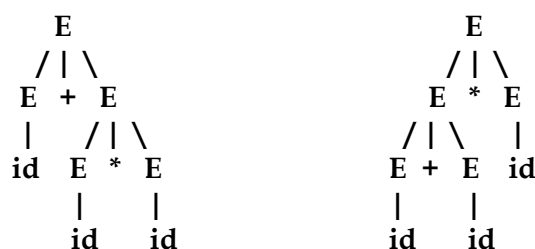
    **EXAMPLE**

    Consider very simple sentence id+ id * id.

| 1st Leftmost Derivation | 2nd Leftmost Derivation |
|---|---|
| E ===> E + E | E ===> E * E |
| ===> id + E | ===> E + E * E |
| ===> id + E * E | ===> id + id * E |
| ===> id + id * E | ===> id + id * E |
| ===> id + id * id | ===> id + id * id |

**1st Parse Tree**      **2nd Parse Tree**

```
      E                       E
    / | \                   / | \
   E  +  E                 E  *  E
   |    / | \             / | \    |
   id  E  *  E           E  +  E   id
       |     |           |     |
       id    id          id    id
```
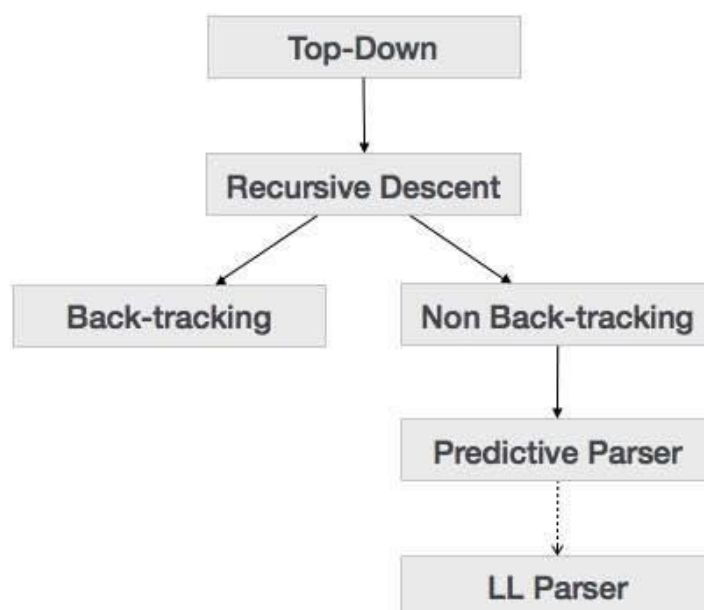
# 2.2 TOP DOWN PARSING

- Parsing is the process of determining if a string of token can be generated by a grammar.

- Mainly 2 parsing approaches:

    - **Top Down Parsing**

    - **Bottom Up Parsing**

- In **top down parsing**, parse tree is constructed from top (root) to the bottom (leaves).

- In **bottom up parsing**, parse tree is constructed from bottom (leaves)) to the top (root).

- It can be viewed as an attempt to construct a parse tree for the input starting from the root and creating the nodes of parse tree in preorder.

- Pre-order traversal means: 1. Visit the root 2. Traverse left subtree 3. Traverse right subtree.

➕ Top down parsing can be viewed as an attempt to find a leftmost derivation for an input string (that is expanding the leftmost terminal at every step).



## 2.2.1 RECURSIVE DESCENT PARSING

It is the most general form of top-down parsing.

It may involve **backtracking**, that is making repeated scans of input, to obtain the correct expansion of the leftmost non-terminal. Unless the grammar is ambiguous or left-recursive, it finds a suitable parse tree
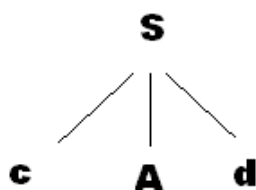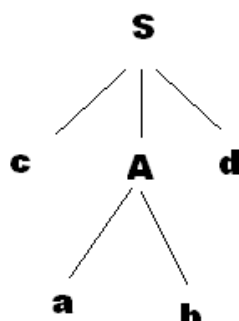
**EXAMPLE**

**Consider the grammar:**

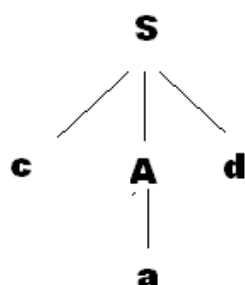**S → cAd**

**A → ab | a**

**and the input string w = cad.**

❖ To construct a parse tree for this string top down, we initially create a tree consisting of a single node labelled **S**.

❖ An input pointer points to **c**, the first symbol of w. **S** has only one production, so we use it to expand **S** and obtain the tree as:

❖ The leftmost leaf, labeled **c**, matches the first symbol of input w, so we advance the input pointer to **a**, the second symbol of **w**, and consider the next leaf, labeled **A**.

❖ Now, we expand **A** using the first alternative **A → ab** to obtain the tree as:



❖ We have a match for the second input symbol, **a**, so we advance the input pointer to **d**, the third input symbol, and compare d against the next leaf, labeled **b**.

❖ Since **b** does not match **d**, we report failure and go back to **A** to see whether there is another alternative for **A** that has not been tried, but that might produce a match.

❖ In going back to **A**, we must reset the input pointer to position 2 , the position it had when we first came to **A**, which means that the procedure for **A** must store the input pointer in a local variable.

❖ The second alternative for **A** produces the tree as:



❖ The leaf **a** matches the second symbol of **w** and the leaf **d** matches the third symbol. Since we have produced a parse tree for **w**, we halt and announce successful completion of parsing. (that is the string parsed completely and the parser stops).

❖ The leaf a matches the second symbol of w and the leaf d matches the third symbol. Since we have produced a parse tree for w, we halt and announce successful completion of parsing. (that is the string parsed completely and the parser stops).

## 2.2.2 PREDICTIVE PARSING

✦ A predictive parsing is a special form of recursive-descent parsing, in which the current input token unambiguously determines the production to be applied at each

step. The goal of predictive parsing is to construct a top-down parser that never backtracks. To do so, we must transform a grammar in two ways:

> **Eliminate left recursion, and**

> **Perform left factoring.**

➕ These rules eliminate most common causes for backtracking although they do not guarantee a completely backtrack-free parsing (called LL(1) as we will see later).

## Left Recursion

➕ A grammar is said to be left –recursive if it has a non-terminal A such that there is a derivation A → Aα, for some string α.

**EXAMPLE**

Consider the grammar

**A → Aα**

**A → β**

❖ It recognizes the regular expression βα*. The problem is that if we use the first production for top-down derivation, we will fall into an infinite derivation chain. This is called left recursion.

❖ Top–down parsing methods cannot handle left recursive grammars, so a transformation that eliminates left-recursion is needed. The left-recursive pair of productions **A → Aα|β** could be replaced by two non-recursive productions.

$$A → β\ A'$$
$$A' →\ α\ A'\ |\ ε$$

➕ Consider The following grammar which generates arithmetic expressions

**E → E + T | T**

**T → T * F | F**

**F → ( E ) | id**

Eliminating the immediate left recursion to the productions for E and then for T, we obtain

**E → T E'**

**E' → + T E' | ε**

**T → F T'**

**T' → * F T' | ε**

**F → ( E ) | id**

- No matter how many A-productions there are, we can eliminate immediate left recursion from them by the following technique. First, we group the A productions as

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \ldots \mid A\alpha_m \mid \beta_1 \mid \beta_2 \mid \ldots \mid \beta_n$$

where no $\beta_i$ begins with an **A**. Then we replace the **A**-productions by

$$A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \ldots \mid \beta_n A'$$

$$A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \ldots \mid \alpha_m A' \mid \varepsilon$$

## Left Factoring

- Left factoring is a grammar transformation that is useful for producing a grammar suitable for predictive parsing.

- The basic idea is that when it is not clear which of two alternative productions to use to expand a non-terminal A, we may be able to rewrite the A-productions to defer the decision until we have seen enough of the input to make the right choice

$$A \rightarrow \alpha \beta_1 \mid \alpha \beta_2$$

- are two A-productions, and the input begins with a non-empty string derived from $\alpha$ we do not know whether to expand **A** to $\alpha \beta_1$ or $\alpha \beta_2$.

- However, we may defer the decision by expanding **A** to $\alpha$**B**. Then, after seeing the input derived from $\alpha$, we may expand **B** to $\beta_1$ or $\beta_2$ .

- The left factored original expression becomes:

$$A \rightarrow \alpha B$$
$$B \rightarrow \beta_1 \mid \beta_2$$

- For the "dangling else "grammar:

  **stmt →if cond then stmt else stmt | if cond then stmt**

  The corresponding left – factored grammar is:

  **stmt → if cond then stmt else_clause**

  **else_clause → else stmt | ε**

## Non Recursive Predictive parser

- It is possible to build a nonrecursive predictive parser by maintaining a stack explicitly, rather than implicitly via recursive calls.

- The key problem during predictive parsing is that of determining the production to be applied for a nonterminal.

- The nonrecursive parser in looks up the production to be applied in a parsing table

### Requirements

1. Stackv

2. Parsing Table

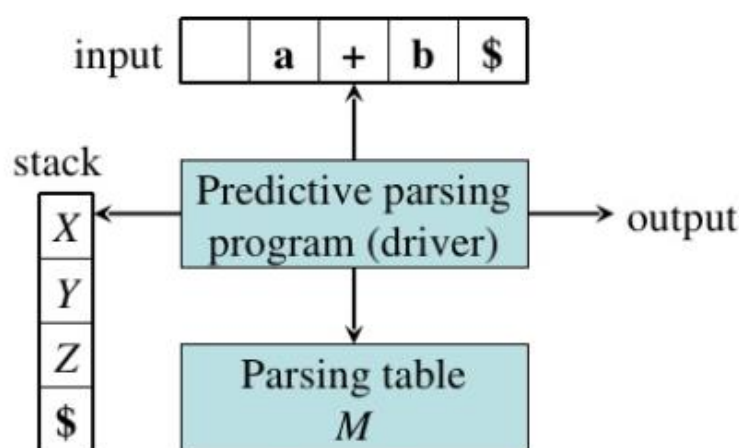3. Input Buffer

4. Parsing



**Figure** : *Model of a nonrecursive predictive parser*

🞣 **Input buffer** - contains the string to be parsed, followed by $(used to indicate end of input string)

🞣 **Stack** – initialized with $, to indicate bottom of stack.

🞣 **Parsing table** - 2 D array M[A,a] where A is a nonterminal and a is terminal or the symbol $

🞣 The parser is controlled by a program that behaves as follows. The program considers **X,** the symbol on top of the stack, and **a** current input symbol. These two symbols determine the action of the parser.

🞣 There are three possibilities,

1. If **X = a = $** , the parser halts and announces successful completion of parsing.

2. If **X = a ≠ $** , the parser pops **X** off the stack and advances the input pointer to the next input symbol,

3. If **X** is a nonterminal, the program consults entry **M|X, a |** of the parsing table **M**. The entry will be either an **X**-production of the grammar or an error entry. If, for example, **M |X, u |= {X → UVW}**, the parser replaces **X** on top of the stack by WVU (with U on top). As output we shall assume that the parser just prints the production.Ton used; any other code could be executed here. If **M|X, a| = error**, the parser calls an error recovery routine.

## Predictive Parsing Algorithm

**INPUT:** A string **w** and a parsing table *M* for grammar *G*.

**OUTPUT:** If **w** is in **L ( G )** , a leftmost derivation of **w**; otherwise, an error indication.

**METHOD**: Initially, the parser is in a configuration in which it has **$S** on the stack with **S**, the start symbol of **G** on top, and **w$** in the input buffer. The program that utilizes the predictive parsing table M to produce a parse for the input is shown below.


set ip to point to the first symbol of **w$**;

**repeat**

let **X** be the lop stack symbol and **a** the symbol pointed to by ip;

if **X** is a terminal or $ then

if **X = a** then

pop **X** from the stack and advance ip

else error ()

else   /* **X** is a nonterminal */

if **M | X, a | = X → Y₁ Y₂** . . . $Y_i$ then begin

pop **X** from the stack;

push **Yₖ, Yₖ₋₁,** . . . , **Yₗ** onto the stack, with **Yₗ** on top;

output the production **X → Y₁ Y₂** . . ..**Yₖ**

end

else error ()

until **X = S** /* stack is empty */


**EXAMPLE**

Consider Grammar:

E → T E′

E' → +T E' | Є

T → F T'

T' → * F T' | Є

F → ( E ) | id

## Construction Of Predictive Parsing Table

- Uses 2 functions:

  - **FIRST()**

  - **FOLLOW()**

- These functions allows us to fill the entries of predictive parsing table

### FIRST

- If 'α' is any string of grammar symbols, then FIRST(α) be the set of terminals that begin the string derived from α . If α==*>ϵ then add ϵ to FIRST(α).First is defined for both terminals and non terminals.

- To Compute First Set

  1. If **X** is **a** terminal , then FIRST(**X**) is {**X**}

  2. If **X→ ϵ** then add **ϵ** to FIRST(**X**)

  3. If **X** is a non terminal and **X→Y₁Y₂Y₃...Yₙ** , then put '**a**' in FIRST(**X**) if for some **i**, **a** is in FIRST(**Yᵢ**) and ϵ is in all of FIRST(**Y₁**),...FIRST(**Yᵢ₋₁**).

**EXAMPLE**

Consider Grammar:

> **E → T E′**
>
> **E' → +T E' | Є**
>
> **T → F T'**
>
> **T' → * F T' | Є**
>
> **F → ( E ) | id**

| Non-terminal | FIRST |
|:---:|:---:|
| E | (, id |
| E' | +, Є |
| T | (, id |
| T' | *, Є |
| F | (, id |

## 2.2.2 PREDICTIVE PARSING

A predictive parsing is a special form of recursive-descent parsing, in which the current input token unambiguously determines the production to be applied at each step. The goal of predictive parsing is to construct a top-down parser that never backtracks. To do so, we must transform a grammar in two ways:

Eliminate left recursion, and Perform left factoring.

These rules eliminate most common causes for backtracking although they do not guarantee a completely backtrack-free parsing (called LL(1) as we will see later)

## Nonrecursive Predictive Parsing:

- A nonrecursive predictive parser can be built by maintaining a stack explicitly, rather than implicitly via recursive calls.

The table-driven parser has an i/p buffer, a stack containing a sequence of grammar symbols, a parsing table and an o/p stream.

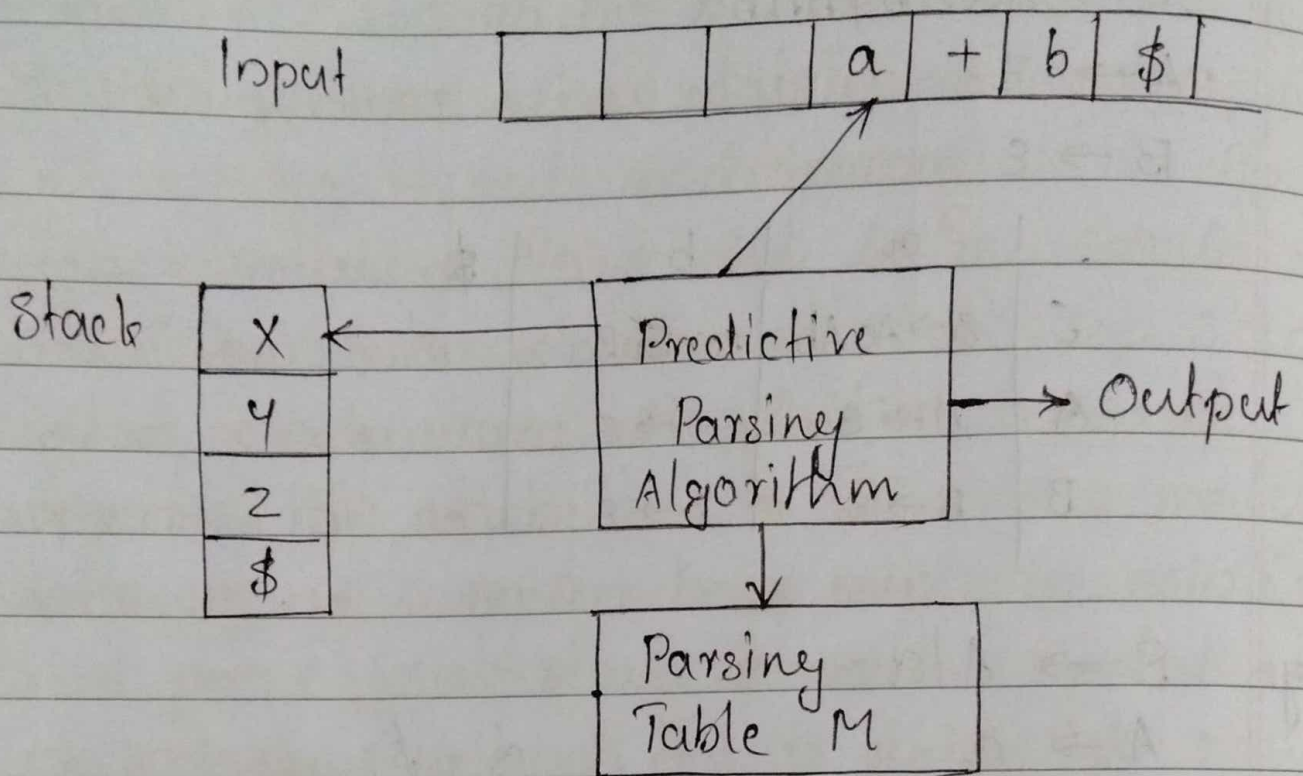- The i/p buffer contains the string to be parsed, followed by the end marker $.

Fig: Model of a predic table-driven predictive Parser

- The symbol $ is reused to mark the bottom of the stack, which initially contains the start symbol of the grammar on top of $.

- The parser is controlled by a program that considers X, the symbol on top of the stack, and a the current i/p symbol. If X is a non-terminal the parser chooses as X- pdn by consulting entry M[X,a] of the parsing table M. otherwise, it checks for a match between the terminal X and current i/p symbol a.

## Algorithm:

**INPUT** : A string $w$ and a parsing table $M$ for grammar $G$.

**OUTPUT** : If $w$ is in $L(G)$, a leftmost derivation of $w$: otherwise, an error indication.

**METHOD** : Initially, the parser is in a configuration with $w\$$ in the i/p buffer and the start symbol $S$ of $G$ on top of the stack, above $\$$.

Set ip to point to the first symbol of $w$; (ie, a)
set x to the top of stack symbol;
while $(x \neq \$)$ /* stack not empty */
{
    If $(x = a)$ pop the stack and advance ip;
    else if $(x$ is a terminal) error();
    else if $(M[x, a]$ is an error entry) error();
    else if $(M[x, a] = x \to y_1 y_2 \cdots y_k)$
    {
        output the production $x \to y_1 y_2 \cdots y_k$;
        pop the stack;
        push $y_k, y_{k-1} \cdots y_1$ on to the stack with $y_1$ on tops;
    }
    set x to the top of stack symbol;
}

| NON Terminal | I/P SYMBOL | | | | | |
|---|---|---|---|---|---|---|
| | id | + | * | ( | ) | $ |
| E | $E \to TE'$ | | | $E \to TE'$ | | |
| E' | | $E' \to +TE'$ | | | $E' \to \varepsilon$ | $E' \to \varepsilon$ |
| T | $T \to FT'$ | | | $T \to FT'$ | | |
| T' | | $T' \to \varepsilon$ | $T' \to *FT'$ | | $T' \to \varepsilon$ | $T' \to \varepsilon$ |
| F | $F \to id$ | | | $F \to (E)$ | | |

eg: Consider the following grammar

$$E \to TE'$$
$$E' \to +TE' \mid \varepsilon$$
$$T \to ET'$$
$$T' \to *FT' \mid \varepsilon$$
$$F \to (E) \mid id.$$

Input : id + id * id

The moves corresponds to the LMD

$$E \underset{lm}{\Rightarrow} TE' \underset{lm}{\Rightarrow} FT'E' \Rightarrow idT'E' \Rightarrow idE' \Rightarrow id+TE' \Rightarrow$$

| Matched | Stack | Input | Action |
|---|---|---|---|
| | E$ | id+id*id$ | |
| | TE'$ | id+id*id$ | output $E \to TE'$ |
| | FT'E'$ | id+id*id$ | output $T \to FT'$ |
| | idT'E'$ | id+id*id$ | output $F \to id$ |
| id | T'E'$ | +id*id$ | match id |
| id | E'$ | +id*id$ | output $T' \to \varepsilon$ |
| id | +TE'$ | +id*id$ | output $E' \to +TE'$ |
| id+ | TE'$ | id*id$ | match + |
| id+ | FT'E'$ | id*id$ | output $T \to FT'$ |
| id+ | idT'E'$ | id*id$ | output $F \to id$ |
| id+id | T'E'$ | *id$ | match id |
| id+id | *FT'E'$ | *id$ | output $T' \to *FT'$ |
| id+id* | FT'E'$ | id$ | match * |
| id+id* | idT'E'$ | id$ | output $F \to id$ |

| id + id * id | T'E'$ | $ | match id |
| id + id * id | E'$ | $ | output T' → ε |
| id + id * id | $ | $ | output E' → ε |

- An error is detected during predictive parsing when the terminal on top of the stack does not match the next i/p symbol or when non-terminal A is on top of the stack, a is the next i/p symbol, and M [A, a] is error (i.e, parsing table entry is empty).

# FIRST and FOLLOW

- The construction ab both top-docon and bottom-cup parsers is aided by two functions, FIRST and FOLLOW, associated with grammar G.

- During top-docon parsing, FIRST and FOLLOW allows us to choose which pdn to apply, based on the next i/p symbol.

- Define FIRST $(\alpha)$, where $\alpha$ is any string ab grammar symbols to be the set ab terminals that begin strings derived from $\alpha$.

eg: If $\alpha \overset{*}{\Rightarrow} \epsilon$, then $\epsilon$ is also in FIRST $(\alpha)$

If $S \Rightarrow \alpha A a \beta \Rightarrow \alpha c \gamma a \beta$,
then $c$ is in FIRST $(A)$

- Consider two A-pdns, $A \rightarrow \alpha | \beta$, where FIRST $(\alpha)$ and FIRST $(\beta)$ are disjoint sets. We can

then choose between these A pdns by looking at the next i/p symbol a, since a can be in atmost one ab FIRST ($\alpha$) and FIRST ($\beta$), not both.

For instance, if a is in FIRST($\beta$) choose the pdn A $\rightarrow$ $\beta$.

To compule FIRST (x) for all grammar symbol x, apply the following rules until no more terminals or e can be added ls any FIRST set.

1. If x is a terminal, then FIRST (x) = $\{x\}$.

2. If x is a non terminal and x $\rightarrow$ $Y_1 Y_2 \cdots Y_k$ is a pdn for some $k \geq 1$, then place a in FIRST (x), if for some i, a is in FIRST ($Y_i$), and e is in all ab FIRST ($Y_1$) $\cdots$ FIRST ($Y_{i-1}$)

   ie, $Y_1 \cdots Y_{i-1} \overset{*}{\Rightarrow} e$

   If e is in FIRST ($Y_j$) for all j = 1, 2 $\cdots$ k, then add e to FIRST (x)

   If $Y_1$ does not derive e, then we add nothing more to FIRST (x), but if $Y_1 \overset{*}{\Rightarrow} \varepsilon$, then we add FIRST ($Y_2$) and so on.

3. If x $\rightarrow$ $\varepsilon$ is a pdn, then add $\varepsilon$ to FIRST (x).

Eg: Consider the following grammar

   $E \rightarrow TE'$

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

$$T \rightarrow FT'$$
$$T' \rightarrow *FT' | \varepsilon$$
$$F \rightarrow (E) | id$$

$$\text{FIRST}(E) = \text{FIRST}(T) = \text{FIRST}(E) = \{ C, id \}$$
$$\text{FIRST}(E') = \{ +, \varepsilon \}$$
$$\text{FIRST}(T') = \{ *, \varepsilon \}$$

- Define <u>FOLLOW (A)</u>, for a nonterminal A, to be the set of terminals a that can appear immediate to the right of A in some sentential forms.

eg: If $S \xrightarrow{*} \propto A a \beta$, Then a is in FOLLOW (A)

- If A is the rightmost symbol in some sentential form, then $ is in FOLLOW(A). ($ is a special end marker symbol that is assumed not to be a symbol of any grammar).

- To compute <u>FOLLOW (A)</u>, for all nonterminals A apply the following rules until nothing can be added to any FOLLOW set.

1. Place $ in 'FOLLOW(S). where s is the start symbol, and $ is the i/p right end marker.

2. If there is a pdn $A \rightarrow \propto B\beta$, then everything

in FIRST ($\beta$), except $\varepsilon$ is in FOLLOW (CB).

3. If there is a pdn A → $\alpha$B, or a pdn A→ $\alpha$B$\beta$ where FIRST ($\beta$) contains $\varepsilon$, then everything in FOLLOW(A) is in FOLLOW (B).

eg: Consider the following grammar

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) | id$


FOLLOW(E) = FOLLOW (E') = { ), $ }.

FOLLOW (T) = FOLLOW (T') = { +, ), $ }

FOLLOW (F) = { *, +, ), $ }.

| Egs:1) | | FIRST | FOLLOW |
|---|---|---|---|
| | S → ABCDE | {a,b,c} | { $ } |
| | A → a/$\varepsilon$ | {a, $\varepsilon$} | {b,c} |
| | B → b/$\varepsilon$ | {b, $\varepsilon$} | {c} |
| | C → c | {c} | {d, e, $} |
| | D → d/$\varepsilon$ | {d, $\varepsilon$} | {e, $} |
| | E → e/$\varepsilon$ | {e, $\varepsilon$} | { $ }. |

// whenever a variable is at the right end of a pdn and if nothing after it then the ~~follow~~ FOLLOW of that variable is FOLLOW of left hand side.

whenever a variable ab[l]y which everything
goes to ε, for ex. in FOLLOW (c), D and E
are going to ε, in this case c is coming to
the right end, so FOLLOW ab right end is equal
to the FOLLOW ab left end //

| | | FIRST | FOLLOW |
|---|---|---|---|
| 2) | $S \rightarrow Bb \mid Cd$ | $\{a, b, c, d\}$ | $\{\$\}$ |
| | $B \rightarrow aB \mid \varepsilon$ | $\{a, \varepsilon\}$ | $\{b\}$ |
| | $C \rightarrow cC \mid \varepsilon$ | $\{c, \varepsilon\}$ | $\{d\}$ |

| | | FIRST | FOLLOW |
|---|---|---|---|
| 3. | $E \rightarrow TE'$ | $\{id, \varepsilon\}$ | $\{\$, )\}$ |
| | $E' \rightarrow +TE' \mid \varepsilon$ | $\{+, \varepsilon\}$ | $\{\$, )\}$ |
| | $T \rightarrow FT'$ | $\{id, \varepsilon\}$ | $\{+, \$, )\}$ |
| | $T' \rightarrow *FT' \mid \varepsilon$ | $\{*, \varepsilon\}$ | $\{+, \$, )\}$ |
| | $F \rightarrow id \mid (E)$ | $\{id, c\}$ | $\{*, +, \$, )\}$ |

// E is the start symbol, so FOLLOW (E)
always contains $. and E is also in the r.h.s
so FOLLOW (E) = $\{\$, )\}$

FOLLOW (E') - E' is at the right end
so FOLLOW (E') = FOLLOW (E).

~~For~~ FOLLOW (T) = FIRST (E')

FIRST(E') contains ε, so replace E' with ε

now T is at the right end so
FOLLOW (T) contains FOLLOW (E)
i.e, FOLLOW(T) = {+, $, )}. //

4).

| | FIRST | FOLLOW |
|---|---|---|
| S → AGB \| CbB \| Ba | {d, g, h, ε, b, a} | {$} |
| A → da \| BC | {d, g, h, ε} | {h, g, $} |
| B → g \| ε | {g, ε} | {$, a, h, g} |
| c → h \| ε | {h, ε} | {g, $, b, h} |

## Construction of a Predictive Parsing Table:

This algm collects information from FIRST and FOLLOW sets into a predictive parsing table $M[A, a]$, a two-dimensional array

- where A is a non-terminal and $a$ is a
  terminal or the symbol \$, the i/p end marker

The algm is based on the following idea.

The production $A \rightarrow \alpha$ is chosen, if the next i/p symbol $a$ is in FIRST($\alpha$). The only complication occurs when $\alpha \Rightarrow \varepsilon$, or more generally $\alpha \overset{*}{\Rightarrow} \varepsilon$

In this case, we should again choose $A \rightarrow \alpha$ if the current i/p symbol is in FOLLOW(A). or if the \$ on the i/p has been reached and \$ is in

FOLLOW (A).

Input : Grammar G

Output: Parsing Table M.

Method: For each production A → α of the
grammar, do the following
1. For each terminal a in FIRST (α), add A → α
to M[A, a]
2. If ε is in FIRST (α), then for each terminal
b in FOLLOW (A), add A → α to M[A, b].
If ε is in FIRST (α) and $ is in FOLLOW (A),
add A → α to M[A, $] as well.

If after performing the above, there is no production
at all in M[A, a], then set M[A, a] to error.
(normally represent by an empty string in the table.

Eg: For the expression grammar

$$E \rightarrow TE'$$
$$E' \rightarrow +TE' | \varepsilon$$
$$T \rightarrow FT'$$
$$T' \rightarrow *FT' | \varepsilon$$
$$F \rightarrow (E) | id.$$

string: id + id * id.

| NON Terminal | id | + | * | ( | ) | $ |
|---|---|---|---|---|---|---|
| E | E → TE' | | | E → TE' | | |
| E' | | E' → +TE' | | | E' → ε | E' → ε |
| T | T → FT' | | | T → FT' | | |
| T' | | T' → ε | T' → *FT' | | T' → ε | T' → ε |
| F | F → id | | | F → (E) | | |

I/P SYMBOL

| | FIRST | FOLLOW |
|---|---|---|
| E → TE' | {id, (} | {$, )} |
| E' → +TE' / ε | {+, ε} | {$, )} |
| T → FT' | {id, (} | {+, ), $} |
| T' → *FT' / ε | {*, ε} | {+, ), $} |
| F → id / (E) | {id, (} | {*, +, ), $} |

$ is included because there is a chance for
There is a $ in the input string.
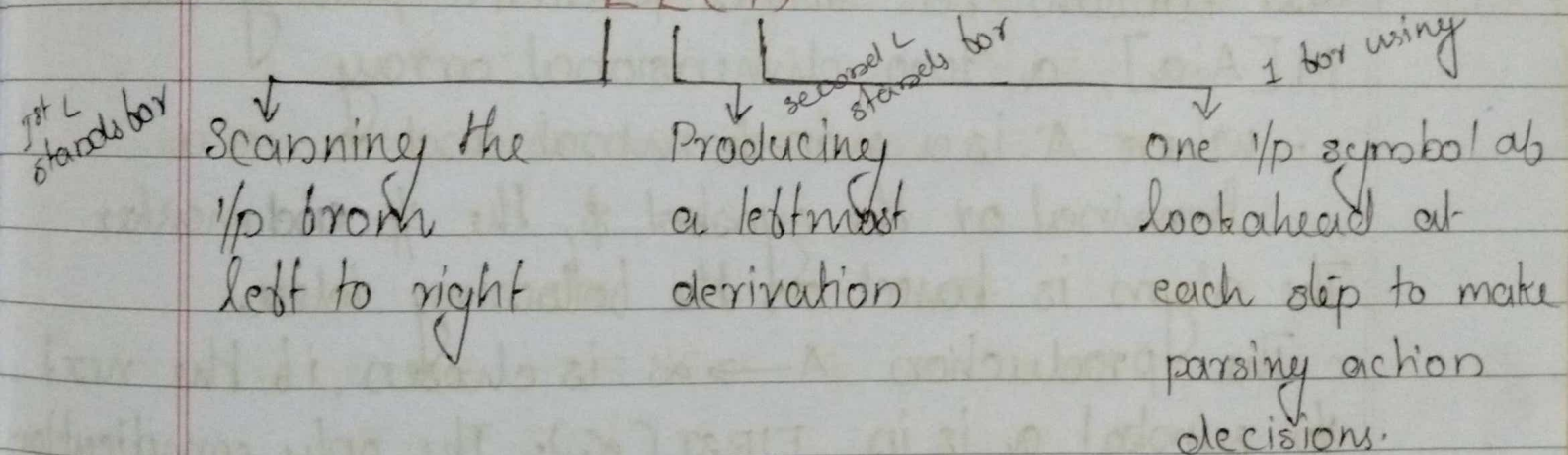E' → ε place the production in the follow of rhs.
Eg: (id E') $
Apply ε production to E' if
the next symbol is )
ie, (id )

## LL(1) Grammar.

Predictive parsers (i.e. recursive descent parsers needing no backtracking) can be constructed for a class of grammars called LL(1).

```
                    LL(1)
        ┌─────┬──────┬─────────────────────────┐
```

1st L stands for | Scanning the i/p from left to right

second L stands for | Producing a leftmost derivation

1 for using | one i/p symbol as lookahead at each step to make parsing action decisions.

A grammar G is LL(1) if and only if whenever A ⟶ $\alpha/\beta$ are two distinct productions of G. the following conditions hold:

1. For no terminal 'a' do both $\alpha$ and $\beta$ derive strings

beginning with a.

2. At most one ab $\alpha$ and $\beta$ can derive the empty string.

3. If $\beta \xRightarrow{*} \varepsilon$, then $\alpha$ does not derive any string beginning with a terminal in FOLLOW(A). Likewise, if $\alpha \xRightarrow{*} \varepsilon$, then $\beta$ does not derive any string beginning with a terminal in FOLLOW(A).

- The first two conditions are equivalent to the statement that FIRST($\alpha$) and FIRST($\beta$) are disjoint sets.
- The third condition is equivalent to stating that if $\varepsilon$ is in FIRST($\beta$), then FIRST($\alpha$) and FOLLOW(A) are disjoint sets.

# 1. Recursive Descent Parsing:

- A recursive descent parsing program consists of a set of procedures, one for each non-terminal.

- Execution begins with the procedure for the start symbol, which halts and announces success if its procedure body scans the entire i/p string.

- General recursive descent may require backtracking (ie, it may require repeated scans

over the I/p).

A typical procedure for a nonterminal in a top-down parser

```
         void A() {
1)          choose an A-pdn, A → x₁x₂ ..... xₖ;
2)          for (i = 1 to k) {
3)              if (xi is a nonterminal)
4)                  call procedure xi();
5)              else if (xi equals the current I/p symbol a)
6)                  advance the I/p to the next symbol;
7)              else /* an error has occurred */;
            }
         }
```

To allow backtracking, the above code need to be modified. It is not possible to choose a unique A-pdn at line 1, so several pdns have to be tried out in some order. Then failure at line 7 is not ultimate failure, but suggest that we need to return to line(1) and try another A-pdn. Only if there are no more A-pdns to try, we say that an error has been found. So inorder to try another A-pdn, I/p pointer has to be reset to the point where it was when it reached line 1. Thus a local variable is needed to store the I/p pointer for further use.

Modify the above pgm code to indicate

include the following:

7) else if A has any alternative pdns
8) choose any one alternative path for A and
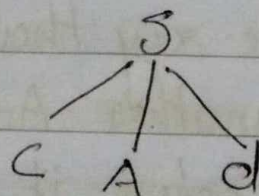            reset i/p pointer.
9) else error();

Eg: Consider the grammar

$$S \rightarrow cAd$$
$$A \rightarrow ab \mid a$$
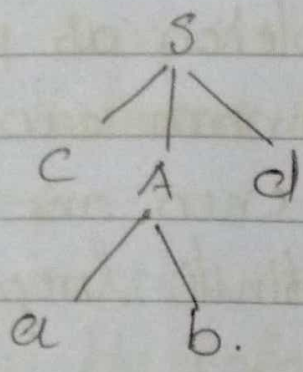
construct a parse tree for the i/p string
$$w = cad.$$

Ans: To construct a parse tree top-down for the i/p
string $w = cad$, begin with a tree consisting of
a single node labeled $S$, and the i/p pointer
pointing to $c$, the first symbol of $w$. $S$ has only
one pdn, so it is used to expand $S$.
Thus the parse tree becomes

```
        S
       /|\
      C A d
```

The leftmost leaf, labeled 'c' matches
the first symbol of input $w$, so i/p pointer is
advanced by 1 position and thus it now point
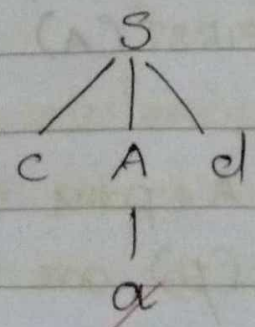
to 'a', the second symbol of w.

- The next leaf is A, when the 1st alternative A → ab is choosen, the tree becomes

$$
\begin{array}{c}
S \\
\diagup \mid \diagdown \\
c \quad A \quad d \\
\diagup \diagdown \\
a \qquad b.
\end{array}
$$

A match is obtained for the second symbol 'a' so the I/p pointer is advanced by 1. and hence, it points to 'd' is compared with the next leaf, labeled 'b'. Since 'b' does not match with 'd', we must go back to A to see whether there is another alternative for A that has not been tried.

In going back to A, we must reset the I/p pointer to position 2. (the position it had, when it first came to A) which means that the procedure for A must store the I/p pointer in a local variable.

The 2nd alternative for A produces the tree

$$
\begin{array}{c}
S \\
\diagup \mid \diagdown \\
c \quad A \quad d \\
\mid \\
a
\end{array}
$$

The leaf 'a' matches 2nd symbol ab w and leaf 'd' matches the 3rd symbol. Since we have produced a parse tree for w, we halt and announce successful completion of parsing.

A left recursive grammar can cause a recursive descent parser even one with back tracking, to go into an infinite loop.