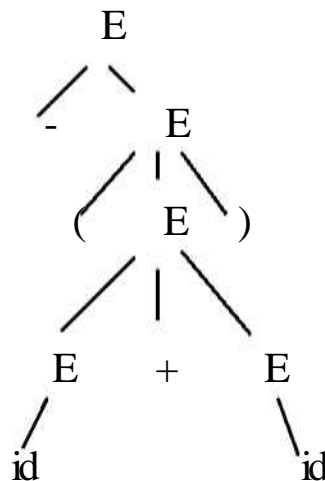\

## Syntactic Analyzer ( Parser )

Every programming language has rules that prescribe the syntactic structure of well formed programs. The syntax of programming language constructs can be described by context free grammars. In our compiler model, the parser obtains a string of tokens from the lexical analyzer, and verifies that the string can be generated by the grammar for the source program. We expect the parser to report any syntax errors in an intelligible fashion. It should also recover from commonly occurring errors so that it can continue processing the remainder of its input.

In syntax analysis we are concerned with groping *tokens* into larger syntactic classes such as *expression* , *statements* , and *procedure*. The syntax analyzer (parser) outputs a *syntax tree*, in which its leaves are the *tokens* and every non-leaf node represents a syntactic *class* type. For example:-

Consider the following grammars:-

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

Then the parse tree for **-(id+id)** is:-



## Syntax Error Handling :-

Often much of the error detection and recovery in a compiler is central around the *parser*. One reason of this is that many errors are syntactic in nature. Errors where the token stream violates the structure of the language are determined by parser, such as an arithmetic expression with unbalanced parentheses.

## Derivations :-

This derivational of view gives a precise description of the *top-down* construction of *parse tree*. The central idea here is that a *production* is treated as rewriting rule in which the *nonterminal* on the left is replaced by the string on the right side of the *production*. For example, consider the following grammar:

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

The *derivation* of the input string **id + id* id** is:

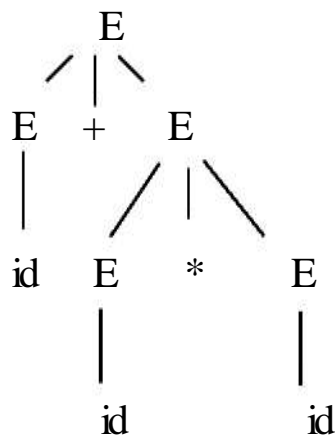| *Left-most derivation* | *Right-most derivation* |
| --- | --- |
| E | E |
| E+E | E+E |
| id +E | E+E*E |
| id+E*E | E+E*id |
| id+id*E | E+id*id |
| id+id*id | id+id*id |

Note:- *parse tree* may by viewed as a graphical representation for a derivation :
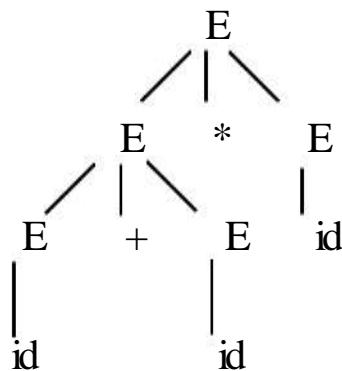
## Ambiguity :-

A grammar that produce more that one parse tree for same sentence is said to be **Ambiguous**. In the another way, by produced more that one *left-most derivation* or more that one *Right-most derivation* for the same sentence.

```
            E
          ╱ │ ╲
        E   +   E
        │     ╱ │ ╲
        id   E   *   E
             │       │
             id      id
```

**(1)**

```
              E
            ╱ │ ╲
          E    *    E
        ╱ │ ╲       │
       E  +   E     id
       │      │
       id     id
```

**(2)**

two parse tree for **id+id\*id**

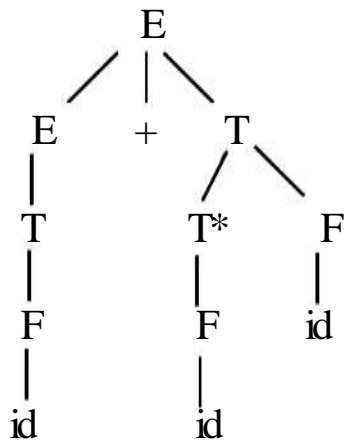| E          |   | E          |
|------------|---|------------|
| E+E        |   | E\*E       |
| id+E       |   | E+E\*E     |
| id+E\*E    |   | id+E\*E    |
| id\*id\*E  |   | id+id\*E   |
| id+id\*i   |   | id+id\*id  |
| d          |   | **(2)**    |
| **(1)**    |   |            |

Two *left-most derivation*s for **id+id\*id**

Sometimes am ambiguous grammar can be rewritten to eliminate the ambiguity. Such as:

| | | |
|---|---|---|
| E → E+E | | E → E+T\|T |
| E → E*E | | E → (E) |
| E → (E) | | E → -E |
| E → -E | | T → T*F\|F |
| E → id | | F → id |

```
              E
            / | \
          E   +   T
          |      / \
          T    T*    F
          |    |     |
          F    F     id
          |    |
          id   id
```

parse tree for **id+id*id**

## Left-Recursion :-

A grammar is left-recursion if has a *nonterminal* **A** such that there is a derivation A ⟶ Aα for some string α . Top-down parser cannot handle *left-recursion* grammars, so a transformation that eliminates left-recursion is needed:

$$A \to A\alpha \mid ß$$

$$A \to ß A'$$

$$A' \to \alpha A' \mid \varepsilon$$

**OR:**

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

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

$A' \rightarrow \alpha_1 A' \mid \alpha_2 A' .. \mid \alpha_m A' \qquad \mid \varepsilon \quad .$

**Example:**

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

---

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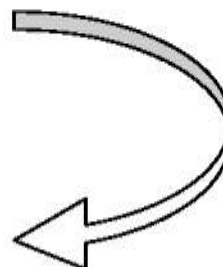$$F \rightarrow (E) \mid id$$

## Left-Factoring :-

The basic idea is that when is not clear which of two alternative production to use to expand a *nonterminal* 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.

$A \rightarrow \alpha \beta_1 \mid \alpha \beta_2$ where $\alpha \neq \varepsilon$

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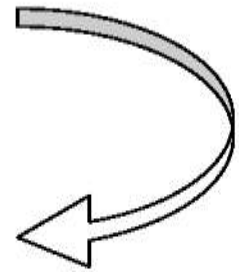
**$A' \rightarrow \beta_1 \mid \beta_2$**

**OR :-**

$A \rightarrow \alpha \beta 1 \quad | \quad \alpha \beta 2 . | \alpha \beta n \quad | \gamma \qquad$ **where** $\alpha \neq \in$

$A \rightarrow \alpha A' \gamma$

$A' \rightarrow \beta 1 \beta 2 .. | \quad | \beta n$

**Example:**

$S \rightarrow iEtS \; iEtSeS \; a \; |$

$E \rightarrow b$

$S \rightarrow iEtSS' a|$

$S' \rightarrow eS \in$

$E \rightarrow b$

*Easy Come, Easy Go*