

# Syntax Analysis(Parser)

A Syntax analyzer is formally defined as :

An Algorithm that Groups the Set of Tokens Sent by the Scanner to Form **Syntax Structures** Such As Expressions, Statements, Blocks,etc.

Simply put, the parser examines if the source code written follows the grammar(production rules) of the language.

The Syntax structure of programming languages and even spoken languages can be expressed in what is called **BNF** notation, which stands for **Bakus Naur Form**.

For example, in spoken English, we can say the following:

```
sentence --> noun-phrase verb-phrase
noun-phrase --> article noun
article --> THE | A | ...
noun --> STUDENT | BOOK | ...
verb-phrase --> verb noun-phrase
verb --> READS | BUYS | ....
```

Note : The BNF Notation uses **different symbols**, for example, a sentence is defined as :

```
< sentence > ::= < noun-phrase > < verb-phrase >
```

But this is very cumbersome, so we use the first notation, since its easier to use.

Now, let us derive a sentence :

```
sentence --> noun-phrase verb-phrase
           --> article noun verb-phrase
           --> THE noun verb-phrase
           --> THE STUDENT verb-phrase
           --> THE STUDENT verb noun-phrase
           --> THE STUDENT READS noun-phrase
           --> THE STUDENT READS article noun
           --> THE STUDENT READS A noun
           --> THE STUDENT READS A BOOK
```

In the same way, the parser tries to **derive** your source program from the starting symbol of the grammar.

Lets say we have these sentences :

```
THE BOOK BUYS A STUDENT
THE BOOK WRITES A DISH
THE DISH TAKES A STROLL
```

Syntax-wise, all of these sentences are correct. However, their meaning is not correct, and they are not useful. What differentiates 2 sentences that are grammatically correct is their meaning or their **semantics**. You and I can agree that the meaning of a grammatically correct sentence is not correct, but how does the computer do it?

# Grammar

A grammar  $G=(V_N, V_T, S, P)$  where:

1.  $V_N$  : A finite set of nonterminals(nonterminals set).
2.  $V_T$  : A finite set of terminals(terminals set).
3.  $S \in V_N$  : The Starting symbol of the grammar.
4.  $P$  = A set of **production rules**(productions).<-- Pending <==> Basically the whole grammar.

Note :

1.  $V_N \cap V_T = \emptyset$ .
2.  $V_N \cup V_T = V$ (the vocabulary of the grammar).

Note : We will use

1. Uppercase Letters A,B,...,Z for non-terminals.
2. Lowercase Letters a,b,...,z for terminals.
3. Greek letters  $\alpha, \beta, \gamma, \dots$  for strings formed from  $V_N$  OR  $V_T = V$ . eg,

if  $V_N = \{S, A, B\}$ ,

$V_T = \{0, 1\}$

then

$\alpha = 0A11B$

$\beta = S110B$

$\gamma = 0010$

## Productions

1. A Production  $\alpha \rightarrow \beta$ (alpha derives beta) is a rewriting rule such that the occurrence of  $\alpha$  can be substituted by  $\beta$  in any string.

Note that  $\alpha$  must contain at least one nonterminal from,  $\in V_N$ .

For example, Assume we have the string  $\gamma\alpha\sigma$ ,

$\gamma\alpha\sigma \rightarrow \gamma\beta\sigma$

2. A Derivation is a sequence of strings  $\alpha_0, \alpha_1, \alpha_2, \alpha_3, \dots, \alpha_n$ , then :

- $\alpha_{i+1} \rightarrow \alpha_i$ ,  $n \geq 0$ .
- $\alpha_{i+1} \xrightarrow{+} \alpha_i$ ,  $n \geq 1$ .

Given a grammar  $G$ , then :

$L(G)$  = Language Generated By the Grammar.

for example, Given the Grammar,  $G = (\{S, B, C\}, \{a, b, c\}, S, P)$

P :

$S \rightarrow aSBC$

$S \rightarrow abC$

$CB \rightarrow BC$

bB → bb

bC → bc

cC → CC

L(G)=?

Lets follow through on the derivations

S → abC → abc(all terminals) ∈ L(G) <--- A sentence

S → aSBC → aabCBC → aabbcBC → blocked, so we try another path

S → aSBC → aabCBC → aabBCC → aabbCC → aabbcC → aabbcc ∈ L(G) <--- A sentence

S → aSBC → ..... → aaabbbccc ∈ L(G) <--- A sentence

Therefore, L(G)={a<sup>n</sup>,b<sup>n</sup>,c<sup>n</sup> | n ≥ 1}

As another Example, we have these productions

E → E+T <-- we can write the productions 1 and 2 as a single production E → E+T | T

E → T

T → T\*F

T → F

F → (E) <-- we can write the productions 5 and 6 as a single production F → (E) | n

F → n

Lets follow through some derivations

E → T → F → n ∈ L(E)

E → E+T → T+T → T+F → T+n → F+n → n+n ∈ L(E)

E → E+T → T+T → F+T → n+T → n+F → n + (E) → n+(T) → n+(T\*F) → n+(F\*F) → n+(n\*F) → n+(n\*n) ∈ L(E)

Therefore, L(G) = {Any arithmetic expression with \* and + operations}, n is an operand here.

Note that, if we add the productions

E → E+T | E-T | T

T → T\*F | T/F | T%F

We would have a language to express all arithmetic expressions with (\*,\,+,-) operations.

Lets Take another Example(things in double quotes are terminals)

Program → block "##"

block → "{" stmt-List "}"

stmt-List → statement "," stmt-List | λ

statement → if-stmt | while-stmt | read-stmt | write-stmt | assignment-stmt | block

if-stmt → "if" condition....

while-stmt → "while" condition....

....

....

read-stmt → "read"

```
write-stmt --> "write"
```

$V_N = \{\text{Program, block, stmt-List, statement, if-stmt, while-stmt, read-stmt, write-stmt, assignment-stmt}\}$

$V_T = \{ \text{"{", "}", "\#", ";", "if", "while", "read", "write"} \}$

Lets Follow through some derivations :

```
Program --> block # --> { stmt-list } # --> { λ } #
```

```
Program --> block # --> {stmt-list} # --> {statement ; stmt-list} # --> {statement ; statement ; stmt-list} # --> {statement ; statement ; λ} # --> {statement ; statement ;} # --> {READ-statement ; statement ;} # --> {READ ; statement ;} # --> {READ ; write-statement ;} # --> {READ ; WRITE ;} #
```

We can write this as

```
{ READ;  
  WRITE;  
}#
```

The language of this language is defined as

$L(G) = \{\text{Set of all programs that can be written in this language}\}.$

This is only a simple example, of a simple language. For something more complex such as C or Pascal, there are hundreds of productions.

## Algorithms for Derivation

A Leftmost derivation is a derivation in which we replace the **leftmost** nonterminal in each derivation step.

A Rightmost derivation is a derivation in which we replace the **rightmost** nonterminal in each derivation step.

For example, given the grammar

```
V --> S R $
```

```
S --> +|-|λ
```

```
R --> .dN | dN.N
```

```
N --> dN | λ
```

$V_N = \{V, R, S, N\}$

$V_T = \{+, -, ., d, \$\}$

Lets follow through on the leftmost derivation

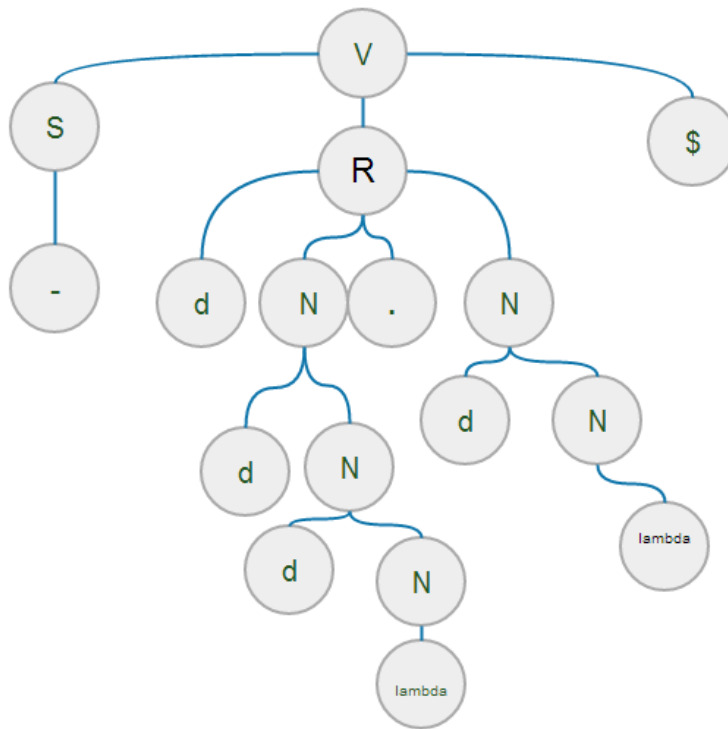
$V \rightarrow SR\$ \rightarrow -R\$ \rightarrow -dN.N\$ \rightarrow -ddN.N\$ \rightarrow -dddN.N\$ \rightarrow -ddd.N\$ \rightarrow -ddd.dN\$ \rightarrow -ddd.d\$ \leftarrow \text{A sentence.}$

Lets follow through on the rightmost derivation

$V \rightarrow SR\$ \rightarrow SdN.N \rightarrow SdN.dN\$ \rightarrow SdN.d\$ \rightarrow sddN.d\$ \rightarrow sdddN.d\$ \rightarrow Sddd.d\$ \rightarrow -ddd.d\$ \leftarrow \text{A sentence.}$

## Derivation Trees

A Derivation Tree is a Tree that displays the derivation of some sentence in the language. For example, lets look at the tree for the previous example



Note that if we traverse the tree in order, recording **only** the leaves, we obtain the sentence.

## Classes of Grammars

According to Chomsky, There are 4 classes of grammars :

1. Unrestricted Grammars : No restrictions whatsoever except the restriction by definition that the left side of the production contains at least one nonterminal from  $V_N$ . This grammar is not practical and we cannot work with it.
2. Context-Sensitive Grammars : For each production  $\alpha \rightarrow \beta$ ,  $|\alpha| \leq |\beta|$ , ie , the **length of alpha( $\alpha$ )** is less than or equal to the **length of Beta( $\beta$ )**. This means that in this class of grammar, there are no  $\lambda$  productions in the form  $A \rightarrow \lambda$ , since  $|\lambda| = 0$  and  $A \geq 1$ .

They Say that Fortran has a context-sensitive grammar.

It is very difficult to work with this class of grammars.

3. Context-Free Grammar(CFG) : Each production in this grammar class is of the form  $A \rightarrow \alpha$  , where  $A \in V_N$  and  $\alpha \in V_T^*$ , that is to say, the left hand side is **only** one nonterminal.

This is the most important class of grammar. Most programming languages's structures are context-free.

We will mostly be working with this class of grammar. Most of the examples we have taken are CFG.

4. Regular Grammar (Regular Expressions) : Each production in this grammar class is of the form  $A \rightarrow aB$  or  $A \rightarrow a$ , where  $A, B \in V_N$  and  $a \in V_T$ , **with the exception** of  $S \rightarrow \lambda$

For example, lets say we have the grammar

$A \rightarrow aA$

$A \rightarrow a$

Therefore, we get

$G(L)=a^+$

However, adding the production

$A \rightarrow \lambda$

Results in the grammar

## Parsing Techniques

There are 2 main parsing techniques used by a compiler.

### Top-Down Parsing

In Top-Down Parsing, the parser builds the derivation tree from the root(S : the starting symbol) down to the leaves(sentence).

In Simple words, the parser tries to derive the sentence using leftmost derivation. For example, say we have this grammar :

```
V --> SR$
S --> + | - | λ
R --> .dN | dN.N
N --> dN | λ
```

Lets examine if the sentence

```
dd.d$
```

is derived from this grammar.

```
V --> SR$ --> +R$ --> dN.N$ --> ddN.N$ --> dd.N$ --> dd.dN$ --> dd.d$
```

Therefore, this sentence is derived from the grammar.

However, this approach is very computationally intensive, and more importantly, this requires knowing the source code in advance. The Parser doesnt know which production it should select in each derivation statement. We will learn how to solve these issues later in the course.

### Bottom-Up Parsing

In Bottom-Up Parsing, the parser builds the derivation tree from the leaves(sentence) up to the root(S : Starting Symbol). This type of tree, built from the leaves to the root, is called a [B-Tree](#).

In Simple words, the parser starts with the given sentence, does **reduction**(opposite of derivation) steps, until the starting symbol is reached.

Note that the string  $\lambda$  is present everywhere in the string, and we can use it wherever we like.

Lets follow the reduction of the example given above.

```
+dd.d$ --> +ddλ.d$ --> +ddN.d$ --> +dN.d$ --> +dN.dλ$ --> +dN.dN$ --> +dN.N$ --> +R$ --> SR$ --> V
```

Which means that the sentence is in the grammar.

Note that we can run into deadlocks here. say we took this path instead :

```
+dd.d$ --> +ddλ.d$ --> +ddN.d$ --> +dN.d$ --> +dN.dλ$ --> +dN.dN$ --> +dNR$ --> +NR$ --> SNR$ --> Deadlock
```

This technique also has a major problem : Which substring should we select to reduce in each reduction step?

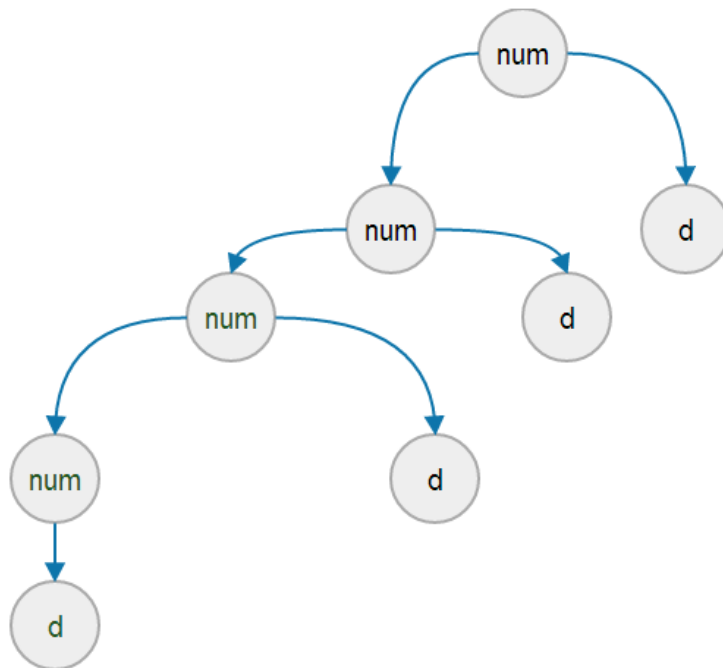
how do we solve this?

## Ambiguity

Given the following grammar :

```
num --> num d
num --> d
```

Let us draw the derivation tree for the sentence dddd



Question : is there another derivation tree that represents the sentence?

The answer is **no**.

If there is only one derivation tree representing the sentence, this means there is only one way to derive the sentence.

Based on this, we can say that :

A Grammar G is said to be ambiguous if there is one sentence with more than one derivation tree.

That is, there is more than one way to derive the sentence.

This means that our algorithm is **non-deterministic**.

Say we have this grammar

$E \rightarrow E + E$

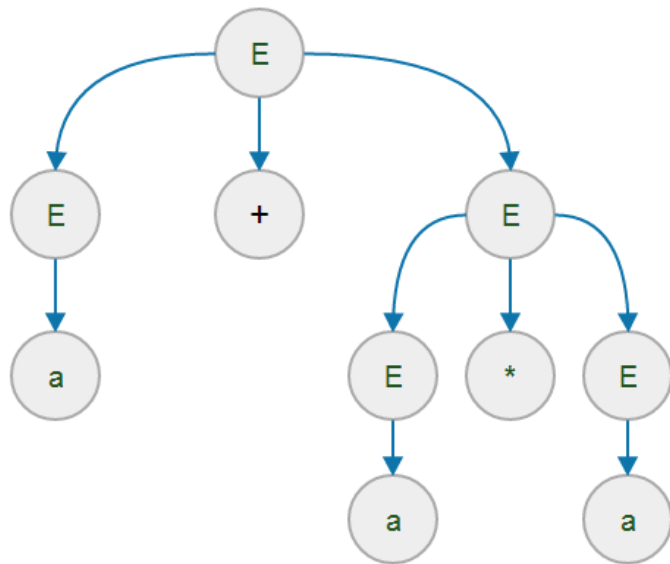
$E \rightarrow E * E$

$E \rightarrow (E) \mid a$

Take the sentence :

$a + a * a$

Lets draw the derivation tree



Due to the fact that we have 2 trees that give the same result, we can say that this grammar is ambiguous.

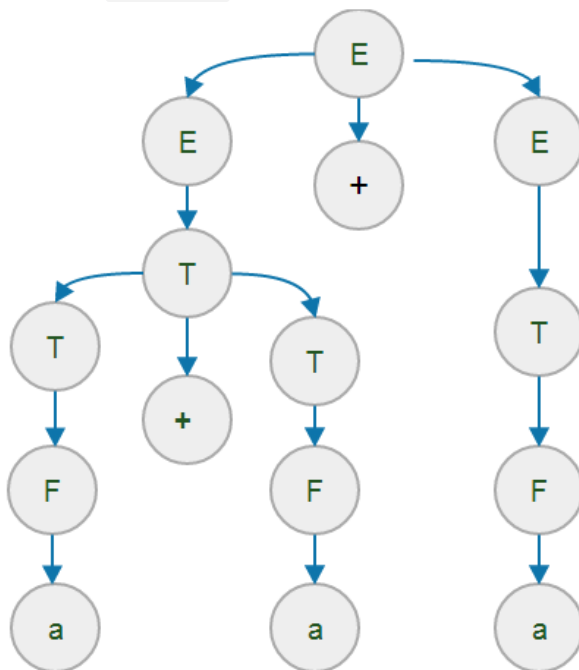
In this case, to enforce the associativity rule, this grammar can be re-written as :

$E \rightarrow E + E \mid T$

$T \rightarrow T * T \mid F$

$F \rightarrow (E) \mid a$

Now, Take the sentence `a + a * a` and find the derivation tree now.



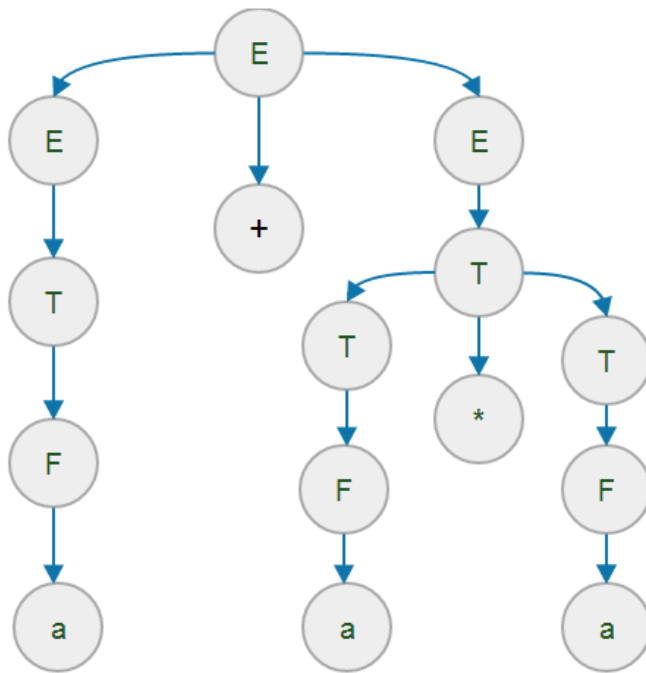
There is only 1 possible derivation tree now. This solves the associativity issue of the grammar before with the `+` and `*` operations.

But lets say we have the sentence :

`a + a + a`

Lets try to find the derivation tree and any alternative trees.





We can see here that there is more than 1 derivation tree, and the language is still ambiguous.

We can solve this if we rewrite the grammar with the **left-associative rule**

$E \rightarrow E + T \mid T$

$T \rightarrow T * F \mid F$

$F \rightarrow (E) \mid a$

The resultant grammar is left-associative.

This grammar solves the problems of :

- ambiguity.
- precedence.
- associativity.

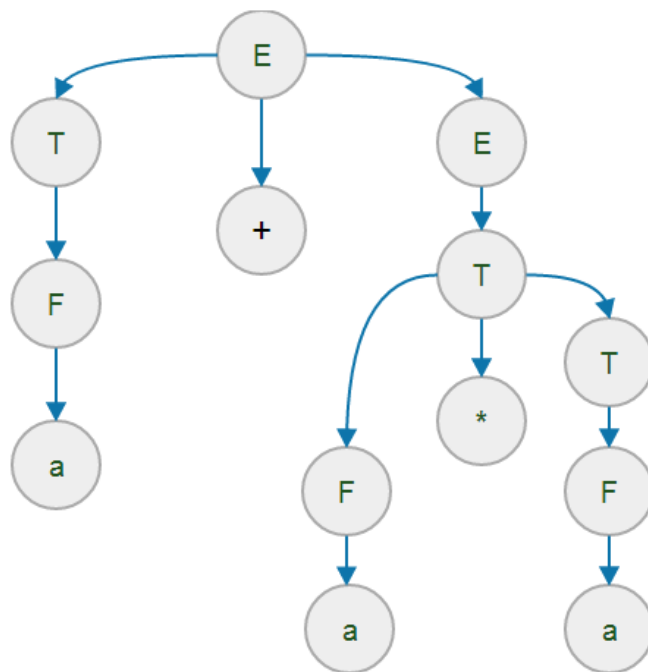
Lets try rewriting it with the **right-associative rule**

$E \rightarrow T + E \mid T$

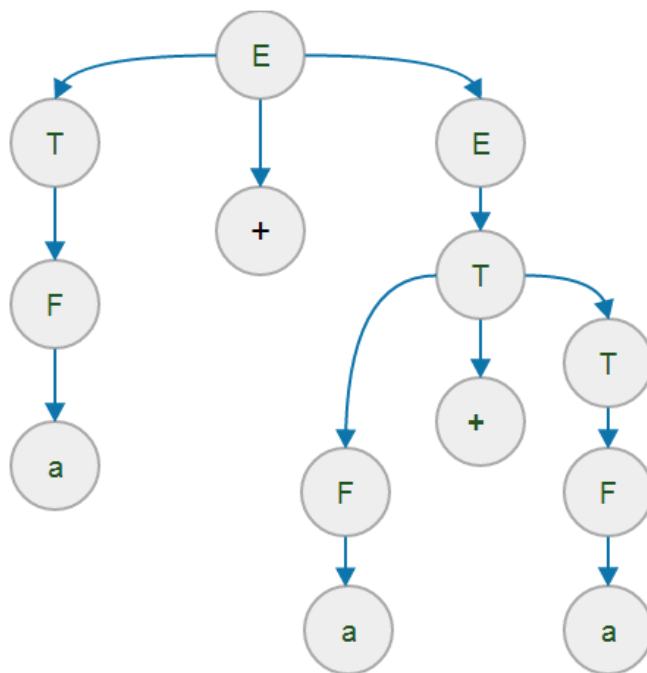
$T \rightarrow F * T \mid F$

$F \rightarrow (E) \mid a$

Lets try creating the derivation tree of `a + a * a`



Now lets draw the derivation tree of `a + a + a`



This new grammar is not ambiguous, however, as we can tell from the derivation trees, there are precedence issues now. It's not technically wrong, but it doesn't follow standard arithmetic rules.

Back to the left-associative grammar now. This grammar is called **left-recursive**. This causes problems when it comes to top-down parsing techniques (we will see why later).

A grammar is said to be left recursive if there is a production of the form:

`A --> Aα`

Conversely, a grammar is right-recursive if there is a production of the form:

`A --> αA`

And causes no problems in top-down parsing.

our grammar has 2 rules of the form

`A --> Aα`

The solution is to transform the grammar to a grammar which is not left-recursive.

This has an algorithm to it.

Given that

$$A \rightarrow A\alpha_1 | A\alpha_2 | A\alpha_3 | \dots | A\alpha_n$$
$$A \rightarrow B_1 | B_2 | B_3 | \dots | B_n$$

To do this, we must introduce a new non-terminal, say  $A'$ .

The grammar now becomes :

$$A \rightarrow B_1 A' | B_2 A' | B_3 A' | \dots | B_n A'$$

and

$$A' \rightarrow \alpha_1 A' | \alpha_2 A' | \alpha_3 A' | \dots | \alpha_n A' | \lambda$$

For example, say we have

$$A \rightarrow Ab$$
$$A \rightarrow a$$
$$L(G) = ab^*$$

Then according to the above

$$A \rightarrow aA'$$
$$A' \rightarrow bA' | \lambda$$

which results in the same grammar.

Lets apply this to the grammar :

$$E \rightarrow E + T | T$$
$$T \rightarrow T * F | F$$
$$F \rightarrow (E) | a$$

This results in :

$$E \rightarrow T E'$$
$$E' \rightarrow + T E' | \lambda$$
$$T \rightarrow F T'$$
$$T' \rightarrow * F T' | \lambda$$
$$F \rightarrow (E) | a$$

This grammar is now **perfect**. It solves all our ambiguity issues, and this is a grammar we can use to construct the production rules for our programming language.

Another ambiguity in programming languages is the `if...else` statement.

Lets take a generic if statement in a generic language:

$$\text{stmt} \rightarrow \text{if-stmt} | \text{while-stmt} | \dots$$
$$\text{if-stmt} \rightarrow \text{IF condition stmt}$$
$$\text{if-stmt} \rightarrow \text{IF condition stmt ELSE stmt}$$
$$\text{condition} \rightarrow C$$

```
stmt --> S
```

This grammar is ambiguous.

Lets take the nested `if...else` statement :

```
IF C
  IF C
    S
  ELSE
    S
```

This statement results in 2 derivations trees.

Both these trees result in the same traversal, but they have different meanings. The first results in the `ELSE` belonging to the first `IF` , while the second results in the `ELSE` belonging to the second `IF` . We as humans know that the `ELSE` belongs to the second `IF` , since we know that the `ELSE` statement follows the nearest `IF` . but how can the compiler know?

There are a bunch of solutions to this problem:

1. Add a delimiter to the `IF` statement, such as `ENDIF` or `END` or `FI` to the end of the statement, resulting in these productions :

```
if-stmt --> IF condition stmt ENDIF
```

```
if-stmt --> IF condition stmt ELSE stmt ENDIF
```

Resulting in this statement :

```
...
IF C
|  IF C
|  |  S
|  |ELSE
|  |  S
|  ENDIF
ENDIF
...
```

The grammar is now unambiguous, since we have to clearly state when an `IF` statement ends. However, this is not a pretty solution, and is extra work for both the programmer and compiler, and results in less readable code.

2. In C and Pascal, the compiler **always** prefers to shift the `ELSE` when it sees it in the source code so it follows the nearest `IF` . We will learn about this in more detail later.

Another thing about this grammar is **left factoring**.

## Left Factoring

Consider the productions :

```
A --> αβ
```

```
A --> αγ
```

Note how the first part of the productions is the same. This grammar can be transformed by introducing a new non-terminal, So what happens now is:

```
A --> αB
```

```
B --> βγ
```

For our grammar, this results in

```
if-stmt --> IF condition stmt
```

```
if-stmt --> IF condition stmt ELSE stmt
```

becoming

```
if-stmt --> IF conditon stmt else-part
```

```
else-part --> ELSE stmt |  $\lambda$ 
```

Does this solve the ambiguity? No, but it helps in removing choices, since the if-stmt is now one production. If we look at the statement :

```
IF C
  IF C
    S
  ELSE
    S
```

It still has 2 derivation trees

## More Ways of Expressing Programming Languages

### Extended BNF Notation

So far, we have been using **BNF Notation**(Production rules) to express languages. However, there is another form to Express a language, which is **Extended BNF Notation**

if there is repetition in the grammar, say in the example of the grammar

```
E --> E + T | T
```

```
T --> T * F | F
```

```
F --> (E) | a
```

which can give us a derivation in the form of

```
E --> E + T --> E + T + T --> E + T + T + T --> T + T + T + T....+T
```

or in the same line,

```
T --> T F --> T F F --> T F F F --> T F F F... F
```

We can express this grammar as :

```
E --> T { + T }
```

```
T --> F { * F }
```

```
F --> (E) | a
```

We know that [x] means that we take x 0 or 1 time only.

However, { x } means we take x zero or any number of times. This is equivalent to (x)\*

We can also express this grammar as:

```
E --> T ( + T )*
```

```
T --> F ( * F )*
```

```
F --> (E) | a
```

### Syntax Diagrams

Another way to express languages are **Syntax Diagrams**. These are used only with Extended-BNF notation.

A square shape represents a nonterminal and an oval shape represents a terminal.

## Parsing Techniques (Continued)

Recall : The parser is an algorithm which accepts or rejects a sentence in the programming language.

Recall : There are 2 kinds of parsers :

1. Top-Down Parsers : In This parsing technique, The parser starts with  $S$  using leftmost derivation to derive the sentence. The Major problem with this parsing technique is that the parser doesn't know which production it should select in each derivation step.
2. Bottom-Up Parsers : The parser in this parsing technique starts from the sentence, doing reduction steps, until it reaches the starting symbol  $S$  of the grammar. The Major problem with this technique is that the parser doesn't know which substring the parser should select in each reduction step.

In Top-Down parsing, we have 2 available algorithms for parsing :

1. Recursive Descent Parsing.
2. LL(1) Predictive Parsing.

In Bottom-Up parsing, we have 2 available algorithms for parsing :

1. LR Parsers.
2. Operator Precedence Parsers --> Uses matrix manipulation.

Before we continue, we need to define a few functions

### The FIRST() Function

Given a string  $\alpha \in V^*$ , then

$$\text{FIRST}(\alpha) = \{ a \mid \alpha \xrightarrow{*} aw, a \in V_T, w \in V^* \}$$

in addition, if  $\alpha \rightarrow \lambda$ , then we add  $\lambda$  to  $\text{FIRST}(\alpha)$ , that is

$$\lambda \in \text{FIRST}(\alpha)$$

That is to say,  $\text{FIRST}(\alpha)$  = Set of all terminals that may begin strings derived from  $\alpha$ .

For example

$$\alpha \xrightarrow{*} cBx$$

$$\alpha \xrightarrow{*} ayD$$

$$\alpha \xrightarrow{*} ab$$

$$\alpha \rightarrow ddd$$

Then

$$\text{FIRST}(\alpha) = \{c, a, d\}$$

Assume as well that

$$\alpha \rightarrow \lambda$$

then

$$\text{FIRST}(\alpha) = \{c, a, d, \lambda\}$$

That is to say,  $\lambda$  appears in the **FIRST()** function.

Lets take an example of this.

### The FOLLOW() Function

We define the FOLLOW() function for **only** nonterminals. That is to say

$\text{FOLLOW}(A), A \in V_N$

Given

$S \xrightarrow{*} uA\beta, u \in V_T^*, A \in V_N, \beta \in V^*$

then

$\text{FOLLOW}(A) = \text{FIRST}(\beta)$

That is to say, FOLLOW(A) = The set of all terminals that may appear after A in the derivation.

$S \xrightarrow{*} aaXdd$

$S \xrightarrow{*} Xa$

$S \xrightarrow{*} BXc$

Then

$\text{FOLLOW}(X) = \{d, a, c\}$

## Rules To Compute FIRST() and FOLLOW() Sets

1.  $\text{FIRST}(\lambda) = \{\lambda\}$ .
2.  $\text{FIRST}(a) = \{a\}$ .
3.  $\text{FIRST}(a\alpha) = \{a\}$ .
4.  $\text{FIRST}(XY) = \text{FIRST}(X) \text{ if } \lambda \notin \text{FIRST}(X) \text{ OR } \text{FIRST}(X) = \lambda \text{ AND } \text{FIRST}(Y)$ .
5. Given the production  $A \rightarrow \alpha X \beta$ , Then : a.  $\text{FIRST}(\beta) \subset \text{FOLLOW}(X)$  if  $\beta \neq \lambda$ . b.  $\text{FOLLOW}(A) \subset \text{FOLLOW}(X)$  if  $\beta = \lambda$ .

Note that the FIRST() and FOLLOW() sets are made of **terminals only**

By these rules, say we have

$A \rightarrow \alpha XB, X \in V_N$

and say we want FOLLOW(X)

Then

$\text{FIRST}(B) \subset \text{FOLLOW}(X)$

We say it is a subset because we can have other productions involving X.

Assuming that  $B = \lambda$ , Things are different.

Say that we have a production that leads to this derivation is

$S \xrightarrow{*} uAy$

and following through this results in this derivation :

$S \xrightarrow{*} uAy \rightarrow u\alpha Xy$

Therefore,

$\text{FOLLOW}(A) \subset \text{FOLLOW}(X)$

This is because whatever follows A can follow X if there is nothing between them.

Notes :

1.  $\lambda$  **may** appear in FIRST() but it doesn't appear in FOLLOW(). We will see this when we define augmented grammars.
2. Generally, we start computing the FIRST() from bottom to top, But follow from top to bottom.

3. When we compute FOLLOW(X), we search for X in the right side of any production.

### Augmented Grammars

Given the grammar  $G=(V_N, V_T, S, P)$ , then the augmented grammar  $G'=(V_N', V_T', S', P')$  can be obtained from G as follows:

1.  $V_N' = V_N \cup \{S'\}$ .
2.  $V_T' = V_T \cup \{ \$ \}$ .
3.  $S'$  = new starting point.
4.  $P = P \cup \{S' \rightarrow SS'\}$

For example :

```
E --> E + T | T
T --> T * F | F
F --> (E) | a
```

Becomes :

```
G --> E$
E --> E + T | T
T --> T * F | F
F --> (E) | a
```

This is because we want to create a FOLLOW() set for S.

Lets take another example :

```
S' --> S$
S --> AB
A --> a | λ
B --> b | λ
```

Lets compute the FIRST() sets for this grammar :

```
FIRST(A) = {a, λ}
FIRST(B) = {b, λ}
FIRST(S) = FIRST(AB) = FIRST(FIRST(A).FIRST(B))
= FIRST({a, λ}, {b, λ})
= FIRST({a, λ, b, λ})
= {a, b, λ}
FIRST(S') = FIRST(S$) = FIRST(FIRST(S).FIRST($))
= FIRST({a, b, λ}. $) = FIRST(a$, b$, $)
= {a, b, $}
```

Now Lets compute the FOLLOW() sets for this grammar :

```
FOLLOW(S) = {$}
FOLLOW(A) = {b, $}
FOLLOW(B) = {$}
```



Lets Take another, slightly more complex example :

```
S' --> S$  
S --> aAc b  
S --> A b c  
A --> b | c |  $\lambda$ 
```

Lets take the FIRST() for this grammar :

```
FIRST(A) = {b,c, $\lambda$ }  
FIRST(S) = FIRST(aAc b)  $\cup$  FIRST(A b c) = {a,}  $\cup$  {b,c}  
= {a,b,c}  
FIRST(S') = FIRST(S$) = FIRST(FIRST(S).FIRST($))  
= FIRST({a,b,c}.{$})  
= {a,b,c}
```

Now lets take the FOLLOW() :

```
FOLLOW(S) = {$}  
FOLLOW(A) = {c,b}
```

---

Say we have the grammar

```
G --> E$  
E --> E + T | T  
T --> T * F | F  
F --> (E) | a
```

Lets calculate FIRST() :

```
FIRST(F) = {(,a}  
FIRST(T) = FIRST(T F)  $\cup$  FIRST(F) = FIRST(T F)  $\cup$  {(,a}  
= {(,a} (Because every T will eventually become an F)  
FIRST(E) = FIRST(E + T)  $\cup$  FIRST(T) = {(,a}  $\cup$  {(,a}  
= {(,a}  
FIRST(G) = FIRST(E$) = {(,a}
```

Now lets Calculate FOLLOW() :

```
FOLLOW(E) = {$,+,)}  
FOLLOW(T) = FOLLOW(E)  $\cup$  {*} = {$,+,*,)}  
FOLLOW(F) = FOLLOW(T) = {$,+,*,)}
```

---

But what makes all this so important?

Well, All of the parsing techniques we are going to learn will heavily rely on FIRST() and FOLLOW().

---

## Top-Down Parsing (Continued)

### Recursive Descent Parsing

Recursive Descent Parsing is very simple. It works like this :

1. Divide the grammar into primitive/simple components

- i. For the token "a" :

```
if(token == "a"){  
    get-next()  
}  
else{  
    report-error()  
}
```

- ii. For  $X = \alpha_1, \alpha_2, \dots, \alpha_n$  :

```
>code(X):\{  
>  
>code(&alpha;<sub>1</sub>);  
>  
>code{&alpha;<sub>2</sub>};  
>...  
>  
>...  
>  
>...  
>  
>code(&alpha;<sub>n</sub>);  
>\}
```

Then

```
> If(token &isin; FIRST(&alpha;<sub>1</sub>))  
>  
>   code(&alpha;<sub>1</sub>)  
>  
> else If(token &isin; FIRST(&alpha;<sub>2</sub>))  
>  
>   code(&alpha;<sub>2</sub>)  
>...  
>  
>...  
>  
>...  
>If(token &isin; FIRST(&alpha;<sub>n</sub>))  
>  
>   code(&alpha;<sub>n</sub>)  
>  
>else{  
>  
> report-error();  
>}
```

If  $\alpha$ ;  $\neq$ ;  $\lambda$ ;

So we say, the code for  $x = \alpha^*$ , we say

```
> while(token &isin; FIRST(&alpha;^*))\{  
>  
> call(&alpha;^*);
```

```
>  
> /}
```

Notes :

1. Every nonterminal has a code(a function).
2. S` in augmented grammar is represented by the function "main".
3. We only start with calling "get-token" in function "main".

For Example, lets say we have

```
G --> E$  
E --> T ( + T ) *  
T --> F ( * F ) *  
F --> ( E ) | a
```

Then we can say :

```
main(){//represents G  
  
get-token;  
  
call E();  
  
if(token!="$")  
{  
    Error;  
}  
else{  
    SUCCESS;  
}  
  
}  
  
function E(){//Represents E -- T ( + T ) *  
  
    call T();  
  
    while(token == "+"){  
  
        get-token();  
  
        call T()  
    }  
  
}  
  
function T(){//T--> F ( * F ) *  
  
    call F();  
    while(token == "*"){  
  
        get-token();  
  
        call F();  
    }  
}  
  
function F(){  
  
    if(token == "(")  
    {  
        get-token();
```

```

    call E();

    if(token == "")
    {
        get-token();
    }
    else
    {
        ERROR;
    }
}
else if(token=="a")
{
    get-token();
}
else
{
    ERROR;
}
}

```

Note that `ERROR` is a function we should write.

---

Lets take another example now.

Given the grammar :

Program  $\rightarrow$  body .

body  $\rightarrow$  Begin stmt ( ; stmt)\* End

stmt  $\rightarrow$  Read | Write | body |  $\lambda$

where we will represent  $\lambda$  as `1` from now on in the example;

and

$V_N = \{ \text{Program, body, stmt, block} \}$

$V_T = \{ ., \text{Begin, }, \text{End, Read, Write} \}$

examples of programs of this language would be

```

Begin
    Read;
    Write;
    Read;
    Write;
End.

```

or

```

Begin
    Read;
End.

```

or

```

Begin
    Read;
    Begin
        Read;
        Write;
    End.
End.

```

```
    Write;  
End.
```

or

```
Begin;  
;  
;  
;  
;  
End.
```

Lets write the code for this programming language.

```
main(){  
    get-token();  
  
    call body();  
  
    if(token != "."){  
        ERROR;  
    }  
    else{  
        SUCCESS;  
    }  
}  
  
function body(){  
    call Begin();  
    if(token == "Begin"){  
        get-token();  
        call stmt();  
        while(token != ";"){  
            get-token();  
            call stmt();  
        }  
        if(token == "End"){  
            get-token();  
        }  
        else{  
            ERROR;  
        }  
    }  
    else{  
        ERROR;  
    }  
}  
  
function stmt(){  
    if(token == "Read"){  
        get-token();  
    }  
    else if(token == "Write"){  
        get-token();  
    }  
    else if(token == "Begin"){  
        call body();  
    }  
    else if(token != ";" || token != "End" ){  
        ERROR();  
    }  
}
```

This Parsing method is a **table-driven** parsing method. The LL(1) parsing table selects which production to choose for the next derivation step.

### Formal Definition of LL(1)

The Formal definition of LL(1) grammars is given by :

Given the Productions :

$A \rightarrow \alpha_1$

$A \rightarrow \alpha_2$

$A \rightarrow \alpha_3$

$A \rightarrow \alpha_n$

then the grammar is LL(1) if :

1.  $FIRST(\alpha_i) \cap FIRST(\alpha_j) = \emptyset$  for all  $i, j$
2. if one of  $\alpha_i$  is  $\lambda$ ,  $\alpha_n = \lambda$ , in addition to 1,  $\forall i < n$

For example, Given the grammar :

$S \rightarrow S\$$

$S \rightarrow aABC$

$A \rightarrow a \mid bbD$

$B \rightarrow a \mid \lambda$

$C \rightarrow b \mid \lambda$

$D \rightarrow C \mid \lambda$

let us see if it is LL(1)

- $FIRST(a) \cap FIRST(bbD) = \emptyset$
- $FIRST(a) \cap FOLLOW(B) = \emptyset$
- $FIRST(b) \cap FOLLOW(C) = \emptyset$
- $FIRST(c) \cap FOLLOW(D) = \emptyset$

Then this grammar is LL(1).

Given another Grammar :

$S' \rightarrow S\$$

$S \rightarrow aAa \mid \lambda$

$A \rightarrow abS \mid \lambda$

- $FIRST(aAa) \cap FOLLOW(S) = \{a\} \cap \{\$, a\} = \{a\} \neq \emptyset$

This grammar is **not** LL(1).

---

Lets assume that we have a grammar that is LL(1). How do we build the LL(1) parsing table?

### LL(1) Parsing Table Building Algorithm

1. For each production  $A \rightarrow \alpha$  in the grammar G,
  - i. Add to the table entry  $T[A, a]$  the production  $A \rightarrow \alpha$ , where  $A \in FIRST(\alpha)$
2. If  $\lambda \in FIRST(\alpha)$ , Add to the table entry  $T[A, b]$  the production  $A \rightarrow \alpha \forall b \in FOLLOW(A)$ .

3. All Remaining Entries are Error Entries.

For example, given the grammar :

$V \rightarrow SR \$^1$

$S \rightarrow +^2 \mid -^3 \mid \lambda^4$

$R \rightarrow dN.N^5 \mid .dN^6$

$N \rightarrow dN^7 \mid \lambda^8$

note that the superscript denotes the production number.

Then the table will look like this

$FIRST(SR \$) = \{+, -, d, .\}$

$FIRST(+)=\{+\}$

$FOLLOW(S)=\{d, .\}$

$FIRST(R)=\{d, .\}$

$FIRST(d)=\{d\}$

$FOLLOW(N)=\{d, ., \$\}$

$V_N \backslash V_T$	+	-	d	.	\$
V	1	1	1	1	
S	2	3	4	4	
R			5	5	
N			7	8	8

There should be no conflict(multiple entries) in the LL(1) table.

$L(G)$  of this grammar = all floating point numbers.

The parser works like this

Stack	Remaining Input	Action
V	-dd.d\$	Production 1
SR\$	-dd.d\$	Production 3
-R\$	-dd.d\$	Pop & advance input
R\$	dd.d\$	Production 5
dN.N\$	dd.d\$	Pop & advance input
N.N\$	d.d\$	Production 7
dN.N\$	d.d\$	Pop & advance input
N.N\$	.d\$	Production 8
.N\$	.d\$	Pop & advance input
N\$	d\$	Production 7
dN\$	d\$	Pop & advance
N\$	\$	Production 8
\$	\$	Pop and Advance
$\lambda$	$\lambda$	Accept

If at any point the parser reaches a place where the input and the stack have 2 different terminal symbols, it throws a syntax error.

Lets Take another example. Let the Grammar be :

program  $\rightarrow$  block \$ <sup>1</sup>

block  $\rightarrow$  { declarations stmts } <sup>2</sup>

decls  $\rightarrow$  D ; decls <sup>3</sup> |  $\lambda$  <sup>4</sup>

stmts  $\rightarrow$  statement ; stmts <sup>5</sup> |  $\lambda$  <sup>6</sup>

statement  $\rightarrow$  if <sup>7</sup> | while <sup>8</sup> | ass <sup>9</sup> | scan <sup>10</sup> | print<sup>11</sup> | block <sup>12</sup> |  $\lambda$  <sup>13</sup>

$V_T = \{\$, \{, \}, D, :, \text{if}, \text{while}, \text{ass}, \text{scan}, \text{print}\}$

$V_N \setminus V_T$	if	while	ass	scan	print	{	}	D	;	\$
Program						1				
block						2				
decls	4	4	4	4	4	4	4	3	4	
stmts	5	5	5	5	5	5	6		5	
statement	7	8	9	10	11	2			13	

No conflict.

Another example is the If....else statement with a delimiter. the grammar looks like this :

$S' \rightarrow S\$$

$S \rightarrow iCSE$

$E \rightarrow S \mid \lambda$

$S \rightarrow a$

$C \rightarrow c$

$V_N \setminus V_T$	i	a	e	c	\$
$S'$	1	1			
S	2	5			
E			3,4		4
C				6	

There is a conflict. To solve this, we can add a delimiter.

$S' \rightarrow S\$$

$S \rightarrow iCSEd$

$E \rightarrow eS \mid \lambda$

$S \rightarrow a$

$C \rightarrow c$

$V_N \setminus V_T$	i	a	e	c	d	\$
$S'$	1	1				
S	2	5				
E			3		4	
C				6		



The grammar is now unambiguous. Alternatively, we can just strike out the transition 4 from the LL(1) table, which is a  $\lambda$  transition. Lets follow through the derivation tree. the resultant table is

$V_N \setminus V_T$	i	a	e	c	\$
S	1	1			
S	2	5			
E			3		4
C				6	

and this works because the natural behavior of the else-part is to follow the nearest if statement.

---

Something Important to note is that **if a grammar is LL(1), then it is unambiguous. However, the opposite is not necessarily true.**

Another thing to note is that in Top-Down parsing, we should avoid a grammar that is not LL(1).

---

## Bottom-Up Parsing (Continued)

Recall that in Bottom-Up parsing, the parser starts from the given sentence, applying reductions until it reaches the starting symbol of the grammar or a deadlock. The major problem with Bottom-Up parsing is which substring we should select in each reduction step.

The answer to the above question is : In each reduction step, we select what is called **the handle**.

The Handle is obtained by a **rightmost derivation in reverse**

For example, Given the grammar :

```
V --> S R $
S --> +|-|λ
R --> .dN | dN.N
N --> dN | λ
```

and the sentence

```
-dd.d$
```

First, we derive the sentence rightmost.

```
V --rm--> SR$ --rm--> SdN.N$ --rm--> SdN.dN$ --rm--> SdN.dd$ --rm--> SddN.d$ --rm--> Sdd.d$ --rm--> -dd.d$
```

So our handles would be :

```
V <-- SR$ <-- SdN.N$ <-- SdN.dN$ <-- SdN.dλ$ <-- SddN.d$ <-- Sddλ.d$ <-- -dd.d$
```

But Compilers dont work like this. We already derived the sentence, why would we go back and do it again?

We cannot build a Bottom-Up parser for every Context-Free Grammar. However, we are fortunate enough that there exist subsets of the Context-Free Grammar for which we can build a deterministic Bottom-Up parser ie, the parser can determine/decide precisely where the handle is in each reduction step.

some of these subsets are

- LR Parsers :
  - SLR(Simple-LR).
  - LALR(Look-Ahead LR).
  - LR.
- Operator Precedence.

We will only be talking about the SLR parser, just to get an idea of how Bottom-Up parsing works.

## SLR Parsing

SLR parsing, and LR parsing in general, is a tabular parsing method.

All LL(1) grammars are a subset of SLR grammars.

All LR parsers contains :

1. A parsing table.
2. A stack.
3. The input string.

As a reminder, the LL(1) parser contains :

1. A parsing table.
2. A stack.
3. The input string.

However, the way we build it is different.

### Building the SLR Parsing Table

An LR(0) item of a grammar G is a production in G with a dot(.) at some position in the right side.

For example, the production

$A \rightarrow aBY$

This production generates the following LR(0) items :

$A \rightarrow .aBY$

$A \rightarrow a.BY$

$A \rightarrow aB.Y$

$A \rightarrow aBY. \rightarrow \text{complete item}$

Note that for  $A \rightarrow \lambda$ , this generates only  $A \rightarrow \lambda.$

Generally speaking, if the right side of the production is L, then there are L+1 resultant LR(0) items.

The LR(0) item

$A \rightarrow aB.Y$

Means that the parser has scanned on the input a string derived from aB and expects to see a string derived from Y .

We need to define the following 2 functions.

### The CLOSURE function

```
function CLOSURE(I)//I is a set of LR(0)items
{
    Repeat
        for(every LR(0) item A-->alpha.Bbeta in I,
            and for every production B-->gamma in G,
            Add the LR(0)item B-->.gamma to I)
        Until no more items to be added;
}
```

Lets apply this to our grammar :

- (1)  $E \rightarrow E+T$
- (2)  $E \rightarrow T$
- (3)  $T \rightarrow T^*F$
- (4)  $T \rightarrow F$
- (5)  $F \rightarrow (E)$
- (6)  $F \rightarrow a$

This grammar is not LL(1) because

$$\text{FIRST}(T^*F) \cap \text{FIRST}(F) = \{(.a) \neq \emptyset$$

We will need to build the LR(0) sets of items.

we start with  $I_0$

- $I_0: E \rightarrow .E$
- $\text{CLOSURE}(I_0)$
- $E \rightarrow \lambda.E\lambda$
- $A \rightarrow \gamma.BB$

we add all productions starting with E and add the . at the start therefore :

$$I_0: \{E \rightarrow .E, E \rightarrow .E+T, E \rightarrow .T\}$$

Lets look at  $E \rightarrow .E+T$

- $E \rightarrow \lambda.E+T$
- $A \rightarrow \gamma.BB$

Therefore, we add all productions starting with T to  $I_0$

$$I_0: \{E \rightarrow .E, E \rightarrow .E+T, E \rightarrow .T, T \rightarrow .T^*F, T \rightarrow .F\}$$

we iterate again, and the resultant set is:

$$I_0: \{E \rightarrow .E, E \rightarrow .E+T, E \rightarrow .T, T \rightarrow .T^*F, T \rightarrow .F, F \rightarrow .(E), F \rightarrow .a\}$$

**The GOTO function**

```
function GOTO(I,X)//I=set of items,X=Grammar symbol
{
    CLOSURE(all items A-->alphaX.Beta Where A-->alpha.XBeta in I)
}
```

Lets apply this to the grammar above. First we separate each grammar symbol production to its own set This results in 4 Item groups:

- $I_1: E \rightarrow .E, E \rightarrow .E+T$
- $I_2: E \rightarrow .T, T \rightarrow .T^*F$
- $I_3: T \rightarrow .F$
- $I_4: F \rightarrow .(E)$
- $I_5: F \rightarrow .a$

and take the CLOSURE for all these sets. The resultant is :

$I_1 : E \rightarrow E., E \rightarrow E.+T$   
 $I_2 : E \rightarrow T., T \rightarrow T.*F$   
 $I_3 : T \rightarrow F.$   
 $I_4 : F \rightarrow (.E), E \rightarrow .E+T, E \rightarrow .T, E \rightarrow .T*F, T \rightarrow .F, F \rightarrow .(E), F \rightarrow .a$   
 $I_5 : F \rightarrow .a$   
 $I_6 : E \rightarrow E+.T, E \rightarrow .T*F, T \rightarrow .F, F \rightarrow .(E), F \rightarrow .a$   
 $I_7 : E \rightarrow T*.F, F \rightarrow .(E), F \rightarrow .a$   
 $I_8 : F \rightarrow (E.), E \rightarrow E.+T$   
 $I_9 : E \rightarrow E+T., T \rightarrow T.*F$   
 $I_{10} : T \rightarrow T*F.$   
 $I_{11} : F \rightarrow (E).$

### Constructing the SLR table

Input : LR(0) sets of items

Output : SLR(1) parsing table

1. For every item  $A \rightarrow \alpha.AB$  in  $I_i$ ,  $\alpha \in V_T$ , and  $GOTO(I_i, a) = I_j$ , then set; then  $ACTION[i, a] = S_j$  (shift and push j on the stack).
2. For item  $A \rightarrow \alpha.$  (complete item) in  $I_i$ ,  $ACTION[i, b] = \text{Reduce by } a \rightarrow A \text{ FOR ALL } b \in FOLLOW(A)$ .
3. For  $S' \rightarrow S$  in  $I_i$ ,  $ACTION[i, \$] = \text{Accept}$ .
4. If  $GOTO(I_i, A) = I_j$  then set,  $GOTO(i, A) = j$ .
5. All remaining entries are error entries.

lets apply this to the example above and generate the table

V	Action	a	+	*	(	)	\$	GOTO	E	T	F
0		S5			S4				1	2	3
1			S6				A				
2			R2	S7		R2	R2				
3			R4	R4		R4	R4				
4		S5			S4				8	2	3
5			R6	R6		R6	R6				
6		S5			S4					9	3
7		S5			S4						10
8			S6			S11					
9			R1		S7	R1	R1				
10			R3	R3		R3	R3				
11			R5	R5		R5	R5				

No conflict  $\rightarrow$  SLR(1) grammar

### Parsing The SLR Table

$E \rightarrow E+T$

```

E --> T
T --> T * F
T --> F
F --> (E)
F --> a

```

Lets examine the sentence

```
a + a $
```

Stack	Remaining	Action
0	a + a \$	S5
0 a 5	+ a \$	R6
0 F 3	+ a \$	R4
0 T 2	+ a \$	R2
0 E 1	+ a \$	S6
0 E 1 + 6	a \$	S5
0 E 1 + 6 a 5	\$	R6
0 E 1 + 6 F 3	\$	R4
0 E 1 + 6 T 9	\$	R1
0 E 1	\$	Accept

But what if the grammar is not SLR?

## LR Parsing Techniques

The main difference between LR and SLR is the CLOSURE function

```

function CLOSURE(I) // I is a set of LR(1) items
{
    Repeat
        for every LR(1) item [A --> alpha.Bbeta, a] in I,
        and for every production B --> gamma in G,
        Add the LR(1) item [B --> .gamma, b] where b belongs to FIRST(Beta a) to I
    Until no more items to be added;
}

```

Where An LR(1) item is an LR(0) item with a **Lookahead Symbol**. For example

$[A \rightarrow \alpha.B, a]$  where  $a$  is the lookahead. The lookahead symbol " $a$ " has no effect whatsoever on an item  $[A \rightarrow \alpha.B, a]$   $B \neq \lambda$  (not complete item) However, if the item is a complete  $[A \rightarrow \alpha., a]$ , this means we reduce by the production  $A \rightarrow \alpha$  on token " $a$ ".

For example, Give the grammar :

```

S --> S'
S --> CC (1)
C --> cC (2)
C --> d (3)

```

$I_0$ :

- $S' \rightarrow .S, \$ I_1$

- $S \rightarrow .CC, \text{FIRST}(\lambda\$) = \text{FIRST}(\$) = \$ I_2$
- $C \rightarrow .cC, \text{FIRST}(C\$) = c, d I_3$
- $C \rightarrow .d, c, d I_4$

$I_1$ :

- $S' \rightarrow .S, \$ \text{Accept}$

$I_2$ :

- $S \rightarrow C.C, \$ I_5$  (if the lookahead is different it goes to a new state)
- $C \rightarrow .cC, \$ I_6$
- $C \rightarrow .d, \$ I_7$

$I_3$ :

- $C \rightarrow c.C, c, d I_8$
- $C \rightarrow .cC, c, d I_3$
- $C \rightarrow .d I_4$

$I_4$ :

- $C \rightarrow d., c, d \text{ Complete}$

$I_5$ :

- $S \rightarrow CC., \$ \text{Complete}$

$I_6$ :

- $C \rightarrow c.C, \$ I_9$
- $C \rightarrow .cC, \$ I_6$
- $C \rightarrow .d, \$ I_7$

$I_7$ :

- $C \rightarrow d., \$ \text{Complete}$

$I_8$ :

- $C \rightarrow cC., c, d \text{ Complete}$

$I_9$ :

- $C \rightarrow cC., \$ \text{Complete}$

V	Action	c	d	\$	GOTO	S	C
0		S3	S4			1	2
1				A			
2		S6	S7				5
3		S3	S4				8
4		R3	R3				
5				R1			
6		S6	S7				9
7				R3			
8		R2	R2				
9				R2			

No conflict, the grammar is an LR grammar.

---

If we look at the above example, we can see that some sets of items have the same core items(LR(0) items), but the lookahead is different. For example  $(I_7, I_4)$ ,  $(I_3, I_6)$ ,  $(I_8, I_9)$ . Lets say we merge the states. What happens is the we add the lookaheads and change the state references for equivalent states . The table becomes this:

V	Action	c	d	\$	GOTO	S	C
0		S3	S4			1	2
1				A			
2		S3	S4				5
3		S3	S4				8
4		R3	R3	R3			
5				R1			
8		R2	R2	R2			

---

This is now a simplified table. if the parsing table after merging has no conflicts(like in the above example), then the grammar is an LALR(1) Grammar.