# RECURSIVELY DESCENT PARSING AND LL PARSING

ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA
DIPARTIMENTO DI
INFORMATICA - SCIENZA E INGEGNERIA

## COSIMO LANEVE

cosimo.laneve@unibo.it

**CORSO 72671 - COMPLEMENTI DI LINGUAGGI DI PROGRAMMAZIONE**
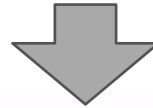
# THIS LECTURE

lexical
analysis

syntactic
analysis

LL-parsing and `ANTLR`

semantic
analysis

**the `SimpLan`
interpreter**

bytecode
generation

# OUTLINE

* recursive descent parsing

    - problems
    - implementations

* predictive parsers

* parsers `LL(1)`

    - **FIRST and FOLLOW sets**
    - `LL(1)` tables
    - ambiguities

* `ANTLR`

**reference**: Torben Morgensen: Basics of Compiler Design, Chapter 3, sections 6—13

# RECURSIVE DESCENT PARSING

analyze the sequence of tokens trying to reconstruct the steps of a **leftmost derivation**

these parsers are called **top-down** because they mimic an anticipated visit of the syntax tree — anticipated = from the root to the leaves

**idea**: the rules for a non-terminal $A$ define a method that recognises $A$

* the right-hand sides of the rules define the structure of the method code

* the sequence of terminals and non-terminals in the rules corresponds to a check that terminals match and to invocations of the methods corresponding to the non-terminal symbols

* the presence of different rules for $A$ is implemented by a `case` or a `if`

# RECURSIVE DESCENT PARSING — EXAMPLE

take the grammar

$$E \rightarrow T + E \mid T$$
$$T \rightarrow (E) \mid (E) * T \mid int \mid int * T$$

the tokens returned by the lexer are

`lpar`    `rpar`    `plus`    `times`    `int(k)` [ k ∈ Nat ]

assume to analyze the token stream:

`int(5) times int(2)`

the parsing starts with the expansion of the initial symbol `E`  and **every rule** of `E` is checked, one at a time `...`

# RECURSIVE DESCENT PARSING — EXAMPLE

try   E → T + E

```
E → T + E | T
T → (E) | (E) * T | int | int * T
```

`int(5) times int(2)`

✳ then you check the **first rule** for `T`:   `T → (E)`

  ● but there is no match with the input token `int(5)`

✳ then you check the **second rule** for `T`:   `T → (E)* T`

  ● but there is no match with the input token `int(5)`

✳ then you check the **third rule** for `T`:   `T → int`

  ● there is match with the input token `int(5)`
  ● but there is no match with the token `plus` after `T` and the token `times` of the input stream

✳ then you check the **forth rule** for `T`:   `T → int * T`

  ● there is match with `int(5)` and then `times` and `int(2)`
  ● but there is no match with the token `plus` because the input stream ends

✳ **we have saturated the choices for  `T`   without succeeding**

  ● backtrack to the other choices for `E`
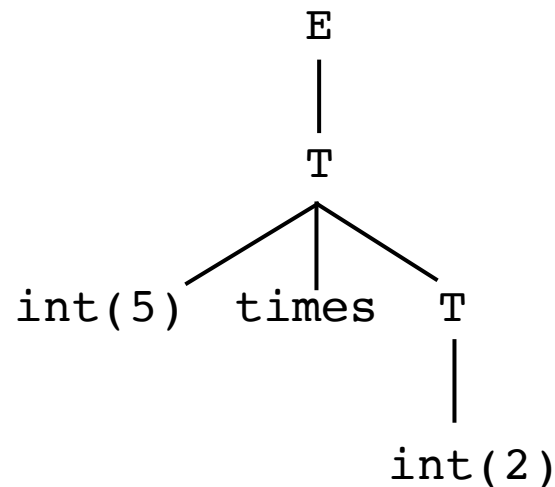
# RECURSIVE DESCENT PARSING — EXAMPLE

```
E → T + E | T
T → (E) | (E) * T | int | int * T
```

`int(5) times int(2)`

then try `E → T` and perform the same steps done for `E → T + E`

✳ the parsing succeeds with the rule `T → int*T` and `T → int`

✳ the returned parse tree is the following one

define a method `boolean` that verifies the matches of the token stream

✳ verify the match with a given **terminal**

```
public boolean term(TOKEN tok){
     TOKEN x = in[next] ;
     next = next + 1 ;
     return x == tok;
}
```

✳ verify the match with a rule of S (say the n-th)

```
public boolean S_n(){ ... }
```

✳ verify the match with a whatever rule of S:

```
public boolean S(){ ... }
```

**note**: every foregoing method increments `next`

# RECURSIVE DESCENT PARSING — IMPLEMENTATION

for the rule $E \rightarrow T + E$

```
public boolean E_1( ){
    return (T() && term(plus) && E()) ;
}
```

this corresponds to
```
boolean B1 = T();
boolean B2 = B1 && term(plus);
return(B2 && E()) ;
```

for the rule $E \rightarrow T$

```
public boolean E_2( ){
    return T();
}
```

for (all) the rules of $E$ (with backtracking)

```
public boolean E() {
    int saved = next ;
    if (E_1()) return true ;
        else { next = saved ; return (E_2()) ; }
}
```

backtrack!

methods for the non-terminal T

```
E → T + E | T
T → (E) | (E) * T | int | int * T
```

```java
public boolean T_1(){
        return ( term(lpar) && E() && term(rpar) );
}

public boolean T_2(){
    return ( term(lpar) && E() && term(rpar) && term(times) && T() );
}

public boolean T_3(){ return ( term(int) ); }

public boolean T_4(){
        return ( term(int) && term(times) && T() ); }

public boolean T(){
        int saved = next;
        if (T_1()) return true ;
        else { next = saved ;
            if (T_2()) return true ;
            else { next = saved ;
                    if (T_3()) return true ;
                    else { next = saved ; return T_4() ; }
            }
        }
}
```

# RECURSIVE DESCENT PARSING — REMARKS

to trigger the parsing

✳ initialize `next` in such a way it points to the first token

✳ invoke `E( )`

✳ assume that a special character `$` represents the end of the input stream in the array `in[]`

✳ the parsing **ends with success** if, at the end of the execution, `in[next] == $`

✳ **remark**: the execution of the recursive descent parsing coincides with the abstract execution computed at the beginning

✳ **other remark**: this is simple to implement (also by hand) **but it does not work, sometimes!**

take the rule $\quad$ S $\rightarrow$ S a

and try to analyze this rule in the recursive descent parsing

✳ **why the process does not work?**

## Definition: left recursive grammar

A grammar $(\mathbf{N}, \mathbf{T}, \rightarrow, \mathbf{S})$ is **left-recursive** if there is $\mathbf{A} \in \mathbf{N}$ such that

$$\mathbf{A} \rightarrow^+ \mathbf{A}\, \gamma \text{ , for some } \gamma$$

the recursive descent parsing **does not work** for lr-grammars

✳ because it performs **an infinite cycle**

**note**: in these cases you need to **change the grammar by removing the left-recursion** (see following slides)

# RECURSIVE DESCENT PARSING — SUMMING UP

the parsing strategy is **extremely simple**

✳ in case you need to **remove the left-recursion** … but this task can be performed **automatically**

it is **not common** because it uses the backtracking

✳ it is very inefficient

✳ in practice, the backtracking may be reduced or eliminated by changing the grammar (**left-factorization**)

it is **good for small grammars**

✳ you need to be careful: the order of productions is important even after left-recursion is eliminated

✳ try to reverse the order of `T → int*T` and `T → int`

✳ what goes wrong? (consider input `int*int`)

# PREDICTIVE PARSERS — MOTIVATIONS

to avoid the backtracking, it would be useful

✷ if the recursive-descent parser **knows the next production to expand**

✷ **idea**: replace the code

```
saved = next ;
if (E_1()) return true;
else {  next = saved; return E_2(); }
```

with

```
switch ( something )  {
      case L1: return E_1();
      case L2: return E_2();
      default: System.out.print("syntax error") ;
}
```

✷ what is the meaning of "`something`", `L1`, `L2` ?
  they are defined by a **lookahead** (analysis of the next tokens)

# PREDICTIVE PARSING AND LEFT FACTORING

the grammar

$$E \rightarrow T + E \mid T$$
$$T \rightarrow (E) \mid (E)*T \mid int \mid int * T$$

is impossible to predict because

✳ the non-terminal `T` has **two productions** that begin with "`(`" and **two productions** that begin with "`int`"

✳ the non-terminal `E` has the two productions that begin with `T` and **it is not evident how to predict**

this grammar **must be left-factorized** before using predictive parsers (see following slides)

| | | |
|---|---|---|
| E → T X | X → + E \| ε |
| T → (E) Y \| int Y | Y → * T \| ε |

# PREDICTIVE PARSERS

they are similar to recursive-descent parsers except that they can **predict** which production to use

* by looking at the next tokens
* without backtracking

predictive parsers accept `LL(k)` grammars

* `L` means "left-to-right" input scan
* `L` means "leftmost derivation"
* `k` means "predict using k tokens of lookahead"

## we study `LL(1)` analysis

* **ANTLR** uses `LL(*)`, a sophisticated technique that consider as many token as needed (this technique is not covered in this course)

# LL(1) LANGUAGES

in recursive-descent parsers, for each non-terminal and input token there may be several possible productions

LL(1) means: for each non-terminal and input token there may be **at most one production** that can be used

LL(1) parsers can be defined by a **2 dimension table**
* one dimension for the non-terminal to expand
* one dimension for the next token
* the table entry contains the production to use

# PARSER `LL(1)`

in practice, instead of using the code

```
switch ( something )  {
          case L1: return E_1();
          case L2: return E_2();
          default: System.out.print("syntax error") ;
  }
```

use a table `LL(1)` and a **parsing stack**

✳ the `LL(1)` table will replace the switch instruction

✳ the parsing stack will replace the call stack

# PARSER `LL(1)` — PARSING TABLE/EXAMPLE

the `LL(1)` parsing table of

```
E → T X                    X → + E | ε
T → (E) Y | int Y          Y → * T | ε
```

|     | int          | *         | +         | (           | )         | $         |
|-----|--------------|-----------|-----------|-------------|-----------|-----------|
| T   | T → int Y    |           |           | T → (E)Y    |           |           |
| E   | E → T X      |           |           | E → T X     |           |           |
| X   |              |           | X → + E   |             | X → ε     | X → ε     |
| Y   |              | Y → * T   | Y → ε     |             | Y → ε     | Y → ε     |

✳ **for the `[E,int]` entry**: when the non-terminal on the stack is `E` and the next token in input is `int`, use the production `E → T X`

✳ **for the `[Y,+]` entry**: when the non-terminal on the stack is `Y` and the next token in input is `+` then remove `Y` (we'll see why)

✳ the empty entries indicate an **error**: example `[E,*]`

# PARSER `LL(1)` — THE PARSING TABLE

the technique is **similar to recursive descent**, but instead of nondeterminism (and the backtrack)

✳ for every non-terminal `S`, look at the next token, say `a`, and the entry `[S,a]` in the table

we use a **stack** in order to record the terminals and non-terminals in the rhs of the production in `[S,a]`

✳ the input is **rejected** when an erroneous state is found (empty entry in the parsing table)

✳ the input is **accepted** when the entry contains **end-of-input** token

# PSEUDO-ALGORITHM OF `LL(1)` PARSING

```
add $ at the end of the array TOKENS ;
next = 0 ;
stack =⟨S $⟩ ;
repeat
   switch (stack){
     case ⟨X rest⟩: if (LL1_TABLE[X,TOKENS[next]] = $\alpha_1\ldots\alpha_n$)
                            stack = ⟨$\alpha_1\ldots\alpha_n$ rest⟩;
                      else  System.out.println("error") ;
                     break ;
           ⟨t rest⟩: if (t == TOKENS[next]) {
                            stack = ⟨rest⟩ ;
                            next = next+1 ;
                     } else  System.out.println("error") ;
                     break ;
    }
until (stack == ⟨ ⟩)
```

# LL(1) PARSING: EXAMPLE

|   | int | * | + | ( | ) | $ |
|---|-----|---|---|---|---|---|
| T | T → int Y | | | T → (E)Y | | |
| E | E → T X | | | E → T X | | |
| X | | | X → + E | | X → ε | X → ε |
| Y | | Y → * T | Y → ε | | Y → ε | Y → ε |

| **Stack** | **Input** | **Action** |
|-----------|-----------|------------|
| E $ | int * int $ | T X |
| T X $ | int * int $ | int Y |
| int Y X $ | int * int $ | terminal |
| Y X $ | * int $ | * T |
| * T X $ | * int $ | terminal |
| T X $ | int $ | int Y |
| int Y X $ | int $ | terminal |
| Y X $ | $ | ε |
| X $ | $ | ε |
| $ | $ | terminal/ACCEPT |

# THE DEFINITION OF THE `LL(1)` PARSING TABLE

let $G = (\mathbf{N}, \mathbf{T}, \rightarrow, S)$, its `LL(1)` table is defined as follows:

1. it has non-terminal symbols in the rows and terminal symbols in the columns

2. for every rule $X \rightarrow \gamma$ in $G$ and for every $t$ such that $\gamma \rightarrow^* t\,\delta$, add the rule $X \rightarrow \gamma$ in the entry $(X, t)$

3. for every rule $X \rightarrow \gamma$ in $G$ such that $\gamma \rightarrow^* \varepsilon$ add the rule $X \rightarrow \gamma$ in the entry $(X, t)$, for every $t$ such that $S \rightarrow^* \delta\, X\, t\, \delta'$

**they seem difficult to compute!**

the `LL(1)` grammars are those with `LL(1)` parsing tables that **do not have multiple entries**

## Definition: the function `NULLABLE`

Let $G = (\mathbf{N}, \mathbf{T}, \rightarrow, S)$ be a context-free grammar. `NULLABLE` is a function on $G$ defined as follows

$$\mathrm{NULLABLE}(G) = \{ \; A \;\; | \;\; A \rightarrow^* \varepsilon \}$$

remark: by definition $\mathrm{NULLABLE}(G) \subseteq \mathbf{N}$

example:

| | | | | | |
|---|---|---|---|---|---|
| E | → | T X | X | → | + E  \| $\varepsilon$ |
| T | → | (E) Y \| Y | Y | → | * T  \| $\varepsilon$ |

then `NULLABLE(`$G$`) = { X, Y, T, E }`. Are you sure about `E` ?

remark: this definition **is not algorithmic**

## Algorithmic definition: the function `NULLABLE`

Let $G = (\mathbf{N}, \mathbf{T}, \rightarrow, S)$ be a context-free grammar. $\texttt{NULLABLE}_i$ are functions on $G$ defined as follows

1. $\texttt{NULLABLE}_0(G) = \{ A \mid A \rightarrow \varepsilon \text{ in } G\}$
2. $\texttt{NULLABLE}_{i+1}(G) = \texttt{NULLABLE}_i(G)$

$$\cup \; \{ A \mid A \rightarrow A_1 \dots A_n \text{ in } G \; \bigwedge \; A_1, \dots, A_n \in \texttt{NULLABLE}_i(G) \}$$

---

✳ it is easy to show that $\texttt{NULLABLE}_i(G) \subseteq \texttt{NULLABLE}_{i+1}(G) \subseteq \mathbf{N}$

✳ therefore there is $\texttt{k}$ such that $\texttt{NULLABLE}_k(G) = \texttt{NULLABLE}_{k+1}(G)$

$$\text{then } \texttt{NULLABLE}(G) = \texttt{NULLABLE}_k(G)$$

# FUNCTION NULLABLE: EXAMPLES

grammar  $\quad$ E $\rightarrow$ T X $\qquad\qquad$ X $\rightarrow$ + E $\mid$ $\varepsilon$

$\qquad\qquad$ T $\rightarrow$ (E)Y $\mid$ int Y $\quad$ Y $\rightarrow$ * T $\mid$ $\varepsilon$

predicate NULLABLE

$\text{NULLABLE}_0(\,G\,) = \{\, X,\ Y\,\}$

$\text{NULLABLE}_1(\,G\,) = \{\, X,\ Y\,\}$

$\text{NULLABLE}(\,G\,) = \text{NULLABLE}_0(\,G\,)$
$\qquad\qquad = \{\, X,\ Y\,\}$

# FUNCTION NULLABLE: EXAMPLES

grammar

$$Z \rightarrow b \mid X\ Y\ Z$$
$$X \rightarrow Y \mid a$$
$$Y \rightarrow \varepsilon \mid c$$

predicate NULLABLE

$\text{NULLABLE}_0(\,G\,) = \{\, Y \,\}$

$\text{NULLABLE}_1(\,G\,) = \{\, X,\ Y \,\}$

$\text{NULLABLE}_2(\,G\,) = \{\, X,\ Y \,\}$

$\text{NULLABLE}(\,G\,) = \text{NULLABLE}_1(\,G\,)$
$= \{\, X,\ Y \,\}$

grammar

$$S \rightarrow a \mid X$$
$$X \rightarrow Y$$
$$Y \rightarrow X$$

predicate NULLABLE

$\text{NULLABLE}_0(\,G\,) = \varnothing$

$\text{NULLABLE}_1(\,G\,) = \varnothing$

$\text{NULLABLE}(\,G\,) = \text{NULLABLE}_0(\,G\,)$
$= \varnothing$

## Algorithmic definition: the function `FIRST`

Let $G = (\mathbf{N}, \mathbf{T}, \rightarrow, S)$ be a context-free grammar. $\texttt{FIRST}_i$ are functions on $\mathbf{N} \cup \mathbf{T}$ that are defined as follows

1. $\texttt{FIRST}_i(t) = \{\, t \,\}$, with $t \in \mathbf{T}$      `// for every i`

2. $\texttt{FIRST}_0(\mathbf{A}) = \begin{cases} \{\, \varepsilon \,\} & \text{if } \mathbf{A} \in \texttt{NULLABLE}(G) \\[1em] \varnothing & \text{if } \mathbf{A} \notin \texttt{NULLABLE}(G) \ \wedge \ \mathbf{A} \in \mathbf{N} \end{cases}$

3. $\texttt{FIRST}_{i+1}(\mathbf{A}) = \texttt{FIRST}_i(\mathbf{A}) \; \bigcup \; \bigcup\limits_{\substack{\mathbf{A} \rightarrow \alpha_1 \cdots \alpha_n \text{ in } G \\ \forall\, i \in 1..k-1\, :\, \alpha_i \in \texttt{NULLABLE}(G)}} \texttt{FIRST}(\alpha_k) \backslash \{\, \varepsilon \,\}$

---

✳ it is easy to show that, for every $\texttt{i}$: $\texttt{FIRST}_i(\mathbf{A}) \subseteq \texttt{FIRST}_{i+1}(\mathbf{A}) \subseteq \mathbf{T} \cup \{\, \varepsilon \,\}$

✳ therefore there is $\texttt{k}$ such that, for every $\mathbf{A}$, $\texttt{FIRST}_k(\mathbf{A}) = \texttt{FIRST}_{k+1}(\mathbf{A})$

$$\text{then } \texttt{FIRST}(\mathbf{A}) = \texttt{FIRST}_k(\mathbf{A})$$

# SETS `FIRST`: EXAMPLES

grammar

$$E \rightarrow \quad T \ X \qquad\qquad X \rightarrow \quad + \ E \quad | \quad \varepsilon$$

$$T \rightarrow \quad (E) \ Y \ | \ \text{int} \ Y \qquad Y \rightarrow \quad * \ T \quad | \quad \varepsilon$$

sets `FIRST`$_i$     `// we only compute FIRST for nonterminals`

$\text{FIRST}_0(X) = \{ \ \varepsilon \ \}$     $\text{FIRST}_0(Y) = \{ \ \varepsilon \ \}$     $\text{FIRST}_0(E) = \varnothing$     $\text{FIRST}_0(T) = \varnothing$

$\text{FIRST}_1(X) = \{ \ +, \ \varepsilon \ \}$     $\text{FIRST}_1(Y) = \{*, \ \varepsilon \ \}$     $\text{FIRST}_1(E) = \varnothing$     $\text{FIRST}_1(T) = \{ \ ( \ , \text{int} \}$

$\text{FIRST}_2(X) = \{ \ +, \ \varepsilon \ \}$     $\text{FIRST}_2(Y) = \{*, \ \varepsilon \ \}$     $\text{FIRST}_2(E) = \{ \ ( \ , \text{int} \}$     $\text{FIRST}_2(T) = \{ \ ( \ , \text{int} \}$

$$\text{FIRST}(X) = \{ \ +, \ \varepsilon \ \} \quad \text{FIRST}(Y) = \{*, \ \varepsilon \ \}$$

$$\text{FIRST}(E) = \{ \ ( \ , \text{int} \} \quad \text{FIRST}(T) = \{ \ ( \ , \text{int} \}$$

grammar

$$S \rightarrow (S) \ S \quad | \quad \varepsilon$$

sets $FIRST_i$

$$FIRST_0( \ S \ ) = \{ \ \varepsilon \ \}$$

$$FIRST_1( \ S \ ) = \{ \ ( \ , \ \varepsilon \ \}$$

$$FIRST( \ S \ ) = FIRST_1( \ S \ ) = \{ \ ( \ , \ \varepsilon \ \}$$

# SETS `FIRST`: EXAMPLES

grammar

$$
\begin{aligned}
Z &\rightarrow b \quad | \quad X\ Y\ Z \\
X &\rightarrow Y \quad | \quad a \\
Y &\rightarrow \varepsilon \quad | \quad c
\end{aligned}
$$

sets $\text{FIRST}_i$:

$\text{FIRST}_0(\,Z\,) = \varnothing \quad \text{FIRST}_0(\,X\,) = \{\ \varepsilon\ \} \quad \text{FIRST}_0(\,Y\,) = \{\ \varepsilon\ \}$

$\text{FIRST}_1(\,Z\,) = \{\,b\,\} \quad \text{FIRST}_1(\,X\,) = \{\,a,\varepsilon\ \} \quad \text{FIRST}_1(\,Y\,) = \{\,c,\varepsilon\ \}$

$\text{FIRST}_2(\,Z\,) = \{\,a,b,c\,\} \quad \text{FIRST}_2(\,X\,) = \{\,a,c,\varepsilon\ \} \quad \text{FIRST}_2(\,Y\,) = \{\,c,\varepsilon\ \}$

$\text{FIRST}(\,X\,) = \{\,a,c,\varepsilon\,\} \qquad \text{FIRST}(\,Y\,) = \{\,c,\varepsilon\,\}$

$\text{FIRST}(\,Z\,) = \{\,a,\ b,\ c\,\}$

# SETS `FIRST`: EXAMPLES

grammar

$$S \rightarrow X \mid XS$$
$$X \rightarrow X \mid \varepsilon$$

sets `FIRST`$_i$:

$\text{FIRST}_0( \, S \, ) = \{ \, \varepsilon \, \}$     $\text{FIRST}_0( \, X \, ) = \{ \, \varepsilon \, \}$

$\text{FIRST}_0( \, S \, ) = \{ \, \varepsilon \, \}$     $\text{FIRST}_0( \, X \, ) = \{ \, \varepsilon \, \}$

$\text{FIRST}( \, S \, ) = \{ \, \varepsilon \, \}$     $\text{FIRST}( \, X \, ) = \{ \, \varepsilon \, \}$

it is easy to compute $\mathbf{FIRST}(\gamma)$ where $\gamma \in (\mathbf{N} \cup \mathbf{T})^*$

`FIRST`$(\varepsilon) = \{\ \varepsilon\ \}$

`FIRST`$(t\gamma) = \{\ t\ \}$,   with  $t \in \mathbf{T}$

`FIRST`$(A\gamma) = $ `FIRST`$(A)$,   with  $A \notin$ `NULLABLE`$(G)$

`FIRST`$(A\gamma) = $ `FIRST`$(A) \backslash \{\ \varepsilon\ \} \cup$ `FIRST`$(\gamma)$,   with  $A \in$ `NULLABLE`$(G)$

the algorithm for computing `FOLLOW` uses this extension

## Algorithmic definition: the function `FOLLOW`

Let $G = (\mathbf{N}, \mathbf{T}, \rightarrow, S)$ be a context-free grammar. $\texttt{FOLLOW}_i$ are functions on $\mathbf{N}$ and defined as follows

1. $\texttt{FOLLOW}_0(S) = \{\, \$ \,\}$ and $\texttt{FOLLOW}_0(A) = \varnothing$

2. $\texttt{FOLLOW}_{i+1}(X) = \texttt{FOLLOW}_i(X) \bigcup_{Z \,\rightarrow\, \delta\, X\, \gamma \text{ in } G} \texttt{FIRST}(\gamma) \backslash \{\varepsilon\}$

$$\bigcup_{Z \,\rightarrow\, \delta\, X\, \gamma \text{ in } G \text{ and } \texttt{NULLABLE}(\gamma)} \texttt{FOLLOW}_i(Z)$$

---

✳ it is easy to show that, for every `i`: $\texttt{FOLLOW}_i(\mathbf{A}) \subseteq \texttt{FOLLOW}_{i+1}(\mathbf{A}) \subseteq \mathbf{T} \cup \{\, \$ \,\}$

✳ therefore there is `k` such that, for every $\mathbf{A}$, $\texttt{FOLLOW}_k(\mathbf{A}) = \texttt{FOLLOW}_{k+1}(\mathbf{A})$

then $\texttt{FOLLOW}(\mathbf{A}) = \texttt{FOLLOW}_k(\mathbf{A})$

remarks: `(1)` when the initial symbol does not appear on the rhs of produc-tions, "`$`" is the unique symbol in its `FOLLOW`

`(2)` `FOLLOW` never contains "$\varepsilon$"

# FOLLOW SETS — EXAMPLES

grammar

$$E \rightarrow T\ X \qquad\qquad X \rightarrow +\ E\ |\ \varepsilon$$
$$T \rightarrow (E)\ Y\ |\ int\ Y \qquad Y \rightarrow *\ T\ |\ \varepsilon$$

sets $FOLLOW_i$

$FOLLOW_0(\,E\,) = \{\,\$\,\}$   $FOLLOW_0(\,T\,) = \varnothing$   $FOLLOW_0(\,X\,) = \varnothing$   $FOLLOW_0(\,Y\,) = \varnothing$

$FOLLOW_1(\,E\,) = \{\,\$,)\}$   $FOLLOW_1(\,T\,) = \{+,\ \$\,\}$   $FOLLOW_1(\,X\,) = \{\,\$\,\}$   $FOLLOW_1(\,Y\,) = \varnothing$

$FOLLOW_2(\,E\,) = \{\,\$,)\}$   $FOLLOW_2(\,T\,) = \{+,\ \$,)\}$   $FOLLOW_2(\,X\,) = \{\$,)\}$   $FOLLOW_2(\,Y\,) = \{+,\ \$\}$

$FOLLOW_3(\,E\,) = \{\,\$,)\}$   $FOLLOW_3(\,T\,) = \{+,\ \$,)\}$   $FOLLOW_3(\,X\,) = \{\$,)\}$   $FOLLOW_3(\,Y\,) = \{+,\ \$,)\}$

$$FOLLOW(\,E\,) = \{\,\$,)\} \qquad\qquad FOLLOW(\,T\,) = \{+,\ \$,)\}$$

$$FOLLOW(\,X\,) = \{\$,)\} \qquad\qquad FOLLOW(\,Y\,) = \{+,\ \$,)\}$$

# FOLLOW SETS — EXAMPLES

grammar

$$S \rightarrow ( S ) S \quad | \quad \varepsilon$$

$\text{FOLLOW}_i$ sets

$\text{FOLLOW}_0( S ) = \{ \$ \}$

$\text{FOLLOW}_1( S ) = \{ \$, ) \}$

$\text{FOLLOW}( S ) = \text{FOLLOW}_1( S ) = \{ \$, ) \}$

# FOLLOW SETS — EXAMPLES

grammar

$$Z \rightarrow \quad b \quad | \quad X\ Y\ Z$$
$$X \rightarrow \quad Y \quad | \quad a$$
$$Y \rightarrow \quad \varepsilon \quad | \quad c$$

$FOLLOW_i$ sets

$FOLLOW_0(\ Z\ ) = \{\ \$\ \} \quad FOLLOW_0(\ X\ ) = \varnothing \quad FOLLOW_0(\ Y\ ) = \varnothing$

$FOLLOW_1(\ Z\ ) = \{\ \$\ \}$

$FOLLOW_1(\ X\ ) = FOLLOW_0(\ X\ ) \cup FIRST(\ Y\ )\backslash\{\varepsilon\} \cup FIRST(Z\ ) = \{\ c\ \} \cup \{a,\ b,\ c\}$

$FOLLOW_1(\ Y\ ) = FOLLOW_0(\ Y\ ) \cup FIRST(\ Z\ ) \cup FOLLOW_0(X\ ) = \{a,\ b,\ c\}$

$$FOLLOW(\ Z\ ) = \{\ \$\ \} \quad FOLLOW(\ X\ ) = \{a,\ b,\ c\}$$

$$FOLLOW(\ Y\ ) = \{a,\ b,\ c\}$$

# DEFINITION OF `LL(1)` PARSING TABLES

the parsing table $\mathtt{LL}^1{}_G$ for a grammar $G$:

for every $\mathtt{A} \to \alpha$ in $G$ do:

1. for every terminal $\mathtt{t} \in \mathtt{FIRST}(\alpha)$ do

   ✳  $\mathtt{LL}^1{}_G\,[\mathtt{A}, \mathtt{t}] = \mathtt{A} \to \alpha$

2. if $\varepsilon \in \mathtt{FIRST}(\alpha)$, for each $\mathtt{t} \in \mathtt{FOLLOW}(\mathtt{A})$ do

   ✳  $\mathtt{LL}^1{}_G\,[\mathtt{A}, \mathtt{t}] = \mathtt{A} \to \alpha$

   [ this rule applies also to $\$$, i.e. when $\$ \in \mathtt{FOLLOW}(\mathtt{A})$:

   if $\varepsilon \in \mathtt{FIRST}(\alpha)$ and $\$ \in \mathtt{FOLLOW}(\mathtt{A})$ do $\mathtt{LL}^1{}_G\,[\mathtt{A}, \$] = \mathtt{A} \to \alpha$

   ]

take the grammar

$$E \to T \ X \qquad\qquad X \to + \ E \ | \ \varepsilon$$
$$T \to (E) \ Y \ | \ \text{int} \ Y \qquad Y \to * \ T \ | \ \varepsilon$$

where in the line of `Y` we put `Y → *T` ?
✳ in the columns of `FIRST(*T)` = { `*` }

where in the line of `Y` we put `Y → ε` ?
✳ in the columns of `FOLLOW(Y)` = { `$, +, )` }

|   | int | * | + | ( | ) | $ |
|---|---|---|---|---|---|---|
| T | T → int Y | | | T → (E)Y | | |
| E | E → T X | | | E → T X | | |
| X | | | X → +E | | X → ε | X → ε |
| Y | | Y → *T | Y → ε | | Y → ε | Y → ε |

# REMARKS ABOUT `LL(1)` TABLES

if any entry is **multiply defined** then $G$ is not LL(1)

in particular when

* $G$ is **left recursive**
* $G$ is **not left-factored**
* $G$ is **ambiguous**
* and in other cases as well

most programming language grammars are not `LL(1)`

* there are tools that build `LL(1)` tables
* the parser generator `ANTLR` uses the `LL` approach

# REMOVING LEFT RECURSION

[see def. slide 12] a grammar si called **left-recursive** if it has a non-terminal $A$ such that $A \Longrightarrow^+ A \gamma$ , for some $\gamma$

case of DIRECT LEFT-RECURSION, i.e. there is $A$ such that

$$A \to A \gamma_1$$
$$\vdots$$
$$A \to A \gamma_m$$
$$A \to \delta_1$$
$$\vdots \qquad \qquad \Big\} \quad \delta_1 \ldots \delta_n \text{ do not start with } A$$
$$A \to \delta_n$$

**remark**: the grammar is equivalent to the regular expression

$$( \delta_1 | \ldots | \delta_n )( \gamma_1 | \ldots | \gamma_m )*$$

# REMOVING DIRECT LEFT RECURSION

$$A \rightarrow A\gamma_1$$
$$\vdots$$
$$A \rightarrow A\gamma_m$$

$$\left.\begin{array}{l} A \rightarrow \delta_1 \\ \vdots \\ A \rightarrow \delta_n \end{array}\right\} \delta_1 \ldots \delta_n \text{ do not start with } A$$

is rewritten into — we use a new non terminal $A'$

$$A \rightarrow \delta_1 A' \qquad\qquad A' \rightarrow \gamma_1 A' \qquad\qquad A' \rightarrow \varepsilon$$
$$\vdots \qquad\qquad\qquad \vdots$$
$$A \rightarrow \delta_n A' \qquad\qquad A' \rightarrow \gamma_m A'$$

## remarks

✳   since the $\delta_i$ do not start with $A$ there is no direct left-recursion anymore

✳   since the $A'$ is a new non-terminal, the $\gamma_i$ *cannot start with it*

✳   there may be **indirect left-recursions** if, for some $i$, $\text{NULLABLE}(\gamma_i)$

# REMOVING DIRECT LEFT RECURSION/EXAMPLE

E → E + F        is rewritten into        E → F E'

E → E - F                                 E' → + F E'

E → F                                     E' → - F E'

F → F * T                                 E' → ε

F → F / T                                 F → T F'

F → T                                     F' → * T F'

T → num                                   F' → / T F'

T → (E)                                   F' → ε

                                          T → num

                                          T →(E)

**exercise**: build the `LL(1)` table

# REMOVING INDIRECT LEFT RECURSION

there are several possibilities

1.  case of  MUTUAL LEFT-RECURSION:

$$A_1 \rightarrow A_2\,\gamma_1$$
$$A_2 \rightarrow A_3\,\gamma_2$$
$$\vdots$$
$$A_{k-1} \rightarrow A_k\,\gamma_{k-1}$$
$$A_k \rightarrow A_1\,\gamma_k$$

break the mutual recursion, e.g. replace

$$A_1 \rightarrow A_2\,\gamma_1$$

with

$$A_1 \rightarrow A_1\,\gamma_k \ldots \gamma_1$$

and solve the direct left recursion

2.  there is a production

$$A \rightarrow \gamma\,A\,\delta \quad \text{where} \quad \texttt{NULLABLE}(\gamma)$$

3.  any combination of 1 and 2

it is always possible to rewrite indirect left recursion into direct one

   ✳   the **process is a bit complex**

# LEFT FACTORING

the grammar

```
E →    T + E | T
T →    (E) | (E)*T  |  int  |  int * T
```

is impossible to predict because

&#10033; the non-terminal `T` has **two productions** that begin with "`(`" and **two productions** that begin with "`int`"

&#10033; the non-terminal `E` has the two productions that begin with `T` and **it is not evident how to predict**

the above grammar **must be left-factored** before using predictive parsers

# LEFT-FACTORING — AN EXAMPLE

the grammar

```
E   →    T + E | T
T   →     (E) | (E)*T | int   |   int * T
```

is left-factored as follows

```
E   →   T E'
E'  →  + E  |  ε
T   →  (E) T'  |  int T'
T'  →  *T  |  ε
```

PROBLEM: left-factoring the standard `if-then-else` statement

```
Stat  →  if Exp then Stat else Stat  |  if Exp then Stat
```

# AMBIGUITY

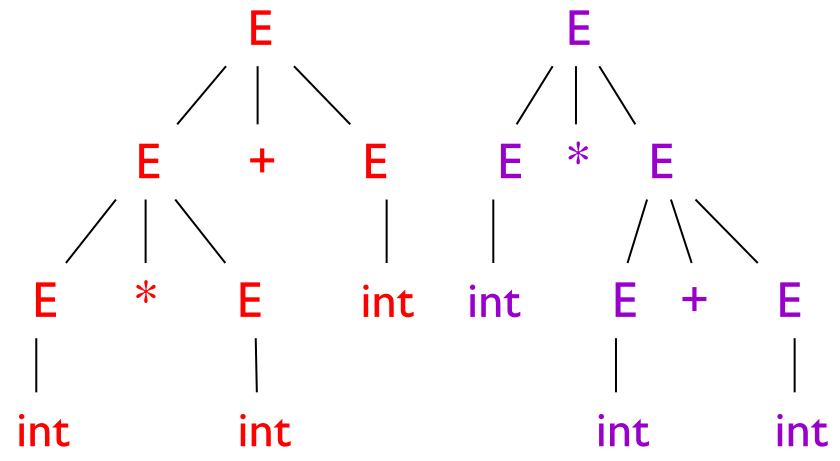a grammar is **ambiguous** if it has **more than one parse tree** for some string

&#x2733; equivalently, there is **more than one rightmost** or **leftmost derivation** for some string

**example**: `E → E + E | E * E | (E) | int`   is ambiguous

because  `int+int+int`   `int*int+int`  have two parse trees



**+ is left-associative**



**\* has higher precedence than +**

47

# AMBIGUITY

ambiguity is **bad**

 ✳  leaves meaning of some programs ill-defined

ambiguity is **common** in programming languages

 ✳  arithmetic expressions
 ✳  `if-then-else`

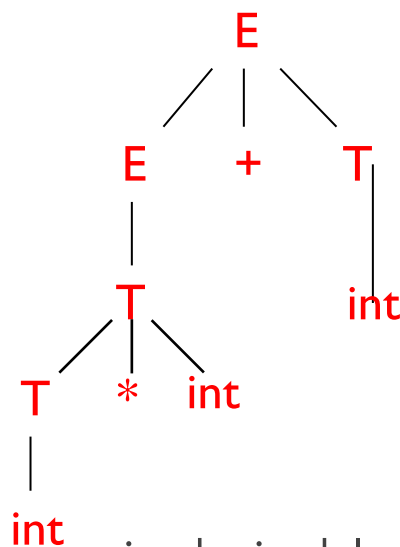in `LL` parsing it is possible to deal with ambiguity by

 ✳  **rewriting grammars**

# DEALING WITH AMBIGUITY
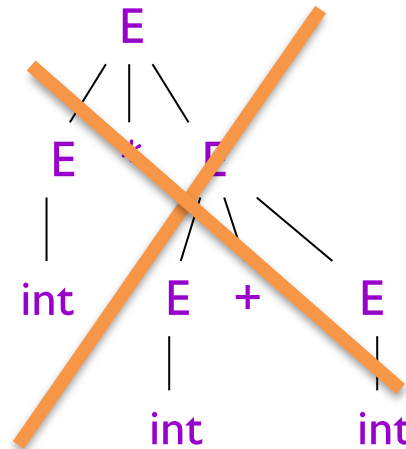
rewrite the grammar for expressions in an unambiguous way:

```
E  →   E + T   |    T
T  →   T * int   |   T * ( E )   |   int   |   ( E )
```

✳ enforces precedence of * over +

✳ enforces left-associativity of + and *



is derivable                is not derivable

the new grammar is neither `LL(1)` nor adequate for rec.descent parsing (left-recursion)
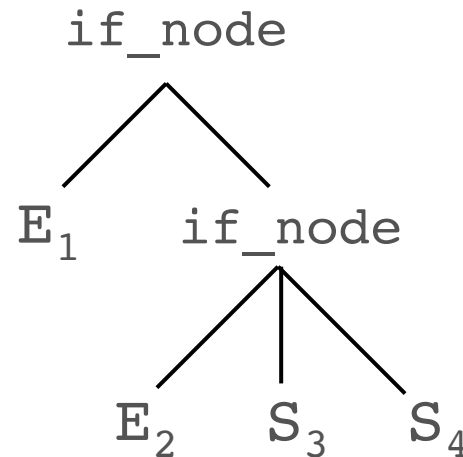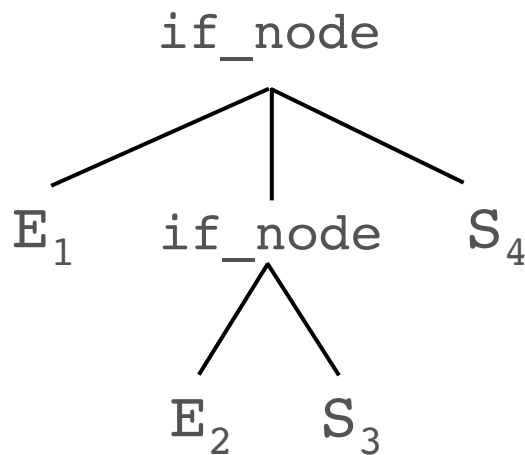
# DEALING WITH AMBIGUITY: THE DANGLING ELSE

the grammar

```
S → ID '=' E | 'if' E 'then' S | 'if' E 'then' S 'else' S
```

is also **ambiguous** because the statement

$$\text{if } E_1 \text{ then if } E_2 \text{ then } S_3 \text{ else } S_4$$

has two (abstract) parse trees



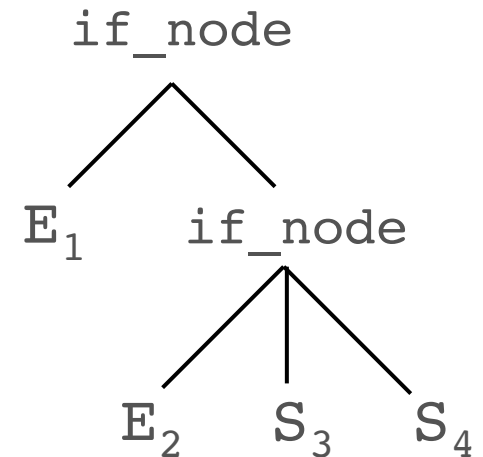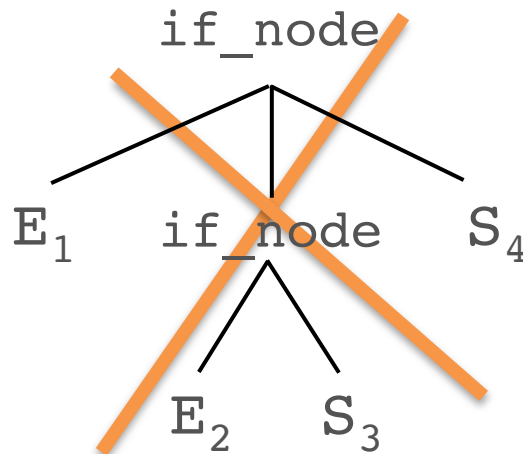in programming languages we want the tree in the right

# THE DANGLING ELSE: A FIX

## else matches the closest unmatched then

we can **factorize** the `if-then` part and **rewrite** the grammar as follows:

```
S     →     ID = E
            | if E then S ELSE

ELSE  →     else S | ε
```

(this new grammar describes the same set of strings and allows the same derivations)

this is a standard hack
to ban this derivation: give
priority to the "`else`" token

# THE DANGLING ELSE: A THEORETICAL FIX

see Gabbrielli-Martini

```
statement              : statementNoIf
                       | ifThenStatement
                       | ifThenElseStatement
                       ;

statementNoIf          : // the statements without if
                       ;
ifThenStatement        : 'if' '(' exp ')' 'then' statement ;

ifThenElseStatement    : 'if' '(' exp ')' 'then' statementNoShortIf 'else'
                           statement ;

statementNoShortIf     : statementNoIf
                       |'if' '(' exp ')' 'then'
                             statementNoShortIf 'else' statementNoShortIf ;
```

the grammar is NOT LL(1): it must be left-factorized!
even if it is accepted in `ANTLR`!

# AMBIGUITY

there is no general techniques for handling ambiguity by transforming grammar

✳   it is always preferable **not to change a grammar**

instead of rewriting the grammar

✳   use the more natural (ambiguous) grammar

✳   along with **disambiguating declarations**

# ANTLR

* start rule

* ambiguity

* left recursion

* non-`LL(*)` decision errors

# ANTLR — START RULE

any grammar needs a so-called start rule

✳ start rule is a rule that is not referenced by another rule

✳ if your grammar does not have such rule, `ANTLR` generator will issue a warning:

```
no start rule (no rule can obviously be followed by EOF)
```

to avoid it, add a dummy start rule to your grammar:

```
start_rule: someOtherRule ;
```

# ANTLR — AMBIGUITY

example:

```
exp : exp '+' exp | exp '*' exp | NUM | '(' exp ')';
NUM: ('0'..'9')+;
```

this grammar should recognize inputs for a simple calculator

✳   `ANTLR v4` behaves badly

✳   to understand the problem, recall that `ANTLR` goes from **left to right** whenever parsing an input

- it first decides which alternative it will use **following the order of the rules**

- then it sticks with the decision

✳   **remark**: left-factorization is solved  automatically by  `ANTLR`

# ANTLR — AMBIGUITY

**example**:

```
exp : exp '+' exp | exp '*' exp | NUM | '(' exp ')';
NUM: ('0'..'9')+;
```

✳ try to simulate it on the input

$$1 + 3 * 4$$

✳ does it match an `exp '+' exp` alternative, or an `exp '*' exp` alternative?

✳ the error suggests to reorder the rules as follows:

```
exp : exp '*' exp | exp '+' exp | NUMBER | '(' exp ')';
NUMBER: ('0'..'9')+;
```

# ANTLR — PROBLEMS WITH ASSOCIATIVITY

**problem**: extend the grammar with "/" and force it to be right-associative

```
exp : exp '/' exp | exp '*' exp | exp '+' exp | '(' exp ')' | NUM   ;
NUM: ('0'..'9')+;
```

**solution:** use a new nonterminal!

```
exp : term | exp '+' term ;
term : factor | factor '/' term | term '*' factor  ;
factor : '(' exp ')' | NUM  ;
```

**better solution:** use the "right-associativity" annotation!

```
exp : <assoc=right> exp '/' exp       // the annotation must be written to
      | exp '+' exp                   // the left
      | exp '*' exp
      | '(' exp ')'
      | NUM                 ;
```

# NEXT LECTURE



lexical analysis

syntactic analysis

semantic analysis

bytecode generation

the `SimpLan` interpreter