# Interpretation and Compilation of Languages
## Master Programme in Computer Science

Mário Pereira     `mjp.pereira@fct.unl.pt`

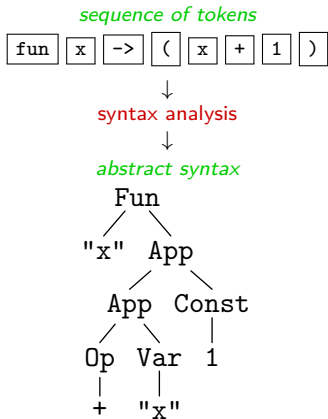Nova School of Science and Technology, Portugal

April 1, 2025

## Lecture 4

based on lectures by Jean-Christophe Filliâtre and Léon Gondelman
previous editions by João Costa Seco, Luís Caires, and Bernardo Toninho

Is to recognize the phrases belonging to the syntax of language.
It's input is the flow of tokens constructed by lexical analysis.
It's output is an abstract syntax tree.

Additionally, syntax analysis must detect syntax errors and

- signal them with a position in the source
- explain them (most often limited to "syntax error" but also "unclosed parenthesis", etc.)
- possibly resume the analysis to discover further errors

To implement syntax analysis, we are using

- a context-free grammar to define the syntax
- a pushdown automaton to recognize it

This is similar to regular expressions / finite automata used for the lexical analysis.

# Today: Syntactic Analysis (Parsing), part 1

1. context-free grammars (recall)

2. top-down parsing algorithm

   - how the algorithm works

   - constructing the expansion table (key ingredient)

## Definition

*A context-free grammar is a tuple $(N, T, S, R)$ where*

- *$N$ is a finite set of nonterminal symbols*
- *$T$ is a finite set of terminal symbols*
- *$S \in N$ is the start symbol (the axiom)*
- *$R \subseteq N \times (N \cup T)^\star$ is a finite set of production rules*

$N = \{E\}$, $T = \{+, *, (,), \texttt{int}\}$, $S = E$,
and $R = \{(E, E\texttt{+}E), (E, E\texttt{*}E), (E, \texttt{(}E\texttt{)}), (E, \texttt{int})\}$

In practice, we write production rules as follows:

$$
\begin{aligned}
E \quad \rightarrow \quad & E + E \\
| \quad & E * E \\
| \quad & ( E ) \\
| \quad & \texttt{int}
\end{aligned}
$$

The terminals are the tokens produced by the lexical analysis.

Here `int` stands for an integer literal token (*i.e.*, its nature, not its value).

Words are of the form $a^n b^n$ with $n \in \mathbb{N}$

$N = \{S\}$, $T = \{a, b\}$ and $R = \{(S, aSb), (S, \epsilon)\}$

$$
\begin{aligned}
S \quad &\to \quad aSb \\
&| \quad \epsilon
\end{aligned}
$$

Words recognized by the grammar: $\epsilon$, *ab*, *aabb*, ...
Words not recognized by the grammar: *a*, *b*, *aab*, *aabbbb*, ...

Words over alphabet $\{a, b\}$ with the same number of $a$'s and $b$'s:

$N = \{S\}$, $T = \{a, b\}$ and
$R = \{(S, \epsilon), (S, aAS), (S, bBS), (A, aAA), (A, b), (B, bBB), (B, a)\}$

$$
\begin{aligned}
S &\rightarrow \epsilon \\
  &\mid a\ A\ S \\
  &\mid b\ B\ S \\
A &\rightarrow a\ A\ A \\
  &\mid b \\
B &\rightarrow b\ B\ B \\
  &\mid a
\end{aligned}
$$

Words recognized by the grammar: $\epsilon$, *ab*, *abbaba*, ...
Words not recognized by the grammar: *a*, *b*, *aab*, *bbba*, ...

### Definition

A word $u \in (N \cup T)^\star$ is *derived* into a word $v \in (N \cup T)^\star$, which we write $u \to v$, if there is a decomposition

$$u = u_1 X u_2$$

with $X \in N$, $X \to \beta \in R$, and

$$v = u_1 \beta u_2$$

Example:

$$\underbrace{E * (}_{u_1} \underbrace{E}_{X} \underbrace{)}_{u_2} \quad \to \quad E * ( \underbrace{E + E}_{\beta} )$$

A sequence $w_1 \rightarrow w_2 \rightarrow \cdots \rightarrow w_n$ is called a derivation.

It is a leftmost derivation (resp. rightmost) if the nonterminal $X$ is the leftmost in each of $w_i$, *i.e.*, $u_1 \in T^\star$ (resp. the rightmost *i.e.* $u_2 \in T^\star$) where $w_i = u_1 X u_2$ for some $u_1$ and $u_2$.

We note $\rightarrow^\star$ the reflexive transitive closure of $\rightarrow$.

$$
\begin{array}{rcl}
E & \rightarrow & E * E \\
  & \rightarrow & \texttt{int} * E \\
  & \rightarrow & \texttt{int} * (\ E\ ) \\
  & \rightarrow & \texttt{int} * (\ E + E\ ) \\
  & \rightarrow & \texttt{int} * (\ \texttt{int} + E\ ) \\
  & \rightarrow & \texttt{int} * (\ \texttt{int} + \texttt{int}\ )
\end{array}
$$

In particular, we get

$$
E \rightarrow^{\star} \texttt{int} * (\ \texttt{int} + \texttt{int}\ )
$$

### Definition

The *language* defined by a context-free grammar $G = (N, T, S, R)$ is the set of words from $T^\star$ that are derived from the axiom, i.e.

$$L(G) = \{\, w \in T^\star \mid S \to^\star w \,\}$$

In our example

$$\texttt{int * ( int + int )} \in L(G)$$

### Definition

*A derivation tree is a tree whose nodes are labeled with grammar symbols, such that*

- *the root is the axiom S*
- *any internal node X is a nonterminal whose subnodes are labeled by $\beta \in (N \cup T)^{\star}$ with $X \to \beta$ a production rule*

Example:

```
          E
        / | \
      E   +   E
      |      / | \
     int   E  *  E
           |     |
          int   int
```

Careful: this is different from the abstract syntax tree.

- Derivation trees describe acceptance of the input by a given grammar and can feature concrete syntax (e.g. parentheses).

- For a derivation tree whose leaves form the word $w$ in infix order, it is clear that there exists $S \to^\star w$.

- Conversely, any derivation $S \to^\star w$ corresponds to a derivation tree whose leaves form the word $w$ in infix order.

- The notion of derivation tree helps to characterize whether a grammar is ambiguous or not, which is important for building parsers.

The leftmost derivation

$$E \rightarrow E + E \rightarrow \text{int} + E \rightarrow \text{int} + E * E \rightarrow \text{int} + \text{int} * E \rightarrow \text{int} + \text{int} * \text{int}$$

Corresponds to the derivation tree



And so is the rightmost derivation

$$E \rightarrow E + E \rightarrow E + E * E \rightarrow E + E * \text{int} \rightarrow E + \text{int} * \text{int} \rightarrow \text{int} + \text{int} * \text{int}$$

### Definition

*A context-free grammar is ambiguous if at least one word accepts several derivation trees.*

Example: the word `int + int * int` accepts two derivation trees



and thus our grammar is ambiguous:

- the parser can't know it should choose the left tree to build the AST
- if it chooses the right tree, it will build an AST with a wrong meaning

It is possible to propose another grammar, that is not ambiguous and that defines the same language.

$$
\begin{aligned}
E &\rightarrow E + T \\
&\mid T \\
T &\rightarrow T * F \\
&\mid F \\
F &\rightarrow ( E ) \\
&\mid \texttt{int}
\end{aligned}
$$

This new grammar reflects the priority of multiplication over addition, and the choice of a left associativity for these two operations.

Now, the word int + int * int * int has a single derivation tree,



corresponding to this leftmost derivation

$$E \to E + T \to T + T \to F + T \to \text{int} + T \to \text{int} + T * F$$
$$\to \text{int} + T * F * F \to \text{int} + F * F * F \to \text{int} + \text{int} * F * F$$
$$\to \text{int} + \text{int} * \text{int} * F \to \text{int} + \text{int} * \text{int} * \text{int}$$

Whether a context-free grammar is ambiguous is not decidable.

(*reminder*: decidable means that we can write a program that, for any
 input, terminates and outputs yes or no)

# Approach

We are going to use decidable sufficient criteria to ensure that a grammar is not ambiguous, and for which we know how to decide membership efficiently (using a pushdown automaton).

The corresponding grammar classes defined by those criteria are called LL(1), LR(0), SLR(1), LALR(1), LR(1), etc.

| | | |
|---|---|---|
| LL | : | Left-to-right, leftmost derivation |
| LR | : | left-to-right, rightmost derivation |
| SLR | : | simple left-to-right, rightmost derivation |
| LALR | : | look-ahead, left-to-right, rightmost derivation |

# Top-down Parsing

Proceed by successive expansions of the leftmost non-terminal (thus constructing a left derivation).

- starting from $S$ and

- using a table $T(X, \overrightarrow{c})$ that, for a non-terminal $X$ to be expanded and the first $k$ characters of the input $\overrightarrow{c}$, indicates the expansion $X \rightarrow \beta$ to be carried out.

(This is referred to as *top-down parsing*).

Suppose $k = 1$ subsequently, and let $T(X, c)$ denote this table.
    That is why we will refer to this technique as LL(1) grammars.

In practice, we assume that a terminal symbol $\#$ denotes the end of the input, and thus the table also indicates the expansion of $X$ when the end of the input is reached.

We use a stack, which is a word from $(N \cup T)^\star$

Initially, the stack is reduced to the start symbol and at each step we scan the top of the stack and the first character $c$ of the input:

- if the stack is empty, we stop and accept the input if and only if $c$ is $\#$
- if the top of the stack is a terminal $a$, then $a$ must be equal to $c$. We pop $a$ from the stack and consume $c$; otherwise, we fail
- if the top of the stack is a non-terminal $X$, then we replace $X$ by the word $\beta = T(X, c)$ on the top of the stack, by pushing the characters of $\beta$ starting from the last one; otherwise, we fail.

Let's rewrite the grammar of the arithmetic expressions and take a look on the following expansion table:

$$
\begin{aligned}
E &\rightarrow T\,E' \\
E' &\rightarrow +\,T\,E' \\
&\mid \epsilon \\
T &\rightarrow F\,T' \\
T' &\rightarrow *\,F\,T' \\
&\mid \epsilon \\
F &\rightarrow (\,E\,) \\
&\mid \texttt{int}
\end{aligned}
$$

|     | +     | *     | (     | )   | int   | #   |
| --- | ----- | ----- | ----- | --- | ----- | --- |
| E   |       |       | $TE'$ |     | $TE'$ |     |
| E'  | $+TE'$ |       |       | $\epsilon$ |       | $\epsilon$ |
| T   |       |       | $FT'$ |     | $FT'$ |     |
| T'  | $\epsilon$ | $*FT'$ |       | $\epsilon$ |       | $\epsilon$ |
| F   |       |       | $(E)$ |     | $\texttt{int}$ |     |

(we'll see in a moment how to construct such tables)

But first let's illustrate the top-down analysis of the input

$$int + int * int$$

Note: the elements in the stack are presented in reverse order.

| stack $\mapsto$ | input |
|---|---|
| $E$ | `int+int*int`$\#$ |
| $E'\,T$ | `int+int*int`$\#$ |
| $E'\,T'F$ | `int+int*int`$\#$ |
| $E'\,T'$`int` | `int+int*int`$\#$ |
| $E'\,T'$ | `+int*int`$\#$ |
| $E'$ | `+int*int`$\#$ |
| $E'\,T+$ | `+int*int`$\#$ |
| $E'\,T$ | `int*int`$\#$ |
| $E'\,T'F$ | `int*int`$\#$ |
| $E'\,T'$`int` | `int*int`$\#$ |
| $E'\,T'$ | `*int`$\#$ |
| $E'\,T'F*$ | `*int`$\#$ |
| $E'\,T'F$ | `int`$\#$ |
| $E'\,T'$`int` | `int`$\#$ |
| $E'\,T'$ | $\#$ |
| $E'$ | $\#$ |
| $\epsilon$ | $\#$ |

| | + | * | ( | ) | int | # |
|---|---|---|---|---|---|---|
| $E$ | | | $TE'$ | | $TE'$ | |
| $E'$ | $+TE'$ | | | $\epsilon$ | | $\epsilon$ |
| $T$ | | | $FT'$ | | $FT'$ | |
| $T'$ | $\epsilon$ | $*FT'$ | | $\epsilon$ | | $\epsilon$ |
| $F$ | | | $(E)$ | | `int` | |

But first let's illustrate the top-down analysis of the input

```
int + int * int
```

Note: the elements in the stack are presented in reverse order.

| stack $\mapsto$ | input |
|---|---|
| $E$ | `int+int*int#` |
| $E'\,T$ | `int+int*int#` |
| $E'\,T'F$ | `int+int*int#` |
| $E'\,T'$`int` | `int+int*int#` |
| $E'\,T'$ | `+int*int#` |
| $E'$ | `+int*int#` |
| $E'\,T+$ | `+int*int#` |
| $E'\,T$ | `int*int#` |
| $E'\,T'F$ | `int*int#` |
| $E'\,T'$`int` | `int*int#` |
| $E'\,T'$ | `*int#` |
| $E'\,T'F*$ | `*int#` |
| $E'\,T'F$ | `int#` |
| $E'\,T'$`int` | `int#` |
| $E'\,T'$ | `#` |
| $E'$ | `#` |
| $\epsilon$ | `#` |

| | + | * | ( | ) | int | # |
|---|---|---|---|---|---|---|
| $E$ | | | $TE'$ | | $TE'$ | |
| $E'$ | $+TE'$ | | | $\epsilon$ | | $\epsilon$ |
| $T$ | | | $FT'$ | | $FT'$ | |
| $T'$ | $\epsilon$ | $*FT'$ | | $\epsilon$ | | $\epsilon$ |
| $F$ | | | $(E)$ | | `int` | |

But first let's illustrate the top-down analysis of the input

```
int + int * int
```

Note: the elements in the stack are presented in reverse order.

| | + | * | ( | ) | int | # |
|---|---|---|---|---|---|---|
| $E$ | | | $TE'$ | | $TE'$ | |
| $E'$ | $+TE'$ | | | $\epsilon$ | | $\epsilon$ |
| $T$ | | | $FT'$ | | $FT'$ | |
| $T'$ | $\epsilon$ | $*FT'$ | | $\epsilon$ | | $\epsilon$ |
| $F$ | | | $(E)$ | | int | |

| stack $\mapsto$ | input |
|---|---|
| $E$ | int+int*int# |
| $E'T$ | int+int*int# |
| $E'T'F$ | int+int*int# |
| $E'T'$int | int+int*int# |
| $E'T'$ | +int*int# |
| $E'$ | +int*int# |
| $E'T+$ | +int*int# |
| $E'T$ | int*int# |
| $E'T'F$ | int*int# |
| $E'T'$int | int*int# |
| $E'T'$ | *int# |
| $E'T'F*$ | *int# |
| $E'T'F$ | int# |
| $E'T'$int | int# |
| $E'T'$ | # |
| $E'$ | # |
| $\epsilon$ | # |

But first let's illustrate the top-down analysis of the input

```
int + int * int
```

Note: the elements in the stack are presented in reverse order.

| | + | * | ( | ) | int | # |
|---|---|---|---|---|---|---|
| $E$ | | | $TE'$ | | $TE'$ | |
| $E'$ | $+TE'$ | | | $\epsilon$ | | $\epsilon$ |
| $T$ | | | $FT'$ | | $FT'$ | |
| $T'$ | $\epsilon$ | $*FT'$ | | $\epsilon$ | | $\epsilon$ |
| $F$ | | | $(E)$ | | int | |

| stack $\mapsto$ | input |
|---|---|
| $E$ | `int+int*int`# |
| $E'T$ | `int+int*int`# |
| $E'T'F$ | `int+int*int`# |
| $E'T'$`int` | `int+int*int`# |
| $E'T'$ | `+int*int`# |
| $E'$ | `+int*int`# |
| $E'T+$ | `+int*int`# |
| $E'T$ | `int*int`# |
| $E'T'F$ | `int*int`# |
| $E'T'$`int` | `int*int`# |
| $E'T'$ | `*int`# |
| $E'T'F*$ | `*int`# |
| $E'T'F$ | `int`# |
| $E'T'$`int` | `int`# |
| $E'T'$ | # |
| $E'$ | # |
| $\epsilon$ | # |

But first let's illustrate the top-down analysis of the input

```
int + int * int
```

Note: the elements in the stack are presented in reverse order.

| | + | * | ( | ) | int | # |
|---|---|---|---|---|---|---|
| $E$ | | | $TE'$ | | $TE'$ | |
| $E'$ | $+TE'$ | | | $\epsilon$ | | $\epsilon$ |
| $T$ | | | $FT'$ | | $FT'$ | |
| $T'$ | $\epsilon$ | $*FT'$ | | $\epsilon$ | | $\epsilon$ |
| $F$ | | | $(E)$ | | int | |

| stack $\mapsto$ | input |
|---|---|
| $E$ | int+int*int# |
| $E'T$ | int+int*int# |
| $E'T'F$ | int+int*int# |
| $E'T'$int | int+int*int# |
| $E'T'$ | +int*int# |
| $E'$ | +int*int# |
| $E'T+$ | +int*int# |
| $E'T$ | int*int# |
| $E'T'F$ | int*int# |
| $E'T'$int | int*int# |
| $E'T'$ | *int# |
| $E'T'F*$ | *int# |
| $E'T'F$ | int# |
| $E'T'$int | int# |
| $E'T'$ | # |
| $E'$ | # |
| $\epsilon$ | # |

But first let's illustrate the top-down analysis of the input

```
int + int * int
```

Note: the elements in the stack are presented in reverse order.

| | + | * | ( | ) | int | # |
|---|---|---|---|---|---|---|
| $E$ | | | $TE'$ | | $TE'$ | |
| $E'$ | $+TE'$ | | | $\epsilon$ | | $\epsilon$ |
| $T$ | | | $FT'$ | | $FT'$ | |
| $T'$ | $\epsilon$ | $*FT'$ | | $\epsilon$ | | $\epsilon$ |
| $F$ | | | $(E)$ | | int | |

| stack $\mapsto$ | input |
|---|---|
| $E$ | int+int*int# |
| $E'T$ | int+int*int# |
| $E'T'F$ | int+int*int# |
| $E'T'$int | int+int*int# |
| $E'T'$ | +int*int# |
| $E'$ | +int*int# |
| $E'T+$ | +int*int# |
| $E'T$ | int*int# |
| $E'T'F$ | int*int# |
| $E'T'$int | int*int# |
| $E'T'$ | *int# |
| $E'T'F*$ | *int# |
| $E'T'F$ | int# |
| $E'T'$int | int# |
| $E'T'$ | # |
| $E'$ | # |
| $\epsilon$ | # |

A top-down parser can be implemented quite easily by introducing a function for each non-terminal of the grammar.

Each function scans the input and, depending on the case as described above, either consumes the input or recursively calls the functions corresponding to other non-terminals, according to the expansion table.

Let's opt for a purely functional programming approach, where the input is a list of tokens of the type:

```
type token = Tplus | Tmult | Tleft | Tright | Tint | Teof
```

we will therefore construct five functions that "consume" the input list:

```
val e : token list -> token list
val e' : token list -> token list
val t : token list -> token list
val t' : token list -> token list
val f : token list -> token list
```

and the recognition of an input can then be done like this

```
let recognize l =
  e l = [Teof]
```

# Implementing a Top-down Parser

The functions proceed by pattern matching on the input and follow the table:

|   | + | * | ( | ) | int | # |
|---|---|---|---|---|-----|---|
| $E$ |   |   | $TE'$ |   | $TE'$ |   |

```
let rec e = function
  | (Tleft | Tint) :: _ as m -> e' (t m)
  | _ -> error ()
```

|   | + | * | ( | ) | int | # |
|---|---|---|---|---|-----|---|
| $E'$ | $+TE'$ |   |   | $\epsilon$ |   | $\epsilon$ |

```
and e' = function
  | Tplus :: m -> e' (t m)
  | (Tright | Teof) :: _ as m -> m
  | _ -> error ()
```

| | + | * | ( | ) | int | # |
|---|---|---|---|---|---|---|
| $T$ | | | $FT'$ | | $FT'$ | |

```
and t = function
  | (Tleft | Tint) :: _ as m -> t' (f m)
  | _ -> error ()
```

| | + | * | ( | ) | int | # |
|---|---|---|---|---|---|---|
| $T'$ | $\epsilon$ | $*FT'$ | | $\epsilon$ | | $\epsilon$ |

```
and t' = function
  | (Tplus | Tright | Teof) :: _ as m -> m
  | Tmult :: m -> t' (f m)
  | _ -> error ()
```

|   | + | * | ( | ) | int | # |
|---|---|---|---|---|-----|---|
| $F$ |   |   | ($E$) |   | int |   |

```
and f = function
  | Tint :: m -> m
  | Tleft :: m -> begin match e m with
      | Tright :: m -> m
      | _ -> error ()
    end
  | _ -> error ()
```

Note that in our implementation

- the expansion table is *implicit*: it is within the code of each function
- the stack is also *not explicit*: it is implemented by the *call stack*
- they could have been made explicit
- alternatively, a more imperative programming approach could have been chosen

```
val next_token : unit -> token
```

One important question remains: how to construct the expansion table systematically for a given grammar?

# Expansion Table

The idea itself is simple: in order to decide whether to perform the expansion $X \to \beta$ when the first character of the input is $c$, we will try to determine if $c$ is among the first characters of the words recognized by $\beta$.

But there is a small catch: a difficulty arises for a production such as $Y \to \epsilon$, and in that case, it's necessary to consider also the set of characters that can follow $Y$ (e.g, it could be that $\beta = Y\beta'$ for some $\beta'$ so the first of $\beta'$ are also the first of $\beta$).

... but for a rule $Z \to \gamma$, it could be that $Z \to^\star \epsilon$ (i.e. indirectly via other rules), so determining the first and follow sets also requires determining if an expansion $\gamma$ can derive $\epsilon$, i.e. have to compute the sets of null.

### Definition (NULL)

Let $\alpha \in (T \cup N)^{\star}$. NULL$(\alpha)$ holds if and only if we can derive $\epsilon$ from $\alpha$ i.e. $\alpha \rightarrow^{\star} \epsilon$.

### Definition (FIRST)

Let $\alpha \in (T \cup N)^{\star}$. FIRST$(\alpha)$ is the set of all terminals starting words derived from $\alpha$, i.e. $\{a \in T \mid \exists w. \alpha \rightarrow^{\star} aw\}$.

### Definition (FOLLOW)

Let $X \in N$. FOLLOW$(X)$ is the set of all terminals that may appear after $X$ in a derivation, i.e. $\{a \in T \mid \exists u, w. S \rightarrow^{\star} uXaw\}$.

To compute $\text{NULL}(\alpha)$, we simply need to compute $\text{NULL}(X)$ for $X \in N$

$\text{NULL}(X)$ holds if and only if

- there exists a production $X \to \epsilon$,
- or there exists a production $X \to Y_1 \ldots Y_m$ where $\text{NULL}(Y_i)$ for all $i$

issue: this is a set of mutually recursive equations

said otherwise, if $N = \{X_1, \ldots, X_n\}$ and if $\vec{V} = (\text{NULL}(X_1), \ldots, \text{NULL}(X_n))$, we look for the least fixpoint to an equation such as

$$\vec{V} = F(\vec{V})$$

## Theorem (existence of a least fixpoint (Tarski))

*Let $A$ be a finite set with an order relation $\leq$ and a least element $\varepsilon$. Any monotonically increasing function $f : A \rightarrow A$, i.e., such that $\forall x, y . x \leq y \Rightarrow f(x) \leq f(y)$, has a least fixpoint.*

(worry not dear students, this is just for you to know that
mathematical magic will take good care of problems)

To compute NULL, we have
$A = \text{BOOL} \times \cdots \times \text{BOOL}$ with $\text{BOOL} = \{\texttt{false}, \texttt{true}\}$

We can equip BOOL with order $\texttt{false} \leq \texttt{true}$ and $A$ with point-wise order

$$(x_1, \ldots, x_n) \leq (y_1, \ldots, y_n) \quad \text{if and only if} \quad \forall i.\, x_i \leq y_i$$

the theorem applies with

$$\varepsilon = (\texttt{false}, \ldots, \texttt{false})$$

since computing $\text{NULL}(X)$ from $\text{NULL}(X_i)$ is monotonic

To compute $\text{NULL}(X_i)$, we thus start with

$$\text{NULL}(X_1) = \texttt{false}, \ldots, \text{NULL}(X_n) = \texttt{false}$$

and we use the equations until we get a fixpoint *i.e.* until the values $\text{NULL}(X_i)$ do not change anymore

$$E \rightarrow T E'$$
$$E' \rightarrow + T E'$$
$$| \ \epsilon$$
$$T \rightarrow F T'$$
$$T' \rightarrow * F T'$$
$$| \ \epsilon$$
$$F \rightarrow ( E )$$
$$| \ \texttt{int}$$

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|-------|-------|-------|-------|-------|
| false | false | false | false | false |

$$
\begin{aligned}
E &\rightarrow T\,E' \\
E' &\rightarrow +\,T\,E' \\
&\mid\ \epsilon \\
T &\rightarrow F\,T' \\
T' &\rightarrow *\,F\,T' \\
&\mid\ \epsilon \\
F &\rightarrow (\,E\,) \\
&\mid\ \text{int}
\end{aligned}
$$

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|-------|-------|-------|-------|-------|
| false | false | false | false | false |
| false | true | false | true | false |

$$
\begin{aligned}
E &\rightarrow T\,E' \\
E' &\rightarrow +\,T\,E' \\
&\mid\ \epsilon \\
T &\rightarrow F\,T' \\
T' &\rightarrow *\,F\,T' \\
&\mid\ \epsilon \\
F &\rightarrow (\ E\ ) \\
&\mid\ \texttt{int}
\end{aligned}
$$

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|-------|-------|-------|-------|-------|
| false | false | false | false | false |
| false | true  | false | true  | false |
| false | true  | false | true  | false |

We have attained a fixpoint,
this is the result for $\text{NULL}(X_i)$.

Why do we seek for a least fixpoint?

$\Rightarrow$ by induction on the number of steps of the fixpoint computation, we show that if $\text{NULL}(X) = \texttt{true}$ then $X \rightarrow^\star \epsilon$

*(soundness)*

$\Leftarrow$ by induction on the number of steps of derivation $X \rightarrow^\star \epsilon$, we show that $\text{NULL}(X) = \texttt{true}$ in the previous computation

*(completeness)*

Similarly, the equations defining FIRST are mutually recursive

$$\text{FIRST}(X) = \bigcup_{X \to \beta} \text{FIRST}(\beta)$$

and

$$
\begin{aligned}
\text{FIRST}(\epsilon) &= \emptyset \\
\text{FIRST}(a\beta) &= \{a\} \\
\text{FIRST}(X\beta) &= \text{FIRST}(X), \quad \text{if } \neg\text{NULL}(X) \\
\text{FIRST}(X\beta) &= \text{FIRST}(X) \cup \text{FIRST}(\beta), \quad \text{if } \text{NULL}(X)
\end{aligned}
$$

Again, we compute a least fixpoint using Tarski's theorem, with
$A = \mathcal{P}(T) \times \cdots \times \mathcal{P}(T)$, point-wise ordered with $\subseteq$, and with $\varepsilon = (\emptyset, \ldots, \emptyset)$

NULL

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|---|---|---|---|---|
| false | true | false | true | false |

$$
\begin{aligned}
E &\rightarrow T\,E' \\
E' &\rightarrow +\,T\,E' \\
&\mid\ \epsilon \\
T &\rightarrow F\,T' \\
T' &\rightarrow *\,F\,T' \\
&\mid\ \epsilon \\
F &\rightarrow (\,E\,) \\
&\mid\ \texttt{int}
\end{aligned}
$$

FIRST

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|---|---|---|---|---|
| $\emptyset$ | $\emptyset$ | $\emptyset$ | $\emptyset$ | $\emptyset$ |

NULL

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|-------|------|-------|------|-------|
| false | true | false | true | false |

$$
\begin{aligned}
E &\rightarrow T E' \\
E' &\rightarrow + T E' \\
&\mid \epsilon \\
T &\rightarrow F T' \\
T' &\rightarrow * F T' \\
&\mid \epsilon \\
F &\rightarrow ( E ) \\
&\mid \text{int}
\end{aligned}
$$

FIRST

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|-----|------|-----|------|-----|
| $\emptyset$ | $\emptyset$ | $\emptyset$ | $\emptyset$ | $\emptyset$ |
| $\emptyset$ | $\{+\}$ | $\emptyset$ | $\{*\}$ | $\{(,\text{int}\}$ |

NULL

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|-------|------|-------|------|-------|
| false | true | false | true | false |

$$
\begin{aligned}
E &\rightarrow T E' \\
E' &\rightarrow + T E' \\
&\mid \ \epsilon \\
T &\rightarrow F T' \\
T' &\rightarrow * F T' \\
&\mid \ \epsilon \\
F &\rightarrow ( E ) \\
&\mid \ \texttt{int}
\end{aligned}
$$

FIRST

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|-----|------|-----------------|-------|-----------------|
| $\emptyset$ | $\emptyset$ | $\emptyset$ | $\emptyset$ | $\emptyset$ |
| $\emptyset$ | $\{+\}$ | $\emptyset$ | $\{*\}$ | $\{(,\texttt{int}\}$ |
| $\emptyset$ | $\{+\}$ | $\{(,\texttt{int}\}$ | $\{*\}$ | $\{(,\texttt{int}\}$ |

NULL

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|-------|------|-------|------|-------|
| false | true | false | true | false |

$$
\begin{aligned}
E &\rightarrow T\,E' \\
E' &\rightarrow +\,T\,E' \\
&\mid \epsilon \\
T &\rightarrow F\,T' \\
T' &\rightarrow *\,F\,T' \\
&\mid \epsilon \\
F &\rightarrow (\,E\,) \\
&\mid \texttt{int}
\end{aligned}
$$

FIRST

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|---------------------|------------|-----------------|------------|-----------------|
| $\emptyset$ | $\emptyset$ | $\emptyset$ | $\emptyset$ | $\emptyset$ |
| $\emptyset$ | $\{\texttt{+}\}$ | $\emptyset$ | $\{\texttt{*}\}$ | $\{\texttt{(,int}\}$ |
| $\emptyset$ | $\{\texttt{+}\}$ | $\{\texttt{(,int}\}$ | $\{\texttt{*}\}$ | $\{\texttt{(,int}\}$ |
| $\{\texttt{(,int}\}$ | $\{\texttt{+}\}$ | $\{\texttt{(,int}\}$ | $\{\texttt{*}\}$ | $\{\texttt{(,int}\}$ |

NULL

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|-------|-------|-------|-------|-------|
| false | true | false | true | false |

$$
\begin{aligned}
E &\rightarrow T\,E' \\
E' &\rightarrow +\,T\,E' \\
&\mid \epsilon \\
T &\rightarrow F\,T' \\
T' &\rightarrow *\,F\,T' \\
&\mid \epsilon \\
F &\rightarrow (\,E\,) \\
&\mid \texttt{int}
\end{aligned}
$$

FIRST

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|----------------|-------|----------------|---------|----------------|
| $\emptyset$ | $\emptyset$ | $\emptyset$ | $\emptyset$ | $\emptyset$ |
| $\emptyset$ | $\{+\}$ | $\emptyset$ | $\{*\}$ | $\{(,\texttt{int}\}$ |
| $\emptyset$ | $\{+\}$ | $\{(,\texttt{int}\}$ | $\{*\}$ | $\{(,\texttt{int}\}$ |
| $\{(,\texttt{int}\}$ | $\{+\}$ | $\{(,\texttt{int}\}$ | $\{*\}$ | $\{(,\texttt{int}\}$ |
| $\{(,\texttt{int}\}$ | $\{+\}$ | $\{(,\texttt{int}\}$ | $\{*\}$ | $\{(,\texttt{int}\}$ |

We have attained a fixpoint,
this is the result for $\text{FIRST}(X_i)$.

Again, the equations defining FOLLOW are mutually recursive

$$\text{FOLLOW}(X) = \bigcup_{Y \to \alpha X \beta} \text{FIRST}(\beta) \cup \bigcup_{Y \to \alpha X \beta, \, \text{NULL}(\beta)} \text{FOLLOW}(Y)$$

We compute a least fixpoint, using the same domain as for FIRST

Remark: we add a special symbol $\#$ in FOLLOW($S$), which we can do directly, or by adding a rule $S' \to S\#$ to the grammar

$$
\begin{aligned}
E &\rightarrow T\,E' \\
E' &\rightarrow +\,T\,E' \\
&\mid \epsilon \\
T &\rightarrow F\,T' \\
T' &\rightarrow *\,F\,T' \\
&\mid \epsilon \\
F &\rightarrow (\,E\,) \\
&\mid \text{int}
\end{aligned}
$$

NULL

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|---|---|---|---|---|
| false | true | false | true | false |

FIRST

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|---|---|---|---|---|
| $\{$(,int$\}$ | $\{$+$\}$ | $\{$(,int$\}$ | $\{$*$\}$ | $\{$(,int$\}$ |

FOLLOW

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|---|---|---|---|---|
| $\{$#$\}$ | $\emptyset$ | $\emptyset$ | $\emptyset$ | $\emptyset$ |

$$
\begin{aligned}
E &\rightarrow T\,E' \\
E' &\rightarrow +\,T\,E' \\
&\mid \epsilon \\
T &\rightarrow F\,T' \\
T' &\rightarrow *\,F\,T' \\
&\mid \epsilon \\
F &\rightarrow (\,E\,) \\
&\mid \texttt{int}
\end{aligned}
$$

NULL

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|-------|-------|-------|-------|-------|
| false | true | false | true | false |

FIRST

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|-------|-------|-------|-------|-------|
| $\{(,\texttt{int}\}$ | $\{\texttt{+}\}$ | $\{(,\texttt{int}\}$ | $\{\texttt{*}\}$ | $\{(,\texttt{int}\}$ |

FOLLOW

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|-------|-------|-------|-------|-------|
| $\{\#\}$ | $\emptyset$ | $\emptyset$ | $\emptyset$ | $\emptyset$ |
| $\{\#,)\}$ | $\{\#\}$ | $\{\texttt{+},\#\}$ | $\emptyset$ | $\{\texttt{*}\}$ |

$$
\begin{aligned}
E &\rightarrow T\,E' \\
E' &\rightarrow +\,T\,E' \\
&\mid \epsilon \\
T &\rightarrow F\,T' \\
T' &\rightarrow *\,F\,T' \\
&\mid \epsilon \\
F &\rightarrow (\,E\,) \\
&\mid \texttt{int}
\end{aligned}
$$

NULL

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|---|---|---|---|---|
| false | true | false | true | false |

FIRST

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|---|---|---|---|---|
| $\{(,\texttt{int}\}$ | $\{+\}$ | $\{(,\texttt{int}\}$ | $\{*\}$ | $\{(,\texttt{int}\}$ |

FOLLOW

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|---|---|---|---|---|
| $\{\#\}$ | $\emptyset$ | $\emptyset$ | $\emptyset$ | $\emptyset$ |
| $\{\#,)\}$ | $\{\#\}$ | $\{+,\#\}$ | $\emptyset$ | $\{*\}$ |
| $\{\#,)\}$ | $\{\#,)\}$ | $\{+,\#,)\}$ | $\{+,\#\}$ | $\{*,+,\#\}$ |

$$\begin{aligned}
E &\to T\,E' \\
E' &\to +\,T\,E' \\
&\mid \epsilon \\
T &\to F\,T' \\
T' &\to *\,F\,T' \\
&\mid \epsilon \\
F &\to (\,E\,) \\
&\mid \texttt{int}
\end{aligned}$$

NULL

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|---|---|---|---|---|
| false | true | false | true | false |

FIRST

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|---|---|---|---|---|
| $\{\texttt{(},\texttt{int}\}$ | $\{\texttt{+}\}$ | $\{\texttt{(},\texttt{int}\}$ | $\{\texttt{*}\}$ | $\{\texttt{(},\texttt{int}\}$ |

FOLLOW

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|---|---|---|---|---|
| $\{\#\}$ | $\emptyset$ | $\emptyset$ | $\emptyset$ | $\emptyset$ |
| $\{\#,\texttt{)}\}$ | $\{\#\}$ | $\{\texttt{+},\#\}$ | $\emptyset$ | $\{\texttt{*}\}$ |
| $\{\#,\texttt{)}\}$ | $\{\#,\texttt{)}\}$ | $\{\texttt{+},\#,\texttt{)}\}$ | $\{\texttt{+},\#\}$ | $\{\texttt{*},\texttt{+},\#\}$ |
| $\{\#,\texttt{)}\}$ | $\{\#,\texttt{)}\}$ | $\{\texttt{+},\#,\texttt{)}\}$ | $\{\texttt{+},\#,\texttt{)}\}$ | $\{\texttt{*},\texttt{+},\#,\texttt{)}\}$ |

$$E \rightarrow T\,E'$$
$$E' \rightarrow +\,T\,E'$$
$$\mid \epsilon$$
$$T \rightarrow F\,T'$$
$$T' \rightarrow *\,F\,T'$$
$$\mid \epsilon$$
$$F \rightarrow (\,E\,)$$
$$\mid \text{int}$$

NULL

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|-------|------|-------|------|-------|
| false | true | false | true | false |

FIRST

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|----------------|-------|----------------|-------|----------------|
| $\{(,\text{int}\}$ | $\{+\}$ | $\{(,\text{int}\}$ | $\{*\}$ | $\{(,\text{int}\}$ |

FOLLOW

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|-----------|-----------|-------------|-------------|---------------|
| $\{\#\}$ | $\emptyset$ | $\emptyset$ | $\emptyset$ | $\emptyset$ |
| $\{\#,)\}$ | $\{\#\}$ | $\{+,\#\}$ | $\emptyset$ | $\{*\}$ |
| $\{\#,)\}$ | $\{\#,)\}$ | $\{+,\#,)\}$ | $\{+,\#\}$ | $\{*,+,\#\}$ |
| $\{\#,)\}$ | $\{\#,)\}$ | $\{+,\#,)\}$ | $\{+,\#,)\}$ | $\{*,+,\#,)\}$ |
| $\{\#,)\}$ | $\{\#,)\}$ | $\{+,\#,)\}$ | $\{+,\#,)\}$ | $\{*,+,\#,)\}$ |

With FIRST and FOLLOW, we build the expansion table $T(X, a)$ as follows.

For each rule $X \to \beta$ of the grammar,

- define $T(X, a) = \beta$ for every $a \in$ FIRST$(\beta)$
- if NULL$(\beta)$, then also define $T(X, a) = \beta$ for every $a \in$ FOLLOW$(X)$

# Example

$$
\begin{aligned}
E &\rightarrow T\,E' \\
E' &\rightarrow +\,T\,E' \\
&\quad | \quad \epsilon \\
T &\rightarrow F\,T' \\
T' &\rightarrow *\,F\,T' \\
&\quad | \quad \epsilon \\
F &\rightarrow (\,E\,) \\
&\quad | \quad \texttt{int}
\end{aligned}
$$

FIRST

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|---|---|---|---|---|
| $\{(,\texttt{int}\}$ | $\{+\}$ | $\{(,\texttt{int}\}$ | $\{*\}$ | $\{(,\texttt{int}\}$ |

FOLLOW

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|---|---|---|---|---|
| $\{\#,)\}$ | $\{\#,)\}$ | $\{+,\#,)\}$ | $\{+,\#,)\}$ | $\{*,+,\#,)\}$ |

| | $+$ | $*$ | $($ | $)$ | $\texttt{int}$ | $\#$ |
|---|---|---|---|---|---|---|
| $E$ | | | $TE'$ | | $TE'$ | |

$$
\begin{aligned}
E &\rightarrow T E' \\
E' &\rightarrow + T E' \\
&\mid \epsilon \\
T &\rightarrow F T' \\
T' &\rightarrow * F T' \\
&\mid \epsilon \\
F &\rightarrow ( E ) \\
&\mid \text{int}
\end{aligned}
$$

FIRST

| E | E' | T | T' | F |
|---|----|---|----|---|
| $\{(,\text{int}\}$ | $\{+\}$ | $\{(,\text{int}\}$ | $\{*\}$ | $\{(,\text{int}\}$ |

FOLLOW

| E | E' | T | T' | F |
|---|----|---|----|---|
| $\{\#,)\}$ | $\{\#,)\}$ | $\{+,\#,)\}$ | $\{+,\#,)\}$ | $\{*,+,\#,)\}$ |

|    | + | * | ( | ) | int | # |
|----|---|---|---|---|-----|---|
| E  |   |   | $TE'$ |   | $TE'$ |   |
| E' | $+TE'$ |   |   | $\epsilon$ |   | $\epsilon$ |

$$
\begin{aligned}
E &\rightarrow T\,E' \\
E' &\rightarrow +\,T\,E' \\
&\mid \epsilon \\
T &\rightarrow F\,T' \\
T' &\rightarrow *\,F\,T' \\
&\mid \epsilon \\
F &\rightarrow (\,E\,) \\
&\mid \texttt{int}
\end{aligned}
$$

FIRST

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|---|---|---|---|---|
| $\{(,\texttt{int}\}$ | $\{+\}$ | $\{(,\texttt{int}\}$ | $\{*\}$ | $\{(,\texttt{int}\}$ |

FOLLOW

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|---|---|---|---|---|
| $\{\#,)\}$ | $\{\#,)\}$ | $\{+,\#,)\}$ | $\{+,\#,)\}$ | $\{*,+,\#,)\}$ |

| | + | * | ( | ) | int | # |
|---|---|---|---|---|---|---|
| $E$ | | | $TE'$ | | $TE'$ | |
| $E'$ | $+TE'$ | | | $\epsilon$ | | $\epsilon$ |
| $T$ | | | $FT'$ | | $FT'$ | |

$$
\begin{array}{rcl}
E & \to & T\,E' \\
E' & \to & +\,T\,E' \\
 & | & \epsilon \\
T & \to & F\,T' \\
T' & \to & *\,F\,T' \\
 & | & \epsilon \\
F & \to & (\,E\,) \\
 & | & \texttt{int}
\end{array}
$$

FIRST

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|---|---|---|---|---|
| $\{(,\texttt{int}\}$ | $\{+\}$ | $\{(,\texttt{int}\}$ | $\{*\}$ | $\{(,\texttt{int}\}$ |

FOLLOW

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|---|---|---|---|---|
| $\{\#,)\}$ | $\{\#,)\}$ | $\{+,\#,)\}$ | $\{+,\#,)\}$ | $\{*,+,\#,)\}$ |

| | $+$ | $*$ | $($ | $)$ | $\texttt{int}$ | $\#$ |
|---|---|---|---|---|---|---|
| $E$ | | | $TE'$ | | $TE'$ | |
| $E'$ | $+TE'$ | | | $\epsilon$ | | $\epsilon$ |
| $T$ | | | $FT'$ | | $FT'$ | |
| $T'$ | $\epsilon$ | $*FT'$ | | $\epsilon$ | | $\epsilon$ |

$$
\begin{aligned}
E &\rightarrow T E' \\
E' &\rightarrow + T E' \\
&\mid \epsilon \\
T &\rightarrow F T' \\
T' &\rightarrow * F T' \\
&\mid \epsilon \\
F &\rightarrow ( E ) \\
&\mid \texttt{int}
\end{aligned}
$$

FIRST

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|---|---|---|---|---|
| $\{(,\texttt{int}\}$ | $\{+\}$ | $\{(,\texttt{int}\}$ | $\{*\}$ | $\{(,\texttt{int}\}$ |

FOLLOW

| $E$ | $E'$ | $T$ | $T'$ | $F$ |
|---|---|---|---|---|
| $\{\#,)\}$ | $\{\#,)\}$ | $\{+,\#,)\}$ | $\{+,\#,)\}$ | $\{*,+,\#,)\}$ |

|  | + | * | ( | ) | int | # |
|---|---|---|---|---|---|---|
| $E$ |  |  | $TE'$ |  | $TE'$ |  |
| $E'$ | $+TE'$ |  |  | $\epsilon$ |  | $\epsilon$ |
| $T$ |  |  | $FT'$ |  | $FT'$ |  |
| $T'$ | $\epsilon$ | $*FT'$ |  | $\epsilon$ |  | $\epsilon$ |
| $F$ |  |  | $(E)$ |  | int |  |

> **Definition (LL(1) Grammar)**
>
> *A grammar is said to be LL(1) if, in the expansion table, there is at most one production in each cell.*

The grammar from the previous slide is LL(1).

LL stands for "Left to right scanning, Leftmost derivation".

$$
\begin{aligned}
E &\rightarrow E + T \\
  &\mid T \\
T &\rightarrow T * F \\
  &\mid F \\
F &\rightarrow ( E ) \\
  &\mid \texttt{int}
\end{aligned}
$$

FIRST

| E | T | F |
|---|---|---|
| $\{(,\texttt{int}\}$ | $\{(,\texttt{int}\}$ | $\{(,\texttt{int}\}$ |

|   | + | * | ( | ) | int | # |
|---|---|---|---|---|-----|---|
| E |   |   | $E+T/T$ |   | $E+T/T$ |   |
| T |   |   | $T*F/F$ |   | $T*F/F$ |   |
| F |   |   | $(E)$ |   | $\texttt{int}$ |   |

A left-recursion grammar, *i.e.*, containing a production of the form

$$X \quad \rightarrow \quad X\alpha$$
$$| \quad \beta$$

will never be LL(1).

Indeed, the FIRST would be the same for these two productions (no matter the word $\beta$).

One needs to suppress the left recursion, for instance

$$
\begin{aligned}
X &\rightarrow \beta\,X' \\
X' &\rightarrow \alpha\,X' \\
&\quad\mid\; \epsilon
\end{aligned}
$$

Also, if a grammar contains

$$X \rightarrow \mathtt{a}\alpha$$
$$| \quad \mathtt{a}\beta$$

it will never be LL(1).

The problem is, again, the FIRST would be the same for both productions.

One needs to suppress productions that start with same terminal
(left factorization)

$$
\begin{aligned}
X &\rightarrow \texttt{a}\,X' \\
X' &\rightarrow \alpha \\
&\mid \beta
\end{aligned}
$$

LL(1) parsers are relatively simple to write.

However, they require writing somewhat unnatural grammars.

We will turn to another solution next week.

Many compilers use top-down hand-written analyzers.

Examples :
- `javac` ($\approx$ 3 kloc of Java code)
- `rustc` ($\approx$ 16 kloc of Rust code)
- `gcc` ($\approx$ 25 kloc of C++ code)

- Compilers, Principles, Techniques, and Tools, Alfred Aho and Monica S. Lam and Ravi Sethi and Jeffrey D. Ullman, Second Edition[Chapter 4.1-4.4], Addison-Wesley (2006) (*"the dragon book"*).