# A!
**Aalto University**
**School of Science**

## CS-C2160 Theory of Computation

Lecture 6: The Parsing Problem, Parse Trees and Recursive-Descent Parsing

Pekka Orponen
Aalto University
Department of Computer Science

---

Topics:

- The parsing problem, canonical derivations and parse trees
- Recursive-descent parsing
- LL(1) grammars
- * Excursion: Attribute grammars
- * Excursion: Parsing tools in the Scala language
- * Supplement: General definition of LL(1) grammars

Material:

- In Finnish: Sections 3.3–3.5 in Finnish lecture notes
- In English: Section 2.1 in the Sipser book, these slides, Wikipedia pages on top-down parsing and attribute grammars.

---

## Recap: Context-free grammars

**Example:**

A (simplified) grammar for arithmetic expressions in a C-like programming language:

$$
\begin{aligned}
E &\rightarrow T \mid E + T \\
T &\rightarrow F \mid T * F \\
F &\rightarrow a \mid (E) \mid f(L) \\
L &\rightarrow L' \mid \varepsilon \\
L' &\rightarrow E \mid E, L'
\end{aligned}
$$

Deriving the string $f(a+a)*a$ in the grammar:

$$
\begin{aligned}
\underline{E} &\Rightarrow \underline{T} &&\Rightarrow \underline{T}*F &&\Rightarrow \underline{F}*F \\
&\Rightarrow f(\underline{L})*F &&\Rightarrow f(\underline{L'})*F &&\Rightarrow f(\underline{E})*F \\
&\Rightarrow f(\underline{E}+T)*F &&\Rightarrow f(\underline{T}+T)*F &&\Rightarrow f(\underline{F}+T)*F \\
&\Rightarrow f(a+\underline{T})*F &&\Rightarrow f(a+\underline{F})*F &&\Rightarrow f(a+a)*\underline{F} \\
&\Rightarrow f(a+a)*a.
\end{aligned}
$$

---

## The Parsing Problem, Canonical Derivations and Parse Trees

## 6.1 The parsing problem and parse trees

We want to solve the following problem:

Given a context-free grammar $G$ and a string $x$. Is $x \in \mathcal{L}(G)$?

A program that solves this problem for a fixed grammar $G$ is called a *parser* for $G$. In this case only the string $x$ is considered as the input.

There are many alternative techniques to design parsers, especially when the grammar $G$ is of some (practically relevant) special form.

**Aalto University**
**School of Science**

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
5/69

---

## Derivations and parse trees

**Example:**

Recall the grammar $G_{\text{expr}}$:

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

Some derivations of string $a + a$ in the grammar are:

$$
\begin{aligned}
\text{(i)} \quad \underline{E} &\Rightarrow \underline{E} + T \Rightarrow \underline{T} + T \Rightarrow \underline{F} + T \\
&\Rightarrow a + \underline{T} \Rightarrow a + \underline{F} \Rightarrow a + a \\
\text{(ii)} \quad \underline{E} &\Rightarrow E + \underline{T} \Rightarrow \underline{E} + F \Rightarrow \underline{T} + F \\
&\Rightarrow \underline{F} + F \Rightarrow F + \underline{a} \Rightarrow a + a \\
\text{(iii)} \quad \underline{E} &\Rightarrow E + \underline{T} \Rightarrow E + \underline{F} \Rightarrow \underline{E} + a \\
&\Rightarrow \underline{T} + a \Rightarrow \underline{F} + a \Rightarrow a + a
\end{aligned}
$$

The underlines denote which non-terminal variable is substituted in which step.

---

- Two canonical ("standard") derivation orders:
- A derivation is a *leftmost* derivation if in each step a rewrite rule is applied to the leftmost available variable. (To emphasise this, we may use the symbol $\underset{\text{lm}}{\Rightarrow}$ instead of $\Rightarrow$.) Derivation (i) on the previous slide is a leftmost derivation.
- *Rightmost derivations* (symbol $\underset{\text{rm}}{\Rightarrow}$) are defined similarly. Derivation (iii) on the previous slide is a rightmost derivation.

**Aalto University**
**School of Science**

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
7/69

---

- Let $G = (V, \Sigma, R, S)$ be a context-free grammar.
- A *parse tree* in $G$ is an ordered tree $\tau$ where:
  - The nodes in $\tau$ are labelled with elements from $V \cup \Sigma \cup \{\varepsilon\}$ so that (i) non-leaf nodes are labeled with elements in $V$ and (ii) the root is labelled with the start variable $S$.
  - If $A$ is the label of a non-leaf node and $X_1, \ldots, X_k$ are the labels of its (ordered) children, then $A \rightarrow X_1 \ldots X_k$ is a production in $R$.
- The string ("sentential form") *represented* by a parse tree is obtained by listing the labels of its leaf nodes in preorder ("from left to right").

**Aalto University**
**School of Science**

CS-C2160 Theory of Computation / Lecture 6
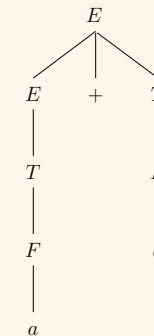Aalto University / Dept. Computer Science
8/69

- Let $G = (V, \Sigma, R, S)$ be a context-free grammar.
- A *parse tree* in $G$ is an ordered tree $\tau$ [1] where:
  - The nodes in $\tau$ are labelled with elements from $V \cup \Sigma \cup \{\varepsilon\}$ so that (i) non-leaf nodes are labeled with elements in $V$ and (ii) the root is labelled with the start variable $S$.
  - If $A$ is the label of a non-leaf node and $X_1, \ldots, X_k$ are the labels of its (ordered) children, then $A \to X_1 \ldots X_k$ is a production in $R$.
- The string ("sentential form") *represented* by a parse tree is obtained by listing the labels of its leaf nodes in preorder ("from left to right").

---
[1] In an ordered tree the children of each node have a fixed left-to-right ordering.

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
9/69

**Example:**

A parse tree for string $a + a$ in grammar $G_{\text{expr}}$:



A derivation for the string:

$$
\begin{aligned}
E &\Rightarrow E + T \Rightarrow T + T \Rightarrow F + T \\
&\Rightarrow a + T \Rightarrow a + F \Rightarrow a + a
\end{aligned}
$$

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
10/69

- A parse tree can be constructed from a derivation

$$S = \gamma_0 \Rightarrow \gamma_1 \Rightarrow \cdots \Rightarrow \gamma_n = \gamma$$
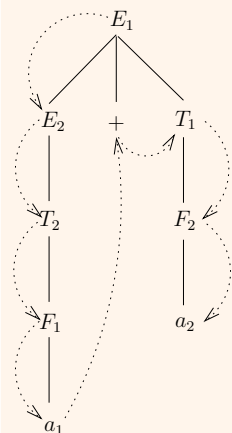
as follows:
  1. The root of the tree is labelled with $S$. If $n = 0$, the tree has no other nodes; otherwise
  2. if the first step in the derivation applies rule $S \to X_1 X_2 \ldots X_k$, the root has $k$ child nodes whose labels from left to right are $X_1, X_2, \ldots, X_k$;
  3. if the next step applies rule $X_i \to Y_1 Y_2 \ldots Y_l$, then the $i$th child node of the root has $l$ children, whose labels from left to right are $Y_1, Y_2, \ldots, Y_l$; and so on.
- We observe that if $\tau$ is the parse tree constructed from derivation $S \Rightarrow^* \gamma$, then the string represented by $\tau$ is $\gamma$.

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
11/69

- Let $\tau$ be a parse tree representing a terminal string $x$.
- We get a leftmost derivation for $x$ by traversing the nodes of $\tau$ in preorder ("from root to leaves, from left to right") and expanding the non-terminal variables encountered as indicated in the tree.
- A rightmost derivation can be obtained similarly by traversing $\tau$ in postorder ("from root to leaves, from right to left").
- By constructing a parse tree from a leftmost derivation and then retrieving the leftmost derivation from the tree, one obtains the original leftmost derivation. The same holds for rightmost derivations.

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
12/69

## Slide 13

**Example:**

Retrieving a leftmost derivation for string $a + a$ from a parse tree.

*Parse tree:*



*Nodes in preorder:*

$$E_1 E_2 T_2 F_1 a_1 + T_1 F_2 a_2$$

*Leftmost derivation:*

$$E \underset{lm}{\Rightarrow} E + T \underset{lm}{\Rightarrow} T + T \underset{lm}{\Rightarrow} F + T$$
$$\underset{lm}{\Rightarrow} a + T \underset{lm}{\Rightarrow} a + F \underset{lm}{\Rightarrow} a + a$$

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
13/69

## Slide 14

### Lemma 6.1

Let $G = (V, \Sigma, R, S)$ be a context-free grammar.

- Each string $\gamma$ that can be derived in $G$ has a parse tree that represents $\gamma$.
- For each parse tree $\tau$ that represents a string $x \in L(G)$ there is a unique leftmost derivation $S \underset{lm}{\Rightarrow}^* x$ and a unique rightmost derivation $S \underset{rm}{\Rightarrow}^* x$.

### Corollary 6.2

Each string $x \in L(G)$ has a leftmost and a rightmost derivation.

That is: parse trees, leftmost derivations and rightmost derivations for words in a language are in one-to-one correspondence.
When solving the parsing problem "Is $x \in L(G)$?", one usually also produces a parse tree (or equivalently a leftmost/rightmost derivation) for $x$ if the answer is "yes".
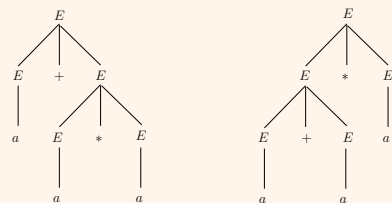
**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
14/69

## Slide 15

### Ambiguity

#### Example

Let us consider the following grammar for simple arithmetic expressions:

$$G'_{\text{expr}} = \{E \to E + E,\ E \to E * E,\ E \to a,\ E \to (E)\}.$$

In this grammar, e.g. string $a + a * a$ has *two* parse trees:



- A context-free grammar $G$ is *ambiguous* if some word $x \in L(G)$ has two different parse trees.
- Otherwise the grammar is *unambiguous*.

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
15/69

## Slide 16

- Ambiguity is usually an unwanted property in computer science because it means that some words have several alternative "interpretations".
- A context-free language for which all the grammars are ambiguous is called an *inherently ambiguous language*.
- As an example, the grammar $G'_{\text{expr}}$ is ambiguous while $G_{\text{expr}}$ is unambiguous. The language $L_{\text{expr}} = L(G'_{\text{expr}})$ is not inherently ambiguous because it also has an unambiguous grammar $G_{\text{expr}}$ generating it.
- On the other hand, e.g. the language

$$\{a^i b^j c^k \mid i = j \text{ or } j = k\}$$

is inherently ambiguous. (The proof of this result is rather complicated and hence omitted here.)

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
16/69

## Recursive-Descent Parsing

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science

17/69

---

## 6.2 Recursive-descent parsing

- One method to search for a leftmost derivation (or parse tree) for a string $x$ in a grammar $G$ is to (i) start from the start variable of $G$ and then (ii) generate systematically and recursively all the possible leftmost derivations (parse trees), (iii) comparing as one proceeds the derived terminal symbols to the ones in the target string $x$.

- If a conflict (= non-match between derived and target symbol) is found, the search backtracks its most recent production rule choice and tries the next available rule.

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science

18/69

---

**Example:**

Let us consider the following grammar $G$:

$$E \rightarrow T + E \mid T - E \mid T$$
$$T \rightarrow a \mid (E)$$

Recursive-descent parsing for the string $a - a$:

$$
\begin{aligned}
E \Rightarrow T+E &\Rightarrow a+T \quad \text{[conflict; backtrack]}\\
&\Rightarrow (E)+T \quad \text{[conflict; backtrack]}\\
\Rightarrow T-E \Rightarrow a-E &\Rightarrow a-T+E \Rightarrow a-a+E\\
&\quad \text{[conflict; backtrack]}\\
&\Rightarrow a-(E)+E\\
&\quad \text{[conflict; backtrack]}\\
&\Rightarrow a-T-E \Rightarrow a-a-E\\
&\quad \text{[conflict; backtrack]}\\
&\Rightarrow a-(E)-E\\
&\quad \text{[conflict; backtrack]}\\
&\Rightarrow a-T \quad \Rightarrow a-a \quad \text{[OK!]}
\end{aligned}
$$

---

- This parsing method can be made efficient if the grammar has the property that at each step the *next symbol* in the input string uniquely determines which rule is to be applied when expanding the leftmost non-terminal variable.
- A grammar that has this property is called an *LL(1) grammar*.
- As an example, we can "factor" the productions of the variable $E$ in the grammar $G$ above and get an equivalent grammar $G'$:

$$
\begin{aligned}
E &\rightarrow T E'\\
E' &\rightarrow +E \mid -E \mid \varepsilon\\
T &\rightarrow a \mid (E)
\end{aligned}
$$

- Parsing the string $a - a$ in $G'$ (at each step, the symbol determining the next rule is marked above the "yields" symbol):

$$E \underset{lm}{\Rightarrow} TE' \overset{a}{\underset{lm}{\Rightarrow}} aE' \overset{-}{\underset{lm}{\Rightarrow}} a-E \underset{lm}{\Rightarrow} a-TE' \overset{a}{\underset{lm}{\Rightarrow}} a-aE' \overset{\varepsilon}{\underset{lm}{\Rightarrow}} a-a.$$

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science

20/69

For an LL(1) grammar, it is easy to write a parser program as a set of recursive procedures. As an example, here is a Python implementation of a parser for the grammar $G'$:

```python
from sys import exit, stdin
def error(s): print(s); exit(1)
def e():
    print("E -> TE'")
    t(); eprime()
def eprime():
    global next
    if next=="+":
        print("E' -> +E")
        next=stdin.read(1)
        e()
    elif next=="-":
        print("E' -> -E")
        next=stdin.read(1)
        e()
    else: print("E ->")
```

Continues on the next slide...

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
21/69

```python
def t():
    global next
    if next=="a":
        print("T -> a")
        next=stdin.read(1)
    elif next=="(":
        print("T -> (E)")
        next=stdin.read(1)
        e()
        if next!=")": error(") expected.")
        next=stdin.read(1)
    else: error("T cannot start with %s"%(next))

next=stdin.read(1)
e()
```

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
22/69

When processing string `a-(a+a)`, the program outputs the following lines:

```
E  -> TE'
T  -> a
E' -> -E
E  -> TE'
T  -> (E)
E  -> TE'
T  -> a
E' -> +E
E  -> TE'
T  -> a
E' ->
E' ->
```

The output corresponds to the leftmost derivation

$$
\begin{aligned}
E \; &\Rightarrow \; TE' \Rightarrow aE' \Rightarrow a-E \Rightarrow a-TE' \\
&\Rightarrow \; a-(E)E' \Rightarrow a-(TE')E' \\
&\Rightarrow \; a-(aE')E' \Rightarrow a-(a+E)E' \\
&\Rightarrow \; a-(a+TE')E' \Rightarrow a-(a+aE')E' \\
&\Rightarrow \; a-(a+a)E' \Rightarrow a-(a+a).
\end{aligned}
$$

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
23/69

## LL(1) Grammars

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
24/69

## 6.3 LL(1) grammars

- Let us next consider the general form of LL(1) grammars.
- LL(1) ≈ "parse input from **L**eft to right and produce a **L**eftmost derivation, using **1** token lookahead".
  Here "1 token lookahead" means that one only considers the next symbol in the target string at a time.
- For instance, the grammar

$$
\begin{aligned}
S &\rightarrow A\,b \mid C\,d \\
A &\rightarrow a\,A \mid \varepsilon \\
C &\rightarrow c\,C \mid \varepsilon
\end{aligned}
$$

  is an LL(1) grammar, even though the right-hand sides of the productions don't always start with a terminal symbol.
- The precise definition of LL(1) grammars is discussed on the supplementary slides at the end of this lecture.

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
25/69

## 6.3 LL(1) grammars

- Let us next consider the general form of LL(1) grammars.
- LL(1) ≈ "parse input from **L**eft to right and produce a **L**eftmost derivation, using **1** token lookahead".
  Here "1 token lookahead" means that one only considers the next symbol in the target string at a time.
- For instance, the grammar

$$
\begin{aligned}
S &\rightarrow A\,b \mid C\,d \\
A &\rightarrow a\,A \mid \varepsilon \\
C &\rightarrow c\,C \mid \varepsilon
\end{aligned}
$$

  is an LL(1) grammar, even though the right-hand sides of the productions don't always start with a terminal symbol.
- The precise definition of LL(1) grammars is discussed on the supplementary slides at the end of this lecture.

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
26/69

## 6.3 LL(1) grammars

- Let us next consider the general form of LL(1) grammars.
- LL(1) ≈ "parse input from **L**eft to right and produce a **L**eftmost derivation, using **1** token lookahead".
  Here "1 token lookahead" means that one only considers the next symbol in the target string at a time. [2]
- For instance, the grammar

$$
\begin{aligned}
S &\rightarrow A\,b \mid C\,d \\
A &\rightarrow a\,A \mid \varepsilon \\
C &\rightarrow c\,C \mid \varepsilon
\end{aligned}
$$

  is an LL(1) grammar, even though the right-hand sides of the productions don't always start with a terminal symbol.
- The precise definition of LL(1) grammars is discussed on the supplementary slides at the end of this lecture.

---

[2] There are also more general notions of "LL($k$)" and "LR($k$)" grammars.

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
27/69

## Left recursion

- Left recursion is a problem for recursive-descent parsing.

### Definition 6.1

A grammar $G = (V, \Sigma, P, S)$ is *left recursive* if one can derive from some variable $A$ with one or more steps the string $A\alpha$, where $\alpha \in (V \cup \Sigma)^\star$.

### Example:

The grammar $G_{\text{expr}}$

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

is left recursive because $E \Rightarrow E + T$ and $T \Rightarrow T * F$.
This kind of left recursion that occurs in a single step is called *immediate left recursion*.

- Left recursion may result in infinite, non-terminating recursion in the parsing process.

**Example:**

In the grammar $G_{\text{expr}}$

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

recursive-descent parsing may start producing the non-terminating derivation

$$
\underline{E} \underset{\text{lm}}{\Rightarrow} \underline{E} + T \underset{\text{lm}}{\Rightarrow} \underline{E} + E + T \underset{\text{lm}}{\Rightarrow} \dots
$$

without ever producing a terminal symbol in the beginning of the derived string.

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
29/69

**Example:**

Also the grammar

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

is left recursive because e.g. $S \Rightarrow A S a \Rightarrow B B S a \Rightarrow B S a \Rightarrow S a$.

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
30/69

# Eliminating immediate left recursion

- Immediate left recursion of form

$$
A \rightarrow A \beta_1 \mid \dots \mid A \beta_n \mid \alpha_1 \mid \dots \mid \alpha_m
$$

can be eliminated by translating it into right recursion

$$
\begin{aligned}
A &\rightarrow \alpha_1 A' \mid \dots \mid \alpha_m A' \\
A' &\rightarrow \beta_1 A' \mid \dots \mid \beta_n A' \mid \varepsilon
\end{aligned}
$$

- Now a derivation of form

$$
A \Rightarrow A \beta_1 \Rightarrow A \beta_2 \beta_1 \Rightarrow \alpha_1 \beta_2 \beta_1
$$

can be "simulated" with the derivation

$$
A \Rightarrow \alpha_1 A' \Rightarrow \alpha \beta_2 A' \Rightarrow \alpha_1 \beta_2 \beta_1 A' \Rightarrow \alpha_1 \beta_2 \beta_1
$$

(Also non-immediate, generic left recursion can be eliminated, see e.g. section 4.3 in the book Aho, Sethi, Ullman: "Compilers — Principles, Techniques, and Tools".)

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
31/69

**Example:**

Eliminating immediate left recursion in the grammar $G_{\text{expr}}$

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

results in the grammar

$$
\begin{aligned}
E &\rightarrow T E' \\
E' &\rightarrow + T E' \mid - T E' \mid \varepsilon \\
T &\rightarrow F T' \\
T' &\rightarrow * F T' \mid \varepsilon \\
F &\rightarrow a \mid (E)
\end{aligned}
$$

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
32/69

# Left factoring

- Another problematic grammar feature for recursive-descent parsing are productions that start with the same symbol.
- As an example, consider statements in the C++ language:

$$\begin{aligned} stmt \;\; &\rightarrow \;\; selection\text{-}stmt \;\mid\; iteration\text{-}stmt \;\mid\; \ldots \\ selection\text{-}stmt \;\; &\rightarrow \;\; \textbf{if} \,(\, expr \,)\, \textbf{then}\, stmt \;\mid\; \\ &\qquad \textbf{if} \,(\, expr \,)\, \textbf{then}\, stmt\, \textbf{else}\, stmt \;\mid\; \\ &\qquad \textbf{switch} \,(\, expr \,)\, stmt \end{aligned}$$

where *iteration-stmt* and others don't start with the **if** symbol.

☞ Based only on the current **if** symbol in the input string, one cannot decide whether the first or the second production for the variable *selection-stmt* should be applied.

**Aalto University**
**School of Science**

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
33/69

---

- Common prefixes of form

$$A \;\;\rightarrow\;\; \alpha\beta_1 \;\mid\; \ldots \;\mid\; \alpha\beta_n \;\mid\; \gamma$$

can be "left factored" as follows:

$$\begin{aligned} A \;\;&\rightarrow\;\; \alpha A' \;\mid\; \gamma \\ A' \;\;&\rightarrow\;\; \beta_1 \;\mid\; \ldots \;\mid\; \beta_n \end{aligned}$$

**Example:**

Left factoring the C++ if-then-else structure

$$\begin{aligned} selection\text{-}stmt \;\;&\rightarrow\;\; \textbf{if} \,(\, expr \,)\, \textbf{then}\, stmt \;\mid\; \\ &\qquad \textbf{if} \,(\, expr \,)\, \textbf{then}\, stmt\, \textbf{else}\, stmt \;\mid\; \\ &\qquad \textbf{switch} \,(\, expr \,)\, stmt \end{aligned}$$

results in

$$\begin{aligned} selection\text{-}stmt \;\;&\rightarrow\;\; \textbf{if} \,(\, expr \,)\, \textbf{then}\, stmt\; selection\text{-}stmt' \;\mid\; \\ &\qquad \textbf{switch} \,(\, expr \,)\, stmt \\ selection\text{-}stmt' \;\;&\rightarrow\;\; \textbf{else}\, stmt \;\mid\; \varepsilon \end{aligned}$$

---

# * Excursion: Attribute Grammars

**Aalto University**
**School of Science**

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
35/69

---

- Attribute grammars are a technique for associating simple semantic rules to context-free grammars.
- Each node in a parse tree, labelled with grammar symbol $X$, is considered an object "of type $X$". The fields in an object of type $X$ are called *attributes* of $X$ and denoted as $X.s$, $X.t$ etc. Each node "object" has its own "instances" of the attribute.
- The productions $A \rightarrow X_1 \ldots X_k$ of the grammar are associated with *evaluation rules* that describe how the values of the respective attribute instances are computed from those in the parent and child nodes.
- The evaluation rules can in principle be arbitrary functions, as long as their parameters only involve locally available information. More precisely, the evaluation rules associated with a production $A \rightarrow X_1 \ldots X_k$ can only mention attributes of the symbols $A, X_1, \ldots, X_k$.

**Aalto University**
**School of Science**

CS-C2160 Theory of Computation / Lecture 6
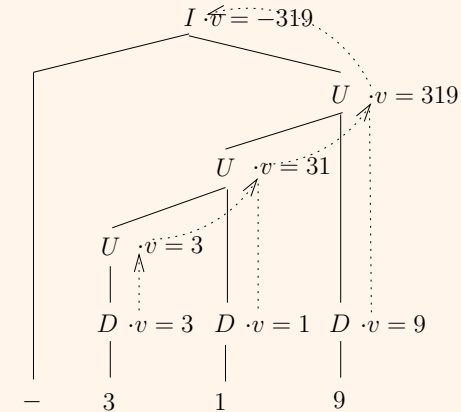Aalto University / Dept. Computer Science
36/69

## Example: Evaluating signed integers

Each node of type $X$ in the parse tree is associated with an attribute instance $X.v$, whose value will be the numeric value of the string derived from $X$. In particular, the value of the instance $v$ in the root node will be the numeric value of the whole string represented by the tree.

*Productions:*      *Evaluation rules:*

$$
\begin{array}{lcl}
I & \to & +U \\
I & \to & -U \\
I & \to & U \\
U & \to & D \\
U & \to & UD \\
D & \to & 0 \\
\cdots \\
D & \to & 9
\end{array}
\qquad
\begin{array}{lcl}
I.v & := & U.v \\
I.v & := & -U.v \\
I.v & := & U.v \\
U.v & := & D.v \\
U_1.v & := & 10 * U_2.v + D.v \\
D.v & := & 0 \\
\\
D.v & := & 9
\end{array}
$$

In the evaluation rule associated with production $U \to UD$, the different instances of variable symbol $U$ are distinguished by the use of indices.

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
37/69

---

The "attributed parse tree" for string "-319":

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
38/69

---

- An attribute $t$ is *synthetic* if the evaluation rule in each production $A \to X_1 \ldots X_k$ mentioning $t$ is of form
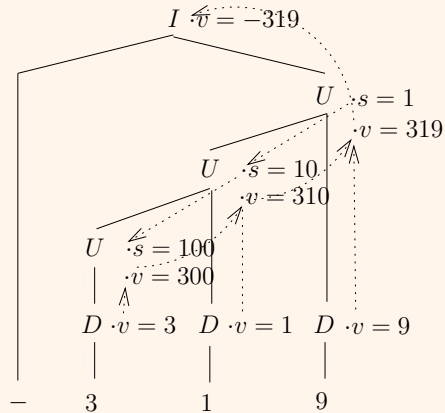
$$A.t := f(A, X_1, \ldots, X_k).$$

- In this case, the value of a $t$ attribute instance depends only on the values of the attribute instances in the node itself and in its child nodes.
- Other forms of attributes are called *inherited*.
- Synthetic attributes are preferable, because their values can be evaluated in a single bottom-up traversal of the parse tree.
- Of course, one can also use inherited attributes, as long as one ensures that there are no dependency cycles in their evaluation rules.

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
39/69

---

## Example

Evaluating signed integers by using an inherited "position multiplier" attribute and a synthetic "value" attribute:

*Productions:*      *Evaluation rules:*

$$
\begin{array}{lcl}
I & \to & +U \\
I & \to & -U \\
I & \to & U \\
U & \to & D \\
U & \to & UD \\
\\
D & \to & 0 \\
\vdots \\
D & \to & 9
\end{array}
\qquad
\begin{array}{rcll}
U.s & := & 1, & I.v := U.v \\
U.s & := & 1, & I.v := -U.v \\
U.s & := & 1, & I.v := U.v \\
& & & U.v := (D.v)*(U.s) \\
U_2.s & := & 10*(U_1.s), \\
U_1.v & := & U_2.v + (D.v)*(U_1.s) \\
& & & D.v := 0 \\
\\
& & & D.v := 9
\end{array}
$$

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
40/69

## Slide 41

For the string "-319" we get the following attributed parse tree:

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
41/69

## Slide 42

- The values of attribute instances can often be computed "on-the-fly" without explicitly constructing the parse tree.

### Example

**Example.** A program that transforms arithmetic expressions from infix notation to postfix notation.

We associate to our grammar $G_{\text{expr}}$ one synthetic, string-valued attribute $pf$. The value of attribute instance $X.pf$ in each parse tree node of type $X$ will be the postfix version of the infix-notation string derived from the node.

| Productions: | | Evaluation rules: | | |
|---|---|---|---|---|
| $E$ | $\to$ $T + E$ | $E_1.pf$ | $:=$ | $(T.pf)\widehat{\ }(E_2.pf)\widehat{\ }('+')$ |
| $E$ | $\to$ $T$ | $E.pf$ | $:=$ | $T.pf$ |
| $T$ | $\to$ $F * T$ | $T_1.pf$ | $:=$ | $(F.pf)\widehat{\ }(T_2.pf)\widehat{\ }('*')$ |
| $T$ | $\to$ $F$ | $T.pf$ | $:=$ | $F.pf$ |
| $F$ | $\to$ $a$ | $F.pf$ | $:=$ | $'a'$ |
| $F$ | $\to$ $(E)$ | $F.pf$ | $:=$ | $E.pf$ |

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
42/69

## Slide 43

A recursive-descent parser for $G_{\text{expr}}$ that also evaluates the values of attribute instances on-the-fly during parsing:

```python
from sys import stdin
def error(s): print(s); exit(1)
def e():            # E -> T + E | T
    global next
    pf1=t()
    if next=="+":
        next=stdin.read(1)
        return pf1+e()+"+"  # E1.pf := T.pf E2.pf +
    else: return pf1         # E.pf  := T.pf
def t():          # T -> F * T | F
    global next
    pf1=f()
    if next=="*":
        next=stdin.read(1)
        return pf1+t()+"*"  # T1.pf := F.pf T2.pf *
    else: return pf1         # T.pf  := F.pf
```

Continues on the next slide...

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
43/69

## Slide 44

```python
def f():           # F -> a | (E)
    global next
    if next=="a":
        next=stdin.read(1)
        return "a"                # F.pf := a
    elif next=="(":
        next=stdin.read(1)
        pf=e()
        if next!=")": error(") expected.")
        next=stdin.read(1)
        return pf                 # F.pf := E.pf
    else: error("F cannot start with this.")

next=stdin.read(1)
print(e())
```

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
44/69

## Slide 45

**\* Excursion: Parsing Tools in the Scala Language**

**Aalto University**
**School of Science**

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
45/69

## Slide 46

- A library suitable for (restricted) recursive-descent parsing is included in the standard Scala language distribution
- Regular expressions can also be included in the grammar rules.

References:

- Chapter 31 in book "Programming in Scala, First Edition"
- Parsers trait
- The abstract Parser class
- The RegexParsers trait

**Aalto University**
**School of Science**

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
46/69

## Slide 47

A parser and evaluator for simple arithmetic expressions:

```scala
import util.parsing.combinator._

object parser extends RegexParsers {
  val integer = "[0-9]+".r  // Regular expression

  def expr: Parser[BigInt] = (  // E -> T + E | T
      term~"+"~expr ^^ {case t~"+"~e => t + e}
    | term ^^ {case t => t} )
  def term: Parser[BigInt] =  // T -> F * T | T
    rep1sep(factor, "*") ^^ { factors => factors.product }
  def factor: Parser[BigInt] = ( // F -> integer | ( E )
      integer ^^ { case intString => BigInt(intString) }
    | "(" ~> expr <~ ")" ^^ { case e => e }  )

  def parse(input: String): (Option[BigInt], String) = {
    parseAll(expr, input) match {
      case Success(value, _) => (Option(value), "success")
      case f => (None, f.toString)
    }
  }
}
```

**Aalto University**
**School of Science**

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
47/69

## Slide 48

- With well-formed input we get the expected result:

```scala
scala> parser.parse("123+4*5")
res0: (Option[BigInt], String) = (Some(143),success)
```

- while an erroneous input gives an error message:

```scala
scala> parser.parse("123++4*5")
res1: (Option[BigInt], String) =
(None,[1.5] failure: '(' expected but '+' found

123++4*5
    ^)
```

**Aalto University**
**School of Science**

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
48/69

- A classic textbook on these topics:
  - Aho, Sethi, Ullman: Compilers — Principles, Techniques, and Tools

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
49/69

# * Supplement: General Definition of LL(1) Grammars

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
50/69

# * Supplement: General definition of LL(1) grammars

- In the following we will formally define LL(1) grammars
- To do this, we need two auxiliary sets
  - FIRST describes which terminal symbols can appear as the first symbols in the strings derivable from a non-terminal variable
  - FOLLOW describes which terminal symbols can follow a non-terminal variable in any of the derivations

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
51/69

- We also need the auxiliary concept of a nullable non-terminal variables

**Definition 6.2**

A non-terminal variable $A$ is *nullable* if $A \Rightarrow^* \varepsilon$.

**Example:**

In the grammar

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

the variables $A$ and $B$ are nullable because

- $\underline{A} \Rightarrow \underline{B}B \Rightarrow \underline{B} \Rightarrow \varepsilon$
- $\underline{B} \Rightarrow \varepsilon$

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
52/69

## FIRST-sets

- For each non-terminal variable $A$ we define the set $\text{FIRST}(A)$ of terminal symbols (incl. $\varepsilon$ is $A$ is nullable) that can be the first symbols in strings derivable from $A$:

$$\text{FIRST}(A) = \{a \in \Sigma \mid A \Rightarrow^* a\gamma \text{ for some } \gamma \in (V \cup \Sigma)^*\} \cup \{\varepsilon \mid A \Rightarrow^* \varepsilon\}$$

**Example:**

In the grammar

$$\begin{aligned} S &\rightarrow Ab \mid Cd \\ A &\rightarrow aA \mid \varepsilon \\ C &\rightarrow cC \mid \varepsilon \end{aligned}$$

- $\text{FIRST}(C) = \{c, \varepsilon\}$ as $C \Rightarrow cC$ ja $C \Rightarrow \varepsilon$
- $\text{FIRST}(A) = \{a, \varepsilon\}$ as $A \Rightarrow aA$ ja $A \Rightarrow \varepsilon$
- $\text{FIRST}(S) = \{a, b, c, d\}$ as
  - $S \Rightarrow Ab \Rightarrow aAb$ and $S \Rightarrow Ab \Rightarrow b$
  - $S \Rightarrow Cd \Rightarrow cCd$ and $S \Rightarrow Cd \Rightarrow d$

---

**Example:**

In the grammar

$$\begin{aligned} S &\rightarrow ASa \mid b \\ A &\rightarrow BB \mid dA \\ B &\rightarrow b \mid \varepsilon \end{aligned}$$

we have

- $\text{FIRST}(S) = \{b, d\}$ as $S \Rightarrow b$ and $S \Rightarrow ASa \Rightarrow dASa$
- $\text{FIRST}(A) = \{b, d, \varepsilon\}$ koska $A \Rightarrow BB \Rightarrow bB$, $A \Rightarrow dA$ and $A \Rightarrow BB \Rightarrow B \Rightarrow \varepsilon$
- $\text{FIRST}(B) = \{b, \varepsilon\}$ as $B \Rightarrow b$ and $B \Rightarrow \varepsilon$.

---

FIRST-sets (for both terminal symbols and non-terminal variables) can be computed inductively:

- If $a$ is a terminal symbol (i.e. $a \in \Sigma$), then $\text{FIRST}(a) = \{a\}$
- If $X \rightarrow \varepsilon$ is a production, then $\varepsilon \in \text{FIRST}(X)$
- If $X \rightarrow X_1 X_2 ... X_k$ is a production, a terminal symbol $a \in \text{FIRST}(X_i)$ for some $1 \le i \le k$ and $\varepsilon \in \text{FIRST}(X_j)$ for all $1 \le j < i$, then $a \in \text{FIRST}(X)$
- If $X \rightarrow X_1 X_2 ... X_k$ is a production and $\varepsilon \in \text{FIRST}(X_j)$ for all $1 \le j \le k$, then $\varepsilon \in \text{FIRST}(X)$

It holds that $\varepsilon \in \text{FIRST}(A)$ if and only if $A$ is nullable.

---

**Example:**

For the grammar

$$\begin{aligned} S &\rightarrow Ab \mid Cd \\ A &\rightarrow aA \mid \varepsilon \\ C &\rightarrow cC \mid \varepsilon \end{aligned}$$

- $\text{FIRST}(a) = \{a\}$, $\text{FIRST}(b) = \{b\}$, $\text{FIRST}(c) = \{c\}$, $\text{FIRST}(d) = \{d\}$
- $\varepsilon \in \text{FIRST}(C)$ as $C \rightarrow \varepsilon$ is a production
- $c \in \text{FIRST}(C)$ as $C \rightarrow cC$ is a production
- $\varepsilon \in \text{FIRST}(A)$ as $A \rightarrow \varepsilon$ is a production
- $a \in \text{FIRST}(A)$ as $A \rightarrow aA$ is a production
- $a \in \text{FIRST}(S)$ as $S \rightarrow Ab$ is a production and $a \in \text{FIRST}(A)$
- $b \in \text{FIRST}(S)$ as $S \rightarrow Ab$ is a production, $\varepsilon \in \text{FIRST}(A)$ and $b \in \text{FIRST}(b)$
- $c, d \in \text{FIRST}(S)$ with similar argumentation

- We expand FIRST to strings over $V \cup \Sigma$ so that we can study which symbols can occur as first ones when deriving strings from the right hand sides of productions
- Let us define this expansion inductively: $\text{FIRST}(X_1...X_k)$ is the smallest subset of $\Sigma \cup \{\epsilon\}$ for which the following conditions hold:
  - $\epsilon \in \text{FIRST}(\epsilon)$
  - $a \in \text{FIRST}(a)$ for each $a \in \Sigma$
  - If $x \in \Sigma$, $x \in \text{FIRST}(X_i)$ for some $1 \le i \le k$ and $\epsilon \in \text{FIRST}(X_j)$ for all $1 \le j < i$, then $x \in \text{FIRST}(X_1...X_k)$
  - If $\epsilon \in \text{FIRST}(X_j)$ for all $1 \le j \le k$, then $\epsilon \in \text{FIRST}(X_1...X_k)$

## Example:

Consider again the grammar

$$
\begin{aligned}
S &\rightarrow Ab \mid Cd \\
A &\rightarrow aA \mid \epsilon \\
C &\rightarrow cC \mid \epsilon
\end{aligned}
$$

Now

- $\text{FIRST}(A) = \{a, \epsilon\}$ and $\text{FIRST}(b) = \{b\}$
- $\text{FIRST}(C) = \{c, \epsilon\}$ and $\text{FIRST}(d) = \{d\}$
- $\text{FIRST}(S) = \{a, b, c, d\}$
- $\text{FIRST}(Ab) = \{a, b\}$
- $\text{FIRST}(Cd) = \{c, d\}$

☞ in the beginning of the parsing, based only on the first symbol in the string, we can decide whether the production $S \rightarrow Ab$ or $S \rightarrow Cd$ should be applied

## Example:

For the grammar

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

we have

- $\text{FIRST}(F) = \{a, (\}$
- $\text{FIRST}(T) = \{a, (\}$
- $\text{FIRST}(E) = \{a, (\}$
- $\text{FIRST}(E + T) = \{a, (\}$

☞ by applying either $E \Rightarrow E + T$ or $E \Rightarrow T$ we can get $a$ to be the first terminal symbol in the derived string

☞ based on the first symbol in the string only, the parser cannot decide which production should be used

## Example:

Grammar:

$$
\begin{aligned}
E &\rightarrow TE' \\
E' &\rightarrow +TE' \mid \epsilon \\
T &\rightarrow FT' \\
T' &\rightarrow *FT' \mid \epsilon \\
F &\rightarrow a \mid (E)
\end{aligned}
$$

FIRST-sets:

- $\text{FIRST}(E) = \{a, (\}$
- $\text{FIRST}(E') = \{+, \epsilon\}$
- $\text{FIRST}(T) = \{a, (\}$
- $\text{FIRST}(T') = \{*, \epsilon\}$
- $\text{FIRST}(F) = \{a, (\}$

Consider parsing the string $a * a + a$

$$
\underline{E} \Rightarrow \underline{T}E' \Rightarrow \underline{F}T'E' \Rightarrow a\underline{T'}E'
$$

Should we now use the production $T' \Rightarrow *FT'$ or $T' \Rightarrow \epsilon$? Why?

## Slide 1

**Example:**

A small "if-then-else" grammar after left factoring:

$$S \rightarrow s \mid \text{if } C \text{ then } S\,S'$$
$$S' \rightarrow \text{else } S \mid \varepsilon$$
$$C \rightarrow c$$

Making leftmost derivation for the string **if** $c$ **then if** $c$ **then** $s$ **else** $s$:

$$\underline{S} \Rightarrow \text{if } \underline{C} \text{ then } S\,S'$$
$$\Rightarrow \text{if } c \text{ then } \underline{S}\,S'$$
$$\Rightarrow \text{if } c \text{ then if } \underline{C} \text{ then } S\,S'\,S'$$
$$\Rightarrow \text{if } c \text{ then if } c \text{ then } \underline{S}\,S'\,S'$$
$$\Rightarrow \text{if } c \text{ then if } c \text{ then } s\,\underline{S'}\,S'$$

Should we now use the production $S' \rightarrow \text{else } S$ or $S' \rightarrow \varepsilon$?

**FIRST-sets:**

- $\text{FIRST}(S) = \{s, \textbf{if}\}$
- $\text{FIRST}(S') = \{\textbf{else}, \varepsilon\}$
- $\text{FIRST}(C) = \{c\}$
- $\text{FIRST}(\textbf{else } S) = \{\textbf{else}\}$
- $\text{FIRST}(\varepsilon) = \{\varepsilon\}$

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
61/69

## Slide 2

### FOLLOW-sets

- For nullable productions the FIRST-set includes the symbol $\varepsilon$
- How should we interpret this when deciding which production to take next?
- Let us define for each non-terminal variable $A$ the set $\text{FOLLOW}(A)$ of terminal symbols (incl. a special symbol \$ describing the end of the string) that *may follow* $A$ in some derivation:
  - $c \in \text{FOLLOW}(A)$ if $c \in \Sigma$ and $S \Rightarrow^* \alpha A c \beta$ for some $\alpha, \beta \in (V \cup \Sigma)^\star$
  - $\$ \in \text{FOLLOW}(A)$ if $S \Rightarrow^* \alpha A$ for some $\alpha \in (V \cup \Sigma)^\star$

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
62/69

## Slide 3

**Example:**

Grammar:

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

Now

- $\text{FOLLOW}(E) = \text{FOLLOW}(E') = \{\$, )\}$ as
  - $E$ is the start variable, $E \Rightarrow^* (E)T'E'$ and $E \rightarrow TE'$
- $\text{FOLLOW}(T) = \text{FOLLOW}(T') = \{+, \$, )\}$ as
  - $E \Rightarrow TE' \Rightarrow T + TE'$, $E \Rightarrow TE' \Rightarrow T$ and $T \Rightarrow FT'$
- $\text{FOLLOW}(F) = \{+, *, \$, )\}$ as
  - $E \Rightarrow TE' \Rightarrow FT'E' \Rightarrow F * FT'E'$
  - $T \Rightarrow FT' \Rightarrow F$ ("what follows $F$, also follows $T$")

**FIRST-sets:**

- $\text{FIRST}(F) = \{a, (\}$
- $\text{FIRST}(T') = \{*, \varepsilon\}$
- $\text{FIRST}(T) = \{a, (\}$
- $\text{FIRST}(E') = \{+, \varepsilon\}$
- $\text{FIRST}(E) = \{a, (\}$

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
63/69

## Slide 4

**Example:**

Grammar:

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

As

- $\text{FOLLOW}(E) = \text{FOLLOW}(E') = \{\$, )\}$
- $\text{FOLLOW}(T) = \text{FOLLOW}(T') = \{+, \$, )\}$
- $\text{FOLLOW}(F) = \{+, *, \$, )\}$

we know that when making the leftmost derivation for $a * a + a$

$$\underline{E} \Rightarrow \underline{T}E' \Rightarrow \underline{F}T'E' \Rightarrow a\underline{T'}E'$$

we should apply the production $T' \Rightarrow *F T'$ instead of $T' \Rightarrow \varepsilon$ because the non-terminal variable $T'$ cannot be followed by the symbol $*$ in any derivation.

**FIRST-sets:**

- $\text{FIRST}(F) = \{a, (\}$
- $\text{FIRST}(T') = \{*, \varepsilon\}$
- $\text{FIRST}(T) = \{a, (\}$
- $\text{FIRST}(E') = \{+, \varepsilon\}$
- $\text{FIRST}(E) = \{a, (\}$

## Example:

A simple "if-then-else" grammar after left factoring:

$$
\begin{aligned}
S &\;\to\; s \mid \textbf{if}\, C\, \textbf{then}\, S\, S' \\
S' &\;\to\; \textbf{else}\, S \mid \varepsilon \\
C &\;\to\; c
\end{aligned}
$$

Building a leftmost derivation for **if** $c$ **then if** $c$ **then** $s$ **else** $s$:

$$
\begin{aligned}
\underline{S} &\;\Rightarrow\; \textbf{if}\, \underline{C}\, \textbf{then}\, S\, S' \\
&\;\Rightarrow\; ... \\
&\;\Rightarrow\; \textbf{if}\, c\, \textbf{then if}\, c\, \textbf{then}\, s\, \underline{S'}\, S'
\end{aligned}
$$

Now

- **else** $\in \mathrm{FIRST}(\textbf{else}\, S)$
- as well as $S' \to \varepsilon$ and **else** $\in \mathrm{FOLLOW}(S')$
- ☞ based only on the first symbol **else**, one cannot decide whether $S' \to \textbf{else}\, S$ or $S' \to \varepsilon$ should be applied

FOLLOW-sets:

- $\mathrm{FOLLOW}(S) = \{\$, \textbf{else}\}$
- $\mathrm{FOLLOW}(S') = \{\$, \textbf{else}\}$
- $\mathrm{FOLLOW}(C) = \{\textbf{then}\}$

---

### Computing FOLLOW-sets inductively

FOLLOW-sets are the smallest sets that fulfill the following conditions:

- If $S$ is the start variable, then $\$ \in \mathrm{FOLLOW}(S)$
- If $A \to \alpha B \beta$ is a production and a terminal symbol $a \in \mathrm{FIRST}(\beta)$, then $a \in \mathrm{FOLLOW}(B)$
- If $A \to \alpha B$ is a production and $a \in \mathrm{FOLLOW}(A)$, then $a \in \mathrm{FOLLOW}(B)$
- If $A \to \alpha B \beta$ is a production, $\varepsilon \in \mathrm{FIRST}(\beta)$ and $a \in \mathrm{FOLLOW}(A)$, then $a \in \mathrm{FOLLOW}(B)$

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
66/69

---

## Example:

Grammar:

$$
\begin{aligned}
E &\;\to\; T\, E' \\
E' &\;\to\; +\, T\, E' \mid \varepsilon \\
T &\;\to\; F\, T' \\
T' &\;\to\; *\, F\, T' \mid \varepsilon \\
F &\;\to\; a \mid (\, E\, )
\end{aligned}
$$

Now

- $\$ \in \mathrm{FOLLOW}(E)$ as $E$ is the start symbol
- $) \in \mathrm{FOLLOW}(E)$ as $F \to (\, E\, )$ is a production
- $\$, ) \in \mathrm{FOLLOW}(E')$ as $E \to TE'$ is a production
- $+ \in \mathrm{FOLLOW}(T)$ as $E' \to +TE'$ and $+ \in \mathrm{FIRST}(E')$
- and so on...

FIRST-sets:

- $\mathrm{FIRST}(F) = \{a, ($\}$
- $\mathrm{FIRST}(T') = \{*, \varepsilon\}$
- $\mathrm{FIRST}(T) = \{a, ($\}$
- $\mathrm{FIRST}(E') = \{+, \varepsilon\}$
- $\mathrm{FIRST}(E) = \{a, ($\}$

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
67/69

---

Finally, we build a two-dimensional *parsing table* $M$, where each set-valued cell $M(A, a)$ includes all those productions that can be applied when the current non-terminal variable is $A$ and the next input string symbol is $a$:

- If $A \to \alpha$ is a production and the terminal symbol $a \in \mathrm{FIRST}(\alpha)$, then $A \to \alpha \in M(A, a)$
- If $A \to \alpha$ is a production, $\varepsilon \in \mathrm{FIRST}(\alpha)$ and $b \in \mathrm{FOLLOW}(A)$, then $A \to \alpha \in M(A, b)$

### Definition 6.3

A grammar is an LL(1) grammar if its parsing table contains at most production in each cell.

**Aalto University**
School of Science

CS-C2160 Theory of Computation / Lecture 6
Aalto University / Dept. Computer Science
68/69

Let us consider again the grammar

$$
\begin{aligned}
E &\rightarrow T E' \\
E' &\rightarrow + T E' \mid \varepsilon \\
T &\rightarrow F T' \\
T' &\rightarrow * F T' \mid \varepsilon \\
F &\rightarrow a \mid (E)
\end{aligned}
$$

The parsing table is

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

**Aalto University**
**School of Science**

**CS-C2160 Theory of Computation / Lecture 6**
Aalto University / Dept. Computer Science
69/69