# Top-down Syntax Analysis

Reinhard Wilhelm, Sebastian Hack, Mooly Sagiv
Saarland University, Tel Aviv University

W2015

Saarland University, Computer Science

# Top-Down Syntax Analysis

input: A sequence of symbols (tokens)

output: A syntax tree or an error message

- Read input from left to right

- Construct the syntax tree in a top-down manner starting with a node labeled with the start symbol

- **until** input accepted (or error) **do**
  - Predict expansion for the actual leftmost nonterminal (maybe using some lookahead into the remaining input) or
  - Verify predicted terminal symbol against next symbol of the remaining input

- Finds leftmost derivations

# Grammar for Arithmetic Expressions

Left factored grammar $G_2$, i.e. left recursion removed.

$$S \rightarrow E$$

$E \rightarrow TE'$      $E$ generates $T$ with a continuation $E'$

$E' \rightarrow +E|\epsilon$    $E'$ generates possibly empty sequence of $+T$s

$T \rightarrow FT'$      $T$ generates $F$ with a continuation $T'$

$T' \rightarrow *T|\epsilon$    $T'$ generates possibly empty sequence of $*F$s

$F \rightarrow \mathbf{id}|(E)$

$G_2$ defines the same language as $G_0$ und $G_1$.

# Recursive Descent Parsing

- parser is a program,

- a procedure $X$ for each non-terminal $X$,
    - parses words for non-terminal $X$,
    - starts with the first symbol read (into variable *nextsym*),
    - ends with the following symbol read (into variable *nextsym*).

- uses one symbol lookahead into the remaining input.

- uses the **FiFo** sets to make the expansion transitions deterministic

$$\textbf{FiFo}(N \rightarrow \alpha) = \textit{FIRST}_1(\alpha) \oplus_1 \textit{FOLLOW}_1(N)$$

# The $FIRST_1$ Sets

- A production $N \rightarrow \alpha$ is applicable for symbols that "begin" $\alpha$

- Example: Arithmetic Expressions, Grammar $G_2$
  - The production $F \rightarrow id$ is applied when the current symbol is **id**
  - The production $F \rightarrow (E)$ is applied when the current symbol is **(**
  - The production $T \rightarrow F$ is applied when the current symbol is **id** or **(**

- Formal definition:

$$FIRST_1(\alpha) = \{1 : w \mid \alpha \stackrel{*}{\Longrightarrow} w, w \in V_T^*\}$$

# The $FOLLOW_1$ Sets

- A production $N \to \epsilon$ is applicable for symbols that "can follow" $N$ in some derivation

- Example: Arithmetic Expressions, Grammar $G_2$
  - The production $E' \to \epsilon$ is applied for symbols $\#$ and $)$
  - The production $T' \to \epsilon$ is applied for symbols $\#$, $)$ and $+$

- Formal definition:

$$FOLLOW_1(N) = \{a \in V_T | \exists \alpha, \gamma : S \stackrel{*}{\Longrightarrow} \alpha N a \gamma\}$$

## Definitions

Let $k \geq 1$

- $k$-prefix of a word $w = a_1 \ldots a_n$

$$k : w = \begin{cases} a_1 \ldots a_n & \text{if } n \leq k \\ a_1 \ldots a_k & \text{otherwise} \end{cases}$$

- $k$-concatenation

$\oplus_k : V^* \times V^* \to V^{\leq k}$, defined by $u \oplus_k v = k : uv$
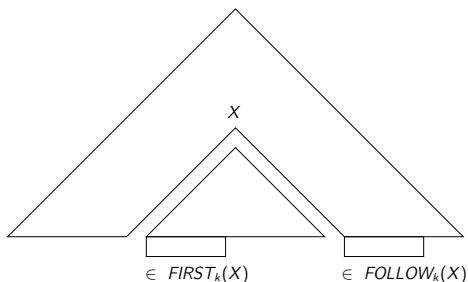
- extended to languages

$$k : L = \{k : w \mid w \in L\}$$

$$L_1 \oplus_k L_2 = \{x \oplus_k y \mid x \in L_1, y \in L_2\}$$

$$V^{\leq k} = \bigcup_{i=1}^{k} V^i \text{ set of words of length at most } k$$

$$V_{T\#}^{\leq k} = V_T^{\leq k} \cup V_T^{k-1}\{\#\} \ldots \text{ possibly terminated by } \#.$$

# $FIRST_k$ and $FOLLOW_k$



■ set of $k$–prefixes of terminal words for $\alpha$

$$FIRST_k : (V_N \cup V_T)^* \to 2^{V_T^{\leq k}}$$
$$FIRST_k(\alpha) = \{ k : u \mid \alpha \overset{*}{\Longrightarrow} u \}$$

■ set of $k$–prefixes of terminal words that may immediately follow $X$

$$FOLLOW_k : V_N \to 2^{V_{T\#}^{\leq k}}$$
$$FOLLOW_k(X) = \{ w \mid S \overset{*}{\Longrightarrow} \beta X \gamma \text{ and } w \in FIRST_k(\gamma) \}$$

## Parser for $G_2$

```
program parser;
var  nextsym: string;
proc scan;
 {reads next input symbol into nextsym}
proc error (message: string);
 {issues error message and stops parser}
proc accept; {terminates successfully}

proc S;
 begin E
 end ;

proc E;
 begin T; E'
 end ;
```

```
proc E';
 begin
   case nextsym in
    {"+"}: if nextsym = "+ "
            then scan
            else error( "+ expected") fi ; E;
    otherwise ;
   endcase
 end ;

proc T;
 begin F; T' end ;

proc T';
 begin
   case nextsym in
    {" *"}: if nextsym = "*"
            then scan
            else error( "* expected") fi ; T;
    otherwise ;
   endcase
 end ;
```

```
proc F;
 begin
   case nextsym in
     {"("}:      if nextsym = "("
                 then scan
                 else error( "( expected") fi ; E;
                 if nextsym = ")"
                 then scan else error(" ) expected") fi;
     otherwise if nextsym ="id"
                 then scan else error("id expected") fi;
   endcase
 end ;
begin
scan; S;
if nextsym = "#" then accept else error("# expected") fi
end .
```

# How to Construct such a Parser Program

- Code was automatically generated from the grammar and the **FiFo** sets.

- The program generating the parser has the functions:

$$
\begin{array}{lcll}
\text{N\_prog} & : & V_N \to \text{code} & \text{nonterminals} \\
\text{C\_prog} & : & (V_N \cup V_T)^* \to \text{code} & \text{concantenations} \\
\text{S\_prog} & : & V_N \cup V_T \to \text{code} & \text{symbols}
\end{array}
$$

## Parser Schema

```
program parser;
    var  nextsym: symbol;
    proc  scan;
        (* reads next input symbol into  nextsym *)
    proc  error (message: string);
        (* issues error message and stops the parser *)
    proc  accept;
        (* terminates parser successfully *)

    N_prog(X_0);                              (* X_0 start symbol *)
    N_prog(X_1);
        .
        .
        .
    N_prog(X_n);
```

```
begin
    scan;
    X_0;
    if nextsym = "#"
        then  accept
        else  error("...")
    fi
end
```

## The Non-terminal Procedures

N = Non-terminal, C = Concatenation, S = Symbol

$N\_prog(X) =$                                    $(* \; X \to \alpha_1|\alpha_2|\cdots|\alpha_{k-1}|\alpha_k \; *)$

       **proc** X;
       **begin**
       **case** nextsym **in**
       $FiFo(X \to \alpha_1):$      $C\_progr(\alpha_1);$
       $FiFo(X \to \alpha_2):$      $C\_progr(\alpha_2);$
               $\vdots$
       $FiFo(X \to \alpha_{k-1}):$    $C\_progr(\alpha_{k-1});$
       **otherwise** $C\_progr(\alpha_k);$
       **endcase**
       **end** ;

$$C\_progr(\alpha_1\alpha_2\cdots\alpha_k) =$$
$$\quad S\_progr(\alpha_1); S\_progr(\alpha_2); \ldots S\_progr(\alpha_k);$$
$$S\_progr(a) =$$
$$\quad \textbf{if } nextsym = a \textbf{ then } scan$$
$$\quad \textbf{else } error (\ "a\ expected")$$
$$\quad \textbf{fi}$$
$$S\_progr(Y) = Y$$

FiFo–sets have to be disjoint (LL(1)–grammar)

# A Generative Solution

Generate the control of a deterministic PDA from the grammar and the **FiFo** sets.

- At compiler–generation time construct a table $M$
  $M\colon V_N \times V_T \to P$
  $M[N, a]$ is the production used to expand nonterminal $N$ when the current symbol is $a$

- For some grammars report that the table cannot be constructed
  The compiler writer can then decide to:
  - change the grammar (but not the language)
  - use a more general parser-generator
  - "Patch" the table (manually or using some rules)

# Creating the table

> Input: cfg $G$, $FIRST_1$ und $FOLLOW_1$ for $G$.
>
> Output: The parsing table $M$ or an indication that such a table cannot be constructed

$M$ is constructed as follows:

- For all $X \to \alpha \in P$ and $a \in FIRST_1(\alpha)$, set $M[X, a] = (X \to \alpha)$

- If $\varepsilon \in FIRST_1(\alpha)$, for all $b \in FOLLOW_1(X)$, set $M[X, b] = (X \to \alpha)$

- Set all other entries of $M$ to *error*

Parser table cannot be constructed if at least one entry is set twice.
Then, $G$ is not LL(1)

# Example – arithmetic expressions

| nonterminal | symbol | Production |
|---|---|---|
| $S$ | $($, $id$ | $S \to E$ |
| $S$ | $+, *, ), \#$ | error |
| $E$ | $($, $id$ | $E \to TE'$ |
| $E$ | $+, *, ), \#$ | error |
| $E'$ | $+$ | $E' \to +E$ |
| $E'$ | $), \#$ | $E' \to \epsilon$ |
| $E'$ | $($, $*, id$ | error |
| $T$ | $($, $id$ | $T \to FT'$ |
| $T$ | $+, *, ), \#$ | error |
| $T'$ | $*$ | $T' \to *T$ |
| $T'$ | $+, ), \#$ | $T' \to \epsilon$ |
| $T'$ | $($, $id$ | error |
| $F$ | $id$ | $F \to id$ |
| $F$ | $($ | $F \to (E)$ |
| $F$ | $+, *, )$ | error |

# LL-Parser Driver (interprets the table $M$)

```
program parser;
    var  nextsym: symbol;
    var  st: stack of item;
    proc  scan;
        (∗ reads next input symbol into  nextsym ∗)
    proc  error (message: string);
        (∗ issues error message and stops the parser ∗)
    proc  accept;
        (∗ terminates parser successfully ∗)
    proc reduce;
        (∗ replaces [X → β.Yγ][Y → α.] by [X → βY.γ]  ∗)
    proc  pop;
        (∗ removes topmost item from st ∗)
    proc  push ( i : item);
        (∗ pushes i onto st ∗)
    proc  replaceby ( i: item);
        (∗ replaces topmost item of st by i ∗)
```

```
begin
    scan; push( [S' → .S] );
    while nextsym ≠ "#" do
    case top in
    [X → β.aγ]:     if nextsym = a
                    then scan; replaceby([X → βa.γ])
                    else error fi ;
    [X → β.Yγ] :    if M[Y, nextsym] = (Y → α)
                    then push([Y → .α])
                    else error fi ;
    [X → α.]:       reduce;
    [S' → S.] :     if nextsym = "#" then accept
                    else error fi
    endcase
    od
end .
```

# Explicit Stack
## Deterministic Pushdown Automaton



$[X \to \alpha.Y\beta]$

$\rho$

$\#$

Stack

w | a | v

Input

Control

Parser–Table

$M$

tree

Output

# LL($k$)-grammar

Goal: formalizing our intuition when the expand-transitions of the Item-Pushdown-Automaton can be made deterministic.

Means: $k$-symbol lookahead into the remaining input.

# LL($k$)-grammar

- Let $G = (V_N, V_T, P, S)$ be a cfg and $k$ be a natural number.
  $G$ is an **LL($k$)-grammar** iff the following holds:

  if there exist two leftmost derivations

  $S \xRightarrow[lm]{*} uY\alpha \xRightarrow[lm]{} u\beta\alpha \xRightarrow[lm]{*} ux$ and

  $S \xRightarrow[lm]{*} uY\alpha \xRightarrow[lm]{} u\gamma\alpha \xRightarrow[lm]{*} uy$, and if $k : x = k : y$,

  then $\beta = \gamma$.

- The expansion of the leftmost non-terminal is always uniquely determined by
  - the consumed part of the input and
  - the next $k$ symbols of the remaining input

## Example 1

Let $G_1$ be the cfg with the productions

$$
\begin{aligned}
STAT \rightarrow \quad & \textbf{if id then } STAT \textbf{ else } STAT \textbf{ fi} \quad | \\
& \textbf{while id do } STAT \textbf{ od} \quad | \\
& \textbf{begin } STAT \textbf{ end} \quad | \\
& \textbf{id} := \textbf{id}
\end{aligned}
$$

## Example 1

Let $G_1$ be the cfg with the productions

$$
\begin{aligned}
STAT \rightarrow \quad & \textbf{if id then } STAT \textbf{ else } STAT \textbf{ fi} \quad | \\
& \textbf{while id do } STAT \textbf{ od} \quad | \\
& \textbf{begin } STAT \textbf{ end} \quad | \\
& \textbf{id} := \textbf{id}
\end{aligned}
$$

$G_1$ is an LL(1)-grammar.

$$
STAT \; \underset{lm}{\overset{*}{\Longrightarrow}} \; w\, STAT\, \alpha \; \underset{lm}{\Longrightarrow} \; w\, \beta\, \alpha \; \underset{lm}{\overset{*}{\Longrightarrow}} \; w\, x
$$

$$
STAT \; \underset{lm}{\overset{*}{\Longrightarrow}} \; w\, STAT\, \alpha \; \underset{lm}{\Longrightarrow} \; w\, \gamma\, \alpha \; \underset{lm}{\overset{*}{\Longrightarrow}} \; w\, y
$$

From $1 : x = 1 : y$ follows $\beta = \gamma$,
e.g., from $1 : x = 1 : y = \textbf{if}$ follows
$\beta = \gamma = "\textbf{if id then } STAT \textbf{ else } STAT \textbf{ fi}"$

# Example 2

Let $G_2$ be the cfg with the productions

$STAT \rightarrow$ **if** id **then** $STAT$ **else** $STAT$ **fi** |
**while** id **do** $STAT$ **od** |
**begin** $STAT$ **end** |
**id** := **id** |
**id**: $STAT$ |               (∗ labeled statem. ∗)
**id** (**id** )              (∗ procedure call ∗)

## Example 2 (cont'd)

$G_2$ is not an LL(1)–grammar.

$$STAT \underset{lm}{\overset{*}{\Longrightarrow}} w \ STAT \ \alpha \underset{lm}{\Longrightarrow} w \ \overbrace{\textbf{id} := \textbf{id}}^{\beta} \ \alpha \underset{lm}{\overset{*}{\Longrightarrow}} w \ x$$

$$STAT \underset{lm}{\overset{*}{\Longrightarrow}} w \ STAT \ \alpha \underset{lm}{\Longrightarrow} w \ \overbrace{\textbf{id} : STAT}^{\gamma} \ \alpha \underset{lm}{\overset{*}{\Longrightarrow}} w \ y$$

$$STAT \underset{lm}{\overset{*}{\Longrightarrow}} w \ STAT \ \alpha \underset{lm}{\Longrightarrow} w \ \overbrace{\textbf{id}(\textbf{id})}^{\delta} \ \alpha \underset{lm}{\overset{*}{\Longrightarrow}} w \ z$$

and $1 : x = 1 : y = 1 : z = $"**id**",
and $\beta$, $\gamma$, $\delta$ are pairwise different.
$G_2$ is an LL(2)–grammar.
$2 : x = $"**id** :=", $2 : y = $"**id** :", $2 : z = $"**id**(" are pairwise different.

# Example 3

Let $G_3$ have the productions

$$
\begin{aligned}
STAT \;\; \rightarrow \;\; & \textbf{if} \;\; \textbf{id} \;\textbf{then} \;\; STAT \;\textbf{else} \;\; STAT \;\textbf{fi} && | \\
& \textbf{while} \;\; \textbf{id} \;\textbf{do} \; STAT \;\textbf{od} && | \\
& \textbf{begin} \;\; STAT \;\textbf{end} && | \\
& VAR := VAR && | \\
& \textbf{id}(\; IDLIST\,) && (\ast \text{ procedure call } \ast) \\
VAR \;\; \rightarrow \;\; & \textbf{id} \;\; | \;\; \textbf{id}\,(IDLIST\,) && (\ast \text{ indexed variable } \ast) \\
IDLIST \;\; \rightarrow \;\; & \textbf{id} \;\; | \;\; \textbf{id},\, IDLIST
\end{aligned}
$$

# Example 3

Let $G_3$ have the productions

$$
\begin{array}{lll}
STAT & \rightarrow & \textbf{if} \ \textbf{id} \ \textbf{then} \ STAT \ \textbf{else} \ STAT \ \textbf{fi} \qquad | \\
 & & \textbf{while} \ \textbf{id} \ \textbf{do} \ STAT \ \textbf{od} \qquad | \\
 & & \textbf{begin} \ STAT \ \textbf{end} \qquad | \\
 & & VAR := VAR \qquad | \\
 & & \textbf{id}(IDLIST) \qquad\qquad\qquad (* \ \text{procedure call} \ *) \\
VAR & \rightarrow & \textbf{id} \ | \ \textbf{id}(IDLIST) \qquad\qquad (* \ \text{indexed variable} \ *) \\
IDLIST & \rightarrow & \textbf{id} \ | \ \textbf{id}, IDLIST \\
\end{array}
$$

$G_3$ is not an LL($k$)–grammar for any $k$.

## Proof

Assume $G_3$ to be $LL(k)$ for a $k > 0$.

Let $STAT \Rightarrow \beta \overset{*}{\underset{lm}{\Longrightarrow}} x$ and $STAT \Rightarrow \gamma \overset{*}{\underset{lm}{\Longrightarrow}} y$ with

$$x = \mathbf{id}\,(\underbrace{\mathbf{id}, \mathbf{id}, \ldots, \mathbf{id}}_{\lceil \frac{k}{2} \rceil \text{ times}}) := \mathbf{id} \text{ and } y = \mathbf{id}\,(\underbrace{\mathbf{id}, \mathbf{id}, \ldots, \mathbf{id}}_{\lceil \frac{k}{2} \rceil \text{ times}})$$

Then $k : x = k : y$,
but $\beta = "VAR := VAR" \neq \gamma = "\mathbf{id}\,(IDLIST)"$.

# Transforming to LL($k$)

Factorization creates an LL(2)–grammar, equivalent to $G_3$.

The productions

$$STAT \rightarrow VAR := VAR \mid \mathbf{id}(IDLIST)$$

are replaced by

$$
\begin{aligned}
STAT &\rightarrow ASSPROC \mid \mathbf{id} := VAR \\
ASSPROC &\rightarrow \mathbf{id}(IDLIST)\ APREST \\
APREST &\rightarrow\ := VAR \mid \varepsilon
\end{aligned}
$$

# A non–LL($k$)–language

Let

$$G_4 = (\{S, A, B\}, \{0, 1, a, b\}, P_4, S)$$

$$P_4 = \left\{ \begin{array}{ccc} S & \rightarrow & A \mid B \\ A & \rightarrow & aAb \mid 0 \\ B & \rightarrow & aBbb \mid 1 \end{array} \right\}$$

$$L(G_4) = \{a^n 0 b^n \mid n \geq 0\} \cup \{a^n 1 b^{2n} \mid n \geq 0\}$$

$G_4$ is not LL($k$) for any $k$. Consider the two leftmost derivations

$$S \xRightarrow[lm]{0} S \xRightarrow[lm]{} A \xRightarrow[lm]{*} a^k 0 b^k$$

$$S \xRightarrow[lm]{0} S \xRightarrow[lm]{} B \xRightarrow[lm]{*} a^k 1 b^{2k}$$

With $u = \alpha = \varepsilon$, $\beta = A$, $\gamma = B$, $x = "a^k 0 b^k"$, $y = "a^k 1 b^{2k}"$ it holds
$k : x = k : y$, but $\beta \neq \gamma$.
Since $k$ can be chosen arbitrarily, we have $G_4$ is not LL($k$) for any $k$.
There even is no LL($k$)-grammar for $L(G_4)$ for any $k$.

# LL($k$)–conditions

**Theorem**

*G is LL(1) iff for different productions $A \to \beta$ and $A \to \gamma$*
*$FIRST_1(\beta) \oplus_1 FOLLOW_1(A) \cap FIRST_1(\gamma) \oplus_1 FOLLOW_1(A) = \emptyset$*

**Corollary**

*G is LL(1) iff for all alternatives $A \to \alpha_1 | \ldots | \alpha_n$:*

1. *$FIRST_1(\alpha_1), \ldots, FIRST_1(\alpha_n)$ are pairwise disjoint; in particular, at most one of them may contain $\varepsilon$*

2. *$\alpha_i \stackrel{*}{\Longrightarrow} \varepsilon$ implies:*

   *$FIRST_1(\alpha_j) \cap FOLLOW_1(A) = \emptyset$ for $1 \leq j \leq n$, $j \neq i$.*

The Theorem was used in the parser construction!

# Further Definitions and Theorems

- $G$ is called a strong LL(k)-grammar if for each two different productions $A \to \beta$ and $A \to \gamma$

$$FIRST_k(\beta) \oplus_k FOLLOW_k(A) \cap FIRST_k(\gamma) \oplus_k FOLLOW_k(A) = \emptyset$$

- A production is called directly left recursive
  if it has the form $A \to A\alpha$

- A non-terminal $A$ is called left recursive if it has a derivation
  $A \overset{+}{\Longrightarrow} A\alpha$.

- A cfg $G$ is called left recursive
  if $G$ contains at least one left recursive non-terminal

**Theorem**

    (a)  *G is not LL(k) for any k if G is left recursive.*

    (b)  *G is not ambiguous if G is LL(k)-grammar.*