

Top-down Syntax Analysis

Sebastian Hack

(based on slides by Reinhard Wilhelm and Mooly Sagiv)

<http://compilers.cs.uni-saarland.de>

Compiler Construction Core Course 2017
Saarland University

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.

$$S \rightarrow E$$

$$E \rightarrow TE' \quad E \text{ generates } T \text{ with a continuation } E'$$

$$E' \rightarrow +E|\epsilon \quad E' \text{ generates possibly empty sequence of } +Ts$$

$$T \rightarrow FT' \quad T \text{ generates } F \text{ with a continuation } T'$$

$$T' \rightarrow *T|\epsilon \quad T' \text{ generates possibly empty sequence of } *Fs$$

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

G_2 defines the same language as G_0 und G_1 .

Grammar for Arithmetic Expressions

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

$$S \rightarrow E$$

$$E \rightarrow TE' \quad E \text{ generates } T \text{ with a continuation } E'$$

$$E' \rightarrow +E|\epsilon \quad E' \text{ generates possibly empty sequence of } +Ts$$

$$T \rightarrow FT' \quad T \text{ generates } F \text{ with a continuation } T'$$

$$T' \rightarrow *T|\epsilon \quad T' \text{ generates possibly empty sequence of } *Fs$$

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

G_2 defines the same language as G_0 and G_1 .

But the parse tree is not so suitable as an abstract syntax tree!

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

$$\mathbf{FiFo}(N \rightarrow \alpha) = \mathit{FIRST}_1(\alpha) \oplus_1 \mathit{FOLLOW}_1(N)$$

The $FIRST_1$ Sets

- A production $N \rightarrow \alpha$ is applicable for symbols that “begin” α
- 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:

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

The $FOLLOW_1$ Sets

- A production $N \rightarrow \epsilon$ is applicable for symbols that “can follow” N in some derivation
- Example: Arithmetic Expressions, Grammar G_2
 - The production $E' \rightarrow \epsilon$ is applied for symbols $\#$ and $)$
 - The production $T' \rightarrow \epsilon$ is applied for symbols $\#,)$ and $+$
- Formal definition:

$$FOLLOW_1(N) = \{a \in V_T \mid \exists \alpha, \gamma : S \xRightarrow{*} \alpha N a \gamma\}$$

Definitions

Let $k \geq 1$

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

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

- k -concatenation

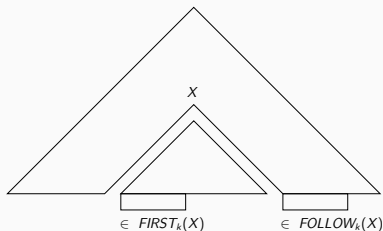
$$\oplus_k : V^* \times V^* \rightarrow V^{\leq k}, \text{ defined by } u \oplus_k v = k : uv$$

- extended to languages

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

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

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



- set of k -prefixes of terminal words for α

$$FIRST_k : (V_N \cup V_T)^* \rightarrow 2^{V_T^{\leq k}}$$

$$FIRST_k(\alpha) = \{k : u \mid \alpha \xRightarrow{*} u\}$$

- set of k -prefixes of terminal words that may immediately follow X

$$FOLLOW_k : V_N \rightarrow 2^{V_T^{\leq k} \#}$$

$$FOLLOW_k(X) = \{w \mid S \xRightarrow{*} \beta X \gamma \text{ and } w \in FIRST_k(\gamma)\}$$

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:

N_prog	:	$V_N \rightarrow$	code	nonterminals
C_prog	:	$(V_N \cup V_T)^*$	\rightarrow code	concatenations
S_prog	:	$V_N \cup V_T$	\rightarrow code	symbols

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 →  $\alpha_1|\alpha_2|\dots|\alpha_{k-1}|\alpha_k$  *)  
  proc X;  
  begin  
    case nextsym in  
      FiFo(X →  $\alpha_1$ ) : C_progr( $\alpha_1$ );  
      FiFo(X →  $\alpha_2$ ) : C_progr( $\alpha_2$ );  
      ⋮  
      FiFo(X →  $\alpha_{k-1}$ ) : C_progr( $\alpha_{k-1}$ );  
    otherwise C_progr( $\alpha_k$ );  
  endcase  
  end ;
```



```

C_progr( $\alpha_1\alpha_2 \cdots \alpha_k$ ) =
    S_progr( $\alpha_1$ ); S_progr( $\alpha_2$ ); ... S_progr( $\alpha_k$ );
S_progr(a) =
    if nextsym = a then scan
    else error ( "a expected" )
    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
 $M: V_N \times V_T \rightarrow P$
 $M[N, a]$ 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 $X \rightarrow \alpha \in P$ and $a \in FIRST_1(\alpha)$, set $M[X, a] = (X \rightarrow \alpha)$
- If $\varepsilon \in FIRST_1(\alpha)$, for all $b \in FOLLOW_1(X)$, set $M[X, b] = (X \rightarrow \alpha)$
- 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	$S \rightarrow E$
S	+, *,), #	error
E	(, id	$E \rightarrow TE'$
E	+, *,), #	error
E'	+	$E' \rightarrow +E$
E'), #	$E' \rightarrow \epsilon$
E'	(, *, id	error
T	(, id	$T \rightarrow FT'$
T	+, *,), #	error
T'	*	$T' \rightarrow *T$
T'	+,), #	$T' \rightarrow \epsilon$
T'	(, id	error
F	id	$F \rightarrow id$
F	($F \rightarrow (E)$
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  $[X \rightarrow \beta.Y\gamma][Y \rightarrow \alpha.]$  by  $[X \rightarrow \beta Y.\gamma]$  *)
  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($[S' \rightarrow \cdot S]$);

while nextsym \neq "#" **do**

case top **in**

$[X \rightarrow \beta \cdot a \gamma]$: **if** nextsym = a
then scan; replaceby($[X \rightarrow \beta a \cdot \gamma]$)
else error fi ;

$[X \rightarrow \beta \cdot Y \gamma]$: **if** $M[Y, \text{nextsym}] = (Y \rightarrow \alpha)$
then push($[Y \rightarrow \cdot \alpha]$)
else error fi ;

$[X \rightarrow \alpha \cdot]$: reduce;

$[S' \rightarrow S \cdot]$: **if** nextsym = "#" **then** accept
else error fi

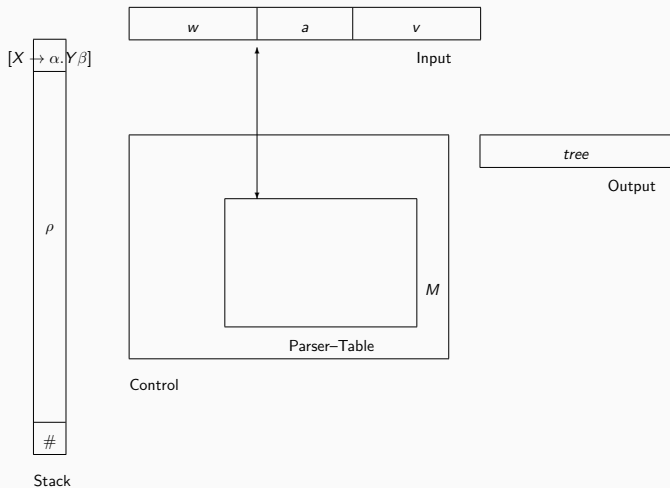
endcase

od

end .

Explicit Stack

Deterministic Pushdown Automaton



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

$$S \xrightarrow[lm]{*} uY\alpha \xrightarrow[lm]{} u\beta\alpha \xrightarrow[lm]{*} ux \quad \text{and}$$

$$S \xrightarrow[lm]{*} uY\alpha \xrightarrow[lm]{} u\gamma\alpha \xrightarrow[lm]{*} uy \quad \text{and}$$

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

```
STAT → if id then STAT else STAT fi |  
      while id do STAT od |  
      begin STAT end |  
      id := id
```

Example 1

Let G_1 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 |

G_1 is an LL(1)-grammar.

$$\begin{array}{l} STAT \xrightarrow[*]{lm} w STAT \alpha \xRightarrow{lm} w \beta \alpha \xrightarrow[*]{lm} w x \\ STAT \xrightarrow[*]{lm} w STAT \alpha \xRightarrow{lm} w \gamma \alpha \xrightarrow[*]{lm} w y \end{array}$$

From $1 : x = 1 : y$ follows $\beta = \gamma$,
e.g., from $1 : x = 1 : y = \mathbf{if}$ follows

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.

$$STAT \xrightarrow[*]{lm} w STAT \alpha \xRightarrow{lm} w \overbrace{\mathbf{id} := \mathbf{id}}^{\beta} \alpha \xrightarrow[*]{lm} w x$$

$$STAT \xrightarrow[*]{lm} w STAT \alpha \xRightarrow{lm} w \overbrace{\mathbf{id} : STAT}^{\gamma} \alpha \xrightarrow[*]{lm} w y$$

$$STAT \xrightarrow[*]{lm} w STAT \alpha \xRightarrow{lm} w \overbrace{\mathbf{id}(\mathbf{id})}^{\delta} \alpha \xrightarrow[*]{lm} w z$$

and $1 : x = 1 : y = 1 : z = \mathbf{id}$,

and β, γ, δ are pairwise different.

G_2 is an LL(2)-grammar.

$2 : x = \mathbf{id} :=$, $2 : y = \mathbf{id} :$, $2 : z = \mathbf{id} ($ are pairwise different.

Example 3

Let G_3 have the productions

<i>STAT</i>	→	if id then <i>STAT</i> else <i>STAT</i> fi	
		while id do <i>STAT</i> od	
		begin <i>STAT</i> end	
		<i>VAR</i> := <i>VAR</i>	
		id(<i>IDLIST</i>)	(* procedure call *)
<i>VAR</i>	→	id id(<i>IDLIST</i>)	(* indexed variable *)
<i>IDLIST</i>	→	id id, <i>IDLIST</i>	

Example 3

Let G_3 have the productions

$STAT$	\rightarrow	if id then $STAT$ else $STAT$ fi	
		while id do $STAT$ od	
		begin $STAT$ end	
		$VAR := VAR$	
		$id(IDLIST)$	(* procedure call *)
VAR	\rightarrow	id $id(IDLIST)$	(* indexed variable *)
$IDLIST$	\rightarrow	id $id, IDLIST$	

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

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

Let $STAT \Rightarrow \beta \xrightarrow[*]{lm} x$ and $STAT \Rightarrow \gamma \xrightarrow[*]{lm} y$ with

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

Then $k : x = k : y$,

but $\beta = \text{"VAR := VAR"} \neq \gamma = \text{"id (IDLIST)"}$.

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

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

A non-LL(k)-language

Let

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

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

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

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 .

Theorem

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

Corollary

G is LL(1) iff for all alternatives $A \rightarrow \alpha_1 | \dots | \alpha_n$:

- $FIRST_1(\alpha_1), \dots, FIRST_1(\alpha_n)$ are pairwise disjoint; in particular, at most one of them may contain ϵ*
- $\alpha_j \xRightarrow{*} \epsilon$ implies:*

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

The Theorem was used in the parser construction!

Further Definitions and Theorems

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

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

- $SLL(1) = LL(1)$
- A production is called **directly left recursive** if it has the form $A \rightarrow A\alpha$
- A non-terminal A is called **left recursive** if it has a derivation $A \xrightarrow{+} A\alpha$.
- 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.*