

Top-down Syntax Analysis

– Wilhelm/Maurer: Compiler Design, Chapter 8 –

Reinhard Wilhelm
Universität des Saarlandes
wilhelm@cs.uni-sb.de
and
Mooly Sagiv
Tel Aviv University
sagiv@math.tau.ac.il

Subjects

- ▶ Functionality and Method
- ▶ Recursive Descent Parsing
- ▶ Using parsing tables
- ▶ Explicit stacks
- ▶ Creating the table
- ▶ $LL(k)$ -grammars
- ▶ Other properties
- ▶ Handling Limitations

Top-Down Syntax Analysis

input: A sequence of symbols (tokens)

output: A syntax tree or an error message

- method**
- ▶ 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.

$$S \rightarrow E$$

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

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

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

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

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

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

$$\begin{aligned}
 \mathbf{FiFo}(N \rightarrow \alpha) &= FIRST_1(\alpha) \oplus_1 FOLLOW_1(N) = \\
 &\begin{cases} FIRST_1(\alpha) \cup FOLLOW_1(N) & \alpha \xRightarrow{*} \epsilon \\ FIRST_1(\alpha) & \text{otherwise} \end{cases}
 \end{aligned}$$

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

Observation: Much redundant code generated. Why this?
Code was automatically generated from the **grammar** and the **FIFO** sets.

Nice application for a **functional programming language!**

Let $G = (V_N, V_T, P, S)$ be a context-free grammar and FiFo be the computed lookahead sets.

The functional program generating the parser would have the functions:

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

$$N_prog(X) = (* X \rightarrow \alpha_1 | \alpha_2 | \dots | \alpha_{k-1} | \alpha_k *)$$

```

proc X;
begin
  case nextsym in
    FiFo(X  $\rightarrow$   $\alpha_1$ ) : C_progr( $\alpha_1$ );
    FiFo(X  $\rightarrow$   $\alpha_2$ ) : C_progr( $\alpha_2$ );
       $\vdots$ 
    FiFo(X  $\rightarrow$   $\alpha_{k-1}$ ) : 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 should 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

Method: 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.
 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($[S' \rightarrow .S]$);

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

case top **in**

$[X \rightarrow \beta.a\gamma]$: **if** nextsym = a

then scan; replaceby($[X \rightarrow \beta a.\gamma]$

else error **fi** ;

$[X \rightarrow \beta.Y\gamma]$: **if** $M[Y, \text{nextsym}] = (Y \rightarrow \alpha)$

then push($[Y \rightarrow .\alpha]$)

else error **fi** ;

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

$[S' \rightarrow S.]$: **if** nextsym = "#"**then** accept

else error **fi**

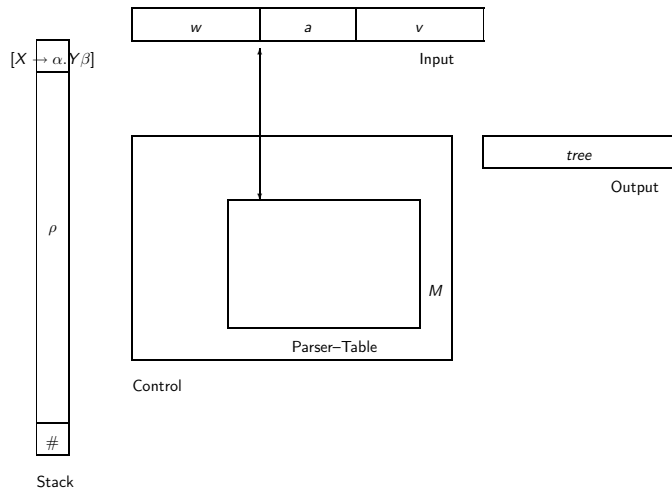
endcase

od

end .

Explicit Stack

Deterministic Pushdown Automaton



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

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

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

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

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

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

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

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

$STAT \rightarrow$ **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
 $\beta = \gamma = \mathbf{"if id then $STAT$ else $STAT$ fi"}$

Example 1

Let G_1 be the cfg with the productions

$$\begin{array}{l} STAT \rightarrow \text{if id then } STAT \text{ else } STAT \text{ fi} \\ \quad \text{while id do } STAT \text{ od} \\ \quad \text{begin } STAT \text{ end} \\ \quad \text{id} := \text{id} \end{array} \quad \left| \begin{array}{l} \\ \\ \\ \end{array} \right.$$

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 $1 : x = 1 : y = \text{if}$ follows
 $\beta = \gamma = \text{"if id then } STAT \text{ else } STAT \text{ fi"}$

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.

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

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

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

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

and β, γ, δ are pairwise different.

G_2 is an LL(2)-grammar.

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

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		
		<i>VAR := VAR</i>		
		<i>id(IDLIST)</i>		(* procedure call *)
VAR	\rightarrow	id <i>id(IDLIST)</i>		(* indexed variable *)
$IDLIST$	\rightarrow	id <i>id, IDLIST</i>		

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

Example 3

Let G_3 have the productions

<p>$STAT \rightarrow$ 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 := VAR</i> <i>id(IDLIST)</i></p> <p>$VAR \rightarrow$ id <i>id(IDLIST)</i></p> <p>$IDLIST \rightarrow$ id <i>id, IDLIST</i></p>	<p> </p> <p> </p> <p> </p> <p> </p> <p>(* procedure call *)</p> <p>(* indexed variable *)</p>
--	---

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

Proof:

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

Let $STAT \Rightarrow \beta \xrightarrow[*]{lm} x$ and $STAT \Rightarrow \gamma \xrightarrow[*]{lm} y$ with
 $x = \text{id} \underbrace{(\text{id}, \text{id}, \dots, \text{id})}_{\lceil \frac{k}{2} \rceil \text{ times}} := \text{id}$ and $y = \text{id} \underbrace{(\text{id}, \text{id}, \dots, \text{id})}_{\lceil \frac{k}{2} \rceil \text{ times}}$

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

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 .

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

Consider the two leftmost derivations

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

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 .

Towards Checkable $LL(k)$ -conditions

Theorem

G is $LL(k)$ -grammar iff the following condition holds:

Are $A \rightarrow \beta$ and $A \rightarrow \gamma$ different productions in P , then

$$FIRST_k(\beta\alpha) \cap FIRST_k(\gamma\alpha) = \emptyset \text{ for all } \alpha \text{ with } S \xrightarrow[lm]{*} wA\alpha$$

Theorem

Let G be a cfg without productions of the form $X \rightarrow \varepsilon$.

G is an $LL(1)$ -grammar iff

for each non-terminal X with the alternatives $X \rightarrow \alpha_1 | \dots | \alpha_n$
the sets $FIRST_1(\alpha_1), \dots, FIRST_1(\alpha_n)$ are pairwise disjoint.

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 $A \rightarrow \alpha_1 | \dots | \alpha_n$:

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

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

The condition of 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 $A \rightarrow \beta$ and $A \rightarrow \gamma$

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

- ▶ A production is called **directly left recursive**, if it has the form $A \rightarrow A\alpha$
- ▶ 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.*

Recursive Descent Parsing

```
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( $X_0 \rightarrow \alpha_0$ );  
  N_prog( $X_1 \rightarrow \alpha_1$ );  
    :  
  N_prog( $X_n \rightarrow \alpha_n$ ); begin  
  scan;  
   $X_0$ ;  
  if nextsym = "#"  
    then accept  
    else error("...")  
  fi  
end
```

```

N_prog( $X \rightarrow \alpha$ ) =
    proc X;
    begin
        progr( $[X \rightarrow \cdot \alpha]$ )
    end ;
progr( $[X \rightarrow \dots (\alpha_1 | \alpha_2 | \dots | \alpha_{k-1} | \alpha_k) \dots]$ ) =
    case nextsym in
    FiFo( $[X \rightarrow \dots (\alpha_1 | \alpha_2 | \dots | \alpha_{k-1} | \alpha_k) \dots]$ ) :
        progr( $[X \rightarrow \dots (\alpha_1 | \alpha_2 | \dots | \alpha_{k-1} | \alpha_k) \dots]$ );
    FiFo( $[X \rightarrow \dots (\alpha_1 \cdot \alpha_2 | \dots | \alpha_{k-1} | \alpha_k) \dots]$ ) :
        progr( $[X \rightarrow \dots (\alpha_1 \cdot \alpha_2 | \dots | \alpha_{k-1} | \alpha_k) \dots]$ );
        :
    FiFo( $[X \rightarrow \dots (\alpha_1 | \alpha_2 | \dots | \alpha_{k-1} | \alpha_k) \dots]$ ) :
        progr( $[X \rightarrow \dots (\alpha_1 | \alpha_2 | \dots | \alpha_{k-1} | \alpha_k) \dots]$ );
    otherwise progr( $[X \rightarrow \dots (\alpha_1 | \alpha_2 | \dots | \alpha_{k-1} | \alpha_k) \dots]$ );
    endcase
progr( $[X \rightarrow \dots (\alpha_1 \alpha_2 \dots \alpha_k) \dots]$ ) =
    progr( $[X \rightarrow \dots (\alpha_1 \alpha_2 \dots \alpha_k) \dots]$ );
    progr( $[X \rightarrow \dots (\alpha_1 \cdot \alpha_2 \dots \alpha_k) \dots]$ );
        :

```

Für $a \in V_T$ ist

$\text{progr}([X \rightarrow \dots .a \dots]) =$

```
  if nextsym = a then scan  
  else error  
  fi
```

Für $Y \in V_N$ ist

$\text{progr}([X \rightarrow \dots .Y \dots]) = Y$