

## Right-Regular Context-Free Parsing

Recurring scheme in programming languages: Lists of sth...

$$S \rightarrow b \quad | \quad S a b$$

Alternative idea: Regular Expressions

$$S \rightarrow (ba)^* b$$

Title: Petter: Compiler Construction (21.05.2020)  
-17: RLL(1) Parsers

Date: Wed May 06 10:42:52 CEST 2020

Duration: 46:07 min

Pages: 18

49/55

## Right-Regular Context-Free Parsing

Recurring scheme in programming languages: Lists of sth...

$$S \rightarrow b \quad | \quad S a b$$

Alternative idea: Regular Expressions

$$S \rightarrow (ba)^* b$$

### Definition: Right-Regular Context-Free Grammar

A right-regular context-free grammar (RR-CFG) is a

4-tuple  $G = (N, T, P, S)$  with:

- $N$  the set of nonterminals,
- $T$  the set of terminals,
- $P$  the set of rules with regular expressions of symbols as rhs,
- $S \in N$  the start symbol

## Right-Regular Context-Free Parsing

Recurring scheme in programming languages: Lists of sth...

$$S \rightarrow b \quad | \quad S a b$$

Alternative idea: Regular Expressions

$$S \rightarrow (ba)^* b$$

### Definition: Right-Regular Context-Free Grammar

A right-regular context-free grammar (RR-CFG) is a

4-tuple  $G = (N, T, P, S)$  with:

- $N$  the set of nonterminals,
- $T$  the set of terminals,
- $P$  the set of rules with regular expressions of symbols as rhs,
- $S \in N$  the start symbol

### Example: Arithmetic Expressions

$$\begin{array}{l} S \rightarrow E \\ E \rightarrow T [(+T)^*] \\ T \rightarrow F [(*F)^*] \\ F \rightarrow (E) \mid \text{name} \mid \text{int} \end{array}$$

49/55

49/55

Idea 1: Rewrite the rules from  $G$  to  $\langle G \rangle$ :

$A$	$\rightarrow$	$\langle \alpha \rangle$	if $A \rightarrow \alpha \in P$
$\langle \alpha \rangle$	$\rightarrow$	$\alpha$	if $\alpha \in N \cup T$
$\langle \epsilon \rangle$	$\rightarrow$	$\epsilon$	
$\langle \alpha^* \rangle$	$\rightarrow$	$\epsilon \sqcup \langle \alpha \rangle \langle \alpha^* \rangle$	if $\alpha \in \text{Regex}_{T,N}$
$\langle \alpha_1 \dots \alpha_n \rangle$	$\rightarrow$	$\langle \alpha_1 \rangle \dots \langle \alpha_n \rangle$	if $\alpha_i \in \text{Regex}_{T,N}$
$\langle \alpha_1   \dots   \alpha_n \rangle$	$\rightarrow$	$\langle \alpha_1 \rangle   \dots   \langle \alpha_n \rangle$	if $\alpha_i \in \text{Regex}_{T,N}$

... and generate the according LL(k)-Parser  $M_{\langle G \rangle}^L$

Idea 1: Rewrite the rules from  $G$  to  $\langle G \rangle$ :

$A$	$\rightarrow$	$\langle \alpha \rangle$	if $A \rightarrow \alpha \in P$
$\langle \alpha \rangle$	$\rightarrow$	$\alpha$	if $\alpha \in N \cup T$
$\langle \epsilon \rangle$	$\rightarrow$	$\epsilon$	
$\langle \alpha^* \rangle$	$\rightarrow$	$\epsilon \sqcup \langle \alpha \rangle \langle \alpha^* \rangle$	if $\alpha \in \text{Regex}_{T,N}$
$\langle \alpha_1 \dots \alpha_n \rangle$	$\rightarrow$	$\langle \alpha_1 \rangle \dots \langle \alpha_n \rangle$	if $\alpha_i \in \text{Regex}_{T,N}$
$\langle \alpha_1   \dots   \alpha_n \rangle$	$\rightarrow$	$\langle \alpha_1 \rangle   \dots   \langle \alpha_n \rangle$	if $\alpha_i \in \text{Regex}_{T,N}$

... and generate the according LL(k)-Parser  $M_{\langle G \rangle}^L$

Example: Arithmetic Expressions cont'd

$S$	$\rightarrow$	$E$
$E$	$\rightarrow$	$T ( + T )^*$
$T$	$\rightarrow$	$F ( * F )^*$
$F$	$\rightarrow$	$( E )   \text{name}   \text{int}$

50/55

50/55

Idea 1: Rewrite the rules from  $G$  to  $\langle G \rangle$ :

$A$	$\rightarrow$	$\langle \alpha \rangle$	if $A \rightarrow \alpha \in P$
$\langle \alpha \rangle$	$\rightarrow$	$\alpha$	if $\alpha \in N \cup T$
$\langle \epsilon \rangle$	$\rightarrow$	$\epsilon$	
$\langle \alpha^* \rangle$	$\rightarrow$	$\epsilon \sqcup \langle \alpha \rangle \langle \alpha^* \rangle$	if $\alpha \in \text{Regex}_{T,N}$
$\langle \alpha_1 \dots \alpha_n \rangle$	$\rightarrow$	$\langle \alpha_1 \rangle \dots \langle \alpha_n \rangle$	if $\alpha_i \in \text{Regex}_{T,N}$
$\langle \alpha_1   \dots   \alpha_n \rangle$	$\rightarrow$	$\langle \alpha_1 \rangle   \dots   \langle \alpha_n \rangle$	if $\alpha_i \in \text{Regex}_{T,N}$

... and generate the according LL(k)-Parser  $M_{\langle G \rangle}^L$

Example: Arithmetic Expressions cont'd

$S$	$\rightarrow$	$E$
$E$	$\rightarrow$	$\langle T ( + T )^* \rangle$
$T$	$\rightarrow$	$F ( * F )^*$
$F$	$\rightarrow$	$( E )   \text{name}   \text{int}$
$\langle T ( + T )^* \rangle$	$\rightarrow$	$T \langle (+ T)^* \rangle$

Idea 1: Rewrite the rules from  $G$  to  $\langle G \rangle$ :

$A$	$\rightarrow$	$\langle \alpha \rangle$	if $A \rightarrow \alpha \in P$
$\langle \alpha \rangle$	$\rightarrow$	$\alpha$	if $\alpha \in N \cup T$
$\langle \epsilon \rangle$	$\rightarrow$	$\epsilon$	
$\langle \alpha^* \rangle$	$\rightarrow$	$\epsilon \sqcup \langle \alpha \rangle \langle \alpha^* \rangle$	if $\alpha \in \text{Regex}_{T,N}$
$\langle \alpha_1 \dots \alpha_n \rangle$	$\rightarrow$	$\langle \alpha_1 \rangle \dots \langle \alpha_n \rangle$	if $\alpha_i \in \text{Regex}_{T,N}$
$\langle \alpha_1   \dots   \alpha_n \rangle$	$\rightarrow$	$\langle \alpha_1 \rangle   \dots   \langle \alpha_n \rangle$	if $\alpha_i \in \text{Regex}_{T,N}$

... and generate the according LL(k)-Parser  $M_{\langle G \rangle}^L$

Example: Arithmetic Expressions cont'd

$S$	$\rightarrow$	$E$
$E$	$\rightarrow$	$\langle T ( + T )^* \rangle$
$T$	$\rightarrow$	$F ( * F )^*$
$F$	$\rightarrow$	$( E )   \text{name}   \text{int}$
$\langle T ( + T )^* \rangle$	$\rightarrow$	$T \langle (+ T)^* \rangle$
$\langle (+ T)^* \rangle$	$\rightarrow$	$\epsilon \sqcup \langle (+ T)^* \rangle$

50/55

50/55

Idea 1: Rewrite the rules from  $G$  to  $\langle G \rangle$ :

$A$	$\rightarrow \langle \alpha \rangle$	if $A \rightarrow \alpha \in P$
$\langle \alpha \rangle$	$\rightarrow \alpha$	if $\alpha \in N \cup T$
$\langle \epsilon \rangle$	$\rightarrow \epsilon$	
$\langle \alpha^* \rangle$	$\rightarrow \epsilon \mid \langle \alpha \rangle \langle \alpha^* \rangle$	if $\alpha \in \text{Regex}_{T,N}$
$\langle \alpha_1 \dots \alpha_n \rangle$	$\rightarrow \langle \alpha_1 \rangle \dots \langle \alpha_n \rangle$	if $\alpha_i \in \text{Regex}_{T,N}$
$\langle \alpha_1 \mid \dots \mid \alpha_n \rangle$	$\rightarrow \langle \alpha_1 \rangle \mid \dots \mid \langle \alpha_n \rangle$	if $\alpha_i \in \text{Regex}_{T,N}$

... and generate the according LL(k)-Parser  $M_{\langle G \rangle}^L$

Example: Arithmetic Expressions cont'd

$S$	$\rightarrow E$
$E$	$\rightarrow \langle T (+T)^* \rangle$
$T$	$\rightarrow \langle F (*F)^* \rangle$
$F$	$\rightarrow (E) \mid \text{name} \mid \text{int}$
$\langle T (+T)^* \rangle$	$\rightarrow T \langle (+T)^* \rangle$
$\langle (+T)^* \rangle$	$\rightarrow \epsilon \mid \langle +T \rangle \langle (+T)^* \rangle$
$\langle +T \rangle$	$\rightarrow +T$

Idea 1: Rewrite the rules from  $G$  to  $\langle G \rangle$ :

$A$	$\rightarrow \langle \alpha \rangle$	if $A \rightarrow \alpha \in P$
$\langle \alpha \rangle$	$\rightarrow \alpha$	if $\alpha \in N \cup T$
$\langle \epsilon \rangle$	$\rightarrow \epsilon$	
$\langle \alpha^* \rangle$	$\rightarrow \epsilon \mid \langle \alpha \rangle \langle \alpha^* \rangle$	if $\alpha \in \text{Regex}_{T,N}$
$\langle \alpha_1 \dots \alpha_n \rangle$	$\rightarrow \langle \alpha_1 \rangle \dots \langle \alpha_n \rangle$	if $\alpha_i \in \text{Regex}_{T,N}$
$\langle \alpha_1 \mid \dots \mid \alpha_n \rangle$	$\rightarrow \langle \alpha_1 \rangle \mid \dots \mid \langle \alpha_n \rangle$	if $\alpha_i \in \text{Regex}_{T,N}$

... and generate the according LL(k)-Parser  $M_{\langle G \rangle}^L$

Example: Arithmetic Expressions cont'd

$S$	$\rightarrow E$
$E$	$\rightarrow \langle T (+T)^* \rangle$
$T$	$\rightarrow \langle F (*F)^* \rangle$
$F$	$\rightarrow (E) \mid \text{name} \mid \text{int}$
$\langle T (+T)^* \rangle$	$\rightarrow T \langle (+T)^* \rangle$
$\langle (+T)^* \rangle$	$\rightarrow \epsilon \mid \langle +T \rangle \langle (+T)^* \rangle$
$\langle +T \rangle$	$\rightarrow +T$

Idea 1: Rewrite the rules from  $G$  to  $\langle G \rangle$ :

$A$	$\rightarrow \langle \alpha \rangle$	if $A \rightarrow \alpha \in P$
$\langle \alpha \rangle$	$\rightarrow \alpha$	if $\alpha \in N \cup T$
$\langle \epsilon \rangle$	$\rightarrow \epsilon$	
$\langle \alpha^* \rangle$	$\rightarrow \epsilon \mid \langle \alpha \rangle \langle \alpha^* \rangle$	if $\alpha \in \text{Regex}_{T,N}$
$\langle \alpha_1 \dots \alpha_n \rangle$	$\rightarrow \langle \alpha_1 \rangle \dots \langle \alpha_n \rangle$	if $\alpha_i \in \text{Regex}_{T,N}$
$\langle \alpha_1 \mid \dots \mid \alpha_n \rangle$	$\rightarrow \langle \alpha_1 \rangle \mid \dots \mid \langle \alpha_n \rangle$	if $\alpha_i \in \text{Regex}_{T,N}$

... and generate the according LL(k)-Parser  $M_{\langle G \rangle}^L$

Example: Arithmetic Expressions cont'd

$S$	$\rightarrow E$
$E$	$\rightarrow \langle T (+T)^* \rangle \cong A$
$T$	$\rightarrow \langle F (*F)^* \rangle \cong B$
$F$	$\rightarrow (E) \mid \text{name} \mid \text{int}$
$A \langle T (+T)^* \rangle$	$\rightarrow T \langle (+T)^* \rangle \cong C$
$C \langle (+T)^* \rangle$	$\rightarrow \epsilon \mid \langle +T \rangle \langle (+T)^* \rangle \cong C$
$\Delta \langle +T \rangle$	$\rightarrow +T$
$B \langle F (*F)^* \rangle$	$\rightarrow F \langle (*F)^* \rangle$
$\langle (*F)^* \rangle$	$\rightarrow \epsilon \mid \langle *F \rangle \langle (*F)^* \rangle$
$\langle *F \rangle$	$\rightarrow *F$

50/55

50/55



Reinold Heckmann

Definition:

An RR-CFG  $G$  is called  $RLL(1)$ ,  
if the corresponding CFG  $\langle G \rangle$  is an  $LL(1)$  grammar.

Discussion

- directly yields the table driven parser  $M_{\langle G \rangle}^L$  for  $RLL(1)$  grammars
- however: mapping regular expressions to recursive productions unnecessarily strains the stack  
→ instead directly construct automaton in the style of Berry-Sethi

50/55

51/55

## Idea 2: Recursive Descent RLL Parsers:

*Recursive descent* RLL(1)-parsers are an alternative to table-driven parsers; apart from the usual function `scan()`, we generate a program frame with the lookahead function `expect()` and the main parsing method `parse()`:

```

void int next;
boolean expect(Set E){
    if ({e, next} ∩ E = ∅){
        cerr << "Expected" << E << "found" << next;
        exit(0);
    }
    return;
}
void parse(){
    next = scan();
    expect(First1(S));
    S();
    expect({EOF});
}

```

52/55

## Idea 2: Recursive Descent RLL Parsers:

$$\begin{aligned}
 \text{generate}(r^*) &= \text{while } (\text{next} \in F_\epsilon(r)) \{ \\
 &\quad \text{generate}(r) \\
 &\} \\
 \text{generate}(r_1 \mid \dots \mid r_k) &= \text{switch}(\text{next}) \{ \\
 &\quad \text{labels}(\text{First}_1(r_1)) \text{ generate}(r_1) \text{ break}; \\
 &\quad \vdots \\
 &\quad \text{labels}(\text{First}_1(r_k)) \text{ generate}(r_k) \text{ break}; \\
 &\} \\
 \text{labels}(\{\alpha_1, \dots, \alpha_m\}) &= \text{label}(\alpha_1); \dots \text{ label}(\alpha_m); \\
 \text{label}(\alpha) &= \text{case } \alpha \\
 \text{label}(\epsilon) &= \text{default}
 \end{aligned}$$

54/55

## Idea 2: Recursive Descent RLL Parsers:

For each  $A \rightarrow \alpha \in P$ , we introduce:

```

void A(){
    generate(α)
}

```

with the meta-program `generate` being defined by structural decomposition of  $\alpha$ :

$$\begin{aligned}
 \text{generate}(r_1 \dots r_k) &= \text{generate}(r_1) \\
 &\quad \text{expect}(\text{First}_1(r_2)); \\
 &\quad \text{generate}(r_2) \\
 &\quad \vdots \\
 &\quad \text{expect}(\text{First}_1(r_k)); \\
 &\quad \text{generate}(r_k) \\
 \text{generate}(e) &= ; \\
 \text{generate}(a) &= \text{next} = \text{scan}(); \\
 \text{generate}(A) &= A()
 \end{aligned}$$

53/55

## Topdown-Parsing

### Discussion

- A practical implementation of an *RLL(1)*-parser via *recursive descent* is a straight-forward idea
- However, *only a subset* of the deterministic contextfree languages can be parsed this way.
- As soon as  $\text{First}_1(\_)$  sets are not disjoint any more,

55/55

## Idea 2: Recursive Descent RLL Parsers:

```
generate( $r^*$ )      = while ( next ∈ F $_e(r)$  ) {  
    generate( $r$ )  
}  
generate( $r_1 \mid \dots \mid r_k$ ) = switch(next) {  
    labels(First $_1(r_1)$ ) generate( $r_1$ ) break ;  
    :  
    labels(First $_1(r_k)$ ) generate( $r_k$ ) break ;  
}  
labels({ $\alpha_1, \dots, \alpha_m$ }) = label( $\alpha_1$ ): ... label( $\alpha_m$ ):  
label( $\alpha$ )           = case  $\alpha$   
label( $\epsilon$ )           = default
```

## Topdown-Parsing

### Discussion

- A practical implementation of an  $RLL(1)$ -parser via recursive descent is a straight-forward idea
- However, only a subset of the deterministic contextfree languages can be parsed this way.
- As soon as  $First_1(\cdot)$  sets are not disjoint any more,

54 / 55

55 / 55

## Topdown-Parsing

### Discussion

- A practical implementation of an  $RLL(1)$ -parser via recursive descent is a straight-forward idea
- However, only a subset of the deterministic contextfree languages can be parsed this way.
- As soon as  $First_1(\cdot)$  sets are not disjoint any more,
  - Solution 1: For many accessibly written grammars, the alternation between right hand sides happens too early. Keeping the common prefixes of right hand sides joined and introducing a new production for the actual diverging sentence forms often helps.
  - Solution 2: Introduce ranked grammars and decide conflicting lookahead always in favour of the higher ranked alternative
    - relation to  $LL$ , parsing not so clear any more
    - not so clear for  $^*$  operator how to decide
  - Solution 3: Going from  $LL(1)$  to  $LL(k)$   
The size of the occurring sets is rapidly increasing with larger  $k$   
*Unfortunately*, even  $LL(k)$  parsers are not sufficient to accept all deterministic contextfree languages. (regular lookahead →  $LL(*)$ )
- In practical systems, this often motivates the implementation of  $k = 1$  only ...

55 / 55