

**Script** generated by TTT

Title: Petter: Compilerbau (24.05.2018)

Date: Thu May 24 14:18:07 CEST 2018

Duration: 89:24 min

Pages: 33

## Right-Regular Context-Free Parsing

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

$S \rightarrow b \mid S a b$

Alternative idea: Regular Expressions

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

## Left Recursion

### Theorem:

Let a grammar  $G$  be reduced and left-recursive, then  $G$  is not  $LL(k)$  for any  $k$ .

### Proof:

Let wlog.  $A \rightarrow A\beta \mid \alpha \in P$   
and  $A$  be reachable from  $S$

Assumption:  $G$  is  $LL(k)$

$\Rightarrow \text{First}_k(\alpha\beta^n\gamma) \cap$   
 $\text{First}_k(\alpha\beta^{n+1}\gamma) = \emptyset$

**Case 1:**  $\beta \rightarrow^* \epsilon$  — Contradiction !!!

**Case 2:**  $\beta \rightarrow^* w \neq \epsilon \implies \text{First}_k(\alpha w^k\gamma) \cap \text{First}_k(\alpha w^{k+1}\gamma) \neq \emptyset$

106 / 288

## Right-Regular Context-Free Parsing

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

$S \rightarrow b \mid S a b$

Alternative idea: Regular Expressions

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

### 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

## 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  $M_{\langle G \rangle}^L$

108 / 288

## 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()`:

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

110 / 288



Reinhold 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

109 / 288

## 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  $\alpha$ :

```
generate(r1...rk) = generate(r1)
                    if (!expect(First1(r2))) exit(0);
                    generate(r2)
                    :
                    if (!expect(First1(rk))) exit(0);
                    generate(rk)
generate(ε)         = ;
generate(a)         = consume();
generate(A)         = A();
```

111 / 288

## Idea 2: Recursive Descent RLL Parsers:

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

112 / 288

## 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,

113 / 288

## Idea 2: Recursive Descent RLL Parsers:

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

112 / 288

## Syntactic Analysis

## Chapter 4: Bottom-up Analysis

114 / 288

# Shift-Reduce Parser



Donald Knuth

## Idea:

We *delay* the decision whether to reduce until we know, whether the input matches the right-hand-side of a rule!

**Construction:** Shift-Reduce parser  $M_G^R$

- The input is shifted successively to the pushdown.
- Is there a **complete right-hand side** (a **handle**) atop the pushdown, it is replaced (**reduced**) by the corresponding left-hand side

# Shift-Reduce Parser

## Example:

$S \rightarrow AB$   
 $A \rightarrow a$   
 $B \rightarrow b$

The pushdown automaton:

**States:**  $q_0, f, a, b, A, B, S;$   
**Start state:**  $q_0$   
**End state:**  $f$

$q_0$	$a$	$q_0 a$
$a$	$\epsilon$	$A$
$A$	$b$	$Ab$
$b$	$\epsilon$	$B$
$AB$	$\epsilon$	$S$
$q_0 S$	$\epsilon$	$f$

# Shift-Reduce Parser

## Observation:

- The sequence of reductions corresponds to a **reverse rightmost-derivation** for the input
- To prove correctness, we have to prove:

$$(\epsilon, w) \vdash^* (A, \epsilon) \quad \text{iff} \quad A \rightarrow^* w$$

- The shift-reduce pushdown automaton  $M_G^R$  is in general also **non-deterministic**
- For a deterministic parsing-algorithm, we have to identify computation-states for reduction

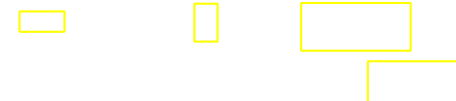
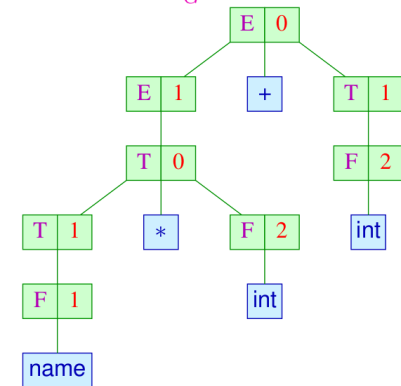
$\implies$  LR-Parsing

# Reverse Rightmost Derivations in Shift-Reduce-Parsers

**Idea:** Observe **reverse rightmost**-derivations of  $M_G^R$ !

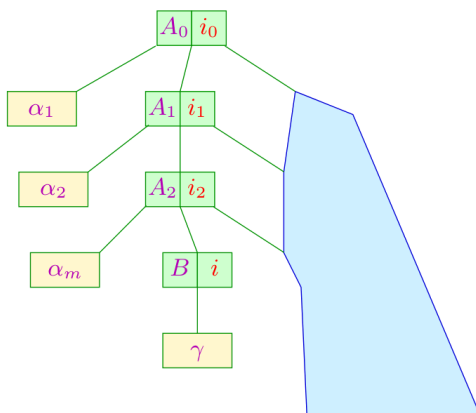
Input:  
 $counter * 2 + 40$

Pushdown:  
 $(q_0)$



## Bottom-up Analysis: Viable Prefix

$\alpha\gamma$  is viable for  $[B \rightarrow \gamma \bullet]$  iff  $S \xrightarrow{*}_R \alpha B v$

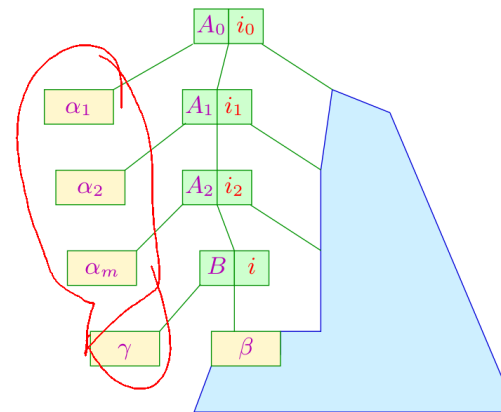


... with  $\alpha = \alpha_1 \dots \alpha_m$

120 / 288

## Bottom-up Analysis: Admissible Items

The item  $[B \rightarrow \gamma \bullet \beta]$  is called **admissible** for  $\alpha'$  iff  $S \xrightarrow{*}_R \alpha B v$  with  $\alpha' = \alpha\gamma$ :



... with  $\alpha = \alpha_1 \dots \alpha_m$

121 / 288

## Characteristic Automaton

### Observation:

The set of viable prefixes from  $(N \cup T)^*$  for (admissible) items can be computed from the content of the **shift-reduce parser's** pushdown with the help of a finite automaton:

States: **Items**

Start state:  $[S' \rightarrow \bullet S]$

Final states:  $\{[B \rightarrow \gamma \bullet] \mid B \rightarrow \gamma \in P\}$

Transitions:

- (1)  $([A \rightarrow \alpha \bullet X \beta], X, [A \rightarrow \alpha X \bullet \beta]), X \in (N \cup T), [A \rightarrow \alpha X \bullet \beta] \in P;$
- (2)  $([A \rightarrow \alpha \bullet B \beta], \epsilon, [B \rightarrow \bullet \gamma]), A \rightarrow \alpha B \beta, [B \rightarrow \bullet \gamma] \in P;$

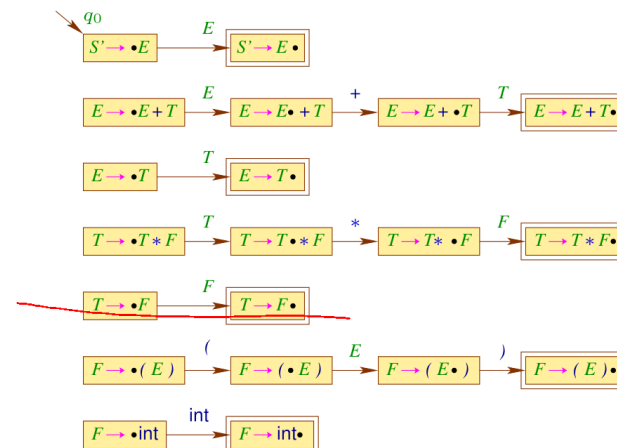
The automaton  $c(G)$  is called **characteristic automaton** for  $G$ .

122 / 288

## Characteristic Automaton

For example:

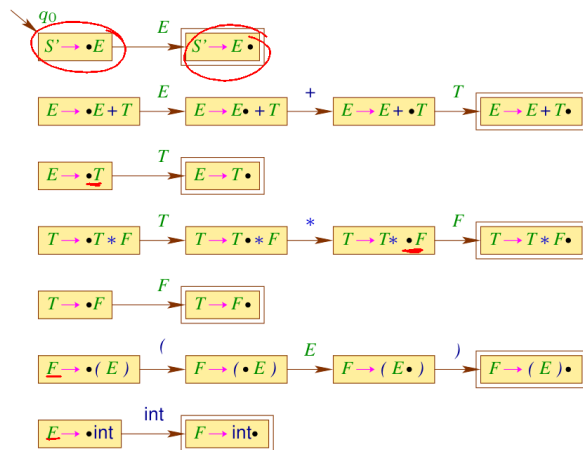
$E \rightarrow E + T$	$T$
$T \rightarrow T * F$	$F$
$F \rightarrow (E)$	$\text{int}$



123 / 288

## Characteristic Automaton

For example:

$$\begin{array}{l} E \rightarrow E+T \quad | \quad T \\ T \rightarrow T*F \quad | \quad F \\ F \rightarrow (E) \quad | \quad \text{int} \end{array}$$


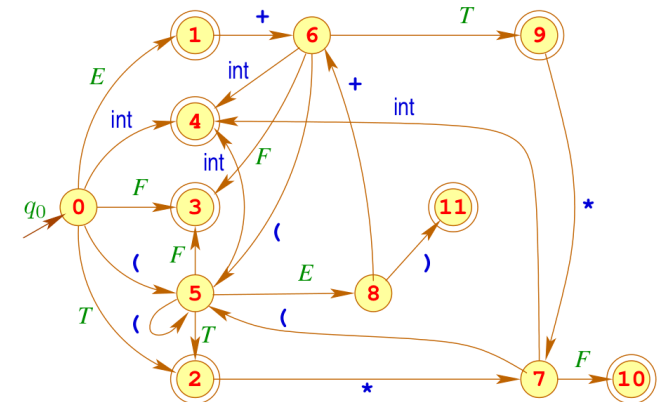
123 / 288

## Canonical LR(0)-Automaton

The canonical  $LR(0)$ -automaton  $LR(G)$  is created from  $c(G)$  by:

- 1 performing arbitrarily many  $\epsilon$ -transitions after every consuming transition
- 2 performing the powerset construction

... for example:



124 / 288

## Canonical LR(0)-Automaton

Observation:

The canonical  $LR(0)$ -automaton can be created **directly** from the grammar.

Therefore we need a helper function  $\delta_\epsilon^*$  ( $\epsilon$ -closure)

$$\delta_\epsilon^*(q) = q \cup \{ [B \rightarrow \bullet \gamma] \mid B \rightarrow \gamma \in P, [A \rightarrow \alpha \bullet B' \beta'] \in q, B' \rightarrow^* B \beta \}$$

We define:

States: Sets of items;

Start state:  $\delta_\epsilon^* \{ [S' \rightarrow \bullet S] \}$

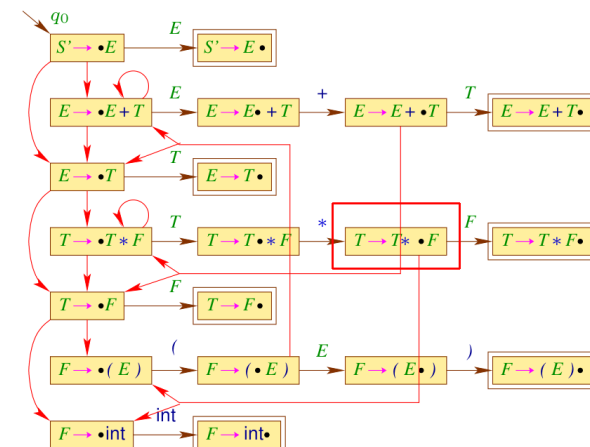
Final states:  $\{ q \mid [A \rightarrow \alpha \bullet] \in q \}$

Transitions:  $\delta(q, X) = \delta_\epsilon^* \{ [A \rightarrow \alpha X \bullet \beta] \mid [A \rightarrow \alpha \bullet X \beta] \in q \}$

125 / 288

## Characteristic Automaton

For example:

$$\begin{array}{l} E \rightarrow E+T \quad | \quad T \\ T \rightarrow T*F \quad | \quad F \\ F \rightarrow (E) \quad | \quad \text{int} \end{array}$$


123 / 288

## Canonical LR(0)-Automaton

### Observation:

The canonical LR(0)-automaton can be created directly from the grammar.

Therefore we need a helper function  $\delta_\epsilon^*$  ( $\epsilon$ -closure)

$$\delta_\epsilon^*(q) = q \cup \{ [B \rightarrow \bullet \gamma] \mid B \rightarrow \gamma \in P, \\ [A \rightarrow \alpha \bullet B' \beta'] \in q, \\ B' \rightarrow^* B \beta' \}$$

We define:

States: Sets of items;

Start state:  $\delta_\epsilon^* \{ [S' \rightarrow \bullet S] \}$

Final states:  $\{ q \mid [A \rightarrow \alpha \bullet] \in q \}$

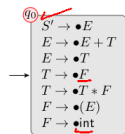
Transitions:  $\delta(q, X) = \delta_\epsilon^* \{ [A \rightarrow \alpha X \bullet \beta] \mid [A \rightarrow \alpha \bullet X \beta] \in q \}$

125 / 288

## Canonical LR(0)-Automaton

For example:

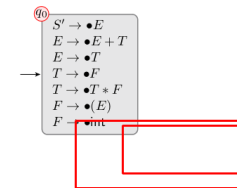
$S' \rightarrow E$	$\rightarrow$	$E$		
$E \rightarrow E + T$	$\rightarrow$	$E + T$		$T$
$T \rightarrow T * F$	$\rightarrow$	$T * F$		$F$
$F \rightarrow ( E )$	$\rightarrow$	$( E )$		$\text{int}$



## Canonical LR(0)-Automaton

For example:

$S' \rightarrow E$	$\rightarrow$	$E$		
$E \rightarrow E + T$	$\rightarrow$	$E + T$		$T$
$T \rightarrow T * F$	$\rightarrow$	$T * F$		$F$
$F \rightarrow ( E )$	$\rightarrow$	$( E )$		$\text{int}$



## LR(0)-Parser

Idea for a parser:

- The parser manages a viable prefix  $\alpha = X_1 \dots X_m$  on the pushdown and uses  $LR(G)$ , to identify reduction spots.
- It can reduce with  $A \rightarrow \gamma$ , if  $[A \rightarrow \gamma \bullet]$  is admissible for  $\alpha$

Optimization:

We push the states instead of the  $X_i$  in order not to process the pushdown's content with the automaton anew all the time.

Reduction with  $A \rightarrow \gamma$  leads to popping the uppermost  $|\gamma|$  states and continue with the state on top of the stack and input  $A$ .

Attention:

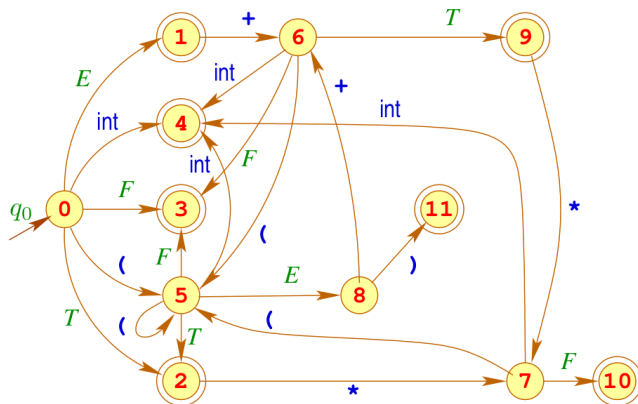
This parser is only deterministic, if each final state of the canonical LR(0)-automaton is conflict free.

## Canonical LR(0)-Automaton

The canonical  $LR(0)$ -automaton  $LR(G)$  is created from  $c(G)$  by:

- performing arbitrarily many  $\epsilon$ -transitions after every consuming transition
- performing the powerset construction

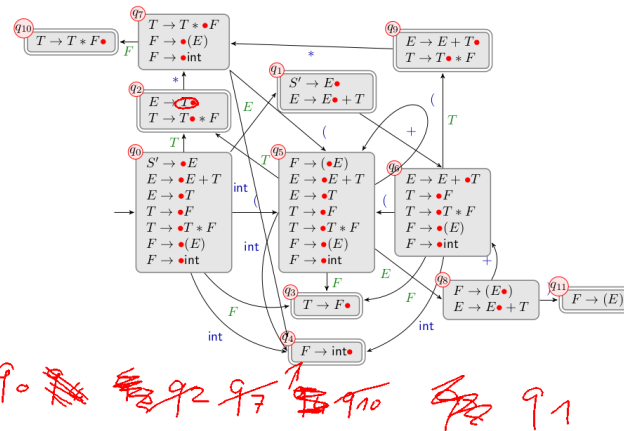
... for example:



124 / 288

## Canonical LR(0)-Automaton

For example:

$$\begin{array}{l} S' \rightarrow E \\ E \rightarrow E + T \quad | \quad T \\ T \rightarrow T * F \quad | \quad F \\ F \rightarrow ( E ) \quad | \quad \text{int} \end{array}$$


126 / 288

## LR(0)-Parser

Idea for a parser:

- The parser manages a viable prefix  $\alpha = X_1 \dots X_m$  on the pushdown and uses  $LR(G)$ , to identify reduction spots.
- It can reduce with  $A \rightarrow \gamma$ , if  $[A \rightarrow \gamma \bullet]$  is admissible for  $\alpha$

Optimization:

We push the **states** instead of the  $X_i$  in order not to process the pushdown's content with the automaton anew all the time. Reduction with  $[A \rightarrow \gamma \bullet]$  leads to popping the uppermost  $|\gamma|$  states and continue with the state on top of the stack and input  $A$ .

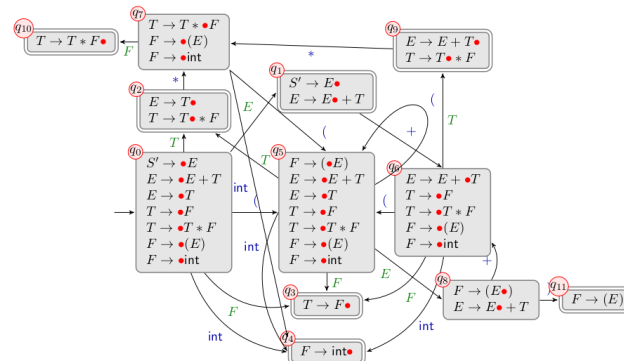
Attention:

This parser is only **deterministic**, if each final state of the canonical  $LR(0)$ -automaton is **conflict free**.

127 / 288

## Canonical LR(0)-Automaton

For example:

$$\begin{array}{l} S' \rightarrow E \\ E \rightarrow E + T \quad | \quad T \\ T \rightarrow T * F \quad | \quad F \\ F \rightarrow ( E ) \quad | \quad \text{int} \end{array}$$


126 / 288



... for example:

$$q_1 = \{[S' \rightarrow E \bullet], [E \rightarrow E \bullet + T]\}$$

$$q_2 = \{[E \rightarrow T \bullet], [T \rightarrow T \bullet * F]\}$$

$$q_3 = \{[T \rightarrow F \bullet]\}$$

$$q_4 = \{[F \rightarrow \text{int} \bullet]\}$$

$$q_9 = \{[E \rightarrow E + T \bullet], [T \rightarrow T \bullet * F]\}$$

$$q_{10} = \{[T \rightarrow T * F \bullet]\}$$

$$q_{11} = \{[F \rightarrow ( E ) \bullet]\}$$

The final states  $q_1, q_2, q_9$  contain more than one admissible item  
 $\Rightarrow$  non deterministic!

### Idea for a parser:

- The parser manages a viable prefix  $\alpha = X_1 \dots X_m$  on the pushdown and uses  $LR(G)$ , to identify reduction spots.
- It can reduce with  $A \rightarrow \gamma$ , if  $[A \rightarrow \gamma \bullet]$  is admissible for  $\alpha$

### Optimization:

We push the **states** instead of the  $X_i$  in order not to process the pushdown's content with the automaton anew all the time. Reduction with  $A \rightarrow \gamma$  leads to popping the uppermost  $|\gamma|$  states and continue with the state on top of the stack and input  $A$ .

### Attention:

This parser is only **deterministic**, if each final state of the canonical  $LR(0)$ -automaton is **conflict free**.