



## LR Parsing, Part 2

### Constructing Parse Tables

Parse table construction  
Grammar conflict handling  
Categories of LR Grammars and Parsers



### An NFA Recognizing Viable Prefixes

- A.k.a. the "characteristic finite automaton" for a grammar G
- States: LR(0) items (= context-free items) of extend. grammar
- Input stream: The grammar symbols on the stack
- Start state:  $[S' \rightarrow -|.S]$  Final state:  $[S' \rightarrow -|S.]$
- Transitions:
  - "move dot across symbol" if symbol found next on stack:
    - $A \rightarrow \alpha.B\gamma$  to  $A \rightarrow \alpha B.\gamma$
    - $A \rightarrow \alpha.b\gamma$  to  $A \rightarrow \alpha b.\gamma$
  - $\epsilon$ -transitions to LR(0)-items for nonterminal productions from items where the dot precedes that nonterminal:
    - $A \rightarrow \alpha.B\gamma$  to  $B \rightarrow |\beta$

### Computing the Closure



For a set  $I$  of LR(0) items compute  $Closure(I)$ :

1.  $Closure(I) := I$
2. If  $\exists [A \rightarrow \alpha.B\beta]$  in  $Closure(I)$  and  $\exists$  production  $B \rightarrow \gamma$  then add  $[B \rightarrow \cdot\gamma]$  to  $Closure(I)$  (if not already there)
3. Repeat Step 2 until no more items can be added to  $Closure(I)$ .

Remarks:

- For  $s=[A \rightarrow \alpha.B\gamma]$ ,  $Closure(s)$  contains all NFA states reachable from  $s$  via  $\epsilon$ -transitions, i.e., starting from which any substring derivable from  $B\beta$  could be recognized. A.k.a.  $\epsilon$ -closure(s).
- Then apply the well-known subset construction to transform Closure-NFA  $\rightarrow$  DFA.
- DFA states will be sets unioning closures of NFA states

### Representing Sets of Items



- Any item  $[A \rightarrow \alpha.\beta]$  can be represented by 2 integers:
  - production number
  - position of the dot within the RHS of that production
- The resulting sets often contain "closure" items (where the dot is at the beginning of the RHS).
  - Can easily be reconstructed (on demand) from other ("kernel") items
    - ▶ **Kernel items:** start state  $[S' \rightarrow -|.S]$ , plus all items where the dot is not at the left end.
  - Store only kernel items explicitly, to save space

### GOTO Function and DFA States



Given: Set  $I$  of items, grammar symbol  $X$

- $GOTO(I, X) := \bigcup_{[A \rightarrow \alpha.X\beta] \in I} Closure(\{[A \rightarrow \alpha X.\beta]\})$ 
  - To become the state transitions in the DFA
- Now do the **subset construction** to obtain the DFA states:
  - $C := Closure(\{[S' \rightarrow -|.S]\})$  //  $C$ : Set of sets of NFA states
  - repeat**
  - for** each set of items  $I$  of  $C$ :
  - for** each grammar symbol  $X$
  - if** ( $GOTO(I,X)$  is not empty and not in  $C$ )
  - add**  $GOTO(I,X)$  to  $C$

### Resulting DFA



- (Example: see whiteboard)
- All states correspond to some viable prefix
- Final states: contain at least one item with dot to the right
  - recognized some handle  $\rightarrow$  reduce *may (must)* follow
- Other states: handle recognition incomplete  $\rightarrow$  shift will follow
- The DFA is also called the GOTO graph (not the same as the LR GOTO Table!!).
- This automaton is deterministic as a FA (i.e., selecting transitions considering only input symbol consumption) but can still be nondeterministic as a pushdown automaton (e.g., in state  $I_1$  above: to reduce or not to reduce?)

## From DFA to parser tables: ACTION



- For each DFA transition  $I_i \rightarrow I_j$  reading a terminal  $a$  in  $\Sigma$  (thus,  $I_i$  contains items of kind  $[X \rightarrow \alpha.a\beta]$ )
  - enter  $S_j$  (shift, new state  $I_j$ ) in ACTION
- For each DFA final state  $I_i$  (containing a complete item  $[X \rightarrow \alpha.]$ )
  - enter  $R_x$  (reduce,  $x$  = prod. rule number for  $X \rightarrow \alpha.$ ) in ACTION[ $i, t$ ] ...
    - LR(0) parser: for all  $t$  in  $\Sigma$  (all entries in row  $i$ )
    - SLR(1) parser: for all  $t$  in  $LA_{SLR}(i, [X \rightarrow \alpha.]) = FOLLOW_1(X)$
    - LALR(1) parser: for all  $t$  in  $LA_{LALR}(i, [X \rightarrow \alpha.])$  (see later)

**ACTION table:**

state		-	,	a	b
0	X	X	S4	S5	
1	A	S2	*	*	
2	X	X	S4	S5	
3	R1	R1	*	*	
4	R3	R3	*	*	
5	R4	R4	*	*	
6	R2	R2	*	*	

## From DFA to parser tables: GOTO Table



- For each DFA transition  $I_i \rightarrow I_j$  reading nonterminal  $A$  (i.e.,  $I_i$  contains an item  $[X \rightarrow \alpha.A\beta]$ )
  - enter  $GOTO[i, A] = j$

**GOTO table:**

state	L	E
0	1	6
1	*	*
2	*	3
3	*	*
4	*	*
5	*	*
6	*	*



## Conflicts and Categories of LR Grammars and Parsers

## Conflict Examples in LR Grammars



- Shift – Reduce conflict:**
  - $E \rightarrow id + E$  (shift +)
  - $| id$  (reduce id)
- Reduce – Reduce conflict:**
  - $E \rightarrow id$  (reduce id)
  - $Pcall \rightarrow id$  (reduce id)
- (Shift – Accept conflict)**
  - $S' \rightarrow L$  (accept)
  - $L \rightarrow L, E$  (shift ,)

## Conflicts in LR Grammars



Observe conflicts in DFA (GOTO graph) kernels or at the latest when filling the ACTION table.

- Shift-Reduce conflict**
  - A DFA accepting state has an outgoing transition, i.e. contains items  $[X \rightarrow \alpha.]$  and  $[Y \rightarrow \beta.Z\gamma]$  for some  $Z$  in  $Nu\Sigma$
- Reduce-Reduce conflict**
  - A DFA accepting state can reduce for multiple nonterminals i.e. contains at least 2 items  $[X \rightarrow \alpha.]$  and  $[Y \rightarrow \beta.]$ ,  $X \neq Y$
- (Shift/Reduce-Accept conflict)**
  - A DFA accepting state containing  $[S' \rightarrow S.|-]$  contains another item  $[X \rightarrow \alpha.S.]$  or  $[X \rightarrow \alpha.S.\beta]$

Only for LR(0) grammars there are no conflicts.

## Handling Conflicts in LR Grammars



(Overview):

- Use lookahead
  - if lucky, the LR(0) states + a few fixed lookahead sets are sufficient to eliminate all conflicts in the LR(0)-DFA
    - SLR(1), LALR(1)
  - otherwise, use LR(1) items  $[X \rightarrow \alpha.\beta, a]$  ( $a$  is look-ahead) to build new, larger NFA/DFA
    - expensive (many items/states  $\rightarrow$  very large tables)
  - if still conflicts, may try again with  $k > 1 \rightarrow$  even larger tables
- Rewrite the grammar (factoring / expansion) and retry...
- If nothing helps, re-design your language syntax
  - Some grammars are not LR( $k$ ) for any constant  $k$  and cannot be made LR( $k$ ) by rewriting either

## Look-Ahead (LA) Sets



- For a LR(0) item  $[X \rightarrow \alpha.\beta]$  in DFA-state  $I_i$ , define **lookahead set**  $LA(I_i, [X \rightarrow \alpha.\beta])$  (a subset of  $\Sigma$ )
- For SLR(1), LALR(1) etc., the LA sets only differ for reduce items
- For SLR(1):**  
 $LA_{SLR}(I_i, [X \rightarrow \alpha.]) = \{a \text{ in } \Sigma: S' \Rightarrow^* \beta X a \gamma\} = FOLLOW_1(X)$   
 for all  $I_i$  with  $[X \rightarrow \alpha.]$  in  $I_i$ 
    - depends on nonterminal  $X$  only, not on state  $I_i$
  - For LALR(1):**  
 $LA_{LALR}(I_i, [X \rightarrow \alpha.]) = \{a \text{ in } \Sigma: S' \Rightarrow^* \beta X a w \text{ and the LR(0)-DFA started in } I_0 \text{ reaches } I_i \text{ after reading } \beta \alpha\}$ 
    - usually a subset of  $FOLLOW_1(X)$ , i.e. of SLR LA set

P. Fritsson, C. Kessler, IDA, Linköping universitet. 13 TDDD16/TDD844 Compiler Construction, 2008

## Made it simple: Is my grammar SLR(1) ?



- Construct the (LR(0)-item) characteristic NFA and its equivalent DFA (= GOTO graph) as above.
- Consider all conflicts in the DFA states:
  - Shift-Reduce:**

$[X \rightarrow \alpha.]$   
 $[Y \rightarrow \beta.b\gamma]$

 $\xrightarrow{b}$  ...  
 Consider all pairs of conflicting items  $[X \rightarrow \alpha.]$ ,  $[Y \rightarrow \beta.b\gamma]$ :  
 If  $b \text{ in } FOLLOW_1(X)$  for any of these  $\rightarrow$  not SLR(1).
  - Reduce-Reduce:**

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

  
 Consider all pairs of conflicting items  $[X \rightarrow \alpha.]$ ,  $[Y \rightarrow \beta.]$ :  
 If  $FOLLOW_1(X)$  intersects with  $FOLLOW_1(Y)$   $\rightarrow$  not SLR(1)
  - (Shift-Accept:** similar to Shift-Reduce)

P. Fritsson, C. Kessler, IDA, Linköping universitet. 14 TDDD16/TDD844 Compiler Construction, 2008

## Example: L-Values in C Language



- L-values on left hand side of assignment. Part of a C grammar:
  - $S' \rightarrow S$
  - $S \rightarrow L = R$
  - $\quad | R$
  - $L \rightarrow *R$
  - $\quad | id$
  - $R \rightarrow L$
- Avoids that  $R$  (for R-values) appears as LHS of assignments
- But  $*R = \dots$  is ok.
- This grammar is LALR(1) but not SLR(1):

P. Fritsson, C. Kessler, IDA, Linköping universitet. 15 TDDD16/TDD844 Compiler Construction, 2008

## Example (cont.)



- LR(0) parser has a shift-reduce conflict in kernel of state  $I_2$ :
- $I_0 = \{ [S' \rightarrow S], [S \rightarrow L=R], [S \rightarrow R], [L \rightarrow *R], [L \rightarrow id], R \rightarrow L \}$
  - $I_1 = \{ [S' \rightarrow S.] \}$
  - $I_2 = \{ [S \rightarrow L=R], [R \rightarrow L.] \}$  Shift = or reduce to R?
  - $I_3 = \{ [S \rightarrow R.] \}$
  - $I_4 = \{ [L \rightarrow *R], [R \rightarrow L], [L \rightarrow *R], [L \rightarrow id] \}$
  - $I_5 = \{ [L \rightarrow id.] \}$
  - $I_6 = \{ [S \rightarrow L=R], [R \rightarrow L], [L \rightarrow *R], L \rightarrow id \}$
  - $I_7 = \{ [L \rightarrow *R.] \}$
  - $I_8 = \{ [R \rightarrow L.] \}$
  - $I_9 = \{ [S \rightarrow L=R.] \}$
- $FOLLOW_1(R) = \{ |-, = \}$   $\rightarrow$  SLR(1) still shift-reduce conflict in  $I_2$  as = does not disambiguate

P. Fritsson, C. Kessler, IDA, Linköping universitet. 16 TDDD16/TDD844 Compiler Construction, 2008

## Example (cont.)



- $I_0 = \{ [S' \rightarrow S], [S \rightarrow L=R], [S \rightarrow R], [L \rightarrow *R], [L \rightarrow id], R \rightarrow L \}$
  - $I_1 = \{ [S' \rightarrow S.] \}$
  - $I_2 = \{ [S \rightarrow L=R], [R \rightarrow L.] \}$
  - $I_3 = \{ [S \rightarrow R.] \}$
  - $I_4 = \{ [L \rightarrow *R], [R \rightarrow L], [L \rightarrow *R], [L \rightarrow id] \}$
  - $I_5 = \{ [L \rightarrow id.] \}$
  - $I_6 = \{ [S \rightarrow L=R], [R \rightarrow L], [L \rightarrow *R], L \rightarrow id \}$
  - $I_7 = \{ [L \rightarrow *R.] \}$
  - $I_8 = \{ [R \rightarrow L.] \}$
  - $I_9 = \{ [S \rightarrow L=R.] \}$
- $LA_{LALR}(I_2, [R \rightarrow L.]) = \{ |-, = \} \rightarrow$  LALR(1) parser is conflict-free as computation path  $I_0 \dots I_2$  does not really allow = following R. = can only occur after R if  $*R$  was encountered before.

P. Fritsson, C. Kessler, IDA, Linköping universitet. 17 TDDD16/TDD844 Compiler Construction, 2008

## LALR(1) Parser Construction



- Method 1:** (simple but not practical)
- Construct the LR(1) items (see later). (If there is already a conflict, stop.)
  - Look for sets of LR(1) items that have the same kernel, and merge them.
  - Construct the ACTION table as for LR(1). If a conflict is detected, the grammar is not LALR(1).
  - Construct the GOTO function:  
 For each merged  $J = I_1 \cup I_2 \cup \dots \cup I_n$   
 the kernels of  $GOTO(I_1, X), \dots, GOTO(I_n, X)$  are identical because the kernels of  $I_1, \dots, I_n$  are identical.  
 Set  $GOTO(J, X) := \bigcup \{ I_i \mid I_i \text{ has the same kernel as } GOTO(I_i, X) \}$
- Method 2:** (details see textbook)
- Start from LR(0) items and construct kernels of DFA states  $I_0, I_1, \dots$
  - Compute lookahead sets by propagation along the  $GOTO(I_i, X)$  edges (fixed point iteration).

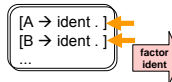
P. Fritsson, C. Kessler, IDA, Linköping universitet. 18 TDDD16/TDD844 Compiler Construction, 2008

## Solve Conflicts by Rewriting the Grammar

- Eliminate Reduce-Reduce Conflict:

### Factoring

$S \rightarrow (A) \mid (B)$   
 $A \rightarrow \text{char} \mid \text{integer} \mid \text{ident}$   
 $B \rightarrow \text{float} \mid \text{double} \mid \text{ident}$

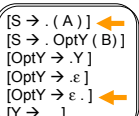


$S \rightarrow (A) \mid (B) \mid (C)$   
 $A \rightarrow \text{char} \mid \text{integer}$   
 $B \rightarrow \text{float} \mid \text{double}$   
 $C \rightarrow \text{ident}$

- Eliminate Shift-Reduce Conflict: (one token lookahead: '(')

### Inline-Expansion

$S \rightarrow (A) \mid \text{OptY} (B)$   
 $\text{OptY} \rightarrow Y \mid \epsilon$   
 $Y \rightarrow \dots$   
 $A \rightarrow \dots$   
 $B \rightarrow \dots$



$S \rightarrow (A) \mid (B) \mid Y(B)$   
 $Y \rightarrow \dots$   
 $A \rightarrow \dots$   
 $B \rightarrow \dots$

## LR(k) Grammar - Formal Definition

p.116

- Let  $G'$  be the augmented grammar for  $G$  (i.e., extended by new start symbol  $S'$  and production rule  $S' \rightarrow S \mid \epsilon$ )

- $G$  is called a **LR(k) grammar** if

- $S'_{rm} \Rightarrow^* \alpha X w_{rm} \Rightarrow \alpha \beta w$  and
- $S'_{rm} \Rightarrow^* \gamma Y x_{rm} \Rightarrow \alpha \beta y$  and
- $w[1:k] = y[1:k]$

imply that  $\alpha = \gamma$  and  $X = Y$  and  $x = y = w$ .

i.e., considering at most  $k$  symbols after the handle, in each rightmost derivation the handle can be localized and the production to be applied can be determined.

Remark:  $w, x, y$  in  $\Sigma^*$   $\alpha, \beta, \gamma$  in  $(N \cup \Sigma)^*$   $X, Y$  in  $N$

Example: see whiteboard

## Some grammars are not LR(k) for any fixed k

- Example:
 
$$S \rightarrow a B c$$

$$B \rightarrow b B b$$

$$B \rightarrow b$$
  - describes language  $\{ a b^{2N+1} c : N \geq 0 \}$

- This grammar is not LR(k) for any fixed k.

**Proof:** As k is fixed (constant), consider for any  $n > k$ :

- $S \Rightarrow^* a b^n B b^n c \Rightarrow a b^n \underline{b} b^n c$

- $S \Rightarrow^* a b^{n+1} B b^{n+1} c \Rightarrow a b^{n+1} \underline{b} b^{n+1} c$

By the LR(k) definition,

- $\alpha = a b^n$   $\beta = b$   $w = b^n c$

The handle cannot be localized with only limited lookahead size k

## No ambiguous grammar is LR(k) for any fixed k

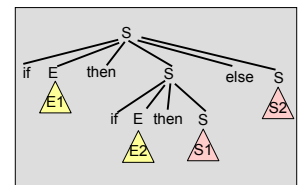
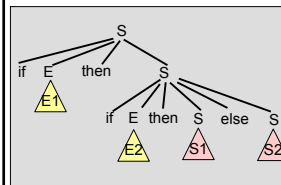
- $$S \rightarrow \text{if } E \text{ then } S$$

$$S \rightarrow \text{if } E \text{ then } S \text{ else } S$$

$$S \rightarrow \text{other statements}$$

is ambiguous – the following statement has 2 parse trees:

if E1 then if E2 then S1 else S2



## (cont.)

- Consider situation (configuration of shift-reduce parser)

... | ... if E then S else ... | ...

- Not clear whether to
  - shift else (following production 2, i.e. if E then S is not handle), or
  - reduce handle if E then S to S (following production 1)

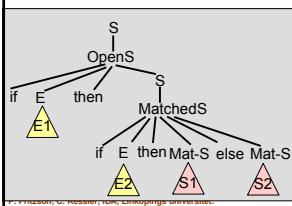
- Any fixed-size lookahead (else and beyond) does not help!

- Suggestion: Rewrite grammar to make it unambiguous

## Rewriting the grammar...

$S \rightarrow \text{MatchedS}$   
 $S \rightarrow \text{OpenS}$   
 $\text{MatchedS} \rightarrow \text{if } E \text{ then MatchedS else MatchedS}$   
 $\text{MatchedS} \rightarrow \text{other statements}$   
 $\text{OpenS} \rightarrow \text{if } E \text{ then } S$   
 $\text{OpenS} \rightarrow \text{if } E \text{ then MatchedS else OpenS}$

is no longer ambiguous



Impossible now to derive any sentential form containing an OpenS nonterminal from a MatchedS

## Some grammars are not LR( $k$ ) for any fixed $k$



- Grammar with productions
$$S \rightarrow a S a \mid \epsilon$$
is unambiguous but not LR( $k$ ) for any fixed  $k$ . (Why?)
- An equivalent LR grammar for the same language is
$$S \rightarrow a a S \mid \epsilon$$

## LR(1) Items and LR( $k$ ) Items



**LR( $k$ ) parser:** Construction similar to LR(0) / SLR(1) parser, but plan for distinguishing between states for  $k > 0$  tokens **lookahead** already from the beginning

- States in the LR(0) GOTO graph may be split up
- LR(1) items:
$$[A \rightarrow \alpha \cdot \beta, a]$$
 for all productions  $A \rightarrow \alpha \beta$  and all  $a$  in  $\Sigma$
- Can be combined for lookahead symbols with equal behavior:
$$[A \rightarrow \alpha \cdot \beta, a|b]$$
 or  $[A \rightarrow \alpha \cdot \beta, L]$  for a subset  $L$  of  $\Sigma$
- Generalized to  $k > 1$ :
$$[A \rightarrow \alpha \cdot \beta, a_1 a_2 \dots a_k]$$

**Interpretation of  $[A \rightarrow \alpha \cdot \beta, a]$  in a state:**

- If  $\beta$  not  $\epsilon$ , ignore second component (as in LR(0))
- If  $\beta = \epsilon$ , i.e.  $[A \rightarrow \alpha \cdot a]$ , reduce **only if next input symbol = a**

## LR(1) Parser



- NFA start state is  $[S' \rightarrow \cdot S, [-]$
- Modify computation of *Closure(I)*, *GOTO(I,X)* and the subset computation for LR(1) items
  - Details see [ASU86, p.232] or [ALSU06, p.261]
- Can have many more states than LR(0) parser
  - Which may help to resolve some conflicts

## Interesting to know...



- For each LR( $k$ ) grammar with some constant  $k > 1$  there exists an equivalent\* grammar  $G'$  that is LR(1).
- For any LL( $k$ ) grammar there exists an equivalent LR( $k$ ) grammar (but not vice versa!)
  - e.g., language  $\{a^n b^n : n > 0\} \cup \{a^n c^n : n > 0\}$  has a LR(0) grammar but no LL( $k$ ) grammar for any constant  $k$ .
- Some grammars are LR(0) but not LL( $k$ ) for any  $k$ 
  - e.g.,  $S \rightarrow A b$   
 $A \rightarrow A a \mid a$  (left recursion, could be rewritten)

\* Two grammars are *equivalent* if they describe the same language.