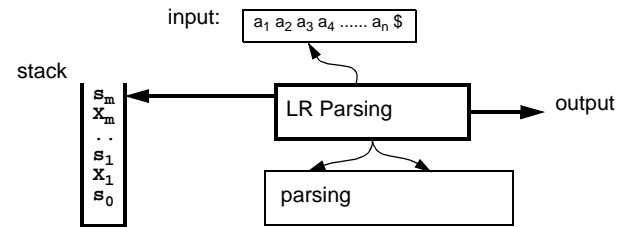


## Parser Generation

- **Main Problem:** given a grammar  $G$ , how to build a **top-down parser** or a **bottom-up parser** for it ?
- **parser** : a program that, given a sentence, reconstructs a derivation for that sentence ---- if done successfully, it "recognize" the sentence
- all parsers read their input **left-to-right**, but construct parse tree differently.
- **bottom-up parsers** --- construct the tree from leaves to root
  - shift-reduce, LR, SLR, LALR, operator precedence
- **top-down parsers** --- construct the tree from root to leaves
  - recursive descent, predictive parsing, LL(1)

## Bottom-Up Parsing

- Construct parse tree "bottom-up" --- from leaves to the root
- Bottom-up parsing always constructs **right-most derivation**
- Important parsing algorithms: **shift-reduce**, **LR parsing**
- **LR parser** components: input, stack (strings of grammar symbols and states), driver routine, parsing tables.



## LR Parsing

- A sequence of new **state** symbols  $s_0, s_1, s_2, \dots, s_m$  ---- each state summarize the information contained in the stack below it.
- **Parsing configurations:** (**stack**, **remaining input**) written as  $(s_0 X_1 s_1 X_2 s_2 \dots X_m s_m, a_i a_{i+1} a_{i+2} \dots a_n \$)$ 
  - next "move" is determined by  $s_m$  and  $a_i$
- **Parsing tables:** ACTION[s,a] and GOTO[s,X]

Table A ACTION[s,a] --- s : state, a : terminal

its entries (1) shift  $s_k$  (2) reduce  $A \rightarrow \beta$   
 (3) accept (4) error

Table G GOTO[s,X] --- s : state, X : non-terminal  
 its entries are states

## Constructing LR Parser

**How to construct the parsing table ACTION and GOTO ?**

- **basic idea:** first construct DFA to recognize handles, then use DFA to construct the parsing tables ! different parsing table yield different LR parsers **SLR(1)**, **LR(1)**, or **LALR(1)**
- **augmented grammar** for context-free grammar  $G = G(T, N, P, S)$  is defined as  $G' = G'(T, N \cup \{S'\}, P \cup \{S' \rightarrow S\}, S')$  ---- adding non-terminal  $S'$  and the production  $S' \rightarrow S$ , and  $S'$  is the new start symbol. When  $S' \rightarrow S$  is reduced, parser accepts.
- **LR(0) item** for productions of a context-free grammar  $G$  ---- is a production with **dot** at some position in the r.h.s.

For  $A \rightarrow XYZ$ , its items are  $A \rightarrow .XYZ$   $A \rightarrow X.YZ$   
 $A \rightarrow XY.Z$   $A \rightarrow XYZ.$

For  $A \rightarrow \epsilon$ , its items are just  $A \rightarrow .$

## LR(0) items and LR(0) DFA

- Informally, item  $A \rightarrow X.YZ$  means a string derivable from  $X$  has been seen, and one from  $YZ$  is expected. **LR(0) items** are used as state names for LR(0) DFA or LR(0) NFA that recognizes **viable prefixes**.
- Viable prefixes** of a CFG are prefixes of right-sentential forms with no symbols to right of the **handle**; we can always add terminals on right to form a right-sentential form.
- Two way to construct the LR(0) DFA:**
  - first construct LR(0) NFA and then convert it to a DFA !
  - construct the LR(0) DFA directly !
- From LR(0) DFA to the Parsing Table** ----- transition table for the DFA is the GOTO table; the states of DFA are states of the parser.

## Example: LR(0) Items

**CFG Grammar:**

$E \rightarrow E + T$		$T$
$T \rightarrow T * F$		$F$
$F \rightarrow ( E )$		$id$

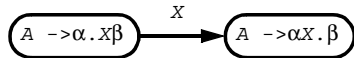
**Augmented Grammar:**

$E' \rightarrow E$		
$E \rightarrow E + T$		$T$
$T \rightarrow T * F$		$F$
$F \rightarrow ( E )$		$id$

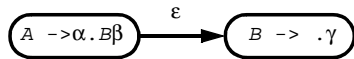
**LR(0) terms:**

$E' \rightarrow . E$		$T \rightarrow . T * F$		$F \rightarrow ( E . )$
$E' \rightarrow E .$		$T \rightarrow T . * F$		$F \rightarrow ( E ) .$
$E \rightarrow . E + T$		$T \rightarrow T * . F$		$F \rightarrow . id$
$E \rightarrow E . + T$		$T \rightarrow T * F .$		$F \rightarrow id .$
$E \rightarrow E + . T$		$T \rightarrow . F$		
$E \rightarrow E + T .$		$T \rightarrow F .$		
$E \rightarrow . T$		$F \rightarrow . ( E )$		
$E \rightarrow T .$		$F \rightarrow ( . E )$		

## From LR(0) NFA to LR(0) DFA

- Construct **LR(0) NFA** with all LR(0) items of  $G$  as states, connect states by moving the **dot**; **final states** are those with dots at the end.
- for each item  $A \rightarrow \alpha.X\beta$ 

- for each pair  $A \rightarrow \alpha.B\beta$ ,  $B \rightarrow .\gamma$ 

(expect to see a string derivable from  $\gamma$ )


- Convert NFA to DFA using **subset construction algorithm**.
- The states of the resulting **LR(0) DFA** ---  $C = \{I_1, I_2, \dots, I_n\}$  are called **canonical LR(0) collection** for grammar  $G'$
- Disadvantage:** the NFA is often huge, and converting from NFA to DFA is tedious and time-consuming.

## Building LR(0) DFA Directly

- Instead of building DFA from NFA, **we can build the LR(0) DFA directly.**
- Given a set of LR(0) items  $I$ , **CLOSURE(I)** is defined as
 

```
repeat
  for each item  $A \rightarrow \alpha.B\beta$  in  $I$  and
  each production  $B \rightarrow \gamma$ 
    add  $B \rightarrow .\gamma$  to  $I$ , if it's not in  $I$ 
until  $I$  does not change
```
- GOTO(I, X)** is defined as
 

```
CLOSURE(all items  $A \rightarrow \alpha X . \beta$  for each  $A \rightarrow \alpha . X \beta$  in  $I$ )
```
- Canonical LR(0) collection** is computed by the following procedure:
 

```
 $I_0 = \text{CLOSURE}(\{S' \rightarrow .S\})$  and  $C = \{I_0\}$ 
repeat
  for each  $I \in C$  and grammar symbol  $X$ 
     $T = \text{GOTO}(I, X)$ ; if  $T \neq \emptyset$  and  $T \notin C$  then  $C = C \cup \{T\}$ ;
until  $C$  does not change
```

Resulting LR(0) DFA: C is the set of states; GOTO is the transition table

## Constructing SLR(1) Parsing Table

- From the LR(0) DFA, we can construct the parsing table ----- SLR(1) parsing table. The parser based on SLR(1) parsing table is called SLR(1) parser. The SLR(1) grammars are those whose SLR(1) parsing table does not contain any conflicts.

- Algorithm --- use  $C = \{I_0, \dots, I_n\}$ , GOTO, FOLLOW:

- If  $A \rightarrow a.a\beta$  is in  $I_i$  and  $GOTO(I_i, a) = I_j$  where  $a$  is a terminal, set  $ACTION[i, a]$  to "shift  $j$ ".
- If  $A \rightarrow \alpha.$  is in  $I_i$ , set  $ACTION[i, a]$  to "reduce  $A \rightarrow \alpha$ " for all terminal  $a$  in  $FOLLOW(A)$ .
- If  $S' \rightarrow S.$  is in  $I_i$ , set  $ACTION[i, \$]$  to "accept"
- If  $GOTO(I_i, A) = I_j$ , set  $GOTO[i, A] = j$
- set all other entries to "error"
- set initial state to be  $I_0$  with  $S' \rightarrow .S$

## Limitation of SLR(1) Parser

- Unfortunately, many unambiguous grammars are not SLR(1) grammars

$S \rightarrow L = R \mid R$   
 $L \rightarrow *R \mid id$   
 $R \rightarrow L$

L means "l-value"  
 R means "r-value"  
 \* means "contents of"

### Canonical LR(0) collection ---

I0: $S' \rightarrow .S$	I3: $S \rightarrow R.$	I6: $S \rightarrow L.=R$
$S \rightarrow .L=R$		$R \rightarrow .L$
$S \rightarrow .R$	I4: $L \rightarrow *.R$	$L \rightarrow .*R$
$L \rightarrow .*R$	$R \rightarrow .L$	$L \rightarrow .id$
$L \rightarrow .id$	$L \rightarrow .*R$	
$R \rightarrow .L$	$L \rightarrow .id$	I7: $L \rightarrow *R.$
I1: $S' \rightarrow S.$	I5: $L \rightarrow id.$	I8: $R \rightarrow L.$
I2: $S \rightarrow L.=R$		I9: $S \rightarrow L=R.$
$R \rightarrow L.$	$FOLLOW(R) = \{=, \dots\}$	

state 2 has a shift/reduce conflict on "=" : shift 6 or reduce R -

## LR(1) Parsing

- Conflict arises because LR(0) states do not encode enough left context - in the previous example, reduction  $R \rightarrow L$  is wrong upon input = because " $R = \dots$ " never appears in right-sentential form.

- Solution: split LR(0) states by adding terminals to states, for example,  $[A \rightarrow \alpha. , a]$  results in reduction only if next symbol is  $a$ .

- An LR(1) term is in the form of  $[A \rightarrow \alpha.\beta , a]$  where  $A \rightarrow \alpha\beta$  is a production and  $a$  is a terminal or  $\$$

- To build LR(1) parsing table --- we first build LR(1) DFA --- then construct the parsing table using the same SLR(1) algorithm except

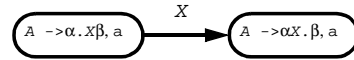
- only if  $[A \rightarrow \alpha. , a]$  is in  $I_i$ , then set  $ACTION[i, a]$  to "reduce  $A \rightarrow \alpha$ "

- To way to build LR(1) DFA ---- from NFA -> DFA or build DFA directly

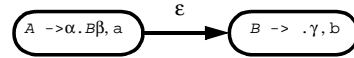
## Building LR(1) DFA

- Construct **LR(1) NFA** with all LR(1) items of  $G$  as states, connect states by moving the **dot**; then convert the NFA to DFA.

1. for each item  $[A \rightarrow \alpha.X\beta, a]$



2. for each pair  $[A \rightarrow \alpha.B\beta, a]$   $B \rightarrow \cdot\gamma$  and  $b$  in  $\text{FIRST}(\beta a)$ .



- Construct the LR(1) DFA directly (see the Dragon book)
- Given a set of LR(1) items  $I$ ,  $\text{CLOSURE}(I)$  is now defined as

```

repeat
  for each item  $[A \rightarrow \alpha.B\beta, a]$  in  $I$  and
    each production  $B \rightarrow \gamma$  and each terminal  $b$  in  $\text{FIRST}(\beta a)$ 
      add  $[B \rightarrow \cdot\gamma, a]$  to  $I$ , if it's not in  $I$ 
until  $I$  does not change
  
```

## Constructing LR(1) Parser

- Canonical LR(1) collection** is computed by the following procedure:

```

I0 = CLOSURE({S' -> .S, $}) and C = {I0}
repeat
  for each I in C and grammar symbol X
    T = GOTO(I, X); if T != empty and T not in C then C = C union { T };
until C does not change
  
```

Resulting LR(1) DFA:  $C$  is the set of states; GOTO is the transition table

- From the LR(1) DFA, we can construct the parsing table ----- LR(1) parsing table. The parser based on LR(1) parsing table is called LR(1) parser. The LR(1) grammars are those whose LR(1) parsing table does **not** contain any conflicts (no duplicate entries).

- Example:**

```

S' -> S
S -> C C
C -> c C | d
  
```

## LALR(1) Parsing

- Bad News:** LR(1) parsing tables are too big; for PASCAL, SLR tables has about 100 states, LR table has about 1000 states.
- LALR (LookAhead-LR) parsing tables have same number of states as SLR, but use lookahead for reductions. The LALR(1) DFA can be constructed from the LR(1) DFA.
- LALR(1) states can be constructed from LR(1) states by merging states with same **core**, or same LR(0) items, and **union** their lookahead sets.

```

Merging I8: C -> c.C, c/d      I9: C -> cC., $
into a new state I89: C -> cC., c/d/$
  
```

```

Merging I3: C -> c.C, c/d      I6: C -> c.C, $
      C -> .cC, c/d          C -> .cC, $
      C -> .d, c/d          C -> .d, $
into a new state I36: C -> c.C, c/d/$
      C -> .cC, c/d/$
      C -> .d, c/d/$I
  
```

## LALR(1) Parsing (cont'd)

- From the LALR(1) DFA, we can construct the parsing table ----- LALR(1) parsing table. The parser based on LALR(1) parsing table is called LALR(1) parser. The LALR(1) grammars are those whose LALR(1) parsing table does **not** contain any conflicts (no duplicate entries).
- LALR(1) DFA and LALR(1) parsing table can be constructed without creating LR(1) DFA --- see Dragon book for detailed algorithm.
- LALR parser makes same number of moves as LR parser on **correct** input.
- On **incorrect** input, LALR parser may make erroneous reductions, but will signal "error" before shifting input, i.e., merging states makes reduce determination "less accurate", but has no effect on shift actions.

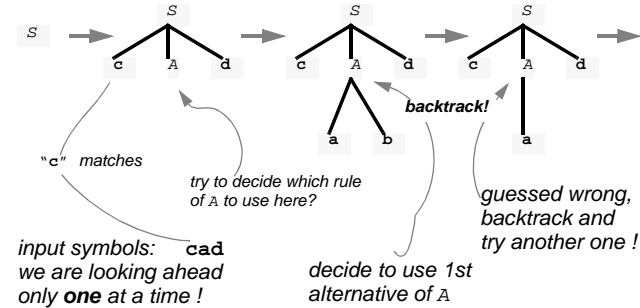
## Summary: LR Parser

- Relation of three LR parsers:  $LR(1) > LALR(1) > SLR(1)$
- Most programming language constructs are **LALR(1)**. The **LR(1)** is unnecessary in practice, but the **SLR(1)** is not enough.
- YACC is an LALR(1) Parser Generator.
- When parsing ambiguous grammars using LR parsers, the parsing table will contain multiple entries. We can specify the **precedence** and **associativity** for terminals and productions to resolve the conflicts. YACC uses this trick.
- **Other Issues** in parser implementation: 1. compact representation of parsing table 2. error recovery and diagnosis.

## Top-Down Parsing

- **Starting from the start symbol** and “guessing” which production to use next step. It often uses **next input token** to guide “guessing”.

example:  $S \rightarrow c A d$   
 $A \rightarrow ab \mid a$



## Top-Down Parsing (cont'd)

- Typical implementation is to write a **recursive procedure** for each **non-terminal** (according to the r.h.s. of each grammar rule)

**Grammar:**

$E \rightarrow T E'$   
 $E' \rightarrow + T E' \mid \epsilon$   
 $T \rightarrow F T'$   
 $T' \rightarrow * F T' \mid \epsilon$   
 $F \rightarrow id \mid ( E )$

```

advance  sets c to next input token
err      reports error message

fun e() = (t(); eprime())
and eprime() = if (c = "+")
  then (advance(); t(); eprime())
and t() = (f(); tprime())
and tprime() = if (c = "*")
  then (advance(); f(); tprime())
and f() =
  (if (c = id) then advance()
   else if (c = "(") then
     (advance(); e();
      if (c = ")") then advance() else err())
   else err())
    
```

## Recursive Descent Parsing

- The previously referred top-down parsing method is often called **recursive descent parsing**!

• **Main challenges:**

1. **back-tracking** is messy, difficult and inefficient  
 (solution: use input “lookahead” to help make the right choice)
2. **more alternatives** --- even if we use one lookahead input char, there are still > 1 rules to choose ---  $A \rightarrow ab \mid a$   
 (solution: rewrite the grammar by **left-factoring**)
3. **left-recursion** might cause infinite loop  
 what is the procedure for  $E \rightarrow E + E$  ?  
 (solution: rewrite the grammar by **eliminating left-recursions**)
4. **error handling** --- errors detected “far away” from actual source.



## Constructing Predictive Parser

- The key is to build the parse table  $M[A, a]$

```

for each production  $A \rightarrow \alpha$  do
  for each  $a \in \text{FIRST}(\alpha)$  do
    add  $A \rightarrow \alpha$  to  $M[A, a]$ 
  if  $\epsilon \in \text{FIRST}(\alpha)$  then
    for each  $b \in \text{FOLLOW}(A)$  do
      add  $A \rightarrow \alpha$  to  $M[A, b]$ 

```

rest of  $M$  is error

- $\text{FIRST}(\alpha)$  is a set of terminals (plus  $\epsilon$ ) that **begin** strings derived from  $\alpha$ , where  $\alpha$  is any string of non-terminals and terminals.
- $\text{FOLLOW}(A)$  is a set of terminals that can **follow**  $A$  in a sentential form, where  $A$  is any non-terminal

## First & Follow

- To compute  $\text{FIRST}(X)$  for any grammar symbol  $X$ :

```

FIRST(X) = {X}, if X is a terminal;
FIRST(X) = FIRST(X) ∪ {a}, if X → aα;
FIRST(X) = FIRST(X) ∪ {ε}, if X → ε; and
FIRST(X) = FIRST(X) ∪ FIRST(Y1Y2...Yk),
if X → Y1Y2...Yk.

```

repeat until nothing new is added to any  $\text{FIRST}$

- $\text{FIRST}(Y_1Y_2...Y_k) = \text{FIRST}(Y_1) - \{\epsilon\}$   
 $\cup \text{FIRST}(Y_2) - \{\epsilon\}$  if  $\epsilon \in \text{FIRST}(Y_1)$   
 $\cup \text{FIRST}(Y_3) - \{\epsilon\}$  if  $\epsilon \in \text{FIRST}(Y_1Y_2)$   
.....  
 $\cup \text{FIRST}(Y_k) - \{\epsilon\}$  if  $\epsilon \in \text{FIRST}(Y_1...Y_{k-1})$   
 $\cup \{\epsilon\}$  if all  $\text{FIRST}(Y_i)_{(i=1,...,k)}$  contain  $\epsilon$

## First & Follow (cont'd)

- To compute  $\text{FOLLOW}(X)$  for any non-terminal  $X$ :

$\text{FOLLOW}(S) = \text{FOLLOW}(S) \cup \{\$, \}$ , if  $S$  is start symbol;

$\text{FOLLOW}(B) = \text{FOLLOW}(B) \cup (\text{FIRST}(\beta) - \{\epsilon\})$ ,  
if  $A \rightarrow \alpha B \beta$  and  $\beta \neq \epsilon$

$\text{FOLLOW}(B) = \text{FOLLOW}(B) \cup \text{FOLLOW}(A)$   
if  $A \rightarrow \alpha B$  or  $A \rightarrow \alpha B \beta$  and  $\epsilon \in \text{FIRST}(\beta)$

- Example:

$\text{FOLLOW}(E) = \{\$, \}$ ,	<b>reason</b>
$= \{\$ \}$ ,	$E$ is start symbol
$\text{FOLLOW}(E') = \text{FOLLOW}(E)$	$F \rightarrow ( E )$
$= \{\$ \}$ ,	$E \rightarrow T E'$

(Read Appel pp 47 - 53 for detailed examples)

## Summary: LL(1) Grammars

- A grammar is **LL(1)** if parsing table  $M[A, a]$  has no duplicate entries, which is equivalent to specifying that for each production

$$A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$$

- All  $\text{FIRST}(\alpha_i)$  are disjoint.
  - At most one  $\alpha_i$  can derive  $\epsilon$ ; in that case,  $\text{FOLLOW}(A)$  must be disjoint from  $\text{FIRST}(\alpha_1) \cup \text{FIRST}(\alpha_2) \cup \dots \cup \text{FIRST}(\alpha_n)$
- Left-recursion and ambiguity grammar lead to **multiple entries** in the parsing table. (try the **dangling-else** example)
  - The **main difficulty** in using (top-down) **predictive parsing** is in rewriting a grammar into an LL(1) grammar. There is no general rule on how to resolve **multiple entries** in the parsing table.