



中山大學
SUN YAT-SEN UNIVERSITY



国家超级计算广州中心
NATIONAL SUPERCOMPUTER CENTER IN GUANGZHOU

Compilation Principle 编译原理

第8讲：语法分析(5)

张献伟

xianweiz.github.io

DCS290, 3/25/2021



中山大學
SUN YAT-SEN UNIVERSITY



Review Questions (1)

- What are the parts of a table-driven predictive parser?

Input buffer, stack, parse table and a driver

- What are the operations on the stack?

Expand the non-terminal, match the terminal

- How to predict the next production to use?

Next input symbol, current nonterminal being processed

- What does LL(k) mean?

L: scans the input from left to right

L: produces a leftmost derivation

k: using k input symbols of lookahead

- How to build the LL(1) parse table?

Two sets: FIRST, FOLLOW

Review Questions (2)

- Which one is typically used, LL(0), LL(1), LL(2) ...? Why not others?

LL(1). LL(0) is too weak, LL(k) has a too large table

- Which are the key differences between top-down and bottom-up parsing?

Top-down is based on leftmost derivation;
bottom-up is the reverse of rightmost derivation.

- What are the key operations of bottom-up parsing?

Shift: pushes a terminal on the stack

Reduce: pops RHS and pushes LHS

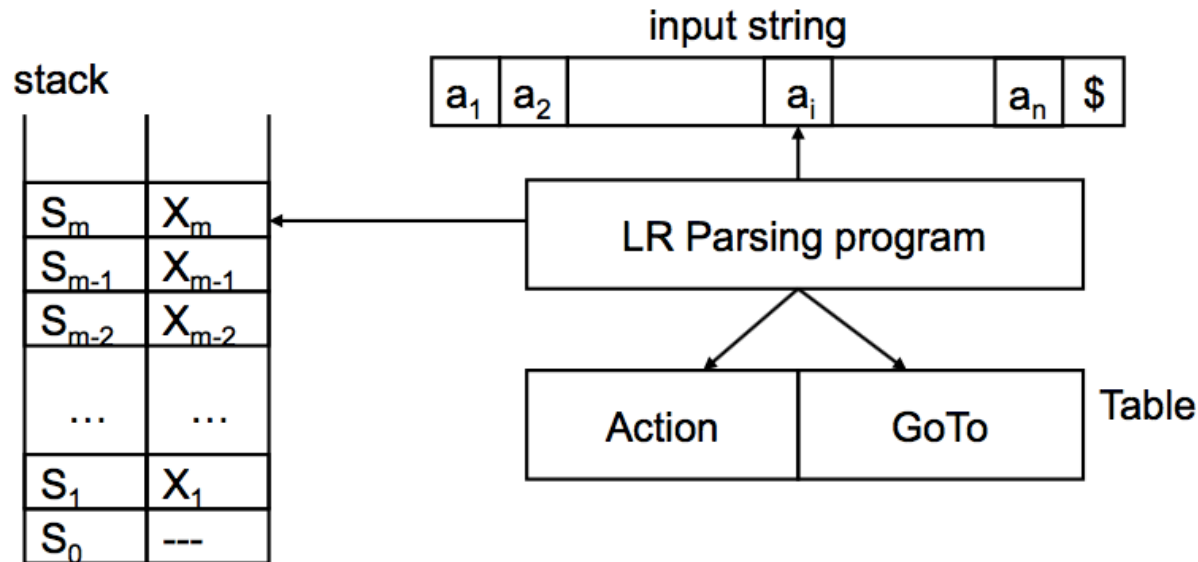
Types of Bottom-Up Parsers

- Types of bottom up parsers
 - Simple precedence parsers
 - Operator precedence parsers
 - Recursive ascent parsers
 - LR family parsers
 - ...
- In this course, we will only discuss **LR family parsers**
 - Efficient, table-driven shift-reduce parsers
 - Most automated tools for bottom-up parsing generate LR family

LR(k) Parser

- **LR(k)**: member of LR family of parsers
 - **L**: scan input from left to right
 - **R**: construct a rightmost derivation in reverse
 - **k**: number of input symbols of lookahead to make decisions
 - $k = 0$ or 1 are of particular interests, is assumed to be 1 when omitted
- Comparison with LL(k) parser
 - Efficient as LL(k)
 - Linear in time and space to length of input (same as LL(k))
 - Convenient as LL(k)
 - Can generate automatically from grammar – YACC, Bison
 - More complex than LL(k)
 - Harder to debug parser when grammar causes conflicting predictions
 - More powerful than LL(k)
 - Handles more grammars: no left recursion removal, left factoring needed
 - Handles more (and most practical) languages: $LL(1) \subset LR(1)$

LR Parser



- The stack holds a sequence of states, $s_0s_1\dots s_m$ (s_m is the top)
 - States are to track where we are in a parse
 - Each grammar symbol X_i is associated with a state s_m
- Contents of stack + input ($X_1X_2\dots X_ma_i\dots a_n$) is a right sentential form
 - If the input string is a member of the language
- Uses $[S_m, a_i]$ to index into parsing table to determine action

Parse Table

- LR parsers use two tables: **action table** and **goto table**
 - The two tables are usually combined
 - Action table specifies entries for terminals
 - Goto table specifies entries for non-terminals
- Action table[动作表]
 - $\text{Action}[s, a]$ tells the parser what to do when the state on top of the stack is s and terminal a is the next input token
 - Possible actions: **shift, reduce, accept, error**
- Goto table[跳转表]
 - $\text{Goto}[s, X]$ indicates the new state to place on top of the stack after a reduction of the non-terminal X while state s is on top of the stack

Possible Actions[可能动作]

- **Shift**

- Transfer the next input symbol onto the top of the stack

- **Reduce**

- If there's a rule $A \rightarrow w$, and if the contents of stack are qw for some q (q may be empty), then we can reduce the stack to qA

- **Accept**

- The special case of reduce: reducing the entire contents of stack to the start symbol with no remaining input
- Last step in a successful parse: have recognized input as a valid sentence

- **Error**

- Cannot reduce, and shifting would create a sequence on the stack that cannot eventually be reduced to the start symbol

Possible Actions (cont.)

- Grammar

$S \rightarrow E$

$E \rightarrow T \mid E + T$

$T \rightarrow \text{id} \mid (E)$

- Input: (id + id)

– $\#(\text{id} + \text{id})\$ \Rightarrow (\text{id}\# + \text{id})\$ \Rightarrow (T\# + \text{id})\$ \Rightarrow (E\# + \text{id})\$ \Rightarrow (E + \text{id}\#)\$ \Rightarrow (E + T\#)\$ \Rightarrow (E\#)\$ \Rightarrow (E)\#\$ \Rightarrow T\#\$ \Rightarrow E\#\$ \Rightarrow S\#\$$

- Input: id+)

– $\#\text{id}+)\$ \Rightarrow \text{id}\#+)\$ \Rightarrow T\#+)\$ \Rightarrow E\#+)\$ \Rightarrow E + \#)\$ \dots$

Example: Parse Table

Grammar:

(1) $S \rightarrow BB$

(2) $B \rightarrow aB$

(3) $B \rightarrow b$

String: bab

State	ACTION			GOTO	
	a	b	\$	S	B
0	s3	s4		1	2
1			acc		
2	s3	s4			5
3	s3	s4			6
4	r3	r3	r3		
5	r1	r1	r1		
6	r2	r2	r2		

- Table entry:

- si : shifts the input symbol and moves to state i (i.e., push state on stack)
- rj : reduce by production numbered j
- acc: accept
- blank: error

Example: Parse Table (cont.)

Grammar:

(1) $S \rightarrow BB$

(2) $B \rightarrow aB$

(3) $B \rightarrow b$

String: bab

State	ACTION			GOTO	
	a	b	\$	S	B
0	s3	s4		1	2
1			acc		
2	s3	s4			5
3	s3	s4			6
4	r3	r3	r3		
5	r1	r1	r1		
6	r2	r2	r2		

state \rightarrow 0 4

b a b

symbol \rightarrow \$ b

b a b \$

Example: Parse Table (cont.)

Grammar:

(1) $S \rightarrow BB$

(2) $B \rightarrow aB$

(3) $B \rightarrow b$

String: **bab**

State	ACTION			GOTO	
	a	b	\$	S	B
0	s3	s4		1	2
1			acc		
2	s3	s4			5
3	s3	s4			6
4	r3	r3	r3		
5	r1	r1	r1		
6	r2	r2	r2		

B

b

a

b

state → 0 4

symbol → \$ B

a b \$

Example: Parse Table (cont.)

Grammar:

(1) $S \rightarrow BB$

(2) $B \rightarrow aB$

(3) $B \rightarrow b$

String: bab

State	ACTION			GOTO	
	a	b	\$	S	B
0	s3	s4		1	2
1			acc		
2	s3	s4			5
3	s3	s4			6
4	r3	r3	r3		
5	r1	r1	r1		
6	r2	r2	r2		

B

b

a

b

state \rightarrow 0 2 3 4

symbol \rightarrow \$ B a b

a b \$

Example: Parse Table (cont.)

Grammar:

(1) $S \rightarrow BB$

(2) $B \rightarrow aB$

(3) $B \rightarrow b$

String: **bab**

State	ACTION			GOTO	
	a	b	\$	S	B
0	s3	s4		1	2
1			acc		
2	s3	s4			5
3	s3	s4			6
4	r3	r3	r3		
5	r1	r1	r1		
6	r2	r2	r2		

$\begin{array}{cc} B & B \\ | & | \\ b & a \end{array}$

state \rightarrow 0 2 3 4
 symbol \rightarrow \$ B a B \$

Example: Parse Table (cont.)

Grammar:

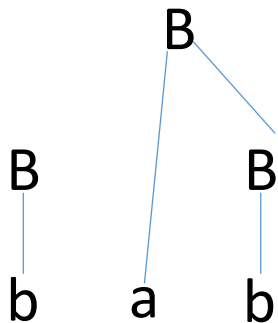
(1) $S \rightarrow BB$

(2) $B \rightarrow aB$

(3) $B \rightarrow b$

String: bab

State	ACTION			GOTO	
	a	b	\$	S	B
0	s3	s4		1	2
1			acc		
2	s3	s4			5
3	s3	s4			6
4	r3	r3	r3		
5	r1	r1	r1		
6	r2	r2	r2		



state \rightarrow 0 2 3 6
 symbol \rightarrow \$ B ~~B~~ B \$

Example: Parse Table (cont.)

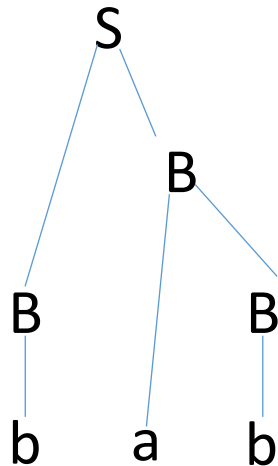
Grammar:

(1) $S \rightarrow BB$

(2) $B \rightarrow aB$

(3) $B \rightarrow b$

String: **bab**



State	ACTION			GOTO	
	a	b	\$	S	B
0	s3	s4		1	2
1			acc		
2	s3	s4			5
3	s3	s4			6
4	r3	r3	r3		
5	r1	r1	r1		
6	r2	r2	r2		

state \rightarrow 0 ~~2~~ 5
 symbol \rightarrow \$ ~~B~~ B

\$

Parser Actions

Initial

s_0	
$\$$	$a_1 a_2 \dots a_n \$$

General

$s_0 s_1 \dots s_m$	
$\$ X_1 \dots X_m$	$a_i a_{i+1} \dots a_n \$$

- If $\text{ACTION}[s_m, a_i] = s_x$, then do shift

- Pushes a_i on stack
 - a_i is removed from input
- Enters state x
 - i.e., pushes state x on stack

$s_0 s_1 \dots s_m x$	
$\$ X_1 \dots X_m a_i$	$a_{i+1} \dots a_n \$$

Parser Actions (cont.)

Initial

s_0

\$

$a_1 a_2 \dots a_n \$$

General

$s_0 s_1 \dots s_m$

$\$X_1 \dots X_m$

$a_i a_{i+1} \dots a_n \$$

- If $\text{ACTION}[s_m, a_i] = r_x$, (i.e., the x^{th} production: $A \rightarrow X_{m-(k-1)} \dots X_m$), then do reduce
 - Pops k symbols from stack
 - Pushes A on stack
 - No change on input
 - $\text{GOTO}[s_{m-k}, A] = y$, then

$s_0 s_1 \dots s_{m-k}$

$\$X_1 \dots X_{m-k} A$

$a_i a_{i+1} \dots a_n \$$



$s_0 s_1 \dots s_{m-k} y$

$\$X_1 \dots X_{m-k} A$

$a_i a_{i+1} \dots a_n \$$

Parser Actions (cont.)

Initial

s_0
 $\$$ $a_1 a_2 \dots a_n \$$

General

$s_0 s_1 \dots s_m$
 $\$X_1 \dots X_m$ $a_i a_{i+1} \dots a_n \$$

- If $\text{ACTION}[s_m, a_i] = \text{acc}$, then parsing is complete
- If $\text{ACTION}[s_m, a_i] = \langle \text{empty} \rangle$, then report error and stop

LR Parsing Program

- **Input:** input string ω and parse table with ACTION/GOTO
- **Output:** reduction steps ω 's bottom-up parse, or error
- **Initial:** s_0 on the stack, $\omega\$$ in the input buffer

```
let  $a$  be the first symbol of  $\omega\$$ 
while (1) { /* repeat forever */
    let  $s$  be the state on top of the stack;
    if (ACTION[ $s,a$ ] = shift  $t$ ) {
        push  $t$  onto the stack;
        let  $a$  be the next input symbol;
    } else if (ACTION[ $s,a$ ] = reduce  $A \rightarrow \beta$ ) {
        pop  $|\beta|$  symbols off the stack;
        let state  $t$  now be on top of the stack;
        push GOTO[ $t,A$ ] onto the stack;
        output the production  $A \rightarrow \beta$ ;
    } else if (ACTION[ $s,a$ ] = accept) break; /* parsing is done */
    else call error-recovery routine;
```

Construct Parse Table

- Construct parsing table: identify the possible states and arrange the transitions among them
- **LR(0)** parsing
 - Simplest LR parsing, only considers stack to decide shift/reduce
 - Weakest, not used much in practice because of its limitations
- **LR(1)** parsing
 - LR parser that considers next token (lookahead of 1)
 - Compared to LR(0), more complex alg and much bigger table
- **SLR(1)** parsing
 - Simple LR, lookahead from first/follow rules derived from LR(0)
 - Keeps table as small as LR(0)
- **LALR(1)** parsing
 - Lookahead LR(1): fancier lookahead analysis using the same LR(0) automaton as SLR(1)

Item[項目]

- An **item** is a production with a “.” somewhere on the RHS
 - Dot indicates extent of RHS already seen in the parsing process
 - The only item for $X \rightarrow \varepsilon$ is $X \rightarrow \cdot$
 - Items are often called “**LR(0) items**” (a.k.a., **configuration**)
- The items for $A \rightarrow XYZ$ are
 - $A \rightarrow \cdot XYZ$
 - Indicates that we hope to see a string derivable from XYZ next on input
 - $A \rightarrow X \cdot YZ$
 - Indicates that we have just seen on the input a string derivable from X and that we hope next to see a string derivable from YZ
 - $A \rightarrow XY \cdot Z$
 - $A \rightarrow XYZ \cdot$
 - Indicates that we have seen the body XYZ and that it may be time to reduce XYZ to A

State[状态]

- Example:
 - Suppose we are currently in this position
 $A \rightarrow X \cdot YZ$
 - We have just recognized X and expect the upcoming input to contain a sequence derivable from YZ (say, $Y \rightarrow u|w$)
 - Y is further derivable from either u or w
 $A \rightarrow X \cdot YZ$
 $Y \rightarrow \cdot u$
 $Y \rightarrow \cdot w$
 - The above three items can be placed into a set, called as **configuration set** of the LR parser
- Parsing tables have one **state** corresponding to each set
 - The states can be modeled as a finite automaton where we move from one state to another via transitions marked with a symbol of the CFG

Augmented Grammar[增广文法]

- We want to start with an item with a dot before the start symbol S and move to an item with a dot after S
 - Represents shifting and reducing an entire sentence of the grammar
 - Thus, we need S to appear on the right side of a production
 - Only one 'acc' in the table
- Modify the grammar by adding the production $S' \rightarrow \cdot S$

Grammar:

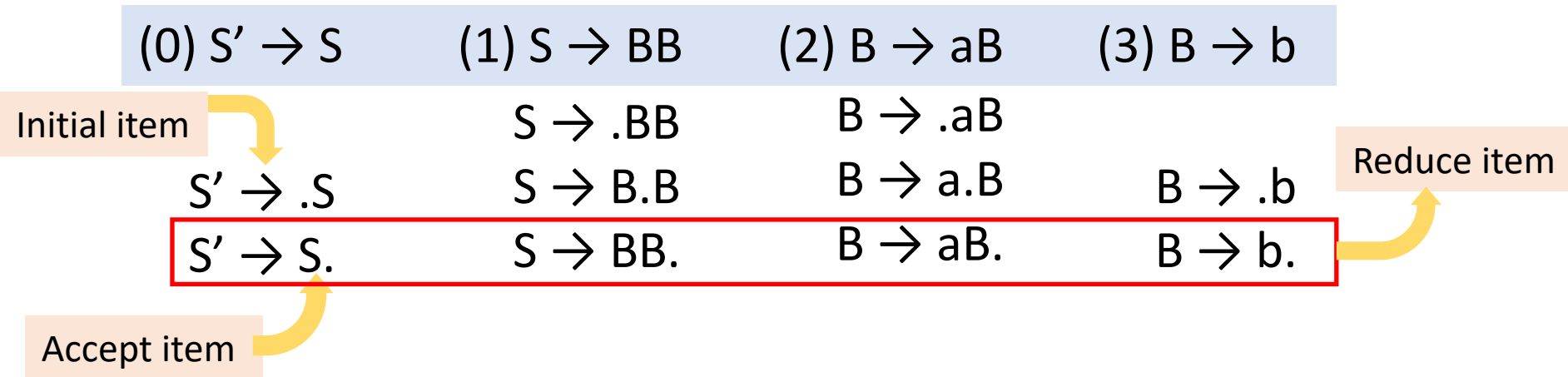
(1) $E \rightarrow E + T$
(2) $E \rightarrow T$
(3) $T \rightarrow T * F$



Augmented grammar:

(0) $E' \rightarrow E$
(1) $E \rightarrow E + T$
(2) $E \rightarrow T$
(3) $T \rightarrow T * F$

Example



- **Closure:** the action of adding equivalent items to a set
 - Example: $S' \rightarrow .S$ $S \rightarrow .BB$ $B \rightarrow .aB$ $B \rightarrow .b$
- Intuitively, $A \rightarrow \alpha.B\beta$ means that we might next see a substring derivable from $B\beta$ ($_sub$) as input. The $_sub$ will have a prefix derivable from B by applying one of the B -productions.
 - Thus, we add items for all the B -productions, i.e., if $B \rightarrow \gamma$ is a production, we add $B \rightarrow .\gamma$ in the closure