

# CSC 416/565: DESIGN AND CONSTRUCTION OF COMPILERS

West Chester University
Dr. Richard Burns
Fall 2023

### **PREVIOUS CLASSES**

- Top-down parsing
  - Recursive descent, requires backtracking
  - Predictive parsing: LL(1), deterministic

### **TODAY**

Bottom-up Parsing

Weakness of LL(k) parsing techniques is that they must *predict* which production to use, after having seen only k tokens.

# **Bottom-up** Parsing Overview

- Builds on ideas of top-down parsing
- More powerful than top-down
- Postpones "reduction" decision until it has seen input tokens corresponding to the entire RHS

# Benefits of Bottom-up Parsing

- (not as restrictive as top-down)
- Left-factoring not necessary
- Elimination of left-recursion unnecessary (because parse tree is built from the leaves up)
  - Right-recursion is not a problem, but may lead to some inefficiencies
- Ambiguous grammars remain problematic
  - (as they should!)

# Bottom-up Parsing

- Most basic bottom-up parsing technique is LR(k) parsing
  - Reading input left-to-right
  - Using a rightmost derivation
  - With k tokens of lookahead
- Idea: bottom-up parsing <u>reduces</u> a string to the start symbol by inverting productions
- *How?* Similar to LL(1) predictive parsing:
  - a parse table is first build based on the grammar
  - At compile time, the driver reads the input tokens, consults the parse table, and creates a parse from the bottom up

```
Grammar: E \rightarrow T + E \mid T T \rightarrow int * T \mid int \mid (E)
```

### Observe what a bottom-up parse looks like:

- How to know which <u>reductions</u> to use, and when to use them?
- Notice: going down-to-up, a <u>rightmost</u> derivation is present
- Eventually reduce, and build complete tree, all the way to the start symbol
- Each step combines subtrees into larger trees

# Shift / Reduce

- Shift / Reduce: the main strategy used by all bottom-up parsers
  - There are a number of LR-parsing variants: (SLR, LALR, ...)
  - All use the same <u>driver</u>
  - Only differ in the generated table

# Parsing Driver for "LR"-family

- Parser has a stack and input.
  - Driver reads the input and consults the parsing table.
- Parsing table contains shift and reduce entries.
- *Initially:* 
  - Stack is empty
  - Parser is at the beginning of <u>input</u>
- Shift: move the next input token onto the top of the stack; bump the input pointer
- Reduce: use a grammar rule "reduction"
- Example: given rule X → ABC, and top of the stack is C B A, pop C B A (the RHS) and push X (the LHS) onto stack
- Note: only ever pushing nonterminals onto the stack in reductions

# Grammar: $E \rightarrow T + E \mid T$ $T \rightarrow int * T \mid int \mid (E)$

### Input String:

```
int * int + int
```

# LR(1) Parse Table Driver Examples

Stack										I	nput							Act
1	a	:=	7	;	b	:=	С	+	(				+	6	,	d	)	\$ shij
1 id4		:=	7	;	b	:=	С	+	(	d	:=	5	+	6	,	d	)	\$ shij
$_{1} id_{4} :=_{6}$			7	;	b	:=	C	+	(	d	:=	5	+	6	,	d	)	\$ shij
$_{1} id_{4} :=_{6} num_{10}$				;	b	:=	C	+	(	d	:=	5	+	6	,	d	)	\$ red
$_{1} id_{4} :=_{6} E_{11}$				;	b	:=	C	+	(	d	:=	5	+	6	,	d	)	\$ red
$_1$ $S_2$				;	b	:=	C	+	(	d	:=	5	+	6	,	d	)	\$ shij
$_{1} S_{2};_{3}$					b	:=	С	+	(	d	:=	5	+	6	,	d	)	\$ shij
$_{1} S_{2};_{3} id_{4}$						:=	С	+	(	d	:=	5	+	6	,	d	)	\$ shij
$_{1} S_{2};_{3} id_{4} :=_{6}$							C	+	(	d	:=	5	+	6	,	d	)	\$ shij
$_{1} S_{2};_{3} id_{4} :=_{6} id_{20}$								+	(	d	:=	5	+	6	,	d	)	\$ red
$_{1} S_{2};_{3} id_{4} :=_{6} E_{11}$								+	(	d	:=	5	+	6	,	d	)	\$ shij
$_{1} S_{2};_{3} id_{4} :=_{6} E_{11}$	+16	j							(	d	:=	5	+	6	,	d	)	\$ shij
$_{1} S_{2};_{3} id_{4} :=_{6} E_{11}$	$+_1$	6 (8								d	:=	5	+	6	,	d	)	\$ shi
$_{1} S_{2};_{3} id_{4} :=_{6} E_{11}$	$+_1$	6 (8	id4	ļ							:=	5	+	6	,	d	)	\$ shij
$_{1} S_{2};_{3} id_{4} :=_{6} E_{11}$	$+_1$	6 (8	id₄	:	=6							5	+	6	,	d	)	\$ shij
$_{1} S_{2};_{3} id_{4} :=_{6} E_{11}$	$+_1$	6 (8	id4	:	=6	num	10						+	6	,	d	)	\$ red
$_{1} S_{2};_{3} id_{4} :=_{6} E_{11}$	$+_1$	6 (8	id4	:	=6	$E_{11}$							+	6	,	d	)	\$ shij
$_{1} S_{2};_{3} id_{4} :=_{6} E_{11}$								16						6	,	d	)	\$ shij
$_{1} S_{2};_{3} id_{4} :=_{6} E_{11}$	$+_1$	6 (8	id4	:	=6	$E_{11}$	+	16	nu	m <sub>10</sub>	0				,	d	)	\$ red
$_{1} S_{2};_{3} id_{4} :=_{6} E_{11}$	$+_1$	6 (8	id4	:	=6	$E_{11}$	+	16	$E_1$	7					,	d	)	\$ red
$_{1} S_{2};_{3} id_{4} :=_{6} E_{11}$	$+_1$	6 (8	id4	:	=6	$E_{11}$									,	d	)	\$ red
$_{1} S_{2};_{3} id_{4} :=_{6} E_{11}$															,	d	)	\$ shij
$_{1} S_{2};_{3} id_{4} :=_{6} E_{11}$	$+_1$	6 (8	$S_1$	2,	18											d	)	\$ shij
$_{1} S_{2};_{3} id_{4} :=_{6} E_{11}$						$id_{20}$											)	\$ red
$_{1} S_{2};_{3} id_{4} :=_{6} E_{11}$																	)	\$ shij
$_{1} S_{2};_{3} id_{4} :=_{6} E_{11}$							)22	2										\$ red
$_{1} S_{2};_{3} id_{4} :=_{6} E_{11}$																		\$ red
$_{1} S_{2};_{3} id_{4} :=_{6} E_{11}$																		\$ red
$_{1} S_{2};_{3} S_{5}$																		\$ red
$_1$ $S_2$																		\$ acc
																		'

Action
shift
shift
shift
$reduce E \rightarrow num$
$reduce S \rightarrow id := E$
shift
shift
shift
shift
reduce $E \rightarrow id$
shift
$reduce E \rightarrow num$
shift
shift
reduce $E \to \text{num}$
reduce $E \to E + E$
$reduce S \rightarrow id := E$
shift
shift
$reduce\ E \rightarrow id$
shift
reduce $E \to (S, E)$
reduce $E \to E + E$
$reduce S \rightarrow id := E$
reduce $S \to S$ ; $S$
accept

id	num	print	;	,	+	:=	(	)	\$	S	$\boldsymbol{E}$	$\boldsymbol{L}$
s4		s7								g2		
			s3						a			
s4		s7								g5		
						s6						
			r1	r1					r1			
s20	s10						s8				g11	
							s9					
s4		s7								g12		
s20	s10						s8				g15	g14
			r5	r5	r5			r5	r5			
			r2	r2	s16				r2			
			s3	s18								
			r3	r3					r3			
				s19				s13				
				r8				r8				
s20	s10						s8				g17	
			r6	r6	s16			r6	r6			
s20	s10						s8				g21	
s20	s10						s8				g23	
			r4	r4	r4			r4	r4			
								s22				
			r7	r7	r7			r7	r7			
				r9	s16			r9				
	s4 s4 s20 s4 s20 s20 s20	s4 s20 s10 s4 s20 s10 s20 s10 s20 s10	s4 s7 s4 s7 s20 s10 s4 s7 s20 s10 s20 s10	s4 s7 s3 s4 s7 r1 s20 s10 s4 s7 s20 s10 r5 r2 s3 r3 r3	s4       s7         s4       s7         s20       s10         s4       s7         s20       s10         r5       r5         r2       r2         s3       s18         r3       r3         s19       r8         s20       s10         s20       s10         r4       r4         r7       r7	s4       s7         s4       s7         s20       s10         s4       s7         s20       s10         r5       r5       r5         r2       r2       s16         s3       s18         r3       r3         s19       r8         s20       s10         s20       s10         s20       s10         r4       r4         r4       r4         r7       r7         r7       r7	s4       s7       s3         s4       s7       s6         r1       r1       r1         s20       s10       r5       r5       r5         r2       r2       r2       s16         s3       s18       r3       r3         s19       r8         s20       s10       r6       r6       s16         s20       s10       r4       r4       r4         r4       r4       r4       r4         r7       r7       r7       r7	s4       s7         s4       s7         s4       s7         s20       s10         s4       s7         s20       s10         s5       r5         r5       r5         r2       r2       s16         s3       s18         r3       r3         s19       r8         s20       s10       s8         s20       r4       r4       r4         r4       r4       r4         r7       r7       r7	s4       s7         s4       s7         s4       s7         s20       s10         s4       s7         s20       s10         s4       s7         s20       s10         r5       r5       r5         r2       r2       r5         s3       s18         r3       r3       r3         s19       s13         r8       r8         s20       s10       s8         r6       r6       s16       r6         s20       s10       s8         s20       s10       s8         r4       r4       r4       r4         r4       r4       r4       r4         r7       r7       r7       r7	s4       s7       s3       a         s4       s7       s6       r1       r1         s20       s10       s8       s9       s4       s7         s20       s10       s8       s9       s8       s9         s4       s7       s5       r5       r3       s13       r3       r3       s13       r8       s8       s8	s4       s7       s3       a       g2         s4       s7       s3       s6       r1       r1         s20       s10       s8       s9       s8       s9       s4       s7       s8       s9       s4       s7       s5       r5       r2       s3       s18       r3       r8       r8       r6       r7       r7	\$4         \$7         \$3         \$6

1 
$$S \rightarrow S$$
;  $S$   
2  $S \rightarrow id := E$   
3  $S \rightarrow print (L)$   
4  $E \rightarrow id$   
5  $E \rightarrow num$   
6  $E \rightarrow E + E$   
7  $E \rightarrow (S, E)$ 

4 
$$E \rightarrow 10$$
  
5  $E \rightarrow \text{num}$   
6  $E \rightarrow E + E$   
9  $L \rightarrow L$ ,  $E$ 

State	а	b	С	d	q	\$	Start	S	Α	В	С	Q
0	3	2	8		8	8	accept	4	1	5		
1			11			4					14	
2		2	8	8	8	8				13		
3		2	8	8						9		
4						8						
5			10		7	10						6
6			6			6						
7			9			9						
8						1						
9			11	4							10	
10				12								
11				3		3						
12			5			5						
13			7	7	7	7						
14						2						

Rule	Derivation
1	$Start \Rightarrow_{rm} S $$
2	$\Rightarrow_{rm} A C $$
3	⇒ <sub>rm</sub> A c \$
5	$\Rightarrow_{rm} a B C d c $ \$
4	$\Rightarrow_{rm} a B d c $$
7	$\Rightarrow_{rm}$ a b B d c \$
7	$\Rightarrow_{rm} a b b B d c$
8	$\Rightarrow_{rm} abbdc$ \$

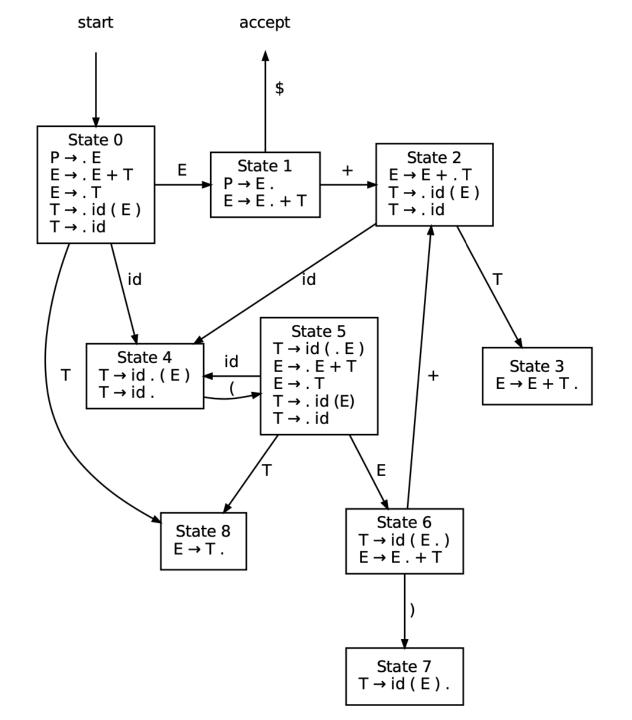
0	_		Initial Configuration	abbdc\$
0	a 3		shift a	bbdc\$
0	a b 2		shift b	bdc\$
	a b 2	b 2	shift b	dc\$
	a b 2	b 2	Reduce $\lambda$ to B	Bdc\$
	a b 2	b B 13	shift B	dc\$
	a b 2		Reduce b B to B	Bdc\$
0	a b 2	B 13	shift B	dc\$
0	a 3		Reduce b B to B	Bdc\$
0	a B 9		shift B	dc\$
0	a B 9		Reduce $\lambda$ to C	Cdc\$
	a B 9	C 10	shift C	dc\$
	a B 3 9	C d 10 12	shift d	c\$
0			Reduce a B C d to A	Ac\$
0	A 1	_	shift A	c\$
0	A c	i	shift c	\$
0	A 1	_	Reduce c to C	C \$
0	A C	1	shift C	\$
0			Reduce A C to S	S\$
0	\$ 4	1	shift S	\$
0	\$ 4		shift\$	\$
0		1	Reduce S \$ to Start	Start \$
0	Start 0		shift Start	\$
			Accept	

State	Go	OTO		ACTION					
	E	T	id	(	)	+	\$		
0	G1	G8	S4						
1						S2	R1		
2		G3	S4						
3					R2	R2	R2		
4				S5	R5	R5	R5		
5	G6	G8	S4						
6					S7	S2			
7					R4	R4	R4		
8					R3	R3	R3		

1. $P \rightarrow E$	
2. $E \rightarrow E +$	Τ
3. $E \rightarrow T$	
4. T $\rightarrow$ id (	E)
5. T $\rightarrow$ id	-

Stack	Symbols	Input	Action
0		id ( id + id) \$	shift 4
0 4	id	(id + id)\$	shift 5
0 4 5	id (	id + id ) \$	shift 4
0454	id ( id	+ id ) \$	reduce $T \rightarrow id$
0458	id ( T	+ id ) \$	reduce $E \rightarrow T$
0456	id (E	+ id ) \$	shift 2
04562	id ( E +	id)\$	shift 4
045624	id ( E + id	)\$	reduce $T \rightarrow id$
045623	id ( E + T	)\$	reduce $E \rightarrow E + T$
0456	id (E	)\$	shift 7
04567	id ( E )	\$	reduce $T \rightarrow id(E)$
0 8	T	\$	$reduce\: E \to T$
0 1	E	\$	accept

# LR(0) Automaton for the Previous Slide



# LR(0) Automaton

- LR(0) automaton represents all possible rules that are currently under consideration by a shift-reduce parser
- Each box represents a *state* in the machine, connected by transitions for both terminals and non-terminals in the grammar
- State contains multiple *items* augmented by a dot (.) to indicate parser's *current position*

### **Example:**

 $E \rightarrow E . + T$ 

Indicates E is currently on the stack and + T is a possible next sequence of tokens

### **Grammar:**

### 1. $P \rightarrow E$ 2. $E \rightarrow E + T$ 3. $E \rightarrow T$ 4. $T \rightarrow id (E)$ 5. $T \rightarrow id$

## Construction of Automaton

- Start production, with a dot (.) at the beginning of the RHS, becomes State 0 (first item is sometimes called the *kernel* of the state)
- Recursively compute the <u>closure</u> of this state and add as additional items in the same state until no more can be added
  - Closure: For each item in the state with a non-terminal X immediately to the right of the dot (.), add all rules in the grammar that have X as the LHS

### Kernel of State 0

$$\mathsf{P} o$$
 .  $\mathsf{E}$ 

### **Closure of State 0**

$$P \rightarrow . E$$
 $E \rightarrow . E + T$ 
 $E \rightarrow . T$ 
 $T \rightarrow . id (E)$ 
 $T \rightarrow . id$ 

# Construction of Automaton (cont.)

### **Closure of State 0**

**Grammar:** 

1. 
$$P \rightarrow E$$
  
2.  $E \rightarrow E + T$   
3.  $E \rightarrow T$   
4.  $T \rightarrow id (E)$   
5.  $T \rightarrow id$ 

- $\begin{array}{c} \mathsf{P} \to . \; \mathsf{E} \\ \mathsf{E} \to . \; \mathsf{E} + \mathsf{T} \\ \mathsf{E} \to . \; \mathsf{T} \\ \mathsf{T} \to . \; \mathsf{id} \; (\; \mathsf{E} \; ) \\ \mathsf{T} \to . \; \mathsf{id} \end{array}$
- All terminals and non-terminals to the right of the dot (.) are possible outgoing transitions
- If the automaton takes that transition it makes a *new state* containing the matching items, with the dot (.) moved one place to the right

Transition on E:

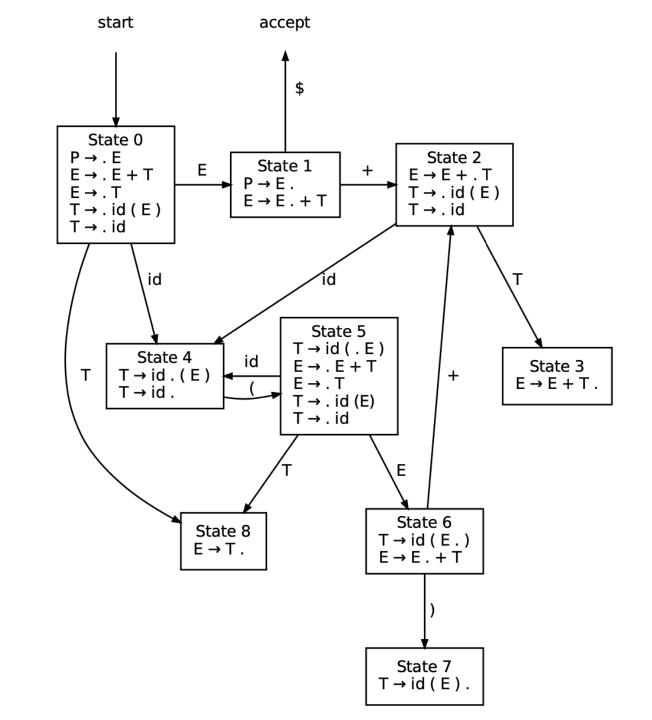
$$P \rightarrow E$$
.  
 $E \rightarrow E$ . + T

Transition on T:

$$\mathsf{E} \to \mathsf{T}$$
 .

Transition on id:

$$T o id . (E)$$
 $T o id .$ 



Shift-Reduce Conflict:

 $\mathsf{T} o \mathsf{id}$  . (  $\mathsf{E}$  )  $\mathsf{T} o \mathsf{id}$  .

# LR(0) Automaton

Reduce-Reduce Conflict:

$$S \rightarrow id (E)$$
.  
 $E \rightarrow id (E)$ .

- LR(0) automata enumerate choices available at any step of the parse
- A state containing an item with a dot (.) at the end indicates a possible reduction
- A transition on a terminal that moves the dot (.) one position to the right indicates a possible shift
- Two types of conflicts can appear in an LR grammar:
  - 1. shift-reduce conflict indicates a choice between a shift action and a reduce action
  - reduce-reduce conflict indicates that two distinct rules have been completely matched

### **SLR Parse Table Creation.**

Given a grammar G and corresponding LR(0) automaton, create tables ACTION[s, a] and GOTO[s, A] for all states s, terminals a, and non-terminals A in G.

For each state s:

For each item like  $A \rightarrow \alpha$ . a  $\beta$ 

ACTION[s, a] = **shift** to state t according to the LR(0) automaton.

For each item like A  $\rightarrow \alpha$  . B  $\beta$ 

GOTO[s, B] = goto state t according to the LR(0) automaton.

For each item like  $A \rightarrow \alpha$ .

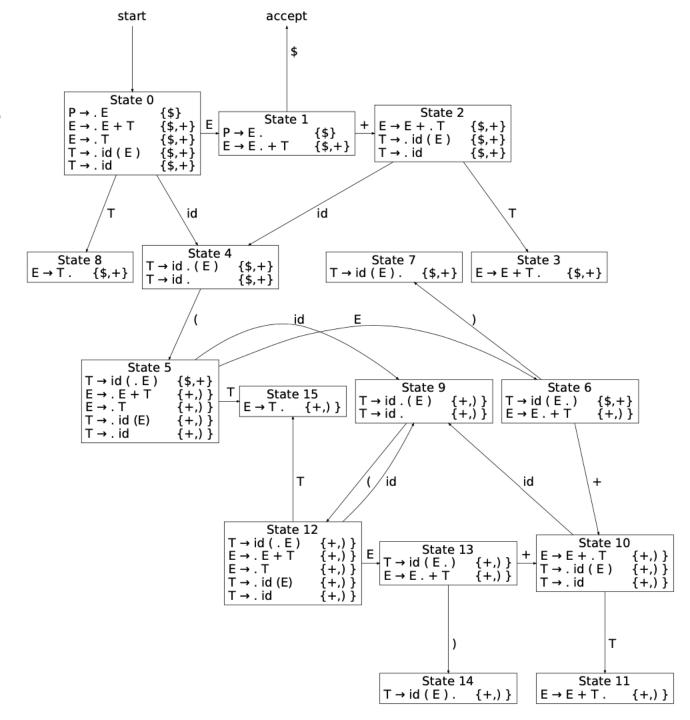
For each terminal a in FOLLOW(A):

ACTION[s, a] = **reduce** by rule A  $\rightarrow \alpha$ 

All remaining states are considered error states.

# Concluding Thoughts

- SLR is a subset of LR(1) and not all LR(1) grammars are SLR
- In practice, a LR(1) automaton is used with an algorithm known as Lookahead LR (LALR)
  - All items in the automaton are augmented with a lookahead of the set of tokens that could potentially follow it



# After the Break

Parsing Tools and an introduction to Programming Assignment #3.

# Parsing Tools

Task of constructing a parser is simple enough to be automated.

# JavaCC – LL(k) Parser Generator

• JavaCC generates a top-down, recursive-descent parser.

Productions written in the form:

```
Assignment \rightarrow Identifier = Expression;
```

Perhaps better written, using tokens for terminals, as:

```
void Assignment() : { } { Identifier() < EQUALS > Expression() < SEMICOLON > }
```

### JavaCC will detect left recursion, which must be eliminated.

### Example:

```
StmList \rightarrow Stm \mid StmList ; Stm
```

### JavaCC syntax:

- 1. curly braces for {Java
   Declarations}
  - Will use these in next part of the course
- 2.{Rule Definitions} used
   to specify the RHS of productions

### Would produce the error:

```
Left Recursion detected:" "StmList ... → StmList ...
```

### JavaCC fragment:

```
void S(): {}
{
    "a" "b" "c"
    |
    "a" "d" "c"
}
```

- Is this grammar LL(1)?
  - It could be after left-factoring.
- Is this grammar LL(2)?
  - Probably yes
  - But LL(1) tables are much smaller and desirable
  - There are no parsing tables in recursive descent anyway!

## Lookahead

### JavaCC Solution:

- LOOKAHEAD directive will allow JavaCC to use more than one token for lookahead
- JavaCC will look at the next two symbols before deciding which rule to use
- Lookahead directives are placed at "choice points"
  - Places in the grammar where there is more than one possible rule that can match

```
void S(): {}
{
    "a" "b" "c"
    |
    "a" "d" "c"
}
```

### **Grammar:**

```
S \rightarrow a(bc|bd)
```

One solution using lookahead of only 2:

```
void S(): {}
{
    "a" (LOOKAHEAD(2)("b""c")|("b""d"))
}
```

NOTE: this wouldn't work

```
void S(): {}
{
    LOOKAHEAD(2) "a" (("b""c")|("b""d"))
}
```

Alternative: could have also factored b out and avoided the lookahead

## Other Notes

- JavaCC can provide a parser for some grammars that are not LL(1)
- the parser that is produced is not guaranteed to correctly parse the language described by the grammar
- a warning will be issued when JavaCC is run on a non-LL(1) grammar
- for a non-LL(1) grammar, the rule that applies first will be used
- LOOKAHEAD suppresses warning messages
  - JavaCC assumes that you know what you are doing

# **Error Recovery**

- What should happen when the parser encounters an error?
- Not only report the *first* error, but ideally recognize *all* of the errors in a program.
- Idea: provide additional grammar rules for error capturing and how far ahead to "skip" to get past error, and to try and resume parsing

### Original Grammar:

```
exp \rightarrow Id
exp \rightarrow exp + exp
exp \rightarrow (exps)
exps \rightarrow exp
exps \rightarrow exp
```

### New Grammar Rule for errors:

```
exp → (error)
exps → error ; exp
```

### *Idea:* skip to the next; or)

- keep discarding input symbols until a lookahead is reached.
  - neither JavaCC nor SableCC supports the <u>error</u> symbol (yacc does)

# Shallow Error Recovery

# Grammar Rule: Stm → IfStm | WhileStm

```
void Stm():
{}
{
    IfStm()
|
    WhileStm()
|
    error_skipto(SEMICOLON)
}
```

- will catch if next token is not <IF> or <WHILE> keyword
- Need to manually define a procedure that will keep skipping tokens until a; and then try to resume

# Deep Error Recovery

```
void Stm() :
    try
            IfStm()
            WhileStm()
      catch (ParserException e) {}
        error skipto(SEMICOLON);
```

- will catch deeper errors
- one advantage of JavaCC (and recursive descent parsers) is that they know which production they are in when the error happened
  - e.g. was in an IfStm()

# Global Error Repair

- finds the smallest set of <u>insertions</u> and <u>deletions</u> that would turn the source string into a syntactically correct string
- one algorithm: "Burke-Fisher error repair"
  - tries every possible single-token insertion, deletion, or replacement, at every point no earlier than K tokens before the point where the parser reported the error
  - Example: if K=10, if parsing engine gets stuck at the 50th token, try *every* possible repair between 40th and 50th token.
  - correction that allows the parser to parse furthest beyond the original reported error is taken
    as best error repair.
- in general, repairs that carry parser 4 tokens beyond the original error are "good enough"
- benefits of technique:
  - 1. grammar is not modified with error productions
  - 2. parsing table is not modified
- Burke-Fisher is not implemented in JavaCC, but that would be an interesting project!