

Lecture VI—Part 2: Syntactic Analysis

Bottom-up Parsing: LR Parsing.

Bottom-Up Parsing

- Bottom-up parsing is more general than top-down parsing
 - And just as efficient
 - Builds on ideas in top-down parsing
 - Preferred method in practice
- Also called LR parsing
 - L means that tokens are read left to right
 - R means that it constructs a rightmost derivation !

An Introductory Example

- LR parsers don't need left-factored grammars and can also handle left-recursive grammars
- Consider the following grammar:

$$E \rightarrow E + (E) \mid \text{int}$$

- Why is this not LL(1)?
- Consider the string: `int + (int) + (int)`

The Idea

- LR parsing *reduces* a string to the start symbol by inverting productions:

str: input string of terminals

repeat

- Identify β in str such that $A \rightarrow \beta$ is a production (i.e., $\text{str} = \alpha \beta \gamma$)
- Replace β by A in str (i.e., str becomes $\alpha A \gamma$)

until str = S

A Bottom-up Parse in Detail (1)

int + (int) + (int)

int + (int) + (int)

A Bottom-up Parse in Detail (2)

int + (int) + (int)

E + (int) + (int)

E
|
int + (int) + (int)

A Bottom-up Parse in Detail (3)

int + (int) + (int)

E + (int) + (int)

E + (E) + (int)

 E E
 | |
int + (int) + (int)

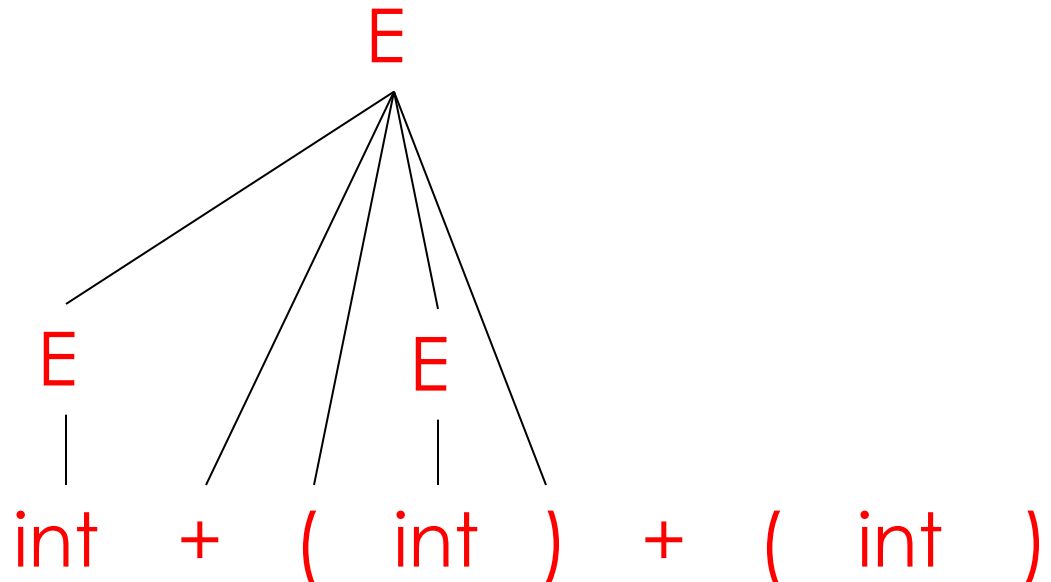
A Bottom-up Parse in Detail (4)

int + (int) + (int)

E + (int) + (int)

E + (E) + (int)

E + (int)



A Bottom-up Parse in Detail (5)

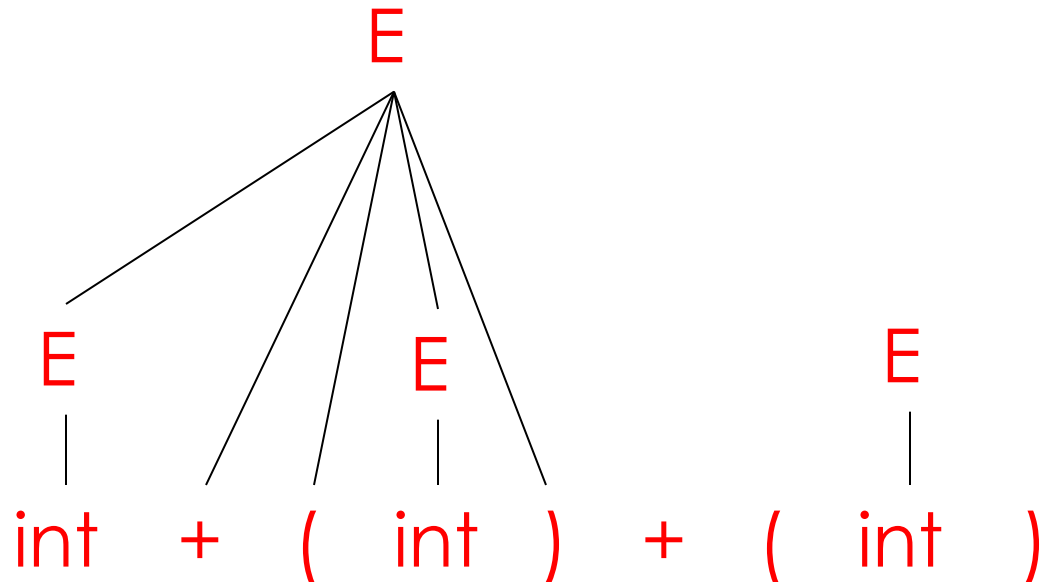
int + (int) + (int)

E + (int) + (int)

E + (E) + (int)

E + (int)

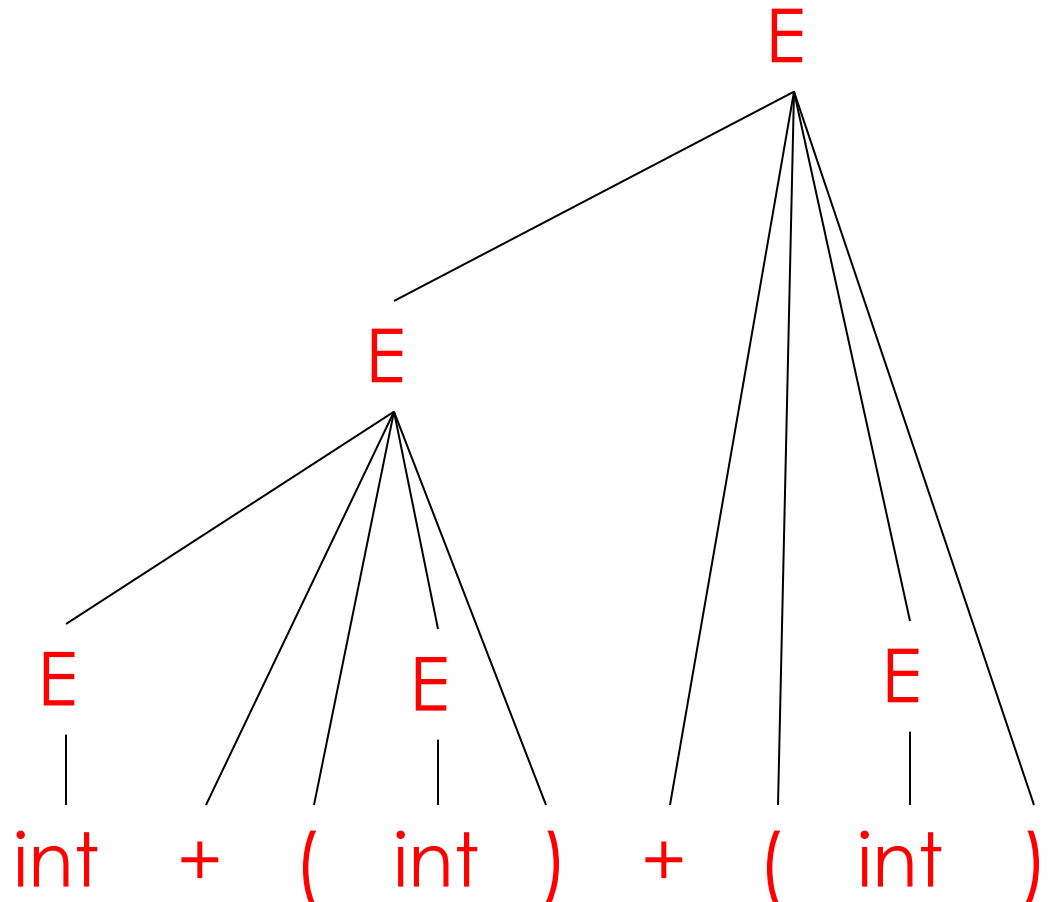
E + (E)



A Bottom-up Parse in Detail (6)

↑
int + (int) + (int)
E + (int) + (int)
E + (E) + (int)
E + (int)
E + (E)
E

A rightmost
derivation in reverse



Important Fact #1

Important Fact #1 about bottom-up parsing:

An LR parser traces a rightmost derivation in reverse

Where Do Reductions Happen

Important Fact #1 has an interesting consequence:

- Let $\alpha\beta\gamma$ be a step of a bottom-up parse
- Assume the next reduction is by $A \rightarrow \beta$
- Then γ is a string of terminals !

Why? Because $\alpha A \gamma \Rightarrow \alpha\beta\gamma$ is a step in a right-most derivation.

Notation

- Idea: Split string into two substrings
 - Right substring (a string of terminals) is as yet unexamined by parser
 - Left substring has terminals and non-terminals
- The dividing point is marked by a |
 - The | is not part of the string
- Initially, all input is unexamined: | $x_1x_2 \dots x_n$

Shift-Reduce Parsing

- Bottom-up parsing uses only two kinds of actions:

Shift

Reduce

Shift

Shift: Move **|** one place to the right

- Shifts a terminal to the left string

$$E + (| \text{int}) \Rightarrow E + (\text{int} |)$$

Reduce

Reduce: Apply an inverse production at the right end of the left string

- If $E \rightarrow E + (E)$ is a production, then

$$E + (\underline{E + (E)} \mid) \Rightarrow E + (\underline{E} \mid)$$

Shift-Reduce Example

| int + (int) + (int)\$ shift

int + (int) + (int)



Shift-Reduce Example

| int + (int) + (int)\$ shift

int | + (int) + (int)\$ red. E \rightarrow int

int + (int) + (int)
↑

Shift-Reduce Example

| int + (int) + (int)\$ shift

int | + (int) + (int)\$ red. E \rightarrow int

E | + (int) + (int)\$ shift 3 times

E
/
int + (int) + (int)
↑

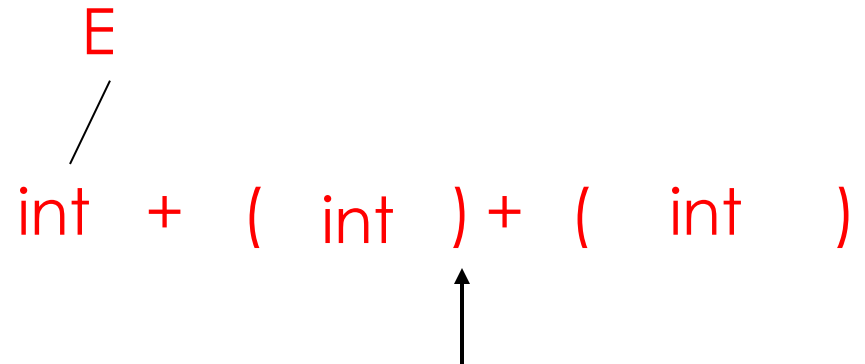
Shift-Reduce Example

| int + (int) + (int)\$ shift

int | + (int) + (int)\$ red. E \rightarrow int

E | + (int) + (int)\$ shift 3 times

E + (int |) + (int)\$ red. E \rightarrow int



Shift-Reduce Example

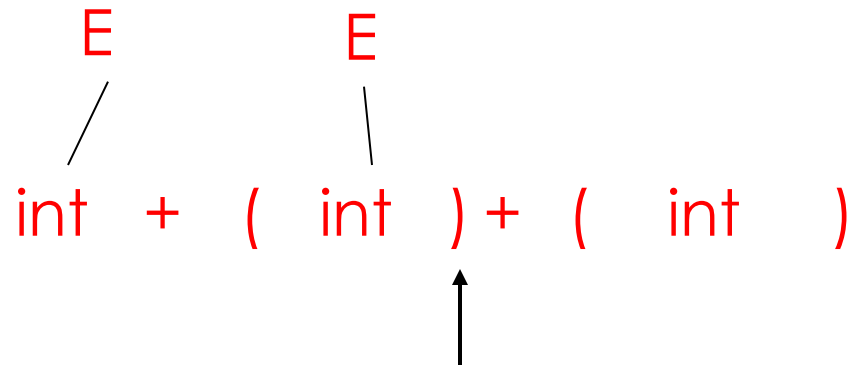
| int + (int) + (int)\$ shift

int | + (int) + (int)\$ red. E → int

E | + (int) + (int)\$ shift 3 times

E + (int |) + (int)\$ red. E → int

E + (E |) + (int)\$ shift



Shift-Reduce Example

| int + (int) + (int)\$ shift

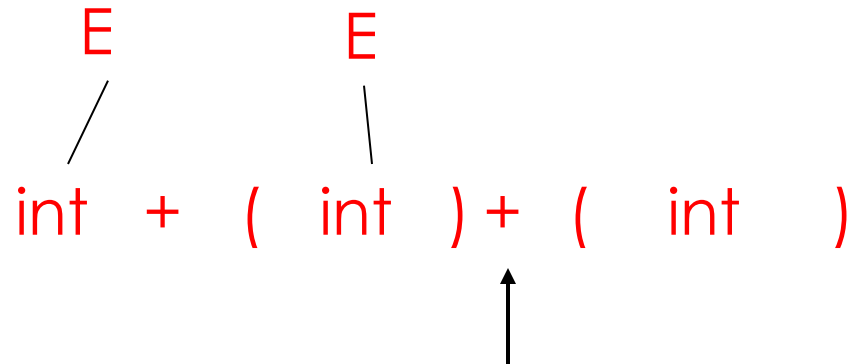
int | + (int) + (int)\$ red. $E \rightarrow \text{int}$

E | + (int) + (int)\$ shift 3 times

E + (int |) + (int)\$ red. $E \rightarrow \text{int}$

E + (E |) + (int)\$ shift

E + (E) | + (int)\$ red. $E \rightarrow E + (E)$



Shift-Reduce Example

| int + (int) + (int)\$ shift

int | + (int) + (int)\$ red. E → int

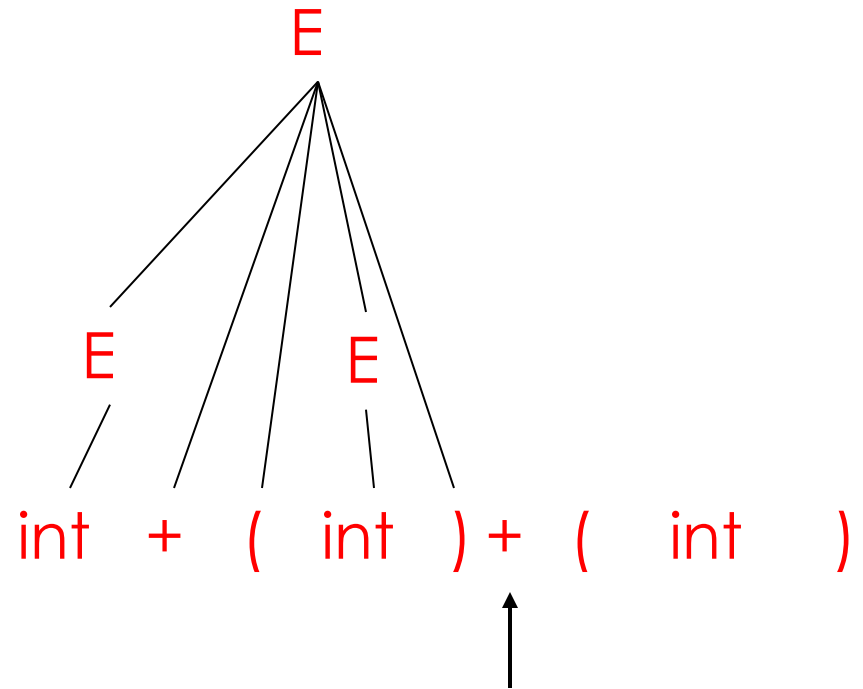
E | + (int) + (int)\$ shift 3 times

E + (int |) + (int)\$ red. E → int

E + (E |) + (int)\$ shift

E + (E) | + (int)\$ red. E → E + (E)

E | + (int)\$ shift 3 times



Shift-Reduce Example

| int + (int) + (int)\$ shift

int | + (int) + (int)\$ red. $E \rightarrow int$

E | + (int) + (int)\$ shift 3 times

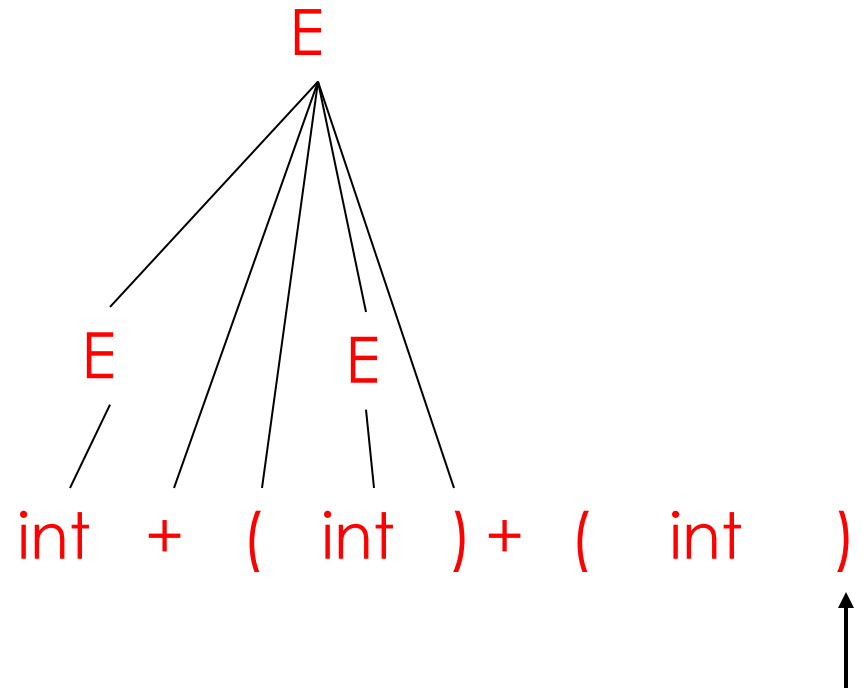
E + (int |) + (int)\$ red. $E \rightarrow int$

E + (E |) + (int)\$ shift

E + (E) | + (int)\$ red. $E \rightarrow E + (E)$

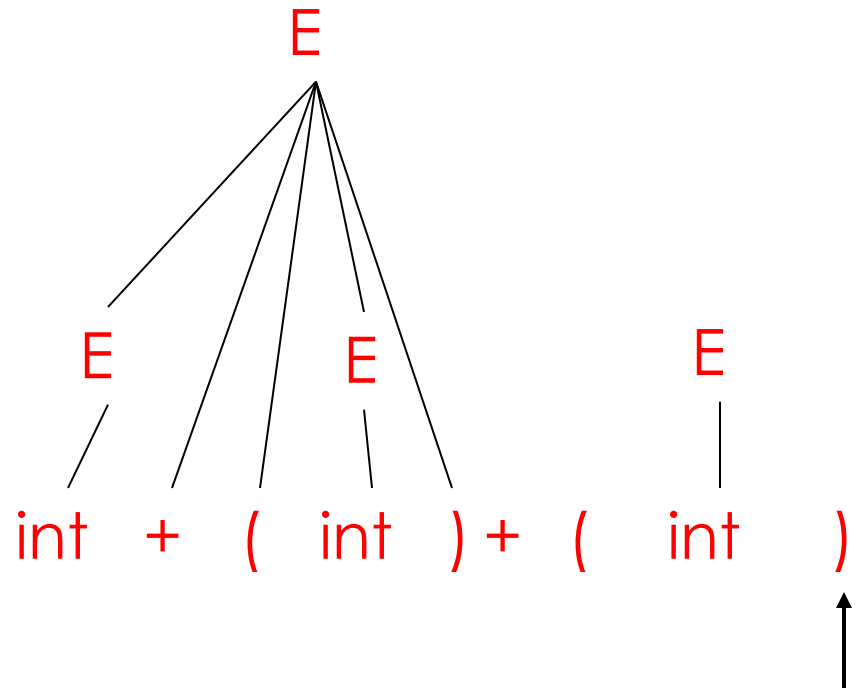
E | + (int)\$ shift 3 times

E + (int |)\$ red. $E \rightarrow int$



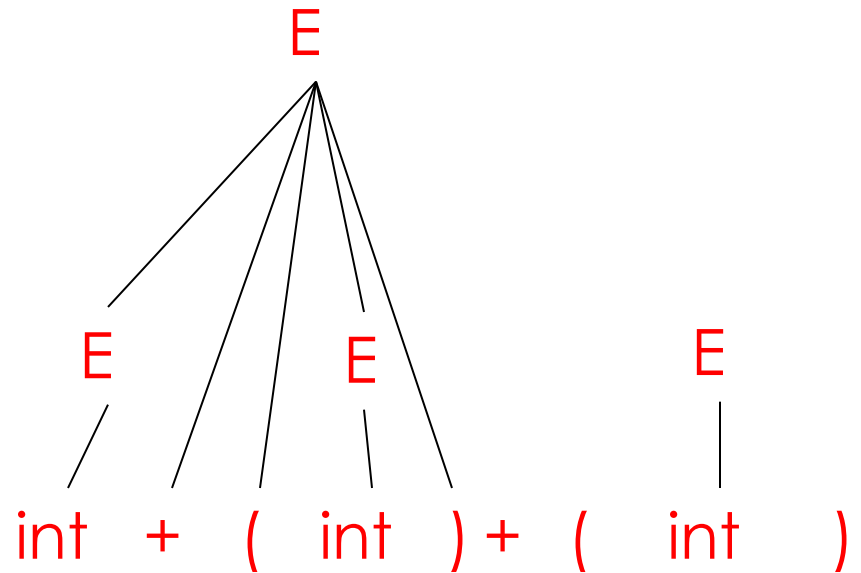
Shift-Reduce Example

$| \text{int} + (\text{int}) + (\text{int})\$$ shift
 $\text{int} | + (\text{int}) + (\text{int})\$$ red. $E \rightarrow \text{int}$
 $E | + (\text{int}) + (\text{int})\$$ shift 3 times
 $E + (\text{int} |) + (\text{int})\$$ red. $E \rightarrow \text{int}$
 $E + (E |) + (\text{int})\$$ shift
 $E + (E) | + (\text{int})\$$ red. $E \rightarrow E + (E)$
 $E | + (\text{int})\$$ shift 3 times
 $E + (\text{int} |)\$$ red. $E \rightarrow \text{int}$
 $E + (E |)\$$ shift



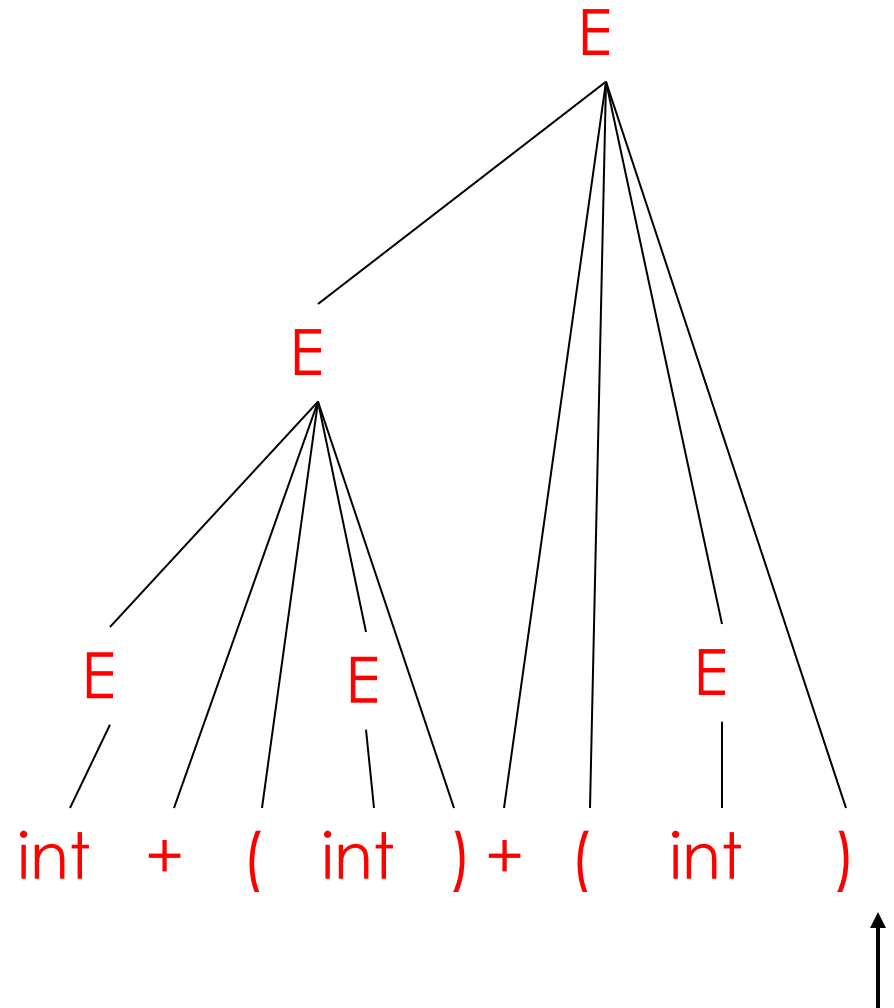
Shift-Reduce Example

$| \text{int} + (\text{int}) + (\text{int})\$$ shift
 $\text{int} | + (\text{int}) + (\text{int})\$$ red. $E \rightarrow \text{int}$
 $E | + (\text{int}) + (\text{int})\$$ shift 3 times
 $E + (\text{int} |) + (\text{int})\$$ red. $E \rightarrow \text{int}$
 $E + (E |) + (\text{int})\$$ shift
 $E + (E) | + (\text{int})\$$ red. $E \rightarrow E + (E)$
 $E | + (\text{int})\$$ shift 3 times
 $E + (\text{int} |) \$$ red. $E \rightarrow \text{int}$
 $E + (E |) \$$ shift
 $E + (E) | \$$ red. $E \rightarrow E + (E)$



Shift-Reduce Example

int + (int) + (int)\$	shift
int + (int) + (int)\$	red. $E \rightarrow \text{int}$
E + (int) + (int)\$	shift 3 times
E + (int) + (int)\$	red. $E \rightarrow \text{int}$
E + (E) + (int)\$	shift
E + (E) + (int)\$	red. $E \rightarrow E + (E)$
E + (int)\$	shift 3 times
E + (int)\$	red. $E \rightarrow \text{int}$
E + (E)\$	shift
E + (E) \$	red. $E \rightarrow E + (E)$
E \$	accept



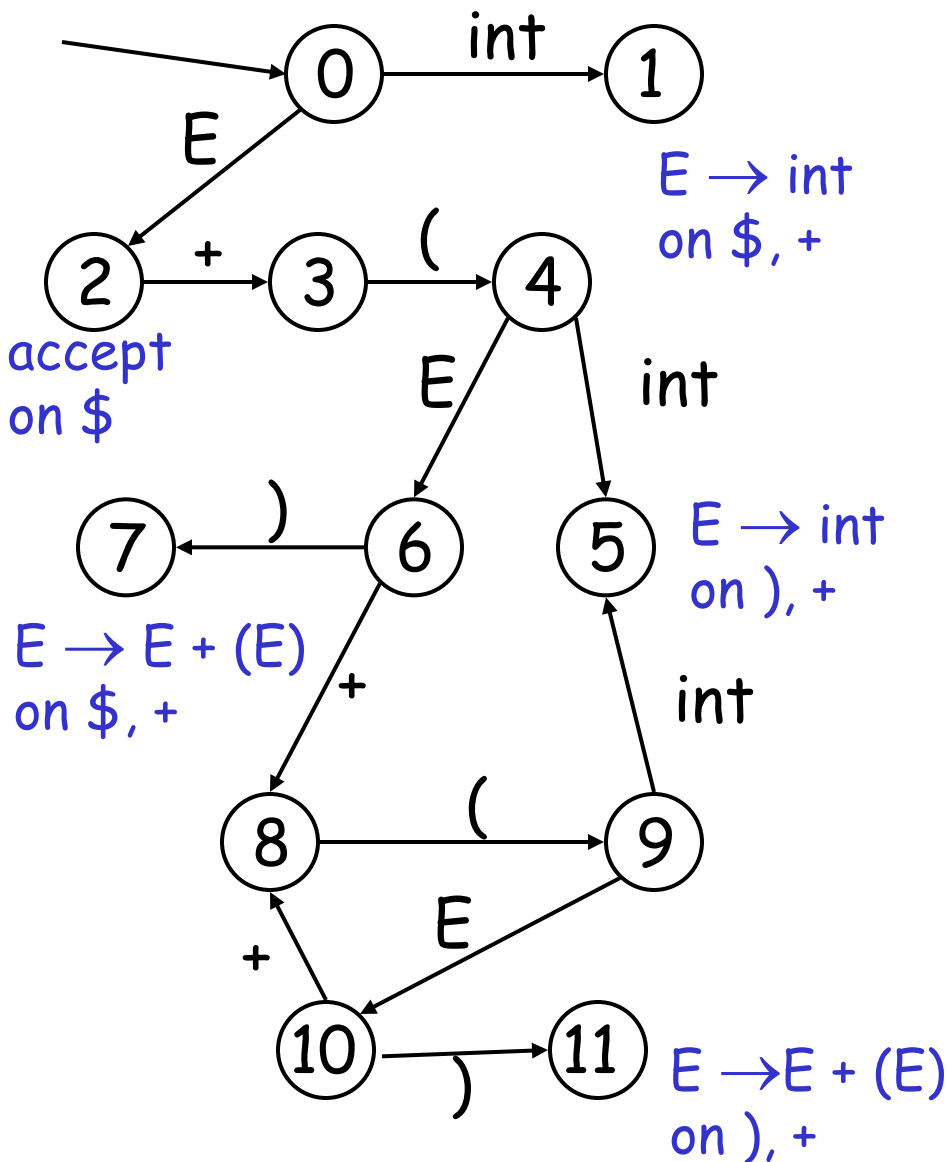
The Stack

- Left string can be implemented by a stack
 - Top of the stack is the |
- Shift pushes a terminal on the stack
- Reduce pops 0 or more symbols off of the stack (production rhs) and pushes a non-terminal on the stack (production lhs)

Key Issue: When to Shift or Reduce?

- Decide based on the stack's content and the lookahead
- Idea: use a finite automaton (DFA) to decide when to shift or reduce
 - The DFA input is the stack's content
 - The alphabet consists of terminals and non-terminals
- We run the DFA on the stack and we examine the resulting state X and the token tok after $|$
 - If X has a transition labeled tok then shift
 - If X is labeled with " $A \rightarrow \beta$ on tok " then reduce

LR(1) Parsing. An Example



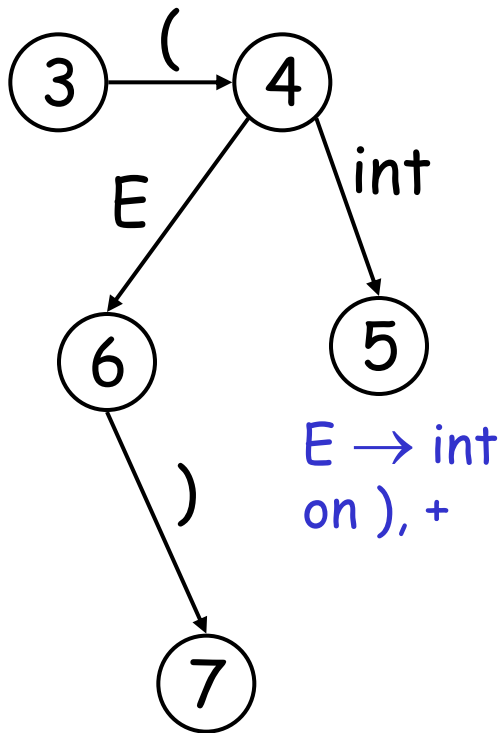
$\text{int} + (\text{int}) + (\text{int})\$$ shift
 $\text{int} | + (\text{int}) + (\text{int})\$$ $E \rightarrow \text{int}$
 $E | + (\text{int}) + (\text{int})\$$ shift(x3)
 $E + (\text{int} |) + (\text{int})\$$ $E \rightarrow \text{int}$
 $E + (E |) + (\text{int})\$$ shift
 $E + (E) | + (\text{int})\$$ $E \rightarrow E + (E)$
 $E | + (\text{int})\$$ shift(x3)
 $E + (\text{int} |)\$$ $E \rightarrow \text{int}$
 $E + (E |)\$$ shift
 $E + (E) | \$$ $E \rightarrow E + (E)$
 $E | \$$ accept

Representing the DFA

- Parsers represent the DFA as a 2D table
- Lines correspond to DFA states
- Columns correspond to terminals and non-terminals
- Typically columns are split into:
 - Those for terminals: **action table**
 - Those for non-terminals: **goto table**

Representing the DFA. Example

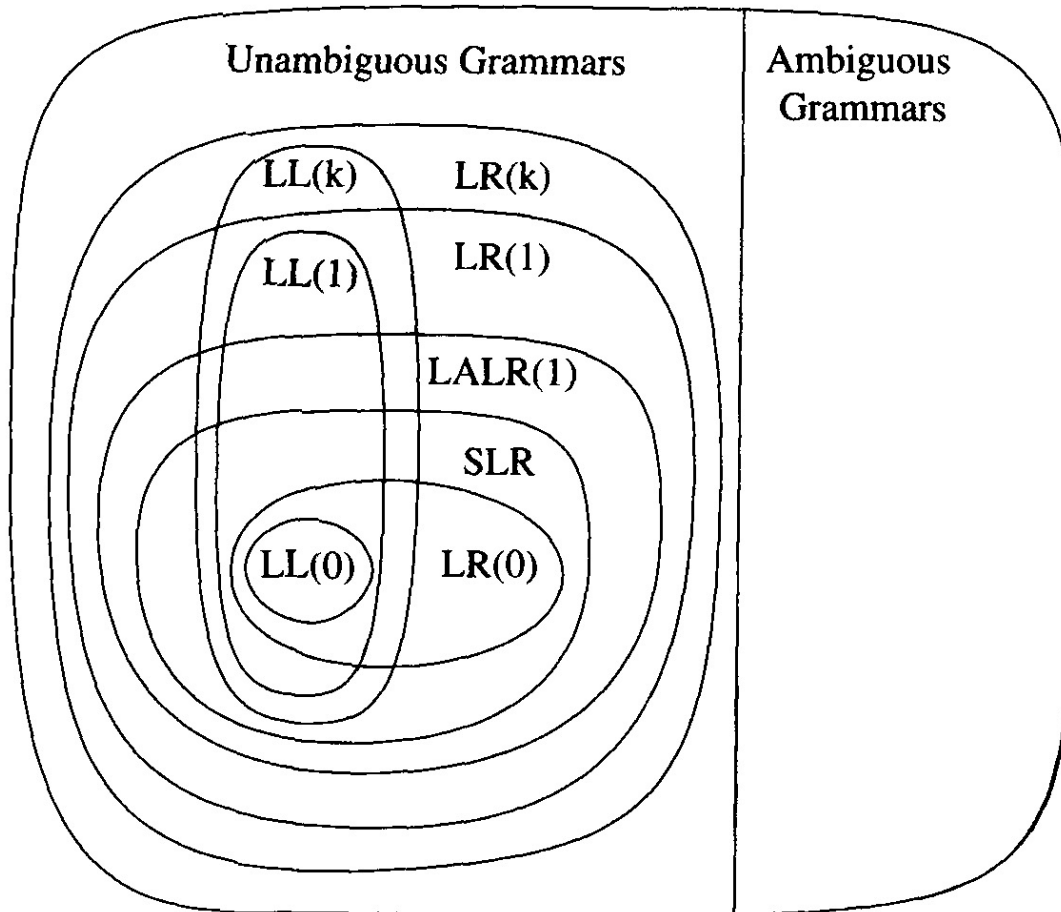
- The table for a fragment of our DFA:



$E \rightarrow E + (E)$
on \$, +

	int	+	()	\$	E
...						
3			s4			
4	s5					g6
5		$r_{E \rightarrow int}$		$r_{E \rightarrow int}$		
6		s8		s7		
7		$r_{E \rightarrow E+(E)}$			$r_{E \rightarrow E+(E)}$	
...						

A Hierarchy of Grammar Classes



From Andrew Appel,
“Modern Compiler
Implementation in Java”