

# Bottom-Up Parsing

CS 132: Compiler Design

January 31, 2011

## REVIEW: LL( $k$ )

A grammar is said to be LL(1) if one can always predict the next grammar production by peeking at one character of input.

A grammar is said to be LL( $k$ ) if we can make the prediction based on the next  $k$  tokens of the input.

A language is said to be LL( $k$ ) if it has an LL( $k$ ) grammar (which generates the right strings, but not necessarily the right parse trees)

# MESSAGING GRAMMARS

$$S \rightarrow E \$$$
$$E \rightarrow n$$
$$| E - n$$
$$\Downarrow$$

Eliminate “left recursion”

$$\Downarrow$$
$$S \rightarrow E \$$$
$$E \rightarrow n$$
$$| n - E$$

- ✓ In this case, a non-LL grammar becomes LL(2).
- ✓ Note: for recursive descent, LL(2) usually isn't that bad!
- ✓ But, this transformation can change the parse tree.

# MESSAGING GRAMMARS

$$S \rightarrow E \$$$

$$E \rightarrow n$$

$$| n - E$$

$$\Downarrow$$

“Left-factor”

$$\Downarrow$$

$$S \rightarrow E \$$$

$$E \rightarrow n E'$$

$$E' \rightarrow \epsilon$$

$$| - E'$$

- ✓ In this case, an LL(2) grammar becomes LL(1).
- ✓ Possibly useful if predictions go via a big lookup table.
- ✓ Parse tree is even less pretty.

## ALTERNATIVE: EBNF

$$S \rightarrow E \$$$

$$E \rightarrow E$$

$$| E - n$$

$$\Downarrow$$

$$S \rightarrow E \$$$

$$E \rightarrow n \{ - n \}$$

or

$$S \rightarrow E \$$$

$$E \rightarrow n ( - n )^*$$

- ✓ Comparatively easy to handle in recursive descent  
(Take one  $n$ , and then loop to gather as many  $- n$ 's as you can.)
- ✓ Result is a "list" of subtractions; can turn these into a proper tree.

# A (CONTRIVED) GRAMMAR

Given the grammar:

$$S \rightarrow a C E \mid b D$$
$$C \rightarrow b c \mid C d$$
$$D \rightarrow d \mid c D e$$
$$E \rightarrow e$$

Draw the parse trees for the inputs  $a b c d e$  and  $b c d e$ .

# LR PARSING

$S \rightarrow a C E \mid b D$

$C \rightarrow b c \mid C d$

$D \rightarrow d \mid c D e$

$E \rightarrow e$

The essence of LR parsing is to build parse trees in a bottom-up fashion:

- ✓ Start with the leaves (  $a b c d e$  and  $b c d e$  respectively)
- ✓ Repeatedly combine subtrees into bigger trees using the productions.

Sadly, we can't naïvely “just run productions backwards”

# SHIFT-REDUCE PARSING

**Stack**

**Input**

**Action** (*Shift/Reduce*)

## SOME TERMINOLOGY

*Assume we have a parse tree.*

*A handle of a parse tree is the leftmost, complete and nontrivial subtree.*

## SOME TERMINOLOGY

*Assume we have a parse tree.*

*A handle of a parse tree is the leftmost, complete and nontrivial subtree.*

*We prune the handle of a parse tree by removing the children/leaves.*

## FROM TREES TO PRODUCTION SEQUENCES

$S \rightarrow a C E \mid b D$

$C \rightarrow b c \mid C d$

$D \rightarrow d \mid c D e$

$E \rightarrow e$

Start with a parse tree (say, for  $a b c d e$ ).

## FROM TREES TO PRODUCTION SEQUENCES

S  $\rightarrow$  a C E | b D

C  $\rightarrow$  b c | C d

D  $\rightarrow$  d | c D e

E  $\rightarrow$  e

Start with a parse tree (say, for a b c d e).

Repeatedly prune the handle.

## FROM TREES TO PRODUCTION SEQUENCES

S  $\rightarrow$  a C E | b D

C  $\rightarrow$  b c | C d

D  $\rightarrow$  d | c D e

E  $\rightarrow$  e

Start with a parse tree (say, for a b c d e).

Repeatedly prune the handle.

How can we get a *rightmost* production sequence out of this?

## FROM TREES TO PRODUCTION SEQUENCES

$S \rightarrow a C E \mid b D$

$C \rightarrow b c \mid C d$

$D \rightarrow d \mid c D e$

$E \rightarrow e$

Start with a parse tree (say, for  $a b c d e$ ).

Repeatedly prune the handle.

How can we get a *rightmost* production sequence out of this?

If only we could find handles in the original string, without precomputing the parse tree...

## VIABLE PREFIXES

Assume that  $\vec{\alpha}$  and  $\vec{\beta}$  are sequences of terminals and non-terminals.

## VIABLE PREFIXES

Assume that  $\vec{\alpha}$  and  $\vec{\beta}$  are sequences of terminals and non-terminals.

We say that  $\vec{\alpha}\vec{\beta}$  is a *viable prefix* if there is a sequence ... of terminals such that

$$\vec{\alpha}\vec{\beta} \dots$$

has a parse tree with  $\vec{\beta}$  as the handle.

## VIABLE PREFIXES

Assume that  $\vec{\alpha}$  and  $\vec{\beta}$  are sequences of terminals and non-terminals.

We say that  $\vec{\alpha}\vec{\beta}$  is a *viable prefix* if there is a sequence ... of terminals such that

$$\vec{\alpha}\vec{\beta} \dots$$

has a parse tree with  $\vec{\beta}$  as the handle.

- ✓ What are some viable prefixes for trees whose root is E? C? D? S?

S  $\rightarrow$  a C E | b D

C  $\rightarrow$  b c | C d

D  $\rightarrow$  d | c D e

E  $\rightarrow$  e

# WHY SHOULD YOU CARE?

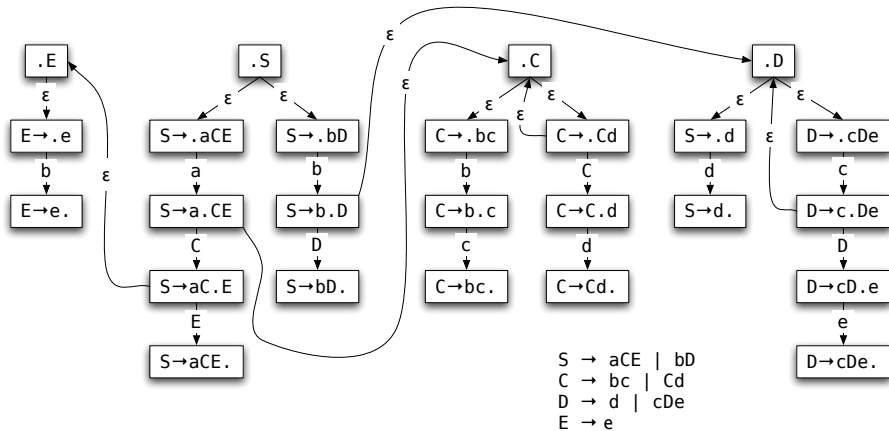
## WHY SHOULD YOU CARE?

- ✓ The viable prefixes of *every* context-free grammar are *regular*!
- ✓ That is, a finite state machine can recognize viable prefixes. It can even tell you what the handle is.

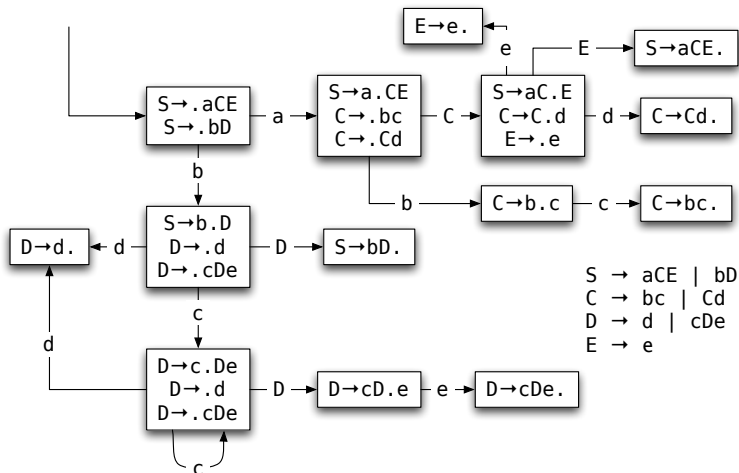
## WHY SHOULD YOU CARE?

- ✓ The viable prefixes of *every* context-free grammar are *regular*!
- ✓ That is, a finite state machine can recognize viable prefixes. It can even tell you what the handle is.
- ✓ Every parseable string (by definition) has a viable prefix.
- ✓ Some grammars have a property that no viable prefix extends another. These are called **LR(0)** grammars.
- ✓ If so, every derivable string starts with a *unique* viable prefix. This gives us a parsing algorithm for **LR(0)** grammars.

## PARSING AUTOMATON: NFA



## BUILDING A PARSING AUTOMATON



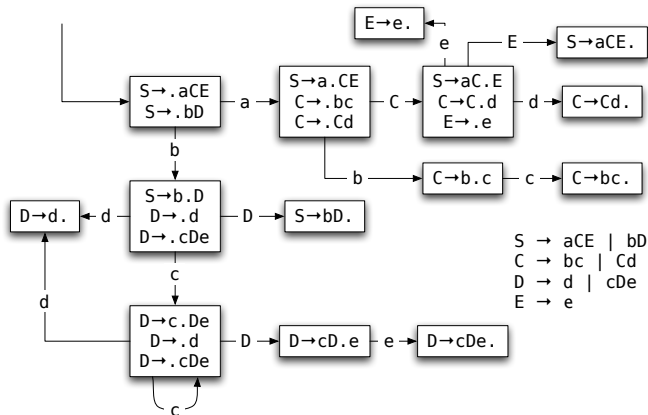
## USING THE PARSING AUTOMATON

Parse *b c d e* and *a b c d e* using the parsing automaton.

<b>Stack</b>	<b>Input</b>	<b>Action</b> (Shift/Reduce)
--------------	--------------	------------------------------

# LR(0)

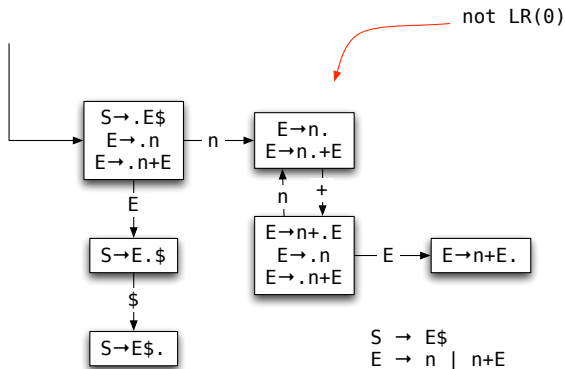
A grammar is LR(0) if the parsing automaton (as described) gives unambiguous instructions.



BTW, why doesn't anyone ever consider LL(0) grammars?

## A NON-LR(0) EXAMPLE

Sadly, few useful grammars are LR(0). For example,



- ✓ Show a parse tree for an input starting with  $n$ , where  $n$  is the handle.
- ✓ Show a parse tree for an input starting with  $n$ , where  $n$  isn't the handle.

Thus, in practice we need to peek past the handle.

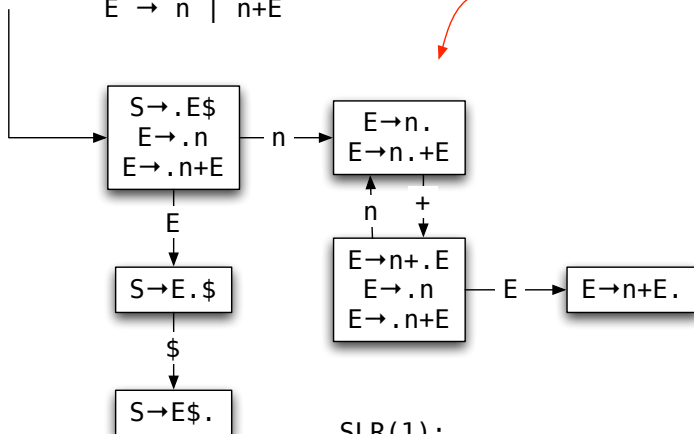
## SLR(1)

Reduce using a rule  $B \rightarrow \beta$  only if the next input token is in  $\text{FOLLOW}(B)$ .

$$S \rightarrow E\$$$

$$E \rightarrow n \mid n+E$$

not LR(0)



SLR(1):

$$\text{FOLLOW}(E) = \{\$, \}$$



## DANGLING ELSE

```
S -> if E then S else S
    | if E then S
    | ...
```

How to parse

```
if a then if b then s1 else s2
```

## DANGLING ELSE

```
S -> if E then S else S
    | if E then S
    | ...
```

How to parse

```
if a then if b then s1 else s2
```

No ambiguous grammar is LR(k), so you'd get a shift-reduce conflict

<pre>S -&gt; if E then S . else S</pre>	<pre>&lt;any&gt;</pre>
<pre>S -&gt; if E then S .</pre>	<pre>else</pre>

# SOLUTIONS

Hack the *grammar*

$S \rightarrow M$

$S \rightarrow U$

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

$M \rightarrow \dots$

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

$U \rightarrow \text{if } E \text{ then } M \text{ else } U$

# SOLUTIONS

Hack the *grammar*

S -> M

S -> U

M -> if E then M else M

M -> ...

U -> if E then S

U -> if E then M else U

Hack the *language*

S -> if E then S else S fi

| if E then S fi

| ...

# SOLUTIONS

Hack the *grammar*

S → M

S → U

M → if E then M else M

M → ...

U → if E then S

U → if E then M else U

Hack the *language*

S → if E then S else S fi

| if E then S fi

| ...

Hack the *parser*:

Whenever you reach this “conflicting” state, just go ahead and shift.

# AMBIGUOUS EXPRESSION GRAMMARS

$$\begin{aligned} E &\rightarrow E + E \mid E - E \\ &\mid E * E \mid E / E \\ &\mid ( E ) \mid n \end{aligned}$$

# AMBIGUOUS EXPRESSION GRAMMARS

$$\begin{aligned}
 E &\rightarrow E + E \mid E - E \\
 &\mid E * E \mid E / E \\
 &\mid ( E ) \mid n
 \end{aligned}$$

All sorts of shift-reduce conflicts:

$E \rightarrow E . + E$	<any>
$E \rightarrow E * E .$	+

$E \rightarrow E . * E$	<any>
$E \rightarrow E + E .$	*

$E \rightarrow E . + E$	<any>
$E \rightarrow E + E .$	+

# SOLUTIONS

Hack the *grammar*

$E \rightarrow E + T \mid E - T \mid T$	$E \rightarrow T (+ T \mid - T)^*$
$T \rightarrow T * F \mid T / F \mid F$	$T \rightarrow F (* F \mid / F)^*$
$F \rightarrow ( E ) \mid n$	$F \rightarrow ( E ) \mid n$

# SOLUTIONS

Hack the *grammar*

$$\begin{array}{ll} E \rightarrow E + T \mid E - T \mid T & E \rightarrow T (+ T \mid - T)^* \\ T \rightarrow T * F \mid T / F \mid F & T \rightarrow F (* F \mid / F)^* \\ F \rightarrow ( E ) \mid n & F \rightarrow ( E ) \mid n \end{array}$$

Hack the *language*

$$E \rightarrow + E E \mid * E E \mid \dots$$

# SOLUTIONS

Hack the *grammar*

$$\begin{array}{ll} E \rightarrow E + T \mid E - T \mid T & E \rightarrow T (+ T \mid - T)^* \\ T \rightarrow T * F \mid T / F \mid F & T \rightarrow F (* F \mid / F)^* \\ F \rightarrow ( E ) \mid n & F \rightarrow ( E ) \mid n \end{array}$$

Hack the *language*

$$E \rightarrow + E E \mid * E E \mid \dots$$

Hack the *parser* to shift or reduce as desired.

# SOLUTIONS

Hack the *grammar*

$$\begin{array}{ll} E \rightarrow E + T \mid E - T \mid T & E \rightarrow T ( + T \mid - T )^* \\ T \rightarrow T * F \mid T / F \mid F & T \rightarrow F ( * F \mid / F )^* \\ F \rightarrow ( E ) \mid n & F \rightarrow ( E ) \mid n \end{array}$$

Hack the *language*

$$E \rightarrow + E E \mid * E E \mid \dots$$

Hack the *parser* to shift or reduce as desired. YACC support:

```
%left PLUS MINUS
```

```
%left STAR SLASH
```

# SOLUTIONS

Hack the *grammar*

$$\begin{array}{ll} E \rightarrow E + T \mid E - T \mid T & E \rightarrow T ( + T \mid - T )^* \\ T \rightarrow T * F \mid T / F \mid F & T \rightarrow F ( * F \mid / F )^* \\ F \rightarrow ( E ) \mid n & F \rightarrow ( E ) \mid n \end{array}$$

Hack the *language*

$$E \rightarrow + E E \mid * E E \mid \dots$$

Hack the *parser* to shift or reduce as desired. YACC support:

```
%left PLUS MINUS
```

```
%left STAR SLASH
```

- ✓ Terminals listed in increasing precedence

# SOLUTIONS

Hack the *grammar*

$$\begin{array}{ll}
 E \rightarrow E + T \mid E - T \mid T & E \rightarrow T (+ T \mid - T)^* \\
 T \rightarrow T * F \mid T / F \mid F & T \rightarrow F (* F \mid / F)^* \\
 F \rightarrow ( E ) \mid n & F \rightarrow ( E ) \mid n
 \end{array}$$

Hack the *language*

$$E \rightarrow + E E \mid * E E \mid \dots$$

Hack the *parser* to shift or reduce as desired. YACC support:

```
%left PLUS MINUS
```

```
%left STAR SLASH
```

- ✓ Terminals listed in increasing precedence
- ✓ Precedence of a rule is that of its last terminal

# SOLUTIONS

Hack the *grammar*

$$\begin{array}{ll} E \rightarrow E + T \mid E - T \mid T & E \rightarrow T (+ T \mid - T)^* \\ T \rightarrow T * F \mid T / F \mid F & T \rightarrow F (* F \mid / F)^* \\ F \rightarrow ( E ) \mid n & F \rightarrow ( E ) \mid n \end{array}$$

Hack the *language*

$$E \rightarrow + E E \mid * E E \mid \dots$$

Hack the *parser* to shift or reduce as desired. YACC support:

```
%left PLUS MINUS
```

```
%left STAR SLASH
```

- ✓ Terminals listed in increasing precedence
- ✓ Precedence of a rule is that of its last terminal
- ✓ Reduce if rule precedence greater than next token, or if same and rule precedence is left-associative. Otherwise shift.

## UNARY MINUS

We want these rules to be in increasing precedence:

```
exp: exp MINUS exp
    | exp TIMES exp
    | MINUS exp
```

But YACC would give the same precedence to the first and third rules.

## SOLUTION: IMAGINARY TOKENS + OVERRIDE RULE PRECEDENCE

```
%left PLUS MINUS
%left STAR SLASH
%left UNARY_MINUS
```

```
exp: exp MINUS exp
    | exp TIMES exp
    | MINUS exp      %prec UNARY_MINUS
```

(The lexer never produces a UNARY\_MINUS token!)