



COMP 412  
FALL 2010

# *Parsing Wrap Up*

## *Comp 412*

Copyright 2010, Keith D. Cooper & Linda Torczon, all rights reserved.

Students enrolled in Comp 412 at Rice University have explicit permission to make copies of these materials for their personal use.

Faculty from other educational institutions may use these materials for nonprofit educational purposes, provided this copyright notice is preserved.



# LR( $k$ ) versus LL( $k$ )

Finding the next step in a derivation

LR( $k$ )  $\Rightarrow$  Each reduction in the parse is detectable with

- $\rightarrow$  the complete left context,
- $\rightarrow$  the reducible phrase, itself, and
- $\rightarrow$  the  $k$  terminal symbols to its right

generalizations of  
LR(1) and LL(1) to  
longer lookaheads

LL( $k$ )  $\Rightarrow$  Parser must select the expansion based on

- $\rightarrow$  The complete left context
- $\rightarrow$  The next  $k$  terminals

Thus, LR( $k$ ) examines more context

*The question is, do languages fall in the gap between LR( $k$ ) and LL( $k$ )?*



## LR(1) versus LL(1)

The following LR(1) grammar has no LL(1) counterpart

- The Canonical Collection has 18 sets of LR(1) Items
  - It is not a simple grammar
  - It is, however, LR(1)

0	<i>Goal</i>	→	<i>S</i>
1	<i>S</i>	→	<i>A</i>
2			<i>B</i>
3	<i>A</i>	→	( <i>A</i> )
4			<u><i>a</i></u>
5	<i>B</i>	→	( <i>B</i> >
6			<u><i>b</i></u>

- It requires an arbitrary lookahead to choose between *A* & *B*
- An LR(1) parser can carry the left context (the '(' s) until it sees *a* or *b*
- The table construction will handle it
- In contrast, an LL(1) parser cannot decide whether to expand *Goal* by *A* or *B*
  - No amount of massaging the grammar will resolve this problem

More precisely, the language described by this LR(1) grammar cannot be described with an LL(1) grammar. In fact, the language has no LL(*k*) grammar, for finite *k*.



# Building LR(1) Tables for Waite's Language

## The Canonical Collection of Sets of LR(1) Items

cc <sub>0</sub>	$[Goal \rightarrow \bullet S, \underline{EOF}], [S \rightarrow \bullet A, \underline{EOF}], [S \rightarrow \bullet B, \underline{EOF}], [A \rightarrow \bullet ( A ), \underline{EOF}],$ $[A \rightarrow \bullet \underline{a}, \underline{EOF}], [B \rightarrow \bullet ( B \succeq, \underline{EOF}], [B \rightarrow \bullet \underline{b}, \underline{EOF}],$
cc <sub>1</sub>	$[Goal \rightarrow S \bullet, \underline{EOF}]$
cc <sub>2</sub>	$[S \rightarrow A \bullet, \underline{EOF}]$
cc <sub>3</sub>	$[S \rightarrow B \bullet, \underline{EOF}]$
cc <sub>4</sub>	$[A \rightarrow ( \bullet A ), \underline{ ) }], [A \rightarrow ( \bullet A ), \underline{ EOF }], [A \rightarrow \bullet \underline{a}, \underline{ ) }], [B \rightarrow \bullet ( B \succeq, \underline{ \succeq }],$ $[B \rightarrow ( \bullet B \succeq, \underline{ EOF }], [B \rightarrow \bullet \underline{b}, \underline{ \succeq }]$
cc <sub>5</sub>	$[A \rightarrow \underline{a} \bullet, \underline{ EOF }]$
cc <sub>6</sub>	$[B \rightarrow \underline{b} \bullet, \underline{ EOF }]$
cc <sub>7</sub>	$[A \rightarrow ( A \bullet ), \underline{ EOF }]$
cc <sub>8</sub>	$[B \rightarrow ( B \bullet \succeq, \underline{ EOF }]$

cc<sub>4</sub> is goto(cc<sub>0</sub>, ( )

0	Goal	→	S
1	S	→	A
2			B
3	A	→	( A )
4			<u>a</u>
5	B	→	( B <u>≥</u>
6			<u>b</u>



# Building LR(1) Tables for Waite's Language

And, the rest of it ...

cc <sub>9</sub>	$[A \rightarrow ( \bullet A ), )], [A \rightarrow ( \bullet A ), ), [A \rightarrow \bullet \underline{a}, ), [B \rightarrow \bullet ( B \succeq, \succeq],$ $[B \rightarrow ( \bullet B \succeq, \succeq], [B \rightarrow \bullet \underline{b}, \succeq]$
cc <sub>10</sub>	$[A \rightarrow \underline{a} \bullet, )]$
cc <sub>11</sub>	$[B \rightarrow \underline{b} \bullet, \succeq]$
cc <sub>12</sub>	$[A \rightarrow ( A ) \bullet, \underline{\text{EOF}}]$
cc <sub>13</sub>	$[B \rightarrow ( B \succeq \bullet, \underline{\text{EOF}}]$
cc <sub>14</sub>	$[A \rightarrow ( A \bullet ), )]$
cc <sub>15</sub>	$[B \rightarrow ( B \bullet \succeq, \succeq]$
cc <sub>16</sub>	$[A \rightarrow ( A ) \bullet, )]$
cc <sub>17</sub>	$[B \rightarrow ( B \succeq \bullet, \succeq]$

0	Goal	→	S
1	S	→	A
2			B
3	A	→	( <u>A</u> )
4			<u>a</u>
5	B	→	( B <u>≥</u>
6			<u>b</u>

# Building LR(1) Tables for Waite's Language



	EOF	(	)	<u>a</u>	<u>≥</u>	<u>b</u>	S	A	B
s <sub>0</sub>		s 4		s 5		s 6	1	2	3
s <sub>1</sub>	acc								
s <sub>2</sub>	r 2								
s <sub>3</sub>	r 3								
s <sub>4</sub>		s 9		s 10		s 11		7	8
s <sub>5</sub>	r 5								
s <sub>6</sub>	r 7								
s <sub>7</sub>			s 12						
s <sub>8</sub>					s 13				
s <sub>9</sub>				s 10		s 11		14	15
s <sub>10</sub>			r 5						
s <sub>11</sub>					r 7				
s <sub>12</sub>	r 4								
s <sub>13</sub>	r 6								
s <sub>14</sub>			s 16						
s <sub>15</sub>					s 17				
s <sub>16</sub>			r 4						
s <sub>17</sub>					r 6				

- Notice how sparse the table is.  
- Goto has 7 of 54  
- Action has 23 of 108
- Notice rows & columns that might be combined.
- Notice  $|CC| > |P|$

0	Goal	→	S
1	S	→	A
2			B
3	A	→	( <u>A</u> )
4			<u>a</u>
5	B	→	( B <u>≥</u>
6			<u>b</u>



# LR(k) versus LL(k)

## Other Non-LL Grammars

0	$B \rightarrow R$
1	$\mid (B)$
2	$R \rightarrow E = E$
3	$E \rightarrow \underline{a}$
4	$\mid \underline{b}$
5	$\mid (E + E)$

Example from D.E Knuth, "Top-Down Syntactic Analysis," *Acta Informatica*, 1:2 (1971), pages 79-110

0	$S \rightarrow \underline{a} A \underline{b}$
1	$\mid \underline{c}$
2	$A \rightarrow \underline{b} S$
3	$\mid B \underline{b}$
4	$B \mid \underline{a} A$
5	$\mid \underline{c}$

Example from Lewis, Rosenkrantz, & Stearns book, "Compiler Design Theory," (1976), Figure 13.1

This grammar is actually LR(0)



# LR( $k$ ) versus LL( $k$ )

Finding the next step in a derivation

LR( $k$ )  $\Rightarrow$  Each reduction in the parse is detectable with

- $\rightarrow$  the complete left context,
- $\rightarrow$  the reducible phrase, itself, and
- $\rightarrow$  the  $k$  terminal symbols to its right

generalizations of  
LR(1) and LL(1) to  
longer lookaheads

LL( $k$ )  $\Rightarrow$  Parser must select the expansion based on

- $\rightarrow$  The complete left context
- $\rightarrow$  The next  $k$  terminals

Thus, LR( $k$ ) examines more context

*"... in practice, programming languages do not actually seem to fall in the gap between LL(1) languages and deterministic languages"*

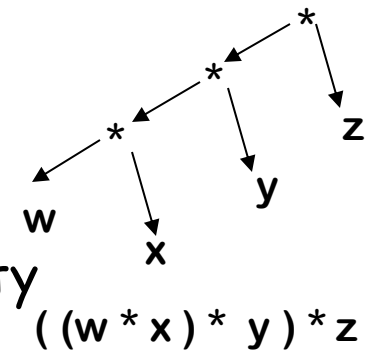
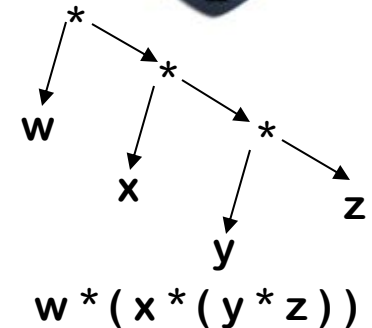
*J.J. Horning, "LR Grammars and Analysers", in Compiler Construction, An Advanced Course, Springer-Verlag, 1976*





# Left Recursion versus Right Recursion

- Right recursion
  - Required for termination in top-down parsers
  - Uses (on average) more stack space
  - Naïve right recursion produces right-associativity
- Left recursion
  - Works fine in bottom-up parsers
  - Limits required stack space
  - Naïve left recursion produces left-associativity
- Rule of thumb
  - Left recursion for bottom-up parsers
  - Right recursion for top-down parsers

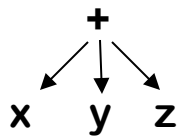




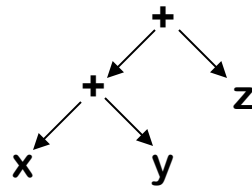
# Associativity

What difference does it make?

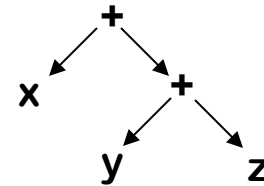
- Can change answers in floating-point arithmetic
- Can change opportunities for optimization
- Consider  $x+y+z$



*Ideal  
operator*



*Left  
association*



*Right  
association*

What if  $y+z$  occurs elsewhere? Or  $x+y$ ? or  $x+z$ ?

The compiler may want to change the "shape" of expressions

- What if  $x = 2$  &  $z = 17$  ? Neither left nor right exposes 19.
- Best shape is function of surrounding context



# Error Detection and Recovery

---

## Error Detection

- Recursive descent
  - Parser takes the last else clause in a routine
  - Compiler writer can code almost any arbitrary action
- Table-driven LL(1)
  - In state  $s_i$  facing word  $x$ , entry is an error
  - Report the error, valid entries in row for  $s_i$  encode possibilities
- Table-driven LR(1)
  - In state  $s_i$  facing word  $x$ , entry is an error
  - Report the error, shift states in row encode possibilities
  - Can precompute better messages from LR(1) items



# Error Detection and Recovery

---

## Error Recovery

- Table-driven LL(1)
  - Treat as missing token, e.g. '}',  $\Rightarrow$  expand by desired symbol
  - Treat as extra token, e.g., 'x - + y',  $\Rightarrow$  pop stack and move ahead
- Table-driven LR(1)
  - Treat as missing token, e.g. '}',  $\Rightarrow$  shift the token
  - Treat as extra token, e.g., 'x - + y',  $\Rightarrow$  don't shift the token

Can pre-compute sets of states  
with appropriate entries...



# Error Detection and Recovery

One common strategy is “hard token” recovery

Skip ahead in input until we find some “hard” token, e.g. ‘;’

- ‘;’ separates statements, makes a logical break in the parse
- Resynchronize state, stack, and input to point after hard token
  - LL(1): pop stack until we find a row with entry for ‘;’
  - LR(1): pop stack until we find a state with a reduction on ‘;’
- Does not correct the input, rather it allows parse to proceed

```
NT ← pop()
repeat until Table[NT,‘;’] ≠ error
  NT ← pop()
token ← NextToken()
repeat until token = ‘;’
  token ← NextToken()
```

*Resynchronizing an LL(1) parser*

Comp 412, Fall 2010

```
repeat until token = ‘;’
  shift token
  shift  $s_e$ 
  token ← NextToken()
reduce by error production
// pops all that state off stack
```

*Resynchronizing an LR(1) parser*

12



# Shrinking the ACTION and GOTO Tables

Three options:

- Combine terminals such as number & identifier, + & -, \* & /
  - Directly removes a column, may remove a row
  - For expression grammar, 198 (vs. 384) table entries
- Combine rows or columns
  - Implement identical rows once & remap states
  - Requires extra indirection on each lookup
  - Use separate mapping for ACTION & for GOTO
- Use another construction algorithm
  - Both LALR(1) and SLR(1) produce smaller tables
    - LALR(1) represents each state with its "core" items
    - SLR(1) uses LR(0) items and the FOLLOW set
  - Implementations are readily available

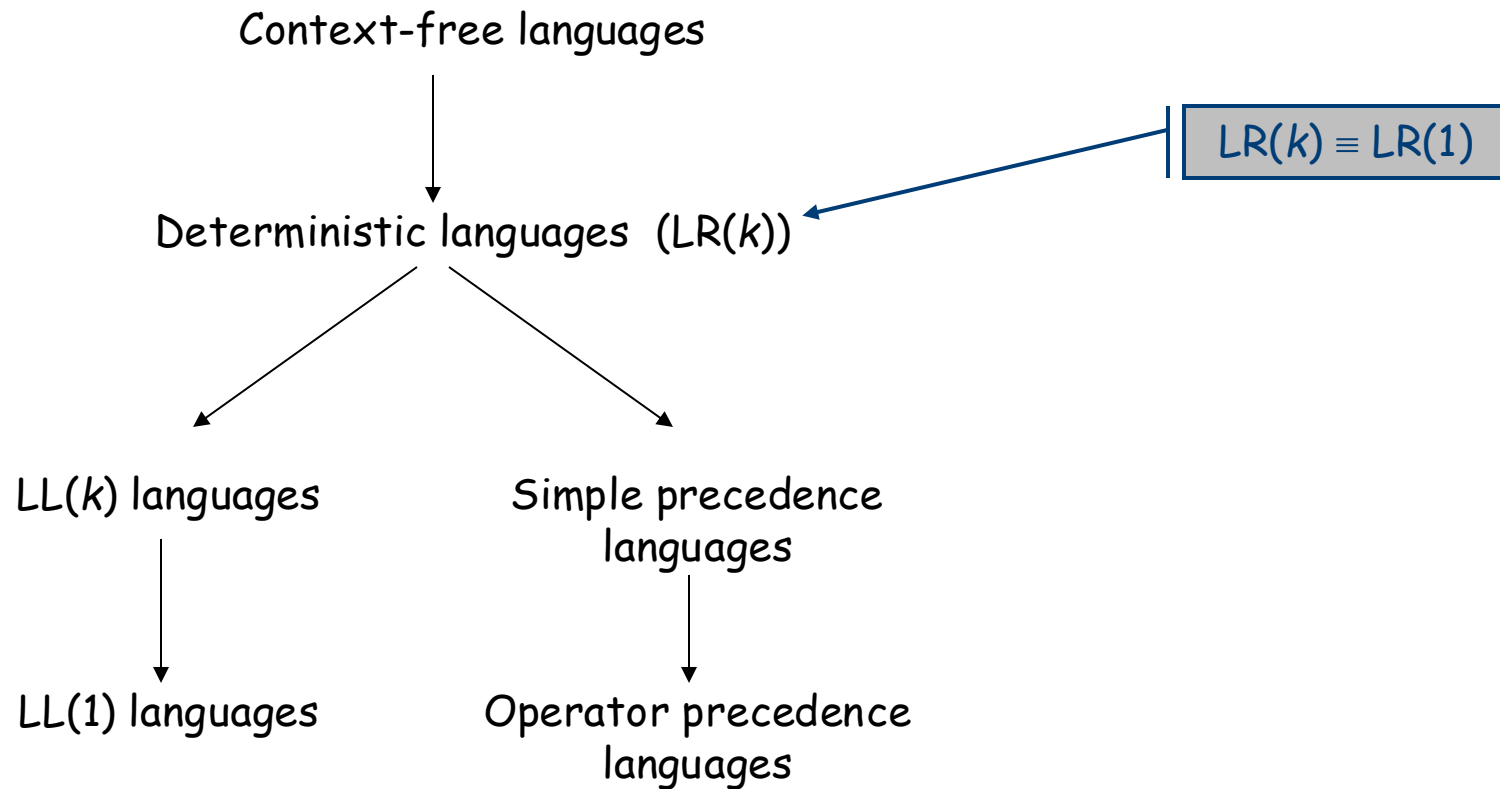
left-recursive expression grammar with precedence, see § 3.7.2 in EAC

classic space-time tradeoff

Fewer grammars, same languages



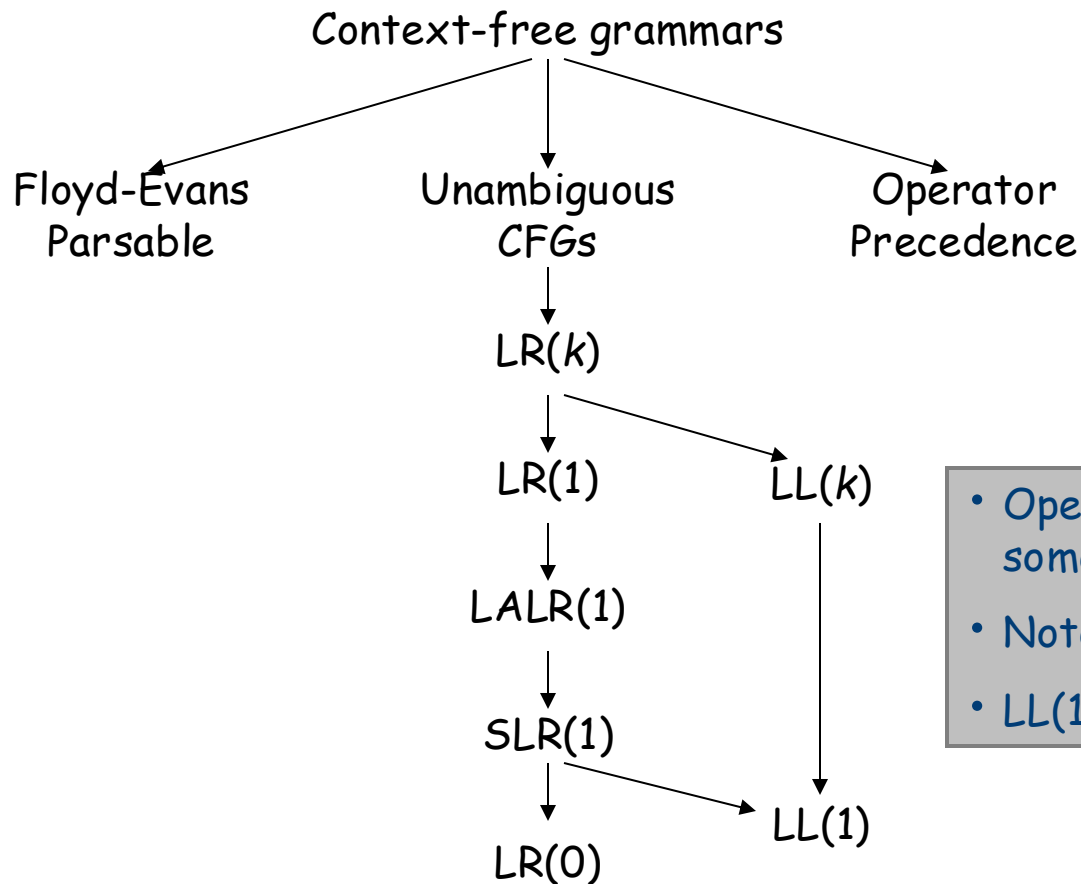
# Hierarchy of Context-Free Languages



## The inclusion hierarchy for context-free languages



# Hierarchy of Context-Free Grammars



- Operator precedence includes some ambiguous grammars
- Note sub-categories of LR(1)
- LL(1) is a subset of SLR(1)

*The inclusion hierarchy for context-free grammars*





# Summary

	<i>Advantages</i>	<i>Disadvantages</i>
<i>Top-down Recursive descent, LL(1)</i>	Fast Good locality Simplicity Good error detection	Hand-coded High maintenance Right associativity
<i>LR(1)</i>	Fast Deterministic langs. Automatable Left associativity	Large working sets Poor error messages Large table sizes