

## Parsing Wrap Up + Lab 2 Overview

### Comp 412

Copyright 2009, 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 Reductions

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 reduction 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

0	$Goal \rightarrow S$	• It requires an arbitrary lookahead to choose between $A$ & $B$
1	$S \rightarrow A$	• An LR(1) parser can carry the left context (the '(' s) until it sees $\underline{a}$ or $\underline{b}$
2	$\quad \quad   B$	• The table construction will handle it
3	$A \rightarrow ( A )$	• In contrast, an LL(1) parser cannot decide whether to expand $Goal$ by $A$ or $B$
4	$\quad \quad   \underline{a}$	→ No amount of massaging the grammar will resolve this problem
5	$B \rightarrow ( B \triangleright$	
6	$\quad \quad   \underline{b}$	

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

More formally, 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  $k$ .

2

## LR(1) versus LL(1)



The Canonical Collection of Sets of LR(1) Items

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

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

0	$Goal \rightarrow S$
1	$S \rightarrow A$
2	$\quad \quad   B$
3	$A \rightarrow ( A )$
4	$\quad \quad   \underline{a}$
5	$B \rightarrow ( B \triangleright$
6	$\quad \quad   \underline{b}$



# LR(1) versus LL(1)

And, the rest of it ...

CC <sub>9</sub>	$[A \rightarrow \bullet(A), ]], [A \rightarrow (\bullet A), ], [A \rightarrow \bullet a, ], [B \rightarrow \bullet(B \geq, \geq), [B \rightarrow (\bullet B \geq, \geq), [B \rightarrow \bullet b, \geq]$
CC <sub>10</sub>	$[A \rightarrow a \bullet, ]]$
CC <sub>11</sub>	$[B \rightarrow b \bullet, \geq]$
CC <sub>12</sub>	$[A \rightarrow (A) \bullet, EOF]$
CC <sub>13</sub>	$[B \rightarrow (B \geq) \bullet, EOF]$
CC <sub>14</sub>	$[A \rightarrow (A) \bullet, )]$
CC <sub>15</sub>	$[B \rightarrow (B \geq) \bullet, \geq]$
CC <sub>16</sub>	$[A \rightarrow (A) \bullet, )]$
CC <sub>17</sub>	$[B \rightarrow (B \geq) \bullet, \geq]$

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

Comp 412, Fall 2009

4

# LR(1) versus LL(1)

	EOF	(	)	a	≥	b	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				

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

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

5



## LR(k) versus LL(k)

### Finding Reductions

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 reduction based on

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

Thus, LR(k) examines more context

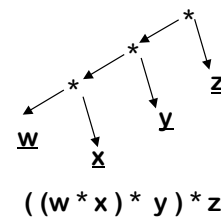
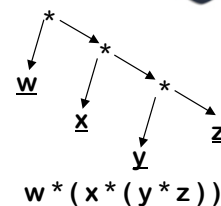
"... in practice, programming languages do not actually seem to fall in the gap between LL(1) languages and deterministic languages"

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
  - Produces right-associative operators
- Left recursion
  - Works fine in bottom-up parsers
  - Limits required stack space
  - Produces left-associative operators
- Rule of thumb
  - Left recursion for bottom-up parsers
  - Right recursion for top-down parsers

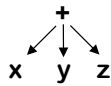




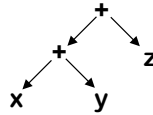
## Associativity

What difference does it make?

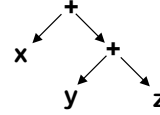
- Can change answers in floating-point arithmetic
- Can change opportunities for optimization
- Consider  $x+y+z$



*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
  - $\rightarrow$  LL(1): pop stack until we find a row with entry for ';'
  - $\rightarrow$  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  $\leftarrow$  pop()
repeat until Table[NT, ';']  $\neq$  error
  NT  $\leftarrow$  pop()
token  $\leftarrow$  NextToken()
repeat until token = ';'
  token  $\leftarrow$  NextToken()

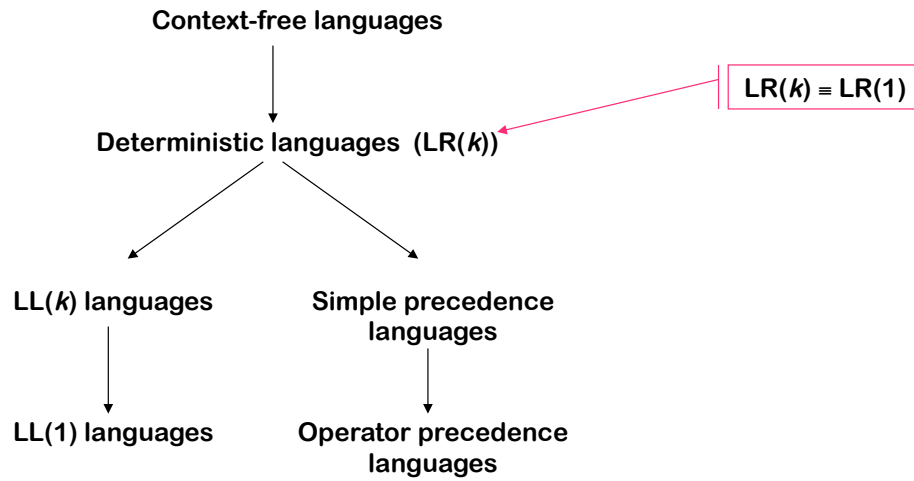
```

```

repeat until token = ';'
  shift token
  shift  $s_e$ 
  token  $\leftarrow$  NextToken()
reduce by error production
// pops all that state off stack

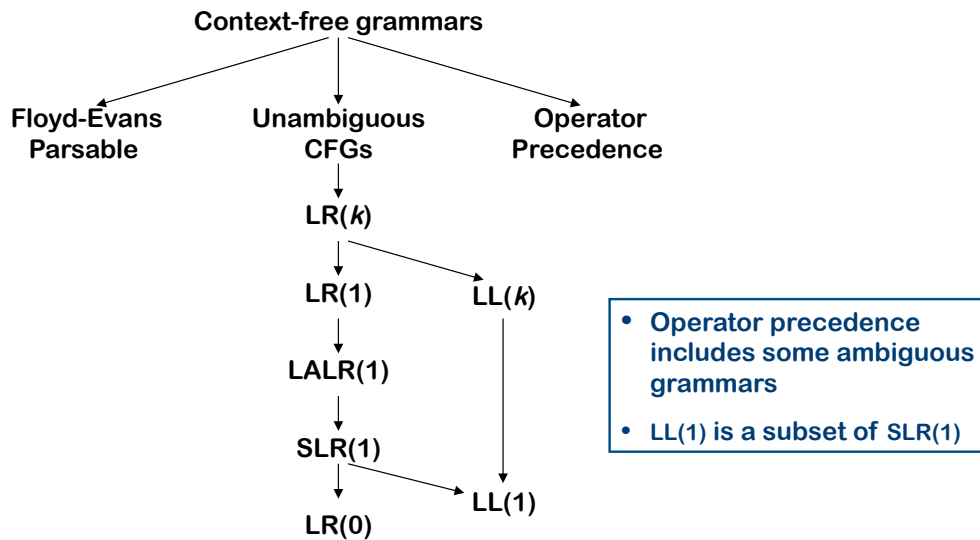
```

# Hierarchy of Context-Free Languages



*The inclusion hierarchy for context-free languages*

# Hierarchy of Context-Free Grammars



*The inclusion hierarchy for context-free grammars*

## Lab 2 Overview: An LL(1) Parser Generator



You will build a program to generate LL(1) parse tables

- Input is a modified form of Backus-Naur Form (MBNF)
- Output is a parse table and some auxiliary data structures
  - Both a pretty (human-readable) version and the C declarations
- You will write:
  - A hand-coded scanner for MBNF
  - A hand-coded, recursive descent parser for MBNF
  - A table generator (First, Follow, First<sup>+</sup> sets)
- Interim deadlines to keep you on track
  - Want you to make progress at a steady rate
- Documents will be online by Wednesday
- Skeleton parser ready in about a week, with its own docs

## Beyond Syntax



There is a level of correctness that is deeper than grammar

```
fie(a,b,c,d)
  int a, b, c, d;
{ ... }
fee() {
  int f[3],g[0],
  h, i, j, k;
  char *p;
  fie(h,i,"ab",j, k);
  k = f * i + j;
  h = g[17];
  printf("<%s,%s>.\n",
  p,q);
  p = 10;
}
```

What is wrong with this C program?  
*(let me count the ways ...)*

To generate code, we need to understand its meaning !



There is a level of correctness that is deeper than grammar

```
fie(a,b,c,d)
  int a, b, c, d;
{ ... }
fee() {
  int f[3],g[0],
  h, i, j, k;
  char *p;
  fie(h,i,"ab",j, k);
  k = f * i + j;
  h = g[17];
  printf("<%s,%s>.\n",
  p,q);
  p = 10;
}
```

What is wrong with this C program?

*(let me count the ways ...)*

- declared g[0], used g[17]
- wrong number of args to fie()
- "ab" is not an int
- wrong dimension on use of f
- undeclared variable q
- 10 is not a character string

All of these are "deeper than syntax"

## Beyond Syntax



To generate code, the compiler needs to answer many questions

- Is "x" a scalar, an array, or a function? Is "x" declared?
- Are there names that are not declared? Declared but not used?
- Which declaration of "x" does each use reference?
- Is the expression "x \* y + z" type-consistent?
- In "a[i,j,k]", does a have three dimensions?
- Where can "z" be stored? *(register, local, global, heap, static)*
- In "f ← 15", how should 15 be represented?
- How many arguments does "fie()" take? What about "printf ()" ?
- Does "\*p" reference the result of a "malloc()" ?
- Do "p" & "q" refer to the same memory location?
- Is "x" defined before it is used?

These cannot be expressed in a CFG