

# CS143: Parsing VI

David L. Dill

Stanford University

# Parsing

- Connections with Formal Language Theory
- Semantic Actions
- Error recovery

# Connections With Formal Language Theory

# Properties of Grammars vs Properties of Languages

" $G$  is  $LL(1)$ " — no conflicts in  $LL(1)$  parse table.

" $L$  is  $LL(1)$ " — there exists an  $LL(1)$  grammar for  $L$ .

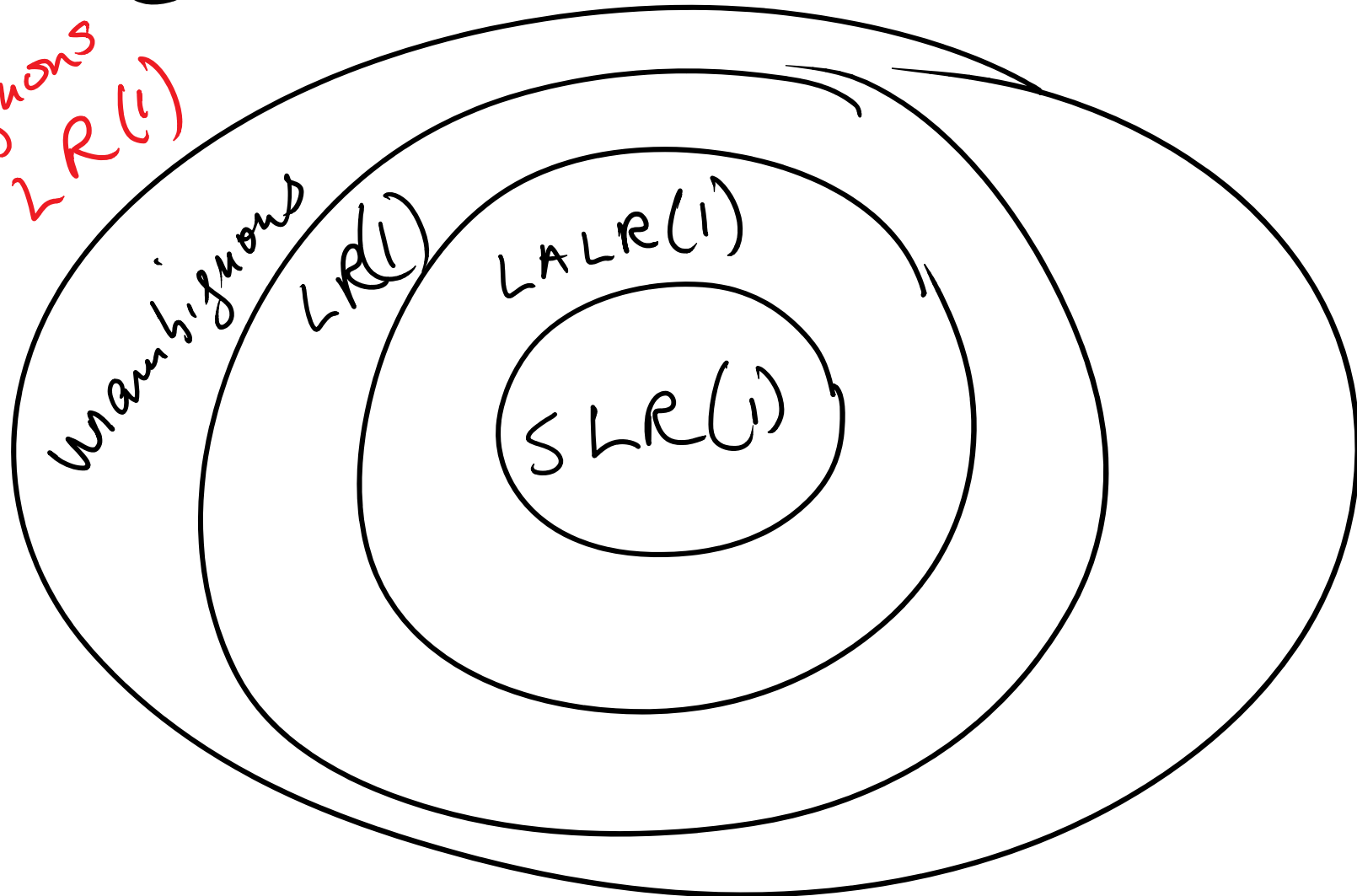
$S \rightarrow S$

$S \rightarrow a$

Is grammar LL(1)?

Is language LL(1)?

# Classes of Context-Free Grammars



$S' \rightarrow S$   
 $S \rightarrow Aa$   
 $S \rightarrow Bb$   
 $S \rightarrow ac$   
 $A \rightarrow a$   
 $B \rightarrow a$

SLR(1)

$S' \rightarrow S$   
 $S \rightarrow Aa$   
 $S \rightarrow Bb$   
 $S \rightarrow bAb$   
 $A \rightarrow a$   
 $B \rightarrow a$

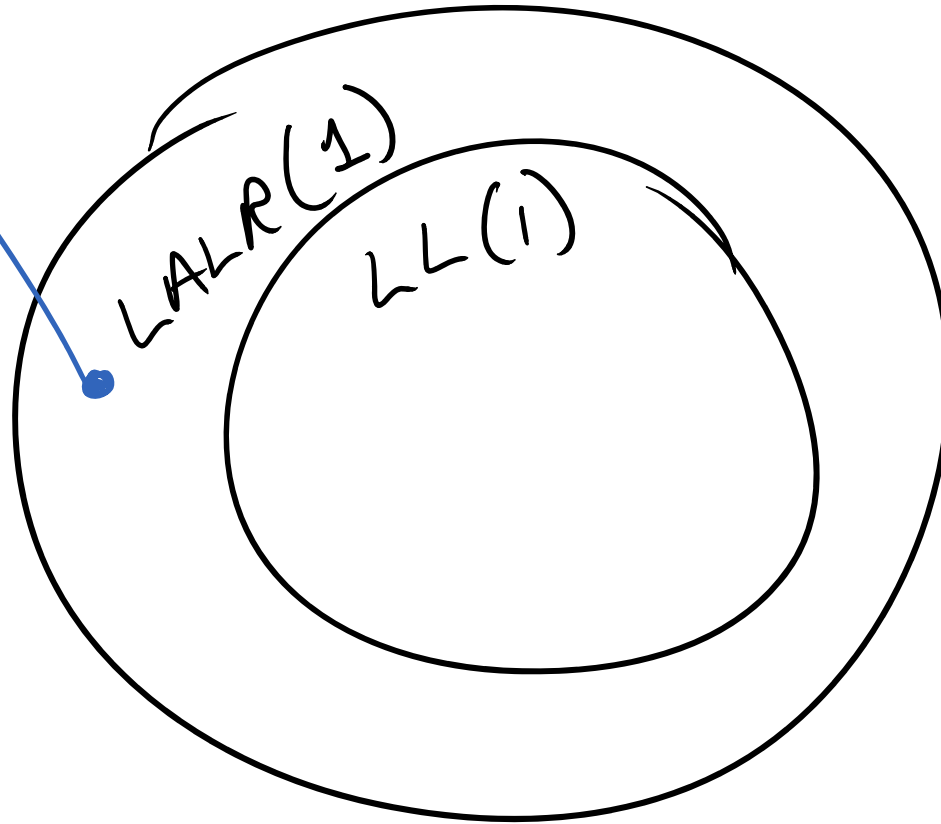
LALR(1)  
(not SLR(1))

$S' \rightarrow S$   
 $S \rightarrow Aa$   
 $S \rightarrow Bb$   
 $S \rightarrow bAb$   
 $S \rightarrow bBa$   
 $A \rightarrow a$   
 $B \rightarrow a$

LR(1)  
(not LALR(1))

# Classes of Grammars

$S \rightarrow aa$   
 $S \rightarrow ab$





Intuition for why LR parsing is more powerful than LL parsing.

$A \rightarrow \alpha$   
 $A \rightarrow \beta$

LL parsing looks at first symbol of RHS of production to decide.

$A \rightarrow \alpha \cdot a$

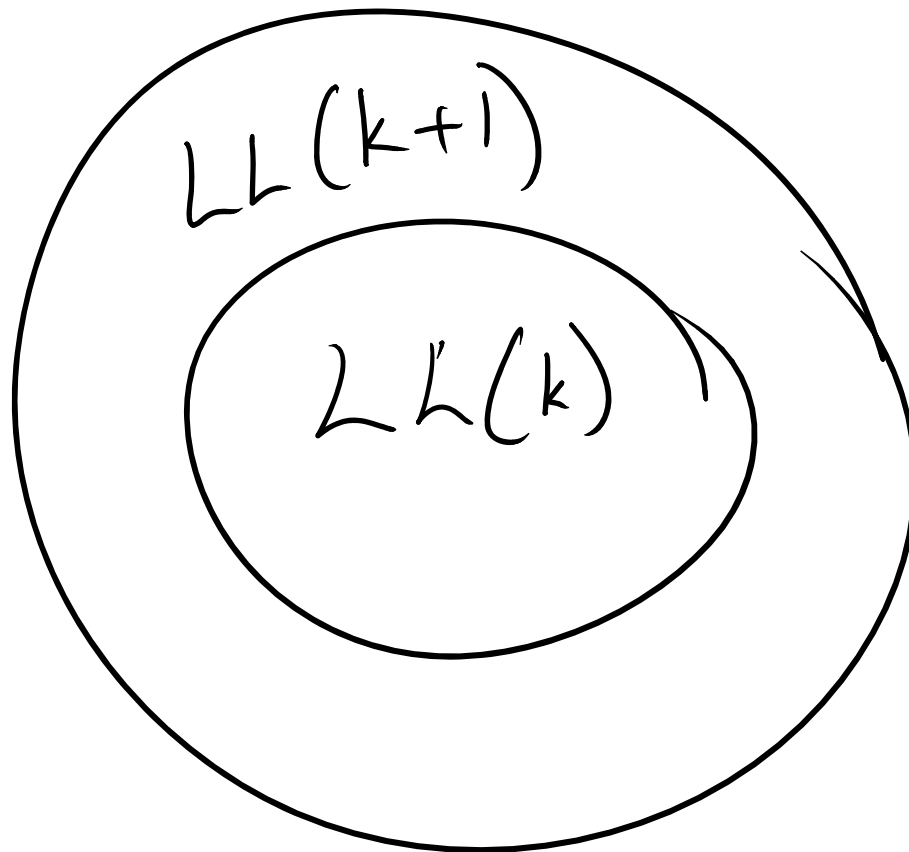
$A \rightarrow \beta \cdot b$

LR parsing uses all of RHS plus following symbol(s) to decide.

# Classes of Grammars

$S \rightarrow aa$   
 $S \rightarrow ab$

↑  
 $LL(2)$   
but not  
 $LL(1)$



# Classes of Grammars

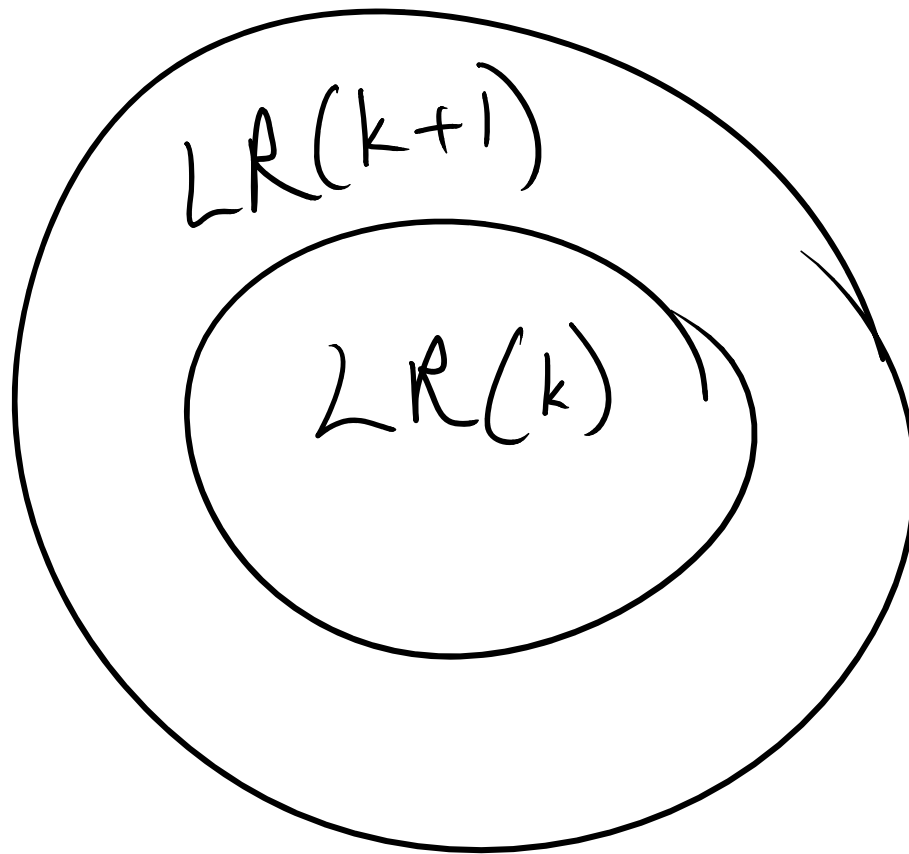
$S \rightarrow Aaa$

$S \rightarrow Bab$

$A \rightarrow a$

$B \rightarrow a$

↑  
 $LR(2)$   
but not  
 $LR(1)$



# Undecidable Problems

Equivalence of CFGs

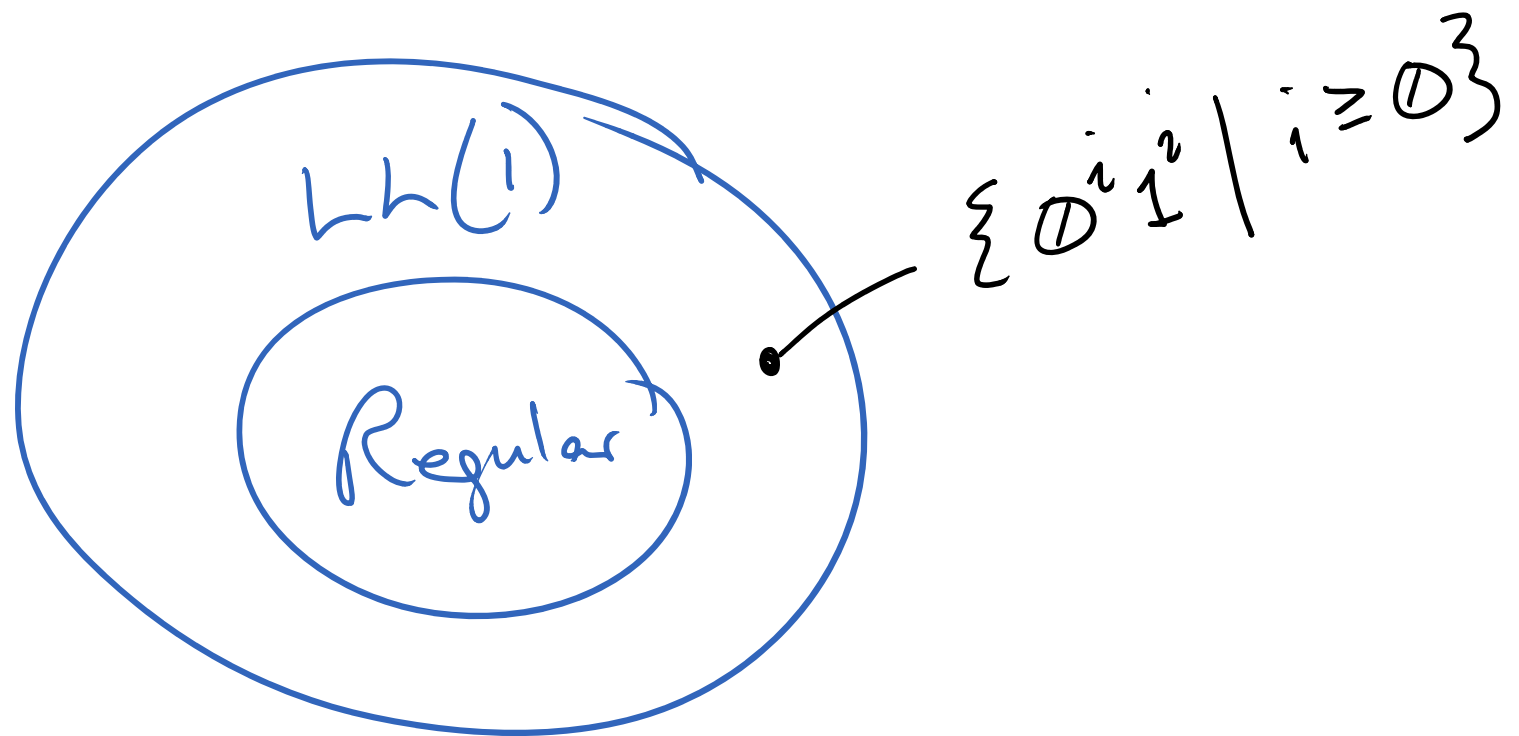
$$L(G_1) = L(G_2)?$$

Ambiguity of a CFG (!)

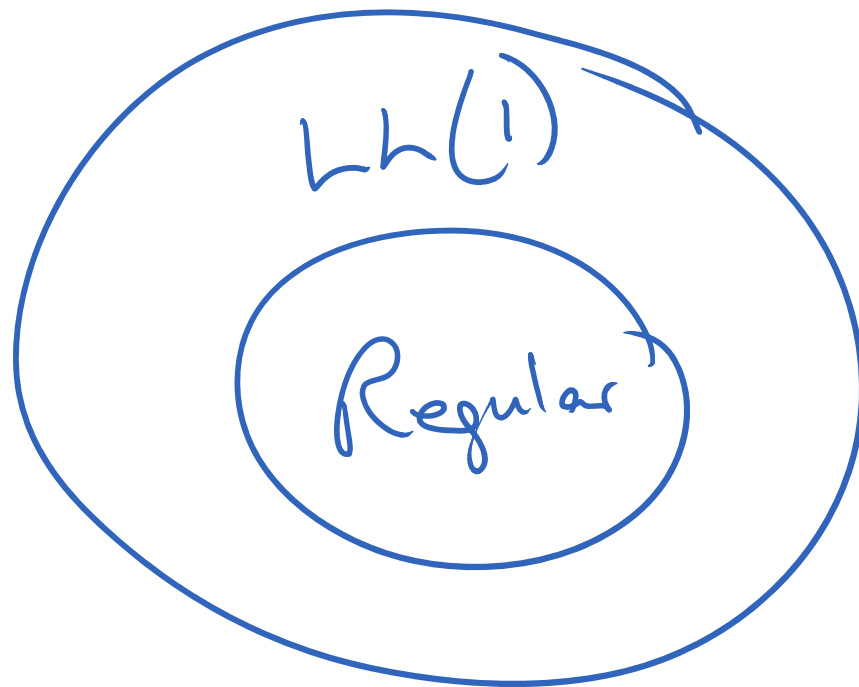
But it is decidable whether  $G$  is  $LL(1)$   $LR(1)$  etc

If  $G$  is  $LL(1)$  or  $LR(1)$  it is not ambiguous.

# Languages



# Languages



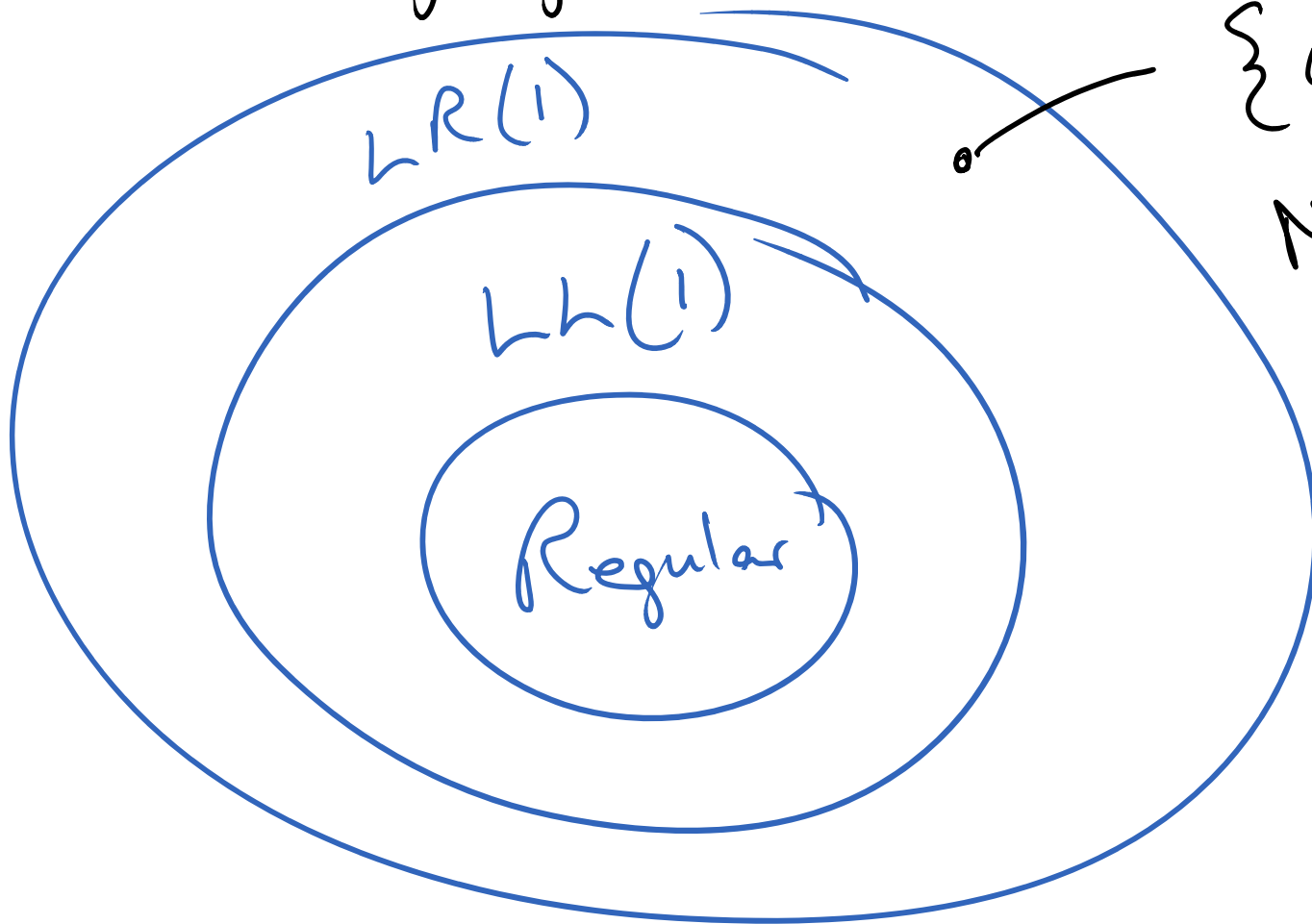
Every regular language can be written as a "right linear CFG"

$$A \rightarrow wB$$

↑

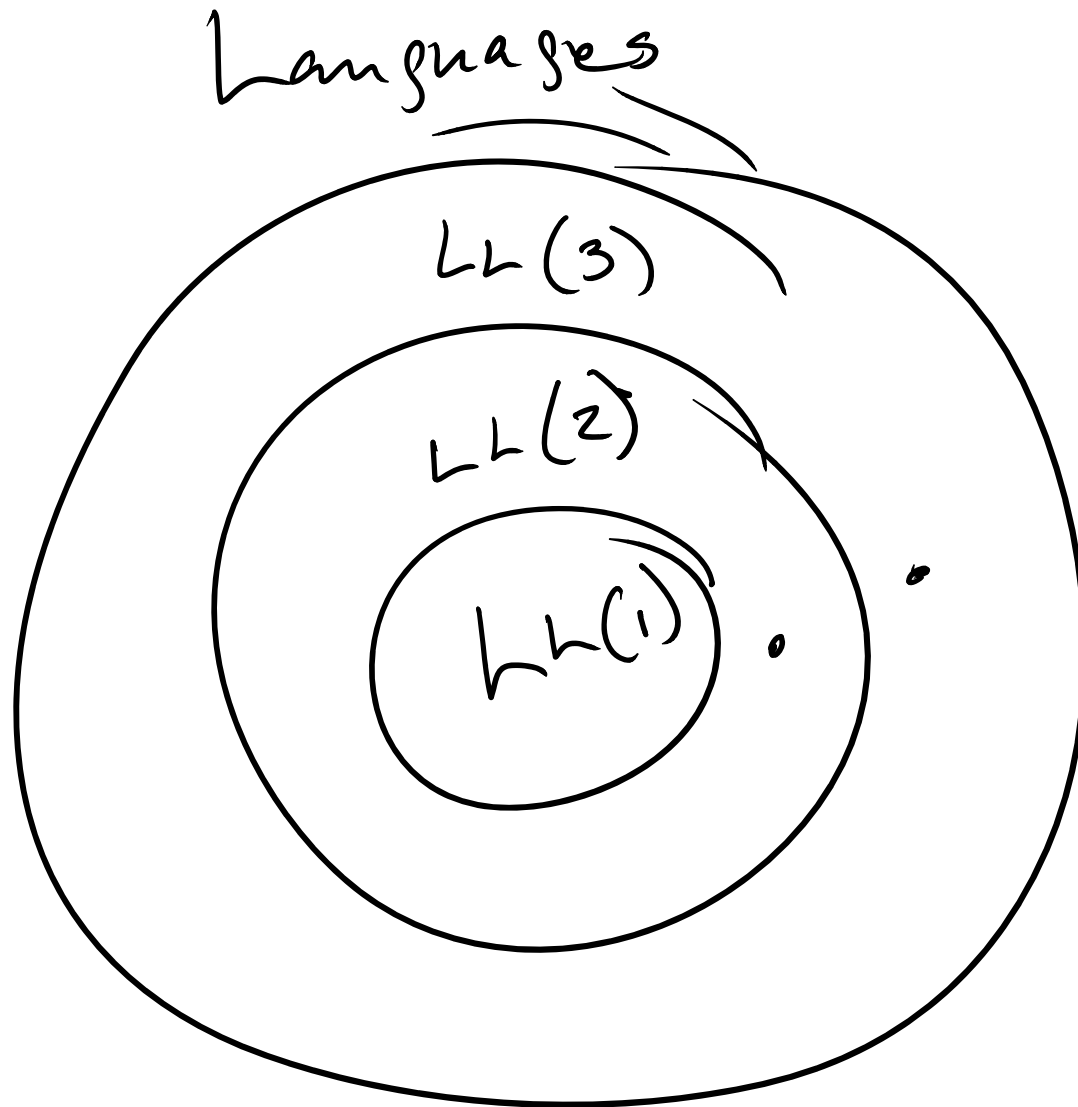
at most one nonterminal, always at right end.

Languages



$\{0^i 1^j \mid i \geq j\}$   
No LL(1) grammar!

↑  
if then else  
has the same  
problem!



← proper subsets.

There is a language that is  $LL(2)$  but not  $LL(1)$ .

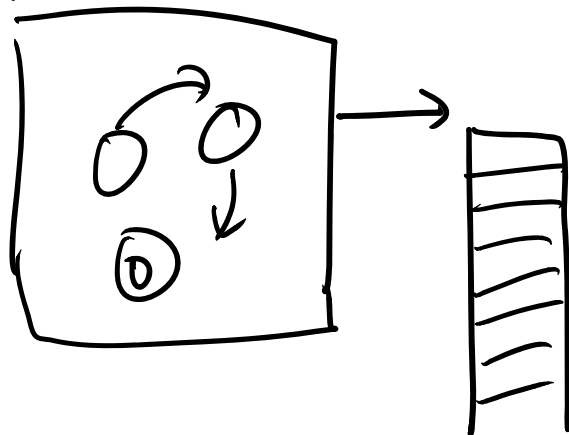


# Languages

$LR(k) = LR(1)$  for all  $k \geq 1$ .

$LR(1) = \text{Deterministic Context-Free Languages.}$

finite control.



Accepted by a  
deterministic push-down  
automaton

$\{ w w^R \mid w \in \{a,b\}^* \}$  — not LR(1)

$\{w w^R \mid w \in \{a, b\}^*\}$  - not LR(1)  
*reverse*

Context-free, but not deterministic  
There is no LR(k) grammar.

$\left( \{w c w \mid w \in \{a, b\}^*\} - \text{is LR}(1) \right)$

Equivalence of deterministic CFLs  
is decidable

Equivalence of two LR(1) CFG's is decidable

But probably not practical.

# Syntax-Directed Translation

# Build Abstract Syntax Trees

$$E \rightarrow E + E \quad \{ \$\$ = \text{plus}(\$1, \$3); \}$$

$$E \rightarrow E * E \quad \{ \$\$ = \text{mul}(\$1, \$3); \}$$

$$E \rightarrow \text{num} \quad \{ \$\$ = \$1; \}$$

Compute Directly

$$E \rightarrow E + E \quad \{ \$\$ = \$1 + \$3; \}$$

$$E \rightarrow E * E \quad \{ \$\$ = \$1 * \$3; \}$$

$$E \rightarrow \text{num} \quad \{ \$\$ = \$1; \}$$

How does this work?

## How Semantic Actions Work

New "value stack" - parallels parse stack.

Shift - value of token is pushed

Reduce - semantic action is executed.

$\$i$  is index  $TOS - (|RHS| - 1)$

$\$\$ = \$1 = \text{new TOS after reduction}$



$$\begin{aligned}
 E &\rightarrow E + E & \{ \$\$ = \$1 + \$3; \} \\
 E &\rightarrow E * E & \{ \$\$ = \$1 * \$3; \} \\
 E &\rightarrow \text{num} & \{ \$\$ = \$1; \}
 \end{aligned}$$

Stack

\$ .

\$ 1

1

\$ E

1

\$ E +

1

\$ E + 2

1 2

input

1 + 2 \* 3 \$

+ 2 \* 3 \$

+ 2 \* 3 \$

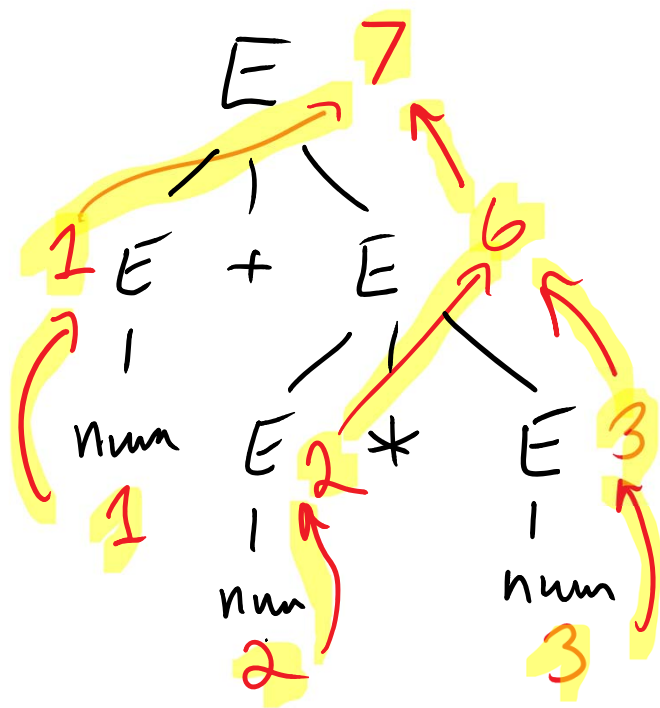
2 \* 3 \$

\* 3 \$

$E \rightarrow E + E \quad \{ \$\$ = \$1 + \$3; \}$   
 $E \rightarrow E * E \quad \{ \$\$ = \$1 * \$3; \}$   
 $E \rightarrow \text{num} \quad \{ \$\$ = \$1; \}$

stack	input
\$ E + E	* 3 \$
1 2	
\$ E + E *	3 \$
1 2	
\$ E + E * 3	\$
1 2 3	
\$ E + E * E	\$
1 2 3	
\$ E + E	\$
1 6	
\$ E	\$
7	

# Bottom-up Computation



# Announcements

- PA2 due today
- Midterm here (NVIDIA auditorium) Thursday – lecture time
- High probability midterm questions
  - Regular expressions, DFAs for something.
  - Lookahead, start conditions.
  - Which symbols/productions are useless (understand the mathematical definition!)
  - Why is CFG ambiguous?
  - What are FNE, FOLLOW, FIRST of a CFG?
  - Given LL(1) parse table and input, what is sequence of stacks/inputs?
  - Generates some LR(0), LR(1), LALR(1) states.
  - Find SLR(1), LALR(1), LR(1) conflicts
  - Given LR(1) parse table and input, what is sequence of stacks/inputs?

# Inherited Values

Values move down the tree.

$E \rightarrow \text{let } x = E \text{ in } E$

$E \rightarrow E + E$

$E \rightarrow E * E$

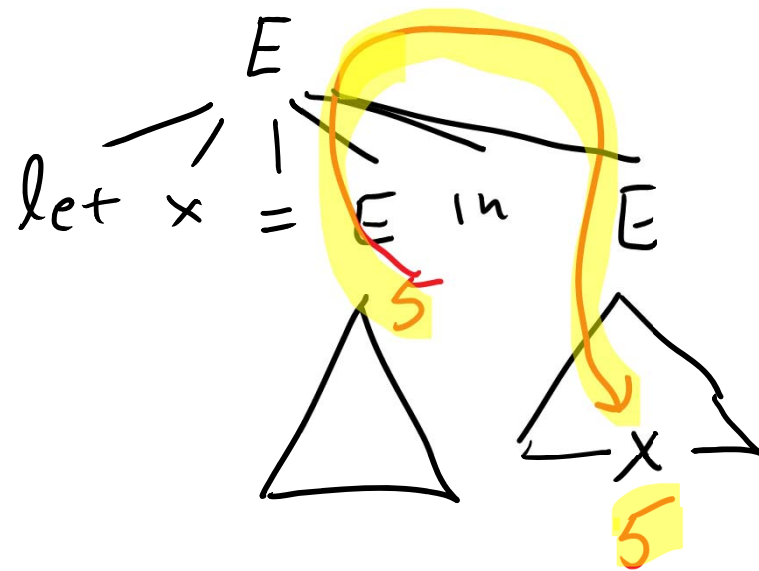
$E \rightarrow \text{num}$

$E \rightarrow x$

$E \rightarrow (E)$

one variable,  $x$ , which  
can be bound to the  
value of an expression.

← and used



Let  $x = 3$  in

$$\left( \text{let } x = 2 \text{ in } x + x \right) * x$$

Handwritten annotations: A red bracket under the expression  $x + x$  with the number 2 written below each  $x$ . A red bracket under the entire expression  $\left( \text{let } x = 2 \text{ in } x + x \right) * x$  with the number 3 written below the  $x$  at the end.

$$= 12$$

Action "in the middle" of a production

$A \rightarrow B \{ \$\$ = \$1 \} C \{ \$\$ = \$2 + \$3 \}$

Bison makes a new nonterminal  
↑ this refers to value of M

$A \rightarrow B M C \{ \$\$ = \$2 + \$3 \}$

$M \rightarrow \epsilon \{ \$\$ = \$1; \}$



Action "in the middle" of a production

$$A \rightarrow B \{ \$\$ = \$1 \} C \{ \$\$ = \$2 + \$3 \}$$

Bison makes a new nonterminal

$$A \rightarrow B M C \{ \$\$ = \$2 + \$3 \}$$
$$M \rightarrow \epsilon \{ \$\$ = \$1; \}$$

↑  
But this refers to "B"  
in the  $A \rightarrow B M C$  production

# Inherited Values *global variable*

$E \rightarrow \text{let } x = E \{ x = \$4; \} \text{ in } E \{ \$\$ = \$7 \}$

$E \rightarrow E + E \quad \{ \$\$ = \$1 + \$3; \}$

$E \rightarrow E * E \quad \{ \$\$ = \$1 + \$3; \}$

$E \rightarrow \text{num} \quad \{ \$\$ = \$1; \}$

$E \rightarrow x \quad \{ \$\$ = x; \}$

$E \rightarrow (E) \quad \{ \$\$ = \$2; \}$

# Inherited Values

\$5

$E \rightarrow \text{let } x = E \{ x = \$4; \} \text{ in } E \{ \$\$ = \$7 \}$

$E \rightarrow E + E \quad \{ \$\$ = \$1 + \$3; \}$

$E \rightarrow E * E \quad \{ \$\$ = \$1 + \$3; \}$

$E \rightarrow \text{num} \quad \{ \$\$ = \$1; \}$

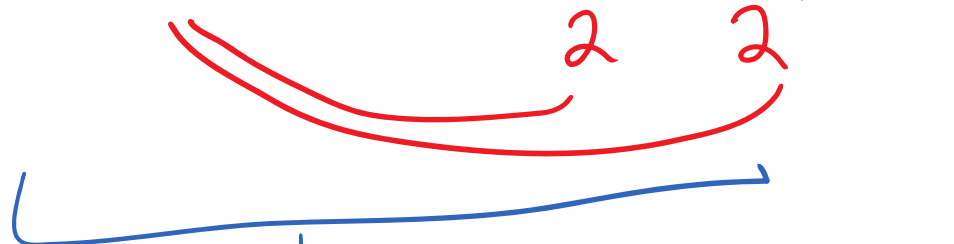
$E \rightarrow x \quad \{ \$\$ = x; \}$

$E \rightarrow (E) \quad \{ \$\$ = \$2; \}$

# Problem

Let  $x = 3$  in

$(\text{let } x = 2 \text{ in } x + x) * x$

  
When this is parsed,  
it sets  $x = 2$

↑ so parser uses  
2, not 3, for  
this.

## Solutions

1. Store  $x$  values in a global stack

$E \rightarrow \text{Let } x = E \{ x_{\text{stack}}, \text{push}(\$4); \}$  in  $E$

$\uparrow$   
new value of  $x$   
goes on top of stack.

## Solutions

1. Store  $x$  values in a global stack

$E \rightarrow \text{Let } x = E \{ xstack, \text{push}(\$4); \}$  in  $E$   
 $\{ xstack.\text{pop}(); \$\$ = \$7; \}$   
 $\uparrow$  restore outer  $x$  value.

## Solutions

1. Store  $x$  values in a global stack

$E \rightarrow \text{Let } x = E \{ xstack.push(\$4); \} \text{ in } E$   
 $\{ xstack.pop(); \$\$ = \$7; \}$

$E \rightarrow x \{ \$\$ = xstack.top(); \}$

get current  $x$  value

# Solutions

2. Use parser value stack

$E \rightarrow \text{Let } x = E \text{ in } \{ \$\$ = \$4; \} E \{ \$\$ = \$7; \}$

$E \rightarrow x \{ \$\$ = \$0 \}$   save x value here



# Solutions

2. Use parser value stack

$E \rightarrow \text{Let } x = E \text{ in } \{ \$\$ = \$4; \} E \{ \$\$ = \$7; \}$

$E \rightarrow x \{ \$\$ = \$0 \}$

↑ save x value here

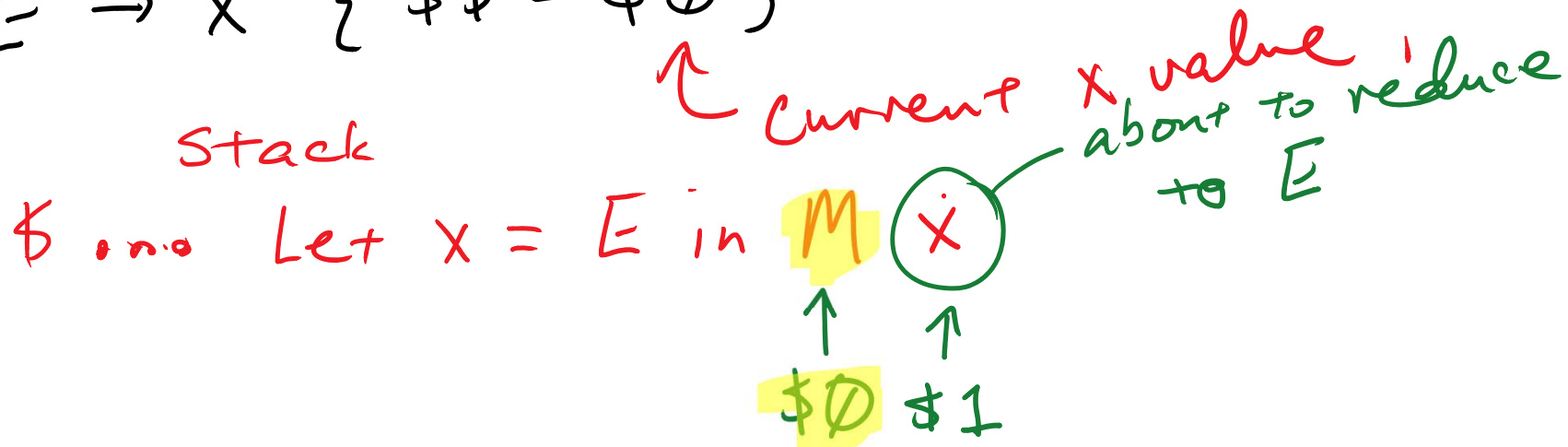
x value gets popped off stack when this is reduced.

# Solutions

2. Use parser value stack

$E \rightarrow \text{Let } x = E \text{ in } \{ \$\$ = \$4; \} E \{ \$\$ = \$7; \}$

$E \rightarrow x \{ \$\$ = \$0 \}$



# Error Recovery

# Challenge

Report errors accurately

"Read user's mind"

Especially hard after the first  
error

# Challenge

Report errors accurately

"Read user's mind"

Especially hard after the first error

Most frequent reason for handwritten parsers

## Panic Mode

Delete tokens until "synchronizing token" is found (e.g. "j")

Resume parsing at that point

error token

YACC and descendants

On parse error

1. Pop stack until top state has an "error" transition
2. Take the transition
3. Delete symbols until it is possible to follow 3 tokens successfully.
4. Resume parsing

Add common errors to lexer / parser

Spelling errors for keywords

Error productions

$E \rightarrow E E$  for omitted  $*$

Problems:

Complexity, conflicts



More Elaborate

Find "closest" syntactically correct program.

"Closest" = minimum edit distance

= minimum number of delete/insert/change moves to get to correct program.

1. Computationally costly.
2. Doesn't read user's mind