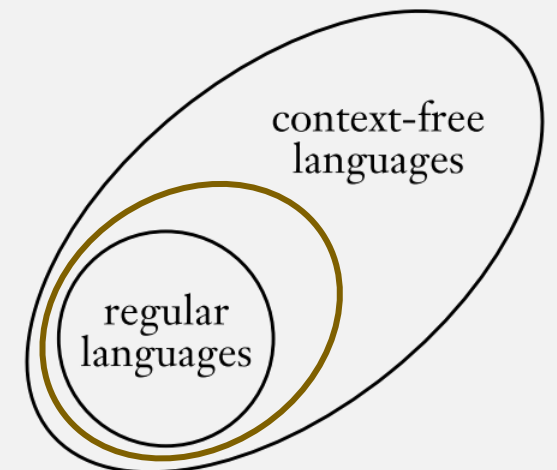


# Regular vs Context-Free Languages (and others?)

Monday, March 20, 2023

UMB CS 420



# *Announcements*

- HW 5 in
  - ~~Due Sunday March 19, 2023 11:59pm EST~~
- HW 6 out
  - Due Sunday March 26, 2023 11:59pm EST

## Quiz Preview

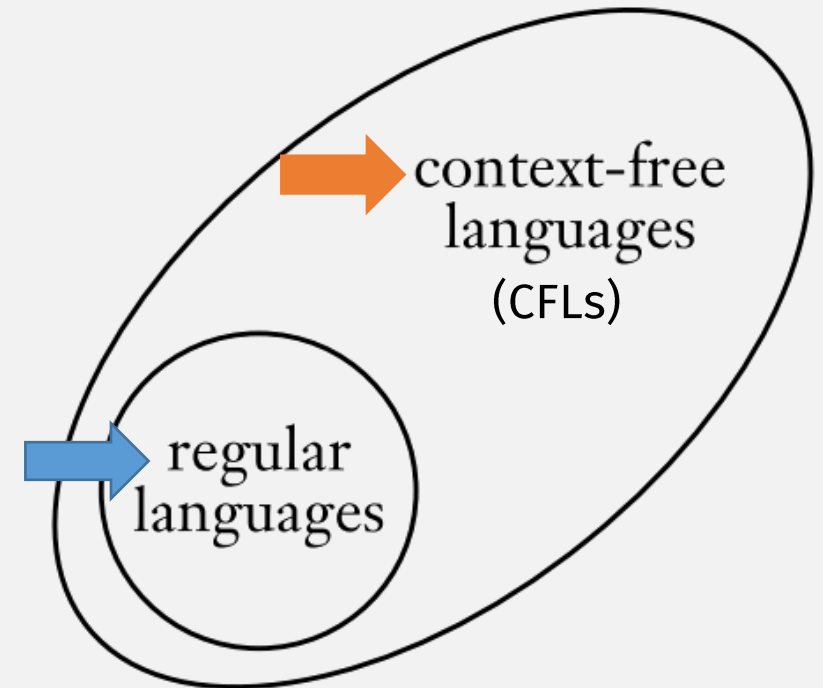
- Is every CFL a regular language?

# Is This Diagram “Correct”?

(What are the statements implied by this diagram?)

➡ 1. Every regular language is a CFL

➡ 2. Not every CFL is a regular language



# How to Prove This Diagram “Correct”?

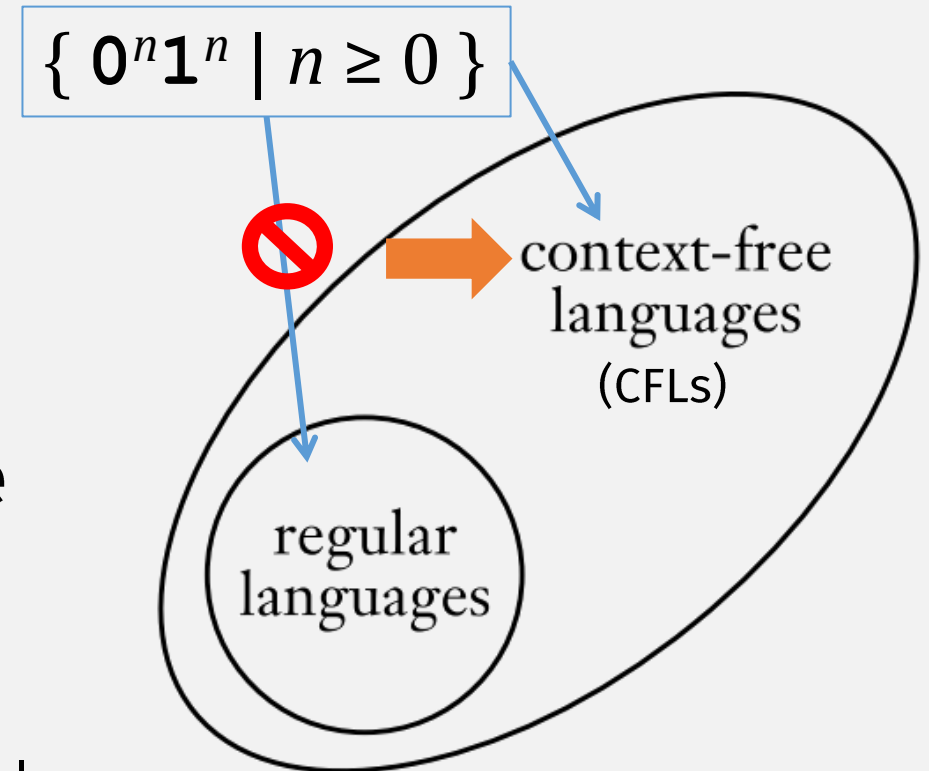
1. Every regular language is a CFL

➔ 2. Not every CFL is a regular language

Find a CFL that is not regular

$\{ 0^n 1^n \mid n \geq 0 \}$

- It's a CFL
  - *Proof:* CFG  $S \rightarrow 0S1 \mid \epsilon$
- It's not regular
  - *Proof:* by contradiction using the Pumping Lemma



# How to Prove This Diagram “Correct”?

➔ 1. Every regular language is a CFL

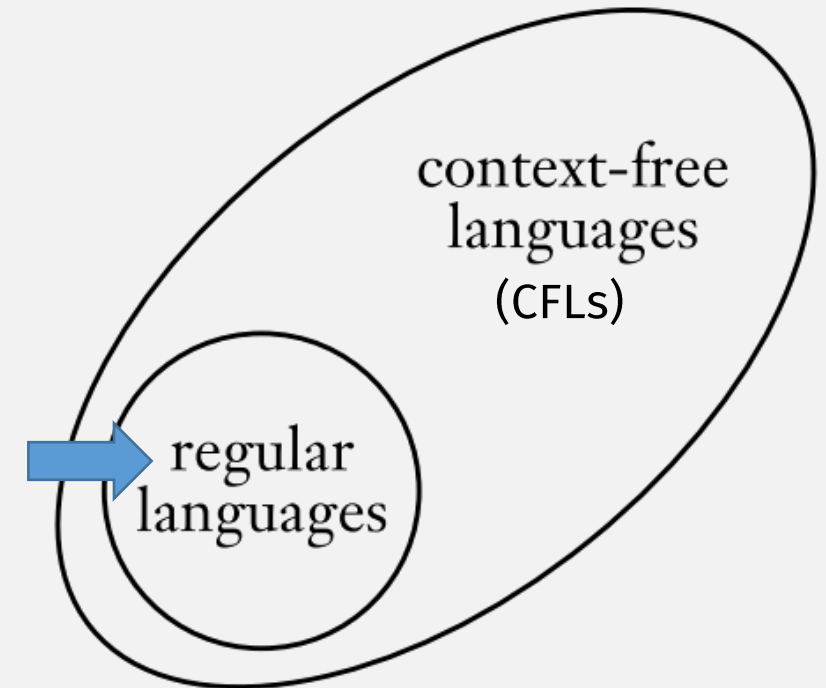
For any regular language  $A$ , show ...

... it has a CFG or PDA

☑ 2. Not every CFL is a regular language

A regular language is represented by a:

- DFA
- NFA
- Regular Expression



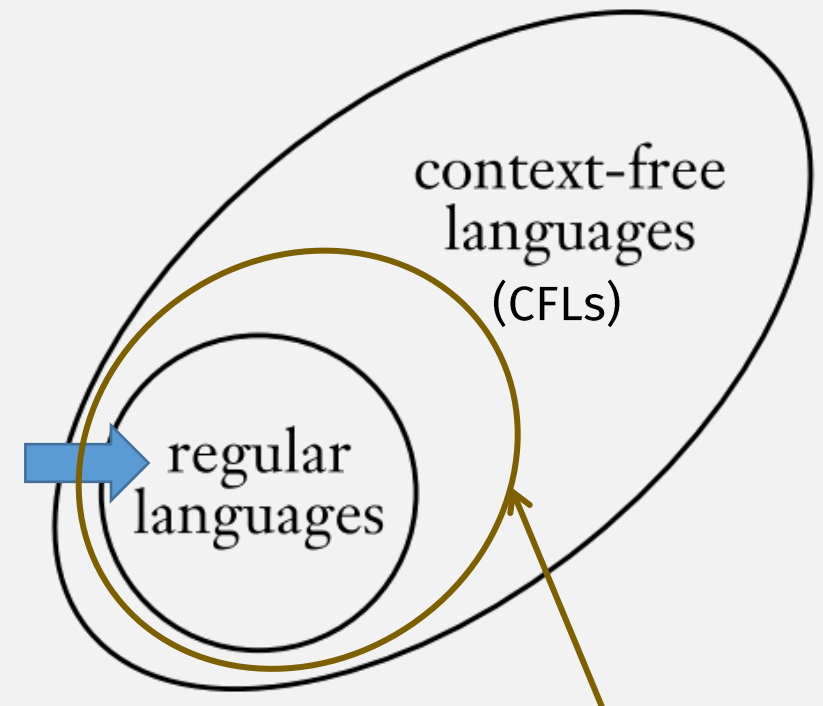
# Regular Languages are CFLs: 3 Ways to Prove

- DFA → CFG or PDA

- NFA → CFG or PDA

See HW 6!

- Regular expression → CFG or PDA



Are there other interesting subsets of CFLs?

# **Deterministic CFLs and DPDAs**

# Previously: Generating Strings

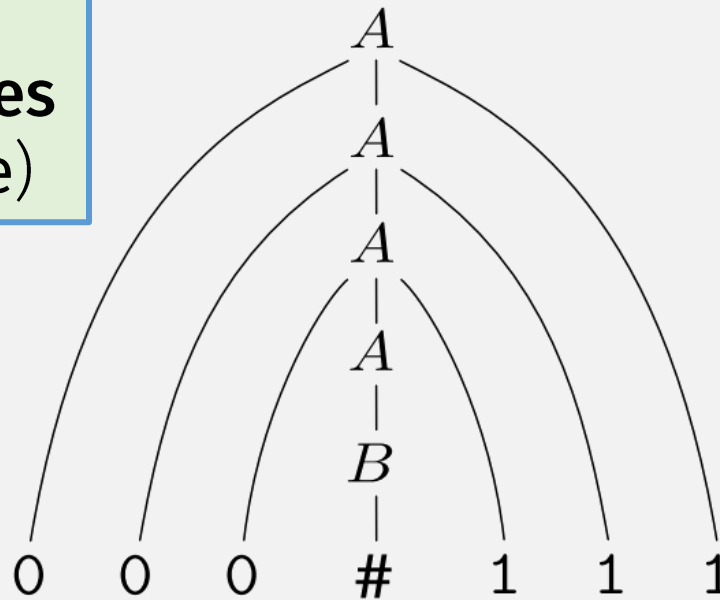
## Generating strings:

1. Start with **start variable**,
2. Repeatedly apply CFG rules to get string (and parse tree)

$$A \rightarrow 0A1$$

$$A \rightarrow B$$

$$B \rightarrow \#$$



$$A \Rightarrow 0A1 \Rightarrow 00A11 \Rightarrow 000A111 \Rightarrow 000B111 \Rightarrow 000\#111$$



# Generating vs Parsing

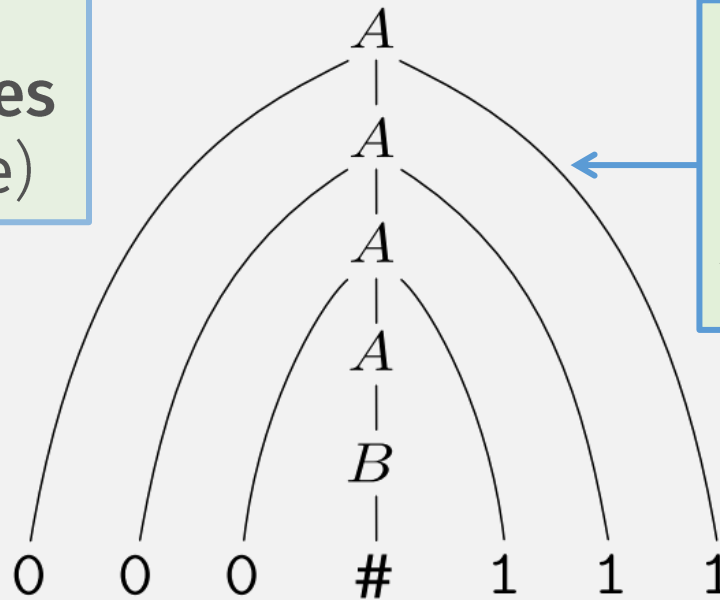
## Generating strings:

1. Start with **start variable**,
2. Repeatedly apply CFG rules to get string (and parse tree)

$$A \rightarrow 0A1$$

$$A \rightarrow B$$

$$B \rightarrow \#$$



In practice, opposite is more interesting:

1. Start with a **string**,
2. Then parse it into **parse tree**

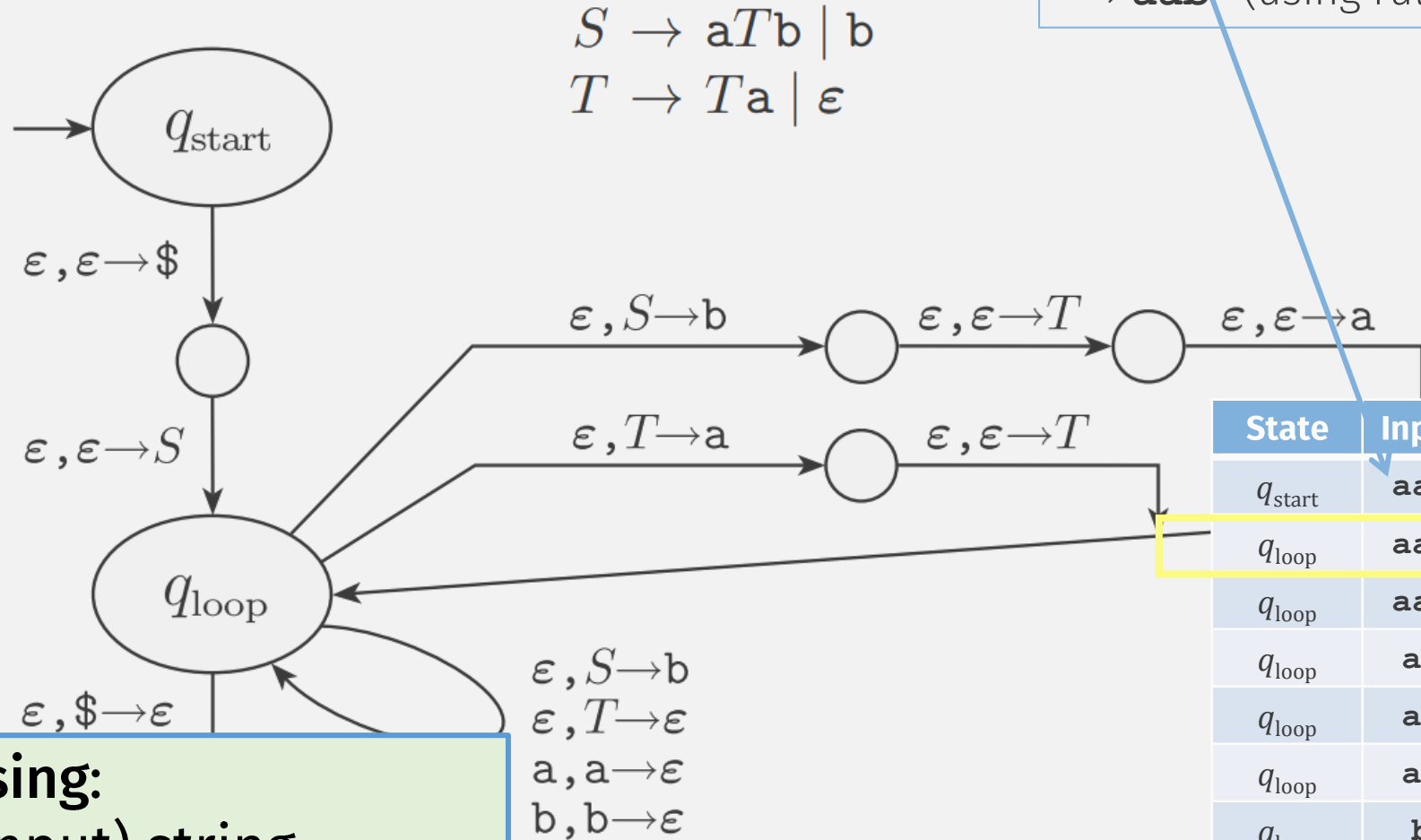
$$A \Rightarrow 0A1 \Rightarrow 00A11 \Rightarrow 000A111 \Rightarrow 000B111 \Rightarrow 000\#111$$

# Generating vs Parsing

- In practice, **parsing** a string more important than **generating** one
  - E.g., a **compiler** (first) parses source code into a parse tree
  - (Actually, *any* program with string inputs must first parse it)

# Previously: Example CFG $\rightarrow$ PDA

Example Derivation using CFG:  
 $S \Rightarrow aTb$  (using rule  $S \rightarrow aTb$ )  
 $\Rightarrow aTab$  (using rule  $T \rightarrow Ta$ )  
 $\Rightarrow aab$  (using rule  $T \rightarrow \epsilon$ )



PDA Example

State	Input	Stack	Equiv Rule
$q_{start}$	aab		
$q_{loop}$	aab	S\$	
$q_{loop}$	aab	aTb\$	$S \rightarrow aTb$
$q_{loop}$	ab	Tb\$	
$q_{loop}$	ab	Tab\$	$T \rightarrow Ta$
$q_{loop}$	ab	ab\$	$T \rightarrow \epsilon$
$q_{loop}$	b	b\$	
$q_{loop}$		\$	
$q_{accept}$			

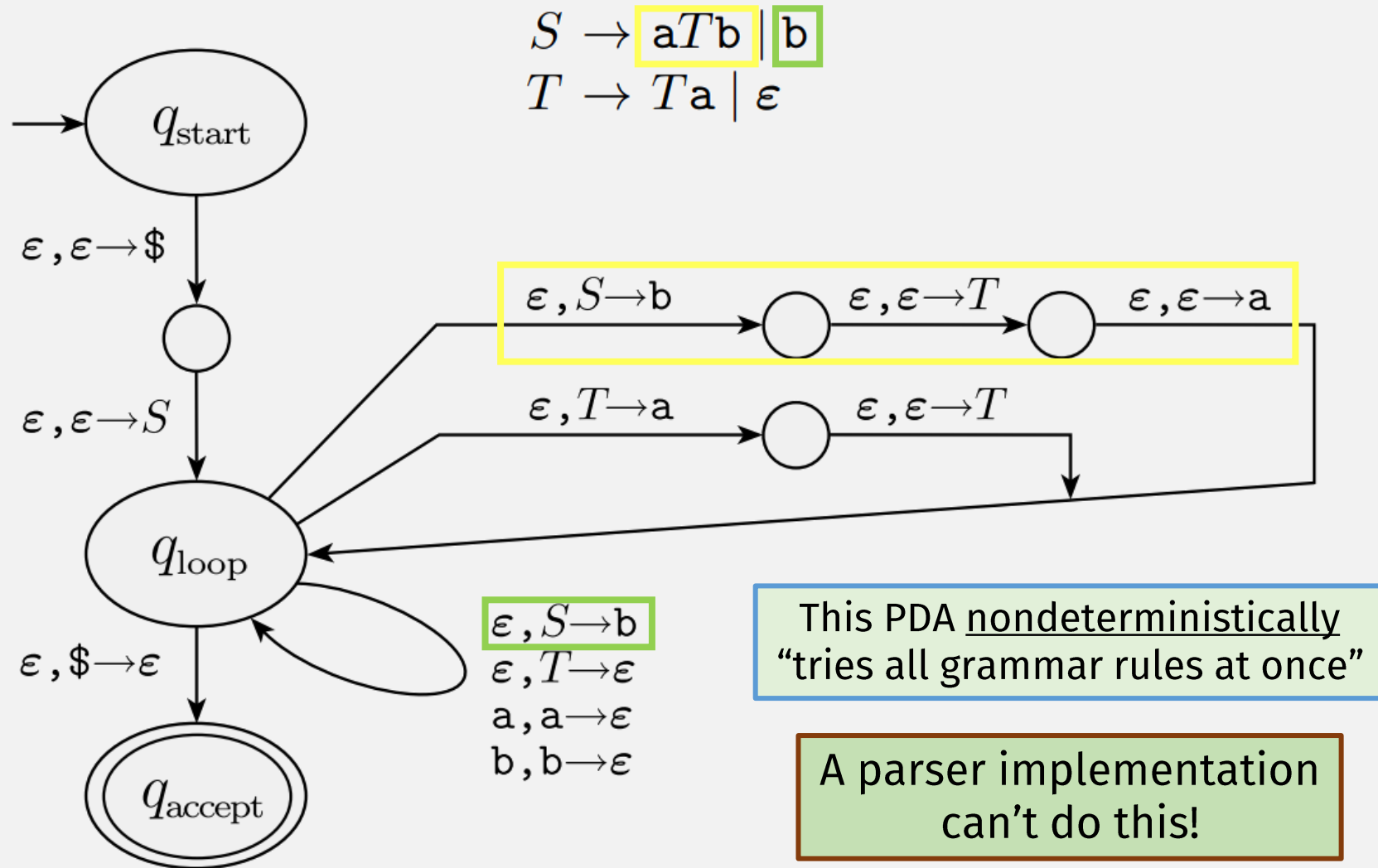
Machine is parsing:

1. Start with (input) string,
2. Find rules that **generate** string

# Generating vs Parsing

- In practice, **parsing** a string more important than **generating** one
  - E.g., a **compiler** (first step) parses source code into a parse tree
  - (Actually, *any* program with string inputs must first parse it)
- But: the PDAs we've seen are non-deterministic (like NFAs)

# Previously: (Nondeterministic) PDA



# Generating vs Parsing

- In practice, **parsing** a string more important than **generating** one
  - E.g., a **compiler** (first step) parses source code into a parse tree
  - (Actually, *any* program with string inputs must first parse it)
- But: the PDAs we've seen are non-deterministic (like NFAs)
- Compiler's parsing algorithm must be deterministic
- So: to model parsers, we need a **Deterministic PDA (DPDA)**

# DPDA: Formal Definition

The language of a DPDA is called a *deterministic context-free language*.

A *deterministic pushdown automaton* is a 6-tuple  $(Q, \Sigma, \Gamma, \delta, q_0, F)$ , where  $Q, \Sigma, \Gamma$ , and  $F$  are all finite sets, and

1.  $Q$  is the set of states,
2.  $\Sigma$  is the input alphabet,
3.  $\Gamma$  is the stack alphabet,
4.  $\delta: Q \times \Sigma_\epsilon \times \Gamma_\epsilon \longrightarrow (Q \times \Gamma_\epsilon) \cup \{\emptyset\}$  is the transition function
5.  $q_0 \in Q$  is the start state, and
6.  $F \subseteq Q$  is the set of accept states.

“do nothing”

A *pushdown automaton* is a 6-tuple

1.  $Q$  is the set of states,
2.  $\Sigma$  is the input alphabet,
3.  $\Gamma$  is the stack alphabet,
4.  $\delta: Q \times \Sigma_\epsilon \times \Gamma_\epsilon \longrightarrow \mathcal{P}(Q \times \Gamma_\epsilon)$
5.  $q_0 \in Q$  is the start state, and
6.  $F \subseteq Q$  is the set of accept states.

Difference: DPDA has only one possible action, for any given state, input, and stack op (similar to DFA vs NFA)

Must take into account  $\epsilon$  reads or stack ops!  
E.g., if  $\delta(q, a, X)$  does “something”, then  $\delta(q, \epsilon, X)$  must “do nothing”

# DPDAs are Not Equivalent to PDAs!

$$\begin{aligned} R &\rightarrow S \mid T \\ S &\rightarrow aSb \mid ab \\ T &\rightarrow aTbb \mid abb \end{aligned}$$

- A PDA can non-deterministically “try all rules” (abandoning failed attempts);
- A DPDA must choose one rule at each step! (cant go back after reading input!)

used  $S$  rule

$aa\underline{abb}$   $\rightsquigarrow$   $aa\underline{S}bb$

**Parsing** = deriving reversed:  
start with string, end with parse tree

used  $T$  rule

$aa\underline{bbbbb}$   $\rightsquigarrow$   $aa\underline{T}bbbb$

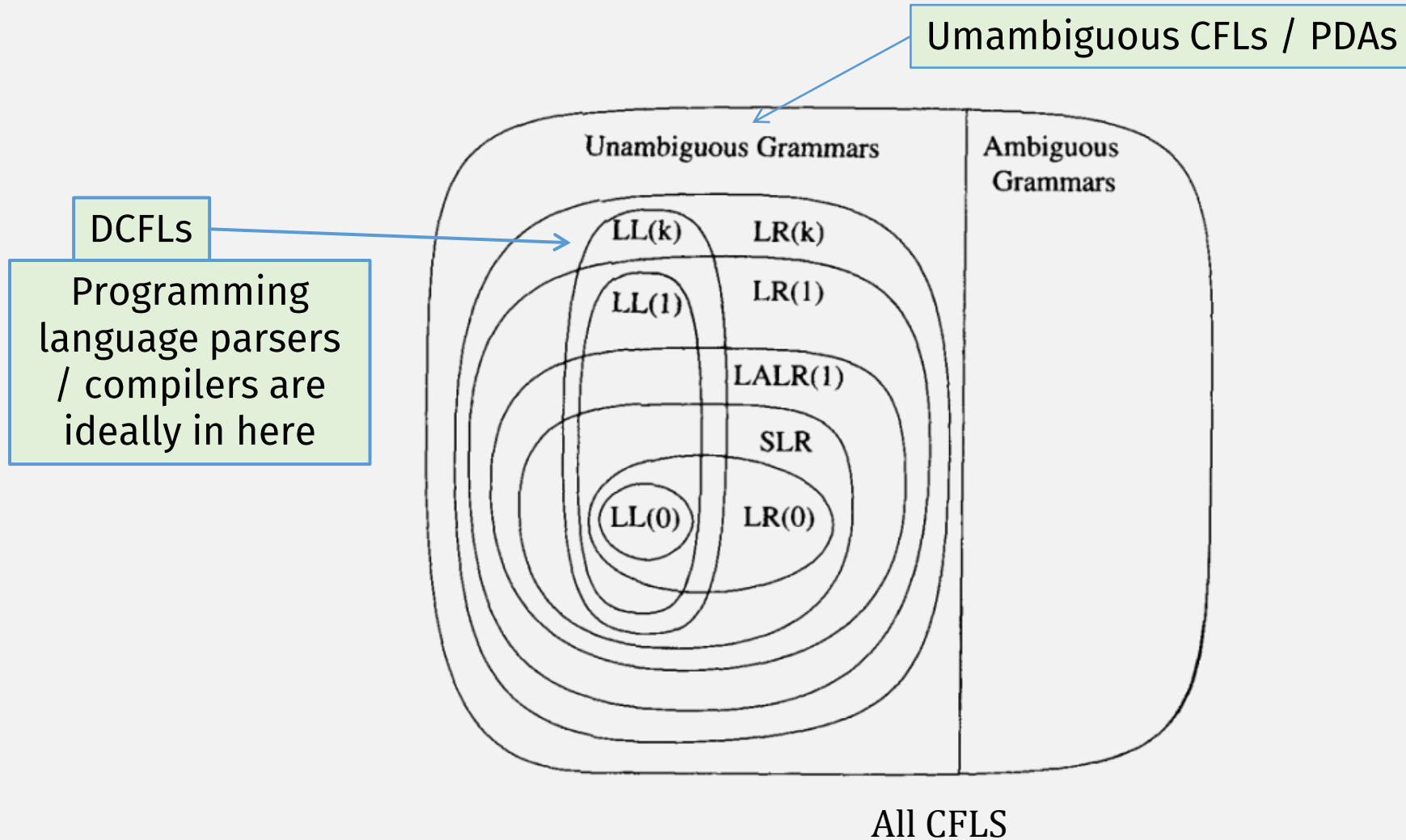
When parsing this string, when does it know which rule was used,  $S$  or  $T$ ?

Choosing “correct” rule depends on rest of the input!

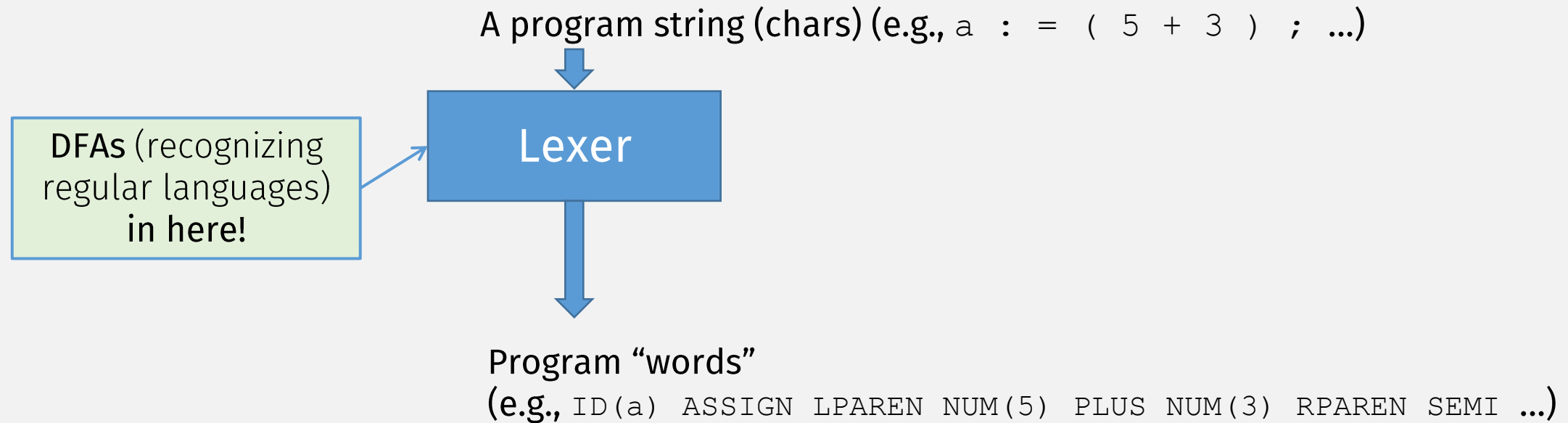
PDAs recognize CFLs, but **DPDAs** only recognize **DCFLs!** (a subset of CFLs)



# Subclasses of CFLs



# Compiler Stages



# A Lexer Implementation

```
%{
/* C Declarations: */
#include "tokens.h" /* definitions of IF, ID, NUM, ... */
#include "errmsg.h"
union {int ival; string sval; double fval;} yylval;
int charPos=1;
#define ADJ (EM_tokPos=charPos, charPos+=yyleng)
}%
/* Lex Definitions: */
digits [0-9]+
%%
/* Regular Expressions and Actions: */
if {ADJ; return IF;}
[a-z][a-z0-9]* {ADJ; yylval.sval=String(yytext);
                return ID;}
{digits} {ADJ; yylval.ival=atoi(yytext);
          return NUM;}
({digits} "." [0-9]*) | ([0-9]* "." {digits}) {ADJ;
                                              yylval.fval=atof(yytext);
                                              return REAL;}
("--" [a-z]* "\n") | (" " | "\n" | "\t")+ {ADJ;}
. {ADJ; EM_error("illegal character");}
```

DFAs  
(represented  
as regular  
expressions)!

Remember our analogy:  
- DFAs are like programs  
- All possible DFA tuples is like  
a programming language

This DFA is a real program!

A “lex” tool converts the  
program:  
- from “DFA Lang” ...  
- to an **equivalent** one in C !

# Compiler Stages

A program (chars) (e.g., `a := ( 5 + 3 ) ; ...`)

Lexer

DFAs (recognizing regular languages) in here!

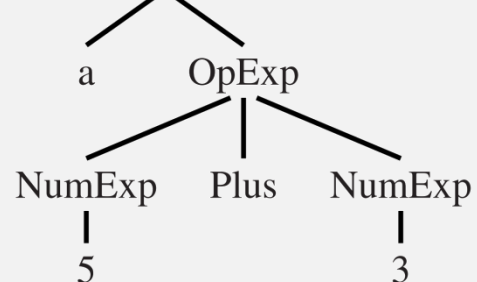
Program "words"

(e.g., `ID(a) ASSIGN LPAREN NUM(5) PLUS NUM(3) RPAREN SEMI ...`)

Parser

DPDAs (recognizing DCFLs) in here!

AssignStm Abstract Syntax tree (AST), i.e., a **parse tree!**



# A Parser Implementation

```
%{
int yylex(void);
void yyerror(char *s) { EM_error(EM_tokPos, "%s", s); }
}%
%token ID WHILE BEGIN END DO IF THEN ELSE SEMI ASSIGN
%start prog
%%

prog: stmlist

stm : ID ASSIGN ID
    | WHILE ID DO stm
    | BEGIN stmlist END
    | IF ID THEN stm
    | IF ID THEN stm ELSE stm

stmlist : stm
        | stmlist SEMI stm
```

Just write  
the CFG!

Remember our analogy:  
CFGs are like **programs**

This CFG is a real program!

A “yacc” tool converts the  
program:  
- from “CFG Lang” ...  
- to an **equivalent** one in C !

# DPDAs are Not Equivalent to PDAs!

Parsing = generating reversed:  
- start with string  
- end with parse tree

$$\begin{aligned} R &\rightarrow S \mid T \\ S &\rightarrow \mathbf{aSb} \mid ab \\ T &\rightarrow \mathbf{aTbb} \mid abb \end{aligned}$$

- **PDA**: can non-deterministically “try all rules” (abandoning failed attempts);  
- **DPDA**: must choose one rule at each step!

Should use *S* rule

$$aa\underline{abb}b \rightsquigarrow aa\underline{S}bb$$

Should use *T* rule

$$aa\underline{abb}bbb \rightsquigarrow aa\underline{T}bbbb$$

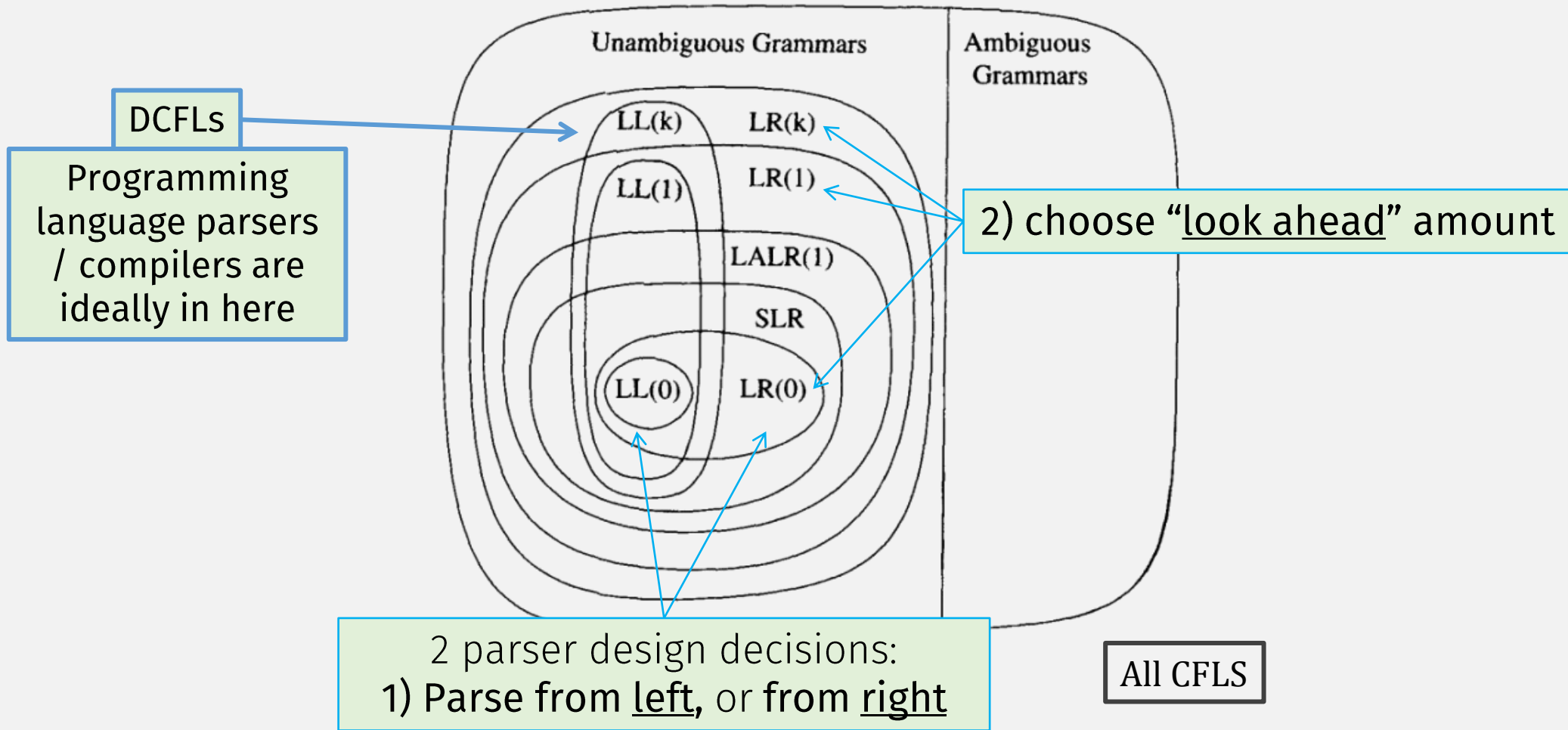
aaa

When parsing reaches this position, does it know which rule, *S* or *T*?

To choose “correct” rule, need to “look ahead” at rest of the input!

PDAs recognize CFLs, but **DPDAs** only recognize **DCFLs!** (a subset of CFLs)

# Subclasses of CFLs



# LL parsing

- **L** = left-to-right
- **L** = leftmost derivation

Game: "You're the Parser":  
Guess which rule applies?

**1**  $S \rightarrow$  if  $E$  then  $S$  else  $S$

**2**  $S \rightarrow$  begin  $S$   $L$

**3**  $S \rightarrow$  print  $E$

**4**  $L \rightarrow$  end

**5**  $L \rightarrow$  ;  $S$   $L$

**6**  $E \rightarrow$  num = num

if 2 = 3 begin print 1; print 2; end else print 0





# LL parsing

- L = left-to-right
- L = leftmost derivation

**1**  $S \rightarrow \text{if } E \text{ then } S \text{ else } S$

**2**  $S \rightarrow \text{begin } S L$

**3**  $S \rightarrow \text{print } E$

**4**  $L \rightarrow \text{end}$

**5**  $L \rightarrow ; S L$

**6**  $E \rightarrow \text{num} = \text{num}$

if 2 ← = 3 begin print 1; print 2; end else print 0



# LL parsing

- **L** = left-to-right
- **L** = leftmost derivation

**1**  $S \rightarrow \text{if } E \text{ then } S \text{ else } S$

**2**  $S \rightarrow \text{begin } S L$

**3**  $S \rightarrow \text{print } E$

**4**  $L \rightarrow \text{end}$

**5**  $L \rightarrow ; S L$

**6**  $E \rightarrow \text{num} = \text{num}$

if 2 = 3 begin print 1; print 2; end else print 0



# LL parsing

- L = left-to-right
- L = leftmost derivation

**1**  $S \rightarrow \text{if } E \text{ then } S \text{ else } S$

**2**  $S \rightarrow \text{begin } S L$

**3**  $S \rightarrow \text{print } E$

**4**  $L \rightarrow \text{end}$

**5**  $L \rightarrow ; S L$

**6**  $E \rightarrow \text{num} = \text{num}$

`if 2 = 3 begin print 1; print 2; end else print 0`

“Prefix” languages (Scheme/Lisp) are easily parsed with LL parsers (zero lookahead)

# LR parsing

- **L** = left-to-right

- **R** = rightmost derivation

1  $S \rightarrow S ; S$

2  $S \rightarrow \text{id} := E$

3  $S \rightarrow \text{print} ( L )$

4  $E \rightarrow \text{id}$

5  $E \rightarrow \text{num}$

6  $E \rightarrow E + E$

a := 7 ;



b := c + ( d := 5 + 6 , d )

When parse is here, can't determine whether it's an assign (:=) or addition (+)

Need to save input (lookahead) to some memory, like a **stack**! this is a job for a (D)PDA!

# LR parsing

- **L** = left-to-right
- **R** = rightmost derivation

$S \rightarrow S ; S$

$E \rightarrow \text{id}$

$S \rightarrow \text{id} := E$

$E \rightarrow \text{num}$

$S \rightarrow \text{print} ( L )$

$E \rightarrow E + E$

a := 7 ;

b := c + ( d := 5 + 6 , d )

Stack	Input	Action
1	a := 7 ; b := c + ( d := 5 + 6 , d ) \$	shift = "push"

push

State name

c

# LR parsing

- **L** = left-to-right
- **R** = rightmost derivation

$$S \rightarrow S ; S$$

$$E \rightarrow \text{id}$$

$$S \rightarrow \text{id} := E$$

$$E \rightarrow \text{num}$$

$$S \rightarrow \text{print} ( L )$$

$$E \rightarrow E + E$$

<i>Stack</i>	<i>Input</i>	<i>Action</i>
1	a := 7 ; b := c + ( d := 5 + 6 , d ) \$	<i>shift</i>
1 id <sub>4</sub>	:= 7 ; b := c + ( d := 5 + 6 , d ) \$	<i>shift</i>
1 id <sub>4</sub> := 6	7 ; b := c + ( d := 5 + 6 , d ) \$	<i>shift</i>



# LR parsing

- **L** = left-to-right
- **R** = rightmost derivation

$$\begin{array}{ll}
 S \rightarrow S ; S & E \rightarrow \text{id} \\
 S \rightarrow \text{id} := E & E \rightarrow \text{num} \\
 S \rightarrow \text{print} ( L ) & E \rightarrow E + E
 \end{array}$$

<i>Stack</i>	<i>Input</i>	<i>Action</i>
1	a := 7 ; b := c + ( d := 5 + 6 , d ) \$	<i>shift</i>
1 id <sub>4</sub>	:= 7 ; b := c + ( d := 5 + 6 , d ) \$	<i>shift</i>
1 id <sub>4</sub> :=6	7 ; b := c + ( d := 5 + 6 , d ) \$	<i>shift</i>
1 id <sub>4</sub> :=6 num <sub>10</sub>	; b := c + ( d := 5 + 6 , d ) \$	<i>reduce E → num</i>



# LR parsing

- L = left-to-right

- R = rightmost derivation

1  $S \rightarrow S ; S$

4  $E \rightarrow id$

2  $S \rightarrow id := E$

5  $E \rightarrow num$

3  $S \rightarrow print ( L )$

6  $E \rightarrow E + E$

Stack	Input	Action
1	a := 7 ; b := c + ( d := 5 + 6 , d ) \$	shift
1 id <sub>4</sub>	:= c + ( d := 5 + 6 , d ) \$	shift
1 id <sub>4</sub> :=6	:= c + ( d := 5 + 6 , d ) \$	shift
1 id <sub>4</sub> :=6 num <sub>10</sub>	; b := c + ( d := 5 + 6 , d ) \$	reduce $E \rightarrow num$

Can determine (rightmost) rule



# LR parsing

- L = left-to-right

- R = rightmost derivation

1  $S \rightarrow S ; S$

4  $E \rightarrow \text{id}$

2  $S \rightarrow \text{id} := E$

5  $E \rightarrow \text{num}$

3  $S \rightarrow \text{print} ( L )$

6  $E \rightarrow E + E$

Stack	Input	Action
1	a := 7 ; b := c + ( d := 5 + 6 , d ) \$	shift
1 id <sub>4</sub>	:= 7 ; b := c + ( d := 5 + 6 , d ) \$	shift
1 id <sub>4</sub> := <sub>6</sub>	= c + ( d := 5 + 6 , d ) \$	shift
1 id <sub>4</sub> := <sub>6</sub> num <sub>10</sub>	= c + ( d := 5 + 6 , d ) \$	reduce $E \rightarrow \text{num}$
1 id <sub>4</sub> := <sub>6</sub> E <sub>11</sub>	; b := c + ( d := 5 + 6 , d ) \$	reduce $S \rightarrow \text{id} := E$

Can determine (rightmost) rule



# LR parsing

- **L** = left-to-right
- **R** = rightmost derivation

$$S \rightarrow S ; S \qquad E \rightarrow id$$

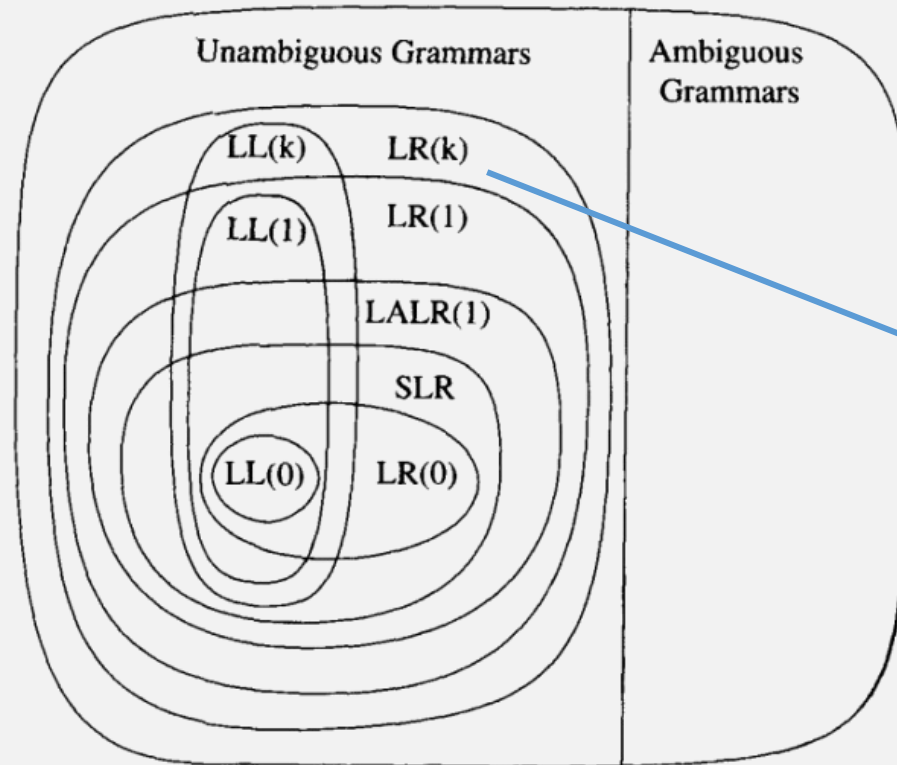
$$S \rightarrow id := E \qquad E \rightarrow num$$

$$S \rightarrow print ( L ) \qquad E \rightarrow E + E$$

Stack	Input	Action
1	a := 7 ; b := c + ( d := 5 + 6 , d ) \$	shift
1 id <sub>4</sub>	:= 7 ; b := c + ( d := 5 + 6 , d ) \$	shift
1 id <sub>4</sub> :=6	7 ; b := c + ( d := 5 + 6 , d ) \$	shift
1 id <sub>4</sub> :=6 num <sub>10</sub>	; b := c + ( d := 5 + 6 , d ) \$	reduce $E \rightarrow num$
1 id <sub>4</sub> :=6 E <sub>11</sub>	; b := c + ( d := 5 + 6 , d ) \$	reduce $S \rightarrow id := E$
1 S <sub>2</sub>	; b := c + ( d := 5 + 6 , d ) \$	shift

LR Parsers also called "Shift-Reduce" Parsers

# To learn more, take a Compilers Class!



A program (string of chars)



Program "words"



Abstract Syntax tree (AST)



This phase needs computation that goes beyond CFLs

# **In-class quiz 3/20**

See gradescope