# Deterministic CFLs, PDAs, and Parsing

Monday, February 28, 2022


(AN UNMATCHED LEFT PARENTHESIS CREATES AN UNRESOLVED TENSION THAT WILL STAY WITH YOU ALL DAY.

# Announcements

- HW 4 in

- HW 5 out
  - Due Sun March 6 11:59pm
  - Problems about PDAs

- Upcoming: Spring Break is week of March 14

# *Previously:* CFLs, CFGs, and Parse Trees

**Generating** strings:
- <u>Start</u> with *start variable,*
- <u>Repeatedly</u> *apply rules* to get a string (and **parse tree**)

$A \to 0A1$
$A \to B$
$B \to \#$



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

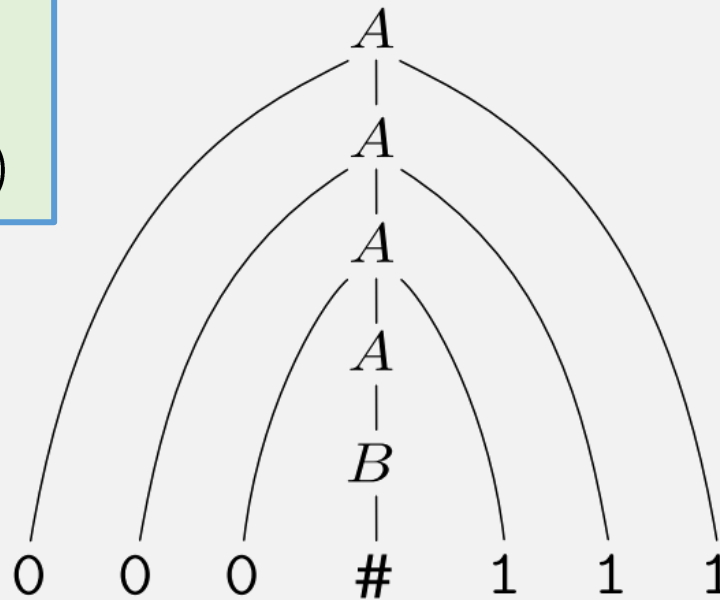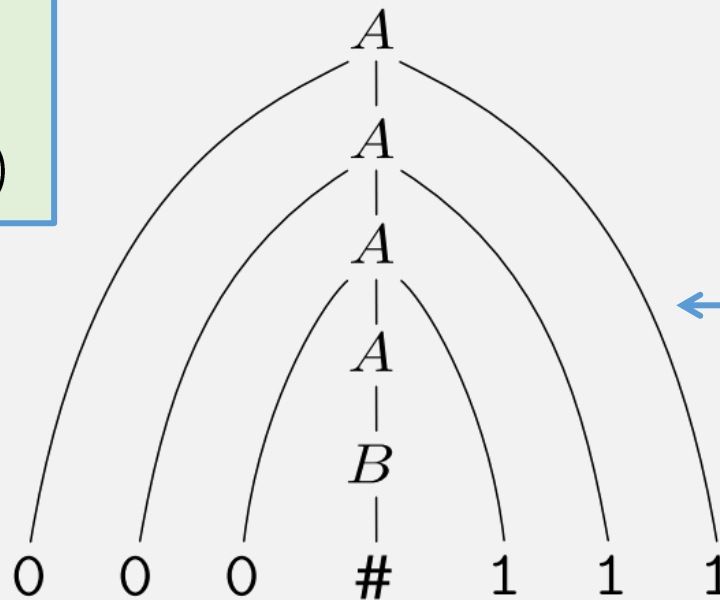# $\mathscr{Today:}$ Generating vs Parsing

Generating strings:
- Start with *start variable*,
- Repeatedly *apply rules* to get a string (and parse tree)

$$A \rightarrow 0A1$$
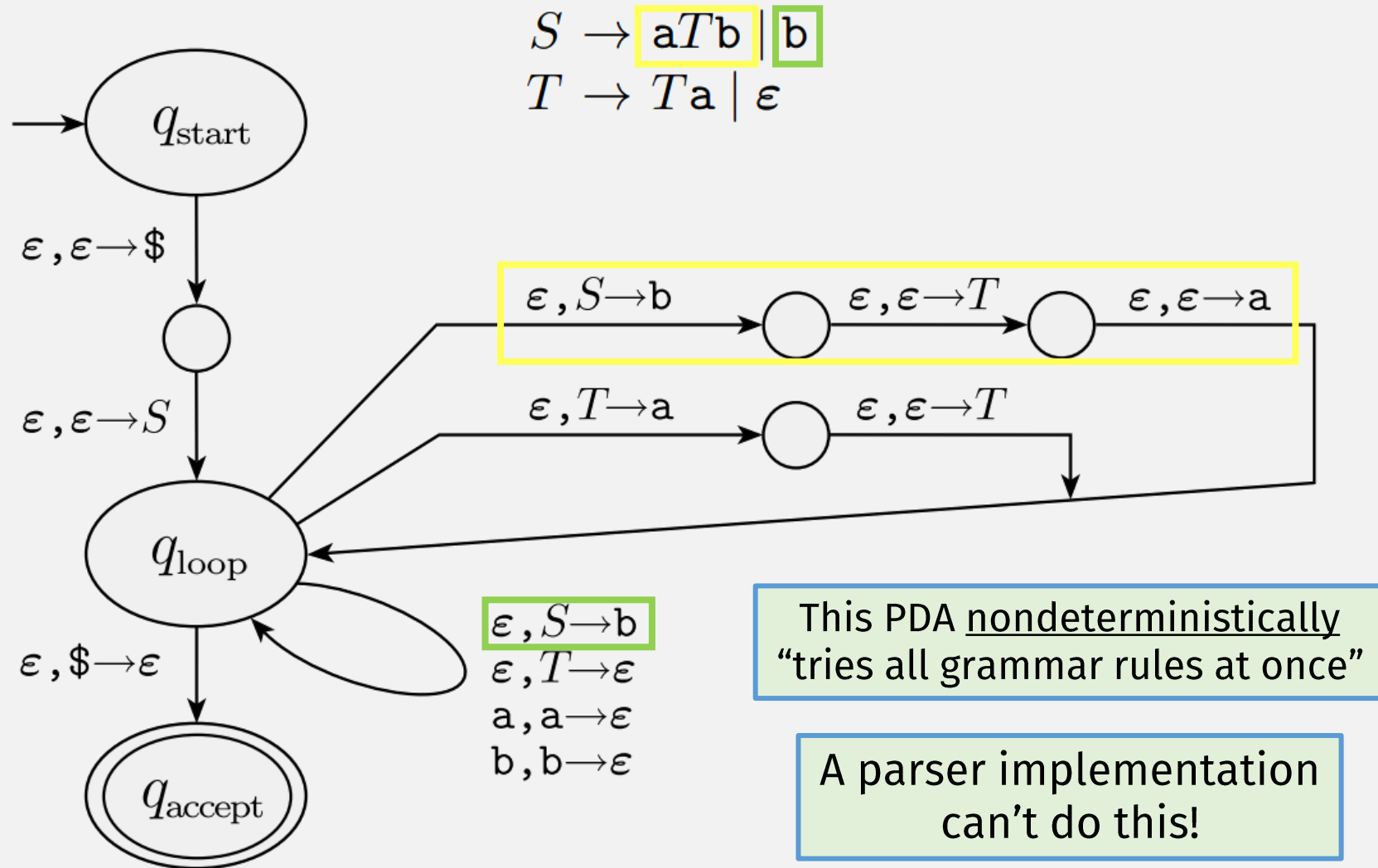$$A \rightarrow B$$
$$B \rightarrow \#$$

In practice, the opposite is more interesting: start with a string, 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 is more important than **generating** one
  - E.g., a **compiler's 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 <u>non-deterministic</u> (like NFAs)

- A compiler's parsing algorithm must be <u>deterministic</u>

- <u>So</u>: to model parsers, we need a **Deterministic PDA** (DPDA)

# *Last time:* (Nondeterministic) PDA

$$S \rightarrow \boxed{\mathbf{a}T\mathbf{b}} \mid \boxed{\mathbf{b}}$$
$$T \rightarrow T\mathbf{a} \mid \varepsilon$$



$\varepsilon, \varepsilon \rightarrow \$$

$\varepsilon, S \rightarrow \mathbf{b}$     $\varepsilon, \varepsilon \rightarrow T$     $\varepsilon, \varepsilon \rightarrow \mathbf{a}$

$\varepsilon, T \rightarrow \mathbf{a}$     $\varepsilon, \varepsilon \rightarrow T$

$\varepsilon, \varepsilon \rightarrow S$

$q_{\mathrm{loop}}$

$\varepsilon, S \rightarrow \mathbf{b}$
$\varepsilon, T \rightarrow \varepsilon$
$\mathbf{a}, \mathbf{a} \rightarrow \varepsilon$
$\mathbf{b}, \mathbf{b} \rightarrow \varepsilon$

$\varepsilon, \$ \rightarrow \varepsilon$

$q_{\mathrm{accept}}$

This PDA <u>nondeterministically</u> "tries all grammar rules at once"

A parser implementation can't do this!

# 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. $\delta: Q \times \Sigma_\varepsilon \times \Gamma_\varepsilon \longrightarrow (Q \times \Gamma_\varepsilon) \cup \{\emptyset\}$ is the transition function
5. $q_0 \in Q$ is the start state, and
6. $F \subseteq Q$ is the set of accept states.

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. $\delta: Q \times \Sigma_\varepsilon \times \Gamma_\varepsilon \longrightarrow \mathcal{P}(Q \times \Gamma_\varepsilon)$
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**)

This must take into account $\varepsilon$ reads or stack ops!
E.g., if $\delta(q, a, X)$ is valid, then $\delta(q, \varepsilon, X)$ must not be

# DPDAs are <u>Not</u> Equivalent to PDAs!

$$R \rightarrow S \mid T$$
$$S \rightarrow \mathbf{a}S\mathbf{b} \mid \mathbf{ab}$$
$$T \rightarrow \mathbf{a}T\mathbf{bb} \mid \mathbf{abb}$$

A PDA non-deterministically "tries all rules" (abandoning failed attempts) but a DPDA must decide on one rule at each step!

Should use $S$ rule

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

$$\mathbf{aa\underline{a}\underline{bb}b} \rightarrowtail \mathbf{a\underline{a}S\underline{b}b}$$

When parsing reaches this input position, which rule should it use, $S$ or $T$?

Should use $T$ rule

$$\mathbf{aa\underline{ab}bbbb} \rightarrowtail \mathbf{a\underline{a}Tbbb}$$
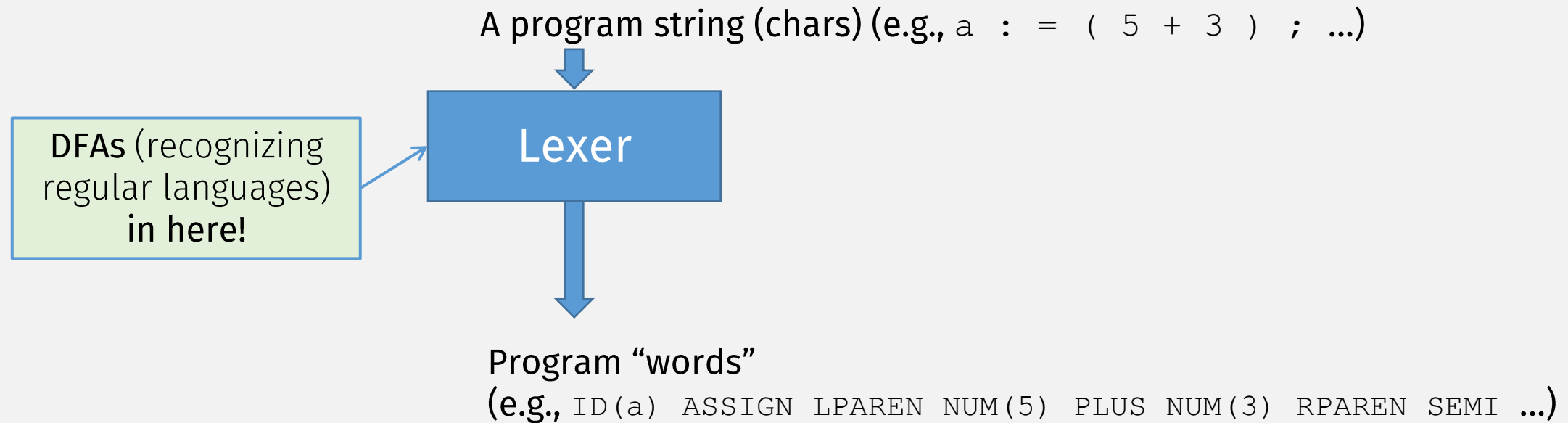
Don't know which rule to use because we can't see rest of the input!

PDAs recognize CFLs, but <u>DPDAs only recognize DCFLs</u>! (a <u>subset</u> of CFLs)

# Subclasses of CFLs

DCFLs

Programming language parsers / compilers are ideally in here

Unambiguous Grammars

Ambiguous Grammars

LL(k)  LR(k)

LL(1)  LR(1)

LALR(1)

SLR

LL(0)  LR(0)

All CFLS

# Compiler Stages

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

DFAs (recognizing regular languages) **in here!**

Lexer

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

# A Lexer Implementation

```
%{
/* C Declarations: */
#include "tokens.h"    /* definitions of IF, ID, NUM, ... */
#include "errormsg.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**)!

A "`lex`" tool translates this to a (C program) implementation of a lexer

# Compiler Stages

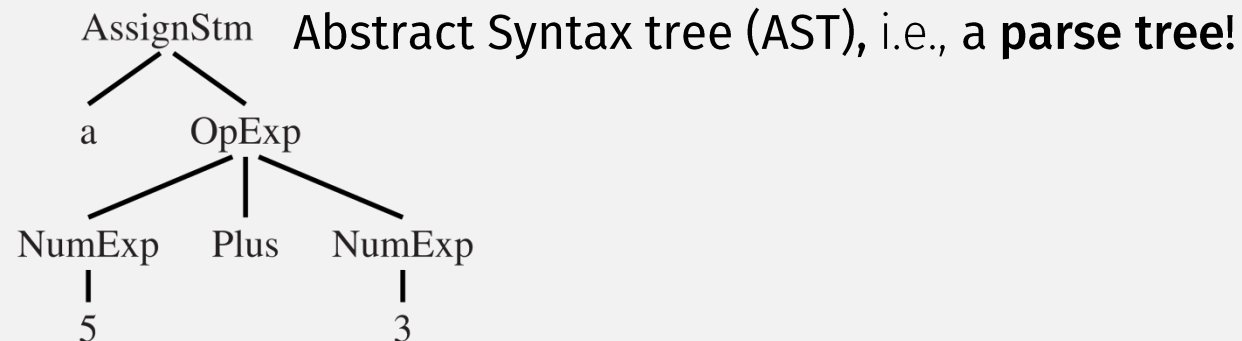A program (chars) (e.g., `a : = ( 5 + 3 ) ; …`)

↓

**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!**

↓

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

```
        AssignStm
        /       \
       a        OpExp
              /   |   \
        NumExp  Plus  NumExp
           |            |
           5            3
```

# 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!

A "`yacc`" tool translates this to a (C program) implementation of a parser

# Parsing

$$R \rightarrow S \mid T$$
$$S \rightarrow aSb \mid ab$$
$$T \rightarrow aTbb \mid abb$$

$$\mathbf{aa\underline{ab}bb} \rightarrowtail \mathbf{aa\underline{S}bb}$$

A parser must be able to <u>choose the one correct rule</u>, when reading input left-to-right

$$\mathbf{aa\underline{ab}bbbb} \rightarrowtail \mathbf{aa\underline{aT}bbbb}$$

# LL parsing

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

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

**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
```

118

# LL parsing

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

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

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

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

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

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

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

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

# LL parsing

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

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

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

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

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

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

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

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

# LL parsing

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

$1\ S \to \text{if } E \text{ then } S \text{ else } S$

$2\ S \to \text{begin } S\ L$

$3\ S \to \boxed{\text{print } E}$

$4\ L \to \text{end}$

$5\ L \to ;\ S\ L$

$6\ E \to \text{num} = \text{num}$

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

"Prefix" languages (like Scheme/Lisp) are easily parsed with LL parsers

# LR parsing

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

$$1 \quad S \rightarrow S\ ;\ S \qquad\qquad 4 \quad E \rightarrow \text{id}$$

$$2 \quad S \rightarrow \text{id} := E \qquad 5 \quad E \rightarrow \text{num}$$

$$3 \quad S \rightarrow \text{print} \ (\ L\ ) \qquad 6 \quad E \rightarrow E\ +\ E$$

```
a := 7;
b := c + (d := 5 + 6, d)
```

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

Need to <u>save</u> input to some temporary memory, like a **stack**: this is a job for a (D)PDA!!

| Stack | Input | Action |
|---|---|---|
| 1 | a := 7 ; b := c + ( d := 5 + 6 , d ) $ | shift *"push"* |
| 1 $\text{id}_4$ | := 7 ; b := c + ( d := 5 + 6 , d ) $ | shift |
| 1 $\text{id}_4$ := $_6$ | 7 ; b := c + ( d := 5 + 6 , d ) $ | shift |
| 1 $\text{id}_4$ := $_6$ $\text{num}_{10}$ | ; b := c + ( d := 5 + 6 , d ) $ | reduce $E \rightarrow \text{num}$ |
| 1 $\text{id}_4$ := $_6$ $E_{11}$ | ; b := c + ( d := 5 + 6 , d ) $ | reduce $S \rightarrow \text{id} := E$ |
| 1 $S_2$ | ; b := c + ( d := 5 + 6 , d ) $ | shift |

push

State name

123

# LR parsing

$$S \to S \; ; \; S \qquad E \to \text{id}$$

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

$$S \to \text{id} := E \qquad E \to \text{num}$$

$$S \to \text{print} \; ( \; L \; ) \qquad E \to E \; + \; E$$

| Stack | Input | Action |
|---|---|---|
| 1 | a := 7 ; b := c + ( d := 5 + 6 , d ) $ | *shift* |
| 1 $\text{id}_4$ | := 7 ; b := c + ( d := 5 + 6 , d ) $ | *shift* |
| 1 $\text{id}_4$ :=$_6$ | 7 ; b := c + ( d := 5 + 6 , d ) $ | *shift* |
| 1 $\text{id}_4$ :=$_6$ $\text{num}_{10}$ | ; b := c + ( d := 5 + 6 , d ) $ | *reduce* $E \to \text{num}$ |
| 1 $\text{id}_4$ :=$_6$ $E_{11}$ | ; b := c + ( d := 5 + 6 , d ) $ | *reduce* $S \to \text{id} := E$ |
| 1 $S_2$ | ; b := c + ( d := 5 + 6 , d ) $ | *shift* |

124

# LR parsing

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

$$S \rightarrow S \; ; \; S \qquad E \rightarrow \mathrm{id}$$
$$S \rightarrow \mathrm{id} := E \qquad E \rightarrow \mathrm{num}$$
$$S \rightarrow \mathrm{print} \; ( \; L \; ) \qquad E \rightarrow E + E$$

| Stack | Input | Action |
|---|---|---|
| 1 | a := 7 ; b := c + ( d := 5 + 6 , d ) \$ | *shift* |
| $1 \; \mathrm{id}_4$ | := 7 ; b := c + ( d := 5 + 6 , d ) \$ | *shift* |
| $1 \; \mathrm{id}_4 :=_6$ | 7 ; b := c + ( d := 5 + 6 , d ) \$ | *shift* |
| $1 \; \mathrm{id}_4 :=_6 \mathrm{num}_{10}$ | ; b := c + ( d := 5 + 6 , d ) \$ | *reduce* $E \rightarrow \mathrm{num}$ |
| $1 \; \mathrm{id}_4 :=_6 E_{11}$ | ; b := c + ( d := 5 + 6 , d ) \$ | *reduce* $S \rightarrow \mathrm{id}:=E$ |
| $1 \; S_2$ | ; b := c + ( d := 5 + 6 , d ) \$ | *shift* |

# LR parsing

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

$$1 \quad S \rightarrow S \; ; \; S \qquad 4 \quad E \rightarrow \mathrm{id}$$

$$2 \quad S \rightarrow \mathrm{id} := E \qquad 5 \quad E \rightarrow \mathrm{num}$$

$$3 \quad S \rightarrow \mathrm{print} \; ( \; L \; ) \qquad 6 \quad E \rightarrow E + E$$

| Stack | Input | Action |
|---|---|---|
| 1 | a := 7 ; b := c + ( d := 5 + 6 , d ) $ | *shift* |
| 1 id$_4$ | := 7 ; b := c + ( d := 5 + 6 , d ) $ | *shift* |
| 1 id$_4$ :=$_6$ | 7 ; b := c + ( d := 5 + 6 , d ) $ | *shift* |
| 1 id$_4$ :=$_6$ num$_{10}$ | ; b := c + ( d := 5 + 6 , d ) $ | *reduce* $E \rightarrow$ num |
| 1 id$_4$ :=$_6$ $E_{11}$ | ; b := c + ( d := 5 + 6 , d ) $ | *reduce* $S \rightarrow$ id:=$E$ |
| 1 $S_2$ | ; b := c + ( d := 5 + 6 , d ) $ | *shift* |

Can determine (rightmost) rule

126

# LR parsing

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

$$1 \quad S \rightarrow S \; ; \; S \qquad 4 \quad E \rightarrow \mathrm{id}$$
$$2 \quad \boxed{S \rightarrow \mathrm{id} := E} \qquad 5 \quad E \rightarrow \mathrm{num}$$
$$3 \quad S \rightarrow \mathrm{print} \; ( \; L \; ) \qquad 6 \quad E \rightarrow E + E$$

| Stack | Input | Action |
|---|---|---|
| 1 | a := 7 ; b := c + ( d := 5 + 6 , d ) $ | *shift* |
| 1 id$_4$ | := 7 ; b := c + ( d := 5 + 6 , d ) $ | *shift* |
| 1 id$_4$ :=$_6$ | = c + ( d := 5 + 6 , d ) $ | *shift* |
| 1 id$_4$ :=$_6$ num$_{10}$ | = c + ( d := 5 + 6 , d ) $ | *reduce* $E \rightarrow$ num |
| 1 id$_4$ :=$_6$ $E_{11}$ | ; b := c + ( d := 5 + 6 , d ) $ | *reduce* $S \rightarrow$ id:=$E$ |
| 1 $S_2$ | b := c + ( d := 5 + 6 , d ) $ | *shift* |

Can determine (rightmost) rule

127

# LR parsing

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

$$S \rightarrow S \; ; \; S \qquad E \rightarrow \text{id}$$
$$S \rightarrow \text{id} := E \qquad E \rightarrow \text{num}$$
$$S \rightarrow \text{print} \; ( \; L \; ) \qquad E \rightarrow E \; + \; E$$

| Stack | Input | Action |
|---|---|---|
| 1 | a := 7 ; b := c + ( d := 5 + 6 , d ) $ | *shift* |
| $1 \; \text{id}_4$ | := 7 ; b := c + ( d := 5 + 6 , d ) $ | *shift* |
| $1 \; \text{id}_4 := _6$ | 7 ; b := c + ( d := 5 + 6 , d ) $ | *shift* |
| $1 \; \text{id}_4 := _6 \text{num}_{10}$ | ; b := c + ( d := 5 + 6 , d ) $ | *reduce $E \rightarrow$ num* |
| $1 \; \text{id}_4 := _6 E_{11}$ | ; b := c + ( d := 5 + 6 , d ) $ | *reduce $S \rightarrow$ id:=E* |
| $1 \; S_2$ | ; b := c + ( d := 5 + 6 , d ) $ | *shift* |

128

# To learn more, take a Compilers Class!

Unambiguous Grammars

LL(k)    LR(k)

LL(1)    LR(1)

LALR(1)

SLR

LL(0)    LR(0)

Ambiguous Grammars

A program (string of chars)

**Lexer**
(DFAs / NFAs)

Program "words"

**Parser**
(DPDAs)

Abstract Syntax tree (AST)

**???**

This phase needs computation that goes beyond CFLs

# Non-CFLs

# *Flashback:* Pumping Lemma for Regular Langs

- The Pumping Lemma <u>describes how strings</u> **repeat**

- Regular language strings can (only) repeat using Kleene pattern
    - But the <u>substrings are independent</u>!

- A non-regular language:

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

Kleene star can't express this pattern:
2nd part <u>depends</u> on (length of) 1st part

- Q: **How do CFLs repeat?**

Repeating pattern

After repeat

Before repeat

Independent

$y$

$q_9$

$x$

$z$

$q_{13}$

$q_1$

# Repetition and Dependency in CFLs

Parts before/after repetition point are linked

Repetition

$$A \to 0A1$$

$$A \to B$$

$$B \to \#$$

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

repetition

$A$

$A$

$A$

$A$

$B$

0   0   0   #   1   1   1

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

# How Do Strings in CFLs Repeat?

- Strings in regular languages repeat states

- Strings in CFLs <u>repeat subtrees</u> in the parse tree

One repeated subtree means that it can be repeated <u>any</u> number of times

Linked parts

# Pumping Lemma for CFLS

**Pumping lemma for context-free languages**  If $A$ is a context-free language, then there is a number $p$ (the pumping length) where, if $s$ is any string in $A$ of length at least $p$, then $s$ may be divided into five pieces $s = uvxyz$ satisfying the conditions

> Now there are <u>two</u> pumpable parts.
> But they must be <u>pumped together</u>!

**1.** for each $i \geq 0$, $uv^i xy^i z \in A,$

**2.** $|vy| > 0$, and

**3.** $|vxy| \leq p.$

**Pumping lemma**  If $A$ is a regular [language then there is a number $p$ (the] pumping length) where if $s$ is any stri[ng] [in $A$ of length at least $p$, then] $s$ may be divided into three pieces, $s = xyz$, sat[isfying the conditions]

**1.** for each $i \geq 0$, $xy^i z \in A,$

**2.** $|y| > 0$, and

**3.** $|xy| \leq p.$



> Two pumpable parts, pumped together

# Non CFL example: $D = \{ww|\ w \in \{0,1\}^*\}$

Previous: $D$ is <u>nonregular</u>: unpumpable <u>counterexample</u> $s$: $0^p10^p1$

<u>Now</u>: **this** $s$ **can** be pumped according to <u>CFL pumping lemma</u>:

$$
\overbrace{\underbrace{000\cdots000}_{u}\ \underbrace{0}_{v}}^{0^p1}\ \underbrace{1}_{x}\ \overbrace{\underbrace{0}_{y}\ \underbrace{000\cdots0001}_{z}}^{0^p1}
$$

Pumping *v* and *y* (together) produces string still in $D$

- CFL Pumping Lemma conditions:
  - ☑ **1.** for each $i \geq 0,\ uv^ixy^iz \in A,$
  - ☑ **2.** $|vy| > 0,$ and
  - ☑ **3.** $|vxy| \leq p.$

**This doesn't prove that the language is a CFL!**
It only means that <u>this attempt</u> to prove that the language is not a CFL failed.

136

# Non CFL example: $D = \{ww|\ w \in \{0,1\}^*\}$

- Need another counterexample string $s$:

If *vyx* is contained in first or second half, then
any pumping will break the match ✖

$$0^p 1^p 0^p 1^p$$

So *vyx* must straddle the middle ✖
But any pumping still breaks the match because order is wrong

- CFL Pumping Lemma conditions:

  **1.** for each $i \geq 0$, $uv^i xy^i z \in A$,

  **2.** $|vy| > 0$, and

  **3.** $|vxy| \leq p.$

Now we have proven that
**this language is not a CFL!**

# CFL Pumping Lemma is Too Restrictive?

**???**

**Pumping lemma for context-free languages**   If $A$ is a context-free language, then there is a number $p$ (the pumping length) where, if $s$ is any string in $A$ of length at least $p$, then $s$ may be divided into five pieces $s = uvxyz$ satisfying the conditions

1. for each $i \geq 0$, $uv^i x y^i z \in A$,
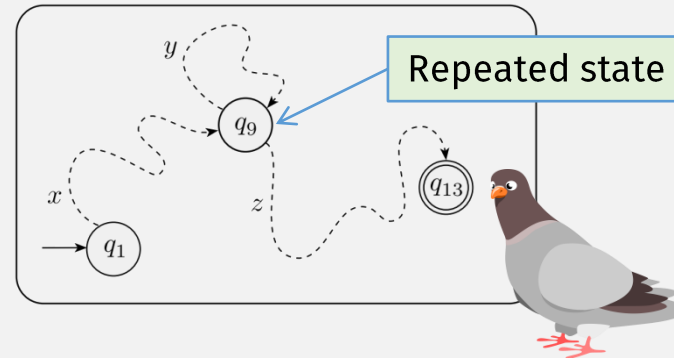2. $|vy| > 0$, and
3. $|vxy| \leq p$.

# *Review:* Regular Language Pumping Lemma

- The <u>pumping length</u> $p$ for a language $L$ is …

  … the # of states in that language's NFA!

**Pumping lemma**     If $A$ is a regular language, then there is a number $p$ (the pumping length) where if $s$ is any string in $A$ of length at least $p$, then $s$ may be divided into three pieces, $s = xyz$, satisfying the following conditions:

   **1.** for each $i \geq 0$, $xy^i z \in A$,
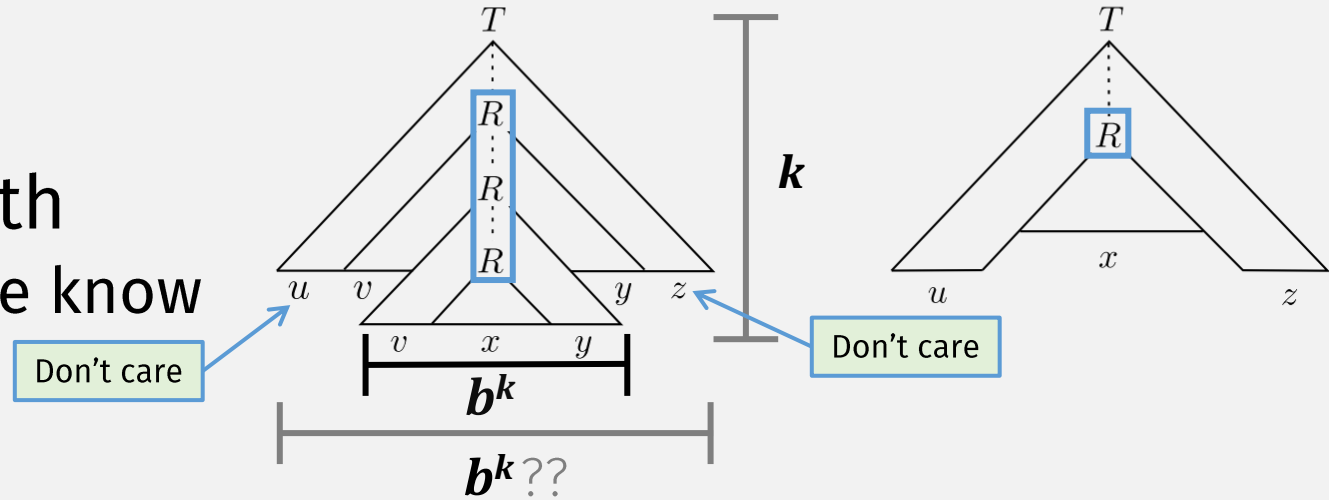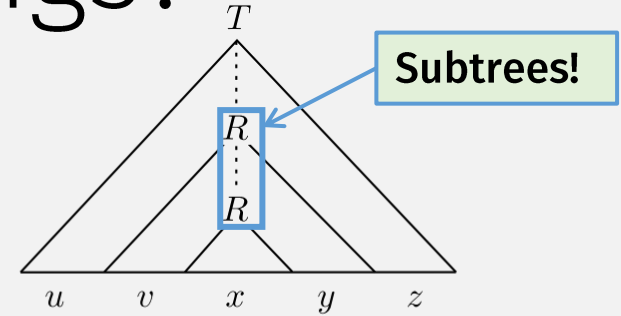   **2.** $|y| > 0$, and
   **3.** $|xy| \leq p$.

- If string length > # of states, then <u>some state must repeat</u>

Repeated state

$q_9$

$q_{13}$

$x$

$q_1$

$y$

$z$

- If a state is <u>repeated once</u>, then it can <u>repeat multiple times</u>

# Repeating Pattern in CFL Strings?

- When are we <u>guaranteed</u> to have a repeated subtree?
  - When <u>height</u> of parse tree > # of rules!

- Let $k$ = # of rules and $b$ = longest rule RHS length
  - Then the length string where we know there's a repeat is $b^k$
  - I.e., pumping length = $b^k$ ???

**Subtrees!**

**Don't care**

**Don't care**

$b^k$

$b^k$ ??

$k$

**Pumping lemma for context-free languages** If $A$ is a context-free language, then there is a number $p$ (the pumping length) where, if $s$ is any string in $A$ of length at least $p$, then $s$ may be divided into five pieces $s = uvxyz$ satisfying the conditions

1. for each $i \geq 0$, $uv^i xy^i z \in A$,
2. $|vy| > 0$, and
3. $|vxy| \leq p$.

**Pumping Length could be too short!**

# A Pumpable Non-CFL?

**Pumping lemma for context-free languages**  If $A$ is a context-free language, then there is a number $p$ (the pumping length) where, if $s$ is any string in $A$ of length at least $p$, then $s$ may be divided into five pieces $s = uvxyz$ satisfying the conditions

1. for each $i \geq 0$, $uv^i xy^i z \in A$,
2. $|vy| > 0$, and
3. $|vxy| \leq p$.

- CFL Pumping Lemma says:
  - "*All CFLs are pumpable*"
  - So if we find a non-pumpable language … it's not a CFL!

- Pumping Lemma does <u>not</u> say:
  - "*All nonCFLs are not pumpable*"
  - (statement != it's inverse)
  - So Pumping Lemma might not be able to prove some non-CFLs!

Example:

$$L = \{ \mathsf{a}^i \mathsf{b}^j \mathsf{c}^k \mathsf{d}^l \mid i = 0 \text{ or } j = k = l \}$$

- For any counterexample, split into *uvxyz* where,
  - $v$ = first char
  - $z$ = remaining chars
  - $u = x = y = \varepsilon$
- If there are **a**s …
  - … it's pumpable bc # of $\mathsf{a}$s is arbitrary
- If there there are no **a**s
  - … it's pumpable bc # of other chars is arbitrary

**This language is pumpable … but not a CFL!**
(can't come up with a CFG)

# Ogden's Lemma (generalizes pumping lemma)

Ogden's lemma is: If $L$ is a CFL, then there is a constant $n$, such that if $z$ is any string of length at least $n$ in $L$, in which we select at least $n$ positions to be *distinguished*, then we can write $z = uvwxy$, such that:

1. $vwx$ has at most $n$ distinguished positions.

2. $vx$ has at least one distinguished position.

3. For all $i$, $uv^i wx^i y$ is in $L$.

> Says that every long enough <u>segment</u> must be pumpable

## Example:

$$L = \{\mathrm{a}^i \mathrm{b}^j \mathrm{c}^k \mathrm{d}^l \mid i = 0 \text{ or } j = k = l\}$$

> **This language is not a CFL because it doesn't satisfy Ogden's Lemma**

## Counterexample: $\mathrm{ab}^n \mathrm{c}^n \mathrm{d}^n$

- $n$ "distinguished" positions must include non-$\mathrm{a}$ character
  - Impossible to pump no matter which $n$ chars are chosen

# A Practical Non-CFL

- **XML**
  - ELEMENT → \<TAG\>CONTENT\</TAG\>
  - Where TAG is any string

- XML also looks like this <u>non-CFL</u>: $D = \{ww | \ w \in \{0,1\}^*\}$

- This means XML is <u>not context-free</u>!
  - <u>Note</u>: HTML *is* context-free because …
  - … there are only a finite number of tags,
  - so they can be embedded into a finite number of rules.

- <u>In practice</u>:
  - XML is <u>parsed</u> as a CFL, with a CFG
  - Then **matching tags checked in a 2$^{nd}$ pass with a more powerful machine** …

# Next Time: A More Powerful Machine ...

$M_1$ accepts its input if it is in language: $B = \{w\#w \mid w \in \{0,1\}^*\}$

$M_1 =$ "On input string $w$:

Infinite memory, initially starts with input

1. Zig-zag across the tape to corresponding positions on either side of the # symbol to check whether these positions contain the same symbol. If they do not, or if no # is found, *reject*. Cross off symbols as they are checked to keep track of which symbols correspond.

Can move to, and read/write from, <u>arbitrary</u> memory locations

# In-class quiz 2/28

See gradescope