

Deterministic CFLs, PDAs, and Parsing

Wednesday, October 6, 2021

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

Announcements

- Reminder: no class next Monday 10/11
- HW4 due Sunday 10/17 11:59pm
 - second Sunday from today

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

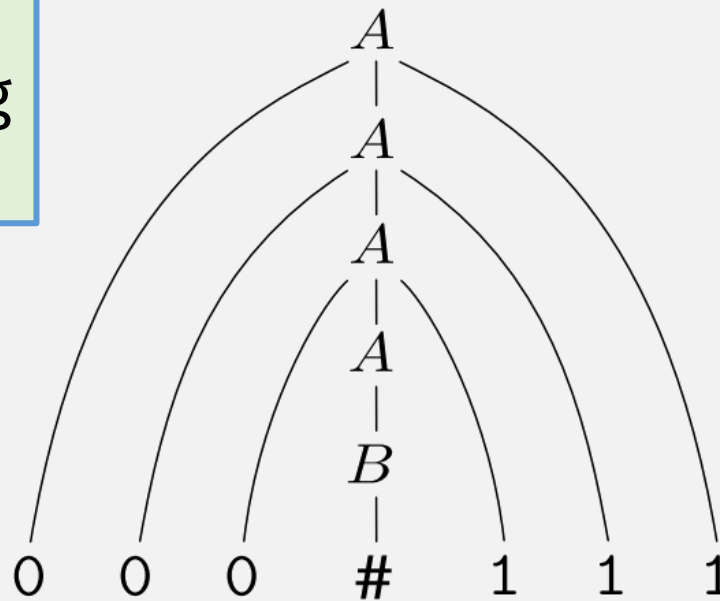
Previously: CFLs, CFGs, and Parse Trees

Generating strings:
start with start variable,
Apply rules to get a string
(and parse tree)

$$A \rightarrow 0A1$$

$$A \rightarrow B$$

$$B \rightarrow \#$$



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

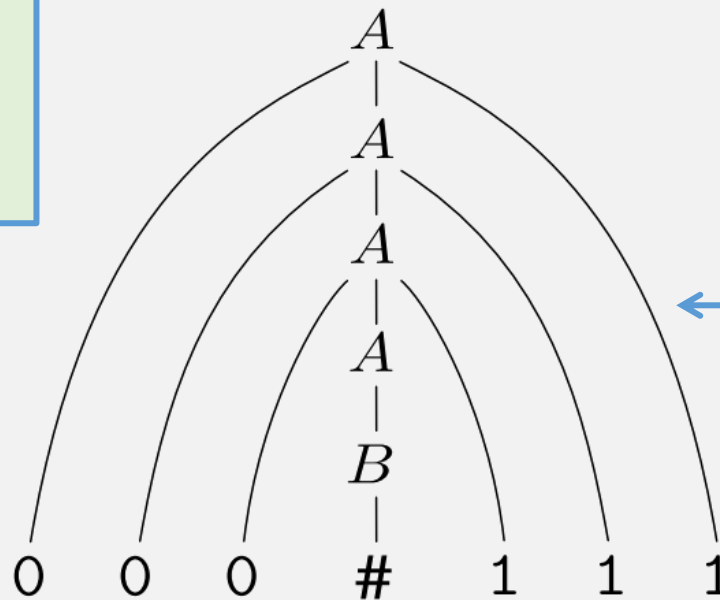
Today: Generating vs Parsing

Generating strings:
start with start variable,
then apply rules to get a
string and parse tree

$$A \rightarrow 0A1$$

$$A \rightarrow B$$

$$B \rightarrow \#$$



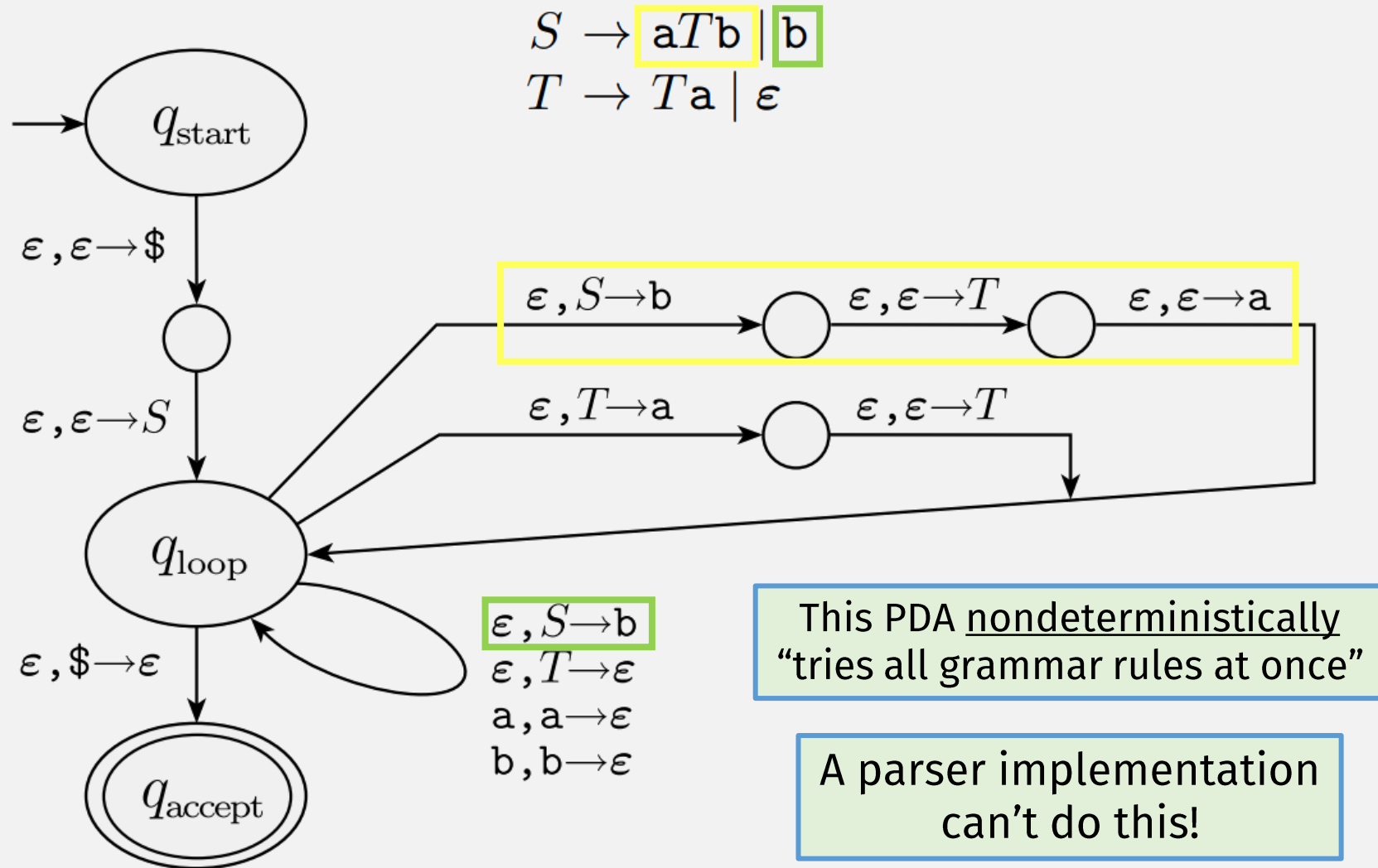
In practice, the
opposite is more interesting:
start with a string,
then **parse** it into parse tree

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

Generating vs Parsing

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

Last time: (Nondeterministic) PDA



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, Σ, Γ , and F are all finite sets, and

1. Q is the set of states,
2. Σ is the input alphabet,
3. Γ 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.

A *pushdown automaton* is a 6-tuple

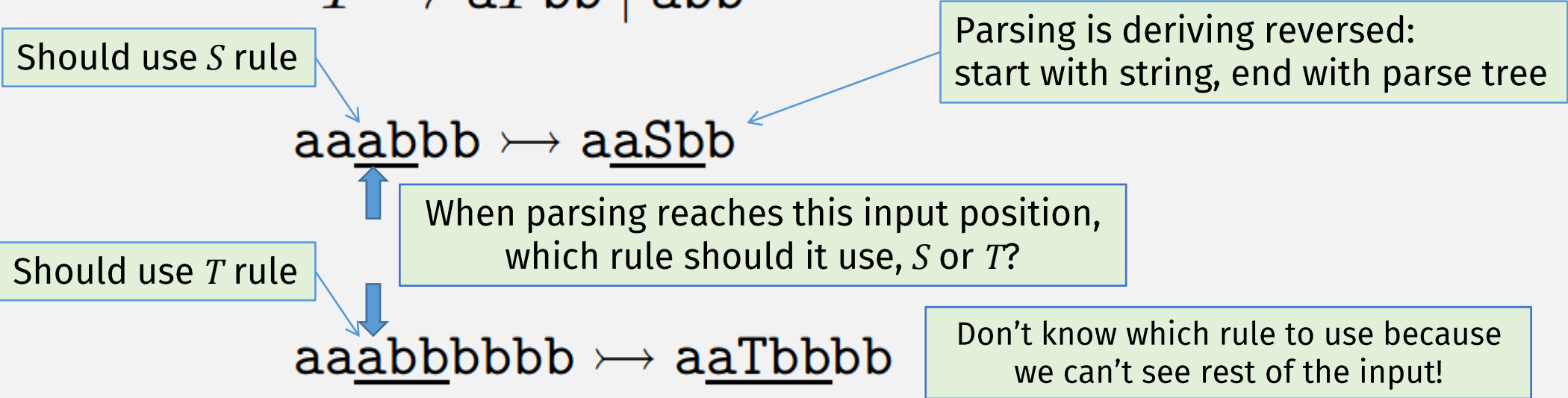
1. Q is the set of states,
2. Σ is the input alphabet,
3. Γ is the stack alphabet,
4. $\delta: Q \times \Sigma_\epsilon \times \Gamma_\epsilon \longrightarrow \mathcal{P}(Q \times \Gamma_\epsilon)$
5. $q_0 \in Q$ is the start state, and
6. $F \subseteq Q$ is the set of accept states.

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

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

DPDAs are Not Equivalent to PDAs!

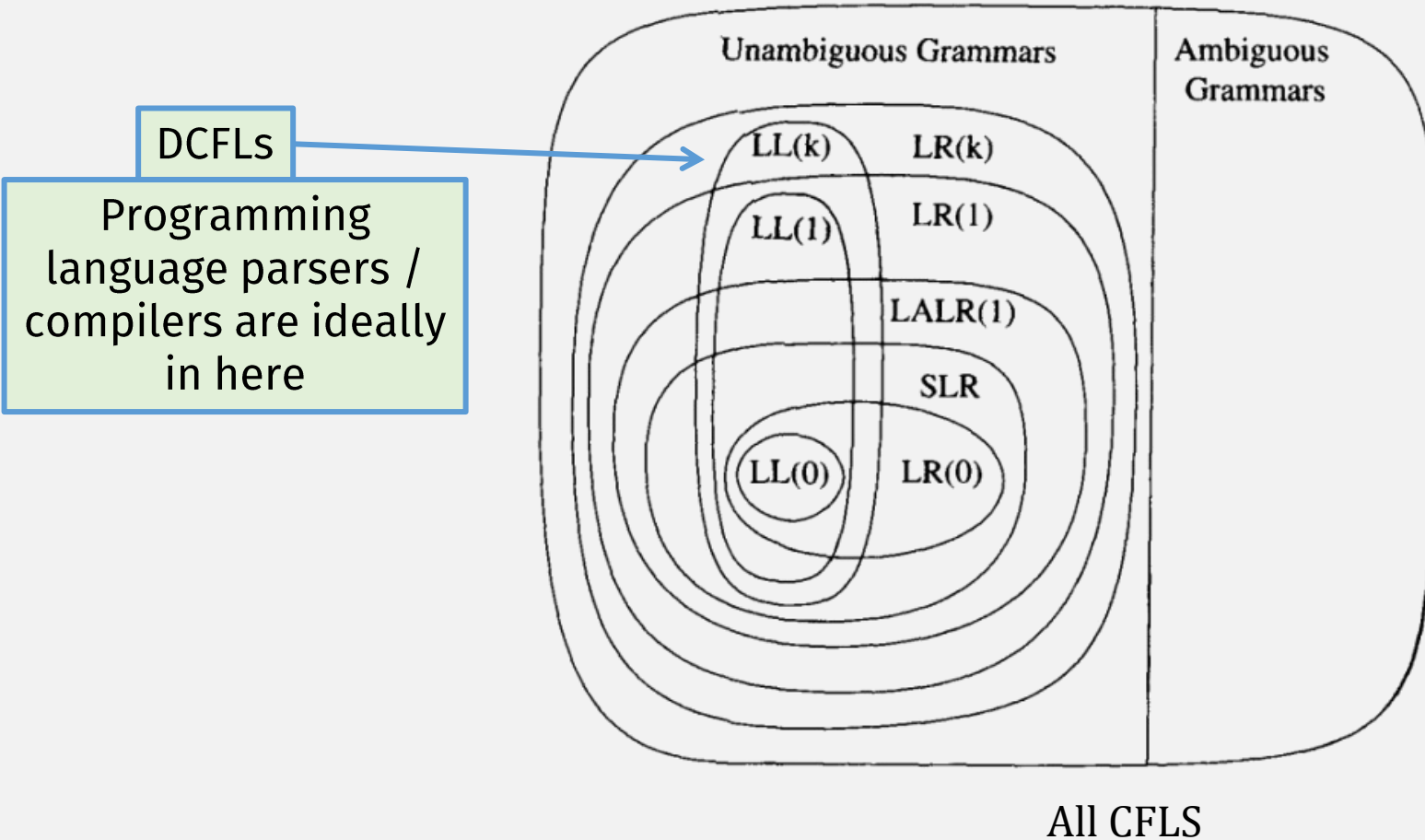
$$\begin{aligned} R &\rightarrow S \mid T \\ S &\rightarrow aSb \mid ab \\ T &\rightarrow aTbb \mid abb \end{aligned}$$



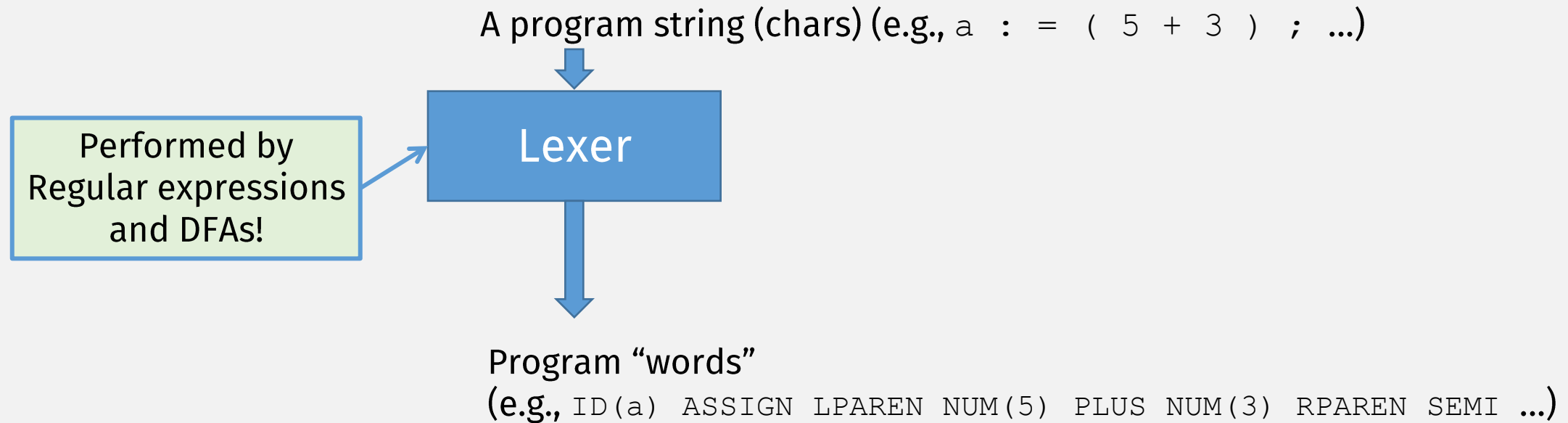
A PDA non-deterministically “tries all rules” (abandons failed attempts) but a DPDA gets only one try!

PDAs recognize CFLs, but a DPDA only recognizes 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");}
```

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

Performed by
Regular expressions
and DFAs!

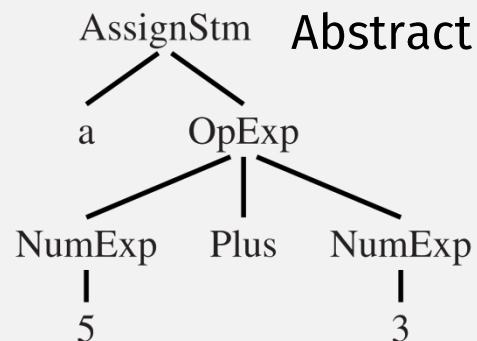
Program "words"

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

Parser

DCFLs and DPDAs

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!



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

Parsing

$$R \rightarrow S \mid T$$

$$S \rightarrow aSb \mid ab$$

$$T \rightarrow aTbb \mid abb$$

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

A parser must be able to choose one correct rule, when reading input left-to-right

$$aa\underline{abbbb}b \rightsquigarrow aa\underline{T}bbbb$$

LL parsing

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

“You’re the Parser” Game:
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 $E \rightarrow \boxed{\text{num} = \text{num}}$

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

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

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 $E \rightarrow E + E$

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

When parse is here, can't determine whether it's an assign or a plus

Need to save input somewhere, like a stack: this is a job for a (D)PDA!!

| 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 |

LR parsing

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

$$S \rightarrow S ; S \quad E \rightarrow \text{id}$$

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

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

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

LR parsing

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

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

| 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 $E \rightarrow \text{num}$ |
| 1 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 |

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 ₄ | := c + (d := 5 + 6 , d) \$ | shift |
| 1 id ₄ := ₆ | := 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



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 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 $E \rightarrow \text{num}$ |
| 1 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 |

Can determine (rightmost) rule



LR parsing

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

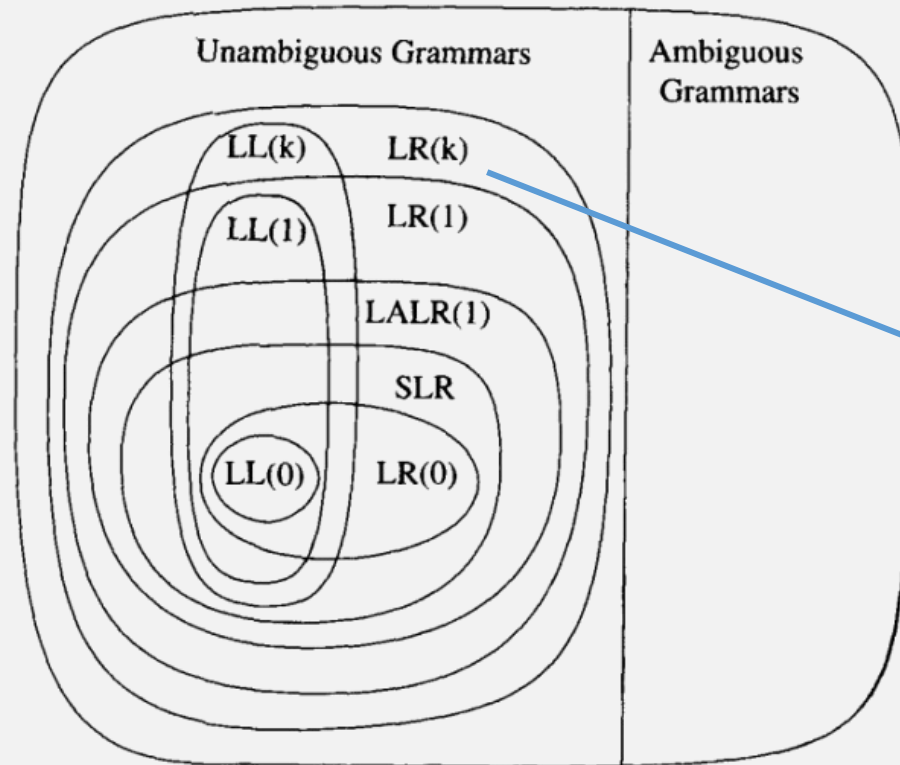
$$S \rightarrow S ; S \qquad E \rightarrow \text{id}$$

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

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

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

To learn more, take a Compilers Class!



A program (string of chars)



Program "words"



Abstract Syntax tree (AST)



Need computation that goes beyond CFLs

Non-CFLs

Flashback: Pumping Lemma for Reg Langs

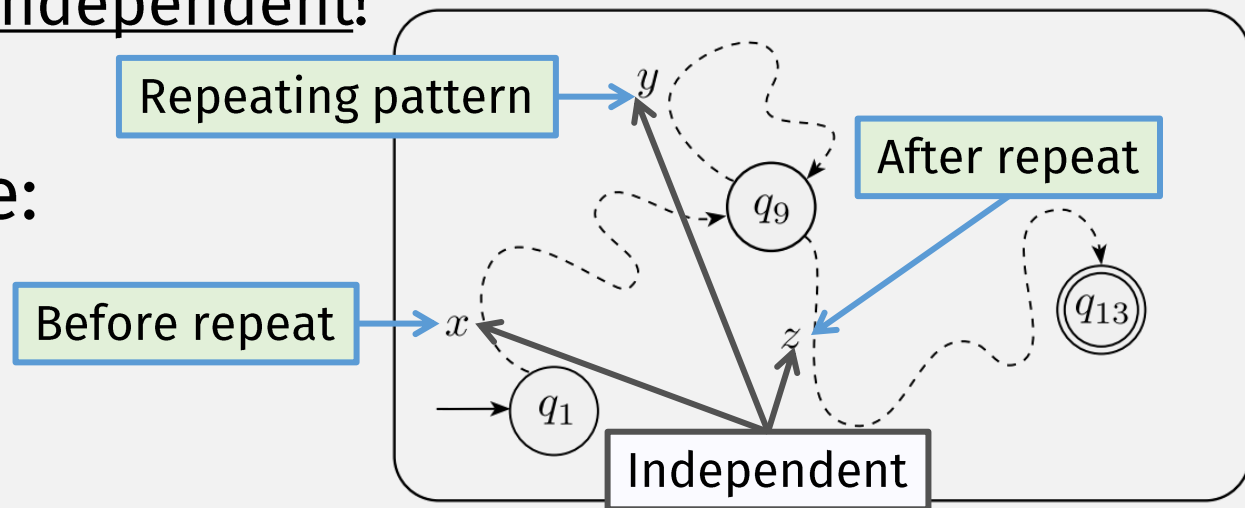
- The Pumping Lemma describes how strings repeat
- Regular language strings can (only) repeat using Kleene pattern
 - But the substrings are independent!

- A non-regular language:

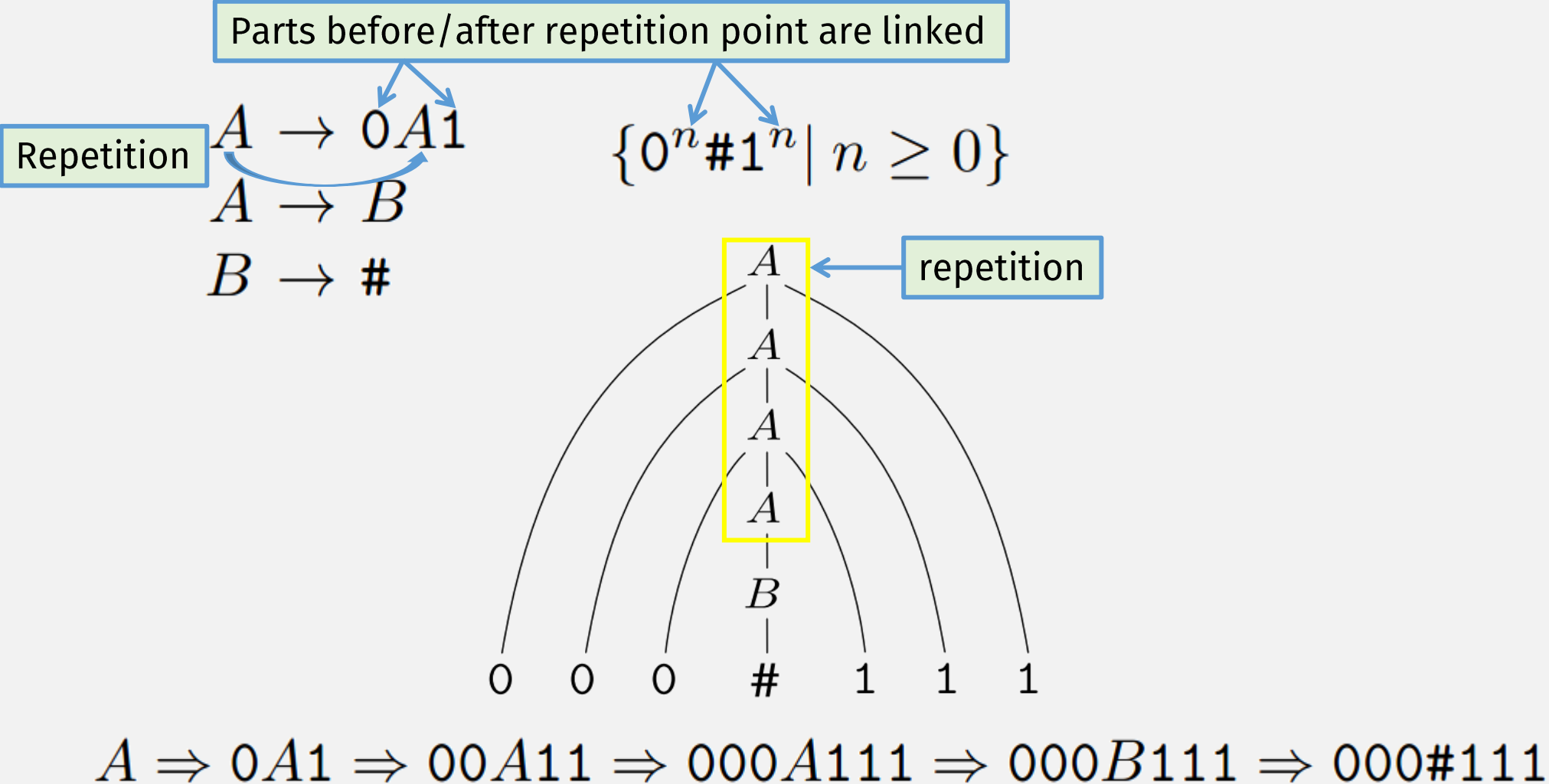
$$\{0^n 1^n \mid n \geq 0\}$$

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

- How do CFLs repeat?



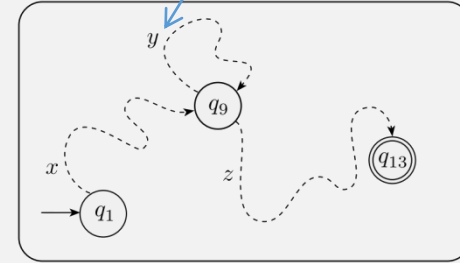
Repetition and Dependency in CFLs



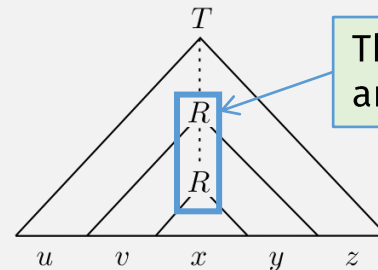
How Do Strings in CFLs Repeat?

NFA can take loop transition any number of times, to process repeated y in input

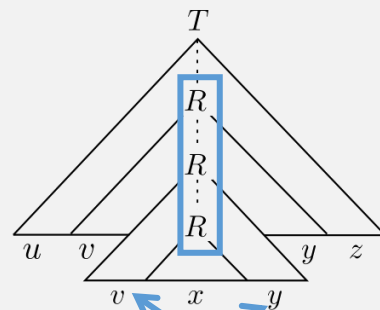
- Strings in regular languages repeat states



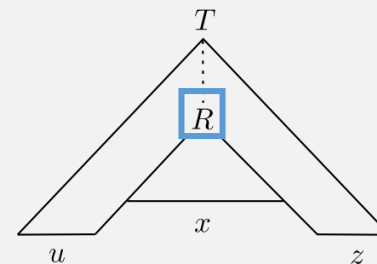
- Strings in CFLs repeat subtrees in the parse tree



This subtree can be repeated any number of times



Linked parts

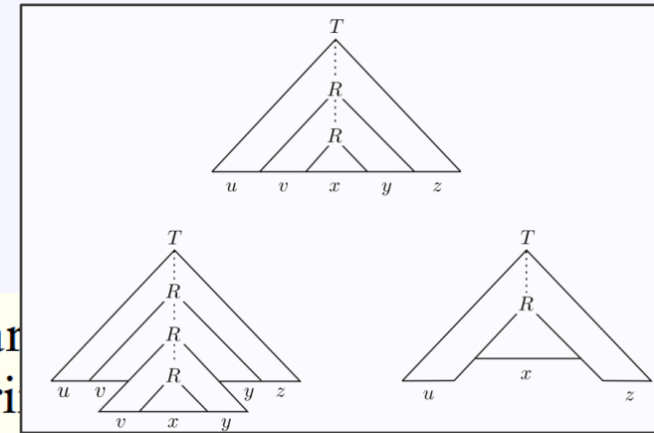


Pumping Lemma for CFLS

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 two pumpable parts. But they must be pumped together!

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 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. $|xy| \leq p$.

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

- Choose another 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. $|vxy| \leq p$.

This language is not a CFL!

CFL Pumping Lemma is Weird?

???

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. $|vxy| \leq p$.

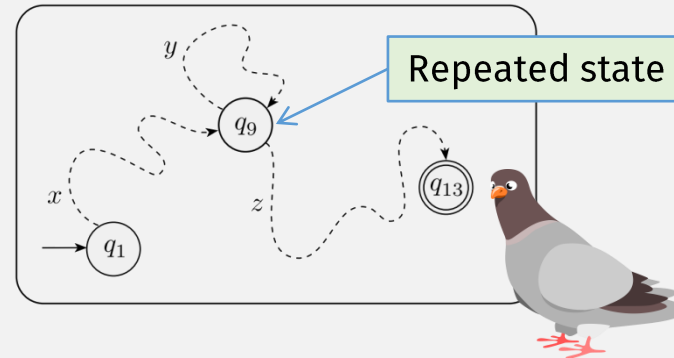
Review: Regular Language Pumping Lemma

- The pumping length 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 $i \geq 0$, $xy^iz \in A$,
2. $|y| > 0$, and
3. $|xy| \leq p$.

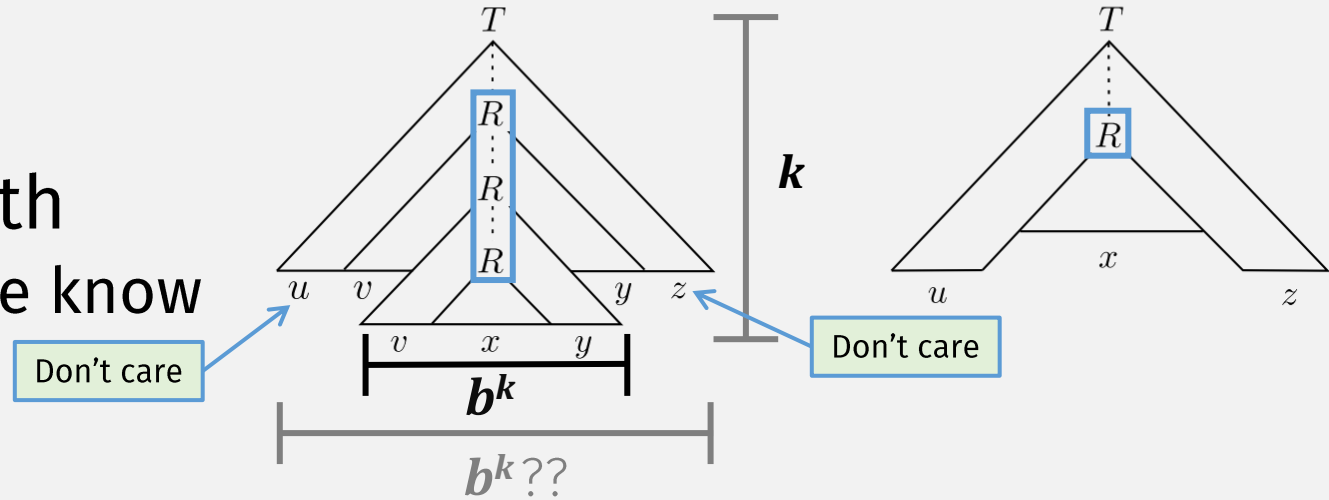
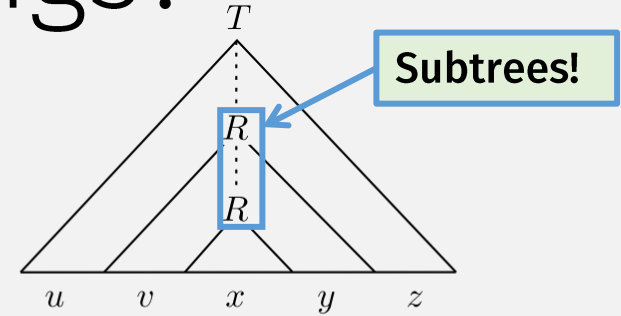
- If string length $>$ # of states, then some state must repeat



- If a state is repeated once, then it can repeat multiple times

Repeating Pattern in CFL Strings?

- When are we guaranteed to have a repeated subtree?
 - When height 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 ???



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 not say:
 - “All nonCFLs are not pumpable”
 - (statement != it’s inverse)
 - So Pumping Lemma might not be able to prove some non-CFLs!

Example:

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

- For any counterexample, split into $uvxyz$ where,
 - v = first char
 - z = remaining chars
 - $u = x = y = \varepsilon$
- If there are **as** ...
 - ... it’s pumpable bc # of **as** is arbitrary
- There there are no **as**
 - ... 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. vw has at most n distinguished positions.
2. vx has at least one distinguished position.
3. For all i , uv^iwx^iy is in L .

Says that every long enough segment must be pumpable

Example:

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

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

Counterexample: $ab^n c^n d^n$

- n "distinguished" positions must include non- a character
 - Impossible to pump no matter which n chars are chosen

A Practical Non-CFL

- **XML**

- ELEMENT \rightarrow \langle TAG \rangle CONTENT \langle /TAG \rangle
- Where TAG is any string

- XML also looks like this non-CFL: $D = \{ww \mid w \in \{0,1\}^*\}$

- This means XML is not context-free!

- Note: HTML is context-free because ...
- ... there are only a finite number of tags,
- so they can be embedded into a finite number of rules.

- In practice:

- XML is parsed as a CFL, with a CFG
- Then matching tags checked in a 2nd 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 :

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.

Infinite memory, initially starts with input

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

In-class quiz 10/6

See gradescope