# CS 420 / CS 620

## CFGs vs PDAs
## subCFLs and DPDAs

Monday October 27, 2025

UMass Boston Computer Science



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

# Announcements

- HW 7
  - Out: Mon 10/20 12pm (noon)
  - Due: Mon 10/27 12pm (noon)

- HW notes
  - Correct Gradescope page assignment of problems is now part of the correctness each submission

- Gradescope note
  - Regrade requests must address a specific deduction

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

# Regular Language vs CFL Comparison

| Regular Languages | Context-Free Languages (CFLs) |
|---|---|
| Regular Expression | Context-Free Grammar (CFG) |
| describes a **Regular Lang** | describes a **CFL** |
| | |
| Deterministic Finite-State Automata (DFA) | **Push-down Automata** (PDA) |
| recognizes a **Regular Lang** | recognizes a **CFL** |
| | |
| Proved: | Must Prove: |
| Regular Lang ⇔ Regular Expr ☑ | CFL ⇔ PDA  **???** |

thm

def

def

thm

# A lang is a CFL **iff** some PDA recognizes it

⇒ If **a language is a** **CFL**, then **a PDA recognizes it**
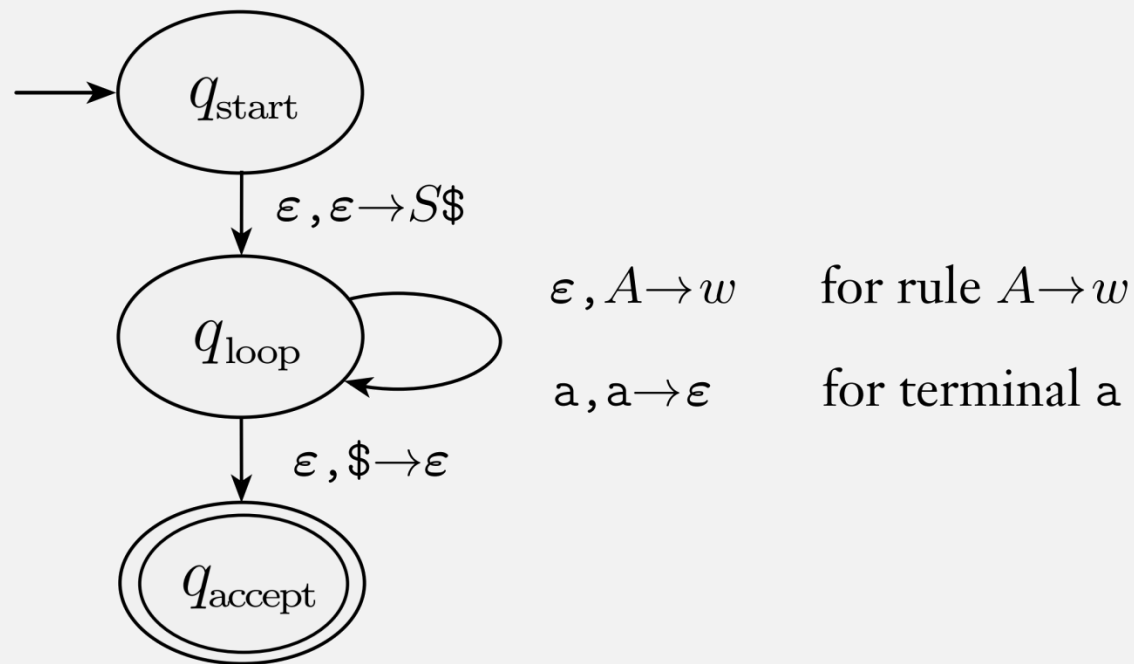
- <u>We know</u>: A **CFL** has a **CFG describing it** (definition of CFL)
- <u>To prove this part, show</u>: the **CFG** has an <u>equivalent</u> PDA

⇐ If **a PDA recognizes a language,** then **it's a CFL**

# **CFG→PDA** (sketch)
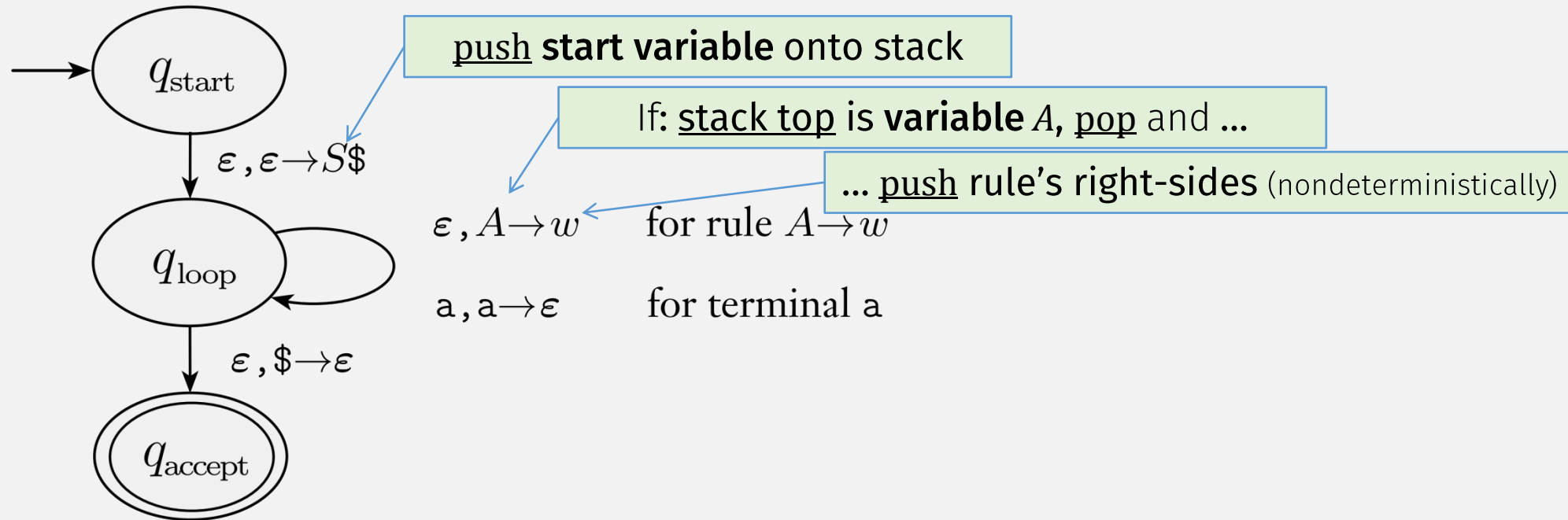
- Construct PDA from CFG such that:
  - PDA accepts input only if CFG generates it

- PDA:
  - simulates generating a string with CFG rules
  - by (nondeterministically) **trying all rules** to find the right ones



$q_{\text{start}}$

$\varepsilon, \varepsilon \to S\$$

$q_{\text{loop}}$

$\varepsilon, A \to w$    for rule $A \to w$

$\text{a}, \text{a} \to \varepsilon$    for terminal a

$\varepsilon, \$ \to \varepsilon$

$q_{\text{accept}}$

# **CFG→PDA** (sketch)
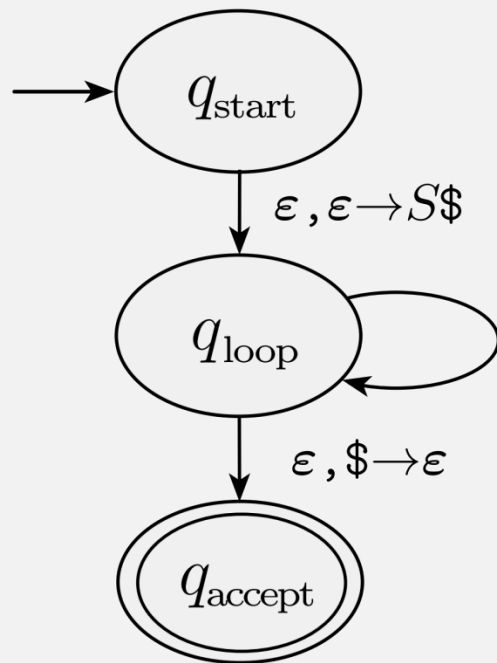
- Construct PDA from CFG such that:
  - PDA accepts input only if CFG generates it

- PDA:
  - simulates generating a string with CFG rules
  - by (nondeterministically) **trying all rules** to find the right ones



push **start variable** onto stack

If: <u>stack top</u> is **variable** $A$, <u>pop</u> and …

… <u>push</u> rule's right-sides (nondeterministically)

$\varepsilon, A \to w$    for rule $A \to w$

$\mathrm{a}, \mathrm{a} \to \varepsilon$      for terminal a

# CFG→PDA (sketch)

- Construct PDA from CFG such that:
  - PDA accepts input only if CFG generates it

- PDA:
  - simulates generating a string with CFG rules
  - by (nondeterministically) **trying all rules** to find the right ones



$q_{start}$

$\varepsilon, \varepsilon \to S\$$

$q_{loop}$

$\varepsilon, \$ \to \varepsilon$

$q_{accept}$

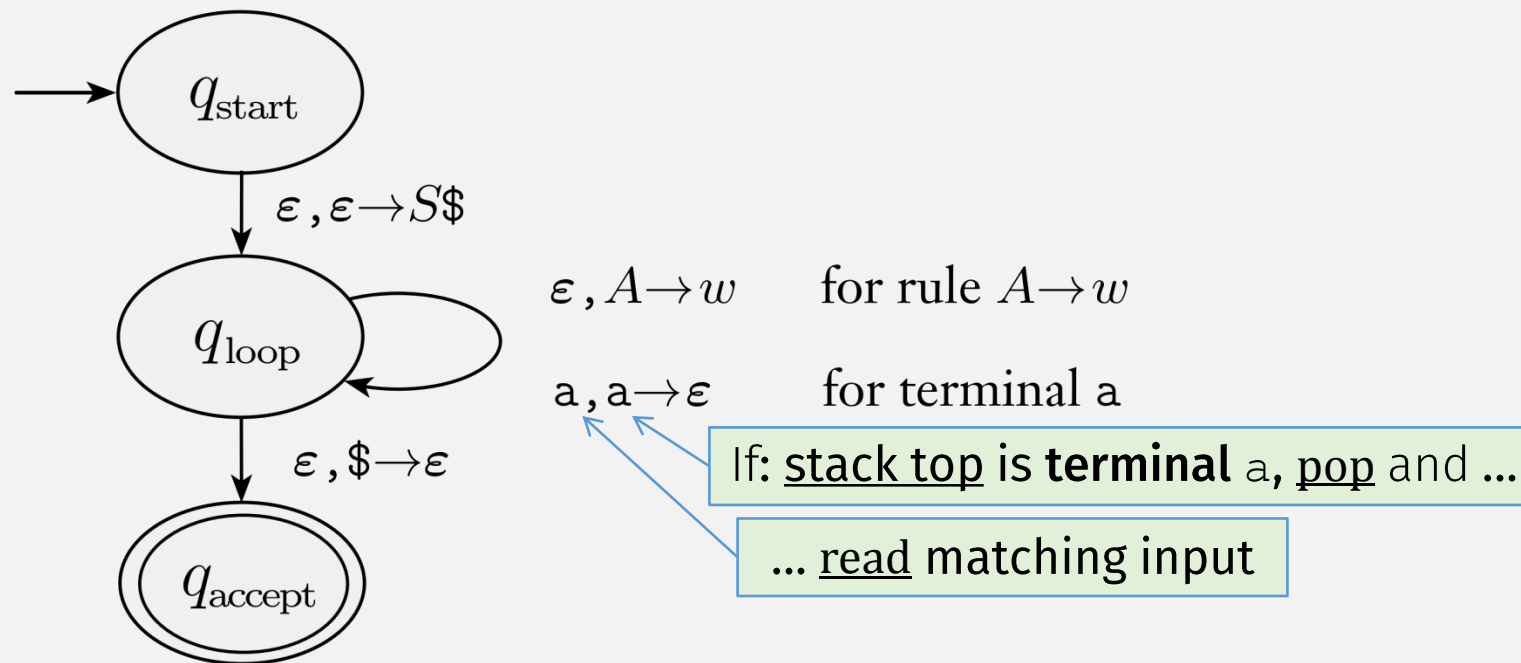Summary: **convert every** <u>CFG rule</u> **to PDA "loop" transition** that:
- Pops LHS variable
- Pushes RHS

(Stack is "workspace" containing intermediate string of vars + terminals)

$\varepsilon, A \to w$     for rule $A \to w$

$a, a \to \varepsilon$     for terminal a

# **CFG→PDA** (sketch)

- Construct PDA from CFG such that:
  - PDA accepts input only if CFG generates it

- PDA:
  - simulates generating a string with CFG rules
  - by (nondeterministically) **trying all rules** to find the right ones

$q_{\text{start}}$

$\varepsilon, \varepsilon \to S\$$

$q_{\text{loop}}$

$\varepsilon, \$ \to \varepsilon$

$q_{\text{accept}}$

$\varepsilon, A \to w$     for rule $A \to w$

$\texttt{a}, \texttt{a} \to \varepsilon$     for terminal a

If: <u>stack top</u> is **terminal** a, <u>pop</u> and ...

... <u>read</u> matching input

# CFG→PDA (sketch)

- Construct PDA from CFG such that:
  - PDA accepts input only if CFG generates it

- PDA:
  - simulates generating a string with CFG rules
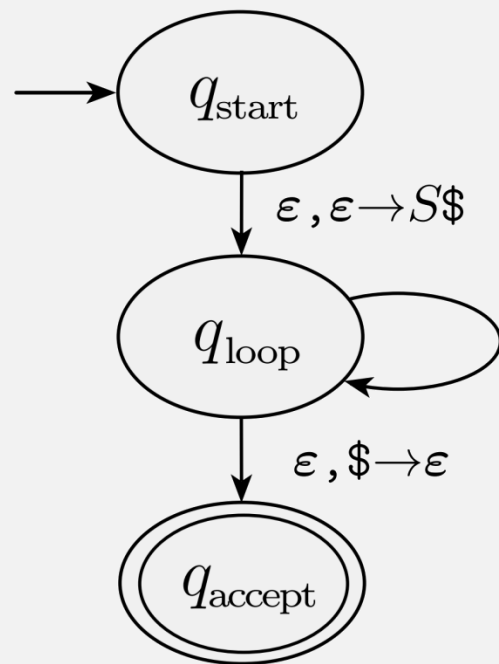  - by (nondeterministically) trying all rules to find the right ones



$$\varepsilon, \varepsilon \rightarrow S\$$$

$$\varepsilon, A \rightarrow w \qquad \text{for rule } A \rightarrow w$$

$$\varepsilon, \$ \rightarrow \varepsilon$$

$$a, a \rightarrow \varepsilon \qquad \text{for terminal a}$$

Summary: **convert every terminal to "loop" transition** that:
- Reads input char
- Pops matching char on stack

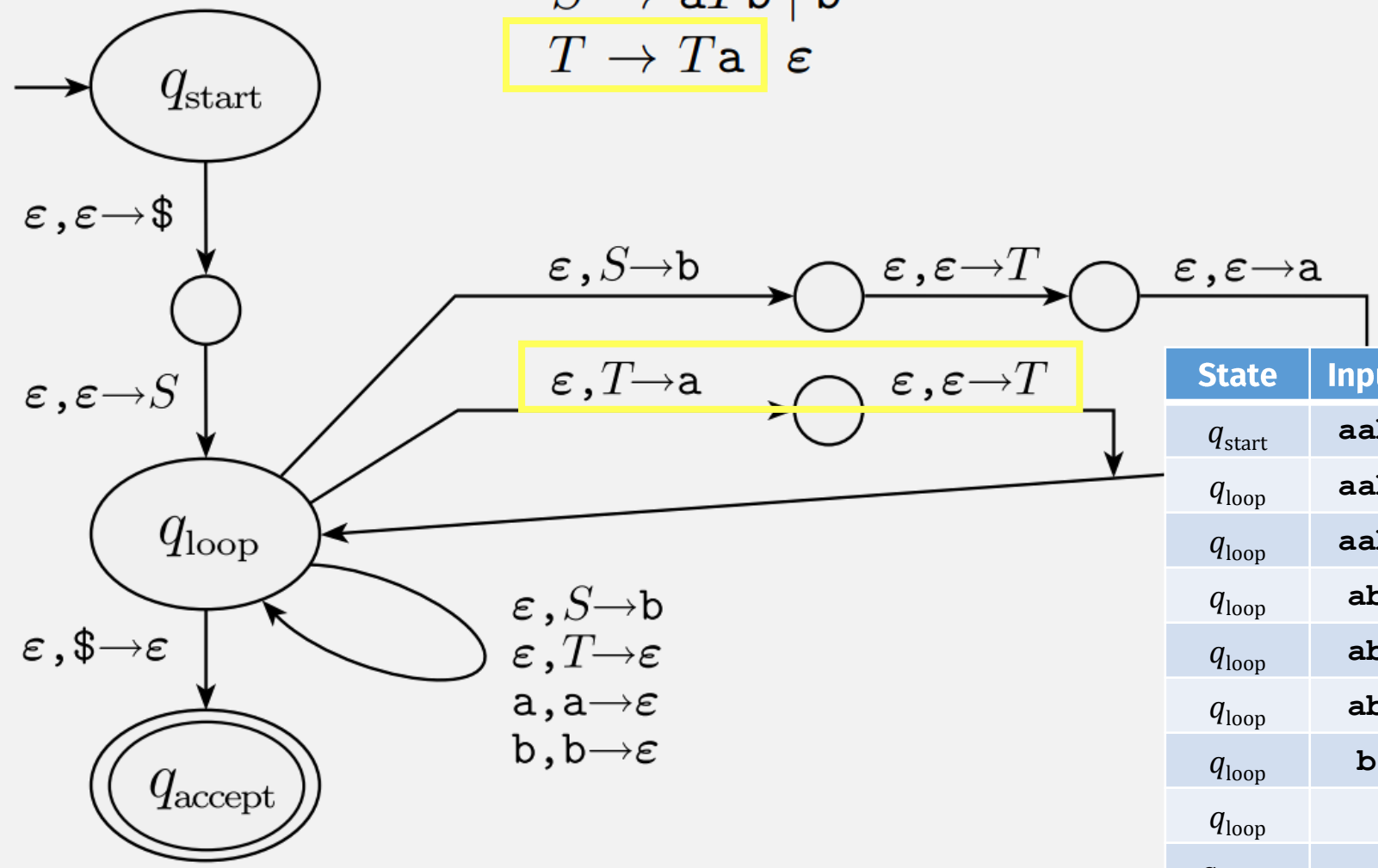(**Read** the **terminals** as they become known)

# Example **CFG→PDA**

Example Derivation using CFG:
$S \Rightarrow \mathbf{a}T\mathbf{b}$ (using rule $S \rightarrow \mathbf{a}T\mathbf{b}$)
$\Rightarrow \mathbf{a}T\mathbf{ab}$ (using rule $T \rightarrow T\mathbf{a}$)
$\Rightarrow \mathbf{aab}$ (using rule $T \rightarrow \varepsilon$)

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



$\varepsilon, \varepsilon \rightarrow \$$

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

$\varepsilon, \varepsilon \rightarrow S$

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

$q_{\text{start}}$

$q_{\text{loop}}$

$\varepsilon, \$ \rightarrow \varepsilon$

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

$q_{\text{accept}}$

PDA Example

| State | Input | Stack | Equiv Rule |
|-------|-------|-------|------------|
| $q_{\text{start}}$ | **aab** | | |
| $q_{\text{loop}}$ | **aab** | $S\$$ | |
| $q_{\text{loop}}$ | **aab** | $\mathbf{a}T\mathbf{b}\$$ | $S \rightarrow \mathbf{a}T\mathbf{b}$ |
| $q_{\text{loop}}$ | **ab** | $T\mathbf{b}\$$ | |
| $q_{\text{loop}}$ | **ab** | $T\mathbf{ab}\$$ | $T \rightarrow T\mathbf{a}$ |
| $q_{\text{loop}}$ | **ab** | $\mathbf{ab}\$$ | $T \rightarrow \varepsilon$ |
| $q_{\text{loop}}$ | **b** | $\mathbf{b}\$$ | |
| $q_{\text{loop}}$ | | $\$$ | |
| $q_{\text{accept}}$ | | | |

# A lang is a CFL **iff** some PDA recognizes it

☑ ⇒ If a language is a CFL, then a PDA recognizes it
- Convert **CFG→PDA**

⇐ If a PDA recognizes a language, then it's a CFL
- <u>To prove this part:</u> show PDA has an equivalent CFG

# **PDA→CFG**: Prelims

Before converting PDA to CFG, <u>modify</u> it so :

**1.** It has a single accept state, $q_{\text{accept}}$.

**2.** It empties its stack before accepting.

**3.** Each transition either pushes a symbol onto the stack (a *push* move) or pops one off the stack (a *pop* move), but it does not do both at the same time.

<u>Important</u>:
This doesn't change the language recognized by the PDA

# PDA $P$ -> CFG $G$ : Transitions and Variables

$$P = (Q, \Sigma, \Gamma, \delta, q_0, \{q_{\text{accept}}\})$$

variables of $G$ are $\{A_{pq} | \ p, q \in Q\}$

- <u>Want</u>: if $P$ goes from **state** $p$ to $q$ **reading input** $x$, then some $A_{pq}$ generates $x$



- <u>So</u>: For **every** <u>pair</u> of **states** $p, q$ in $P$, **add variable** $A_{pq}$ to $G$

- <u>Then</u>: **connect** the **variables together** by,
  - **Add rules:** $A_{pq} \rightarrow A_{pr}A_{rq}$, for **each state** $r$



  - These rules allow: **grammar** to **simulate every possible transition**
  - (We haven't added **input read/generated terminals** yet)

The Key IDEA

- <u>To add terminals</u>: **pair up stack pushes** and **pops** (essence of a CFL)

# PDA $P$ -> CFG $G$ : Generating Strings

$$P = (Q, \Sigma, \Gamma, \delta, q_0, \{q_{\text{accept}}\})$$

variables of $G$ are $\{A_{pq} \mid p, q \in Q\}$

• <u>The key</u>: **pair up stack pushes and pops** (essence of a CFL)

if $\delta(p, a, \varepsilon)$ contains $(r, u)$ and $\delta(s, b, u)$ contains $(q, \varepsilon)$,

put the rule $A_{pq} \to aA_{rs}b$ in $G$

# PDA $P$ -> CFG $G$ : Generating Strings

$$P = (Q, \Sigma, \Gamma, \delta, q_0, \{q_{\text{accept}}\})$$

variables of $G$ are $\{A_{pq} | p, q \in Q\}$

- <u>The key</u>: **pair up stack pushes and pops** (essence of a CFL)

if $\delta(p, a, \varepsilon)$ contains $(r, u)$ and $\delta(s, b, u)$ contains $(q, \varepsilon)$,

put the rule $A_{pq} \rightarrow a A_{rs} b$ in $G$

# PDA $P$ -> CFG $G$ : Generating Strings

$$P = (Q, \Sigma, \Gamma, \delta, q_0, \{q_{\text{accept}}\})$$

variables of $G$ are $\{A_{pq} | \; p, q \in Q\}$

- <u>The key</u>: **pair up stack pushes and pops** (essence of a CFL)

if $\delta(p, a, \varepsilon)$ contains $(r, u)$ and $\delta(s, b, u)$ contains $(q, \varepsilon)$,

put the rule $A_{pq} \rightarrow a A_{rs} b$ in $G$

# A language is a CFL ⇔ A PDA recognizes it

☑ ⇒ If **a language is a CFL,** then **a PDA recognizes it**
   • Convert **CFG→PDA**

☑ ⇐ If **a PDA recognizes a language,** then **it's a CFL**
   • Convert **PDA→CFG**

# Regular Language vs CFL Comparison

| Regular Languages | Context-Free Languages (CFLs) |
|---|---|
| Regular Expression | Context-Free Grammar (CFG) |
| <u>describes</u> a **Regular Lang** | <u>describes</u> a **CFL** |
| | |
| Deterministic Finite-State Automata (DFA) | **Push-down Automata** (PDA) |
| <u>recognizes</u> a **Regular Lang** | <u>recognizes</u> a **CFL** |
| | |
| <u>Proved:</u> | <u>Proved:</u> |
| Regular Lang ⇔Regular Expr ☑ | CFL ⇔ PDA ☑ |

thm

def

def

thm

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

# 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

context-free
languages
(CFLs)

regular
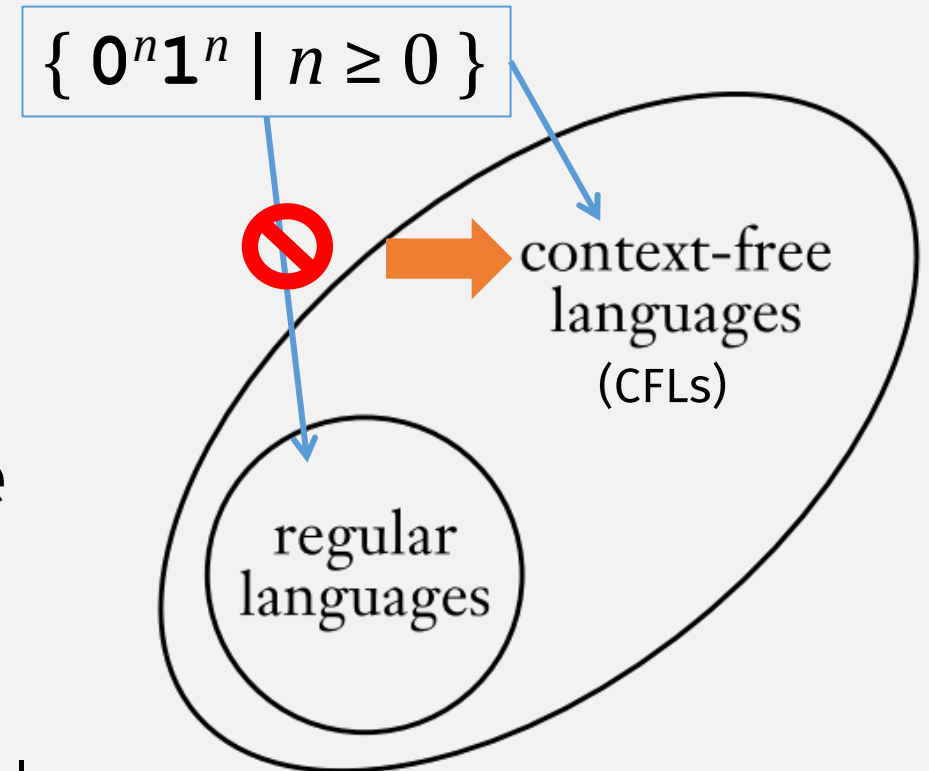languages

# How to <u>Prove</u> This Diagram "Correct"?

1. Every regular language is a CFL

   ➡️ 2. **Not every CFL is a regular language**

   Find a <u>counterexample</u> CFL that is **not regular**

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

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

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

🚫 ➡️ context-free languages (CFLs)

regular languages

# How to <u>Prove</u> This Diagram "Correct"?

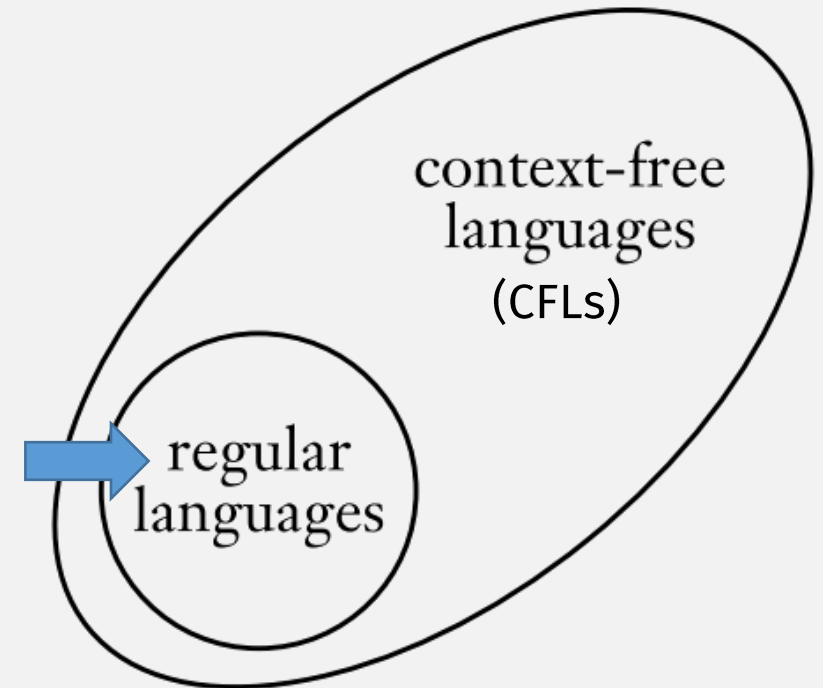➡️ 1. **Every regular language is a CFL**

> 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
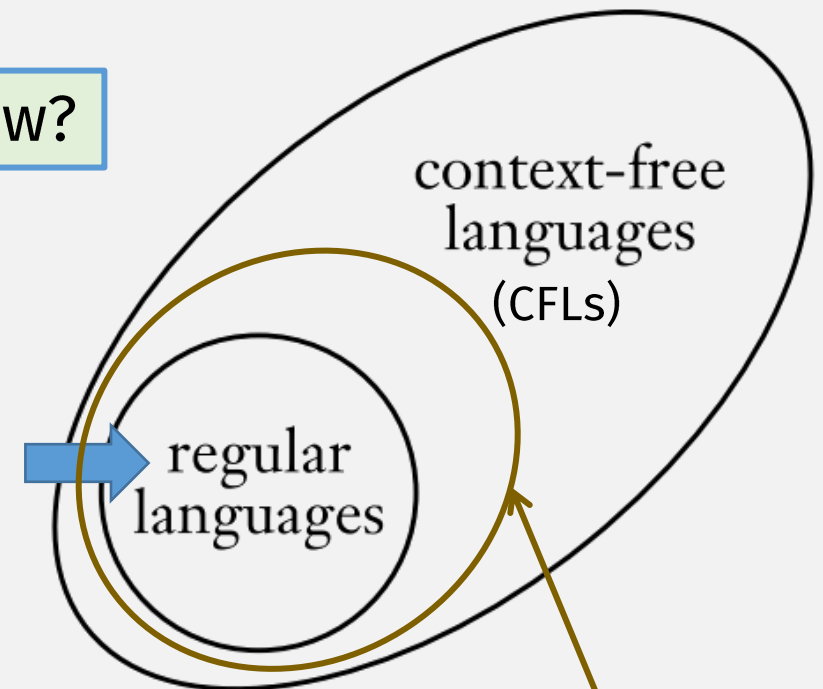
context-free
languages
(CFLs)

regular
languages

# Regular Languages are CFLs: 3 Ways to Prove

- DFA → CFG or PDA

Coming soon to a future hw?

- NFA → CFG or PDA

- Regular expression → CFG or PDA

context-free
languages
(CFLs)

regular
languages

Are there other interesting
subsets of CFLs?

# Deterministic CFLs and DPDAs

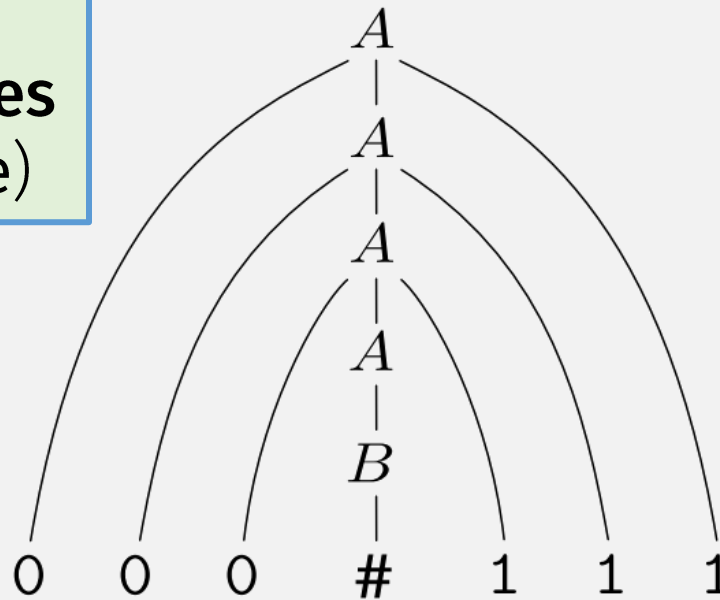# *Previously:* Generating Strings

**Generating** strings:
1. <u>Start</u> with **start variable,**
2. <u>Repeatedly</u> *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)

In practice, opposite is more interesting:
1. Start with **string**,
2. Then **parse** into **parse tree**

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



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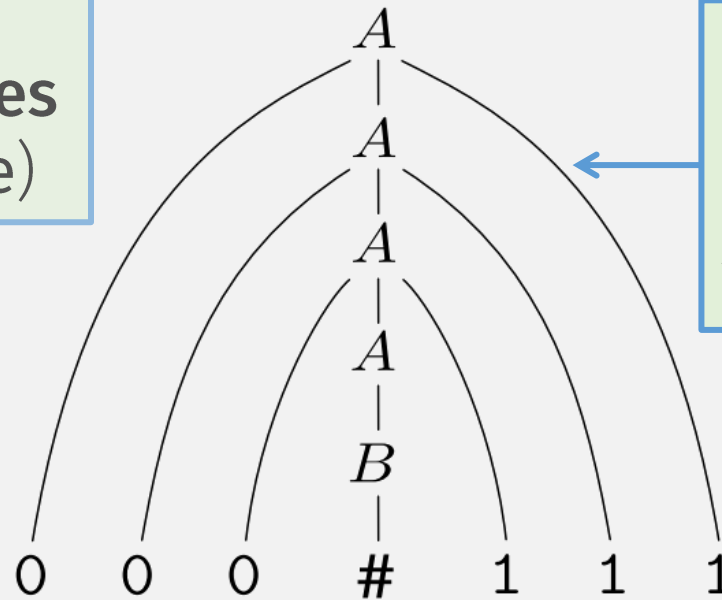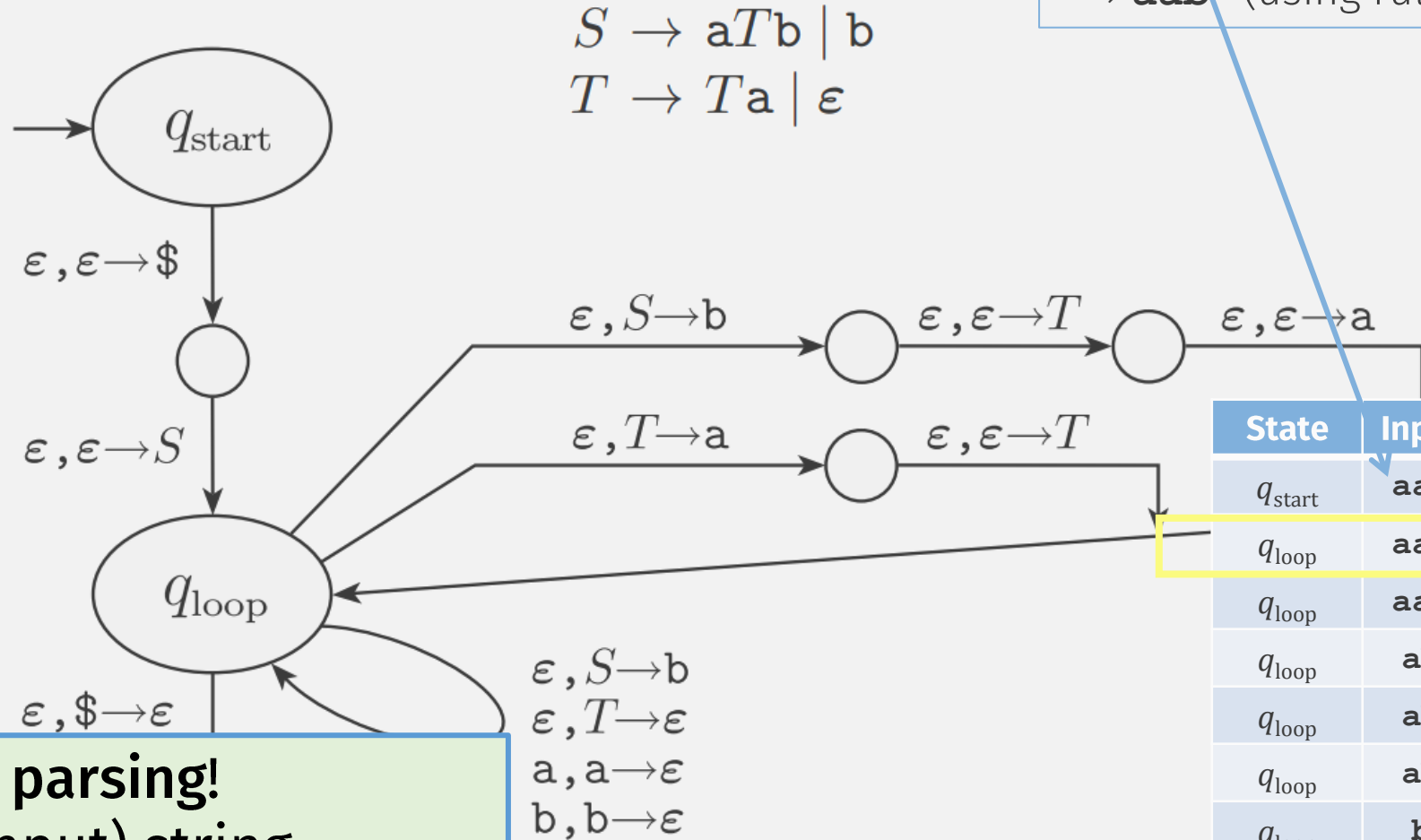
# Generating vs Parsing

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

*Previously:* # Example **CFG→PDA**

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



PDA Example

| State | Input | Stack | Equiv Rule |
|---|---|---|---|
| $q_{\text{start}}$ | **aab** | | |
| $q_{\text{loop}}$ | **aab** | $S\$$ | |
| $q_{\text{loop}}$ | **aab** | $\mathbf{a}T\mathbf{b}\$$ | $S \rightarrow \mathbf{a}T\mathbf{b}$ |
| $q_{\text{loop}}$ | **ab** | $T\mathbf{b}\$$ | |
| $q_{\text{loop}}$ | **ab** | $T\mathbf{ab}\$$ | $T \rightarrow T\mathbf{a}$ |
| $q_{\text{loop}}$ | **ab** | $\mathbf{ab}\$$ | $T \rightarrow \varepsilon$ |
| $q_{\text{loop}}$ | **b** | $\mathbf{b}\$$ | |
| $q_{\text{loop}}$ | | $\$$ | |
| $q_{\text{accept}}$ | | | |

This Machine is **parsing!**
1. <u>Start</u> with (input) string,
2. <u>Find</u> **rules** that **generate** string

# Generating vs Parsing

- In practice, **parsing** a string more important than **generating**
  - E.g., **a compiler** (first step) **parses source code string** 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**)

# *Previously:* (Nondeterministic) PDA

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



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

A parser implementation can't do this!

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

"do nothing"

Not power set

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)

Must consider: $\varepsilon$ reads or stack ops!
E.g., if $\delta(q, a, X)$ does "something", then $\delta(q, \varepsilon, X)$ must "do nothing"

# 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 can non-deterministically "<u>try all</u> rules" (abandoning failed attempts)

- A DPDA must <u>choose one</u> rule at each step! (cant go back after reading input!)

**Parsing** = deriving reversed:
<u>start</u> with **string**, <u>end</u> with **parse tree**

used $S$ rule

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

used $T$ rule

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

$$\mathbf{aa\underline{a}bbbbb} \rightarrowtail \mathbf{aa\underline{T}bbb}$$

Choosing "correct" rule depends on rest of the input!

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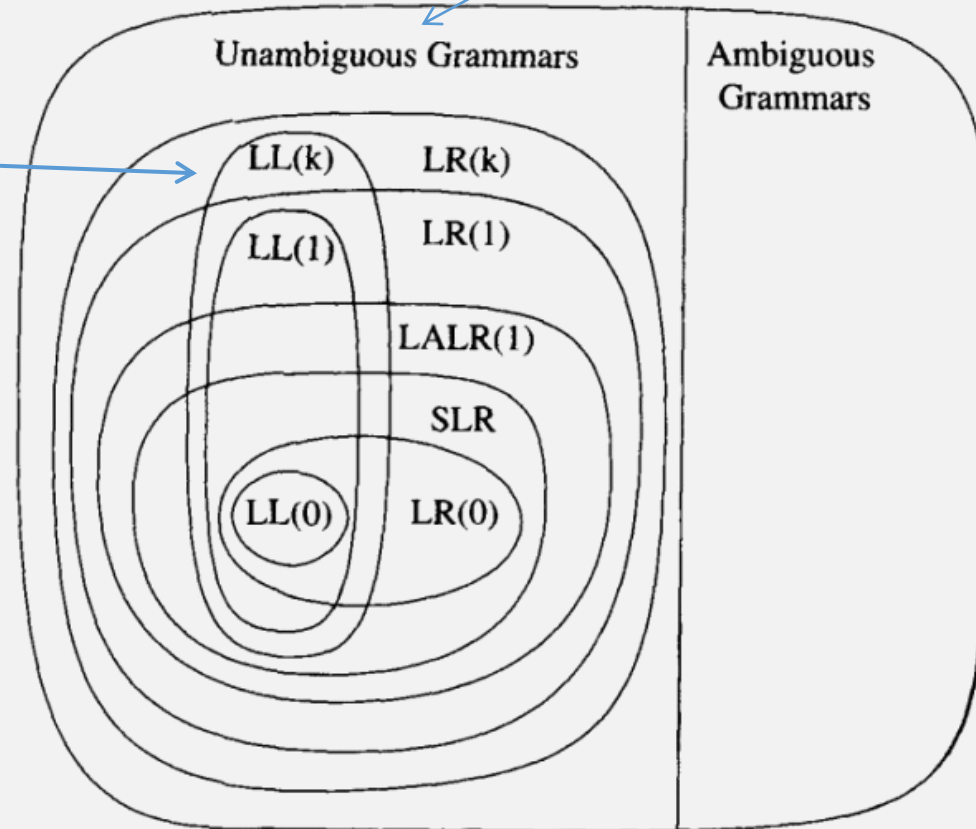
# Subclasses of CFLs

Umambiguous CFLs / PDAs

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 ) ; …`)

↓

**Lexer**

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

↓

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**)!

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

It's more than an analogy!

<u>This</u> DFA <u>is</u> a real program!

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

# Compiler Stages

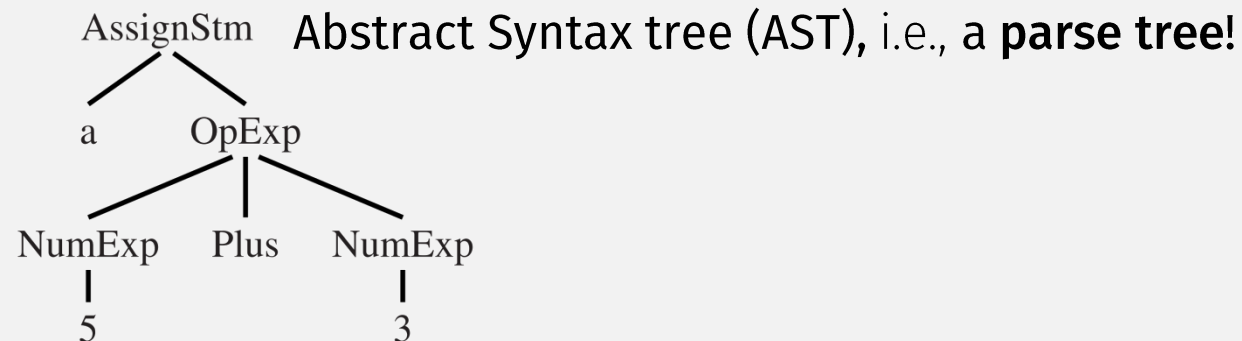A program (chars) (e.g., `a : = ( 5 + 3 ) ; …`)

**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!

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

```
AssignStm
├── a
└── OpExp
    ├── NumExp
    │   └── 5
    ├── Plus
    └── NumExp
        └── 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
```

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

It's more than an analogy!

This CFG is a real program!

Just write the CFG!

A "`yacc`" tool converts the program:
- from "CFG Lang" …
- to an **equivalent** one in C !

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

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

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

- **PDA**: can non-deterministically "<u>try all</u> rules" (abandoning failed attempts);
- **DPDA**: must <u>choose one</u> rule at each step!

Should use $S$ rule

$$\mathbf{aa\underline{a}bbb} \rightarrowtail \mathbf{a\underline{a}Sbb}$$

$$\mathbf{aa\underline{a}}$$

Should use $T$ rule

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

$$\mathbf{aa\underline{a}bbbbb} \rightarrowtail \mathbf{a\underline{a}Tbbb}$$

To choose "correct" rule, need to **"look ahead"** at rest of the input!

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

2) choose "look ahead" amount

2 parser design decisions:
1) Parse from left, or from right

All CFLS

# LL parsing

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

**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

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

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

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

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

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

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

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

# LL parsing

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

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

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

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

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

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

**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
```

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

# LR parsing

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

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

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

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

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

# LR parsing

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

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

```
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

# 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 id$_4$ | := 7 ; b := c + ( d := 5 + 6 , d ) \$ | shift |
| 1 id$_4$ :=$_6$ | 7 ; b := c + ( d := 5 + 6 , d ) \$ | shift |

c

# 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 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 |

# LR parsing

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

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

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

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

| 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 |

Can determine (rightmost) rule

# LR parsing

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

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

| 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:=$E$ |

Can determine (rightmost) rule

# LR parsing

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

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

| 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 Parsers** also called "Shift-Reduce" Parsers

# To learn more, take a Compilers Class!



Unambiguous Grammars

Ambiguous Grammars

LL(k)    LR(k)
LL(1)    LR(1)
         LALR(1)
         SLR
LL(0)    LR(0)

A program (string of chars)

**Lexer**
(DFAs / NFAs)

Program "words"

**Parser**
(DPDAs)

Abstract Syntax tree (AST)

**???**

This phase needs computation that goes beyond CFLs
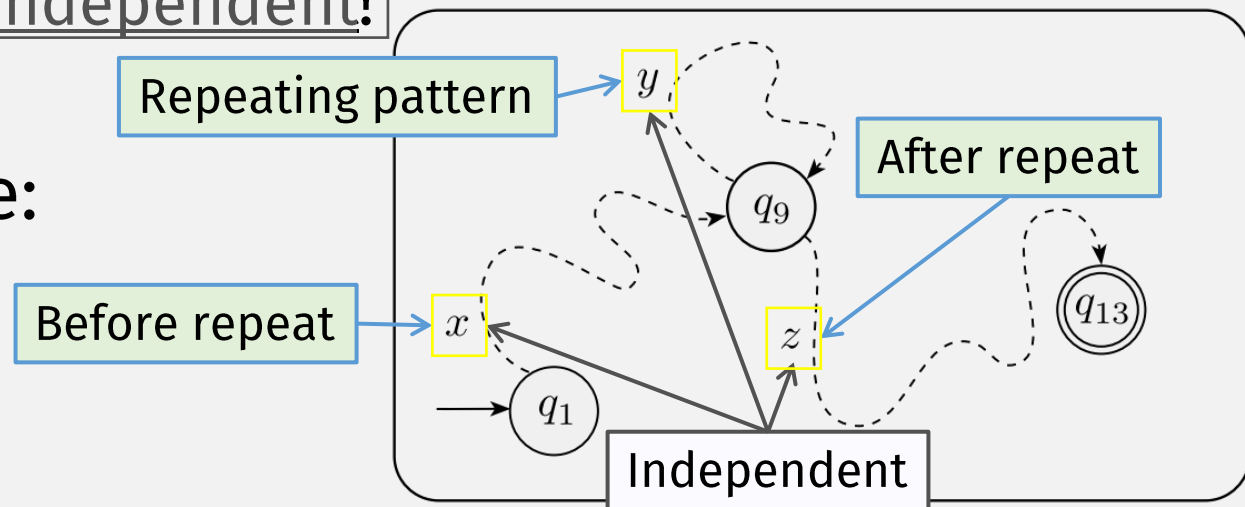
# *Flashback:* Pumping Lemma for Regular Langs

- **Pumping Lemma** describes how strings repeat

- **Regular language** strings repeat using **Kleene star** operation
  - Key: 3 substrings $x\,y\,z$ independent!

- A non-regular language:

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

Repeating pattern

After repeat

Before repeat

$y$

$q_9$

$q_{13}$

$x$

$z$

$q_1$

Independent

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

- Q: How do CFLs repeat?

# Repetition and Dependency in CFLs

Parts before/after repetition point underline{linked} (not independent)

Repetition

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

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

repetition

$$
\begin{array}{c}
A \\
| \\
A \\
| \\
A \\
| \\
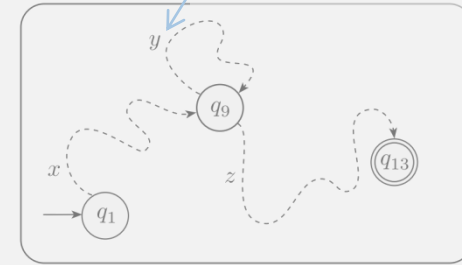A \\
| \\
B \\
| \\
0 \quad 0 \quad 0 \quad \# \quad 1 \quad 1 \quad 1
\end{array}
$$

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

# How Do Strings in CFLs Repeat?

- Strings in regular languages <u>repeat states</u>

- 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

5 substrings

Linked parts

Linked parts repeat together

# 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

> Two pumpable parts.
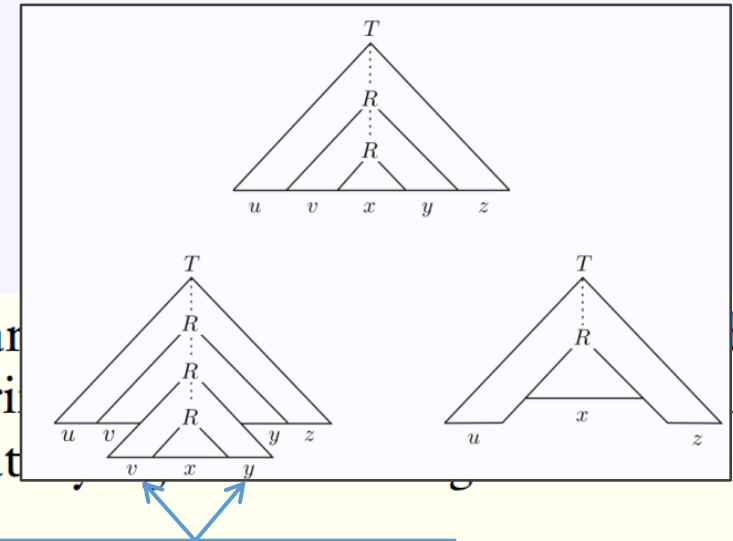> But they must be pumped together!

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

**Pumping lemma** If $A$ is a regular [...] ber $p$ (the pumping length) where if $s$ is any stri[...] s may be divided into three pieces, $s = xyz$, sat[...]

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

One pumpable part

Two pumpable parts, pumped together

# A Non CFL example

$$\text{language } B = \{\mathrm{a}^n\mathrm{b}^n\mathrm{c}^n \mid n \geq 0\} \text{ is not context free}$$

Intuition
- **Strings in CFLs** can have <u>two parts</u> that are **"pumped" together**
- **Language** $B$ **requires** <u>three parts</u> to be **"pumped" together**
- So **it's not a CFL!**

Proof?

Want to prove: $a^n b^n c^n$ **is not** a 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. $|vxy| \leq p$.

Reminder: CFL Pumping lemma says: all strings $a^n b^n c^n \geq$ length $p$ are splittable into $uvxyz$ where $v$ and $y$ are pumpable

Proof (by contradiction):

Now we must find a contradiction ...

- Assume: $a^n b^n c^n$ **is** a CFL
  - So it must satisfy the pumping lemma for CFLs
  - I.e., all strings $\geq$ length $p$ are pumpable
- Counterexample = $a^p b^p c^p$

Contradiction if:
- A string in the language ☑
- $\geq$ length $p$ ☑
- Is **not splittable** into $uvxyz$ where $v$ and $y$ are pumpable

**???**

$p$ as      $p$ bs      $p$ bs

a ... b ... c ...

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

Want to prove: $a^n b^n c^n$ **is not** a CFL

# Possible Splits

<u>Proof</u> (by contradiction):

*contradiction*

- <u>Assume:</u> $a^n b^n c^n$ **is** a CFL
  - So it must satisfy the pumping lemma for CFLs
  - I.e., all strings ≥ length $p$ are pumpable

- <u>Counterexample</u> = $a^p b^p c^p$
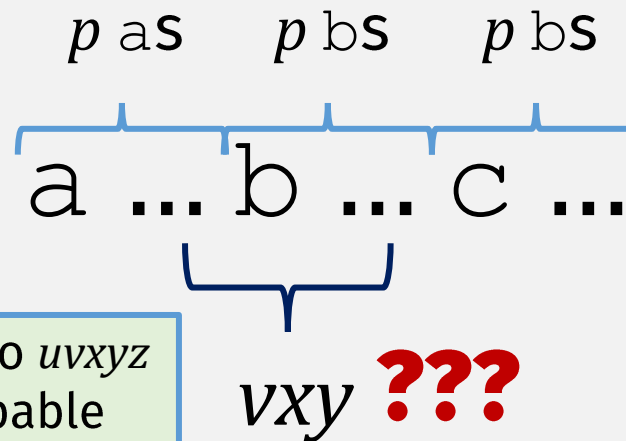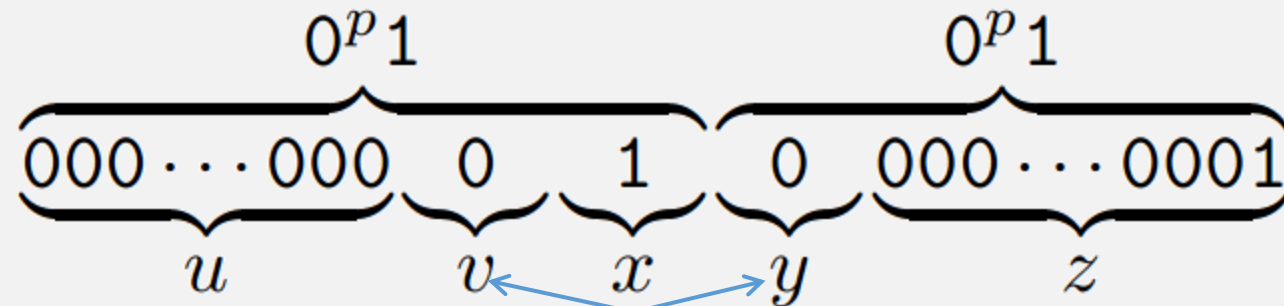
Contradiction if:
- A string in the language
- ≥ length $p$
- Is **not splittable** into $uvxyz$ where $v$ and $y$ are pumpable

- Possible Splits (using condition # 3: $|vxy| \leq p$)

Not pumpable

  ☒ • $vxy$ is all $a$s
  ☒ • $vxy$ is all $b$s
  ☒ • $vxy$ is all $c$s
  ☒ • $vxy$ has $a$s and $b$s
  ☒ • $vxy$ has $b$s and $c$s
   • ($vxy$ cannot have $a$s, $b$s, and $c$s)

$p$ $a$s    $p$ $b$s    $p$ $b$s

a ... b ... c ...

$vxy$ **???**

$a^p b^p c^p$ **cannot be split** into $uvxyz$ where $v$ and $y$ are pumpable

So $a^n b^n c^n$ **is not a CFL**

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

Be careful when choosing **counterexample** $s$: $\boxed{0^p 1 0^p 1}$

This $s$ **can** be pumped according to **CFL pumping lemma**:



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

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

So **this attempt** to **prove** that
the **language** is **not** a CFL **failed**.
(It **doesn't prove** that the language **is a CFL**!)

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

- Need <u>another counterexample</u> 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 $i \geq 0$, $uv^i x y^i z \in A$,

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

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

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

# A Practical Non-CFL

- **XML**
  - ELEMENT $\rightarrow$ <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 <u>finite</u> number of tags,
  - so **they can be embedded** into a <u>finite</u> 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<sup>nd</sup> pass with a <u>more powerful machine </u>…

# *Next:* 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 (initial contents are the input string)

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 arbitrary memory locations!