

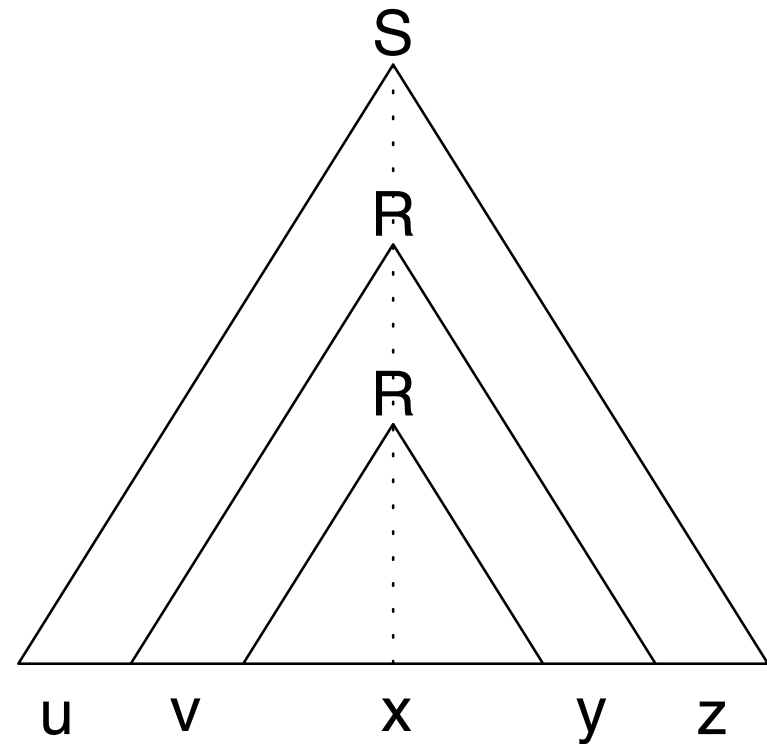
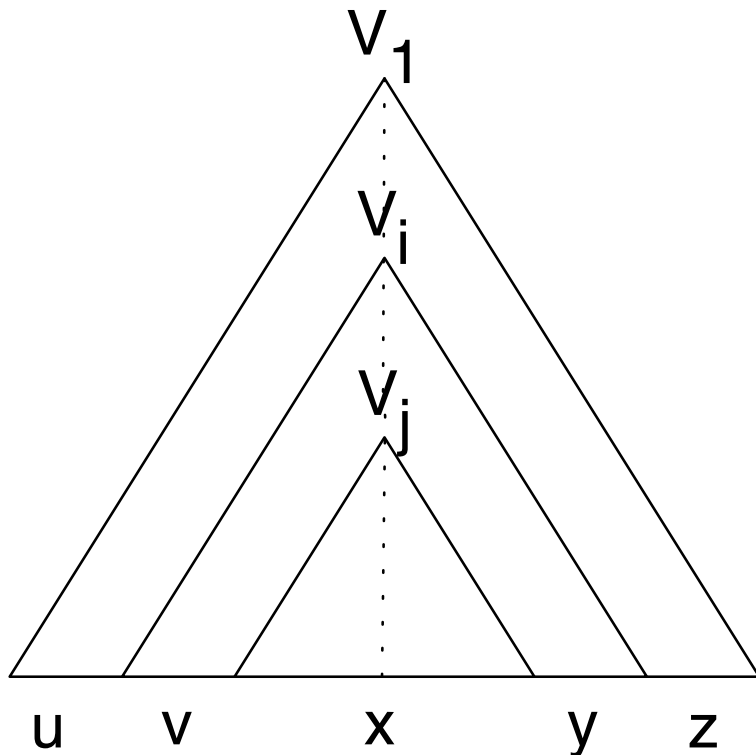
# Parsing & PDA

No, the other PDA.  
No, not that one either.

November 4, 2010

# Applying Pigeonhole

- ✓ Assume we have a string in  $L$  whose shortest parse tree has height  $\geq |V|+1$ . [height = edges]



# Pumping

Final questions:

(1) How do we know  $v$  and  $y$  aren't both  $\epsilon$ ?

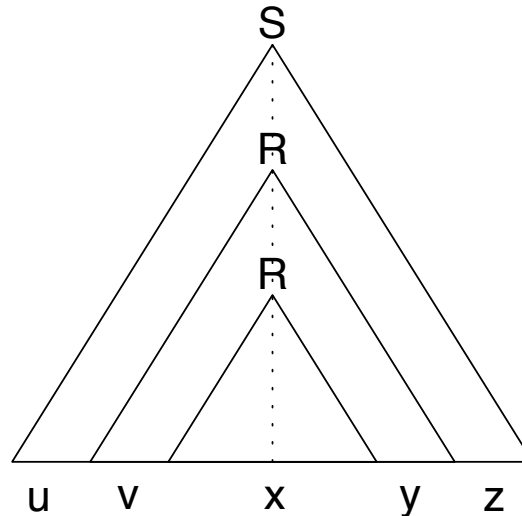
E.g.,

$R \Rightarrow P$

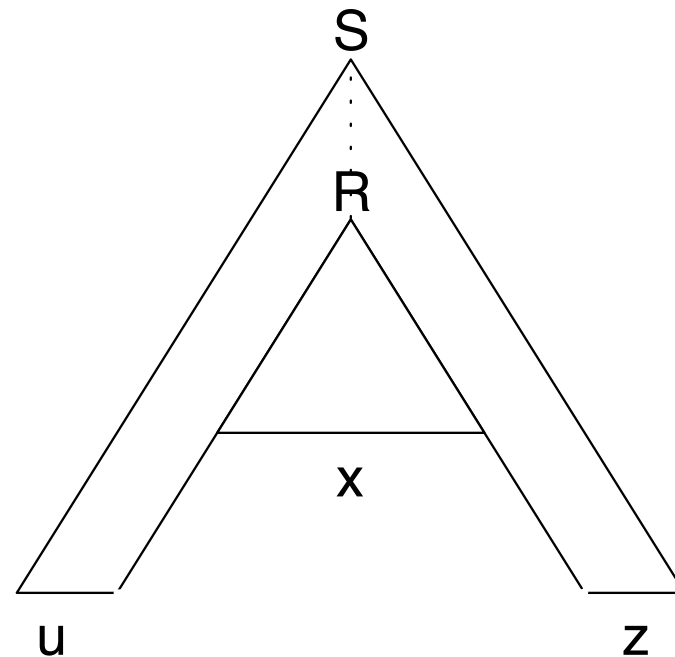
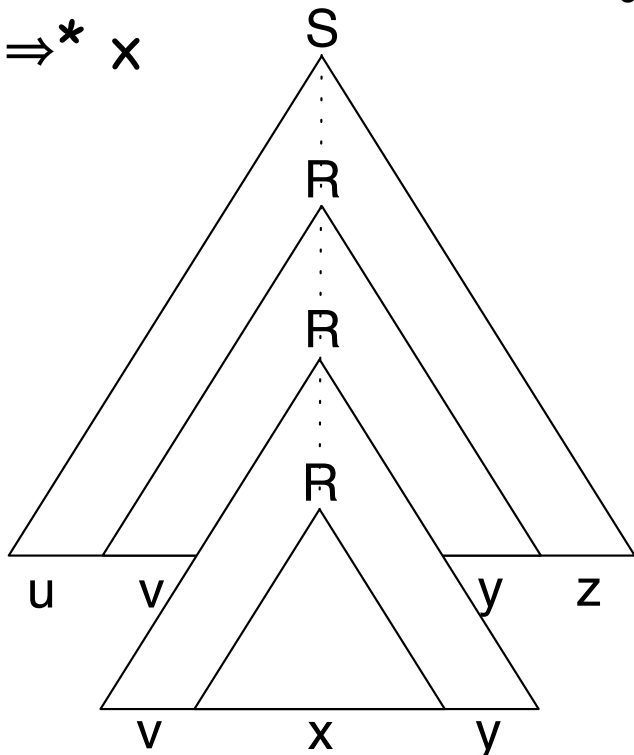
$P \Rightarrow Q$

$Q \Rightarrow R$

$R \Rightarrow^* x$



(2) Is there a point beyond which all strings have "tall" parse trees?



# Regular Grammars

- ✓ A grammar is said to be **regular** if its rules are of the following forms:

$$X \rightarrow a$$

$$X \rightarrow \varepsilon$$

$$X \rightarrow aY$$

# Example

$$S \rightarrow 1B$$

$$B \rightarrow 1B$$

$$B \rightarrow 0C$$

$$C \rightarrow 0S$$

$$S \rightarrow 0S$$

$$C \rightarrow 1B$$

$$C \rightarrow \varepsilon$$

# Context-Free Grammars

- ✓ An unrestricted grammar consists of
  1. A set  $V$  of **variables** (a.k.a. **nonterminals**)
  2. A disjoint set  $\Sigma$  (of **terminals**)
  3. A set of rules of the form  $\text{LEFT} \rightarrow \text{RIGHT}$  where  
 $\text{LEFT} \in (V \cup \Sigma)^+$  and  $\text{RIGHT} \in (V \cup \Sigma)^*$
  4. One designated  $S \in V$ , called the **start variable** (a.k.a. **start symbol**)
  
- ✓ If every  $\text{LEFT}$  is a single variable, the grammar is said to be **context-free**.

# Rewriting Strings

- ✓ If we have a rule  $\text{LEFT} \rightarrow \text{RIGHT}$ , we can replace  $\text{LEFT}$  by  $\text{RIGHT}$  inside any string:

$$\alpha\text{LEFT}\beta \Rightarrow \alpha\text{RIGHT}\beta$$

- ✓ The language of a grammar  $G$  is the set

$$L(G) = \{ w \in \Sigma^* \mid S \Rightarrow^* w \}$$

- ✓ The same language  $L$  might be the language of many different grammars.
  - ✓  $L$  is said to be a **context-free language** if it can be generated by **at least one** context-free grammar.

# CFG Example

$$S \rightarrow \varepsilon$$

$$S \rightarrow 0B$$

$$S \rightarrow 1A$$

$$A \rightarrow 0S$$

$$A \rightarrow 1AA$$

$$B \rightarrow 1S$$

$$B \rightarrow 0BB$$

- ✓ What is the start symbol?
- ✓ What are the terminals and nonterminals?
- ✓ What strings does this grammar generate?

# The Parsing Problem

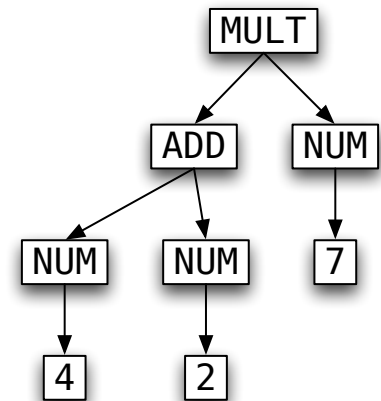
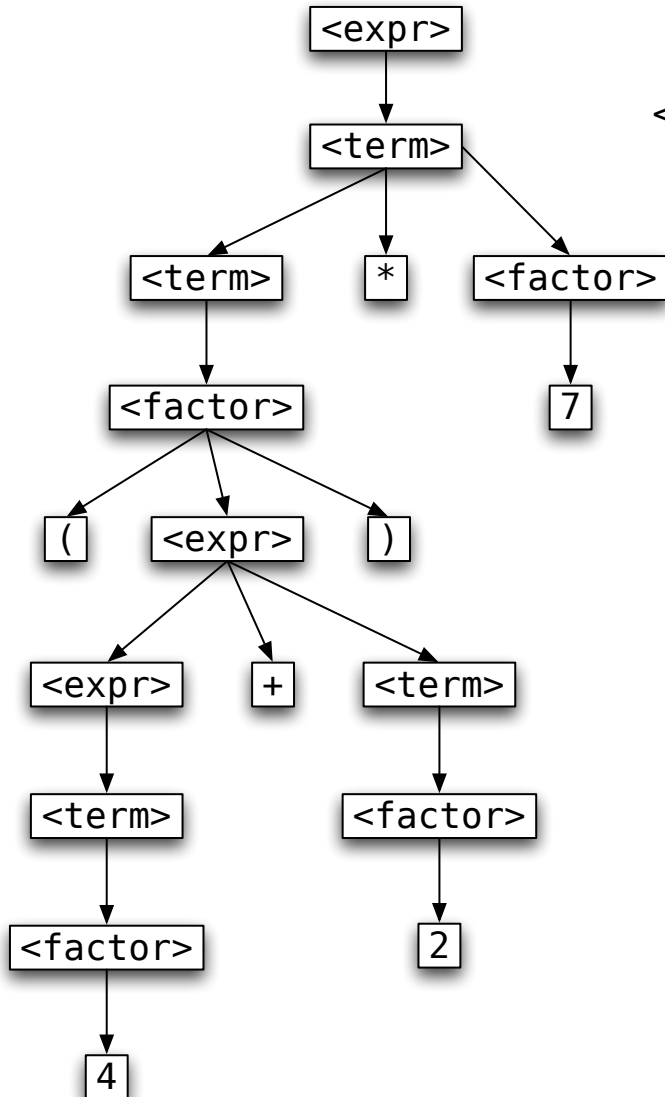
- ✓ Given a grammar and a string
  - ✓ Is the string in the language.
  - ✓ Usually, produce a “proof”
    - ✓ Parse tree or some abstraction thereof.

# Parse Tree vs. AST

$\langle \text{expr} \rangle ::= \langle \text{term} \rangle$   
                  |  $\langle \text{expr} \rangle + \langle \text{term} \rangle$

$\langle \text{term} \rangle ::= \langle \text{factor} \rangle$   
                  |  $\langle \text{term} \rangle * \langle \text{factor} \rangle$

$\langle \text{factor} \rangle ::= \langle \text{int} \rangle$   
                  | (  $\langle \text{expr} \rangle$  )



# Naive Parsing

- ✓ Backtracking search.
  - ✓ Try all ways to derive the string.
- ✓ Inefficient.
- ✓ Harder to implement than you might think...

# Bogus Backtracking

$S \rightarrow Aa \mid Ba$   
 $A \rightarrow a \mid c \mid ac$   
 $B \rightarrow Bb \mid b$

Consume\_S():

try Consume\_A(), then consume a  
if this fails, try Consume\_B(), then consume a

Consume\_A():

try consume a  
if this fails, try consume c  
if this fails, try consume a then consume c

Consume\_B():

try Consume\_B(), then consume b  
if this fails, try consume b

# LL(k) Grammars

- ✓ If each Consume function “knew” the right choice
  - ✓ we’d never need to backtrack
  - ✓ we’d never get tangled in infinite loops
  - ✓ It would be easy to write correct Consume functions
  - ✓ Our parser would run in linear time.

$S \rightarrow aA \mid bB$

- ✓ We say that a grammar is **LL(k)** if, by “peeking ahead” no more than k tokens, we can guarantee a decision that is
  - ✓ correct
  - ✓ unique (unambiguous)
- ✓ A language is LL(k) if there exists at least one LL(k) grammar.

# LL(k) Grammars?

$$S \rightarrow E \$$$
$$E \rightarrow n$$
$$| \text{ plus } E E$$
$$| \text{ times } E E$$
$$S \rightarrow A$$
$$| B$$
$$A \rightarrow a$$
$$| x A$$
$$B \rightarrow b$$
$$| y B$$
$$S \rightarrow A$$
$$| B$$
$$A \rightarrow a$$
$$| x A$$
$$B \rightarrow b$$
$$| x B$$
$$S \rightarrow E \$$$
$$E \rightarrow n$$
$$| n + E$$
$$S \rightarrow E \$$$
$$E \rightarrow n$$
$$| E + n$$
$$S \rightarrow E \$$$
$$E \rightarrow E + E$$
$$| E * E$$
$$| n$$

# Recursive Descent

✓ Naive backtracking works for LL(k) grammars!

$$\begin{array}{l} S \rightarrow s \mid a B \$ \\ B \rightarrow d \mid c B B \end{array}$$

Consume\_S():

try to consume s

if this fails, consume a, Consume\_B(), then consume \$

Consume\_B():

try to consume d

if this fails, consume c, Consume\_B(), then Consume\_B()

# Predictive Parsing

- ✓ Generalization of Recursive Descent.
- ✓ Stack of predictions, initially S.
- ✓ At each step, predict or match.
- ✓ Succeed if we run out of input & predictions at the same time.

# Digression

- ✓ Q: If regular languages are recognized by finite-state machines, what abstract machines recognize the context-free languages?
- ✓ A: Pushdown Automata (PDAs)
  - ✓ Finite state machine + stack
  - ✓ Transitions
    - ✓ Depend on the state and input symbol and top of stack!
    - ✓ Changes state and removes/replaces/ top of stack.
  - ✓ Accepting states as before (or accept on empty stack)
  - ✓ In general, can be nondeterministic.

# Example

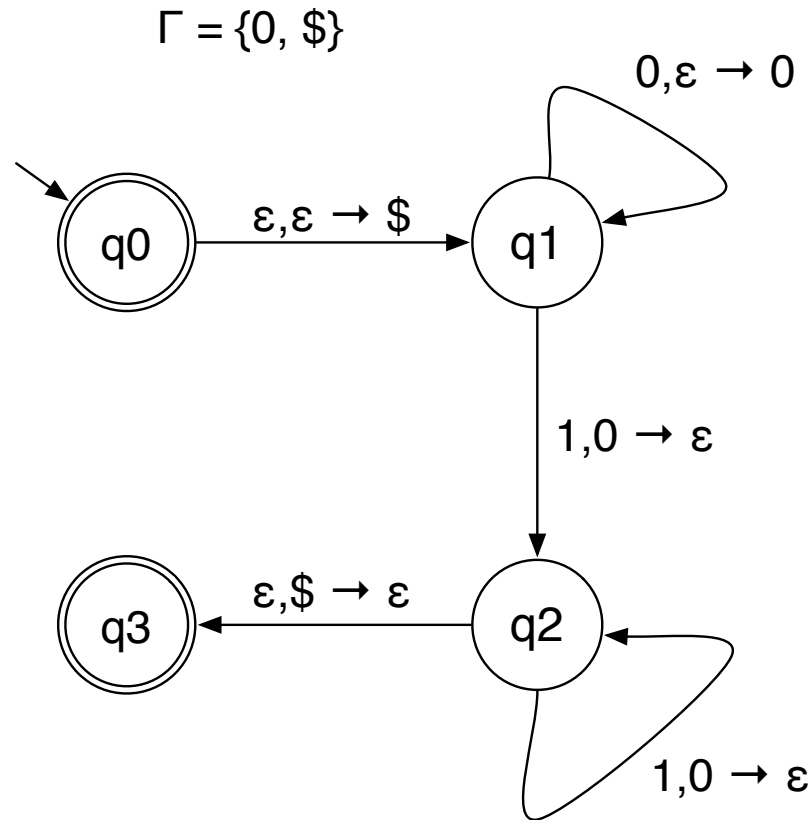
✓ Using a state machine and a stack, how can we recognize  $\{ 0^n 1^n \mid n \geq 0 \}$ ?

# Official Definition

- ✓ A PDA is a tuple  $(Q, \Sigma, \Gamma, q_0, F, \delta)$ 
  - ✓  $Q$  is a finite set of states
  - ✓  $\Sigma$  is a finite alphabet
  - ✓  $\Gamma$  is a finite “stack” alphabet
  - ✓  $q_0 \in Q$  is a start state
  - ✓  $F \subseteq Q$  is a set of accepting states
  - ✓  $\delta : Q \times (\Sigma \cup \{\varepsilon\}) \times (\Gamma \cup \{\varepsilon\}) \rightarrow \mathcal{P}(Q \times (\Gamma \cup \{\varepsilon\}))$

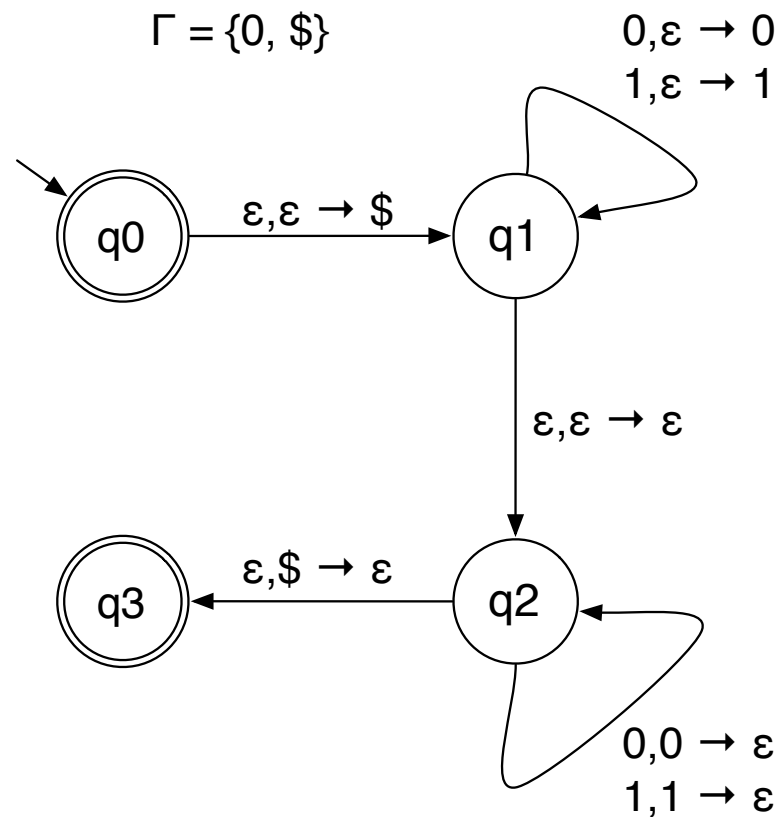
# Example

- ✓ Using a state machine and a stack, how can we recognize  $\{ 0^n 1^n \mid n \geq 0 \}$ ?



# Example

✓ Using a state machine and a stack, how can we recognize  $\{ ww^R \mid w \in \Sigma^* \}$ ?



# PDA vs CFG

- ✓ PDAs recognize all context-free languages.
  - ✓ Given a grammar, construct a PDA to do predict-match parsing
  - ✓ Use nondeterminism to **always** guess the correct prediction!  
(No backtracking officially required)
- ✓ PDAs recognize only context-free languages.
  - ✓ Turn the PDA into a grammar that simulates it.
  - ✓ See the book for details.
  - ✓ Basically, for each pair of states  $(p,q)$ , the nonterminal  $A_{pq}$  produces strings that get you from  $p$  to  $q$  starting and ending with an empty stack.