

# Parsing and PDAs

April 6, 2011

CS 81: Computability and Logic

# REVIEW

- ✓ Grammars
- ✓ Context-Free Grammars
- ✓ Ambiguous and Unambiguous Grammars
- ✓ A Pumping Lemma for CFLs ( $uvxyz$ )
  - ▶  $\{a^n b^n c^n \mid n \geq 0\}$  is not context free.
  - ▶  $\{ww \mid w \in \{a, b\}^*\}$  is not context free

# CHOMSKY HIERARCHY

- ✓ Type 0: Unrestricted Grammars
- ✓ Type 1: Context-Sensitive Grammars
- ✓ Type 2: Context-Free Grammars
- ✓ Type 3: Regular Grammars

## BY POPULAR DEMAND...

What is the language of the following  
*context-sensitive (monotone) grammar?*

$$S \rightarrow abc$$

$$S \rightarrow aSQ$$

$$bQc \rightarrow bbcc$$

$$cQ \rightarrow Qc$$

## STRINGS WITH COMPOSITE (NON-PRIME) LENGTHS

 $S \rightarrow HaC$  $C \rightarrow aC$  $C \rightarrow aT$  $T \rightarrow ZU$  $AZ \rightarrow ZA$  $aZ \rightarrow ZAa$  $HZ \rightarrow HY$  $YA \rightarrow AY$  $Ya \rightarrow aY$  $YU \rightarrow T$  $T \rightarrow X$  $aX \rightarrow Xaa$  $AX \rightarrow Xa$  $HX \rightarrow \varepsilon$

## THE PARSING PROBLEM

Given a grammar and a string,

1. Is the string in the language?

## THE PARSING PROBLEM

Given a grammar and a string,

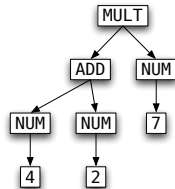
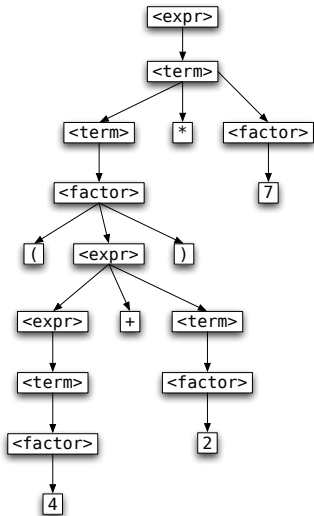
1. Is the string in the language?
2. Why? (Parse tree or other evidence)

## EVIDENCE

```
<expr> ::= <term>
         | <expr> + <term>
```

```
<term> ::= <factor>
         | <term> * <factor>
```

```
<factor> ::= <int>
          | ( <expr> )
```



# NAIVE PARSING

How about backtracking search? All ways to derive the string from  $S$ .

# NAIVE PARSING

How about backtracking search? All ways to derive the string from  $S$ .

- ✓ Inefficient
- ✓ Harder to implement than you might think...

# BOGUS BACKTRACKING

---

$S \rightarrow Ac \mid Bc$

$A \rightarrow a \mid c \mid ac$

$B \rightarrow Bb \mid b$

---

Consume\_S():

try Consume\_A(), then consume c

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

# BOGUS BACKTRACKING

---

S -> Ac | Bc

A -> a | c | ac

B -> Bb | b

---

Consume\_S():

try Consume\_A(), then consume c

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

Consume\_A():

try consume a

if fails, try consume c

if fails, try consume a then consume c

[What's wrong?]

# BOGUS BACKTRACKING

---

S -> Ac | Bc

A -> a | c | ac

B -> Bb | b

---

Consume\_S():

try Consume\_A(), then consume c

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

Consume\_A():

try consume a

if fails, try consume c

if fails, try consume a then consume c

[What's wrong?]

Consume\_B():

try Consume\_B(), then consume b

if fails, try consume b

[What's wrong?]

# LL(k) GRAMMARS

$S \rightarrow Aa \mid Ba$

If each `Consume` function always “knew” which right-hand-side was correct, **we would never need to backtrack or and never get tangled in infinite loops.**

# LL(k) GRAMMARS

$S \rightarrow Aa \mid Ba$

If each `Consume` function always “knew” which right-hand-side was correct, **we would never need to backtrack or and never get tangled in infinite loops.**

Then

- ✓ It would be easy to write correct `Consume` functions
- ✓ Our parser would run in linear time.

# LL(k) GRAMMARS

$S \rightarrow Aa \mid Ba$

If each **Consume** function always “knew” which right-hand-side was correct, **we would never need to backtrack or and never get tangled in infinite loops.**

Then

- ✓ It would be easy to write correct **Consume** functions
- ✓ Our parser would run in linear time.

We say that a grammar is **LL(k)** if, by “peeking ahead” no more than **k** tokens, we can guarantee a decision that is

1. correct
2. unique

# WHEN WILL TOP-DOWN PARSING WORK?

Are these grammars LL(k) for some k?

$$\begin{aligned} S &::= E \$ \\ E &::= n \\ &\quad | \text{ plus } E E \\ &\quad | \text{ times } E E \end{aligned}$$


---


$$\begin{aligned} S &::= A \\ &\quad | B \\ A &::= a \\ &\quad | x A \\ B &::= b \\ &\quad | y B \end{aligned}$$

$$\begin{aligned} S &::= A \\ &\quad | B \\ A &::= a \\ &\quad | x A \\ B &::= b \\ &\quad | x B \end{aligned}$$


---


$$\begin{aligned} S &::= E \$ \\ E &::= n \\ &\quad | n + E \end{aligned}$$

$$\begin{aligned} S &::= E \$ \\ E &::= n \\ &\quad | E + n \end{aligned}$$


---


$$\begin{aligned} S &::= E \$ \\ E &::= E + E \\ &\quad | E * E \\ &\quad | n \end{aligned}$$

# MACHINES

Question: If FSMs recognize regular languages, what machines recognize CFLs?

# MACHINES

Question: If FSMs recognize regular languages, what machines recognize CFLs?

Answer: Pushdown Automata (PDAs)

- ✓ Finite state machine + stack
- ✓ Transitions
  - ▶ Depend on the state and/or input symbol
  - ▶ Change state + add/remove/replace top-of-stack
- ✓ In general, can be nondeterministic.

## EXAMPLE

Using a state machine and a stack, how could we recognize

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

## OFFICIAL DEFINITION

A PDA consists of

- ✓ A finite set  $Q$  of states
- ✓ A finite alphabet  $\Sigma$
- ✓ A finite “stack alphabet”  $\Gamma$
- ✓ A start state  $q_0 \in Q$
- ✓ Accepting states  $F \subseteq Q$
- ✓ Transitions

$$\delta : Q \times (\Sigma \cup \{\varepsilon\}) \times (\Gamma \cup \{\varepsilon\}) \rightarrow \mathcal{P}(Q \times (\Gamma \cup \{\varepsilon\}))$$

## PDA EXAMPLE

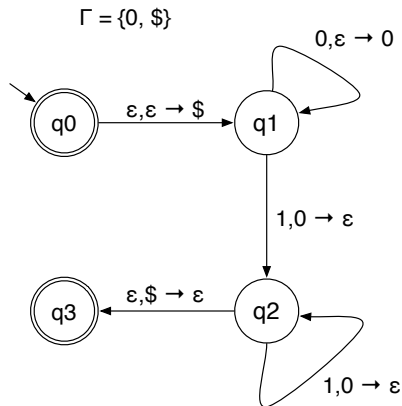
Using a state machine and a stack, how could we recognize

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

# PDA EXAMPLE

Using a state machine and a stack, how could we recognize

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



## PDA EXAMPLE

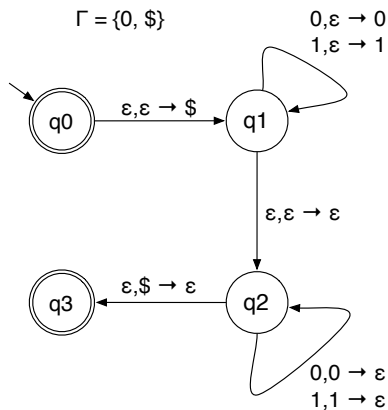
Using a state machine and a stack, how could we recognize

$$\{ ww^R \mid w \in \{0, 1\}^* \}$$

# PDA EXAMPLE

Using a state machine and a stack, how could we recognize

$$\{ ww^R \mid w \in \{0, 1\}^* \}$$



## PDA's vs CFGs: SUMMARY

### 1. PDAs can recognize any CFL

- ▶ Given a grammar, construct a PDA for top-down parsing
- ▶ Use nondeterminism to always guess the correct prediction  
(No backtracking required)

### 2. PDAs recognize only CFLs

- ▶ Turn a PDA into a grammar that simulates it
- ▶ See the book: for each pair of states  $\langle p, q \rangle$ , define a nonterminal  $A_{pq}$  (strings that get you from  $p$  to  $q$ , starting and ending with empty stack).

## CLOSURE PROPERTIES

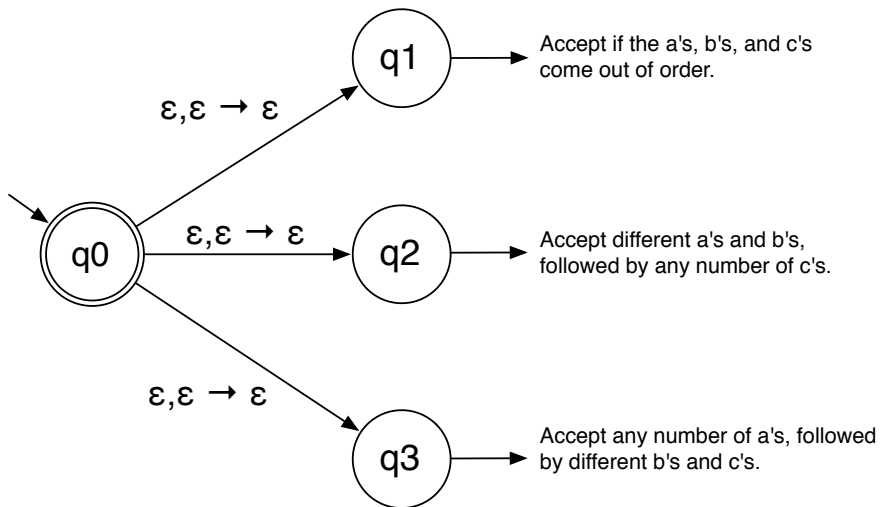
Context-Free Languages are:

- ✓ Closed under star, union, concatenation.
- ✓ Not closed under intersection, complement.

Key issue: nondeterministic PDAs are strictly more powerful than deterministic PDAs!

(No analogy to the subset construction)

$$\Sigma^* \setminus \{ a^n b^n c^n \mid n \geq 0 \}$$



# PRACTICAL PARSING

## Recursive Descent

- ✓ Form of code follows the grammar.
- ✓ Efficient and correct for **LL** grammars.

## Another very practical method: Shift-Reduce Parsing

- ✓ Efficient and correct for **LR** grammars
- ✓ Parser code is usually computer-generated!

But neither of these handle ALL CFGs...

# CYK (CKY) ALGORITHM

- ✓ Works for any CFG.
- ✓ Parses inputs in  $O(n^3)$  ( $n$  = length of input)
- ✓ Requires a grammar in “Chomsky Normal Form”
- ✓ Key ideas:
  - ▶ For every substring, what nonterminals produce it?
  - ▶ Dynamic programming for efficiency