

# MA/CSSE 474

## Theory of Computation

### Context-Free Grammars

### Context-Free Grammars

A **context-free grammar** (a.k.a. CFG)  $G$  is a quadruple,  $(V, \Sigma, R, S)$ , where:

- $V$  is the **rule alphabet (vocabulary)**, which contains nonterminals and terminals.
- $\Sigma$  (the set of **terminals**) is a subset of  $V$ ,
- $R$  (the set of **rules**) is a finite subset of  $(V - \Sigma) \times V^*$ ,
- $S$  (the **start symbol**) is an element of  $V - \Sigma$ .

**Example:**

$(\{S, a, b\}, \{a, b\}, \{S \rightarrow aSb, S \rightarrow \epsilon\}, S)$

Rules are also called **productions**.

**Note:** Some authors say that a CFG is  $(N, \Sigma, R, S)$ , where  $N$  is the set of nonterminal symbols. In that case  $V = N \cup \Sigma$ . I may sometimes use  $N$  in this way.  $N = V - \Sigma$ .

Q1-2

## Derivations

$x \Rightarrow_G y$  iff  $x = \alpha A \beta$

$\downarrow$  and  $A \rightarrow \gamma$  is in  $R$   
 $y = \alpha \gamma \beta$

$w_0 \Rightarrow_G w_1 \Rightarrow_G w_2 \Rightarrow_G \dots \Rightarrow_G w_n$  is a **derivation** in  $G$ .

Let  $\Rightarrow_G^*$  be the reflexive, transitive closure of  $\Rightarrow_G$ .

Then the **language generated by  $G$** , denoted  $L(G)$ , is:

$$\{w \in \Sigma^* : S \Rightarrow_G^* w\}.$$

A language  $L$  is **context-free** if there is some context-free grammar  $G$  such that  $L = L(G)$ .

Q3

## An Example Derivation

Example:

Let  $G = (\{S, a, b\}, \{a, b\}, \{S \rightarrow a S b, S \rightarrow \epsilon\}, S)$

$S \Rightarrow a S b \Rightarrow aa S bb \Rightarrow aaa S bbb \Rightarrow aaabbb$

So we can write  $S \Rightarrow^* aaabbb$

## Regular Grammars

A brief side-trip into Chapter 7

## Regular Grammars

In a regular grammar, every rule in  $R$  must have a right-hand side that is:

- $\epsilon$ , or
- a single terminal, or
- a single terminal followed by a single nonterminal.

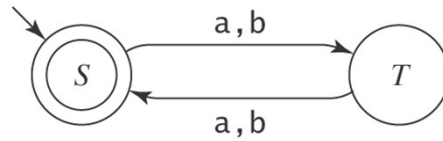
Regular:  $S \rightarrow a$ ,  $S \rightarrow \epsilon$ , and  $T \rightarrow aS$

Not regular:  $S \rightarrow aSa$  and  $aSa \rightarrow T$

Q4

## Regular Grammar Example

$$L = \{w \in \{a, b\}^* : |w| \text{ is even}\} \quad ((aa) \cup (ab) \cup (ba) \cup (bb))^*$$



$S \rightarrow \epsilon$   
 $S \rightarrow aT$   
 $S \rightarrow bT$   
 $T \rightarrow a$   
 $T \rightarrow b$   
 $T \rightarrow aS$   
 $T \rightarrow bS$

## Regular Languages and Regular Grammars

**Theorem:** A language is regular iff it can be defined by a regular grammar.

**Proof:** By two constructions (first one on next slide)

## Regular Languages and Regular Grammars

**Regular grammar**  $\rightarrow$  **FSM:**

$\text{grammartofsm}(G = (V, \Sigma, R, S)) =$

1. Create in  $M$  a separate state for each nonterminal in  $V$ .
2. Start state is the state corresponding to  $S$ .
3. If there are any rules in  $R$  of the form  $X \rightarrow w$ , for some  $w \in \Sigma$ , create a new state labeled #.
4. For each rule of the form  $X \rightarrow w Y$ , add a transition from  $X$  to  $Y$  labeled  $w$ .
5. For each rule of the form  $X \rightarrow w$ , add a transition from  $X$  to # labeled  $w$ .
6. For each rule of the form  $X \rightarrow \epsilon$ , mark state  $X$  as accepting.
7. Mark state # as accepting.

**FSM**  $\rightarrow$  **Regular grammar:** Similar.

## Example - Even Length Strings

$S \rightarrow \epsilon$	$T \rightarrow a$
$S \rightarrow aT$	$T \rightarrow b$
$S \rightarrow bT$	$T \rightarrow aS$
	$T \rightarrow bS$

Construct the FSM

End of regular grammars side-trip

Q5

## Recursive Grammar Rules

- A rule is **recursive** iff it is  $X \rightarrow w_1 Y w_2$ , where:  
 $Y \Rightarrow^* w_3 X w_4$  for some  $w_1, w_2, w_3$ , and  $w_4$  in  $V^*$ .
- A grammar  $G$  is recursive iff  $G$  contains at least one recursive rule.

- Examples:**

$S \rightarrow (S)$	$S \rightarrow (T)$
	$T \rightarrow (S)$

**In general, non-recursive grammars  
are boring!**

## Self-Embedding Grammar Rules

- A rule in a grammar  $G$  is **self-embedding** iff it is :  
 $X \rightarrow w_1 Y w_2$ , where  $Y \Rightarrow^* w_3 X w_4$  and  
 both  $w_1 w_3$  and  $w_4 w_2$  are in  $\Sigma^+$ .
- A grammar is self-embedding iff it contains at least one self-embedding rule.

- Examples:  $S \rightarrow a S a$     self-embedding

$S \rightarrow a S$     recursive but not self-embedding

$S \rightarrow a T$

$T \rightarrow S a$     self-embedding

What is the difference between self-embedding and recursive?

**Q6**

## Where Context-Free Grammars Get Their Power

- If a grammar  $G$  is not self-embedding then  $L(G)$  is regular.
- If a language  $L$  has the property that every grammar that defines it is self-embedding, then  $L$  is not regular.

## Equal Numbers of a's and b's

Let  $L = \{w \in \{a, b\}^* : \#_a(w) = \#_b(w)\}$ .

Find a CFG  $G$  such that  $L = L(G)$

Q7

## Arithmetic Expressions

$G = (V, \Sigma, R, E)$ , where  
 $V = \{+, *, (, ), id, E\}$ ,  
 $\Sigma = \{+, *, (, ), id\}$ ,  
 $R = \{$   
 $\quad E \rightarrow E + E$   
 $\quad E \rightarrow E * E$   
 $\quad E \rightarrow (E)$   
 $\quad E \rightarrow id$   
 $\quad \}$

## BNF

A notation for writing practical context-free grammars

- The symbol  $|$  should be read as “or”.

Example:  $S \rightarrow aSb \mid bSa \mid SS \mid \epsilon$

- Allow a nonterminal symbol to be any sequence of characters surrounded by angle brackets.

Examples of nonterminals:

$\langle \text{program} \rangle$   
 $\langle \text{variable} \rangle$



## BNF for a Java Fragment

```

<block>      ::= {<stmt-list>} |
                {}

<stmt-list> ::= <stmt> |
                <stmt-list> <stmt>

<stmt>       ::= <block> |
                while (<cond>) <stmt> |
                if (<cond>) <stmt> |
                do <stmt> while (<cond>); |
                <assignment-stmt>; |
                return |
                return <expression> |
                <method-invocation>;
  
```

## Spam Generation

```

<spc> → space | . | - | _ | = | : | * | / | | | empty

<Word> → <V> <spc> <l> <spc> <A> <spc> <G> <spc> <R> <spc> <A>
<V> → V | v | \
<l> → I | i | ! | j | : | i | í | ï | î | ï | í | ÿ | ŷ | ŷ | ! | l | l
<A> → A | a | / | \ | @ | ^ | Á | Â | Ã | Ä | Å | á | ä | å | à | á | ä
<G> → G | g | & | 6 | 9
<R> → R | r | ®
  
```

Example production:

```

<spc> → -
<V> → v
<l> ::= !
<A> ::= ä
<G> ::= G
<R> ::= ®
<A> ::= ^
<Word> → v-!-ä-G-®-^
  
```

These production rules yield 1,843,200 possible spellings.

How Many Ways Can You Spell V1@gra? By [Brian Hayes](#)

**American Scientist**, July-August 2007

<http://www.americanscientist.org/template/AssetDetail/assetid/55592>

## HTML

```
<ul>
  <li>Item 1, which will include a sublist</li>
    <ul>
      <li>First item in sublist</li>
      <li>Second item in sublist</li>
    </ul>
  <li>Item 2</li>
</ul>
```

A grammar:

*/\* Text is a sequence of elements.*

$HTMLtext \rightarrow Element\ HTMLtext \mid \epsilon$

$Element \rightarrow UL \mid LI \mid \dots$  (and other kinds of elements that are allowed in the body of an HTML document)

*/\* The <ul> and </ul> tags must match.*

$UL \rightarrow \langle ul \rangle HTMLtext \langle /ul \rangle$

*/\* The <li> and </li> tags must match.*

$LI \rightarrow \langle li \rangle HTMLtext \langle /li \rangle$

## English

$S \rightarrow NP\ VP$

$NP \rightarrow the\ Nominal \mid a\ Nominal \mid Nominal \mid$   
 $ProperNoun \mid NP\ PP$

$Nominal \rightarrow N \mid Adjs\ N$

$N \rightarrow cat \mid dogs \mid bear \mid girl \mid chocolate \mid rifle$

$ProperNoun \rightarrow Chris \mid Fluffy$

$Adjs \rightarrow Adj\ Adjs \mid Adj$

$Adj \rightarrow young \mid older \mid smart$

$VP \rightarrow V \mid V\ NP \mid VP\ PP$

$V \rightarrow like \mid likes \mid thinks \mid shoots \mid smells$

$PP \rightarrow Prep\ NP$

$Prep \rightarrow with$

## Designing Context-Free Grammars

- Generate related regions together.

$$A^n B^n$$

- Generate concatenated regions:

$$A \rightarrow BC$$

- Generate outside in:

$$A \rightarrow aAb$$

## Concatenating Independent Languages

Let  $L = \{a^n b^n c^m : n, m \geq 0\}$ .

The  $c^m$  portion of any string in  $L$  is completely independent of the  $a^n b^n$  portion, so we should generate the two portions separately and concatenate them together.

$G = (\{S, N, C, a, b, c\}, \{a, b, c\}, R, S)$  where:

$$R = \left\{ \begin{array}{l} S \rightarrow NC \\ N \rightarrow aNb \\ N \rightarrow \epsilon \\ C \rightarrow cC \\ C \rightarrow \epsilon \end{array} \right\}.$$

$$L = \{ a^{n_1} b^{n_1} a^{n_2} b^{n_2} \dots a^{n_k} b^{n_k} : k \geq 0 \text{ and } \forall i (n_i \geq 0) \}$$

Examples of strings in L:  $\epsilon$ , abab, aabbbaabbbabab

Note that  $L = \{a^n b^n : n \geq 0\}^*$ .

$G = (\{S, M, a, b\}, \{a, b\}, R, S)$  where:

$$R = \{ \begin{array}{l} S \rightarrow MS \\ S \rightarrow \epsilon \\ M \rightarrow aMb \\ M \rightarrow \epsilon \end{array} \}.$$

### Another Example: Unequal a's and b's

$$L = \{a^n b^m : n \neq m\}$$

$$G = (V, \Sigma, R, S), \text{ where}$$

$$V = \{a, b, S, \quad \},$$

$$\Sigma = \{a, b\},$$

$$R =$$

## Simplifying Context-Free Grammars

*Remove non-productive and unreachable non-terminals.*

## Unproductive Nonterminals

*removeunproductive*( $G$ : CFG) =

1.  $G' = G$ .
2. Mark every nonterminal symbol in  $G'$  as unproductive.
3. Mark every terminal symbol in  $G'$  as productive.
4. Until one entire pass has been made without any new symbol being marked do:
  - For each rule  $X \rightarrow \alpha$  in  $R$  do:
    - If every symbol in  $\alpha$  has been marked as productive and  $X$  has not yet been marked as productive then:
      - Mark  $X$  as productive.
5. Remove from  $G'$  every unproductive symbol.
6. Remove from  $G'$  every rule that contains an unproductive symbol.
7. Return  $G'$ .

## Unreachable Nonterminals

*removeunreachable*( $G$ : CFG) =

1.  $G' = G$ .
2. Mark  $S$  as reachable.
3. Mark every other nonterminal symbol as unreachable.
4. Until one entire pass has been made without any new symbol being marked do:
  - For each rule  $X \rightarrow \alpha A \beta$  (where  $A \in V - \Sigma$ ) in  $R$  do:
    - If  $X$  has been marked as reachable and  $A$  has not then:
      - Mark  $A$  as reachable.
5. Remove from  $G'$  every unreachable symbol.
6. Remove from  $G'$  every rule with an unreachable symbol on the left-hand side.
7. Return  $G'$ .

## Proving the Correctness of a Grammar

$$A^n B^n = \{a^n b^n : n \geq 0\}$$

$$G = (\{S, a, b\}, \{a, b\}, R, S),$$

$$R = \left\{ \begin{array}{l} S \rightarrow a S b \\ S \rightarrow \epsilon \end{array} \right\}$$

- Prove that  $G$  generates only strings in  $L$ .
- Prove that  $G$  generates all the strings in  $L$ .

## Structure

Context free languages:

We care about structure.

