# COMS W3261 – Lecture 5.2:

Context-Free Grammars.

__Idea__:  Introduce a new, more powerful way of describing languages

__Example__:

$$A \longrightarrow 0A1$$
$$A \longrightarrow B$$
$$B \longrightarrow \#$$

__Variables__: things that can be substituted $(A, B)$. Often written as capital letters.

__Terminals__: symbols in the final string, cannot be substituted. $(0, 1, \#)$.

How to generate a string.
1. Writing down the start variable (top left)
2. Replacing any variable using a substitution rule
3. Repeat step 2 until only terminals remain.

$$A \longrightarrow 0A1 \longrightarrow 00A11 \longrightarrow 00B11 \longrightarrow 00\#11.$$
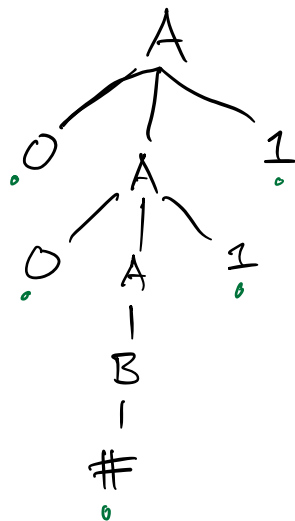
$$A \longrightarrow B \longrightarrow \#.$$

$$A \longrightarrow 0A1 \longrightarrow 0B1 \longrightarrow 0\#1.$$

__Def__. A sequence of substitutions used to create a string of terminals is called a __derivation__.

We can represent a derivation pictorially with a __parse tree__.
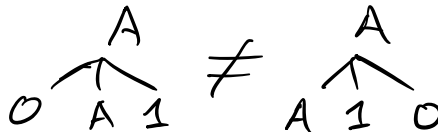
Ex. Parse tree for 00#11

A

O A 1

O A 1

B

#

= 00#11.

(Substituting each symbol according to the rule we use as we move downward.)

Read a parse tree by concatenating symbols left to right.

$A \to OA1$.

$$\underset{O \quad A \quad 1}{A} \neq \underset{A \quad 1 \quad O}{A}$$

Def. The language $L(G)$ of a grammar $G$ is the set of all strings that can be produced by derivation

$$G: \begin{array}{l} A \to OA1 \\ A \to B \\ B \to \# \end{array} \qquad L(G) = \{0^n \# 1^n \mid n \geq 0\}$$

Def. The set of all languages produced by a context-free (CFG) grammar is called the Context-Free Languages. (CFL)s.

Example: A fragment of English.

$$S \to \langle NP \rangle \langle VP \rangle$$

$$\langle NP \rangle \to A \, N$$

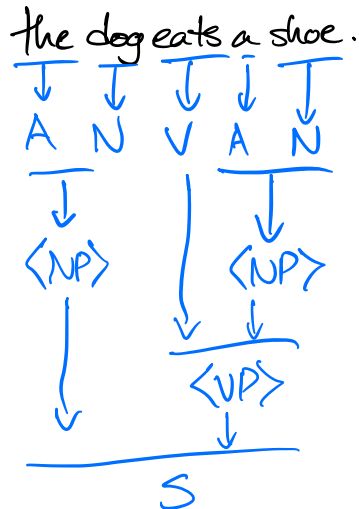// using $\langle \; \rangle$ to denote one single variable symbol.

$\langle VP \rangle \longrightarrow V \mid V \langle NP \rangle$

$N \longrightarrow$ dog $\mid$ cat $\mid$ car $\mid$ shoe $\mid$ person

$V \longrightarrow$ smells $\mid$ sees $\mid$ is $\mid$ eats

$A \longrightarrow$ a $\mid$ the

// the bar '|' abbreviates multiple rules as one

$N \rightarrow$ dog
$N \rightarrow$ cat
$N \rightarrow \cdots$

$S \longrightarrow \langle NP \rangle \langle VP \rangle \longrightarrow AN \langle VP \rangle \longrightarrow aN \langle VP \rangle$

$aNV \longrightarrow$ a shoe $V \longrightarrow$ a shoe is

$S \longrightarrow \langle NP \rangle \langle VP \rangle \longrightarrow \langle NP \rangle V \langle NP \rangle \longrightarrow ANV \langle NP \rangle \longrightarrow$

$ANVAN \longrightarrow \cdots \longrightarrow$ the cat sees the person.

the dog eats a shoe.



**Def.** (CFG, formally.) A context-free grammar is a 4-tuple, $(V, \Sigma, R, S)$, where:

$V$ is a finite set called the variables

$\Sigma$ is a finite set called the terminals (disjoint from $V$)

$R$ is a finite set of _rules_, where each rule maps 1 variable to a sequence of variables and terminals. (e.g., $A \longrightarrow 01A$)

$S \in V$ is the start variable.

For any strings of variables and terminals $u$, $v$, and $w$, if $A \rightarrow w$ is a rule of the grammar, then we say that $uAv$ yields $uwv$, where 'yields' is written $uAv \Rightarrow uvw$.

For any strings of variables + terminals $u$ and $v$, we say $u$ _derives_ $v$, written $u \overset{*}{\Rightarrow} v$, if $u = v$ or if there exists a sequence $u_1, u_2, \dots u_k$, $k \geq 0$, such that $u \Rightarrow u_1 \Rightarrow u_2 \cdots \Rightarrow u_k \Rightarrow v$.

(The language of a grammar $G$ is $\{w \in \Sigma^* \mid S \overset{*}{\Rightarrow} w\}$.)

_Example._ $G_4 = (V, \Sigma, R, \langle Expr \rangle)$ where

$V = \{\langle Expr \rangle, \langle Term \rangle, \langle Factor \rangle\}$,
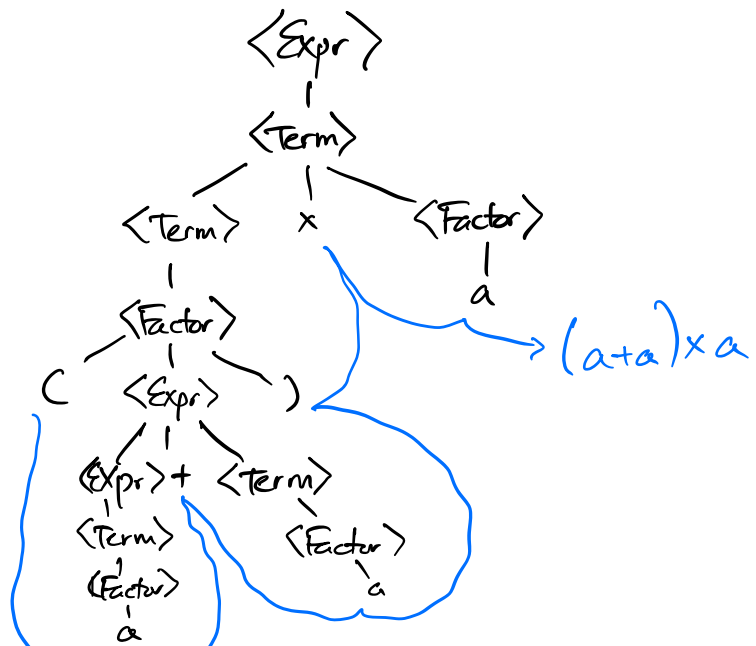
$\Sigma = \{a, +, \times, (, )\}$

"times symbol"   $a \times b.$

$R = \qquad \langle Expr \rangle \longrightarrow \langle Expr \rangle + \langle Term \rangle \mid \langle Term \rangle$

$\qquad \langle Term \rangle \longrightarrow \langle Term \rangle \times \langle Factor \rangle \mid \langle Factor \rangle$

$\qquad \langle Factor \rangle \longrightarrow (\langle Expr \rangle) \mid a$

$(a + a) \times a.$



$(a+a) \times a$

# Building CFGs — two techniques.

Suppose we want to build a CFG for $\{0^n 1^n \mid n \geq 0\}$

$$L = \cup \{1^n 0^n \mid n \geq 0\}$$

$$S_1 \to 0 S_1 1 \mid \varepsilon \qquad \text{// recognizes } \{0^n 1^n \mid n \geq 0\}$$

$$S_2 \to 1 S_2 0 \mid \varepsilon \qquad \text{// recognizes } \{1^n 0^n \mid n \geq 0\}$$
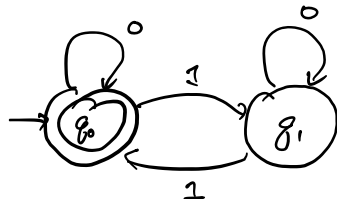
To take the 'union' of these CFGs, add a new start symbol.

$$G: \quad \begin{aligned} S &\to S_1 \mid S_2 \\ S_1 &\to 0 S_1 1 \mid \varepsilon \\ S_2 &\to 1 S_2 0 \mid \varepsilon \end{aligned} \qquad L(G) = L$$

(Aside: we made a CFG for a nonregular language$^{(!!)}$)

## Technique 2. Converting a DFA to a CFG.

Goal: CFG for $L = \{w \mid w \in \{0,1\}^*, w$ has an even number of ones$\}$.



1. Make a variable $R_i$ for each state $q_i$ of the DFA.
2. For each transition $\delta(q_i, a) = q_j$, add the rule $R_i \to a R_j$

$$\delta(q_0, 0) = 0$$

3. Add $R_i \to \varepsilon$ for each accept state.
4. $R_0$ is the start symbol.

($\Sigma$ is implicitly the same)

$$\Sigma = \Sigma$$
$$V = \{R_0, R_1\}$$
$$R = \begin{aligned} R_0 &\to 0 R_0 \quad * \\ R_0 &\to 1 R_1 \quad \bigstar \\ R_1 &\to 0 R_1 \quad - \\ R_1 &\to 1 R_0 \quad - \\ R_0 &\to \varepsilon \end{aligned}$$

0101 : On my DFA, I start in $q_0$ and compute to, $q_0, q_1, q_1, q_0$ and then accept.

In DFA: $R_0 \longrightarrow 0R_0 \longrightarrow 01R_1 \longrightarrow 010R_1 \longrightarrow 0101R_0$
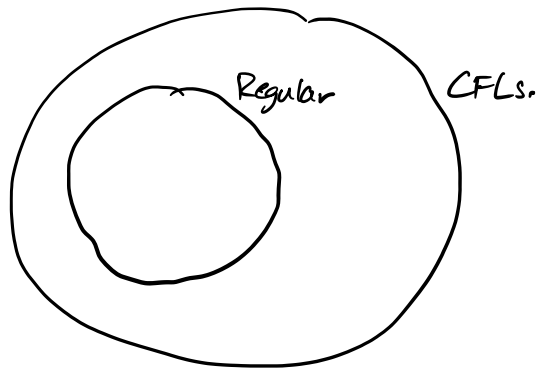$\longrightarrow 0101.$

Claim: this (non-rigorous) construction shows that any regular language derives from some CFG (!!
$$( \text{Reg. lang.} \longrightarrow \text{DFA} \longrightarrow \text{CFG.})$$

Regular Languages $\subseteq$ Context-Free Languages.

CFLs $\nsubseteq$ Reg. Languages.

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


Regular    CFLs.

**Def.** If a grammar generates the same string in ways corresponding to different parse trees, it is ambiguous.
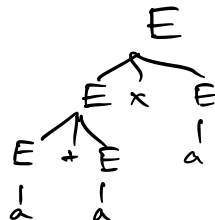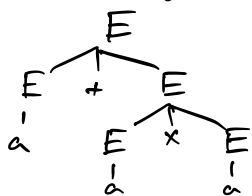
$\begin{pmatrix} S \to AA \\ A \to 1 \end{pmatrix}$

$S \to AA \to 1A \to 11$
$S \to AA \to A1 \to 11$

Example: $G_5 = (V, \Sigma, R, E)$,
where $V = \{E\}$, $\Sigma = \{+, \times, (,), a\}$,
$R$ as follows:

$$E \to E \times E \mid E + E \mid (E) \mid a$$

Now: Derivation of $a + a \times a$ ?

**Def.** A _leftmost_ derivation is one in which we always replace the _leftmost_ variable. Formally, a grammar is ambiguous if it generates a string with at least two different _leftmost_ derivations.

(leftmost #1) $E \rightarrow E+E \rightarrow a+E \rightarrow a+ E\times E \rightarrow a+a\times E \rightarrow a+a\times a$

(not leftmost) $E \rightarrow E+E \rightarrow E+E\times E \rightarrow a+E\times E \rightarrow a+a\times E \rightarrow a+a\times a$

(leftmost) #2 $E \rightarrow E\times E \rightarrow E+E\times E \rightarrow a+E\times E \rightarrow a+a\times E \rightarrow$
$a+a\times a.$

**Next time:** normal forms for grammars, pushdown automata!

Reminder: HW #3 due Monday @ 11:59 EST

Review: Section 2.1.