



CS311: Computational Theory

Lecture 6: CONTEXT-FREE GRAMMARS – Ch 2

Lecture Learning Objectives

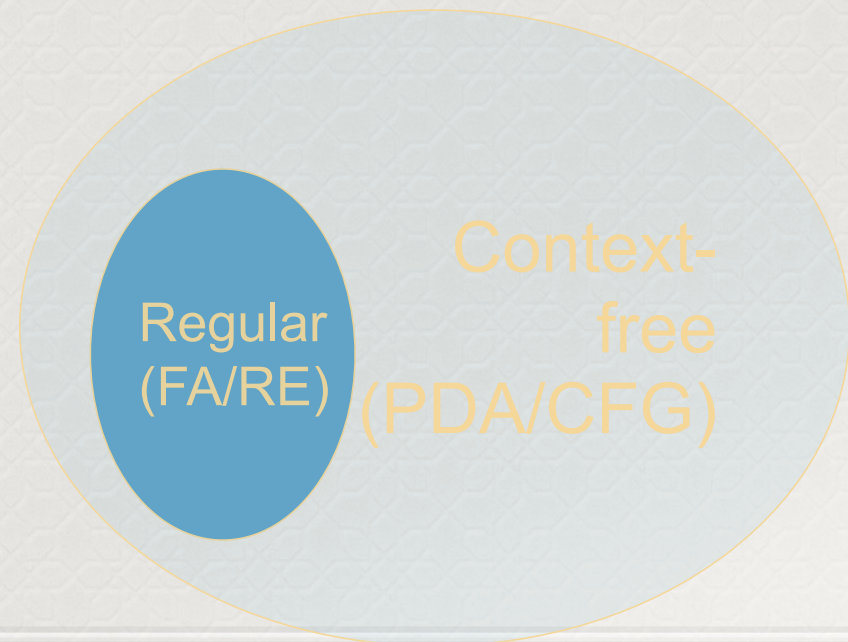
1. *Understand Context-Free Grammars and their applications.*

Not all languages are regular

- *So what happens to the languages which are not regular?*
- *Can we still come up with a language recognizer?*
 - *i.e., something that will accept (or reject) strings that belong (or do not belong) to the language?*

Context-Free Languages

- *A language class larger than the class of regular languages*
- *Supports natural, recursive notation called “context-free grammar”*
- *Applications:*
 - *Parse trees, compilers*
 - *XML*



An Example

- A palindrome is a word that reads identical from both ends
 - E.g., madam, redivider, malayalam, 010010010
- Let $L = \{ w \mid w \text{ is a binary palindrome} \}$
- Is L regular?
 - No.
 - Proof:
 - Let $w = 0^N 1 0^N$ (assuming N to be the p/l constant)
 - By Pumping lemma, w can be rewritten as xyz , such that xy^kz is also L (for any $k \geq 0$)
 - But $|xy| \leq N$ and $y \neq \epsilon$
 - $\implies y = 0^+$
 - $\implies xy^kz$ will NOT be in L for $k=0$
 - \implies Contradiction

But the language of palindromes...

is a CFL, because it supports recursive substitution (in the form of a CFG)

- This is because we can construct a “grammar” like this:

1. $A \Rightarrow \epsilon$
2. $A \Rightarrow 0$
3. $A \Rightarrow 1$
4. $A \Rightarrow 0A0$
5. $A \Rightarrow 1A1$

Productions

Terminal

Variable or non-terminal

Same as:

$A \Rightarrow 0A0 \mid 1A1 \mid 0 \mid 1 \mid \epsilon$

How does this grammar work?

How does the CFG for palindromes work?

An input string belongs to the language (i.e., accepted) iff it can be generated by the CFG

- Example: $w=01110$
- G can generate w as follows:

1. $A \Rightarrow 0A0$
2. $\Rightarrow 01A10$
3. $\Rightarrow 01110$

$G:$
 $A \Rightarrow 0A0 \mid 1A1 \mid 0 \mid 1 \mid \epsilon$

Generating a string from a grammar:

1. Pick and choose a sequence of productions that would allow us to generate the string.
2. At every step, substitute one variable with one of its productions.

Context-Free Grammar: Definition

- A context-free grammar $G=(V,\Sigma,R,S)$, where:
 - V : set of variables or non-terminals
 - T : set of terminals (= alphabet $\cup \{\epsilon\}$)
 - P : set of productions, each of which is of the form
$$V \implies \alpha_1 \mid \alpha_2 \mid \dots$$
 - Where each α_i is an arbitrary string of variables and terminals
 - $S \implies$ start variable

CFG for the language of binary palindromes:

$G=(\{A\},\{0,1\},P,A)$

$P: A \implies 0A0 \mid 1A1 \mid 0 \mid 1 \mid \epsilon$

More examples

- *Parenthesis matching in code*
- *Syntax checking*
- *In scenarios where there is a general need for:*
 - *Matching a symbol with another symbol, or*
 - *Matching a count of one symbol with that of another symbol, or*
 - *Recursively substituting one symbol with a string of other symbols*

Example #2

- *Language of balanced paranthesis*
e.g., ()(((())((())))....
- *CFG?*

```
G:  
S => (S) | SS | ε
```

How would you “interpret” the string “(((())(()))” using this grammar?

Example #3

- A grammar for $L = \{0^m 1^n \mid m \geq n\}$
- CFG?

```
G:  
S => 0S1 | A  
A => 0A | ε
```

How would you interpret the string "00000111" using this grammar?

Example #4

A program containing if-then(-else) statements

if Condition then Statement else Statement

(Or)

if Condition then Statement

CFG?

Example #5

$\langle \text{SENTENCE} \rangle \rightarrow \langle \text{NOUN-PHRASE} \rangle \langle \text{VERB-PHRASE} \rangle$
 $\langle \text{NOUN-PHRASE} \rangle \rightarrow \langle \text{CMPLX-NOUN} \rangle \mid \langle \text{CMPLX-NOUN} \rangle \langle \text{PREP-PHRASE} \rangle$
 $\langle \text{VERB-PHRASE} \rangle \rightarrow \langle \text{CMPLX-VERB} \rangle \mid \langle \text{CMPLX-VERB} \rangle \langle \text{PREP-PHRASE} \rangle$
 $\langle \text{PREP-PHRASE} \rangle \rightarrow \langle \text{PREP} \rangle \langle \text{CMPLX-NOUN} \rangle$
 $\langle \text{CMPLX-NOUN} \rangle \rightarrow \langle \text{ARTICLE} \rangle \langle \text{NOUN} \rangle$
 $\langle \text{CMPLX-VERB} \rangle \rightarrow \langle \text{VERB} \rangle \mid \langle \text{VERB} \rangle \langle \text{NOUN-PHRASE} \rangle$
 $\langle \text{ARTICLE} \rangle \rightarrow \text{a} \mid \text{the}$
 $\langle \text{NOUN} \rangle \rightarrow \text{boy} \mid \text{girl} \mid \text{flower}$
 $\langle \text{VERB} \rangle \rightarrow \text{touches} \mid \text{likes} \mid \text{sees}$
 $\langle \text{PREP} \rangle \rightarrow \text{with}$

Example #5 Derivation

a boy sees

<SENTENCE> ⇒ <NOUN-PHRASE><VERB-PHRASE>
⇒ <COMPLEX-NOUN><VERB-PHRASE>
⇒ <ARTICLE><NOUN><VERB-PHRASE>
⇒ a <NOUN><VERB-PHRASE>
⇒ a **by** <VERB-PHRASE>
⇒ a **by** <COMPLEX-VERB>
⇒ a **by** <VERB>
⇒ a **by** sees

More examples

- $L_1 = \{0^n \mid n \geq 0\}$
- $L_2 = \{0^n \mid n \geq 1\}$
- $L_3 = \{0^i 1^j 2^k \mid i=j \text{ or } j=k, \text{ where } i, j, k \geq 0\}$
- $L_4 = \{0^i 1^j 2^k \mid i=j \text{ or } i=k, \text{ where } i, j, k \geq 1\}$

Applications of CFLs & CFGs

- *Compilers use parsers for syntactic checking*
- *Parsers can be expressed as CFGs*
 1. *Balancing paranthesis:*
 - $B \Rightarrow BB \mid (B) \mid \text{Statement}$
 - $\text{Statement} \Rightarrow \dots$
 2. *If-then-else:*
 - $S \Rightarrow SS \mid \text{if Condition then Statement else Statement} \mid \text{if Condition then Statement} \mid \text{Statement}$
 - $\text{Condition} \Rightarrow \dots$
 - $\text{Statement} \Rightarrow \dots$
- *C paranthesis matching { ... }*
- *Pascal begin-end matching*
- *YACC (Yet Another Compiler-Compiler)*

More applications

- *Markup languages*
 - *Nested Tag Matching*
 - *HTML*
 - ◇ `<html> ...<p> </p> ... </html>`
 - *XML*
 - ◇ `<PC> ... <MODEL> ... </MODEL> .. <RAM> ... </RAM> ... </PC>`

Tag-Markup Languages

Roll \Rightarrow **<ROLL>** *Class Students* **</ROLL>**

Class \Rightarrow **<CLASS>** *Text* **</CLASS>**

Text \Rightarrow *Char Text* | *Char*

Char \Rightarrow **a** | **b** | ... | **z** | **A** | **B** | .. | **Z**

Students \Rightarrow *Student Students* | ϵ

Student \Rightarrow **<STUD>** *Text* **</STUD>**

Here, the left hand side of each production denotes one non-terminals (e.g., “Roll”, “Class”, etc.)

Those symbols on the right hand side for which no productions (i.e., substitutions) are defined are terminals (e.g., ‘a’, ‘b’, ‘|’, ‘<’, ‘>’, “ROLL”, etc.)

Structure of a production



The above is same as:

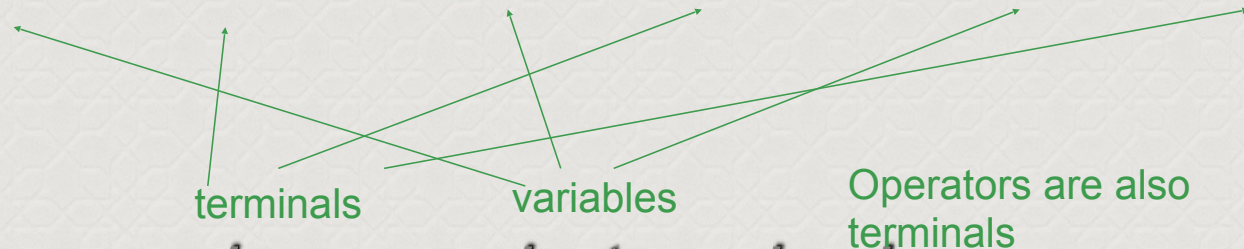
1. $A \Longrightarrow \alpha_1$
2. $A \Longrightarrow \alpha_2$
3. $A \Longrightarrow \alpha_3$
- ...
- K. $A \Longrightarrow \alpha_k$

CFG conventions

- *Terminal symbols $\Leftarrow a, b, c \dots$*
- *Non-terminal symbols $\Leftarrow A, B, C, \dots$*
- *Terminal or non-terminal symbols $\Leftarrow X, Y, Z$*
- *Terminal strings $\Leftarrow w, x, y, z$*
- *Arbitrary strings of terminals and non-terminals $\Leftarrow \alpha, \beta, \gamma, \dots$*

Syntactic Expressions in Programming Languages

*result = a*b + score + 10 * distance + c*



Regular languages have only terminals

- *Reg expression = $[a-z][a-z0-1]^*$*
- *If we allow only letters a & b, and 0 & 1 for constants (for simplification)*
 - *Regular expression = $(a+b)(a+b0+1)^*$*

String membership

How to say if a string belong to the language defined by a CFG?

1. *Derivation*

- *Head to body*

2. *Recursive inference*

- *Body to head*

Example:

- $w = 01110$
- *Is w a palindrome?*

Both are equivalent forms

G:
 $A \Rightarrow 0A0 \mid 1A1 \mid 0 \mid 1 \mid \epsilon$

$A \Rightarrow 0A0$
 $\Rightarrow 01A10$
 $\Rightarrow 01110$

Simple Expressions...

- We can write a CFG for accepting simple expressions
- $G = (V, \Sigma, R, S)$
 - $V = \{E, F\}$
 - $\Sigma = \{0, 1, a, b, +, *, (,)\}$
 - $S = \{E\}$
 - R :
 - $E \Rightarrow E + E \mid E * E \mid (E) \mid F$
 - $F \Rightarrow aF \mid bF \mid 0F \mid 1F \mid a \mid b \mid 0 \mid 1$

Generalization of derivation

- *Derivation is head \implies body*
- $A \implies X$ (*A derives X in a single step*)
- $A \implies^*_G X$ (*A derives X in a multiple steps*)
- *Transitivity:*
*IF $A \implies^*_G B$, and $B \implies^*_G C$, THEN $A \implies^*_G C$*

Context-Free Language

- *The language of a CFG, $G=(V,\Sigma,R,S)$, denoted by $L(G)$, is the set of terminal strings that have a derivation from the start variable S .*
 - $L(G) = \{ w \text{ in } T^* \mid S \Rightarrow^*_G w \}$

Left-most & Right-most Derivation Styles

G:

$E \Rightarrow E + E \mid E * E \mid (E) \mid F$
 $F \Rightarrow aF \mid bF \mid 0F \mid 1F \mid \varepsilon$

Derive the string $a^*(ab+10)$ from G:

$E \stackrel{*}{\Rightarrow}_G a^*(ab+10)$

Left-most derivation:

Always substitute leftmost variable

$\cdot E$
 $\Rightarrow E * E$
 $\Rightarrow F * E$
 $\Rightarrow a * E$
 $\Rightarrow a * (E)$
 $\Rightarrow a * (E + E)$
 $\Rightarrow a * (F + E)$
 $\Rightarrow a * (aF + E)$
 $\Rightarrow a * (abF + E)$
 $\Rightarrow a * (ab + E)$
 $\Rightarrow a * (ab + F)$
 $\Rightarrow a * (ab + 1F)$
 $\Rightarrow a * (ab + 10F)$
 $\Rightarrow a * (ab + 10)$

$\cdot E$
 $\Rightarrow E * E$
 $\Rightarrow E * (E)$
 $\Rightarrow E * (E + E)$
 $\Rightarrow E * (E + F)$
 $\Rightarrow E * (E + 1F)$
 $\Rightarrow E * (E + 10F)$
 $\Rightarrow E * (E + 10)$
 $\Rightarrow E * (F + 10)$
 $\Rightarrow E * (aF + 10)$
 $\Rightarrow E * (abF + 0)$
 $\Rightarrow E * (ab + 10)$
 $\Rightarrow F * (ab + 10)$
 $\Rightarrow aF * (ab + 10)$
 $\Rightarrow a * (ab + 10)$

Right-most derivation:

Always substitute rightmost variable

Leftmost vs. Rightmost derivations

Q1) *For every leftmost derivation, there is a rightmost derivation, and vice versa. True or False?*

True - will use parse trees to prove this

Q2) *Does every word generated by a CFG have a leftmost and a rightmost derivation?*

Yes – easy to prove (reverse direction)

Q3) *Could there be words which have more than one leftmost (or rightmost) derivation?*

Yes – depending on the grammar

How to prove that your CFGs are correct?

(using induction)

CFG & CFL

G_{pal}
 $A \Rightarrow 0A0 \mid 1A1 \mid 0 \mid 1 \mid \epsilon$

- Theorem: A string w in $(0+1)^*$ is in $L(G_{pal})$, if and only if, w is a palindrome.
- Proof:
 - Use induction
 - on string length for the IF part
 - On length of derivation for the ONLY IF part

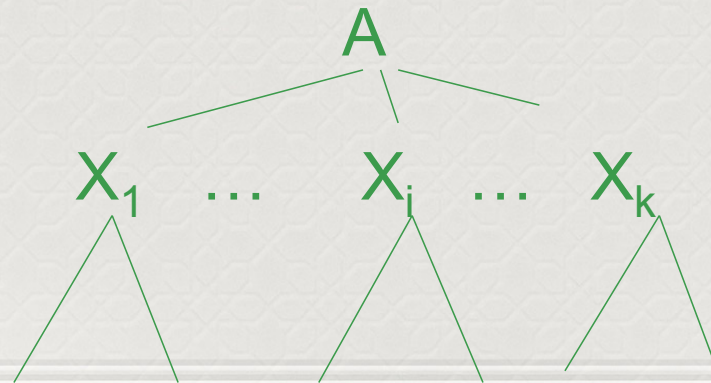
Parse trees

Parse Trees

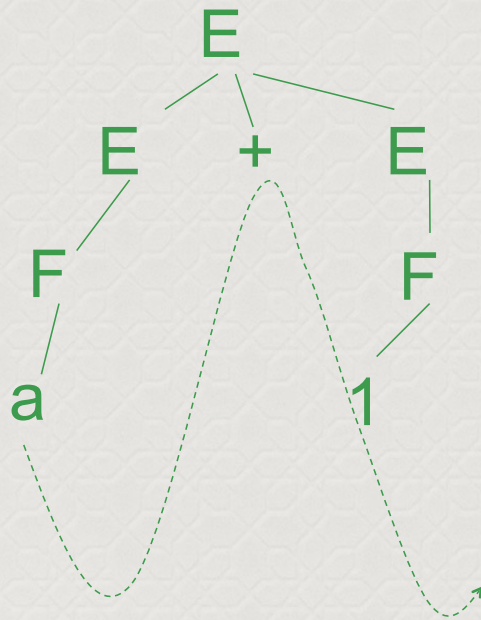
- Each CFG can be represented using a parse tree:
 - Each internal node is labeled by a variable in V
 - Each leaf is terminal symbol
 - For a production, $A \Rightarrow X_1 X_2 \dots X_k$, then any internal node labeled A has k children which are labeled from X_1, X_2, \dots, X_k from left to right

Parse tree for production and all other subsequent productions:

$A \Rightarrow X_1 \dots X_i \dots X_k$

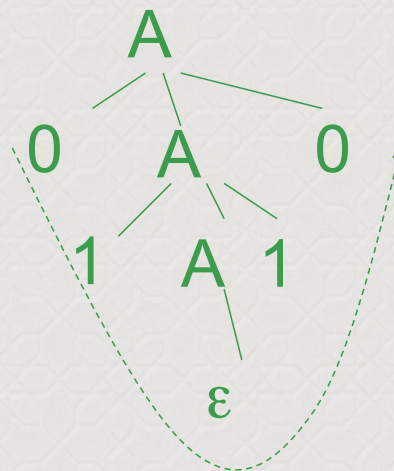


Examples



Parse tree for $a + 1$

Recursive inference



Parse tree for 0110

Derivation

G:

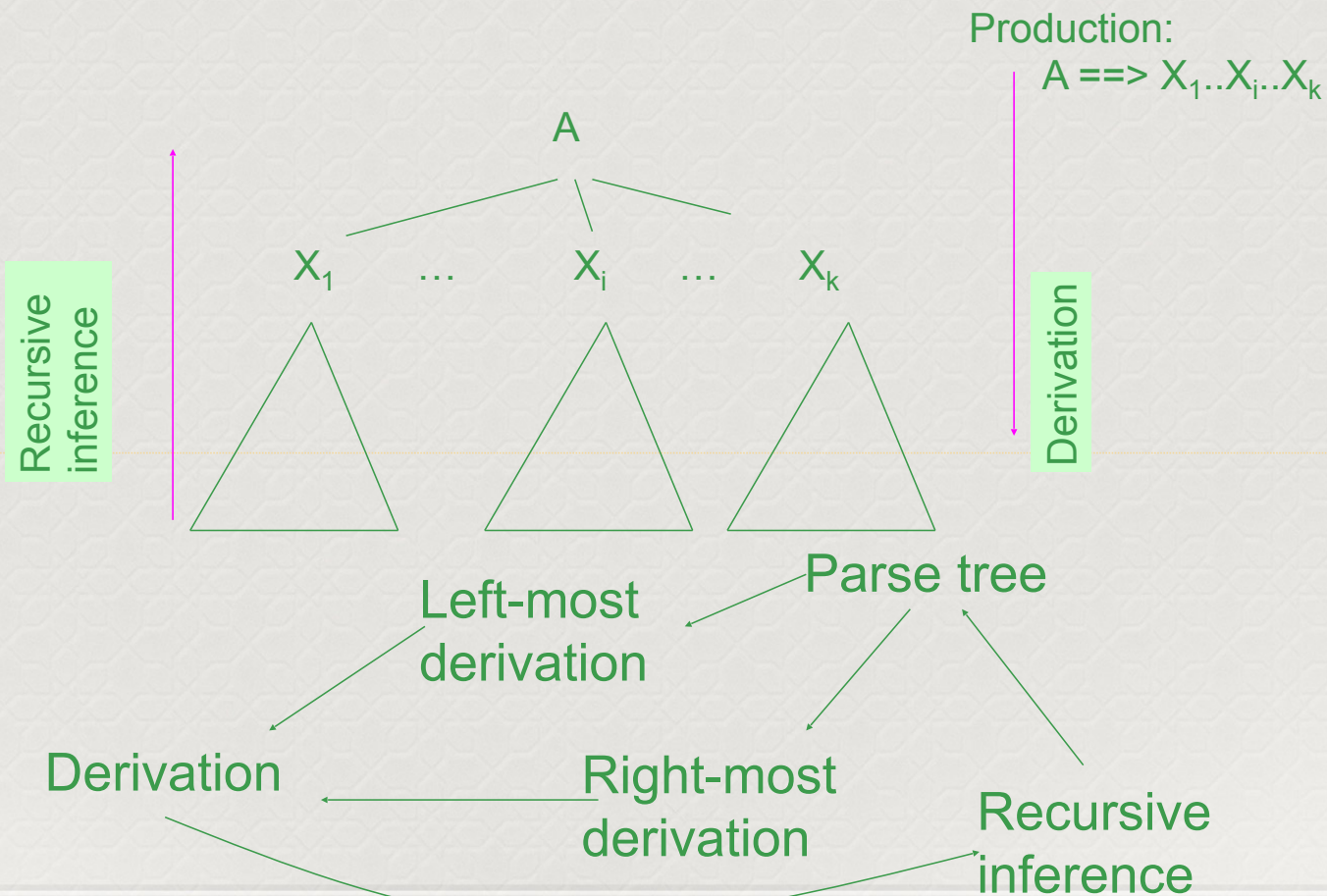
$E \Rightarrow E + E \mid E * E \mid (E) \mid F$

$F \Rightarrow aF \mid bF \mid 0F \mid 1F \mid 0 \mid 1 \mid a \mid b$

G:

$A \Rightarrow 0A0 \mid 1A1 \mid 0 \mid 1 \mid \epsilon$

Parse Trees, Derivations, and Recursive Inferences



Interchangeability of different CFG representations

- *Parse tree \implies left-most derivation*
 - *DFS left to right*
- *Parse tree \implies right-most derivation*
 - *DFS right to left*
- *\implies left-most derivation \iff right-most derivation*
- *Derivation \implies Recursive inference*
 - *Reverse the order of productions*
- *Recursive inference \implies Parse trees*
 - *bottom-up traversal of parse tree*



Ambiguity in CFGs and CFLs

Ambiguity in CFGs

- A CFG is said to be **ambiguous** if there exists a string which has more than one left-most derivation

Example:

$S \Rightarrow AS \mid \varepsilon$

$A \Rightarrow A1 \mid 0A1 \mid 01$

Input string: 00111

Can be derived in two ways

LM derivation #1:

$S \Rightarrow AS$
 $\Rightarrow 0A1S$
 $\Rightarrow 0A11S$
 $\Rightarrow 00111S$
 $\Rightarrow 00111$

LM derivation #2:

$S \Rightarrow AS$
 $\Rightarrow A1S$
 $\Rightarrow 0A11S$
 $\Rightarrow 00111S$
 $\Rightarrow 00111$

Why does ambiguity matter?

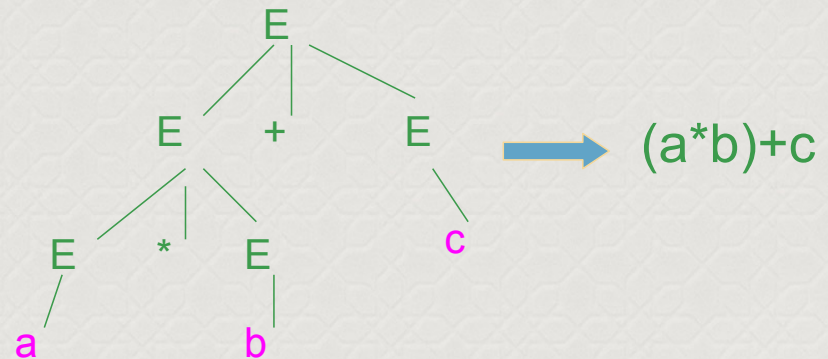
$E \Rightarrow E + E \mid E * E \mid (E) \mid a \mid b \mid c \mid 0 \mid 1$

Values are different !!!

*string = a * b + c*

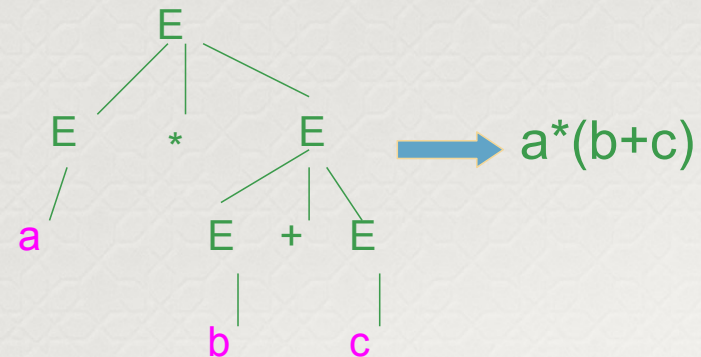
- LM derivation #1:

• $E \Rightarrow E + E \Rightarrow E * E + E$
 $\Rightarrow * a * b + c$



- LM derivation #2

• $E \Rightarrow E * E \Rightarrow a * E \Rightarrow$
 $a * E + E \Rightarrow * a * b + c$



The calculated value depends on which of the two parse trees is actually used.

Connection between CFLs and RLs

Removing Ambiguity in Expression Evaluations

- *It MAY be possible to remove ambiguity for some CFLs*
 - *E.g., in a CFG for expression evaluation by imposing rules & restrictions such as precedence*
 - *This would imply rewrite of the grammar*

• Precedence: $()$, $*$, $+$

Modified unambiguous version:

Ambiguous version:

$E \Rightarrow E + E \mid E * E \mid (E) \mid a \mid b \mid c \mid 0 \mid 1$

$E \Rightarrow E + T \mid T$
 $T \Rightarrow T * F \mid F$
 $F \Rightarrow I \mid (E)$
 $I \Rightarrow a \mid b \mid c \mid 0 \mid 1$

How will this avoid ambiguity?

Inherently Ambiguous CFLs

- However, for some languages, it may not be possible to remove ambiguity
- A CFL is said to be *inherently ambiguous* if every CFG that describes it is ambiguous

Example:

- $L = \{ a^n b^n c^m d^m \mid n, m \geq 1 \} \cup \{ a^n b^m c^m d^n \mid n, m \geq 1 \}$
- L is inherently ambiguous
- Why?

Input string: $a^n b^n c^n d^n$

Chomsky Normal Form

- *A context-free grammar $G = (V, \Sigma, R, S)$ is in Chomsky normal form if every rule is of the form*
- *$A \rightarrow BC$ or $A \rightarrow x$*
- *with variables $A \in V$ and $B, C \in V \setminus \{S\}$, and $x \in \Sigma$ For the start variable S we also allow the rule $S \rightarrow \epsilon$*
- *Advantage: Grammars in this form are far easier to analyse.*

Theorem 2.9

- *Every context-free language can be described by a grammar in Chomsky normal form.*
- *Outline of Proof:*
 - *We rewrite every CFG in Chomsky normal form.*
 - *We do this by replacing, one-by-one, every rule that is not ‘Chomsky’.*
 - *We have to take care of: Starting Symbol, ϵ symbol, all other violating rules.*

Example of Chomsky NF

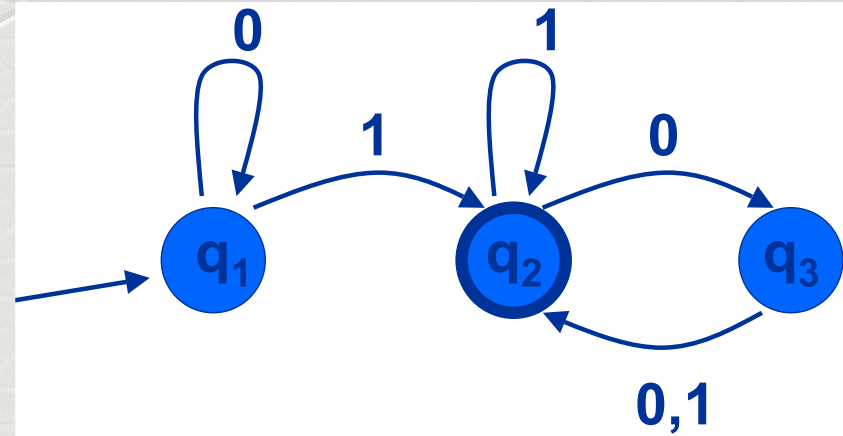
- *Initial grammar: $S \rightarrow aSb \mid \epsilon$*
- *In Chomsky normal form:*
 - $S_0 \rightarrow \epsilon \mid T_a T_b \mid T_a X$
 - $X \rightarrow S T_b$
 - $S \rightarrow T_a T_b \mid T_a X$
 - $T_a \rightarrow a$
 - $T_b \rightarrow b$

RL \subseteq CFL

- *Every regular language can be expressed by a context-free grammar.*
- *Proof Idea:*
 - *Given a DFA $M = (Q, \Sigma, \delta, q_0, F)$, we construct a corresponding CF grammar $GM = (V, \Sigma, R, S)$ with $V = Q$ and $S = q_0$*
 - *Rules of GM:*
 - $q_i \rightarrow x \delta(q_i, x)$ for all $q_i \in V$ and all $x \in \Sigma$
 - $q_i \rightarrow \varepsilon$ for all $q_i \in F$

Example $RL \subseteq CFL$

- *The DFA*



- *leads to the context-free grammar*
- $G_M = (Q, \Sigma, R, q_1)$ with the rules
 - $q_1 \rightarrow 0q_1 \mid 1q_2$
 - $q_2 \rightarrow 0q_3 \mid 1q_2 \mid \epsilon$
 - $q_3 \rightarrow 0q_2 \mid 1q_2$

Summary

- *Context-free grammars*
- *Context-free languages*
- *Productions, derivations, recursive inference, parse trees*
- *Left-most & right-most derivations*
- *Ambiguous grammars*
- *Removing ambiguity*
- *CFL/CFG applications*
 - *parsers, markup languages*