

COMPSCI 501: Formal Language Theory

Lecture 7: Context-Free Grammars

Marius Minea
 marius@cs.umass.edu

University of Massachusetts Amherst

6 February 2019

Grammars describe programming languages

(6.8.4) *selection-statement*:

```

if ( expression ) statement
if ( expression ) statement else statement
switch ( expression ) statement
    
```

(6.8.5) *iteration-statement*:

```

while ( expression ) statement
do statement while ( expression ) ;
for ( expressionopt ; expressionopt ; expressionopt ) statement
for ( declaration expressionopt ; expressionopt ) statement
    
```

Grammars describe natural language

A good student reads books.

S → NP VP	noun phrase + verb phrase
NP → Noun	
NP → Det NP	determiner: article, adjective
VP → Verb	
VP → Verb NP	verb + (direct) object
Det → a good	
Noun → books student	
Verb → reads	

Terminology

$A \rightarrow 0A1$
 $A \rightarrow B$
 $B \rightarrow \#$

- ▶ Grammar = set of **productions** (substitution **rules**)
- ▶ Left-hand side: **variable (nonterminal)**
- ▶ Right-hand side: string of symbols (variables and **terminals**)
- ▶ A **start variable**

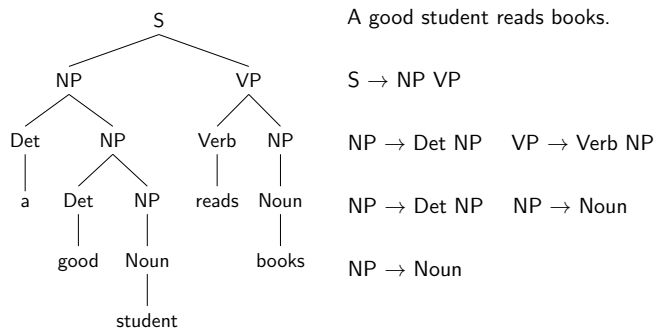
Derivation: apply substitution rules from start variable, until no variables remain.

$A \Rightarrow 0A1 \Rightarrow 00A11 \Rightarrow 000A111 \Rightarrow 000B111 \Rightarrow 000\#111$

Or represent as **parse tree**.

Parse Tree for a Sentence

Hierarchical representation of derivation
 Each symbol on RHS linked as child of LHS variable



Sentence: terminals (leaves) left to right

Formal Definition

A context-free grammar is a 4-tuple (V, Σ, R, S) , where

- ▶ V is a finite set, called **variables** (also called nonterminals)
- ▶ Σ is a finite set of **terminals**; $\Sigma \cap V = \emptyset$
- ▶ R is a finite set of **rules**, $R \subseteq V \times (\Sigma \cup V)^*$
- ▶ $S \in V$ is the start variable

Denote $uAv \Rightarrow uvv$ (**yields**) if $A \rightarrow w$ is a rule

u **derives** v ($u \xRightarrow{*} v$) if u yields v in ≥ 0 steps, i.e.,
 $u = v$ or $u \Rightarrow u_1 \Rightarrow u_2 \dots \Rightarrow u_k \Rightarrow v$

Language of a grammar: all strings (of **terminals**) generated from start symbol
 $L(G) = \{w \in \Sigma^* \mid S \xRightarrow{*} w\}$

Context-free language: generated by a context-free grammar

Designing Grammars

$0^n 1^n$

“count” by matching individual symbols, use nonterminals for rest

Match *first* 0 with *last* 1

expose the same (but smaller) pattern : recursion
(direct or indirect: circular dependencies between variables)

$S \rightarrow 0S1|\epsilon$

Balanced Parentheses:

first symbol is (and must be matched by) somewhere
string in between and string after are balanced

$S \rightarrow (S)S|\epsilon$

Designing Grammars

$L = \{ww^R | w \in \Sigma^*\}$

match first symbol with last symbol

$S \rightarrow \epsilon | aSa | bSb | \dots$ for all symbols in Σ

Other tips:

split into sublanguages, use union
recognize recursive structures

Regular Languages are Context-Free

Can we convert a regular expression to a grammar ?

Grammars have concatenation
and union (multiple right-hand sides)

Kleene star? $A_star \rightarrow \epsilon | AA_star$

Can also easily convert a DFA to a CFG:

- ▶ make variable R_i for each state q_i
- ▶ for transition $q_i \xrightarrow{a} q_j$, add rule $R_i \rightarrow aR_j$
- ▶ for each accept state q_i , add rule $R_i \rightarrow \epsilon$
- ▶ starting state q_0 gives start variable R_0

Ambiguity

Can we have multiple derivations for a parse tree?

Yes, order of expanding children ($S \rightarrow AB$, $A \rightarrow a$, $B \rightarrow b$)

Not a problem.

Leftmost derivation: replace leftmost variable at each step
unique for each parse tree

Multiple **parse trees** for one string = **ambiguity**.

(= different **rules** to derive the same sentence)

Is a problem: meaning usually associated with production rule.

“the girl touches *the boy with the flower*”

“the girl *touches* the boy *with the flower*”

(propositional phrase part of noun phrase or verb phrase)

Def: A string is derived **ambiguously** if it has more than one leftmost derivation.

Ambiguity: If-Then-Else

(Backus-Naur-Form: mix of grammar and regex notation)

$Stmt ::= ExpStmt | IfStmt | WhileStmt | Block$

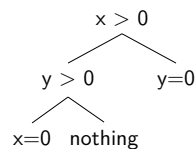
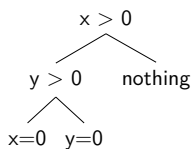
$ExpStmt ::= expr ;$

$IfStmt ::= if (expr) Stmt else Stmt \quad | \quad if (expr) Stmt$

$WhileStmt ::= while (expr) Stmt \quad Block ::= \{ Stmt^* \}$

Problem: which **if** is matched by the **else** ?

if ($x > 0$) **if** ($y > 0$) $x = 0$; **else** $y = 0$;



Disambiguating If-Then-Else

Redesign grammar; distinguish between:

balanced if (with **else**)

unbalanced if (without **else**)

Other statements are included in the “balanced” category

$Stmt ::= BalancedStmt | UnBalancedIf$

$BalancedStmt ::= ExpStmt | WhileStmt | Block | BalancedIf$

$ExpStmt ::= expr ;$

$WhileStmt ::= while (expr) Stmt \quad Block ::= \{ Stmt^* \}$

$BalancedIf ::= if (expr) BalancedStmt else Stmt$

$UnBalancedIf ::= if (expr) Stmt$

Disambiguating Expressions

$\text{Expr} \rightarrow \text{Expr} + \text{Expr} \mid \text{Expr} * \text{Expr} \mid (\text{Expr}) \mid \text{num}$

Two parse trees for: $\text{num} + \text{num} * \text{num}$

Disambiguate: introduce one nonterminal for each precedence level

$\text{Expr} \rightarrow \text{Term} \mid \text{Expr} + \text{Term}$

$\text{Term} \rightarrow \text{Factor} \mid \text{Term} * \text{Factor}$

$\text{Factor} \rightarrow (\text{Expr}) \mid \text{num}$

Chomsky Normal Form

A CFG is in **Chomsky normal form** if every rule is of the form

$$\begin{aligned} A &\rightarrow BC \\ A &\rightarrow a \\ S &\rightarrow \varepsilon \end{aligned}$$

where B, C may not be the start variable

Theorem: Any context-free language can be generated by a CFG in Chomsky Normal Form

Proof: By Construction

Conversion to Chomsky Normal Form

- ▶ Add new start variable $S_0 \rightarrow S$, if S appears on RHS.
- ▶ Eliminate rules $A \rightarrow \epsilon$:
for each rule with A on RHS, add a copy eliminating A must do for each occurrence:
 $R \rightarrow AuA$ will yield $R \rightarrow uA|Au|u$
for $R \rightarrow A$ we add $R \rightarrow \varepsilon$ unless it was previously removed (ensures termination)
- ▶ Eliminate unit rules $A \rightarrow B$:
for any $B \rightarrow u$, add $A \rightarrow u$, unless this is a previously removed unit rule
- ▶ Eliminate rules with ≥ 3 symbols on RHS:
for $SA \rightarrow u_1u_2 \dots u_k$, add rules:
 $A \rightarrow u_1A_1, A_1 \rightarrow u_2A_2, \dots, A_{k-2} \rightarrow u_{k-1}u_k$.
If any u_i are terminals, add rule $U_i \rightarrow u_i$

Outlook: Chomsky Hierarchy

Type-3: **regular** grammar: generate regular languages

$A \rightarrow a, A \rightarrow \varepsilon, A \rightarrow aB$ (right-regular), OR

$A \rightarrow a, A \rightarrow \varepsilon, A \rightarrow Ba$ (left-regular); **don't** combine!

recognized by deterministic automata

Type-2: **context-free** grammar

$A \rightarrow \gamma$ left: nonterminal; right: any string

recognized by (nondeterministic) pushdown automata

Type-1: **context-sensitive** grammar

$\alpha A \beta \rightarrow \alpha \gamma \beta$ with $|\gamma| \geq |A|$

(also $S \rightarrow \varepsilon$)

recognized by linear-bounded nondeterministic Turing machine

Type-0: **recursively enumerable** languages

$\alpha \rightarrow \beta$ (α contains some nonterminal)

recognized by Turing machine