

CS321 Week 4: Context-Free Grammars

Jingke Li, Portland State University

Winter 2014

Today's Topics

- ▶ Context-free grammars
- ▶ Derivations and parse trees
- ▶ Ambiguity and recursions
- ▶ Equivalence and transformations

Next Two Weeks:

- ▶ Top-down parsing
- ▶ Bottom-up parsing

Language Specifications

To process a language, we need specifications for all its components: words, sentence structures, semantics, etc.

While specifications can come in many different forms, to compilers, we want the specifications be both *precise* and *concise*.

Formal systems for programming language specifications:

- ▶ lexical tokens — regular expressions
- ▶ syntax structures — context-free grammars
- ▶ semantics — denotation, abstract machine, or logic axioms

Formal Grammars

Grammar is a system of rules that defines the syntactical structure of a language. (*Merriam-Webster.com*)

A *formal grammar* is used to define the syntactical structure of a formal language. It consists of the following components:

- ▶ A finite set of *terminals* — symbols that form sentences.
- ▶ A finite set of *nonterminals* — variables that can be extended to sentences by the grammar rules.
 - ▶ One of the nonterminals is designated as the *start symbol*.
- ▶ A finite set of *productions* — grammar rules of the form "*lhs* \rightarrow *rhs*", where *lhs* and *rhs* are sequences of terminals and nonterminals, with *lhs* containing at least one nonterminal. The productions specify the relationship between terminals and nonterminals, and how they should be combined to form sentences.

Chomsky Hierarchy

- ▶ *Unrestricted* Grammars: $\alpha A \beta \rightarrow \gamma$
 - ▶ Generate *recursively enumerable* languages
- ▶ *Context-Sensitive* Grammars: $\alpha A \beta \rightarrow \alpha \gamma \beta$
 - ▶ Generate *context-sensitive* languages
- ▶ *Context-Free* Grammars: $A \rightarrow \gamma$
 - ▶ Generate *context-free* languages (which include most programming languages)
- ▶ *Regular* Grammars: $A \rightarrow a$ or $A \rightarrow aB$
 - ▶ Generate *regular* languages, i.e. regular expressions

Context-Free Grammars (CFGs)

A CFG is defined by

- ▶ A finite set of terminals V_t
 - ▶ For a PL, terminals are just tokens recognized by the lexer
- ▶ A finite set of nonterminals V_n , including a start symbol S
- ▶ A finite set of productions of the form
$$A \rightarrow x_1 \cdots x_m \quad \text{where } A \in V_n, x_i \in V_n \cup V_t$$

Backus-Naur Form (BNF):

BNF is a particular notation for CFGs. It was first used in the formal description of Algol 60. The notation convention has evolved through the years; current convention include:

- ▶ Nonterminals begin with upper-case letters.
- ▶ Terminals are quoted. However, terminals that do not correspond to literal characters or strings can be represented as they are.

Extended BNF

EBNF is *any* extension of BNF, usually with the following features:

- ▶ A vertical bar, |, represents a choice,
- ▶ Parentheses, (and), represent grouping,
- ▶ Square brackets, [and], represent an optional construct,
- ▶ Curly braces, { and }, represent zero or more repetitions,

EBNF does not have extra power than BNF. In fact, any grammar in EBNF can be transformed into one in proper BNF.

An Original (Algol 60-Style) BNF Example

```

(program) ::= begin (statement list) end
(statement list) ::= (statement) (statement list)
(statement list) ::= (statement)
(statement) ::= (assignment)
(statement) ::= (read statement)
(statement) ::= (write statement)
(assignment) ::= ID := (expression) ;
(read statement) ::= read ( (id list) ) ;
(write statement) ::= write ( (expr list) ) ;
(id list) ::= ID , (id list)
(id list) ::= ID
(expr list) ::= (expression) , (expr list)
(expr list) ::= (expression)
(expression) ::= (primary) + (expression)
(expression) ::= (primary) - (expression)
(expression) ::= (primary)
(primary) ::= ( (expression) )
(primary) ::= ID
(primary) ::= INTEGER
    
```

A Corresponding Modern Version

```

Program → "begin" StmtList "end"
StmtList → Statement StmtList
StmtList → Statement
Statement → Assignment
Statement → ReadStmt
Statement → WriteStmt
Assignment → id ":=" Expression ";"
ReadStmt → "read" "(" IdList ")" ";"
WriteStmt → "write" "(" ExprList ")" ";"
IdList → id "," IdList
IdList → id
ExprList → Expression "," ExprList
ExprList → Expression
Expression → Primary "+" Expression
Expression → Primary "-" Expression
Expression → Primary
Primary → "(" Expression ")"
Primary → id
Primary → integer
    
```

A Corresponding EBNF

```

Program → "begin" StmtList "end"
StmtList → Statement { Statement }
Statement → Assignment | ReadStmt | WriteStmt
Assignment → id ":=" Expression ";"
ReadStmt → "read" "(" IdList ")" ";"
WriteStmt → "write" "(" ExprList ")" ";"
IdList → id { "," id }
ExprList → Expression { "," Expression }
Expression → Primary { "+" | "-" } Primary
Primary → "(" Expression ")" | id | integer
    
```

Transforming EBNF to BNF

Rules:

$$A \rightarrow \alpha [\beta] \gamma \Rightarrow \begin{array}{l} A \rightarrow \alpha B \gamma \\ B \rightarrow \beta \\ B \rightarrow \epsilon \end{array}$$

$$A \rightarrow \alpha \{\beta\} \gamma \Rightarrow \begin{array}{l} A \rightarrow \alpha B \gamma \\ B \rightarrow \beta B \\ B \rightarrow \epsilon \end{array}$$

Example:

$$\text{IdList} \rightarrow \text{id} \{ \text{" , " id} \} \Rightarrow \begin{array}{l} \text{IdList} \rightarrow \text{id IdTail} \\ \text{IdTail} \rightarrow \text{" , " id IdTail} \\ \text{IdTail} \rightarrow \epsilon \end{array}$$

Language Defined by a Grammar

The set of sentences *derived* from the start symbol S of a grammar G by repeatedly applying the production rules defines the language of the grammar, denoted $L(G)$: $L(G) = \{x \in V_t^* \mid S \xRightarrow{*} x\}$.

A *derivation* step is to replace the lhs of a production by its rhs.

Example:

$$\begin{array}{ll} (1) E \rightarrow E + T & E \xRightarrow{(1)} E + T \xRightarrow{(2)} T + T \\ (2) E \rightarrow T & \xRightarrow{(3)} T * P + T \xRightarrow{(4)} P * P + T \\ (3) T \rightarrow T * P & \xRightarrow{(5)} \text{id} * P + T \xRightarrow{(5)} \text{id} * \text{id} + T \\ (4) T \rightarrow P & \xRightarrow{(4)} \text{id} * \text{id} + P \xRightarrow{(5)} \text{id} * \text{id} + \text{id} \\ (5) P \rightarrow \text{id} & \end{array}$$

Notations:

\Rightarrow one derivation step
 $\xRightarrow{+}$ one or more steps
 $\xRightarrow{*}$ zero or more steps

Note: When context is clear, we omit quotes around terminal symbols in our examples.

Derivations

A sentence can be derived by different sequences of productions.

Example (Revisit):

(1) $E \rightarrow E + T$	$E \xRightarrow{(1)} E + T \xRightarrow{(2)} T + T$
(2) $E \rightarrow T$	$\xRightarrow{(3)} T * P + T \xRightarrow{(4)} P * P + T$
(3) $T \rightarrow T * P$	$\xRightarrow{(5)} id * P + T \xRightarrow{(5)} id * id + T$
(4) $T \rightarrow P$	$\xRightarrow{(4)} id * id + P \xRightarrow{(5)} id * id + id$
(5) $P \rightarrow id$	

Two alternative derivation sequences:

$E \xRightarrow{(1)} E + T \xRightarrow{(2)} T + T$	$E \xRightarrow{(1)} E + T \xRightarrow{(4)} E + P$
$\xRightarrow{(3)} T * P + T \xRightarrow{(4)} P * P + T$	$\xRightarrow{(5)} E + id \xRightarrow{(2)} T + id$
$\xRightarrow{(4)} P * P + P \xRightarrow{(5)} P * P + id$	$\xRightarrow{(3)} T * P + id \xRightarrow{(4)} P * P + id$
$\xRightarrow{(5)} P * id + id \xRightarrow{(5)} id * id + id$	$\xRightarrow{(5)} id * P + id \xRightarrow{(5)} id * id + id$

Derivations of Compiler Interests

- **Left-Most Derivation** ($\Rightarrow_{lm}, \xRightarrow{+}_{lm}, \xRightarrow{*}_{lm}$):
Always choose the left-most nonterminal to expand.
- **Right-Most Derivation** ($\Rightarrow_{rm}, \xRightarrow{+}_{rm}, \xRightarrow{*}_{rm}$):
Always choose the right-most nonterminal to expand.

Example:

$S \rightarrow S ; S \mid id := E \mid print (E)$
 $E \rightarrow id \mid num$
 Sentence: $id := num ; print (id)$

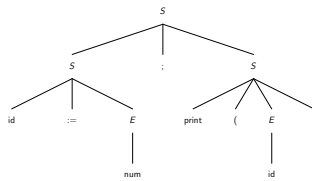
$\underline{S} \Rightarrow_{lm} \underline{S} ; S$	$\underline{S} \Rightarrow_{rm} S ; \underline{S}$
$\Rightarrow_{lm} id := \underline{E} ; S$	$\Rightarrow_{rm} S ; print (\underline{E})$
$\Rightarrow_{lm} id := num ; \underline{S}$	$\Rightarrow_{rm} \underline{S} ; print (id)$
$\Rightarrow_{lm} id := num ; print (\underline{E})$	$\Rightarrow_{rm} id := \underline{E} ; print (id)$
$\Rightarrow_{lm} id := num ; print (id)$	$\Rightarrow_{rm} id := num ; print (id)$

Parse Tree

A tree representation of the derivation from the start symbol to a sentence, *a.k.a. concrete syntax tree*.

- **root** — the start symbol
- **internal nodes** — nonterminals
- **leaf nodes** — terminals

Both the left-most derivation and the right-most derivation are linear representations of a parse tree.



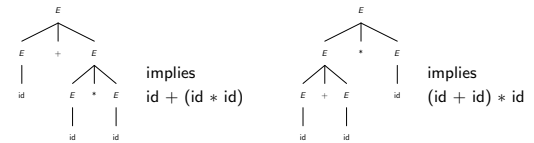
Ambiguity in Grammars

Sometimes, multiple parse trees can be derived from a grammar for the same sentence.

Example 1: Arithmetic expressions

Grammar: $E \rightarrow E + E \mid E * E \mid (E) \mid - E \mid id$
 Sentence: $id + id * id$

Parse Trees:



A grammar is said to be *ambiguous* if it can derive some sentence with more than one parse tree.

- Parsers can not be built with an ambiguous grammar!

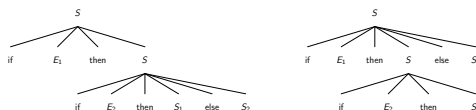
Ambiguity in Grammars (cont.)

Example 2: The "Dangling Else"

Grammar: $S \rightarrow \text{if } E \text{ then } S \mid \text{if } E \text{ then } S \text{ else } S \mid \text{other}$

Sentence: $\text{if } E_1 \text{ then if } E_2 \text{ then } S_1 \text{ else } S_2$

Parse Trees:



Implied semantics:

$\text{if } E_1 \text{ then } \{ \text{if } E_2 \text{ then } S_1 \text{ else } S_2 \}$ $\text{if } E_1 \text{ then } \{ \text{if } E_2 \text{ then } S_1 \} \text{ else } S_2$

Recursion in Grammars

Recursion in productions is necessary for any non-trivial CFGs. (Why?)

Recursions can appear in different positions in a production:

$P \rightarrow a P b \mid P c d \mid e P \mid f$
 $S \rightarrow A a \mid b$
 $A \rightarrow A c \mid S d \mid e$

Of Compiler Interests:

- **Left-recursion:** $E \xRightarrow{+} E \alpha$ for some string α .
 - Immediate form: $E \rightarrow E + id$
 - General form: $E \rightarrow F, F \rightarrow E + id$
- **Right-recursion:** $E \xRightarrow{+} \alpha E$ for some string α .
 - Immediate form: $E \rightarrow id + E$
 - General form: $E \rightarrow G, G \rightarrow id + E$

Grammars with left-recursive productions are not suitable for the predictive parsing method (*i.e.* top-down parsing).

Equivalence of Grammars

As shown earlier, each grammar G defines a single language $L(G)$.

However, the reverse is not true. For a given language L , there can be multiple grammars that define it. I.e., we can have G_1, G_2, G_3, \dots that

$$L = L(G_1) = L(G_2) = L(G_3) = \dots$$

Example 1:

$$G_1: S \rightarrow Sx \mid \epsilon$$

$$G_2: S \rightarrow xS \mid \epsilon$$

Both define the language x^* .

Example 2:

$$G_1: E \rightarrow E + E \mid E - E \mid \text{id} \text{ (ambiguous!)}$$

$$G_2: E \rightarrow E + \text{id} \mid E - \text{id} \mid \text{id}$$

Both define the set of binary expressions with $+$ and $-$.

Grammar Transformations

For writing parsers, it is often necessary to transform a given grammar to an equivalent, but more parser-friendly grammar.

Examples:

- Elimination of ambiguity
- Elimination of left recursion
- Left factoring

Elimination of Ambiguity

Approach 1: Rewrite grammar with operator precedence.

Example:

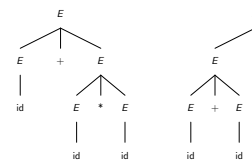
Precedence (from high to low): unary $-$, $*$, $+$.

$$\begin{array}{ll} 1. E \rightarrow E + E & 1. E \rightarrow E + T \\ 2. E \rightarrow E * E & 2. E \rightarrow T \\ 3. E \rightarrow -E & 3. T \rightarrow T * P \\ 4. E \rightarrow (E) & 4. T \rightarrow P \\ 5. E \rightarrow \text{id} & 5. P \rightarrow -P \\ & 6. P \rightarrow (E) \\ & 7. P \rightarrow \text{id} \end{array} \Rightarrow$$

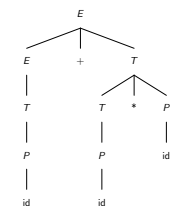
Operator Precedence Example (cont.)

Sentence: $\text{id} + \text{id} * \text{id}$

Before: multiple parse trees



After: unique parse tree



Eliminating Ambiguity (cont.)

Approach 2: Rewrite grammar with more syntactic constraints.

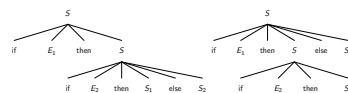
Example:

$$\begin{aligned} S &\rightarrow \text{if } E \text{ then } S \mid \text{if } E \text{ then } S \text{ else } S \mid \text{other} \\ \Rightarrow \\ S &\rightarrow m_S \mid unm_S \\ m_S &\rightarrow \text{if } E \text{ then } m_S \text{ else } m_S \mid \text{other} \\ unm_S &\rightarrow \text{if } E \text{ then } S \mid \text{if } E \text{ then } m_S \text{ else } unm_S \end{aligned}$$

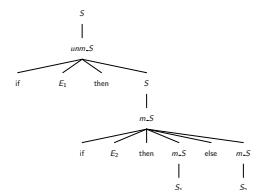
Syntactic Constraints Example (cont.)

Sentence: $\text{if } E_1 \text{ then if } E_2 \text{ then } S_1 \text{ else } S_2$

Before: multiple parse trees



After: unique parse tree



Syntactic Constraints Example (cont.)

Details are important. If we change the transformation slightly:

$$\begin{aligned}
 S &\rightarrow \text{if } E \text{ then } S \mid \text{if } E \text{ then } S \text{ else } S \mid \text{other} \\
 \Rightarrow \\
 S &\rightarrow m.S \mid unm.S \\
 m.S &\rightarrow \text{if } E \text{ then } m.S \text{ else } S \mid \text{other} \\
 unm.S &\rightarrow \text{if } E \text{ then } S
 \end{aligned}$$

then the resulting grammar is still ambiguous!

Elimination of Left Recursions

Case 1. Simple lists.

Approach: Convert directly.

- ▶ $A \rightarrow A\alpha \mid \epsilon$
Sentence patterns: $L = \alpha^*$
Equivalent productions: $A \rightarrow \alpha A \mid \epsilon$
With EBNF operators: $A \rightarrow \{\alpha\}$
- ▶ $A \rightarrow A\alpha \mid \beta$
Sentence patterns: $L = \beta\alpha^*$
Equivalent productions: $A \rightarrow \beta A'$
 $A' \rightarrow \alpha A' \mid \epsilon$
With EBNF operators: $A \rightarrow \beta \{\alpha\}$

Elimination of Left Recursions (cont.)

Case 2. Immediate left recursions.

Approach: Collect all productions for a given nonterminal; group them into left-recursive ones, and non-left-recursive ones.

- ▶ $A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$
Grouping and factoring:
 $A \rightarrow A(\alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_m) \mid (\beta_1 \mid \beta_2 \mid \dots \mid \beta_n)$
Rewriting:
 $A \rightarrow AB \mid C$
 $B \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_m$
 $C \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$
Sentence patterns: $L = CB^*$
Equivalent productions:
 $A \rightarrow (\beta_1 \mid \dots \mid \beta_n)A'$
 $A' \rightarrow (\alpha_1 \mid \dots \mid \alpha_m)A' \mid \epsilon$

Elimination of Left Recursions (cont.)

Case 3. General non-immediate left recursions.

Approach: Inline productions to turn general non-immediate left recursions into immediate left recursions; then apply the replacement transformation.

- ▶ $A \rightarrow B\alpha_1 \mid C\alpha_2 \mid \alpha_3$
 $B \rightarrow A\beta_1 \mid \beta_2$
 $C \rightarrow A\gamma_1 \mid \gamma_2$
Inlining: $A \rightarrow A\beta_1\alpha_1 \mid \beta_2\alpha_1 \mid A\gamma_1\alpha_2 \mid \gamma_2\alpha_2 \mid \alpha_3$

Example:

$$\begin{aligned}
 E &\rightarrow F \mid T \\
 F &\rightarrow E + x \mid E - y \mid P \\
 \Rightarrow \\
 E &\rightarrow E + x \mid E - y \mid P \mid T \\
 F &\rightarrow E + x \mid E - y \mid P \\
 \Rightarrow \\
 E &\rightarrow PE' \mid TE' \\
 E' &\rightarrow +xE' \mid -yE' \mid \epsilon \\
 F &\rightarrow E + x \mid E - y \mid P
 \end{aligned}$$

Left Factoring

Productions with common prefix cause problems for predictive parsers:

$$\begin{aligned}
 S &\rightarrow \text{if } E \text{ then } S \text{ end if ;} \\
 S &\rightarrow \text{if } E \text{ then } S \text{ else } S \text{ end if ;}
 \end{aligned}$$

An easy fix is to perform left factoring:

$$\begin{aligned}
 S &\rightarrow \text{if } E \text{ then } S \text{ } T \\
 T &\rightarrow \text{end if ;} \\
 T &\rightarrow \text{else } S \text{ end if ;}
 \end{aligned}$$

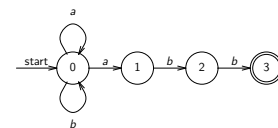
When the common prefix is not explicit, it is hard to fix.

CFGs vs. Regular Expressions

Every set that can be described by a regular expression can also be described by a context-free grammar.

Example: the R.E. $(a|b)^*abb$ can be described by

$$\begin{aligned}
 A_0 &\rightarrow aA_0 \mid bA_0 \mid aA_1 \\
 A_1 &\rightarrow bA_2 \\
 A_2 &\rightarrow bA_3 \\
 A_3 &\rightarrow \epsilon
 \end{aligned}$$



Regular Grammars are just context-free grammars with their productions limited to the following two forms,

$$\begin{aligned}
 A &\rightarrow aB \\
 C &\rightarrow \epsilon
 \end{aligned}$$

Chomsky Hierarchy Revisit

The Chomsky hierarchies are properly nested, *i.e.*

Regular languages \subset CFG languages \subset CSG languages \subset R.E. languages

Examples:

- ▶ The language $\{ [^i]^i \mid i \geq 1 \}$ is not regular, but it can be defined by a context-free grammar:

$$\begin{aligned} S &\rightarrow [T] \\ T &\rightarrow [T] \mid \epsilon \end{aligned}$$

- ▶ The following language is not describable by CFG:

$$L = \{ w c w \mid w \text{ is in } (a|b)^* \}$$

In contrast, the following is a CFG:

$$L' = \{ w c w^R \mid w \text{ is in } (a|b)^* \}$$