

03 - Grammars

Dr. Robert Lowe

Division of Mathematics and Computer Science
Maryville College

Outline

- 1 Grammar and Metalanguages
- 2 Types of Languages
- 3 Ledgard

Metalanguages and Subject Languages

- We need a way to make a formal specification of a language.
- Languages which describe other languages are called **metalanguages**. Examples:
 - Set Builder Notation
 - BNF - Backus Naur Form
 - Regular Expressions
 - YACC
- Languages described by a metalanguage are called **subject languages**.

Languages and Grammars

- A **language** L is the set of all texts expressible within that language.
- Languages are typically infinite sets.
- An expressible text in some language L is referred to as a **sentence**.
- The set of rules used to verify membership in a language is called a **grammar**.
- A grammar is sometimes also called a **syntax**.

Grammars

A grammar G is defined as follows:

$$G = \langle T, N, S, P \rangle$$

- T A finite set of basic symbols (called an **alphabet**). These symbols are called **terminal symbols** of the language.
- N A finite alphabet of **non-terminal symbols**. Members of N label parts of sentences in L .
- S $S \in N$ is the **distinguished symbol** or **start symbol**. This is the name of entire sentences.
- P A set of productions. Set of rules which transform strings of terminal and non-terminal symbols into a valid sentence.

Example Language: Simplified Email Addresses

Example Sentence:

robert.lowe@maryvillecollege.edu

$$T = \{a - z, A - Z, 0 - 9, ., @\}$$

$$N = \{< \text{email} >, < \text{user} >, < \text{subdomain} >, \\ < \text{domain} >, < \text{tld} >\}$$

$$S = < \text{email} >$$

Email Syntax Productions (P)

Productions (P)

$$\begin{aligned} \langle \text{email} \rangle &\Rightarrow \langle \text{user} \rangle @ \langle \text{domain} \rangle \\ \langle \text{user} \rangle &\Rightarrow \{a - z, A - Z, 0 - 9, .\} \langle \text{user} \rangle \\ \langle \text{domain} \rangle &\Rightarrow \{ \langle \text{subdomain} \rangle . \langle \text{tld} \rangle, \\ &\quad \langle \text{subdomain} \rangle . \langle \text{domain} \rangle \} \\ \langle \text{tld} \rangle &\Rightarrow \{ \text{com}, \text{org}, \text{edu} \} \end{aligned}$$

Synthesis of an Email Sentence

< email > ⇒ < user > @ < domain >
⇒ robert.lowe@ < domain >
⇒ robert.lowe@ < subdomain > . < tld >
⇒ robert.lowe@maryvillecollege. < tld >
⇒ robert.lowe@maryvillecollege.edu

Some Properties of Grammars and Their Compilers

- T and N must be disjoint.
- Elements in N must be distinguishable in some way.
- Each production must be distinguishable.
- The work of a compiler with these definitions in mind is:
 - 1 Verify that every symbol in a program is in T .
 - 2 Reduce the program string to a string over N .
 - 3 Find a series of productions in P which produce S .
 - 4 If no such series of productions is possible, this is not a valid program!
 - 5 On success, generate code.
- **Discuss:** Does this process have anything to do with closure?

Rules for Productions

- Each production is a pair of strings (p, q) .
- p and q contain members of T or N (or both).
- A production rule is a rule where replacing p by q is allowed.
- We often write $p \rightarrow q$ or $p ::= q$ for each production.
- The p must contain at least one non-terminal. **Discuss** Why? (HINT: The word “terminal” is not just clever phrasing!)
- The rule of substitution:
 - Suppose $V \Rightarrow W$ where $V, W \in (N \cup T)^*$
 - This holds if we can decompose V and W as:
 - $V = XV'Y$
 - $W = XW'Y$
 - Where there exists production $V' \rightarrow W'$.
- Let's do this for the email grammar!

Chomsky Hierarchy

- Let $a \in T$
- Let $A, B \in N$
- Let $\alpha, \beta, \gamma \in (N \cup T)^*$ where $\gamma \neq \lambda$

Type-0 Recursively Enumerable

$$\alpha A \beta \rightarrow \beta$$

Type-1 Context-Sensitive

$$\alpha A \beta \rightarrow \alpha \gamma \beta$$

Type-2 Context-Free

$$A \rightarrow \alpha$$

Type-3 Regular

$$A \rightarrow a \text{ and } A \rightarrow aB$$

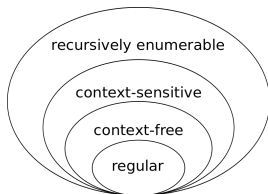


Image Source: https://en.wikipedia.org/wiki/Chomsky_hierarchy

Grammars for Programming Languages

- Recursively Enumerable (Type-0) grammars are generally too powerful to describe programming languages. (They would also be too difficult to compile!)
- Most programming languages do contain at least a handful of context-sensitive attributes, though writing true context sensitive grammars is not needed in most cases.
- Context-Free (Type-2) grammars will be our focus. Even in context-sensitive languages, the syntax is often expressed first as a context free grammar with additional constraints applied by the compiler.
- Regular Grammars are too limited to express general programming languages. They are typically useful for searching and general pattern matching.

Backus Naur Form (BNF)

- BNF is a metalanguage which describes context-free grammars.
- Invented by Naur, edited by Backus, first used to describe Algol-60.
- By far, BNF is the most popular metalanguage notation.
- BNF non-terminals are described using meaningful phrases enclosed in angle brackets $\langle \rangle$.
- Replacement arrows are $::=$
- Alternate replacements are separated with $|$.
- Terminal strings are enclosed in single quotes.

Example Grammar

$\langle S \rangle ::= \langle \text{Expression} \rangle$

$\langle \text{Expression} \rangle ::= \langle \text{Term} \rangle \mid \langle \text{Expression} \rangle '+' \langle \text{Term} \rangle$

$\langle \text{Term} \rangle ::= \langle \text{Factor} \rangle \mid \langle \text{Term} \rangle '*' \langle \text{Factor} \rangle$

$\langle \text{Factor} \rangle ::= \langle \text{Unit} \rangle \mid (\langle \text{Expression} \rangle)$

$\langle \text{Unit} \rangle ::= '0' \mid '1' \mid '2' \mid '3' \mid '4' \mid '5' \mid '6' \mid '7' \mid '8' \mid '9'$

The Ledgard Programming Language

- Ledgard is a teaching language.
- Designed by Lee Wittenberg in order to teach students how to write compilers.
- Named in honor of Henry Ledgard, author of *Programming Language Landscapes*.
- Essentially contains “just enough” of the elements of programming languages to explore compiler creation.

Ledgard Syntax

$\langle \textit{program} \rangle ::= \text{'program'} \langle \textit{decl-list} \rangle \text{'begin'} \langle \textit{stmt-list} \rangle \text{'end'} \text{';'}$

$\langle \textit{decl-list} \rangle ::= \langle \textit{declaration} \rangle \mid \langle \textit{decl-list} \rangle \langle \textit{declaration} \rangle$

$\langle \textit{declaration} \rangle ::= \langle \textit{identifier-list} \rangle \text{':'} \langle \textit{type} \rangle \text{';'}$

$\langle \textit{identifier-list} \rangle ::= \langle \textit{identifier} \rangle \mid \langle \textit{identifier-list} \rangle \text{';' } \langle \textit{identifier} \rangle$

$\langle \textit{type} \rangle ::= \langle \textit{simple-type} \rangle \mid \langle \textit{array-type} \rangle$

$\langle \textit{simple-type} \rangle ::= \text{'integer'} \mid \text{'boolean'}$

$\langle \textit{array-type} \rangle ::= \text{'array'} \text{'['} \langle \textit{bounds} \rangle \text{']' } \text{'of'} \langle \textit{type} \rangle$

$\langle \textit{bounds} \rangle ::= \langle \textit{integer-literal} \rangle \text{'..' } \langle \textit{integer-literal} \rangle$

Ledgard Syntax (continued)

$\langle \textit{stmt-list} \rangle ::= \langle \textit{statement} \rangle \mid \langle \textit{stmt-list} \rangle \langle \textit{statement} \rangle$

$\langle \textit{statement} \rangle ::= \langle \textit{assignment-stmt} \rangle \mid \langle \textit{exchange-stmt} \rangle \mid$
 $\langle \textit{if-stmt} \rangle \mid \langle \textit{loop-stmt} \rangle \mid \langle \textit{input-stmt} \rangle \mid \langle \textit{output-stmt} \rangle$

$\langle \textit{assignment-stmt} \rangle ::= \langle \textit{variable} \rangle \textit{' := ' } \langle \textit{expression} \rangle \textit{' ; '}$

$\langle \textit{exchange-stmt} \rangle ::= \langle \textit{variable} \rangle \textit{' :=: ' } \langle \textit{variable} \rangle \textit{' ; '}$

$\langle \textit{if-stmt} \rangle ::= \textit{' if ' } \langle \textit{expression} \rangle \textit{' then ' } \langle \textit{stmt-list} \rangle \textit{' end ' } \textit{' if ' } \textit{' ; '}$
 $\mid \textit{' if ' } \langle \textit{expression} \rangle \textit{' then ' } \langle \textit{stmt-list} \rangle \textit{' else ' } \langle \textit{stmt-list} \rangle \textit{' end ' } \textit{' if '}$
 $\textit{' ; '}$

$\langle \textit{loop-stmt} \rangle ::= \textit{' while ' } \langle \textit{expression} \rangle \textit{' loop ' } \langle \textit{stmt-list} \rangle \textit{' end '}$
 $\textit{' loop ' } \textit{' ; '}$

Ledgard Syntax (continued)

$\langle \textit{input-statement} \rangle ::= \textit{'input'} \langle \textit{variable-list} \rangle \textit{' ;'}$

$\langle \textit{output-statement} \rangle ::= \textit{'output'} \langle \textit{variable-list} \rangle \textit{' ;'}$

$\langle \textit{variable-list} \rangle ::= \langle \textit{variable} \rangle \mid \langle \textit{variable-list} \rangle \textit{' ,'} \langle \textit{variable} \rangle$

$\langle \textit{expression} \rangle ::= \langle \textit{operand} \rangle \mid \langle \textit{operand} \rangle \langle \textit{operator} \rangle \langle \textit{operand} \rangle$

$\langle \textit{operand} \rangle ::= \langle \textit{variable} \rangle \mid \langle \textit{integer-literal} \rangle \mid \langle \textit{boolean-literal} \rangle \mid \textit{'('} \langle \textit{expression} \rangle \textit{')'} \mid \textit{'not'} \langle \textit{operand} \rangle$

$\langle \textit{variable} \rangle ::= \langle \textit{variable} \rangle \mid \langle \textit{variable} \rangle \textit{'['} \langle \textit{expression} \rangle \textit{']'}$

Ledgard Syntax (continued)

$\langle \textit{boolean-literal} \rangle ::= \text{'true'} \mid \text{'false'}$

$\langle \textit{operator} \rangle ::= \text{'<'} \mid \text{'<='} \mid \text{'=='} \mid \text{'<>'} \mid \text{'>='} \mid \text{'>'} \mid \text{'+'} \mid \text{'-'} \mid \text{'*'} \mid \text{'/'} \mid$
 $\text{'and'} \mid \text{'or'}$

- An integer-literal is just a string of digits 0-9.
- Comments begin with – and continue to the end of the line.