



# Programming Languages

G22.2110  
Summer 2010

Introduction



# Introduction

The main themes of programming language design and use:

- Paradigm (Model of computation)
- Expressiveness
  - ◆ control structures
  - ◆ abstraction mechanisms
  - ◆ types and their operations
  - ◆ tools for programming in the large
- Ease of use: Writeability / Readability / Maintainability

# Language as a tool for thought

- Role of language as a communication vehicle among programmers is more important than ease of writing
- All general-purpose languages are *Turing complete* (They can compute the same things)
- But languages can make expression of certain algorithms difficult or easy.
  - ◆ Try multiplying two Roman numerals
- Idioms in language A may be useful inspiration when writing in language B.

# Idioms

- Copying a string `q` to `p` in C:

```
while (*p++ = *q++) ;
```

- Removing duplicates from the list `@xs` in Perl:

```
my %seen = ();  
@xs = grep { ! $seen{$_}++; } @xs;
```

- Computing the sum of numbers in list `xs` in Haskell:

```
foldr (+) 0 xs
```

Is this natural?     *It is if you're used to it*

# Course Goals

- *Intellectual*: help you understand benefit/pitfalls of different approaches to language design, and how they work.
- *Practical*:
  - ◆ you will probably design languages in your career (at least small ones)
  - ◆ understanding how to use a programming paradigm can improve your programming even in languages that don't support it
  - ◆ knowing how feature is implemented helps us understand time/space complexity
- *Academic*: good start on core exam

# Compilation overview

Major phases of a compiler:

1. lexer: text  $\longrightarrow$  tokens
2. parser: tokens  $\longrightarrow$  parse tree
3. intermediate code generation
4. optimization
5. target code generation
6. optimization

# Programming paradigms

- *Imperative (von Neumann)*: **Fortran, Pascal, C, Ada**
  - ◆ programs have mutable storage (state) modified by assignments
  - ◆ the most common and familiar paradigm
- *Functional (applicative)*: **Scheme, Lisp, ML, Haskell**
  - ◆ functions are first-class values
  - ◆ *side effects* (e.g., assignments) discouraged
- *Logical (declarative)*: **Prolog, Mercury**
  - ◆ programs are sets of assertions and rules
- *Object-Oriented*: **Simula 67, Smalltalk, C++, Ada95, Java, C#**
  - ◆ data structures and their operations are bundled together
  - ◆ inheritance
- Functional + Logical: **Curry**
- Functional + Object-Oriented: **O'Caml, O'Haskell**

# Genealogy

- **FORTRAN** (1957) ⇒ **Fortran90**, **HP**
- **COBOL** (1956) ⇒ **COBOL 2000**
  - ◆ still a large chunk of installed software
- **Algol60** ⇒ **Algol68** ⇒ **Pascal** ⇒ **Ada**
- **Algol60** ⇒ **BCPL** ⇒ **C** ⇒ **C++**
- **APL** ⇒ **J**
- **Snobol** ⇒ **Icon**
- **Simula** ⇒ **Smalltalk**
- **Lisp** ⇒ **Scheme** ⇒ **ML** ⇒ **Haskell**

with lots of cross-pollination: e.g., **Java** is influenced by **C++**, **Smalltalk**, **Lisp**, **Ada**, etc.



# Predictable performance vs. ease of writing

- Low-level languages mirror the physical machine:
  - ◆ **Assembly, C, Fortran**
- High-level languages model an abstract machine with useful capabilities:
  - ◆ **ML, Setl, Prolog, SQL, Haskell**
- Wide-spectrum languages try to do both:
  - ◆ **Ada, C++, Java, C#**
- High-level languages have garbage collection, are often interpreted, and cannot be used for real-time programming. The higher the level, the harder it is to determine cost of operations.

Modern imperative languages (e.g., Ada, C++, Java) have similar characteristics:

- large number of features (grammar with several hundred productions, 500 page reference manuals, ...)
- a complex type system
- procedural mechanisms
- object-oriented facilities
- abstraction mechanisms, with information hiding
- several storage-allocation mechanisms
- facilities for concurrent programming (not C++)
- facilities for generic programming (new in Java)

# Language libraries

The programming environment may be larger than the language.

- The predefined libraries are *indispensable* to the proper use of the language, *and its popularity*.
- The libraries are defined in the language itself, but they have to be internalized by a good programmer.

Examples:

- C++ standard template library
- Java Swing classes
- Ada I/O packages

# Language definition

- Different users have different needs:
  - ◆ *programmers*: tutorials, reference manuals, programming guides (idioms)
  - ◆ *implementors*: precise operational semantics
  - ◆ *verifiers*: rigorous axiomatic or natural semantics
  - ◆ *language designers and lawyers*: all of the above
- Different levels of detail and precision
  - ◆ but none should be sloppy!

# Syntax and semantics

- Syntax refers to external representation:
  - ◆ Given some text, is it a well-formed program?
- Semantics denotes meaning:
  - ◆ Given a well-formed program, what does it mean?
  - ◆ Often depends on context.

The division is somewhat arbitrary.

- Note: It *is* possible to fully describe the syntax and semantics of a programming language by syntactic means (e.g., Algol68 and W-grammars), but this is highly impractical.

Typically use a grammar for the context-free aspects, and different method for the rest.

- Similar looking constructs in different languages often have subtly (or not-so-subtly) different meanings

A *grammar*  $G$  is a tuple  $(\Sigma, N, S, \delta)$

- $N$  is the set of *non-terminal* symbols
- $S$  is the distinguished non-terminal: the root symbol
- $\Sigma$  is the set of *terminal* symbols (alphabet)
- $\delta$  is the set of rewrite rules (productions) of the form:

$$ABC \dots ::= XYZ \dots$$

where  $A, B, C, D, X, Y, Z$  are terminals and non terminals.

- The *language* is the set of sentences containing **only** terminal symbols that can be generated by applying the rewriting rules starting from the root symbol (let's call such sentences *strings*)

# The Chomsky hierarchy

- Regular grammars (Type 3)
  - ◆ all productions can be written in the form:  $N ::= TN$
  - ◆ one non-terminal on left side; at most one on right
- Context-free grammars (Type 2)
  - ◆ all productions can be written in the form:  $N ::= XYZ$
  - ◆ one non-terminal on the left-hand side; mixture on right
- Context-sensitive grammars (Type 1)
  - ◆ number of symbols on the left is no greater than on the right
  - ◆ no production shrinks the size of the sentential form
- Type-0 grammars
  - ◆ no restrictions

# Regular expressions

An alternate way of describing a regular language is with regular expressions.

We say that a regular expression  $R$  denotes the language  $\llbracket R \rrbracket$ .

Recall that a language is a set of strings.

Basic regular expressions:

- $\epsilon$  denotes  $\emptyset$
- a character  $x$ , where  $x \in \Sigma$ , denotes  $\{x\}$
- (sequencing) a sequence of two regular expressions  $RS$  denotes  $\{\alpha\beta \mid \alpha \in \llbracket R \rrbracket, \beta \in \llbracket S \rrbracket\}$
- (alternation)  $R|S$  denotes  $\llbracket R \rrbracket \cup \llbracket S \rrbracket$
- (Kleene star)  $R^*$  denotes the set of strings which are concatenations of zero or more strings from  $\llbracket R \rrbracket$
- parentheses are used for grouping

Shorthands:

- $R^? \equiv \epsilon|R$
- $R^+ \equiv RR^*$



# Regular grammar example

A grammar for floating point numbers:

$$\text{Float} ::= \text{Digits} \mid \text{Digits} . \text{Digits}$$
$$\text{Digits} ::= \text{Digit} \mid \text{Digit Digits}$$
$$\text{Digit} ::= 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$$

A regular expression for floating point numbers:

$$(0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9)^+ (.(0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9)^+)?$$

Perl offer some shorthands:

$$[0-9]^+(\.[0-9]^+)?$$

or

$$\backslash d^+(\.\backslash d^+)?$$

Lexical: formation of words or tokens.

- Described (mainly) by regular grammars
- Terminals are characters. Some choices:
  - ◆ character set: ASCII, Latin-1, ISO646, Unicode, etc.
  - ◆ is case significant?
- Is indentation significant?
  - ◆ Python, Occam, Haskell

Example: identifiers

$$\text{Id} ::= \text{Letter IdRest}$$
$$\text{IdRest} ::= \epsilon \mid \text{Letter IdRest} \mid \text{Digit IdRest}$$

Missing from above grammar: limit of identifier length

# BNF: notation for context-free grammars

(BNF = Backus-Naur Form) Some conventional abbreviations:

- alternation:  $\text{Symb} ::= \text{Letter} \mid \text{Digit}$
- repetition:  $\text{Id} ::= \text{Letter} \{ \text{Symb} \}$   
or we can use a Kleene star:  $\text{Id} ::= \text{Letter} \text{Symb}^*$   
for one or more repetitions:  $\text{Int} ::= \text{Digit}^+$
- option:  $\text{Num} ::= \text{Digit}^+ [ . \text{Digit}^* ]$
- abbreviations do not add to expressive power of grammar
- need convention for metasymbols – what if “|” is in the language?

# Parse trees

A parse tree describes the grammatical structure of a sentence

- root of tree is root symbol of grammar
- leaf nodes are terminal symbols
- internal nodes are non-terminal symbols
- an internal node and its descendants correspond to some production for that non terminal
- top-down tree traversal represents the process of generating the given sentence from the grammar
- construction of tree from sentence is *parsing*

# Ambiguity

If the parse tree for a sentence is not unique, the grammar is *ambiguous*:

$$E ::= E + E \mid E * E \mid \text{Id}$$

Two possible parse trees for “A + B \* C”:

- $((A + B) * C)$
- $(A + (B * C))$

One solution: rearrange grammar:

$$\begin{aligned} E &::= E + T \mid T \\ T &::= T * \text{Id} \mid \text{Id} \end{aligned}$$

Harder problems – disambiguate these (courtesy of Ada):

- $\text{function\_call} ::= \text{name} (\text{expression\_list})$
- $\text{indexed\_component} ::= \text{name} (\text{index\_list})$
- $\text{type\_conversion} ::= \text{name} (\text{expression})$

# Dangling else problem

Consider:

$$S ::= \text{if } E \text{ then } S$$
$$S ::= \text{if } E \text{ then } S \text{ else } S$$

The sentence

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

is ambiguous (Which then does else  $S_2$  match?)

Solutions:

- Pascal rule: else matches most recent if
- grammatical solution: different productions for balanced and unbalanced if-statements
- grammatical solution: introduce explicit end-marker

The general ambiguity problem is unsolvable