

GA.2110-003 Programming Languages - Fall 2022

Class 1

Thomas Wies

New York University

Course Information and Resources

Course web page (general info, syllabus, etc.)

<https://cs.nyu.edu/wies/teaching/pl-fa22/>

Brightspace (announcements, course-related discussions, grades)

<https://brightspace.nyu.edu/d2l/home/222148>

You should already be enrolled.

Github (class notes and code, homework submission)

<https://github.com/nyu-pl-fa22>

Important Dates

Class Meetings

Tue 4:55-6:55pm in Silv 408

Office hours: Thomas Wies, Wed 4-5pm (60FA 403)

Recitations

Fri 5:55-6:45pm (two parallel sessions: 60FA 150 and online)

Office hours: Elaine Li, Fri 1-2pm (60FA 418)

Office hours: Nisarg Patel, Mon 11am-12pm (60FA 418)

Graders

- ▶ Vaibhav Mavi
- ▶ Rajat Narlawar
- ▶ Adithya Viswanathan

Grader office hours will be announced soon.

Midterm Exam

Tue Nov 1, 4:55-6:55pm in Silv 408

Final Exam

Wed Dec 21, 4-6pm

Grading

Grading Key

- ▶ 30% for homework assignments
- ▶ 30% for midterm exam
- ▶ 40% for final exam

Submission Policies

- ▶ All work must be your own.
- ▶ Solutions must be submitted before the announced date and time deadline for full credit.
- ▶ For every 24 hours late you lose 10%.
- ▶ Late solutions will not be accepted after the late deadline. (usually one week).
- ▶ If you turn in a solution to a programming assignment that does not compile, it will not be accepted. You can resubmit according to the above rules.

TODO Items

Check out the welcome announcement on Brightspace and complete the TODOs:

- ▶ Fill out the Github consent form.
- ▶ Follow the setup instructions for the Scala tool chain.

Textbooks

- ▶ No required textbooks!
- ▶ Though, course web page lists several good books that I recommend.
- ▶ Also, see the syllabus for class notes and pointers for recommended reading.

What is a Program?

Programs

What is a Program?

- ▶ A sequence of characters
- ▶ An abstract syntax tree (AST) obtained from a sequence of characters
- ▶ A mathematical interpretation of an AST
- ▶ The implementation of an algorithm
- ▶ A fulfillment of a customer's requirements
- ▶ A set of files on a drive
- ▶ Ones and zeroes
- ▶ Electromagnetic states of a machine

Programs

What is a Program?

- ▶ A sequence of characters
- ▶ An abstract syntax tree (AST) obtained from a sequence of characters
- ▶ A mathematical interpretation of an AST
- ▶ The implementation of an algorithm
- ▶ A fulfillment of a customer's requirements
- ▶ A set of files on a drive
- ▶ Ones and zeroes
- ▶ Electromagnetic states of a machine

In this class, we will typically take the view that a program is a sequence of characters, the AST obtained from parsing this sequence, respectively, its mathematical interpretation.

Though we recognize the close association to the other possible meanings.

Programs

What is a Program?

- ▶ A sequence of characters
- ▶ An abstract syntax tree (AST) obtained from a sequence of characters
- ▶ A mathematical interpretation of an AST
- ▶ The implementation of an algorithm
- ▶ A fulfillment of a customer's requirements
- ▶ A set of files on a drive
- ▶ Ones and zeroes
- ▶ Electromagnetic states of a machine

In this class, we will typically take the view that a program is a sequence of characters, the AST obtained from parsing this sequence, respectively, its mathematical interpretation.

Though we recognize the close association to the other possible meanings.

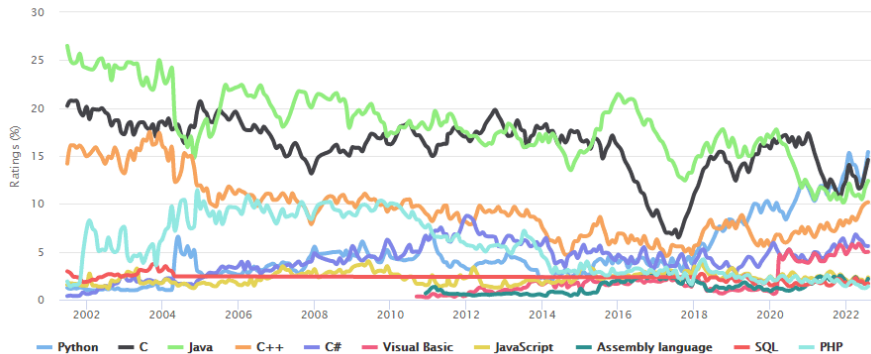
A *Programming Language* describes what sequences are allowed and give valid syntax trees (the syntax) and what these trees mean (the semantics).

Why Study Programming Languages?

Why Study Programming Languages?

TIOBE Programming Community Index

Source: www.tiobe.com



Why Study Programming Languages?

1. *Understand foundations*: understanding notation and terminology helps to learn complicated features and new languages quickly
2. *Make good choices when writing code*: understanding benefits and costs helps you make good choices when programming
3. *Better debugging and profiling*: sometimes you need to understand what's going on under the hood
4. *Simulate useful features in languages that lack them*: being exposed to different idioms and paradigms broadens your set of tools and techniques
5. *Make better use of language technology*: parsers, analyzers, optimizers appear in many contexts
6. *Leverage extension languages*: many tools are customizable via specialized languages (e.g. Google Apps Scripts, Emacs Lisp, ...)

A Brief History of Programming Languages

The first computer programs were written in *machine language*.

Machine language is just a sequence of ones and zeroes.

The computer interprets sequences of ones and zeroes as *instructions* that control the *central processing unit* (CPU) of the computer. The length and meaning of the sequences depends on the CPU.

Example

On the 6502, an 8-bit microprocessor used in the Apple II computer, the following bits add 1 and 1: 10101001000000010110100100000001.

Or, using base 16, a common shorthand: A9016901.

Programming in machine language requires an extensive understanding of the low-level details of the computer and is extremely tedious if you want to do anything non-trivial.

But it *is* the most straightforward way to give instructions to the computer: no extra work is required before the computer can run the program.

A Brief History of Programming Languages

Before long, programmers started looking for ways to make their job easier. The first step was *assembly language*.

Assembly language assigns meaningful names to the sequences of bits that make up instructions for the CPU.

A program called an *assembler* is used to translate assembly language into machine language.

Example

The assembly code for the previous example is:

LDA #\$01

ADC #\$01

Question: *How do you write an assembler?*

A Brief History of Programming Languages

Before long, programmers started looking for ways to make their job easier. The first step was *assembly language*.

Assembly language assigns meaningful names to the sequences of bits that make up instructions for the CPU.

A program called an *assembler* is used to translate assembly language into machine language.

Example

The assembly code for the previous example is:

```
LDA #$01
```

```
ADC #$01
```

Question: *How do you write an assembler?*

Answer: in machine language (at least the first time)!

A Brief History of Programming Languages

As computers became more powerful and software more ambitious, programmers needed more efficient ways to write programs.

This led to the development of *high-level* languages, the first being Fortran.

High-level languages have features designed to make things much easier for the programmer.

In addition, they are largely *machine-independent*: the same program can be run on different machines without rewriting it.

But high-level languages require a *compiler*. The compiler's job is to convert high-level programs into machine language. Alternatively, the program can be interpreted by another program, an *interpreter*, which is already compiled to machine language. More on this later...

Question: *How do you write a compiler?*

A Brief History of Programming Languages

As computers became more powerful and software more ambitious, programmers needed more efficient ways to write programs.

This led to the development of *high-level* languages, the first being Fortran.

High-level languages have features designed to make things much easier for the programmer.

In addition, they are largely *machine-independent*: the same program can be run on different machines without rewriting it.

But high-level languages require a *compiler*. The compiler's job is to convert high-level programs into machine language. Alternatively, the program can be interpreted by another program, an *interpreter*, which is already compiled to machine language. More on this later...

Question: *How do you write a compiler?*

Answer: in assembly language (at least the first time)

Language Design

Language design is influenced by various viewpoints

- ▶ *Programmers*: Desire expressive features, predictable performance, supportive development environment
- ▶ *Implementors*: Prefer languages to be simple
- ▶ *Verifiers/testers*: Languages should have rigorous semantics and discourage unsafe constructs

The interplay of design and implementation is important to understand limitations and idiosyncrasies of programming language features and we will return to this theme periodically.

Classifying Programming Languages

Major Programming Paradigms

- ▶ *Imperative (von Neumann)*: Fortran, Pascal, C, Ada, Rust
 - ▶ programs have mutable storage (state) modified by assignments
 - ▶ the most common and familiar paradigm
- ▶ *Functional (applicative)*: Scheme, Lisp, SML, OCaml, Haskell
 - ▶ based on lambda calculus
 - ▶ functions are first-class objects
 - ▶ *side effects* (e.g., assignments) discouraged
- ▶ *Object-Oriented*: Simula 67, Smalltalk, Java, Python, C#, Scala, JavaScript
 - ▶ data structures and their operations are bundled together
 - ▶ inheritance, information hiding
- ▶ *Logical (declarative)*: Prolog, Mercury
 - ▶ programs are sets of assertions and rules

Classifying Programming Languages

Multi-Paradigm Languages

Most modern general-purpose languages support multiple paradigms.

- ▶ *Object-Oriented + Imperative + Functional*: C++, Java, C#, Objective-C
- ▶ *Object-Oriented + Functional + Imperative*: Scala
- ▶ *Functional + Imperative + Object-Oriented*: OCaml, F#
- ▶ *Logical + Functional*: Curry, Oz

Concurrent Programming

- ▶ Not really a category of programming languages
- ▶ Usually implemented as extensions within existing languages
 - ▶ Threads (e.g. Pthreads in C)
 - ▶ Actors (e.g. Erlang)
 - ▶ Channels (e.g. Go)
 - ▶ Futures and promises
 - ▶ ...
- ▶ Some exceptions (e.g. *dataflow* languages)

Classifying Programming Languages

Compared to machine or assembly language, all others are high-level.

But within high-level languages, there are different levels as well.

Somewhat confusingly, these are also referred to as low-level and high-level.

- ▶ *Low-level* languages give the programmer more control (at the cost of requiring more effort) over how the program is translated into machine code.
 - ▶ C, Fortran
- ▶ *High-level* languages hide many implementation details, often with some performance cost
 - ▶ Basic, Python, Java, Lisp, Scheme, SML, OCaml, Prolog, Haskell
- ▶ *Wide-spectrum* languages try to do both:
 - ▶ Ada, C++, Rust
- ▶ High-level languages typically have garbage collection and are often interpreted.
- ▶ The higher the level, the harder it is to predict performance (bad for real-time or performance-critical applications)

Programming Idioms

- ▶ All general-purpose languages have essentially the same capabilities (they are Turing-complete)
- ▶ But different languages can make the same task difficult or easy
 - ▶ Try multiplying two Roman numerals
- ▶ Idioms in language A may be useful inspiration when writing in language B.
- ▶ One goal of this class is for you to become comfortable with many different idioms.

Characteristics of Modern Languages

Modern general-purpose languages 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
- ▶ facilities for generic programming
- ▶ development support including IDEs, libraries, compilers, build systems

We will discuss many of these in detail this semester.

Compilation vs Interpretation

Compilation

- ▶ Translates a program into machine code (or something close to it, e.g. byte code)
- ▶ Thorough analysis of the input language
- ▶ Nontrivial translation

Interpretation

- ▶ Executes a program one statement at a time using a virtual machine
- ▶ May employ a simple initial translator (preprocessor)
- ▶ Most of the work done by the virtual machine
- ▶ Often involves *just-in-time compilation* to machine code at run-time

Compilation overview

Major phases of a compiler:

1. *Lexer*: Text \longrightarrow Tokens
2. *Parser*: Tokens \longrightarrow Parse Tree
3. *Intermediate code generation*: Parse Tree \longrightarrow Intermed. Representation (IR)
4. *Optimization I*: IR \longrightarrow IR
5. *Target code generation*: IR \longrightarrow assembly/machine language
6. *Optimization II*: target language \longrightarrow target language

Syntax and Semantics

Syntax refers to the structure of the language, i.e. what sequences of characters are programs.

- ▶ Formal specification of syntax requires a set of rules
- ▶ These are often specified using *grammars*

Semantics denotes meaning:

- ▶ Given a program, what does it mean?
- ▶ Meaning may depend on context

We will not be covering semantic analysis (this is covered in the compilers and abstract interpretation courses). Though, I can provide you pointers if you are interested.

We now look at grammars in more detail.

Grammars

A *grammar* G is a tuple (Σ, N, P, S) , where:

- ▶ N is a set of *non-terminal* symbols
- ▶ $S \in N$ is a distinguished non-terminal: the *root* or *start* symbol
- ▶ Σ is a set of *terminal* symbols, also called the *alphabet*. We require Σ to be disjoint from N (i.e. $\Sigma \cap N = \emptyset$).
- ▶ P is a set of rewrite rules (productions) of the form:

$$ABC \dots \rightarrow XYZ \dots$$

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

Any sequence consisting of terminals and non-terminals is called a *word*.

The *language* defined by a grammar is the set of words containing *only* terminal symbols that can be generated by applying the rewriting rules starting from S .

Grammars Example

- ▶ $N = \{S, X, Y\}$
- ▶ $S = S$
- ▶ $\Sigma = \{a, b, c\}$
- ▶ P consists of the following rules:
 - ▶ $S \rightarrow b$
 - ▶ $S \rightarrow XbY$
 - ▶ $X \rightarrow a$
 - ▶ $X \rightarrow aX$
 - ▶ $Y \rightarrow c$
 - ▶ $Y \rightarrow Yc$

Some sample derivations:

- ▶ $S \rightarrow b$
- ▶ $S \rightarrow XbY \rightarrow abY \rightarrow abc$
- ▶ $S \rightarrow XbY \rightarrow aXbY \rightarrow aaXbY \rightarrow aaabY \rightarrow aaabc$

The Chomsky hierarchy

- ▶ Regular grammars (Type 3)

- ▶ All productions must be of one of the following shapes:

$$X ::= \epsilon \quad X ::= a \quad X ::= Y \quad X ::= aY$$

- ▶ Recognizable by finite state automaton
 - ▶ Used in *lexers*

- ▶ Context-free grammars (Type 2)

- ▶ All productions have a single non-terminal on the left
 - ▶ Right side of productions can be any word
 - ▶ Recognizable by non-deterministic pushdown automaton
 - ▶ Used in *parsers*

The Chomsky hierarchy

- ▶ Context-sensitive grammars (Type 1)
 - ▶ Each production is of the form $\alpha A \beta \rightarrow \alpha \gamma \beta$,
 - ▶ A is a non-terminal, and α, β, γ are arbitrary words (α and β may be empty, but not γ)
 - ▶ Recognizable by linear bounded automaton
- ▶ Recursively-enumerable grammars (Type 0)
 - ▶ No restrictions
 - ▶ Recognizable by turing machine

Regular expressions

An alternate way of describing a regular language over an alphabet Σ 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 words).

Regular expressions over alphabet Σ :

- ▶ ϵ 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) $R \mid S$ denotes $\llbracket R \rrbracket \cup \llbracket S \rrbracket$
- ▶ (Kleene star) R^* denotes the set of words which are concatenations of zero or more words from $\llbracket R \rrbracket$
- ▶ parentheses are used for grouping
- ▶ $R^? \equiv \epsilon \mid R$
- ▶ $R^+ \equiv RR^*$

Regular grammar example

A grammar for floating point numbers:

Float \rightarrow Digits | Digits . Digits
Digits \rightarrow Digit | Digit Digits
Digit \rightarrow 0|1|2|3|4|5|6|7|8|9

A regular expression for floating point numbers:

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

The same thing in Perl:

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

or

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

Tokens

Tokens are the basic building blocks of programs:

- ▶ keywords (`begin`, `end`, `while`).
- ▶ identifiers (`myVariable`, `yourType`)
- ▶ numbers (`137`, `6.022e23`)
- ▶ symbols (`+`, `-`)
- ▶ string literals (`"Hello world"`)
- ▶ described (mainly) by regular grammars

Example

: identifiers

$$\text{Id} \rightarrow \text{Letter IdRest}$$
$$\text{IdRest} \rightarrow \epsilon \mid \text{Letter IdRest} \mid \text{Digit IdRest}$$

Other issues: international characters, case-sensitivity, limit of identifier length

Backus-Naur Form

Backus-Naur Form (BNF) is a notation for context-free grammars:

- ▶ alternation: $\text{Symb} ::= \text{Letter} \mid \text{Digit}$
- ▶ repetition: $\text{Id} ::= \text{Letter Symb}^*$
for one or more repetitions: $\text{Int} ::= \text{Digit}^+$
- ▶ option: $\text{Num} ::= \text{Digit}^+ [. \text{Digit}^*]$

Note that these abbreviations do not add to the expressive power of context-free grammars.

Parse trees

A parse tree describes the way in which a word in the language of a grammar is derived:

- ▶ root of tree is start 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 word from the grammar
- ▶ construction of tree from word is *parsing*

Ambiguity

If the parse tree for a word 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:

$$E ::= E + T \mid T$$

$$T ::= T * \text{Id} \mid \text{Id}$$

Ambiguity

If the parse tree for a word 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:

$$E ::= E + T \mid T$$

$$T ::= T * \text{Id} \mid \text{Id}$$

Why is ambiguity bad?

Dangling else problem

Consider:

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

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

The word

$\text{if } E1 \text{ then if } E2 \text{ then } S1 \text{ else } S2$

is ambiguous (Which **then** does **else** $S2$ match?)

Dangling else problem

Consider:

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

The word

$\text{if } E1 \text{ then if } E2 \text{ then } S1 \text{ else } S2$

is ambiguous (Which **then** does **else** $S2$ 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

Lexing and Parsing

Lexer

- ▶ Reads sequence of characters of input program
- ▶ Produces sequence of tokens (identifiers, keywords, numbers, ...)
- ▶ Specified by regular expressions

Parser

- ▶ Reads sequence of tokens
- ▶ Produces parse tree
- ▶ Specified by context-free grammars

Both lexers and parsers can be automatically generated from their grammars using tools such as lex/flex respectively yacc/bison.