# CS 4120 / 4121
# CS 5120 / 5121

Introduction to Compilers
Fall 2013
Ross (Tate)

Lecture 1: Overview

CS 4120 Introduction to Compilers

---

# Outline

- About this course
- Introduction to compilers
  - What are compilers?
  - Why should we learn about them?
  - Anatomy of a compiler
- Introduction to lexical analysis
  - Text stream to tokens

CS 4120 Introduction to Compilers

2

---

# Course Information

- MWF 1:30- 2:15$_{PM}$ in Olin 245
- Instructor: Ross Tate
- Teaching Assistants:
  Fabian Mühlböck, Sam Hopkins, Lee Gao
- E-mail: cs4120-l@cs.cornell.edu
- Web page:
  http://www.cs.cornell.edu/courses/cs4120
- Newsgroup:
  https://piazza.com/class#fall2013/cs4120

3

---

# 4 = 5 & 0 = 1

- CS 4120 and 5120 are really the same course
  - same lectures
  - same assignments or nearly so
  - 5120 is for MEng students, 4120 for others
- CS 4121 (5121) is required!
  - most coursework is in the project

CS 4120 Introduction to Compilers

4

---

# Textbooks

- Lecture notes provided; no required textbook

- On reserve in Uris Library:
  - **Compilers—Principles, Techniques and Tools.** Aho, Lam, Sethi and Ullman (The Dragon Book)
    (strength: parsing)
  - **Modern Compiler Implementation in Java.** Andrew Appel.
    (strength: translation)
  - **Advanced Compiler Design and Implementation.** Steve Muchnick.
    (strength: optimization)

CS 4120 Introduction to Compilers

5

---

# Work

- Homeworks: 5, 35% total
  - 6-8% each

- Programming Assignments: 7, 65%
  - 7-13% each

- Exams: None

CS 4120 Introduction to Compilers

6

## Academic integrity

• Taken seriously.

• Do your own (or your group's) work.

• Report who you discussed homework with (whether student in class or not).

• Feel free to share test inputs on Piazza.

7

## Homeworks

• Three assignments in first half of course; two homeworks in second half

• **Not** done in groups—you may discuss with others but do your own work
– Report who you discussed homework with

CS 4120 Introduction to Compilers

8

## Projects

• Seven programming assignments
• Implementation language: Java
  – or anything you can compile to an executable jar
• Groups of 3-4 students
  – same group for entire class (ordinarily)
  – same grade for all (ordinarily)
  – workload and success in this class depend on working and planning well with your group. Be a good citizen.
  – tell us **early** if you are having problems.
• End of this class: some time to form groups
  – create your group on CMS for PA1.
  – contact us if you are having trouble finding a group.

CS 4120 Introduction to Compilers

9

## Assignments

• Due at midnight on due date
• Late homeworks, programming assignments increasingly penalized
  – 1 day: 5%, 2 days: 15%, 3 days: 30%, 4 days: 50%
  – weekend = 1 day
  – Extensions often granted, but must be approved 2 days in advance
• Projects submitted via CMS
• Solutions available via CMS

CS 4120 Introduction to Compilers

10

## Why take this course?

• CS 4120 is an elective course
• Expect to learn:
  – practical applications of theory, algorithms, data structures
  – parsing
  – deeper understanding of what code is
  – how high-level languages are implemented
  – a little programming-language semantics
  – Intel x86 architecture, Java
  – how programs really execute on computers
  – how to be a better programmer (esp. in groups)

CS 4120 Introduction to Compilers

11

## What are Compilers?

• Translators from one representation of program code to another
• Old: high-level source code to machine language (object code)
• Modern:
  – High-level to mid-level (Java to bytecode)
  – Mid-level to low-level (bytecode to x86)

CS 4120 Introduction to Compilers

12

2

## Source Code

- Source code: optimized for human readability
  - expressive: matches human notions of grammar
  - redundant to help avoid programming errors
  - computation possibly not fully determined by code

```
int expr(int n)
{
    int d;
    d = 4 * n * n * (n + 1) * (n + 1);
    return d;
}
```

CS 4120 Introduction to Compilers

13

## Machine code

- Optimized for hardware
  - Redundancy, ambiguity reduced
  - Information about intent and reasoning lost
  - Assembly code ≈ machine code

```
expr:
    pushl   %ebp                    55
    movl    %esp, %ebp              89 e5
    subl    $4, %esp                83 ec 04
    movl    8(%ebp), %eax           8b 45 08
    movl    %eax, %edx              89 c2
    imull   8(%ebp), %edx           0f af 55 08
    movl    8(%ebp), %eax           8b 45 08
    incl    %eax                    40
    imull   %eax, %edx              0f af d0
    movl    8(%ebp), %eax           8b 45 08
    incl    %eax                    40
    imull   %edx, %eax              0f af c2
    sall    $2, %eax                c1 e0 02
    movl    %eax, -4(%ebp)          89 45 fc
    movl    -4(%ebp), %eax          8b 45 fc
    leave                           c9
            ret                     c3
```

14

## Example (Output assembly code)

Unoptimized Code              Optimized Code

```
expr:                         expr:
    pushl   %ebp                  pushl   %ebp
    movl    %esp, %ebp            movl    %esp, %ebp
    subl    $4, %esp              movl    8(%ebp), %edx
    movl    8(%ebp), %eax         movl    %edx, %eax
    movl    %eax, %edx            imull   %edx, %eax
    imull   8(%ebp), %edx         incl    %edx
    movl    8(%ebp), %eax         imull   %edx, %eax
    incl    %eax                  imull   %edx, %eax
    imull   %eax, %edx            sall    $2, %eax
    movl    8(%ebp), %eax         leave
    incl    %eax                  ret
    imull   %edx, %eax
    sall    $2, %eax
    movl    %eax, -4(%ebp)
    movl    -4(%ebp), %eax
    leave
    ret
```

CS 4120 Introduction to Compilers

15

## How to translate?

- Source-code and machine-code mismatch

- Goals:
  - source-level expressiveness for task
  - best performance for concrete computation
  - reasonable translation efficiency (< $O(n^3)$)
  - maintainable compiler code

CS 4120 Introduction to Compilers

16

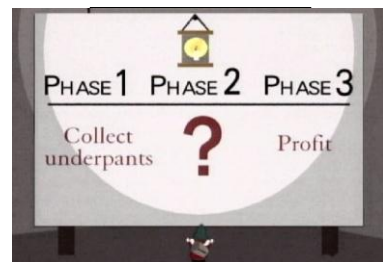## How to translate correctly?

- Programming languages describe computation precisely
- Therefore: translation can be precisely described (a compiler can be correct)
- Correctness is very important!
  - hard to debug programs with broken compiler…
  - non-trivial: programming languages are expressive
  - implications for development cost, security
  - this course: techniques for building correct compilers
  - some compilers have been **proven** correct!
    [X. Leroy, Formal Certification of a Compiler Back End, POPL '06]

CS 4120 Introduction to Compilers
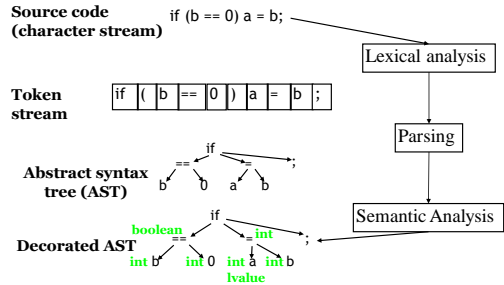
17

## How to translate effectively?



PHASE 1   PHASE 2   PHASE 3

Collect
underpants      ?      Profit

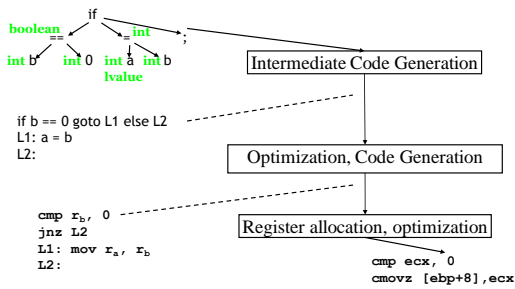CS 4120 Introduction to Compilers

18

3

## Idea: translate in steps

- Compiler uses a series of different **program representations**.
- Intermediate representations that are good for program manipulations of various kinds (analysis, optimization, code generation).
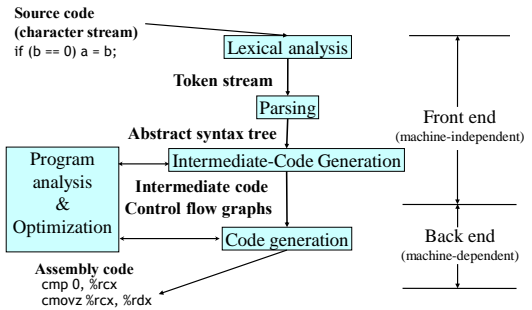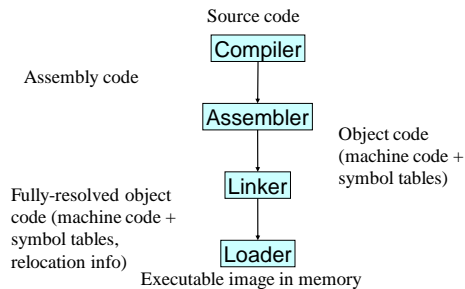
CS 4120 Introduction to Compilers

19

## Compilation in a Nutshell 1

**Source code (character stream)**  if (b == 0) a = b;



Lexical analysis

**Token stream**   `if ( b == 0 ) a = b ;`

Parsing

**Abstract syntax tree (AST)**

Semantic Analysis

**Decorated AST**

CS 4120 Introduction to Compilers

20

## Compilation in a Nutshell 2

Intermediate Code Generation

```
if b == 0 goto L1 else L2
L1: a = b
L2:
```

Optimization, Code Generation

```
cmp r_b, 0
jnz L2
L1: mov r_a, r_b
L2:
```

Register allocation, optimization

```
cmp ecx, 0
cmovz [ebp+8],ecx
```

CS 4120 Introduction to Compilers

21

## Simplified Compiler Structure

**Source code (character stream)**
if (b == 0) a = b;

Lexical analysis

**Token stream**

Parsing

**Abstract syntax tree**

Intermediate-Code Generation

Program analysis & Optimization

**Intermediate code Control flow graphs**

Code generation

**Assembly code**
cmp 0, %rcx
cmovz %rcx, %rdx

Front end (machine-independent)

Back end (machine-dependent)

22

## Even bigger picture

Source code

Compiler

Assembly code

Assembler

Object code (machine code + symbol tables)

Linker

Fully-resolved object code (machine code + symbol tables, relocation info)

Loader

Executable image in memory

CS 4120 Introduction to Compilers

23

## CubeX

```
interface List<E> {
    fun elements() : Iterable<E>;
}
class Nil() extends List<Nothing> {
    fun elements() : Iterable<Nothing>
        { return []; }
}
class Cons<E>(E head, List<E> tail)
        extends List<Nothing> {
    elems := [head] ++ tail.elements();
    fun elements() : Iterable<E> {return elems;}
}
```

24

## CubeX

- Object-Oriented

- Generics

- Pure except for non-termination

- Memory managed

25

## Project

- Compile CubeX to C
  – With optimizations
- Choice of Extension, such as
  – Compile to x86
  – Variance and inference
  – Recursive inheritance
    and F-bounded polymorphism
  – Continuation-based custom iterables
  – Iterable comprehensions and liftings

26