# CS 4120 / 4121

Introduction to Compilers

Fall 2009

Andrew Myers

Lecture 1: Overview

---

# 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

---

# Course Information

- MWF 1:25- 2:15$_{PM}$ in Phillips 203
- Instructor: Andrew Myers
- Teaching Assistants: Anthony Jawad
- E-mail: cs4120-l@cs.cornell.edu
- Web page: http://www.cs.cornell.edu/courses/cs4120
- Newsgroup: cornell.class.cs4120

---

# Academic integrity

- Taken seriously.
- Do your own (or your group's) work.
- Report who you discussed homework with (whether student in class or not).

# CS 4121 is required!

– most coursework is in the project

# Textbooks

- Required text
  - **Modern Compiler Implementation in Java.** Andrew Appel.
  - on reserve in Engineering Library
- Optional texts
  - **Compilers—Principles, Techniques and Tools.** Aho, Lam, Sethi and Ullman (The Dragon Book)
  - **Advanced Compiler Design and Implementation.** Steve Muchnick.

# Work

- Homeworks: 4, 20% total
  - 5/5/5/5
- Programming Assignments: 6, 50%
  - 5/7/8/10/10/10
- Exams: 2 prelims, 30%
  - 15/15
  - No final exam

# Homeworks

- Three assignments in first half of course; one homework in second half
- **Not** done in groups—you may discuss with others but do your own work
  - Report who you discussed homework with

# Projects

- Six programming assignments
- Implementation language: Java
  - talk to us if your group wants to use something else (e.g., OCaml)
- 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 beginning of class
- 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

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
- Typically: high-level source code to machine language (object code)
- Not always:
  - Java compiler: Java to interpretable bytecodes
  - Java JIT: bytecode to executable image

CS 4120 Introduction to Compilers

12

# 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;
}
```

# 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
```

# Example (Output assembly code)

Unoptimized Code

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

Optimized Code

```
expr:
    pushl   %ebp
    movl    %esp, %ebp
    movl    8(%ebp), %edx
    movl    %edx, %eax
    imull   %edx, %eax
    incl    %edx
    imull   %edx, %eax
    imull   %edx, %eax
    sall    $2, %eax
    leave
    ret
```
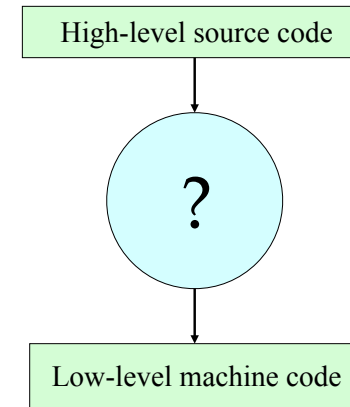
# How to translate?

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

## 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!
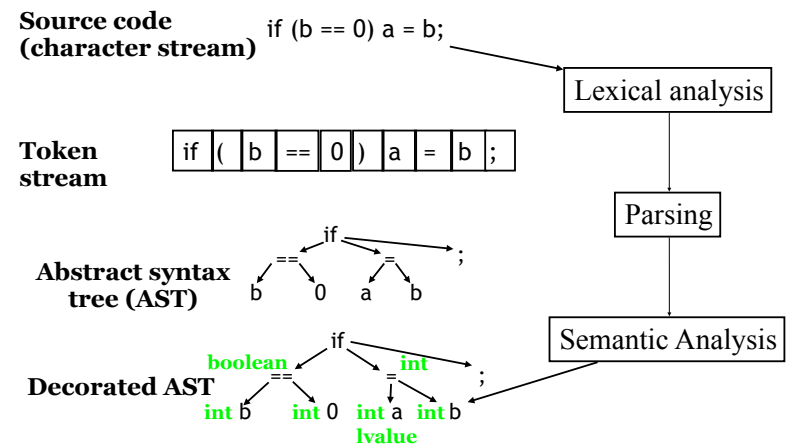    [X. Leroy, Formal Certification of a Compiler Back End, POPL '06]

CS 4120 Introduction to Compilers

17

## How to translate effectively?



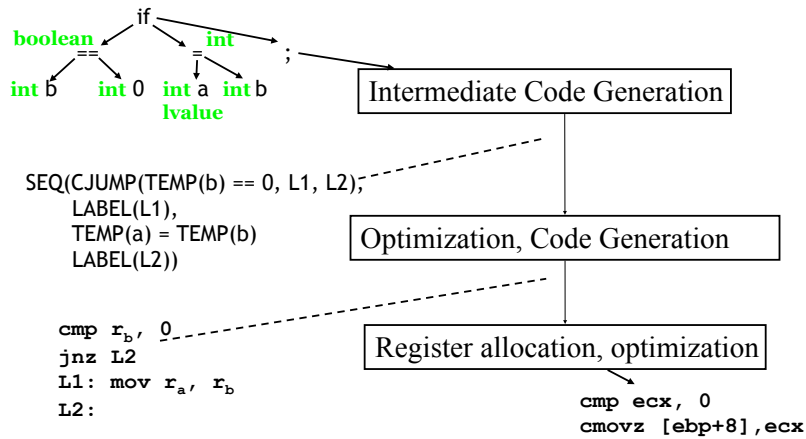CS 4120 Introduction to Compilers

18

## 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).
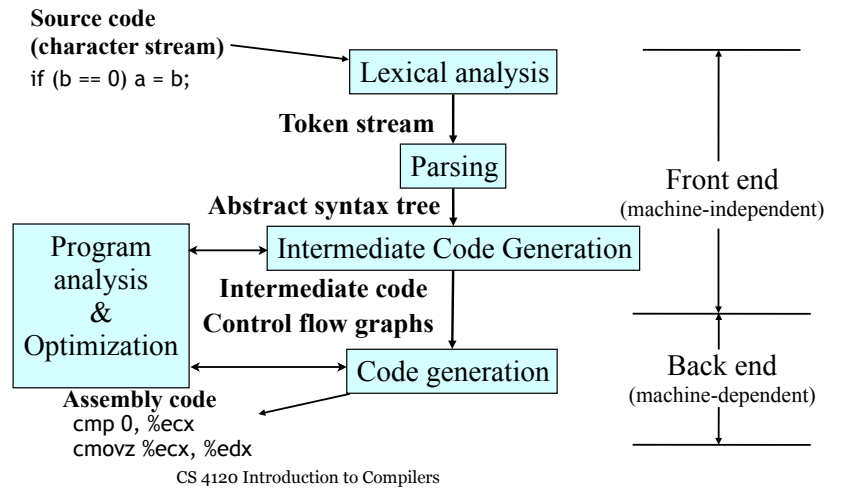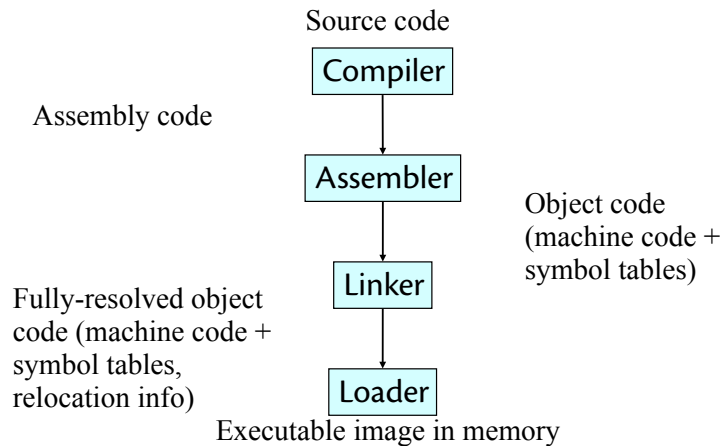
CS 4120 Introduction to Compilers

19

## Compilation in a Nutshell 1



CS 4120 Introduction to Compilers

20

## Compilation in a Nutshell 2

```
           if
boolean  ==   int    ;
int b int 0  int a int b
         =
        lvalue
```

Intermediate Code Generation

```
SEQ(CJUMP(TEMP(b) == 0, L1, L2);
    LABEL(L1),
    TEMP(a) = TEMP(b)
    LABEL(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
```

---

## Simplified Compiler Structure

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

→ Lexical analysis

**Token stream**

Parsing

**Abstract syntax tree**

Program analysis & Optimization ↔ Intermediate Code Generation

**Intermediate code**
**Control flow graphs**

Code generation

**Assembly code**
cmp 0, %ecx
cmovz %ecx, %edx

Front end
(machine-independent)

Back end
(machine-dependent)

---

## 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

---

## Schedule

- Detailed schedule on web page, with links

| | |
|---|---|
| Lexical analysis and parsing: | 6 |
| Semantic analysis: | 5 |
| Intermediate code: | 4 |
| Prelim #1 | |
| Code generation: | 3 |
| Separate compilation and objects: | 4 |
| Optimization: | 8 |
| Prelim #2 | |
| Run-time, link-time support: | 2 |
| Advanced topics: | 7 |

# First step: Lexical Analysis

Source code
(character stream)

**Lexical analysis**

Token stream

Parsing

Abstract syntax tree

Intermediate Code Generation

Intermediate code

Code generation

Assembly code

---

# What is Lexical Analysis?

- Converts character stream to token stream of pairs $\langle token\ type, attribute \rangle$

```
if  (x1 * x2<1.0) {
    y = x1;
}
```

| i | f | | ( | x | 1 | | * | | x | 2 | < | 1 | . | 0 | ) | { | \n |

if  (  Id: **x1**  *  Id: **x2**  <  Num: *1.0*  )  {  Id: **y**

---

# Token stream

- Gets rid of whitespace, comments
- Only $\langle$ Token type, attribute $\rangle$:
  - $\langle$ Id, "x" $\rangle$, $\langle$ Float, *1.0e0* $\rangle$
- Token location preserved for debugging, run-time/compile-time error messages (source file, line number, character posn...)
  - $\langle$ Id, "x", "Main.java", 542 $\rangle$

- Issues:
  - how to specify tokens
  - how to implement tokenizer/lexer