


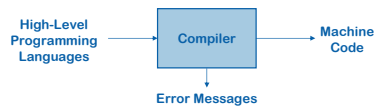
# Compilers: Summary and Wrapup

**CS 471**  
**December 5, 2007**



## What You've Learned This Semester

**What happens when you compile your code**



**How to implement a compiler**

- What is challenging, time consuming
- Available tools and their other applications

**How to apply theory to practice (with less direction)**  
**How to break a large problem down into subcomponents**

## Goals of a Compiler

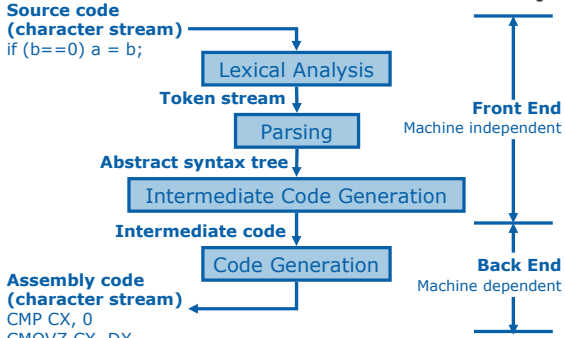
**A compiler's job is to**

- Lower the abstraction level
- Eliminate overhead from language abstractions
- Map source program onto hardware efficiently
  - Hide hardware weaknesses, utilize hardware strengths
- Equal the efficiency of a good assembly programmer

**Optimizing compilers should *improve* the code**

- Performance\*
- Code size
- Security
- Reliability
- Power consumption

## Simplified Compiler Structure



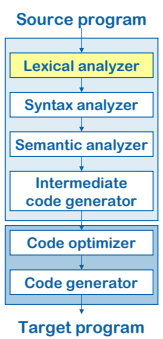
**Source code (character stream)**  
if (b==0) a = b;

**Front End**  
Machine independent

**Back End**  
Machine dependent

**Assembly code (character stream)**  
CMP CX, 0  
CMOVZ CX, DX

## Phases of a Compiler



### Lexical Analyzer

- Group sequence of characters into lexemes (keywords, identifiers, constants)
- Makes use of the theory of regular languages and finite state machines
- Lex and Flex are tools that construct lexical analyzers from regular expression specifications

## Regular expressions

**Language** – set of strings  
**String** – finite sequence of symbols  
**Symbols** – taken from a finite alphabet

**Specify languages using regular expressions**

Symbol	a	one instance of a
Epsilon	$\epsilon$	empty string
Alternation	$R \mid S$	string from either L(R) or L(S)
Concatenation	$R \cdot S$	string from L(R) followed by L(S)
Repetition	$R^*$	0 or more strings from L(R)

## Finite Automata

**Automaton (DFA) can be represented as:**

- A transition table
 

	"	non-"
0	1	error
1	2	1
2	error	error
- A graph

## Language

**Each string is accepted or rejected**

- Starting in the start state
- Automaton follows one edge for every character (edge must match character)
- After n-transitions for an n-character string, if final state then accept

**Language: set of strings that the FSA accepts**

## Lex: A Lexical Analyzer Generator

- Lex produces a C program from a lexical specification
- <http://www.epaperpress.com/lexandyacc/index.html>

```

%%
DIGITS [0-9]+
ALPHA [A-Za-z]
CHARACTER {ALPHA}|_
IDENTIFIER {ALPHA}{CHARACTER}{DIGITS}*
%%
if
{IDENTIFIER}          {return IF; }
{DIGITS}              {return ID; }
{[0-9]+"."[0-9]*}|{[0-9]*"."[0-9]+} {return NUM; }
.                    {error(); }
  
```

## Phases of a Compiler

Source program

Parser

- Convert a linear structure - sequence of tokens - to a hierarchical tree-like structure - an AST
- The parser imposes the syntax rules of the language
- Work should be linear in the size of the input (else unusable) -> type consistency cannot be checked in this phase
- Deterministic context free languages and pushdown automata for the basis
- Bison and yacc allow a user to construct parsers from CFG specifications

## Context-Free Grammar Terminology

- Terminals**
  - Token or  $\epsilon$
- Non-terminals**
  - Syntactic variables
- Start symbol**
  - A special nonterminal is designated (S)
- Productions**
  - Specify how non-terminals may be expanded to form strings
  - LHS: single non-terminal, RHS: string of terminals or non-terminals
- Vertical bar is shorthand for multiple productions

$$S \rightarrow (S) S$$

$$S \rightarrow \epsilon$$

## Shift-Reduce Parsing

**Bottom-up parsing uses two kinds of actions: Shift and Reduce**

**Shift:** Move  $|$  one place to the right

- Shifts a terminal to the left string

$$E + ( | int ) \Rightarrow E + ( int | )$$

**Reduce:** Apply an inverse production at the right end of the left string

- If  $E \rightarrow E + ( E )$  is a production, then

$$E + ( E + ( E ) | ) \Rightarrow E + ( E | )$$



## Phases of a Compiler

**Intermediate Code Generation**

- Makes it easy to port compiler to other architectures (e.g. Pentium to MIPS)
- Can also be the basis for interpreters (such as in Java)
- Enables optimizations that are not machine specific

18 CS 471 - Fall 2007

## Intermediate Code

- Abstract machine code – (Intermediate Representation)
- Allows machine-independent code generation, optimization

19 CS 471 - Fall 2007

## Phases of a Compiler

**Intermediate Code Optimization**

- Constant propagation, dead code elimination, common sub-expression elimination, strength reduction, etc.
- Based on dataflow analysis – properties that are independent of execution paths

20 CS 471 - Fall 2007

## Analysis and Transformation

**Most optimizations require some global understanding of program flow**

- Moving, removing, rearranging instructions

**Achieve understanding by discovering the control flow of the procedure**

- What blocks follow/are reachable from other blocks
- Where loops exist (focus optimization efforts)
- We call this **Control-Flow Analysis**

**Also, find and connect definitions and uses of variables**

- We call this **Data-Flow Analysis**

**All approaches to control-flow analysis require**

- Identification of basic blocks
- Construction of a flow graph of the procedure

21 CS 471 - Fall 2007

## When to Apply Optimization

HIR	AST	Inlining Specialization
	IR	Constant folding Constant propagation Value numbering
MIR	Canonical IR	Dead code elimination Loop-invariant code motion Common sub-expression elimination
	Abstract Assembly	Strength reduction Branch prediction/optimization
LIR	Assembly	Register allocation Loop unrolling Cache optimization

22 CS 471 - Fall 2007

## Scope of Optimization

**Local (or single block)**

- Confined to straight-line code
- Simplest to analyze

**Intraprocedural (or global)**

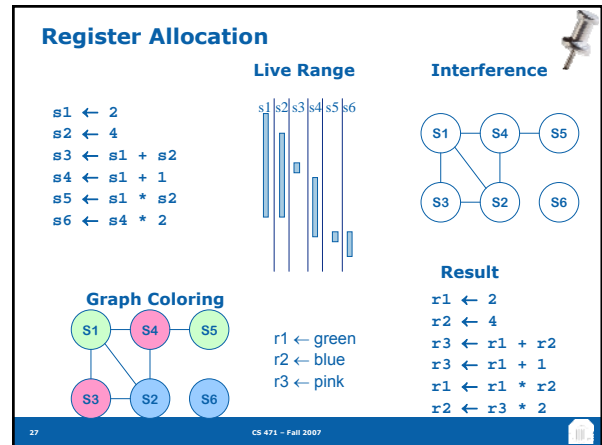
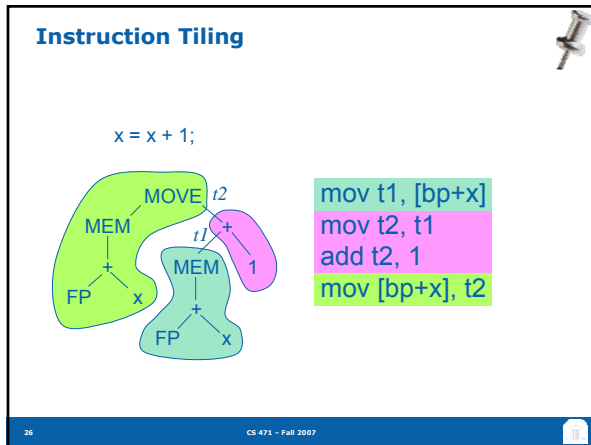
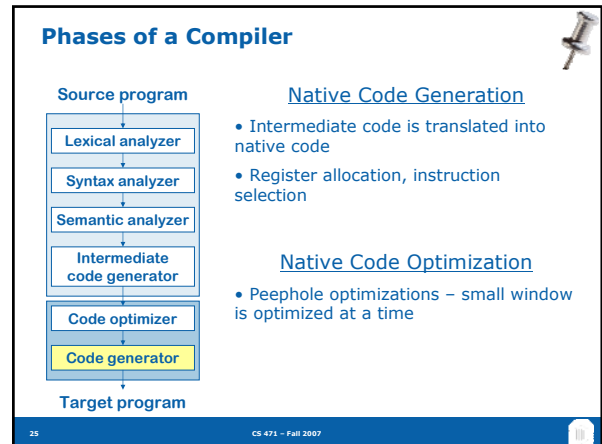
- Consider the whole procedure

**Interprocedural (or whole program)**

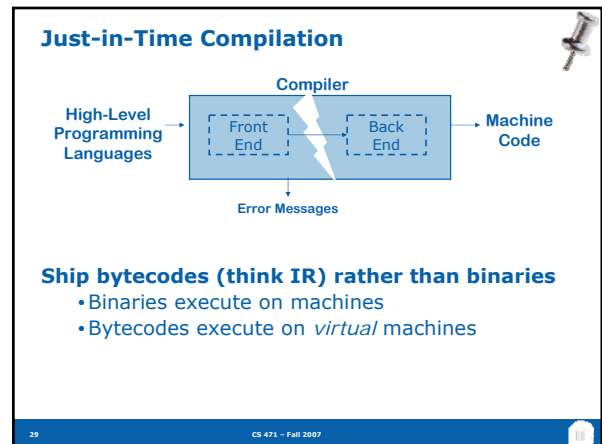
- Consider the whole program

23 CS 471 - Fall 2007

Flag	Description
-O0	Barely any transformations, just code generation. Target code can be debugged with no information loss.
-O1	Some transformations that preserve execution ordering. Debuggability of the generated code is hardly affected. User variables should not disappear and function inlining is not done.
-O2	More aggressive transformations. May affect execution ordering and usually provide faster code. Debuggability may be somewhat compromised by disappearing user variables and function bodies.
-O3	Very aggressive transformations that may or may not provide better code. Some semantics may be modified. Execution order completely distorted. Debuggability is seriously compromised.
-Os	Optimize for size. Enables transformations that reduce generated code size. May sometimes improve application performance because there is less code to execute. May lead to reduced memory footprints which may produce fewer page faults.
-O4	There is no -O4. Anything above -O3 is treated as -O3.



- ### Alternatives to the Traditional Model
- Static Compilation**  
All work is done "ahead-of-time"
  - Just-in-Time Compilation**  
Postpone some compilation tasks
  - Multiversioning and Dynamic Feedback**  
Include multiple options in binary
  - Dynamic Binary Optimization**  
Traditional compilation model  
Executables can adapt



## The Big Picture

### We now know:

- All of the components of a compiler
- What needs to be done statically vs. dynamically
- The potential impact of language or architecture changes
- Why Java moved the "back-end" to run time

30

CS 471 - Fall 2007

## Compiler Wars: Deliverables and Rules

CS 471  
December 10, 2007



## 1. Robustifying your Compiler

(For real this time)

Sample outputs are online

[www.cs.virginia.edu/kim/courses/cs471/project/](http://www.cs.virginia.edu/kim/courses/cs471/project/)

### Contains:

- test[0-9]+.tig.out
- merge.tig.out
- queens.tig.out

32

CS 471 - Fall 2007

## 2. Generate Test Cases

- A test case that may break others' lexical analyzers  
**kh7fe.lex.tig**
- A test case that may break others' parsers  
**kh7fe.parse.tig**
- A test case that may break others' typecheckers  
**kh7fe.type.tig**
- A test case that may break others' irGen  
**kh7fe.irgen.tig**
- A README file - describe expected output for each test file; what you expect to break in each case; why it's unique  
**kh7fe.README.tests**

33

CS 471 - Fall 2007

## 3. Submit your Compilers

- Your tiger.lex file (We will build your lexical analyzer)  
**kh7fe.lex**
- Your grammar (with semantic actions)  
**kh7fe.grm**
- Your typechecker (Input: file; Output: type errors)  
**kh7fe.typecheck**
- Your IR generator (Input: file; Output: type errors and (if applicable) IR tree)  
**kh7fe.irgen**
- README file describing all of the issues you corrected between PA7 and compiler wars  
**kh7fe.README.compiler**
- A tar.gz file containing everything in your PA7 directory (assuming it has been updated since PA7)  
**kh7fe.finalCompiler.tar.gz**

34

CS 471 - Fall 2007

## Then What?

Submit everything

**BEFORE SUNDAY DEC 9 @ 9PM FIRM**  
(enforced by toolkit)

Show up here on Monday Dec 10 @ 2PM for ...

COMPILER WARS!



35

CS 471 - Fall 2007

## How Will It Work?

### Bracketed progression (like the ACC tournament)

- Round 1: Lexical Analysis
- Round 2: Parsing & AST Generation
- Round 3: Type Checking
- Round 4: IR Generation

**Each round: Full bracket progression to determine a winner (single elimination)**

**Each game: Compare test cases passed with your opponent**

- 51 given test cases
- 2 custom cases

**Tie breaker: Cleanest code (as determined by judges)**

36

CS 471 - Fall 2007

## Prize structure

### Prizes for:

- Best lexical analyzer
- Best parser
- Best typechecker
- Best ir generator
- Best overall
- Most improved compiler
- Best test case

### Prizes:

- Corresponding PA late day reprieve (1 day)
- Grab bag

37

CS 471 - Fall 2007

## PA8 Grade

**50%**

**submitting all PA8 materials *on time* and showing up for compiler wars**

**50%**

**completeness of the final product (PA2-PA7)**

**Extra Office Hours  
Sunday 6-7 PM Olsson 209**

38

CS 471 - Fall 2007

## Rest of Class Time Today...

1) Please fill out the official evaluations

2) Please send me anonymous feedback if you have detailed suggestions for doing things differently the next time I teach this course

- Topic emphasis, e.g. less time on X more time on Y
- Homework vs. projects vs. tests
- Lab timeslots?
- Group vs. individual projects?

**Thanks for a Great Semester!**

39

CS 471 - Fall 2007