Overview

Course Objectives

To learn the process of translating a modern high-level language to executable code.

- Learn the fundamental techniques from lectures, text book and exercises from the book.
- Apply these techniques in practice to construct *a fully working compiler* for a non-trivial subset of Java called "Decaf".

In the end, you should be able to compile small Java-like programs with your compiler, and see it actually work!

Compilers	Introduction	CSE 304/504	1 / 33
	Overview		
What is a Compiler?			

- Programming problems are easier to solve in *high-level languages*
 - High-level languages are closer to the problem domain
 - E.g. Java, Python, SQL, Tcl/Tk, ...
- Solutions have to be executed by a machine
 - Instructions to a machine are specified in a language that reflects to the cycle-by-cycle working of a processor
- Compilers are the bridges:
 - Software that *translates* programs written in high-level languages to efficient executable code.

An Example

int gcd(int m, int n) gcd: { pushl %ebp if (m == 0)movl %esp,%ebp return n; cmpl \$0,8(%ebp) else if (m > n)jne .L2 return gcd(n, m); movl 12(%ebp),%eax else jmp .L1 return gcd(n%m, m); .align 16 } jmp .L3 .align 16 .L2: movl 8(%ebp),%eax cmpl %eax,12(%ebp) jge .L4 movl 8(%ebp),%eax pushl %eax . . .

Compilers Introduction CSE 304/504 3 / 33

Overview

Example (contd.)

gcd:			cltd
	pushl %ebp		idivl %esi
	movl %esp,%ebp		movl %edx,%ebx
	pushl %esi	.L14:	
	pushl %ebx		movl %esi,%ecx
	movl 8(%ebp),%esi		movl %ebx,%esi
	movl 12(%ebp),%ebx		movl %ecx,%ebx
.L11:			jmp .L11
	testl %esi,%esi		.align 16
	jne .L8	.L13:	
	movl %ebx,%eax		leal -8(%ebp),%esp
	jmp .L13		popl %ebx
	.align 16		popl %esi
.L8:	-		movl %ebp,%esp
	cmpl %ebx,%esi		popl %ebp
	jg .L14		ret
	movl %ebx,%eax		

Requirements

- In order to translate statements in a language, one needs to understand both
 - the *structure* of the language: the way "sentences" are constructed in the language, and
 - the *meaning* of the language: what each "sentence" stands for.
- Terminology:
 - Structure \equiv **Syntax**
 - Meaning = Semantics

Compilers	Introduction	CSE 304/504	5 / 33
	Overview		
Translation Strategy			

Classic Software Engineering Problem

- **Objective:** Translate a program in a high level language into *efficient* executable code.
- **Strategy:** Divide translation process into a series of phases Each phase manages some particular aspect of translation.

Interfaces between phases governed by specific intermediate forms.

Translation Process



Translation Steps

- Syntax Analysis Phase: Recognizes "sentences" in the program using the *syntax* of the language
- Semantic Analysis Phase: Infers information about the program using the *semantics* of the language
- Intermediate Code Generation Phase: Generates "abstract" code based on the syntactic structure of the program and the semantic information from Phase 2.
- **Optimization Phase:** Refines the generated code using a series of *optimizing* transformations.
- Final Code Generation Phase: Translates the abstract intermediate code into specific machine instructions.

Structure of a Compiler: an Analogy

Syntax-Directed Translation: the *structure* (syntax) of a sentence in a language is used to give it a *meaning* (semantics).

- Die CD-ROM Informatikunterricht liefert alle wichtigen Programme und Bücher, die für den Unterricht oder zur Bearbeitung von Hausaufgaben notwendig sind.
- Green wire connect after first cut white also red wire.
- He sailed the coffee out of the leaf.
- This sentence has four words.

Compilers	Introduction	CSE 304/504	9 / 33
	Syntax-Directed Translation		
Syntax			

Defining and Recognizing Sentences in a Language

- Layered approach
- Alphabet: defines allowed symbols
- Lexical Structure: defines allowed words
- Syntactic Structure: defines allowed sentences

We will later associate *meaning* with sentences (semantics) based on their syntactic structure.

Formal Language Specification

Solid theoretical results applied to a practical problem.

- Declarative vs. Operational Notations
- Declarative notation is used to define a language
 - Defines precisely the set of allowed objects (words/sentences)
 - Examples: Regular expressions, Grammars.
- Operational notation is used to recognize statements in a language
 - Defines an algorithm for determining whether or not a given word/sentence is in the language
 - Example: Automata
- Results from theory on converting between the two notations.

Compilers	Introduction	CSE 304/504	11 / 33
Syr	ntax-Directed Translation		
Formal Languages	S		

A *language* is a set of strings over a set of symbols.

- The set of symbols of a language is called its alphabet (usually denoted by $\boldsymbol{\Sigma}.$
- Each string in the language is called a sentence.
- Parts of sentences are called phrases.

Context-Free Grammars

A well-studied notation for defining formal languages.

- A Context Free Grammar (CFG, or "grammar" unless otherwise qualified) is defined over an alphabet, called terminal symbols.
- A CFG is defined by a set of productions.
- Each production is of the form

$$X \longrightarrow \beta$$

where

- X is a single non-terminal symbol
 - representing a set of phrases in the language, and
- β is a sequence of terminal and non-terminal symbols
- Example:

```
Stmt \longrightarrow while Expr do Stmt
```

- A unique non-terminal, called the start symbol, represents the set of all sentences in the language.
- The language defined by a grammar G is denoted by $\mathcal{L}(G)$.

Compilers	Introduction	CSE 304/504	13 / 33

Syntax-Directed Translation

Example Grammar

"List of digits separated by + and - signs" (Example 2.1 in book):

$$L \longrightarrow L + L$$

$$L \longrightarrow L - L$$

$$L \longrightarrow D$$

$$D \longrightarrow 0|1|...|9$$

Derivation of 9-5+2 from L:

$$L \implies L - L$$
$$\implies D - L$$
$$\implies 9 - L$$
$$\implies 9 - L + L$$
$$\implies 9 - D + L$$
$$\implies 9 - 5 + L$$
$$\implies 9 - 5 + D$$
$$\implies 9 - 5 + 2$$

Parse Trees

Pictorial representation of derivations



Note: one parse tree may correspond to multiple derivations!

Compilers	Introduction	CSE 304/504	15 / 33
	Syntax-Directed Translation		

Ambiguity

A grammar is ambiguous if some sentence in the language has more than one parse tree.



Associativity and Precedence

- $9-5+2 \equiv (9-5)+2$
- $9-5-2 \equiv (9-5)-2$
- $9+5+2 \equiv (9+5)+2$
- "+" and "-" usually have the same precedence and are left-associative.

i.e. the second parse tree in the previous slide is the "correct" one

• The grammar can be changed to reflect the associativity and precedence:



Compilers

Introduction

CSE 304/504 17 / 33

Syntax-Directed Translation

Syntax-Directed Translation Schemes

- A notation that attaches "program fragments" (also called actions) to productions in a grammar.
- The intuition is, whenever a production is used in recognizing a sentence, the corresponding action will be taken.
- Example:

$$L \longrightarrow L + D \{add\}$$

$$L \longrightarrow L - D \{sub\}$$

$$L \longrightarrow D$$

$$D \longrightarrow 0 \{push \ 0\}$$

$$D \longrightarrow 1 \{push \ 1\}$$

$$\vdots$$

Syntax-Directed Translation

- Actions can be seen as "additional leaves" introduced into a parse tree.
- Reading the actions left-to-right in the tree gives the "translation".

Example:



Syntax-Directed Translation

Grammars for Language Specification

- The language (i.e. set of allowed strings) of most programming languages can be specified using CFGs.
- The grammar notation may be tedious for some aspects of a language.
- For instance, an integer is defined by a grammar of the following form:

$$I \longrightarrow P | + P | - P$$

$$P \longrightarrow D P$$

$$P \longrightarrow D$$

$$D \longrightarrow 0|1|...|9$$

- For simpler fragments, the notation of regular expressions may be used.
- I = (+|-)?[0-9]+

Syntax Analysis in Practice

- Usually divided into Lexical Analysis followed by Parsing.
- Lexical Analysis:
 - A lexical analyzer converts a stream of characters into a stream of tokens.
 - Each token has a *name* (associated with terminal symbols) and a *value* (also called its attribute).
 - A lexical analyzer is specified by a set of regular expression patterns and actions that are executed when the patterns are matched.
- Parsing:
 - A parser converts a stream of tokens into a tree.
 - Parsing uncovers the *structure* of a sentence in the language.
 - Parsers are specified by grammars (actually, by translation schemes which are sets of productions associated with actions).

Compilers	Introduction	CSE 304/504	21 / 33
	Phases of Translation		
Syntax Analysi	is		
Source Program	Lexical Analysis Token Stream		

Parsing

Abstract Syntax Tree

Lexical Analysis

First step of syntax analysis

- **Objective:** Convert the *stream of characters representing input program* into a sequence of *tokens*.
- Tokens are the "words" of the programming language.
- Examples:
 - The sequence of characters "static int" is recognized as two tokens, representing the two words "static" and "int".
 - The sequence of characters "*x++" is recognized as three tokens, representing "*", "x" and "++".

	Compilers	Introduction	CSE 304/504	23 / 33
_		Phases of Translation		

Parsing

Second step of syntax analysis

- **Objective:** Uncover the *structure* of a sentence in the program from a stream of *tokens*.
- For instance, the phrase "x = +y", which is recognized as four tokens, representing "x", "=" and "+" and "y", has the structure =(x, +(y)), i.e., an assignment expression, that operates on "x" and the expression "+(y)".
- **Output:** A *tree* called *abstract syntax tree* that reflects the structure of the input sentence.

Abstract Syntax Tree (AST)

- Represents the syntactic structure of the program, hiding a few details that are irrelevent to later phases of compilation.
- For instance, consider a statement of the form: "if (m == 0) S1 else S2" where S1 and S2 stand for some block of statements. A possible AST for this statement is:





Type Checking

A instance of "Semantic Analysis"

- **Objective:** Decorate the AST with semantic information that is necessary in later phases of translation.
- For instance, the AST



Intermediate Code Generation

- **Objective:** Translate each sub-tree of the decorated AST into *intermediate code*.
- Intermediate code hides many machine-level details, but has instruction-level mapping to many assembly languages.
- Main motivation for using an intermediate code is *portability*.

Phases of Translation

Intermediate Code Generation, an Example





Code Optimization

- **Objective:** Improve the time and space efficiency of the generated code.
- Usual strategy is to perform a series of transformations to the intermediate code, with each step representing some efficiency improvement.
- *Peephole optimizations*: generate new instructions by combining/expanding on a small number of consecutive instructions.
- *Global optimizations*: reorder, remove or add instructions to change the structure of generated code.

Phases of Translation

Code Optimization, an Example

loadint m loadimmed 0 intequal jmpz .L1 jmp .L2 .L1: code for S1 jmp .L3 .L2: code for S2 jmp .L3 .L3:





Final Code Generation

- **Objective:** Map instructions in the intermediate code to specific machine instructions.
- Supports standard object file formats.
- Generates sufficient information to enable symbolic debugging.

Phases of Translation

Final Code Generation, an Example

 loadint m
 ⇒
 m

 jmpnz .L2
 t

 .L1:
 j

 code for S1
 .L1:

 jmp .L3
 .L1:

 .L2:
 j

 code for S2
 .L2:

 .L3:
 .L3:

```
movl 8(%ebp), %esi
testl %esi, %esi
jne .L2
.L1:
    .... code for S1
jmp .L3
.L2:
    .... code for S2
.L3:
```

Compilers

Introduction

CSE 304/504 33 / 33