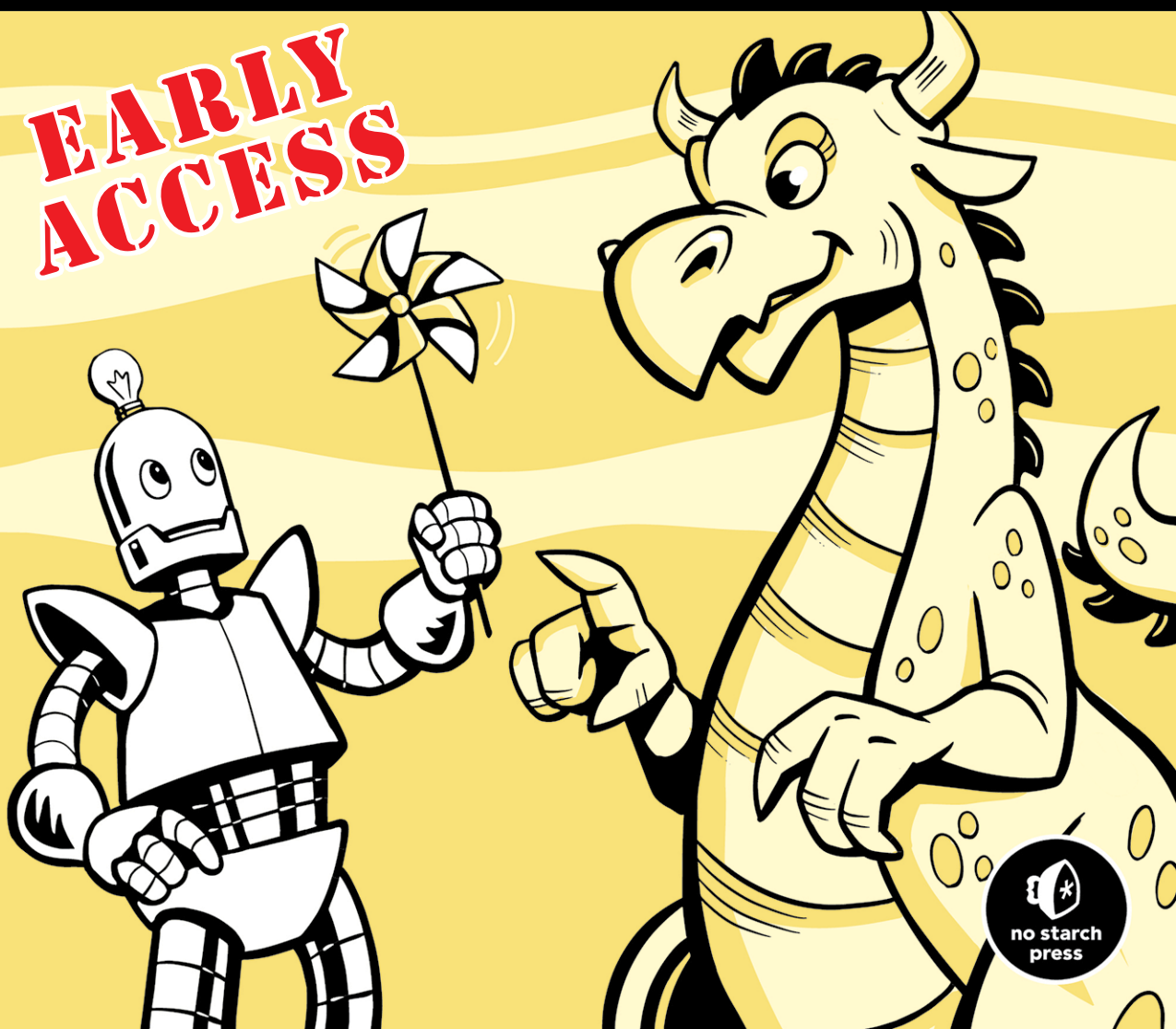


# WRITING A C COMPILER

NORA SANDLER

**EARLY  
ACCESS**



# **NO STARCH PRESS EARLY ACCESS PROGRAM: FEEDBACK WELCOME!**

The Early Access program lets you read significant portions of an upcoming book while it's still in the editing and production phases, so you may come across errors or other issues you want to comment on. But while we sincerely appreciate your feedback during a book's EA phase, please use your best discretion when deciding what to report.

At the EA stage, we're most interested in feedback related to content—general comments to the writer, technical errors, versioning concerns, or other high-level issues and observations. As these titles are still in draft form, we already know there may be typos, grammatical mistakes, missing images or captions, layout issues, and instances of placeholder text. No need to report these—they will all be corrected later, during the copyediting, proof-reading, and typesetting processes.

If you encounter any errors (“errata”) you’d like to report, please fill out [this Google form](#) so we can review your comments.

# **WRITING A C COMPILER**

## **NORA SANDLER**

Early Access edition, 3/25/22

Copyright © 2022 by Nora Sandler.

ISBN 13: 978-1-7185-0042-6 (print)

ISBN 13: 978-1-7185-0043-3 (ebook)

Publisher: William Pollock

Managing Editor: Jill Franklin

Production Manager: Rachel Monaghan

Developmental Editor: Alex Freed

Production Editor: Paula Williamson

Cover Illustrator: James L. Barry

Interior Design: Octopod Studios

Compositor: Happenstance Type-O-Rama

No Starch Press and the No Starch Press logo are registered trademarks of No Starch Press, Inc. Other product and company names mentioned herein may be the trademarks of their respective owners. Rather than use a trademark symbol with every occurrence of a trademarked name, we are using the names only in an editorial fashion and to the benefit of the trademark owner, with no intention of infringement of the trademark.

All rights reserved. No part of this work may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopying, recording, or by any information storage or retrieval system, without the prior written permission of the copyright owner and the publisher.

The information in this book is distributed on an “As Is” basis, without warranty. While every precaution has been taken in the preparation of this work, neither the author nor No Starch Press, Inc. shall have any liability to any person or entity with respect to any loss or damage caused or alleged to be caused directly or indirectly by the information contained in it.

# CONTENTS

Introduction

## **PART I: THE BASICS**

Chapter 1: Introduction to Compilers

Chapter 2: Returning an Integer

Chapter 3: Unary Operators

Chapter 4: Binary Operators

Chapter 5: Logical and Relational Operators

Chapter 6: Local Variables

Chapter 7: If Statements and Conditional Expressions

Chapter 8: Compound Statements

Chapter 9: Loops

Chapter 10: Functions

Chapter 11: Static Variables

## **PART II: IMPLEMENTING TYPES**

Chapter 12: Long Integers

Chapter 13: Unsigned Integers

Chapter 14: Floating-Point Numbers

Chapter 15: Pointers

Chapter 16: Arrays and Pointer Arithmetic

Chapter 17: Characters and Strings

Chapter 18: Supporting Dynamic Memory Allocation

Chapter 19: Structures

## **PART III: OPTIMIZATIONS**

Chapter 20: Optimizing TACKY Programs

Chapter 21: Register Allocation

Conclusion: Next Steps

The chapters in **red** are included in this Early Access PDF.

# 2

## RETURNING AN INTEGER

In this chapter, you'll write a tiny compiler that can only handle the simplest possible C programs. You'll learn how to read a simple assembly program, and you'll implement four basic compiler passes that you'll keep building on for the rest of the book. Let's start by looking at the four compiler passes you'll build in this chapter.

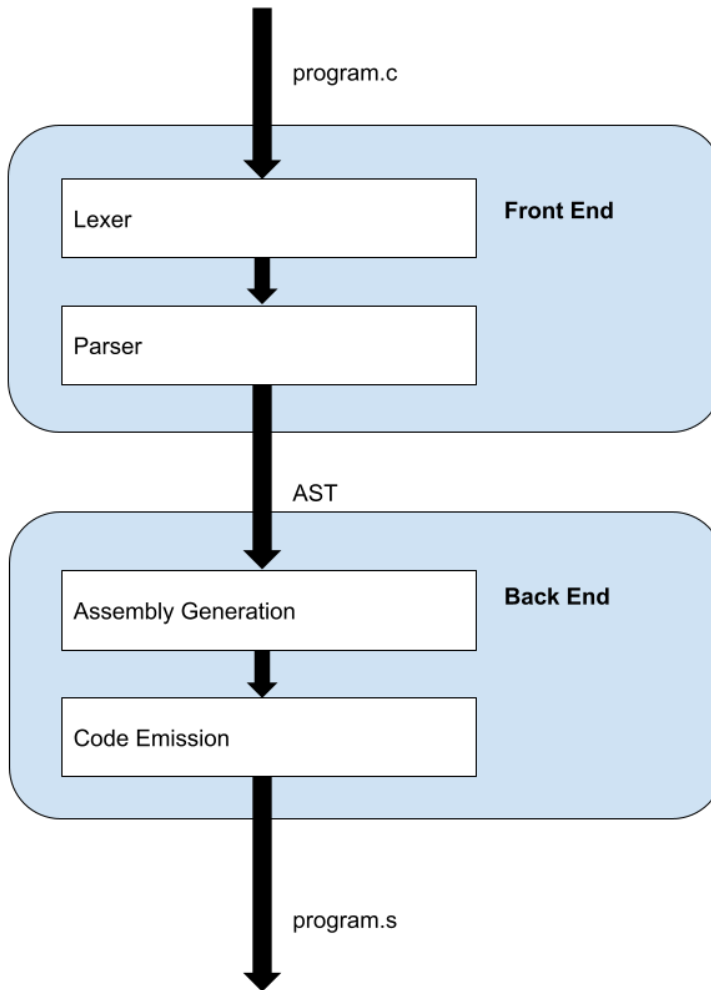


Figure 2-1: Stages of the compiler

## The Four Compiler Passes

The compiler you write in this chapter will process source code in four stages:

The *lexer* breaks up the source code into a list of *tokens*. Tokens are the smallest syntactic units of a program, and include things like delimiters, arithmetic symbols, keywords, and identifiers. If a program is like a book, tokens are like individual words.

The *parser* converts the list of tokens into an *abstract syntax tree (AST)*,

which represents the program in a form that we can easily traverse and analyze.

The *code generation* pass converts the AST into assembly. At this stage, we still represent the assembly instructions in a data structure that the compiler can understand, not as text.

The *code emission* pass writes the assembly to a file so the assembler and linker can turn it into an executable.

This is a pretty normal way of structuring a compiler, although the exact stages and intermediate representations vary. It's also overkill for this chapter; the programs you'll handle here could be compiled in just one pass! But setting up this structure now will make it easier to expand your compiler in future chapters. As you implement more language features, you'll extend these compiler stages and add a few new ones. Each chapter in the book starts with a diagram of the compiler's architecture in that chapter, including the stages you've already implemented and any you'll need to add. Figure 2-1 shows the four stages you'll implement in this chapter.

Before you start coding, let's take a quick look at how to compile C to assembly with GCC, and how to read assembly programs.

## Hello, Assembly!

The simplest possible C program looks like this:

```
1 int main() {
2     return 2;
}
```

Listing 2-1

A simple program that returns the number 2.

This program consists of a single function **1**, `main`, containing a single return statement **2**, which returns an integer—in this case, **2** **3**. Let's translate the code in Listing 2-1 into assembly using GCC:

```
$ gcc -S -O -fno-asynchronous-unwind-tables -fcf-
protection=none return_2.c
```

These GCC options produce fairly readable assembly:

**-S** Don't run the assembler or linker. This makes GCC emit assembly instead of a binary file.

**-O** Optimize the code. This eliminates some instructions we don't care about right now. When you inspect GCC output in later chapters, you'll usually want to turn optimization off so you can more clearly see how code generation works.

**-fno-asynchronous-unwind-tables** Don't generate the unwind table,

which is used for debugging. We don't care about it.

**-fcf-protection=none** Disable control-flow protection. This is a security feature that adds extra instructions that we don't care about. Control-flow protection might already be disabled by default on your system, in which case this option won't do anything.

The result, stored in *return\_2.s*, should basically look like this:

```

1  .globl main
2  main:
3      movl    $2, %eax
4      ret

```

Listing 2-2      The program from Listing 2-1 translated into assembly.

**NOTE** All the assembly listings in this book use AT&T syntax. Elsewhere, you'll sometimes see x64 assembly written in Intel syntax. They're just two different notations for the same language; the biggest difference is that they put instruction operands in different order.

Your *.s* file might contain a few other assembler directives, but you can safely ignore them for now. The four lines in Listing 2-2 are a complete assembly program. Assembly programs have several kinds of statements. The first line, `.globl main` **1**, is an *assembler directive*, a statement that provides directions for the assembler. Assembler directives always starts with a period. Here, `main` is a *symbol*, a placeholder for a memory address. An assembly instruction can include a symbol when it needs to refer to the address of a particular function or variable, but the compiler doesn't know where that function or variable will end up in memory. Later, after the linker has combined the different object files that make up the executable, it can associate each symbol with a memory address; this process is called *symbol resolution*. Then the linker will update every place that uses a symbol to use the corresponding address instead; this is called *relocation*.

The `.globl main` directive tells the assembler that `main` is a *global* symbol. By default, a symbol can only be used in the same assembly file (and therefore the same object file) where it's defined. But because `main` is global, other object files can refer to it too. The assembler will record this fact in a section of the object file called the *symbol table*. The symbol table contains information about all the symbols in an object file or executable. The linker relies on the symbol table during symbol resolution. If the symbol table doesn't list `main` as a global symbol, but another object file tries to refer to it, linking will fail.

Next, we use `main` **2** as a *label* for the code that follows it. Labels consist of a string or number followed by a colon. This label marks the location that the symbol `main` refers to. For example, the instruction `jmp`



`main` should cause the program to jump to the instruction at line 3. But the label can't indicate the final location of `main`; like I mentioned earlier, we won't know that until link time. Instead, it defines `main` as an offset from the start of the current *section* in this object file. (An object file includes different sections for machine instructions, global variables, debug information, and so on, which are loaded into different parts of the program's address space at runtime. The object file produced from Listing 2-2 will only have one section: the text section, which contains machine instructions.) Because 3 is the very first machine instruction in this file, the offset of `main` will be 0. The assembler will record this offset in the symbol table so the linker can use it to determine the final address of `main` during symbol resolution.

## FURTHER READING ON LINKERS

The last couple of paragraphs really oversimplified how linking works! If I included a totally accurate explanation of linkers, this chapter would be 90% about linkers and 10% about your actual compiler. But you should go read more about linkers, because you need to understand them in order to really get what's going on in a running program. Here are some blog posts on linkers that I like:

- “Beginner's Guide to Linkers,” by David Drysdale, is a good starting point. (<http://www.lurklurk.org/linkers/linkers.html>)
- Ian Lance Taylor's 20-part essay on linkers goes into a lot more depth. The first post is at <https://www.airs.com/blog/archives/38>, and there's a table of contents at <https://lwn.net/Articles/276782/>.
- “Position Independent Code (PIC) in shared libraries,” a blog post by Eli Bendersky, provides an overview of how compilers, linkers, and assemblers work together to produce position-independent code, focusing on 32-bit machines (<https://eli.thegreenplace.net/2011/11/03/position-independent-code-pic-in-shared-libraries/>).
- “Position Independent Code (PIC) in shared libraries on x64,” also by Eli Bendersky, builds on the previous article, focusing on 64-bit systems (<https://eli.thegreenplace.net/2011/11/11/position-independent-code-pic-in-shared-libraries-on-x64/>).

Next, we have `movl` 3, an example of a *machine instruction*, which is an instruction that appears in the final executable. The `movl` instruction in Listing 2-2 moves the value 2 into a *register*—a very small, very fast storage slot that has its own name and sits right on the CPU. Here, we move 2 into the register named EAX, which can hold 32 bits. According to our platform's calling convention, return values are passed to the caller in EAX (or RAX, the 64-bit equivalent, depending on the return value's type). Since the caller also knows about this convention, it can retrieve the return value from EAX after the function returns. The `l` suffix in `movl` indicates that the operands to this instruction are long integers. In assembly, unlike most

modern implementations of C, “long” means 32 bits. A `movq` instruction operates on *quadwords*, which is how x64 assembly refers to 64-bit integers. I’ll just write `mov` when I want to refer to this instruction without specifying its size.

Finally, we have another machine instruction, `ret 4`, which returns control to the caller. You might see `retq` here instead of `ret`, since this instruction implicitly operates on a 64-bit return address. I’m skipping a lot of details, like what calling conventions are, who decides on them, or how `ret` knows where the caller is. I’ll come back to those details when we add function calls in chapter 10.

At this point, it’s fair to ask who the caller is, since `main` is the only function in this program. It’s also fair to wonder why we need the `.globl main` directive, since there don’t seem to be any other object files that could contain references to `main`. The answer is that the linker adds a tiny bit of wrapper code called `crt0` to handle setup before `main` runs, and teardown after it exits. (The `crt` stands for “C Runtime.”) This wrapper code basically does the following:

1. Makes a function call to `main`. This is why `main` needs to be globally visible; if it’s not, `crt0` can’t call it.
2. Retrieves the return value from `main`.
3. Invokes the `exit` system call, passing it the return value from `main`. Then `exit` handles whatever work needs to happen inside the operating system to terminate the process and turn the return value into an exit code.

The bottom line is that you don’t need to worry about process startup or teardown; you can treat `main` like a normal function.

To verify that the assembly in Listing 2-2 works correctly, you can assemble and link it, run it, and check the exit code with the `$?` shell operator:

```
$ gcc return_2.s -o return_2
$ ./return_2
$ echo $?
2
```

Note that you can pass an assembly file to GCC just like a regular source file. GCC assumes any input files with a `.s` extension contain assembly, so it will just assemble and link those files without trying to compile them first.

## Writing the Compiler Driver

As we saw in the last chapter, a compiler isn't very useful on its own. To turn a source file into an executable, you'll need to write a compiler driver that invokes the preprocessor, compiler, assembler, and linker. It's a good idea to write a compiler driver that works with *test\_compiler* before starting on the compiler itself, so you can validate each compiler stage against the test suite as you go. The compiler driver should do the following:

1. Preprocess a source file:  

```
| gcc -E -P INPUT_FILE -o PREPROCESSED_FILE
```
2. By convention, the preprocessed file should have a *.i* file extension.
3. Compile the preprocessed source file, and output an assembly file with a *.s* extension. You'll have to stub out this step, since you haven't written your compiler yet.
4. Assemble and link the assembly file to produce an executable:  

```
| gcc ASSEMBLY_FILE -o OUTPUT_FILE
```

To work with *test\_compiler*, your compiler driver must be a command-line program that accepts a path to a C source file as its only argument. If this command succeeds, it must produce an executable in the same directory as the input file, with the same name (minus the file extension). In other words, if you run `./YOUR_COMPILER /path/to/program.c`, it should produce an executable at `/path/to/program` and terminate with an exit code of zero. If your compiler fails, the compiler driver should return a non-zero exit code, and should not write any assembly or executable files; that's how *test\_compiler* verifies that your compiler catches errors in invalid programs. Finally, your compiler driver should support a `--lex` option that directs it to just perform the lexing pass, as well as a `--parse` option that directs it to just run the lexer and parser but stop before code generation. Neither of these options should produce any output files.

Once you've written the compiler driver, you're ready to start working on the actual compiler.

You need to implement the four compiler passes I listed at the beginning of the chapter: the lexer, which produces a list of tokens; the parser, which turns those tokens into an abstract syntax tree; the code generator, which converts the abstract syntax tree into assembly, and the assembly emitter, which writes that assembly to a file. Let's look at each of those passes in more detail.

## Writing the Lexer

The lexer should read in a source file and return a list of tokens. Before

you can start writing the lexer, you need to know what tokens you might encounter. Here are all the tokens in Listing 2-1:

`int`: a keyword  
`main`: an identifier, whose value is “main”  
`(` : an open parenthesis  
`)` : a close parenthesis  
`{` : an open brace  
`return`: a keyword  
`2`: a constant, whose value is “2”  
`;` : a semicolon  
`}` : a close brace

I’ve used two lexer-specific terms here. An *identifier* is an ASCII letter followed by a mix of letters and digits; identifiers are case sensitive. An (integer) *constant* consists of one or more digits. (C supports hexadecimal and octal integer constants too, but you can ignore them to keep things simple. We’ll add character and floating-point constants in part II.)

Note that identifiers and constants have values in the list of tokens above, but the other types of tokens don’t. There are many possible identifiers (`foo`, `variable1`, or `my_cool_function`), so each identifier token produced by the lexer needs to retain its specific name. Likewise, each constant token needs to hold an integer value. By contrast, there’s only one possible `return` keyword, so a `return` keyword token doesn’t need to store any extra information. Even though `main` is the only identifier right now, it’s a good idea to build the lexer in a way that can support arbitrary identifiers later on. Also note that there are no whitespace tokens. If we were compiling a language like Python, where whitespace is significant, we’d need to include whitespace tokens.

You can define each token type with a regular expression. Table 2-1 gives the corresponding regular expression for each token in PCRE syntax:

Table 2-1      Tokens

Token	Regular Expression
Identifier	<code>[a-zA-Z_]\w*\b</code>
Constant	<code>[0-9]+\b</code>
Int keyword	<code>int\b</code>
Return keyword	<code>return\b</code>

Open parenthesis	<code>\ (</code>
Close parenthesis	<code>\ )</code>
Open brace	<code>{</code>
Close brace	<code>}</code>
Semicolon	<code>;</code>

The process of tokenizing a program then looks roughly like this:

```

while input isn't empty:
    find longest match at start of input for any regex in
    Table 2-1
    convert matching substring into a token
    remove matching substring from start of input
    trim whitespace from start of input
    if no valid token can be created, raise an error

```

Listing 2-3

Converting a string to a sequence of tokens

Note that identifiers and constants must end at word boundaries. For example, the first three digits of `123;bar` match the regular expression for a constant, and can be converted into the constant `123`. That's because `;` isn't in the `\w` character class, so the boundary between `3` and `;` is a word boundary.

However, the first three digits of `123bar` do not match the regular expression for a constant, because those digits are followed by more characters in the `\w` character class instead of a word boundary. If your lexer sees a string like `123bar` it should raise an error, because the start of the string doesn't match the regular expression for any token.

You can assume that your C source file only contains ASCII characters. The C standard provides a mechanism called *universal character names* to include non-ASCII characters in identifiers, but we won't implement them. Many C implementations let you use Unicode characters directly, but you don't need to support that either.

## Testing the Lexer

You can test your lexer against all the programs in `tests/chapter_2`. The sample programs in `tests/chapter_2/invalid_lex` all contain invalid tokens, so they should all cause the lexer to fail with an appropriate error message. The sample programs in `tests/chapter_2/invalid_parse` and `tests/chapter_2/valid` only contain valid tokens, so the lexer should be able to process them successfully. You can use the following command to test that your program

fails on the programs in [tests/chapter\\_2/invalid\\_lex](#) and succeeds on everything else:

```
$ ./test_compiler /path/to/your_compiler --chapter 2 --stage lex
```

This command just tests whether the lexer succeeds or fails. You may want to write your own tests to validate that it produces the correct list of tokens for valid programs and emits an appropriate error for invalid ones.

## Implementation Tips

**Treat keywords like other identifiers.** The regex for identifiers also matches keywords. Don't try to simultaneously find the end of the next token and figure out whether it's a keyword or not. First, find the end of the token. Then, if it looks like an identifier, check whether it matches any of the keywords.

**Don't split on whitespace.** It might seem like a good idea to start by splitting the string on whitespace, but it's not. It will just complicate things, because whitespace isn't the only boundary between tokens. For example, `main()` has three tokens and no whitespace.

## Writing the Parser

Now that you have a list of tokens, the next step is to figure out how those tokens are grouped together into language constructs. In most programming languages, including C, this grouping is hierarchical: each language construct in the program is composed of several simpler constructs. Individual tokens represent the most basic constructs, like variables, constants, and arithmetic operators. Tree data structures are a natural way to express this hierarchical relationship. A tree representation of a program is called an *abstract syntax tree*, or *AST*. Most compilers use ASTs internally, and yours will too. Your parser will accept the list of tokens produced by the lexer and generate an AST. Then your code generation stage will traverse that AST to figure out what assembly code to emit.

There are two basic approaches to writing a parser. One option is to handwrite the code for your parser. The other option is to use a *parser generator* like Bison or ANTLR to produce your parsing code automatically. Using a parser generator is less work than hand-writing a parser, but this book uses a handwritten parser for a few reasons. Most importantly, hand-writing a parser will give you a solid understanding of how your parser works. It's easy to use a parser generator without really understanding the code it produces. Many parser generators also have a steep learning curve, and you're better off learning general techniques like recursive descent parsing *before* you spend a lot of time mastering specific

tools.

Handwritten parsers also have some practical advantages over those produced by parser generators; they can be faster and easier to debug, and provide better support for error handling. In fact, both GCC and Clang use handwritten parsers. So writing a parser by hand isn't just an academic exercise.

That said, if you'd rather use a parser generator, that's fine too! It all depends on what you're hoping to get out of the book. But I won't talk about how to use them, so you'll have to figure that out on your own. If you decide to go that route, make sure to research what parsing libraries are available in your implementation language of choice.

Whichever option you choose, the first step is designing the abstract syntax tree you want your compiler to produce. It might help to see an example of an AST first.

## ***An Example Abstract Syntax Tree***

Let's take a look at the AST for this code snippet:

```
if (a < b) {
    return 2 + 2;
}
```

Listing 2-4

A simple if statement

This is an `if` statement, so we'll label the root of the AST `if`. The `if` node will have two children:

1. The condition, `a < b`
2. The "then" clause, `return 2 + 2;`

Each of these constructs can be broken down further. For example, the condition is a binary operation with three children:

3. The left operand, variable `a`
4. The operator, `<`
5. The right operand, variable `b`

Figure 2-2 shows the whole AST for this code snippet:

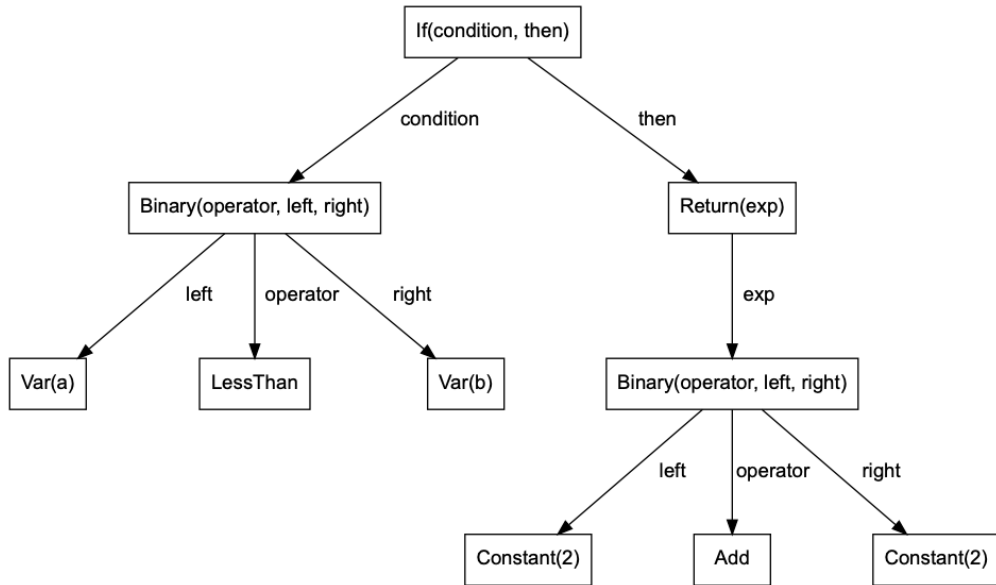


Figure 2-2 An AST for a simple if statement.

The AST in Figure 2-2 contains all the same information as Listing 2-4. By looking at it, we can tell what actions the program will take, and in what order. But, unlike Listing 2-4, this AST presents that information in a way that your compiler can easily work with. In later stages the compiler will traverse the tree, performing a different action at each type of node it encounters. We'll use this general strategy to accomplish a bunch of different tasks, from resolving variable names to generating assembly.

Now that we understand what ASTs look like in general, we can also construct a much simpler AST for the C program from Listing 2-1:



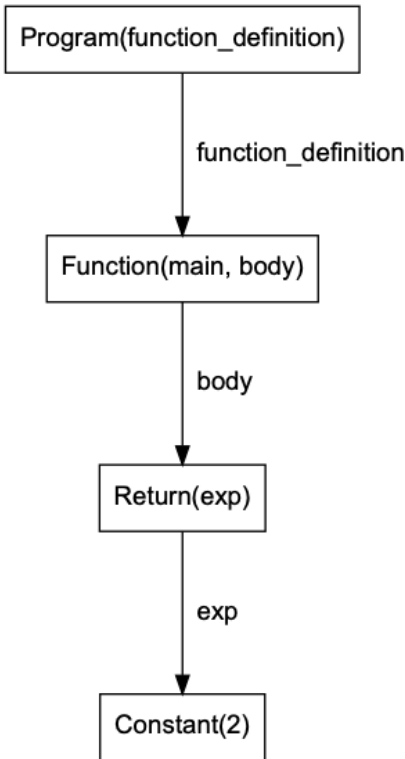


Figure 2-3      The AST for Listing 2-1

Next, you'll define the necessary data structures to construct ASTs like Figure 2-3 in code.

## Defining the AST

This book gives AST descriptions in a language designed for specifying ASTs, the *Zephyr Abstract Syntax Description Language (ASDL)*. I'm using ASDL here as convenient, programming-language-neutral notation. You won't use ASDL directly in your compiler; instead, you'll define equivalent data structures in your chosen programming language. The next few paragraphs include a very brief overview of ASDL. The original paper, which describes the whole language, is listed in the Further Reading section below.

Here's the ASDL description that covers the tiny subset of C you'll implement in this chapter (programs like Listing 2-1):

```

1 program = Program(function_definition)
2 function_definition = Function(3 identifier name, statement
   body)

```

```

4 statement = Return(exp)
5 exp = Constant(int)

```

Listing 2-5

Abstract syntax tree definition for this chapter

Each of the four lines in Listing 2-5 describes how to build one type of AST node. Note that every AST node in Figure 2-2 has a corresponding definition in ASDL. The root of this AST is the `program` node [1](#). At the moment, this node can have exactly one child, of type `function_definition`. A function definition has two children: a function name of type `identifier`, and a function body of type `statement` [2](#). Right now, a function consists of a single statement and has no arguments. Later, you'll add support for function arguments and more complex function bodies. Note that `name` and `body` in this definition are *field names*, human-friendly labels that don't change the structure of the AST. Field names are optional in ASDL. When a field name is present, it comes immediately after the field type, like in `identifier name` [3](#). I'll add field names when I think they'll make things more readable.

In ASDL, `identifier` is a built-in type that represents function and variable names; they're basically strings, but we want to distinguish them from string literals like `"Hello, World!"` because they appear in different parts of an AST. Since `identifier` is a built-in type, it has no children. The other child of the `function_definition` node is `statement`. [4](#) Right now, the only kind of statement is a return statement. A return statement has one child: its return value, of type `exp`, short for expression. The only `exp` at the moment is a constant integer [5](#). `int` is another built-in ASDL type, so our tree is finished.

Of course, return statements aren't the only statements in C, and constants aren't the only expressions. In later chapters, we'll add new constructors to represent the other kinds of statements and expressions. For example, we'll add an `If` constructor to `statement` to represent if statements:

```

statement = Return(exp) | If(exp condition, statement then,
statement? else)

```

The `statement?` type indicates an optional statement, since if statements don't always have an else clause. The `|` symbol separates constructors. Here, it tells us that a `statement` can be either a return statement, defined by the `Return` constructor, or an if statement, defined by the `If` constructor.

Now that you understand the AST definition in Listing 2-5, you need to implement it in whatever language you're using to write your compiler. The standard way to represent ASTs varies a lot between programming languages. If you're implementing your compiler in a functional language like F#, ML, or Haskell, it's easy to define your AST using algebraic

datatypes. Enums in Rust are basically algebraic datatypes, so they're also a good way to represent ASTs. If you're using an object-oriented language like Java, you can define an abstract class for each type of node, and define classes that extend or inherit from those abstract classes for each constructor. For example, you could define an `Exp` abstract class, and `Constant` and `BinaryExp` classes that extend it.

## FURTHER READING ON AST DEFINITIONS

If you're still not sure how to write an AST definition in your implementation language of choice, here are a couple papers that might help:

- "Abstract Syntax Tree Implementation Idioms," by Joel Jones, provides a good overview of how to implement ASTs in various programming languages (<https://hillside.net/plop/plop2003/Papers/Jones-ImplementingASTs.pdf>).
- "The Zephyr Abstract Syntax Description Language," the original paper on ASDL, includes examples of AST definitions in a few different languages (<https://www.cs.princeton.edu/~appel/papers/asdl97.pdf>).

## Defining the Formal Grammar

An AST has all the information you'll need in later stages of the compiler. It does not, however, tell you exactly what tokens make up each language construct. For example, nothing in the AST description in Listing 2-5 tells us that a return statement needs to end with a semicolon, or a function body needs to be enclosed in braces. (This is why it's called an *abstract* syntax tree—by contrast, a *concrete* syntax tree would include every token from the original input). Once you have an AST, those specific details are irrelevant, so it's convenient to leave them out. When you're parsing a sequence of tokens to construct your AST, however, those details are extremely important, because they indicate where each language construct begins and ends.

So, in addition to an AST description, you'll need a set of rules about how to build a language construct from a list of tokens. This set of rules is called a *formal grammar*, and it will correspond closely to the AST description. Here's the formal grammar for C programs like Listing 2-1:

```
<program> ::= <function>
<function> ::= "int" <identifier> "(" ")" "{" <statement>
  "}"
<statement> ::= "return" <exp> ";"
<exp> ::= <int>
<identifier> ::= ? An identifier token ?
<int> ::= ? A constant token ?
```

Listing 2-6

Formal grammar for this chapter

The grammar in Listing 2-6 is in a notation called *Backus-Naur Form*

(BNF). Each line of this grammar is a *production rule* that defines how a language construct can be formed from a sequence of other language constructs and tokens. Every symbol that appears on the left-hand side of a production rule (like `<function>`) is called a *non-terminal symbol*. Individual tokens, like keywords, identifiers, and punctuation, are called *terminal symbols*. All non-terminal symbols are wrapped in angle brackets, and specific tokens (like `;`) are wrapped in quotes. The `<identifier>` and `<int>` symbols are special. They represent individual identifier and constant tokens, respectively, but these tokens aren't set strings like the other terminal symbols. Since there's not an easy way to define those symbols in Backus-Naur form, we describe each of them using a *special sequence*—that's just a plain English description of the symbol, wrapped in question marks.

Listing 2-6 looks a lot like the AST definition in Listing 2-5. In fact, it has exactly the same structure—every AST node in Listing 2-5 corresponds to a non-terminal symbol in Listing 2-6. The only difference is that Listing 2-6 specifies exactly which tokens we'll find at each node of the tree. This helps us figure out when we need to start processing a new node at the next level down in the AST, and when we've finished processing a node and can go back up to its parent on the level above.

Just like later chapters will introduce multiple constructors for some AST nodes, they'll introduce multiple production rules for the corresponding symbols. For example, here's how we'll add a production rule for `<statement>` to support if statements:

```
<statement> ::= "return" <exp> ";" | "if" "(" <exp> ")" <statement> [
    "else" <statement> ]
```

Note that brackets in BNF indicate that something is optional, just like questions marks in ASDL.

You'll need to refer to this formal grammar while writing the parser, but you don't need to explicitly define the grammar rules anywhere in your compiler.

## Recursive Descent Parsing

Now that you have an AST definition and a formal grammar, let's talk about how to actually write the parser. We'll use a straightforward technique called *recursive descent parsing*. A recursive descent parser uses a different function to parse each non-terminal symbol and return the corresponding AST node. For example, when the parser expects to encounter the `<statement>` symbol we defined in Listing 2-6, it will call a function to parse that symbol and return the `statement` AST node we defined in Listing 2-5. To parse an entire program, you'll call the function that parses the `<program>` symbol. With each function call to handle a new symbol,

the parser descends to a lower level in the tree. That’s where the *descent* in recursive descent comes from. (It’s called *recursive* descent because the grammar rules are often recursive, in which case the functions to process them will be too. For example, the operand of an expression could be another expression—we’ll see an example of that in the next chapter.)

Because the process of parsing a symbol is easier to explain in code, let’s look at some pseudocode for parsing a `<statement>` symbol. Once you understand it, you can write your own code to process `<statement>` and all the other non-terminal symbols in Listing 2-6.

```

1 parse_statement(tokens):
2     expect("return", tokens)
3     return_val = parse_exp(tokens)
4     expect(";", tokens)
5     return Return(return_val)

expect(expected, tokens):
    actual = take_token(tokens)
    if actual != expected:
        fail()

```

Listing 2-7

#### Parsing a statement

We’ll call the `parse_statement` function **1** when we expect the list of remaining tokens to start with a `<statement>`. According to Listing 2-6, a `<statement>` consists of three symbols: the `return` keyword, an `<exp>` symbol, and a “`;`” token. First, we call a helper function, `expect` **2**, to verify that the first token really is a `return` keyword. If it is, `expect` just discards it so we can move on to the next token. If it isn’t, there’s a syntax error in the program. Next, the grammar tells us that the `return` keyword should be followed by an `<exp>` symbol. We need to turn this symbol into an `exp` AST node so we can construct the return statement. Since this is a different non-terminal symbol, it should be handled by a separate function, `parse_exp`, which I haven’t defined here. Once we’ve gotten the AST node representing the return value back from `parse_exp` **3**, we just need to verify that it’s followed by the last token, a semicolon. We handle this with `expect` **4**, just like we handled the `return` keyword at **2**. At this point, we know the statement is syntactically valid, so we can return an AST node **5**.

Note that the `parse_statement` function removes all the tokens that made up the statement from the `tokens` list. After `parse_statement` returns, its caller will keep processing the remaining tokens in `tokens`. If there are any tokens left after parsing the entire program, that’s a syntax error.

The other thing to note is that this pseudocode is written a very imperative way. Functional languages (the sorts of languages I

recommended in the last chapter!) generally won't let you modify the input list like I'm doing here. So the details of how your parser passes tokens around will probably differ from Listing 2-7, but the overall structure will be the same.

Right now, each symbol in our formal grammar has only one production rule. In later chapters, when some symbols have multiple production rules, your parser will need to figure out which production rule to use. It can do that by looking at the first few tokens in the list without removing them. Recursive descent parsers that look ahead a few tokens to figure out which production rule to use are called *predictive parsers*. The alternative to predictive parsing is *recursive descent with backtracking*—trying each production rule in turn until you find one that works.

## Testing the Parser

Your parser should fail on the programs in [tests/chapter\\_2/invalid\\_parse](#) and succeed on the programs in [tests/chapter\\_2/valid](#). To test the parser, run:

```
$ ./test_compiler /path/to/your_compiler --chapter 2 --stage
parse
```

This command only tests whether the parser succeeds or fails, so you may want to write your own tests to confirm that it produces the correct AST for valid programs and emits an appropriate error for invalid ones.

## Implementation Tips

**Write a pretty-printer.** A pretty-printer is a function that prints out your AST in a human-readable way. This will make debugging your parser a lot easier. For example, a pretty-printed AST for the program in Listing 2-1 might look something like this:

```
Program(
  Function(
    name="main"
    body=Return(
      Const(2)
    )
  )
)
```

**Give informative error messages.** This will also help you debug your parser (and if anyone ever wants to use your compiler, it will help them too). An error message like `Expected ";" but found "return"` is a lot more helpful than `Fail`.

## Writing the Code Generator

The code generation pass should convert the AST into x64 assembly. It should traverse the AST in roughly the order the program executes, producing the appropriate assembly instructions for each node. First, you need to define an appropriate data structure to represent the assembly program, just like you needed to define a data structure to represent the AST when you wrote the parser. You’re adding yet another data structure, instead of writing assembly to a file right away, so that you can modify the assembly code after you’ve generated it. You won’t need to rewrite any assembly in this chapter, but in later chapters you will.

I’ll use ASDL again to describe the structure we’ll use to represent assembly. Here’s the definition:

```

1program = Program(function_definition)
2function_definition = Function(identifier name, instruction*
   instructions)
3instruction = Mov(operand src, operand dst) | Ret
4operand = Imm(int) | Register

```

Listing 2-8

ASDL definition of an assembly program

This looks a lot like the AST definition from the last section! In fact, this *is* an AST—but it’s an AST that represents an assembly program, not a C program. Every node corresponds to a construct in assembly, like a single instruction, rather than a construct in C, like a statement. I’ll refer to the data structure defined in Listing 2-8 as the “assembly AST” to distinguish it from the AST defined in Listing 2-5.

Let’s walk through Listing 2-8. An assembly program still consists of a single function **1**, which has a name and a list of instructions **2**. The *\** in *instruction\** indicates that it’s a list of instructions, not just one. The two instructions **3** that can appear in our very simple assembly programs are *mov* and *ret*, which we saw in Listing 2-2. The *mov* instruction has two operands: its copies the first operand, the source, to the second operand, the destination. The *ret* instruction doesn’t have any operands. The two possible operands **4** to an instruction are a register and an *immediate value*, which is a value included in the instruction – in other words, a constant. For now, you don’t need to specify which register to operate on, because your generated code will only use EAX. You’ll need to refer to other registers in later chapters.

This pass will have a similar structure to the parser, so I won’t write out the pseudocode here. You’ll need a function to handle each type of AST node, which will call other functions to handle that node’s children. Here’s the equivalent assembly you need for each AST node:

Table 2-2      AST nodes and equivalent assembly

AST Node	Assembly Construct
<code>Program(function_definition)</code>	<code>Program(function_definition)</code>
<code>Function(name, body)</code>	<code>Function(name, instructions)</code>
<code>Return(exp)</code>	<code>Mov(exp, Register)</code> <code>Ret</code>
<code>Constant(int)</code>	<code>Imm(int)</code>

This translation is pretty straightforward, but there are a couple things to note. The first is that a single statement results in multiple assembly instructions. The second is that this translation only works if an expression can be represented as a single assembly operand. That's true right now, because the only expression is a constant integer. But it won't be true once we encounter unary operators in the next chapter. At that point, your compiler will need to generate multiple instructions to calculate an expression, and then figure out where that expression is stored in order to copy it into EAX.

## Testing the Code Generator

To test the code generation stage, run:

```
$ ./test_compiler /path/to/your_compiler --chapter 2 --stage
codegen
```

The `--stage codegen` option will run your whole compiler. Like the equivalent `lex` and `parse` options, it will just check whether the compiler succeeds or fails. Your code generation stage should be able to handle any program that parses successfully, so this just tests that your code generation stage can handle every valid program. You may want to write your own tests to confirm that your compiler generates the assembly you expect.

## Implementation Tips

**Plan ahead for Part II!** If you go on to do Part II of the book, you'll need to update the `Mov` instruction, and many of the other instructions we'll add in the next few sections, to store type information. You might want to define your assembly AST in a way that makes it easy to add more fields to each constructor later on. If there's no good way to do that in your implementation language, that's okay; it just means you'll have a little extra refactoring to do in part II.

## Writing the Code Emitter

Now that your compiler can generate assembly instructions, the last step is writing those instructions to a file. Here's how to print each assembly



construct:

Table 2-3      Formatting assembly

Assembly Construct	Output	
Top-level Constructs		
Program(function_definition)	(just print out the function definition)	
Function(name, instructions)		<pre>.globl &lt;name&gt; &lt;name&gt;:     &lt;instructions&gt;</pre>
Instructions		
Mov(src, dst)		<pre>movl &lt;src&gt;, &lt;dst&gt;</pre>
Ret		<pre>ret</pre>
Operands		
Register		<pre>%eax</pre>
Imm(int)		<pre>\$&lt;int&gt;</pre>

Note that there must be a line break between instructions, just like in Listing 2-2. The code emission step will need to traverse the assembly AST, just like the code generation stage traverses the AST from Listing 2-5. Because the assembly AST corresponds so closely to the final assembly program, the code emission stage will be very simple, even as you add more functionality to the compiler in later chapters.

**NOTE** If you're compiling on macOS, you need an underscore in front of the function name. For example, if you compile a function called `main`, the label in the resulting assembly should be `_main`. If you're on any other system, don't include an underscore.

## Testing the Whole Compiler

To test the whole compiler from lexing to code generation, run:

```
$ ./test_compiler /path/to/your_compiler --chapter 2
```

This will compile each program in `tests/chapter_2/valid` with your compiler and GCC, run both executables, and verify that they produce the same exit code. It will also validate that all the invalid programs fail, but you should have already confirmed that during the earlier stages.

## Implementation Tips

**Generate readable assembly.** When you debug your compiler, you'll spend a lot of time reading the assembly it produces. Your life will be easier if that assembly is nicely formatted. You can indent every line except for labels, like GCC does, to make your assembly more readable. (That's also how I've formatted Listing 2-2.) If you like, you can also include comments in your assembly programs. A `#` symbol in an

assembly program comments out the rest of the line—it works just like `//` in C.

## Summary

In this chapter, you wrote a compiler that can transform a complete C program into an executable that runs on your computer. You learned how to interpret a program written in x64 assembly, a formal grammar in Backus-Naur form, and an AST definition in ASDL. The skills and concepts you learned in this chapter—and the four compiler stages you implemented—are the foundation for everything you’ll do in the rest of the book.

In the next chapter, you’ll add support for unary operators to your compiler. Along the way, you’ll learn about how assembly programs manage the stack, and we’ll introduce a new way to represent the programs you compile to make them easier to analyze, transform, and optimize.

# 3

## UNARY OPERATORS

C has several *unary operators*, which operate on a single value. In this chapter, you'll extend your compiler to handle two unary operators: negation and bitwise complement. You'll update the lexer, parser, and code emission pass to handle these new operators, but the code generation stage will require the biggest changes. Between parsing and code emission, you'll need to transform complex, nested expressions into simple operations that can be expressed in assembly. Instead of performing this transformation in a single compiler pass, we'll introduce a new intermediate

representation between the AST produced by the parser and the assembly program produced by the code generation pass. We'll also break up code generation into several smaller passes. To get started, let's look at a C program that uses our new unary operators, and the corresponding assembly we want to generate.

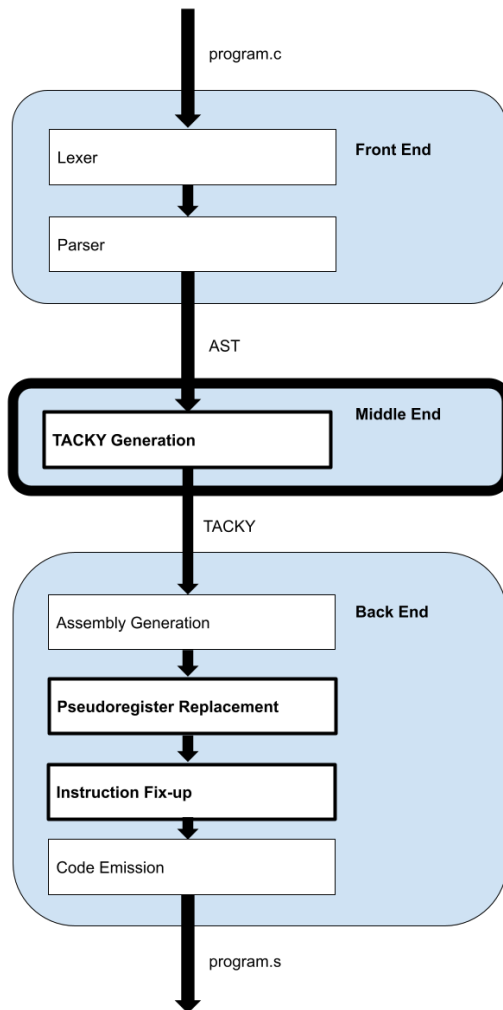


Figure 3-1

Stages of the compiler

## Negation and Bitwise Complement in Assembly

In this chapter, you'll learn how to compile programs like this one:

```
int main() {
    return ~(-2);
}
```

Listing 3-1 A C program with negation and bitwise complement

This program uses both of the unary operations we'll introduce in this chapter. It also includes a nested expression. If you implement your compiler the way I suggest, it will produce the assembly listing below from the program in Listing 3-1:

```

    .globl main
main:
    pushq    %rbp
    movq     %rsp, %rbp
    subq     $8, %rsp
1   movl     $2, 2-4(%rbp)
3   negl     -4(%rbp)
4   movl     -4(%rbp), %r10d
5   movl     %r10d, -8(%rbp)
6   notl     -8(%rbp)
    movl     -8(%rbp), %eax
7   movq     %rbp, %rsp
    popq     %rbp
    ret
```

Listing 3-2 Assembly code for Listing 3-1

The first three instructions after `main` are the *function prologue*, which set up the current stack frame; I'll cover them more when I talk about the stack in detail below. After the function prologue, we'll calculate the intermediate result, -2, and then final result, 1, storing each of them at unique memory address. The resulting assembly isn't very efficient; we waste a lot of instructions copying values from one address to another. But this approach sets us up to generate more efficient assembly later on. In Part III, you'll see how to store as many intermediate results as possible in registers, instead of memory, which will speed things up and eliminate a lot of unneeded copies.

**NOTE** If you compile Listing 3-1 to assembly using GCC, or any other production C compiler, it won't look anything like Listing 3-2, because those compilers evaluate constant expressions at compile time, even when optimizations are disabled! The basis for this seems to be section 6.6 of the C standard, which states that "[a] constant expression can be evaluated during translation rather than runtime,

and accordingly may be used in any place that a constant may be." Evaluating all constant expressions at compile time is an easy way to implement this part of the standard.

The first `movl` instruction `1` stores `2` at an address in memory. The operand `-4(%rbp)` `2` means “the value stored in the RBP register, minus 4.” The value in RBP is a memory address on the stack (more on that below), so `2` refers to another memory address four bytes lower. That address is where instruction `1` will store `2`. Then we negate the value at this address with the `neg` instruction `3`, so `-4(%rbp)` now contains the value `-2`. (Just like `mov`, `neg` has an `l` suffix to indicate that it’s operating on a 32-bit value.)

Next, we need to handle the outer bitwise complement expression. The first step is copying the source value, stored in `-4(%rbp)`, to the destination address at `-8(%rbp)`. We can’t do this in a single instruction, because the `mov` instruction can’t have memory addresses as both source and destination operands. At least one operand to `mov` needs to be a register or an immediate value. We’ll get around this by copying `-2` from memory into a scratch register, R10D `4`, and from there to the destination memory address `5`. We then take the bitwise complement of `-2` with the `not` instruction `6`, so memory address `-8(%rbp)` now contains value we want to return:  $\sim(-2)$ , which comes out to `1`. To return this value, we have to move it into EAX. The next three instructions make up the *function epilogue*, which tears down the stack frame and then returns from the function `7`.

## REPRESENTING SIGNED INTEGERS IN TWO'S COMPLEMENT

All modern computers use a *two's complement* representation of signed integers. A firm grasp on two's complement will help you understand and debug the assembly code your compiler generates. If you aren't already familiar with two's complement, or you need a refresher, here are a couple helpful resources:

"Two's Complement", by Thomas Finley, covers how and why two's complement representations work. (<https://www.cs.cornell.edu/~tomf/notes/cps104/twoscomp.html>)

The second chapter of *The Elements of Computing Systems*, by Noam Nisan and Shimon Schocken, covers similar material from a more hardware-focused perspective. This is the companion book for the Nand to Tetris project. This chapter is freely available at <https://www.nand2tetris.org/course>; click on the book icon under "Project 2: Boolean Arithmetic".

## The Stack

There are still two unanswered questions about the code in Listing 3-2:

what the function prologue and epilogue do, and why we refer to stack addresses relative to a value in the RBP register. To answer both those questions, we need to talk more about the segment of program memory called the *stack*. The address of the top of the stack is stored in the RSP register, which is also called the *stack pointer*. (By convention, RSP points to the last used stack slot, rather than the first free one.) Like you'd expect with any stack data structure, you can push things onto the stack and pop values off of it; the `push` and `pop` assembly instructions do exactly that.

The stack grows towards lower memory addresses. When you push something onto the stack, you decrement RSP. Whenever I say “top of the stack”, I mean the address stored in RSP, which is the lowest address on the stack. Note that the stack diagrams in this book, unlike other diagrams you may have seen, are oriented with lower memory addresses at the top. To help you remember how these diagrams are laid out, you can think of memory addresses like line numbers in a code listing; the top of the listing is at line 0, and line numbers increase as you go down. That means the top of the stack is on top of the diagram.

An instruction like `push $3` does two things:

1. Write the value being pushed (in this example, `3`) to the next empty spot on the stack. The `push` and `pop` instructions adjust the stack pointer in 8-byte increments, and the top value on the stack is currently at the address stored in RSP, so the next empty spot is `RSP - 8`.
2. Decrement RSP by eight. The new address in RSP is now the top of the stack, and the value at that address is `3`.

Figure 3-2 illustrates the effect of a `push` instruction on the stack and RSP register.

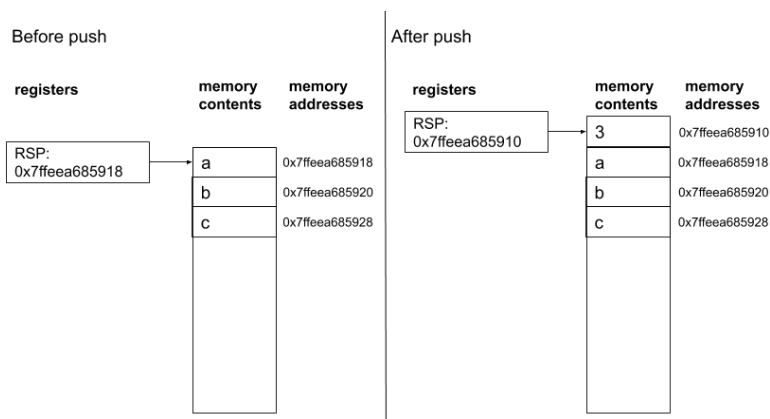


Figure 3-2 Effect of `push $3` instruction on memory and RSP.

The `pop` instruction does the opposite. For example, `pop %rax` would copy the value at the top of the stack into the RAX register, and then add eight to RSP.

Since the `push` instruction decrements the stack pointer by eight bytes, it has to push an 8-byte value. Likewise, the `pop` instruction can only pop 8-byte values off the stack. Because x64 memory addresses are eight bytes, we can use `push` and `pop` to put them on and take them off the stack. But the `int` type is only four bytes. If you want to put a 4-byte value on the stack, like the literal `2` from Listing 3-1, you can't use `push`, so you have to use `mov` instead. (On 32-bit architectures, the reverse is true; you can push and pop 4-byte values but not 8-byte values. In either case, it's also possible to push and pop 2-byte values, but as far as I know you'd never want to do that.)

The stack isn't just an undifferentiated chunk of memory; it's divided into sections called *stack frames*. Whenever a function is called, it allocates some memory at the top of the stack by decreasing the stack pointer. This memory is the function's stack frame. Just before the function returns, it deallocates its stack frame, restoring the stack pointer to its previous value. We'll store the base of the current stack frame in the RBP register, which is the usual approach. We can refer to data in the current stack frame relative to the address stored in RBP. That way we don't need absolute addresses, which we can't know in advance. Since the stack grows toward lower memory addresses, any address in the current stack frame will be lower than the address stored in RBP, which is why we refer to local variables with operands like `-4(%rbp)`. We can also refer to data in the caller's stack frame, like function arguments, relative to RBP. We'll need to do that later when we implement function calls. (Alternatively, we could refer to local variables and parameters relative to RSP, and not bother with RBP at all; some compilers do this as an optimization. We'll stick with RBP-relative addressing because the resulting assembly is easier to understand.)

So, the first thing a function needs to do is set up a new stack frame, and the last thing it needs to do before it returns is restore the caller's stack frame. The function prologue sets up the stack frame in three instructions, as shown in Figure 3-3:

1. `pushq %rbp` saves the current value of RBP, the address of the base of the caller's stack frame, onto the stack. We save it because we'll need it to restore the caller's stack frame later. This value will be at the bottom of the new stack frame established by the next instruction.
2. `movq %rsp, %rbp` makes the top of the stack the base of the new



stack frame. At this point, the top and bottom of the current stack frame are the same. The current stack frame holds exactly one value, which both RSP and RBP point to: the base of the caller's stack frame, which we saved in the previous instruction.

3. `subq $n, %rsp` decrements the stack pointer by `n` bytes. The stack frame now has `n` bytes available to store local and temporary variables. In Figure 3-3, this instruction allocates 24 bytes, enough space for six 4-byte integers. It would also work to just push values onto the stack as needed, instead of allocating space for all of them up front, but most compilers don't do that. One problem with that approach is that you can only push 8-byte values. That's inconvenient when you want to store 4-byte integers, like we do right now.

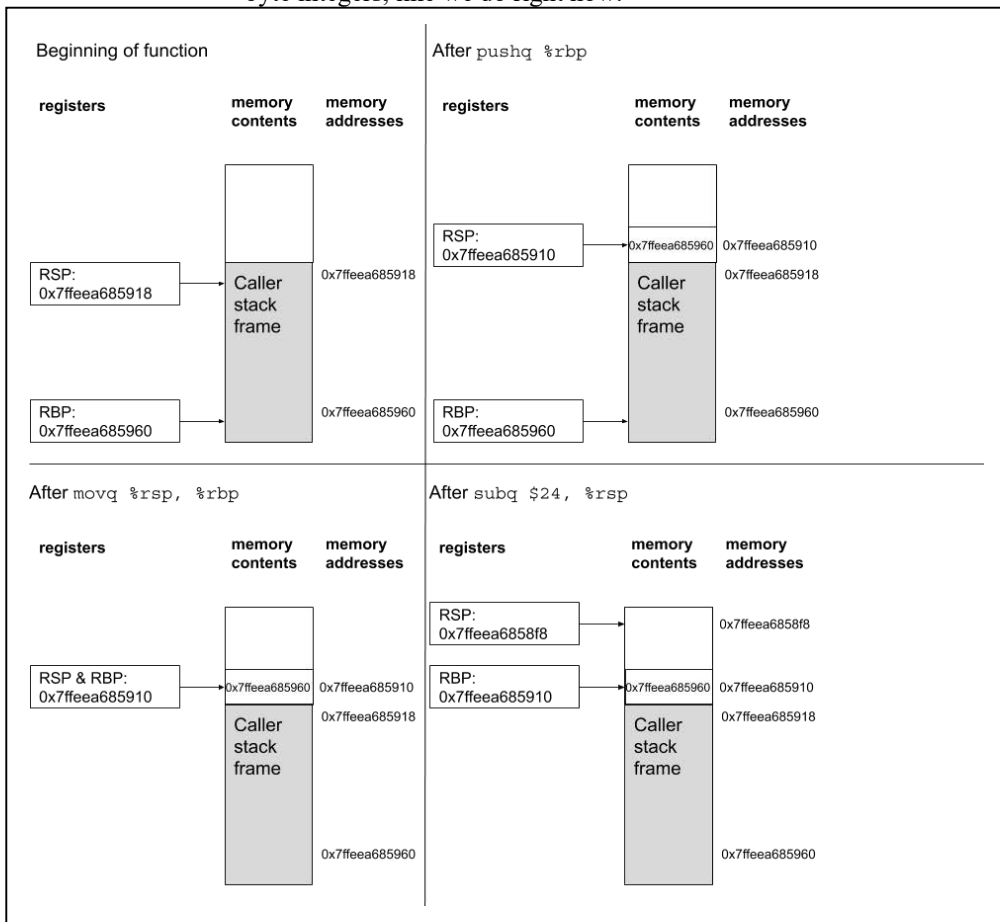


Figure 3-3 State of the stack at each point in the function prologue

The function epilogue needs to restore the caller's stack frame; that means RSP and RBP need to have the same values they did before the function prologue. This requires two instructions, as shown in Figure 3-4:

1. `movq %rbp, %rsp` puts us back where we were after the second instruction of the function prologue: both RSP and RBP point to bottom of the current stack frame, which holds the caller's value for RBP.
2. `popq %rbp` reverses the first instruction of the function prologue and restores the caller's values for both the RSP and RBP registers. It restores RBP because the value at the top of the stack was the base address of the caller's stack frame that we stored in the first instruction of the prologue. It restores RSP because it pops the last value in this stack frame off the stack, leaving RSP pointing to the top of the caller's stack frame.

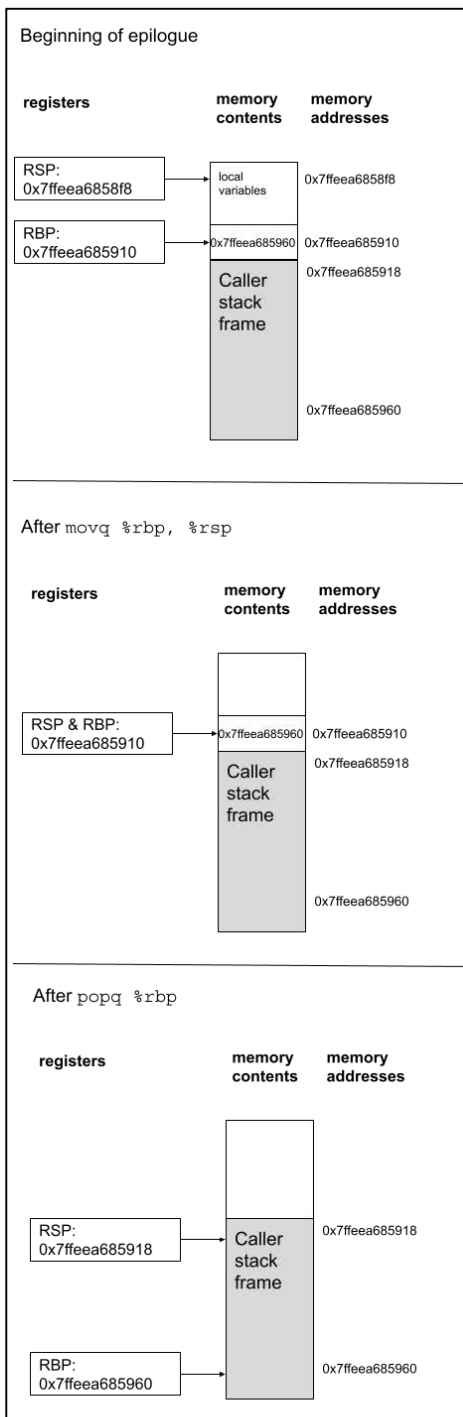


Figure 3-4 State of the stack at each point in the function epilogue

Now that we know what output our compiler should produce, we're ready to continue coding. Let's start by extending the lexer and parser.

## Extending the Lexer

You'll need to extend the lexer to recognize three new tokens:

`~` : a tilde, the bitwise complement operator

`-` : a hyphen, the negation operator

`--` : two hyphens, the decrement operator

You won't implement the decrement operator in this chapter, but you still need to add a token for it. Otherwise, your compiler will accept programs it should reject, like this one:

```
int main() {
    return --2;
}
```

Listing 3-3      An invalid C program using the decrement operator

This shouldn't compile, because you can't decrement a constant. But if your compiler doesn't know that `--` is a distinct token, it will think Listing 3-3 is equivalent to:

```
int main() {
    return -(-2);
}
```

Listing 3-4      A valid C program with two negation operators in a row

which is a perfectly valid program. Your compiler should reject language features you haven't implemented—it shouldn't compile them incorrectly. That's why your lexer needs to know that `--` is a single token, not just two negation operators in a row. (On the other hand, the lexer should lex `~~` as two bitwise complement operators in a row. Expressions like `~~2` are perfectly valid.)

You can process the new tokens exactly the same way you handled punctuation like `;` and `(` in the previous chapter. First, you'll need to define a regular expression for each new token—the regular expressions here will just be the strings `~`, `-`, and `--`. Next, have your lexer check the input against these new regexes, as well as the regexes from the previous chapter, every time it tries to produce a token. Remember that when the start of the input stream matches more than one possible token, you should always

choose the longest one. So, if your input stream ever starts with `--`, you'll parse it as a decrement operator rather than two negation operators.

## Testing the Lexer

The lexer should successfully lex all the test cases for this chapter, including the valid test programs in [tests/chapter\\_3/valid](#) and the invalid test programs in [tests/chapter\\_3/invalid\\_parse](#). To test your lexer against all the test cases so far, run:

```
$ ./test_compiler /path/to/your_compiler --chapter 3 --stage
lex
```

This will also run all the lexing test cases from Chapter 2, to make sure your lexer can still handle them.

## Extending the Parser

To parse the new operators in this chapter, we first need to extend the AST and formal grammar we defined in Chapter 2. Let's look at the AST first. Since unary operations are expressions, we'll represent them with a new constructor for the `exp` AST node. Here's the updated AST definition, with new parts bolded:

```
program = Program(function_definition)
function_definition = Function(identifier name, statement
body)
statement = Return(exp)
exp = Constant(int) | Unary(unary_operator, exp)
unary_operator = Complement | Negate
```

Listing 3-5

Abstract syntax tree with unary operations

The updated rule for `exp` indicates that an expression can be either a constant integer or a unary operation. A unary operation consists of one of the two unary operators, `Complement` or `Negate`, applied to an inner expression. Notice that the definition of `exp` is recursive: the `Unary` constructor for an `exp` node contains another `exp` node. That lets us construct arbitrarily deeply nested expressions, like `-(~(-~-(-4)))`.

We also need to make the corresponding changes to the grammar:

```
<program> ::= <function>
<function> ::= "int" <identifier> "(" ")" "{" <statement>
"}"
<statement> ::= "return" <exp> ";"
<exp> ::= <int> | <unop> <exp> | "(" <exp> ")"
```

```

<unop> ::= "~" | "~"
<identifier> ::= ? An identifier token ?
<int> ::= ? A constant token ?

```

Listing 3-6

Formal grammar with unary operations

Listing 3-6 includes a new production rule for unary expressions, and a new `<unop>` symbol to represent the two unary operators. Those changes correspond exactly with the addition to the AST in Listing 3-5. We've also added a third production rule for the `<exp>` symbol, which doesn't correspond to anything in Listing 3-5. This rule just indicates that if you wrap an expression in parentheses, the result is still an expression. It doesn't have a corresponding constructor in the AST because the rest of the compiler doesn't need to distinguish between an expression wrapped in parentheses and the same expression without parentheses. The expressions `1`, `(1)`, and `(( (1) ))` should all be represented by the same AST node: `Constant(1)`.

The decrement operator `--`, doesn't show up anywhere in this grammar. That means your parser should fail if it encounters a `--` token.

To update the parsing stage, you first need to modify your compiler's AST data structure to match Listing 3-5. Then you need to update your recursive descent parsing code to reflect the changes in Listing 3-6. Parsing an expression is a bit more complicated than it was in the previous chapter, because the `<exp>` symbol has three different production rules and you need to figure out which one to apply. This pseudocode sketches out how to parse an expression:

```

parse_exp(tokens):
1   next_token = peek(tokens)
2   if next_token is an int:
      --snip--
3   else if next_token is "~" or "-":
4       operator = parse_unop(tokens)
5       inner_exp = parse_exp(tokens)
6       return Unary(operator, inner_exp)
7   else if next_token == "(":
      take_token(tokens)
      inner_exp = parse_exp(tokens)
      expect(")", tokens)
8   return inner_exp
   else:
       fail()

```

Listing 3-7

Pseudocode for parsing an expression

The first step is looking at the next token in the input to figure out

which production rule to apply. We call `peek` [1](#) to look at this token without removing it from the input stream. Once we know which production rule to use, we'll want to process the whole input, including that first token, using that rule. That's why we don't want to consume this token from the input just yet. (Like I mentioned in the last chapter, you might not consume tokens from the input stream exactly as I've described here, so you might look at the first token without actually calling a `peek` function.)

If the expression we're about to parse is valid, `next_token` should be an integer, a unary operator, or an open parenthesis. If it's an integer [2](#), we can parse it exactly the same way we did in the previous chapter. If it's a unary operator [3](#), we need to apply the second production rule for `<exp>` from Listing 3-6 to construct a unary expression. This rule is `<unop>` `<exp>`, so we'll parse the unary operator and then the inner expression. The `<unop>` symbol is a single token, `next_token`, which we've already inspected. In Listing 3-7, we handle `<unop>` in a separate function [4](#) (`parse_unop`, whose definition I've omitted). In practice, it might be unnecessary to define a separate function to parse just one token. Either way, we'll end up with a very simple AST node representing the appropriate unary operator. The operator should be followed by an `<exp>` symbol, which we'll process with a recursive call to `parse_exp` [5](#). (This is the recursive part of "recursive descent.") That call should return an `exp` AST node representing the operand of the unary expression. Now we have AST nodes for both the operator and the operand, so we can return the AST node for the whole unary expression [6](#).

If `next_token` is an open parenthesis [7](#), it should be immediately followed by a valid expression, so we remove the parenthesis from the input stream and call `parse_exp` recursively to handle the expression that follows. The inner expression should be followed by a closing parenthesis to balance out the opening parenthesis we already processed. We call `expect` to remove that closing parenthesis or throw a syntax error if it's missing. Since the AST doesn't need to indicate that there were parentheses, we can just return the inner expression as-is [8](#).

If `next_token` isn't an integer, a unary operator, or an open parenthesis, the expression must be malformed, so we throw a syntax error.

## Testing the Parser

The parser should be able to handle every valid test case in [tests/chapter\\_3/valid](#), and raise an error on every invalid test case in [tests/chapter\\_3/invalid\\_parse](#). It should also continue to handle valid and invalid test cases from the last chapter correctly. To test your parser against

the test cases from this chapter and last chapter, run:

```
$ ./test_compiler /path/to/your_compiler --chapter 3 --stage
parse
```

## TACKY: A New Intermediate Representation

Converting the AST to assembly isn't as straightforward as it was in the last chapter. C expressions can have nested sub-expressions, and assembly instructions can't. A single expression like `-(~2)` needs to be broken up into two assembly instructions: one to apply the inner bitwise complement operation, and one to apply the outer negation operation.

We'll bridge the gap between C and assembly using a new intermediate representation (IR), *three-address code (TAC)*. In three-address code, the operands of each instruction must be constants or variables, not nested expressions. That means each instruction uses at most three values: two operands and a destination. (It would be more accurate to call this two-address code until we implement binary operators in the next chapter.) To rewrite nested expressions in three-address code, we often need to introduce new temporary variables. For example, `return 1 + 2 * 3` would become:

```
tmp0 = 2 * 3
tmp1 = 1 + tmp0
return tmp1
```

Listing 3-8

Three-address code for `return 1 + 2 * 3`

There are two main reasons to use three-address code instead of converting an AST directly to assembly. The first is that it lets us handle major structural transformations, like removing nested expressions, separately from the details of assembly language, like figuring out which operands are valid for which instructions. This lets us keep each compiler pass small, instead of having one really big compiler pass handling all those concerns. The second reason is that three-address code is well-suited to several of the optimizations we'll implement in Part III. It has a simple, uniform structure, which makes it easy to answer questions like “is the result of this expression ever used?” or “will this variable always have same value?” The answers to those questions will determine what optimizations are safe to perform.

### MULTIPLE LANGUAGES, MULTIPLE TARGETS

Intermediate representations like three-address code are useful for one other reason, although it isn't relevant to this project. An intermediate representation can provide a common target for



multiple source languages and a common starting point for assembly generation for multiple target architectures. The LLVM compiler framework is a great example of this: it supports several frontends and backends using a single intermediate representation. If you want to compile a new programming language, you can just compile it to the LLVM IR, and then LLVM can do all the work of optimizing that IR and producing machine code for a bunch of different CPU architectures. Or, if you want to run software on some exotic new CPU architecture, you can just write a backend that converts the LLVM IR into machine code for that architecture, and you'll automatically be able to compile any language with an LLVM frontend for that architecture.

It's pretty normal for compilers to use some sort of three-address code internally, but the details vary. I've decided to name the intermediate representation in this book *TACKY*. (Naming your intermediate representations is, in my opinion, one of the best parts of compiler design.) I made up TACKY for this book, but it's similar to three-address code in other compilers.

## Defining TACKY

We can define TACKY in ASDL, just like our other intermediate representations. It looks almost, but not quite, like the AST definition from Listing 3-4:

```
program = Program(function_definition)
function_definition = Function(identifier, linstruction*
body)
instruction = Return(val) | Unary(unary_operator, val src,
val dst)
val = Constant(int) | Var(identifier)
unary_operator = Complement | Negate
```

Listing 3-9

Definition of TACKY

In TACKY, a function body consists of a list of instructions<sup>1</sup>, not just a single statement. In this respect, it's similar to the assembly AST we defined in the previous chapter. For now, TACKY has two instructions. *Return* returns a value. *Unary* performs some unary operation on *src*, the source value for the expression, and stores the result in *dst*, the destination. Both of these instructions operate on *vals*, which can be either constant integers (*Constant*) or temporary variables (*Var*).

TACKY makes a couple of assumptions that aren't explicit in Listing 3-9. The first assumption is that the *dst* of a unary operation will be a temporary *Var*, not a *Constant*. Trying to assign a value to a constant wouldn't make sense. The second assumption is that you'll always assign a value to a temporary before you use it. Right now, the only way to assign a value to a variable is by making it the *dst* of a unary operation. There are

two ways to use a variable: by returning it, or by using it as the `src` of a unary operation. Because we're generating TACKY from an AST that we know is valid, we can guarantee that both of those assumptions hold.

You'll need to define a data structure for TACKY, just like you did for the AST and assembly AST. It can be similar to the data structures you used for the AST and assembly AST. For example, if you defined a separate algebraic datatype or abstract class for each node in the assembly AST, you'll want to take the same approach here. Once you have your data structure, you're ready to write the IR generation stage, which converts the AST from Listing 3-5 into TACKY.

## Generating TACKY

Your IR generation pass needs to take an AST in the form defined in Listing 3-5, and return a TACKY AST in the form defined in Listing 3-9. The tricky part is turning an expression into a list of instructions; once you have that figured out, handling all the other AST nodes is easy. Table 3-1 lists a few examples of ASTs and the resulting TACKY:

Table 3-1 TACKY representations of unary expressions

AST	TACKY
<code>Return (Constant (3))</code>	<code>Return (Constant (3))</code>
<code>Return (Unary (Complement, Constant (2)))</code>	<code>Unary (Complement, Constant (2), Var (tmp0)) Return (Var (tmp0))</code>
<code>Return (Unary (Negate, Unary (Complement, Unary (Negate, Constant (8))))</code>	<code>Unary (Negate, Constant (8), Var (tmp0)) Unary (Complement, Var (tmp0), Var (tmp1)) Unary (Negate, Var (tmp1), Var (tmp2)) Return (Var (tmp2))</code>

In the examples above, we convert each unary operation into a `Unary` TACKY instruction, starting with the innermost expression and working our way out. We store the result of each `Unary` instruction in a temporary variable, which we then use in the outer expression or return statement. The

pseudocode in Listing 3-10 describes how to generate these TACKY instructions.

```

emit_tacky(e, instructions):
1   match e with
    | 2 Constant(c) -> return 3 Constant(c)
    | Unary(op, inner) ->
4       src = emit_tacky(inner, instructions)
5       dst_name = make_temporary()
       dst = Var(dst_name)
       tacky_op = convert_unop(op)
6       instructions.append(Unary(tacky_op, src, dst))
       return dst

```

Listing 3-10

Pseudocode to convert an expression into a list of TACKY instructions

This pseudocode emits the instructions needed to calculate an expression by appending them to the `instructions` argument. It also returns a TACKY `val` that represents the result of the expression, which we'll use when we translate the outer expression or statement.

The `match` statement in Listing 3-10 checks which type of expression we're translating, then runs the clause to handle that kind of expression <sup>1</sup>. If the expression is a constant, we'll just return the equivalent TACKY `Constant` without generating any new instructions. Note that the `Constant` constructs at <sup>2</sup> and <sup>3</sup> are different; <sup>2</sup> is a node in the original AST, while <sup>3</sup> is a node in the TACKY AST. (The same is true for the two `Unary` constructs that appear in the next clause.)

If `e` is a unary expression, we'll construct TACKY values for the source and destination. First, we'll call `emit_tacky` recursively on the source expression to get the corresponding TACKY value <sup>4</sup>. This will also generate the TACKY instructions to calculate that value. Then, we'll create a new temporary variable for the destination <sup>5</sup>. The `make_temporary` helper function will generate a unique name for this variable. We'll use another helper function, `convert_unop`, to convert the unary operator to its TACKY equivalent. I won't provide pseudocode for either of these helper functions, since they're very simple. Once we have our source, destination, and unary operator, we'll construct the `Unary` TACKY instruction and append it to the `instructions` list <sup>6</sup>. Finally, we'll return `dst` as the result of the whole expression.

## Testing the TACKY Generator

The TACKY generator should be able to handle every valid test case from this chapter and the one before. To test this stage, we'll run the whole compiler and check whether it succeeds or fails, without inspecting its

output. You can run those tests with:

```
$ ./test_compiler /path/to/your_compiler --chapter 3 --stage
codegen
```

The TACKY stage shouldn't encounter any invalid test cases, because the lexer and parser should catch them first.

## Implementation Tips

**Generate globally unique names.** In the TACKY examples in Table 3-1, it's clear that giving two temporary variables the same name would be an error. In later chapters, we'll want to guarantee that no two temporaries share the same name, even if they're in different functions. That makes IR generation easier, because you don't have to think about whether it's safe for two temporaries to have the same name or not. So, you'll want a convenient way to generate unique names. One easy solution is to maintain a global counter; to generate a unique name, just increment the counter and use its new value (or its new value plus some descriptive string) as your variable name. Because these names won't appear in assembly, they don't need to follow any particular naming convention; they just have to be unique.

**Handle expressions in general, not return statements in particular.** Right now, expressions only appear in return statements, but they'll show up in other kinds of statements in later chapters. Make sure your solution can be extended to handle expressions that aren't in return statements.

## Assembly Generation

TACKY is closer to assembly, but it still doesn't specify exactly which assembly instructions we need. The next step is converting the program from TACKY into the assembly AST we defined in the last chapter. We'll do this in three small compiler passes. First, we'll produce an assembly AST, but still refer to temporary variables directly. Next, we'll replace those variables with concrete addresses on the stack. That step will result in some invalid instructions, because many x64 assembly instructions can't use memory addresses for both operands. So, in the last compiler pass, we'll rewrite the assembly AST to fix any invalid instructions.

## Converting TACKY to Assembly

First, we need to extend the assembly AST we defined in the last chapter. We need a way to represent the `neg` and `not` instructions that we

used in Listing 3-2. We also need to decide how, or whether, we'll represent the function prologue and epilogue in the assembly AST.

We have a few different options for handling the prologue and epilogue. We could go ahead and add the `push`, `pop`, and `sub` instructions to the assembly AST. We could add high-level instructions that correspond to the entire prologue and epilogue, instead of maintaining a 1-1 correspondence between assembly AST constructs and assembly instructions. Or we could omit the function prologue and epilogue entirely, and add them during code emission. The assembly AST below just includes an instruction for decrementing the stack pointer (the third instruction in the function prologue) so we can record how many bytes we need to subtract. Because the rest of the prologue and epilogue are completely fixed, we can easily add them during code emission even if they're not included in the assembly AST. That said, the other approaches to representing the function prologue and epilogue can also work, so feel free to choose whichever seems best to you.

We'll also introduce *pseudoregisters* to represent temporary variables. We'll be able to use pseudoregisters exactly the same way as real registers in the assembly AST; the only difference is that they don't correspond to hardware registers, so we have an unlimited supply of them. Because they aren't real registers, they can't appear in the final assembly program; they'll need to be replaced by real registers or memory addresses in a later compiler pass. For now, we'll assign every pseudoregister to its own address in memory. In Part III, we'll write a *register allocator*, which will speed up the program by assigning as many pseudoregisters as possible to hardware registers instead of memory.

Here's the updated assembly AST, with new parts bolded:

```
program = Program(function_definition)
function_definition = Function(identifier name, instruction*
instructions)
instruction = Mov(operand src, operand dst)
              | Unary(unary_operator, operand)
              | AllocateStack(int)
              | Ret
unary_operator = Neg | Not
operand = Imm(int) | Reg(reg) | Pseudo(identifier) |
Stack(int)
reg = AX | R10
```

Listing 3-11

Assembly definition with unary operators

The `instruction` node has a couple of new constructors to represent our new assembly instructions. We'll represent both new unary instructions with the `Unary` constructor. Since this constructor represents a single `not`

or `neg` instruction, it takes just one operand that's used as both source and destination. The `AllocateStack` constructor represents the third instruction in the function prologue, `subq $n, %rsp`. Its one child, an integer, indicates the number of bytes we'll subtract from RSP.

We also have several new instruction operands. The `Reg` operand can represent either of the two hardware registers we've seen so far: EAX and R10D. The `Pseudo` operand lets us use an arbitrary identifier as a pseudoregister. We'll use this to refer to the temporary variables we produced while generating TACKY. Ultimately, we need to replace every pseudoregister with a location on the stack; we'll represent those with the `Stack` operand, which indicates the stack address at the given offset from RBP. For example, in Listing 3-2 we used `-4(%rbp)` as an operand. We'd represent this as `Stack(-4)` in the assembly AST.

**NOTE** Every hardware register has several aliases, depending on how many bytes of the register you need. EAX refers to the lower 32 bits of the 64-bit RAX register, and R10D refers to the lower 32 bits of the 64-bit R10 register. The names AX and R10B refer to the lower 8 bits of RAX and R10, respectively. Register names in the assembly AST are size-agnostic, so `AX` in Listing 3-11 can refer to the register alias RAX, EAX, or AX, depending on context.

Now we can write a straightforward conversion from TACKY to assembly, given in Table 2-2 below:

Table 3-2 Conversion from TACKY to Assembly

TACKY	Assembly
<b>Top-level constructs</b>	
<code>Program(function_definition)</code>	<code>Program(function_definition)</code>
<code>Function(name, instructions)</code>	<code>Function(name, instructions)</code>
<b>Instructions</b>	
<code>Return(val)</code>	<code>Mov(val, Reg(AX))</code> <code>Ret</code>
<code>Unary(unary_operator, src, dst)</code>	<code>Mov(src, dst)</code> <code>Unary(unary_operator, dst)</code>
<b>Operators</b>	
<code>Complement</code>	<code>Not</code>

Negate	Neg
<b>Operands</b>	
Constant(int)	Imm(int)
Var(identifier)	Pseudo(identifier)

Since our new assembly instructions use the same operation for the source and destination, we just copy the source value into the destination before issuing the `neg` or `not` instruction. Note that we’re not using the `AllocateStack` instruction yet; we’ll add it in the very last stage before code emission, once we know how many bytes we need to allocate. We’re also not using any `Stack` operands; we’ll replace every `Pseudo` operand with a `Stack` operand in the next compiler pass. And we’re not using the R10D register; we’ll introduce it when we rewrite invalid instructions.

## Replacing Pseudoregisters

Next, we’ll write a compiler pass to replace each `Pseudo` operand with a `Stack` operand, leaving the rest of the assembly AST unchanged. In Listing 3-2, we used two stack locations: `-4(%rbp)` and `-8(%rbp)`. We’ll stick with that pattern: the first temporary variable we assign a value to will be at `Stack(-4)`, the next will be at `Stack(-8)`, and so on. We’ll subtract four for each new variable, since every temporary variable is a 4-byte integer. You’ll need to maintain a map from identifiers to offsets as you go, so you can replace each pseudoregister with the same address on the stack every time it appears. For example, if you were processing the following list of instructions:

```
Mov(Imm(2), Pseudo(A))
Unary(Neg, Pseudo(A))
```

you would need to make sure that `Pseudo(A)` was replaced with the same `Stack` operand in both instructions.

This compiler pass should also return the stack offset of the final temporary variable, because that tells us how many bytes to allocate on the stack in the final compiler pass.

## Fixing Up Instructions

Now we need to traverse the assembly AST one more time and make two small fixes. The first fix is inserting the `AllocateStack` instruction at the very beginning of the instruction list in the

`function_definition`. The integer argument to `AllocateStack` should be the stack offset of the last temporary variable we allocated in the previous compiler pass. That way, we'll allocate enough space on the stack to accommodate every address we use. For example, if we replace three temporary variables, replacing the last one with `-12(%rbp)`, we'll insert `AllocateStack(12)` at the front of the instruction list.

The second fix is rewriting invalid `Mov` instructions. When we replaced a bunch of pseudoregisters with stack addresses, we may have ended up with `Mov` instructions where both the source and destination are `Stack` operands. In particular, this will happen if the unary expression in your program has at least one level of nesting. But `mov`, like many other instructions, can't have memory addresses in both the source and the destination. If you try to assemble a program with an instruction like `movl -4(%rbp), -8(%rbp)`, the assembler will reject it. Whenever you encounter an invalid `mov` instruction, you'll need to rewrite it to first copy from the source address into R10D, and then copy from R10D to the destination. For example,

```
| movl -4(%rbp), -8(%rbp)
```

would become

```
| movl -4(%rbp), %r10d
| movl %r10d, -8(%rbp)
```

I've chosen R10D as a scratch register because it doesn't serve any other special purpose. Some registers are used by particular instructions; for example, the `idiv` instruction, which performs division, requires the dividend to be stored in EAX. Other registers are used to pass arguments during function calls. Using any of these registers for scratch at this stage could cause conflicts later. For example, you might copy a function argument into the correct register, but then accidentally overwrite it while using that register to transfer a different value between memory addresses. But because R10D doesn't have any special purpose, we don't have to worry about that kind of conflict.

## Testing Code Generation

Once you've implemented the assembly generation passes, you can test them exactly the same way as the TACKY generator:

```
| $ ./test_compiler /path/to/your_compiler --chapter 3 --stage
codegen
```

## Implementation Tips

**Plan ahead for Part II.** The `Unary` instruction, like `Mov`, will



eventually need to record type information; consider defining it in a way that will make it easier to add type information later on.

**Define a register datatype.** It might seem easiest to store register names as strings, but I think you're better off defining a new datatype to represent them. Like I mentioned earlier, registers in our assembly AST are size-agnostic, but register names in the final assembly program are not. Your code will be clearer if you distinguish between registers in the assembly AST, which aren't yet associated with a particular integer size, and the registers names that will appear in the final assembly program.

## Extending the Code Emitter

Finally, we need to extend the code emission stage to handle our new constructs and print out the function prologue and epilogue. Here's how to print out each construct, with new and changed constructs bolded:

Table 3-3      Formatting assembly

Assembly Construct	Output
<b>Top-level constructs</b>	
Program(function_definition)	(just print out the function definition)
<b>Function(name, instructions)</b>	<pre> .global &lt;name&gt;  &lt;name&gt;:      pushq    %rbp     movq     %rsp, %rbp     &lt;instructions&gt; </pre>
<b>Instructions</b>	
Mov(src, dst)	<pre> movl &lt;src&gt;, &lt;dst&gt; </pre>
<b>Ret</b>	<pre> movq    %rbp, %rsp popq    %rbp ret </pre>
<b>Unary(unary_operator, operand)</b>	<pre> &lt;unary_operator&gt; &lt;operand&gt; </pre>

<b>AllocateStack(int)</b>		subq     \$<int>, %rsp
<b>Operators</b>		
<b>Neg</b>		negl
<b>Not</b>		notl
<b>Operands</b>		
<b>Reg(AX)</b>		%eax
<b>Reg(R10)</b>		%r10d
<b>Stack(int)</b>		<int>(%rbp)
<b>Imm(int)</b>		\$<int>

We'll always insert the function prologue right after the function's label. We'll also emit the whole function epilogue whenever we encounter a single `ret` instruction. Because RBP and RSP contain memory addresses, which are eight bytes, we'll operate on them using quadword instructions, which have a `q` prefix. Note that the program now includes two versions of the `mov` instruction: `movl` and `movq`. They're identical apart from the size of their operands.

## Testing the Whole Compiler

Once you've updated the code emission stage, your compiler should produce correct assembly for all the test cases in this chapter. To test it out, run:

```
$ ./test_compiler /path/to/your_compiler --chapter 3
```

Just like in the previous chapter, this will compile all the valid examples, run them, and verify the return code. It also runs the invalid examples, but those should already fail at the parsing stage.

## Summary

In this chapter, you extended your compiler to implement negation and bitwise complement. You also implemented a new intermediate representation, wrote a couple different compiler passes that transform assembly code, and learned how stack frames are structured. Next, you'll implement binary operations like addition and subtraction. The changes to the backend in the next chapter will be pretty simple; the tricky part will be getting the parser to respect operator precedence and associativity.

# 4

## BINARY OPERATORS

In this chapter, you'll implement five new operators: addition, subtraction, multiplication, division, and the modulo operator. These are all *binary operators*, which take two operands. This chapter won't require any new compiler stages; you'll just need to extend each of the stages you've already written. In the parsing stage, we'll see why recursive descent parsing doesn't work well for binary operators. You'll learn about a different technique, *precedence climbing*, that will be easier to build on in later chapters. Precedence climbing is the last major parsing technique we'll

need. Once it's in place, we'll be able to add new syntax with relatively little effort for the rest of the book. In the code generation stage, we'll introduce several new assembly instructions that perform binary operations. As usual, we'll start with the lexer.

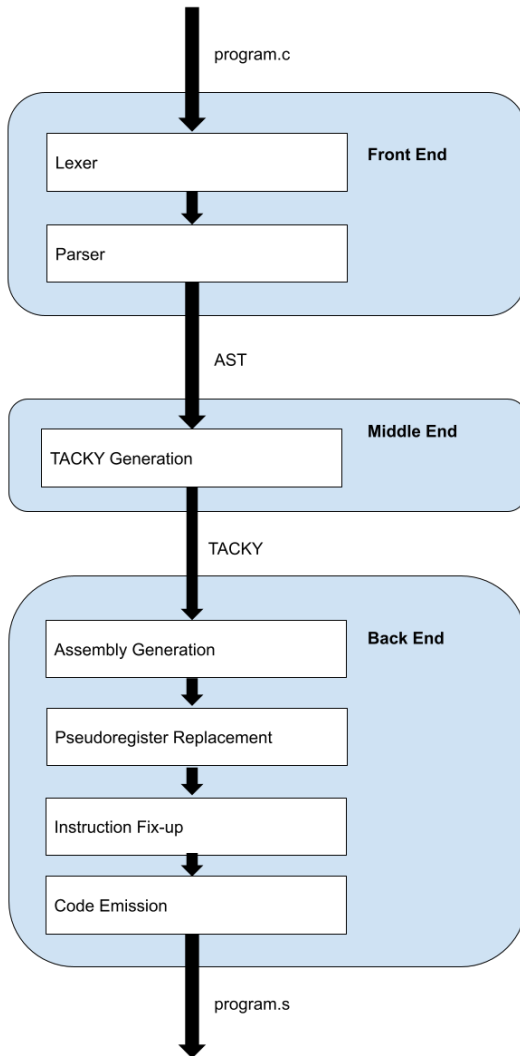


Figure 4-1

Stages of the compiler

## Extending the Lexer

The lexer will need to recognize four new tokens:

- `+` : a plus sign, the operator for addition
- `*` : an asterisk, the operator for multiplication
- `/` : a forward slash, the division operator
- `%` : a percent sign, the modulo operator

This list doesn't include the `-` token, because you already added it in the last chapter. The lexing stage doesn't distinguish between negation and subtraction; it should produce the same token either way.

You can implement these tokens the same way you did the single-character tokens in earlier chapters.

### Testing the Lexer

You know the drill. Your lexer shouldn't fail on any of the test cases in this chapter.

```
$ ./test_compiler /path/to/your_compiler --chapter 4 --stage
lex
```

## Extending the Parser

In this chapter, we'll need to add another kind of expression to the AST: binary operations. Listing 4-1 gives the updated AST definition:

```
program = Program(function_definition)
function_definition = Function(identifier name, statement
body)
statement = Return(exp)
exp = Constant(int)
    | Unary(unary_operator, exp)
    | Binary(binary_operator, exp, exp)
unary_operator = Complement | Negate
binary_operator = Add | Subtract | Multiply | Divide | Mod
```

Listing 4-1

Abstract syntax tree with binary operations

There are a couple things to note about this AST definition. The first is that the parser, unlike the lexer, distinguishes between negation and subtraction. A `-` token will be parsed as either `Negate` or `Subtract`, depending on where it appears in an expression.

The second point is that the structure of the AST determines the order in which we evaluate nested expressions. Let's look at a couple examples to see how the AST's structure controls the order of operations. The AST in Figure 4-2 represents the expression  $1 + (2 * 3)$ , which evaluates to 7.

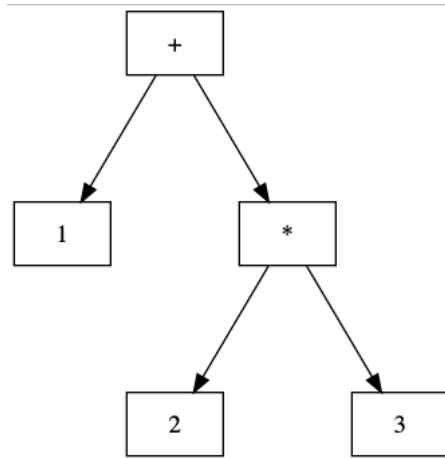


Figure 4-2 AST for  $1 + (2 * 3)$

The + operation has two operands: 1 and  $(2 * 3)$ . If you were going to evaluate this expression, you would need to calculate  $2 * 3$  first, and then add 1 to the result. The AST in Figure 4-3, on the other hand, represents the expression  $(1 + 2) * 3$ , which evaluates to 9:

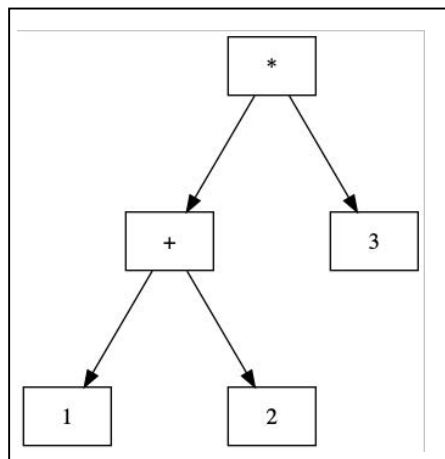


Figure 4-3 AST for  $(1 + 2) * 3$

In this case, you would need to evaluate  $1 + 2$  first, then multiply by

3. As a general rule, before evaluating an AST node you need to evaluate both of its children. This pattern, where you need to process a node's children before you process the node itself, is called *post-order traversal*. (Note that any tree data structure can be traversed in post-order, not just ASTs.)

Your compiler will traverse the AST to generate code, not to evaluate expressions, but the idea is the same. When you convert the AST for a binary expression to TACKY, you need to generate instructions to calculate both operands, then generate instructions for the operator itself. (We also used post-order traversal to process unary operations in the last chapter.)

The point of all this is that it's very important for your parser to group nested expressions correctly, because if you try to parse `1 + (2 * 3)` but end up with the AST from Figure 4-3, you'll end up compiling the program incorrectly.

The examples we just looked at used parentheses to explicitly group nested expressions. Some expressions, like `1 + 2 * 3`, don't parenthesize every nested expression. In those cases, we group expressions based on the *precedence* and *associativity* of the operators. Operators with higher precedence are evaluated first; since `*` has higher precedence than `+`, you'd parse `1 + 2 * 3` as `1 + (2 * 3)`. Associativity tells you how to handle operators at the same precedence level. If an operation is *left-associative*, you apply the operator on the left first, and if it's *right-associative*, you apply the operator on the right first. For example, since addition and subtraction are left-associative, `1 + 2 - 3` would be parsed as `(1 + 2) - 3`. All the new operators in this chapter are left-associative, and there are two precedence levels: `*`, `/`, and `%` have higher precedence, and `+` and `-` have lower precedence.

## The Trouble with Recursive Descent Parsing

It's surprisingly tricky to write a recursive descent parser that handles operator precedence and associativity correctly. To see why, let's try extending the grammar rule for expressions. The obvious rule would look something like this:

```
<exp> ::= 1 <int> | <unop> <exp> | "(" <exp> ")" | 2 <exp>
<binop> <exp>
```

Listing 4-2

A deceptively simple but unworkable grammar rule

The rule in Listing 4-2 corresponds to the AST definition: we added a new constructor to the AST, so it stands to reason that we can just add another production rule 2 to the grammar. But this production rule won't

work with the recursive-descent parsing algorithm we’ve used up to now. If we try to use it, we’ll run into two problems.

The first problem is that rule 2 is *ambiguous*: sometimes it allows you to parse a list of tokens in more than one way. Based on this rule, Figures 4-2 and 4-3 are equally valid parses of `1 + 2 * 3`. We need to know about the relative precedence of `+` and `*` to decide which parse to use, but the rule in Listing 4-2 doesn’t capture that information.

The second problem is that rule 2 is *left-recursive*. That means the left-most symbol in this production rule for `<exp>` is, itself, `<exp>`. You can’t parse a left-recursive rule with a recursive descent parser. First of all, it’s impossible to decide which production rule to apply. Let’s say your input starts with an `<int>` token. Maybe your expression is a single `<int>`, so you should apply production rule 1. Or maybe it’s a more complex expression and the `<int>` is just the first operand of the first sub-expression, so you should apply production rule 2. There’s no way to tell until you’ve parsed some of the input.

Even if you could determine which production rule to use, processing rule 2 would lead to unbounded recursion. The first symbol in this rule is an `<exp>`, so `parse_exp` would have to process that symbol by calling itself recursively. But, because `parse_exp` would be calling itself with exactly the same input, since it didn’t consume any tokens before the recursive call, it would never terminate.

There are a couple of ways to solve these two problems. If you want a pure recursive descent parser, you can refactor the grammar to remove the ambiguity and left-recursion. Since that approach has some drawbacks, we’ll use an alternative to recursive descent parsing called precedence climbing. However, it’s helpful to take a look at the pure recursive-descent solution first.

## The Adequate Solution: Refactoring the Grammar

If we refactor the grammar, we’ll end up with one grammar rule for each precedence level:

```
<exp> ::= <term> { ("+" | "-") <term> }
<term> ::= <factor> { ("*" | "/" | "%") <factor> }
<factor> ::= <int> | <unop> <factor> | "(" <exp> ")"
```

Listing 4-3

A recursive descent-friendly grammar for binary operations

Using the grammar in Listing 4-3, there’s only one way to parse `1 + 2 * 3`, and there’s no left recursion. The curly braces indicate repetition, so a



single `<exp>`, for example, can contain any number of `<term>`s. It might be a `<term>`, or `<term> + <term>`, or `<term> - <term> + <term>`, and so on. The parser then groups that long string of terms or factors into a left-associative tree. (Note that we can't use a rule like `<exp> ::= <term> "+" <exp>` because it would result in a right-associative tree.)

This approach works, but it gets more unwieldy as you add more precedence levels. We have three precedence levels now, if you count `<factor>`; we'll add four more when we introduce logical and relational operators in the next chapter. If we went with this approach, we'd need to add a new symbol to the grammar—and a corresponding function to our parser—for each precedence level we add. That's a lot of boilerplate, since the functions to parse all the different binary expressions will be almost identical.

## ***The Better Solution: Precedence Climbing***

Precedence climbing is a simpler way to parse binary expressions; it can handle production rules like `<exp> <binop> <exp>`. The basic idea is that every operator will have a numeric precedence level, and `parse_exp` will take a minimum precedence level as an argument. That lets you specify the appropriate precedence level for whatever sub-expression you're parsing. For example, let's say you just saw a `+` token, and now want to parse what comes next as the right-hand side of an addition expression—you would specify that it should only include operations that are higher-precedence than `+`. This solution makes it easy to add new operators; you have to assign those operators an appropriate numeric precedence level, but otherwise your parsing code doesn't need to change.

## **Mixing Precedence Climbing with Recursive Descent**

Luckily, we can use precedence climbing here without rewriting the recursive descent parsing code we wrote earlier. We'll write a hybrid parser that uses precedence climbing for binary expressions, and recursive descent for everything else. Remember that in a recursive descent parser, we define one parsing function to handle each symbol. That makes it straightforward to mix the two approaches: we can just use precedence climbing in the `parse_exp` function, and recursive descent in the functions that parse all the other symbols. The `parse_exp` function will remove tokens from the input stream and return an `exp` AST node, just like a recursive descent-based parsing function would. But it will use a different strategy to get that result.

Since we already know how to parse unary and parenthesized expressions with recursive descent, let's represent those with a separate symbol from binary operations. That will make it easier to parse the two types of expressions using different techniques. Here's the resulting grammar:

```
<program> ::= <function>
<function> ::= "int" <identifier> "(" ")" "{" <statement>
}"
<statement> ::= "return" <exp> ";"
<exp> ::= <factor> | <exp> <binop> <exp>
<factor> ::= <int> | <unop> <factor> | "(" <exp> ")"
<unop> ::= "-" | "~"
<binop> ::= "-" | "+" | "*" | "/" | "%"
<identifier> ::= ? An identifier token ?
<int> ::= ? A constant token ?
```

Listing 4-4

The final grammar to handle binary operations

A `<factor>` (which we were calling an `<exp>` in the last chapter) can be parsed with the usual recursive descent approach. (We'll keep calling this symbol a "factor," like we do in Listing 4-3, since it can appear as a term in a multiplication, division, or modulo expression.) It looks almost exactly like last chapter's rule for `<exp>`, except that we now allow binary operations as well as factors inside parentheses. That means `(1 + 2)` is a factor, because `"(" <exp> ")"` is a production rule for `<factor>`. However, `-1 + 2` is not, because `"-" <exp>` is not a production rule for `<factor>`. Because the rules for `<exp>` and `<factor>` refer to each other, the functions to parse those symbols will be mutually recursive. An `<exp>` is either a binary operation, defined in the obvious way, or it's just a factor.

The pseudocode to parse factors also looks almost the same as last chapter:

```
parse_factor(tokens):
    next_token = peek(tokens)
    if next_token is an int:
        --snip--
    else if next_token is "~" or "-":
        operator = parse_unop(tokens)
1[] inner_exp = parse_factor(tokens)
        return Unary(operator, inner_exp)
    else if next_token == "(":
2[] take_token(tokens)
        inner_exp = parse_exp(tokens)
        expect(")", tokens)
        return inner_exp
    else:
        fail()
```

Listing 4-5

Pseudocode for parsing a factor

The only difference is that we call `parse_factor` where we expect a `<factor>` **1**, and `parse_exp` where we expect an `<exp>` **2**; before, we just called `parse_exp` in both places.

## Making Operators Left-Associative

Next, we need to figure out what `parse_exp` looks like. First, let's make the problem simpler by only considering the `+` and `-` operators, which are both at the same precedence level. To handle these operators, `parse_exp` needs to group expressions in a left-associative way, but it doesn't need to handle multiple precedence levels yet.

In this simple case, we'll encounter inputs like `factor1 + factor2 - factor3 + factor4`. These should always be parsed in a left-associative way to produce expressions like `((factor1 + factor2) - factor3) + factor4`. As a result, the right operand of every expression, including sub-expressions, will be a single factor. For example, the right operand of `(factor1 + factor2)` is `factor2`, and the right operator of `((factor1 + factor2) - factor3)` is `factor3`.

Once we realize that the right operand of an expression is always a single factor, we can write pseudocode to parse these expressions:

```

parse_exp(tokens):
1  left = parse_factor(tokens)
   next_token = peek(tokens)
2  while next_token is "+" or "-":
   operator = parse_binop(tokens)
3   right = parse_factor(tokens)
4  left = Binary(operator, left, right)
   next_token = peek(tokens)
return left

```

Listing 4-6

Parsing left-associative expressions without considering precedence level

In Listing 4-6, we start by parsing a single factor **1**. This factor will be either the whole expression or the left operand of a larger expression. Then, we check if the next token is a binary operator **2**. If it is, we consume it from the input and convert it to an AST node. Then we construct a binary expression where the left operand is everything we've parsed so far and the right operand is the next factor **4**, which we get by calling `parse_factor` **3**. We repeat this process until we see a token other than `+` or `-` after a factor; that means there are no binary expressions left to construct, so we're done.

## Dealing with Precedence

Listing 4-6 lets us parse left-associative binary operators, but it doesn't handle different precedence levels. Now let's extend it to handle `*`, `/`, and `%`. These operators are also left-associative, but they're at a higher precedence level than `+` and `-`.

Once we add these operators, the right operand of every expression can be either a single factor, or a sub-expression involving only the new, higher-precedence operators. For example, `1 + 2 * 3 + 4` would be parsed as `(1 + (2 * 3)) + 4`. The right operand of the whole expression is a single factor, 4. The right operand of the inner sub-expression, `1 + (2 * 3)`, is a product, `2 * 3`.

We can be even more precise. If the outermost expression is a `+` or `-` operation, its right operand will only contain factors, `*`, `/`, and `%`. But if the outermost expression is itself a `*`, `/`, or `%` operation, its right operand must be single factor.

To generalize: whenever we're parsing an expression of the form `e1 <op> e2`, all the operators in `e2` should be higher-precedence than `<op>`. We can achieve this by tweaking the code from Listing 4-6:

```

parse_exp(tokens, min_prec):
    left = parse_factor(tokens)
    next_token = peek(tokens)
    while next_token is a binary operator and
precedence(next_token) >= min_prec:
        operator = parse_binop(tokens)
        right = parse_exp(tokens, precedence(next_token) +
1)
        left = Binary(operator, left, right)
        next_token = peek(tokens)
    return left

```

Listing 4-7

Parsing left-associative operators with precedence climbing

This pseudocode is our entire precedence climbing algorithm. The `min_prec` argument lets us state that all operators in the sub-expression we're currently parsing need to exceed some precedence level. For example, we could include only operators that are higher-precedence than `+`. We enforce this by comparing the precedence of the current operator to `min_prec` at each iteration of the while loop; we exclude the operator and anything that follows it from the current expression if its precedence is too low. Then, when we process the right-hand side of an operation, we set the minimum precedence higher than the precedence of the current operator. This guarantees that higher-precedence operators will be evaluated first.

Since operators at the same precedence level as the current operator won't be included in the right-hand expression, the resulting AST will be left-associative.

When you're calling `parse_exp` from any other function (including from `parse_factor`, to handle parenthesized expressions), you'll start with a minimum precedence of zero, so the result includes operators at every precedence level.

The code in Listing 4-7 requires us to assign every binary operator a precedence value; the values I've assigned are listed in Table 4-1.

Table 4-1      Precedence Values of Binary Operators

Operator	Precedence
<code>*</code>	50
<code>/</code>	50
<code>%</code>	50
<code>+</code>	45
<code>-</code>	45

The exact precedence values don't matter, as long as higher-precedence operators have higher values. The numbers I've chosen here give us plenty of room to add new lower-precedence operators in the next chapter.

## Precedence Climbing in Action

Let's walk through an example where we parse the following expression:

```
1 * 2 - 3 * (4 + 5)
```

We'll trace the execution of our precedence-climbing code (Listing 4-7) as it parses this expression. In each code snippet below, I've added a level of indentation inside each function call, to make it easier to track how deep we are in the call stack.

We'll start by calling `parse_exp` on the whole expression with a minimum precedence of zero:

```
parse_exp("1 * 2 - 3 * (4 + 5)", 0):
```

Inside `parse_exp`, we'll start by parsing the first factor:

```
left = parse_factor("1 * 2 - 3 * (4 + 5)")
      = Constant(1)
next_token = ""
```

This first call to `parse_factor` will just parse the token `1`, returning `Constant(1)`. Next, we peek at the token that follows, which is `*`. This token is a binary operator, and its precedence is greater than zero, so we enter the `while` loop.

The first iteration of the loop looks like this:

```
// loop iteration #1
operator = parse_binop("* 2 - 3 * (4 + 5)")
          = ""
right = parse_exp("2 - 3 * (4 + 5)", 51)
        left = parse_factor("2 - 3 * (4 + 5)")
              = Constant(2)
        next_token = "-"
        // precedence(next_token) < 51
        = Constant(2)
left = Binary(*, Constant(1), Constant(2))
next_token = ""
```

Inside the loop, `parse_binop` consumes `next_token` from the input, which leaves `2 - 3 * (4 + 5)`. Next, we need to call `parse_exp` recursively to get the right-hand side of this product. Since the precedence of `*` is 50, the second argument to `parse_exp` will be 51. In the recursive call, we again get the next factor (`2`) and the token that follows it (`-`). The `-` token is a binary operator, but its precedence is only 45; it doesn't meet the minimum precedence of 51, so we don't enter the while loop. Instead, we return `Constant(2)`.

Back in the outer call to `parse_exp`, we use `Binary` to construct the AST node for `1 * 2` from the values we've parsed so far. Then, we check the next token to see whether we have more sub-expressions to process. The next token is `-`; we peeked at it, but didn't remove it from the input, inside the recursive call to `parse_exp`. Because `-` is a binary operator, and it exceeds our minimum precedence of zero, we jump back to the beginning of the `while` loop to parse the next sub-expression:

```
// loop iteration #2
operator = parse_binop("- 3 * (4 + 5)")
          = "-"
right = parse_exp("3 * (4 + 5)", 46)
        left = parse_factor("3 * (4 + 5)")
              = Constant(3)
        next_token = ""
        // loop iteration #1
```

```

        operator = parse_binop("* (4 + 5)")
        = "*"
    right = parse_exp("(4 + 5)", 51)
        left = parse_factor("(4 + 5)")
            parse_exp("4 + 5", 0)
            --snip--
            = Binary(+, Constant(4),
Constant(5))
            = Binary(+, Constant(4),
Constant(5))
            = Binary(+, Constant(4), Constant(5))
        left = Binary(*, Constant(3), Binary(+,
Constant(4), Constant(5)))
        = Binary(*, Constant(3), Binary(+, Constant(4),
Constant(5)))
        left = Binary(-,
            Binary(*, Constant(1), Constant(2)),
            Binary(*, Constant(3), Binary(+,
Constant(4), Constant(5))))

```

The second time through the loop, we consume `-` from the input and make a recursive call to `parse_exp`. This time, because the precedence of `-` is 45, the second argument to `parse_exp` will be 46.

Following our usual routine, we get the next factor (3) and the next token (\*). Since the precedence of \* exceeds the minimum precedence, we need to parse another sub-expression. We consume \*, leaving (4 + 5), then make yet another recursive call to `parse_exp`.

In this next call to `parse_exp`, we start by calling `parse_factor` as usual. This call will consume the rest of our input and return the AST node for 4 + 5. To handle that parenthesized expression, `parse_factor` will need to recursively call `parse_exp` with the minimum precedence reset to zero, but we won't step through that here. At this point, there are no tokens left in our expression. Let's assume this is a valid C program and the next token is a semicolon. Since the next token isn't a binary operator, we just return the expression we got from `parse_factor`.

At the next level up, we construct the AST node for 3 \* (4 + 5) from the sub-expressions we've processed in this call. Once again, we peek at the next token, see that it isn't a binary operator, and return.

Finally, back in the original call to `parse_exp`, we construct the final expression from the left operand that we constructed in the first loop iteration (1 \* 2), the current value of `next_token` (-), and the right operand that was just returned from the recursive call (3 \* (4 + 5)). For the last time, we check the next token, see that it isn't a binary operator, and

return.

Now that we've seen how to parse binary expressions with precedence climbing, you're ready to extend your own parser. Remember that you'll use precedence climbing to parse binary expressions, and recursive descent to parse all the other symbols in the grammar, including factors.

## FURTHER READING ON PRECEDENCE CLIMBING

These blog posts helped me understand precedence climbing, and how it relates to similar algorithms that solve the same problem. You might find them helpful too.

- "Parsing expressions by precedence climbing," a blog post by Eli Bendersky, provides a good overview of the precedence climbing algorithm. It also covers right-associative operators, which I didn't discuss here. (<https://eli.thegreenplace.net/2012/08/02/parsing-expressions-by-precedence-climbing>)
- "Some problems of recursive descent parsers," also by Eli Bendersky, goes into more detail about how to handle binary expressions with a pure recursive descent parser. (<https://eli.thegreenplace.net/2009/03/14/some-problems-of-recursive-descent-parsers>)
- Andy Chu has written two useful blog posts on precedence climbing. The first, "Pratt Parsing and Precedence Climbing are the Same Algorithm" explores the fundamental similarities between these two approaches (<https://www.oilshell.org/blog/2016/11/01.html>). The second, "Precedence Climbing is Widely Used," discusses their differences (<https://www.oilshell.org/blog/2017/03/30.html>). These posts clarify some of the confusing terminology around different parsing algorithms.

## Testing the Parser

The parser should be able to handle every valid test case in [tests/chapter\\_4/valid](#), and raise an error on every invalid test case in [tests/chapter\\_4/invalid\\_parse](#). To test your parser against the test cases from this chapter and the ones before it, run:

```
$ ./test_compiler /path/to/your_compiler --chapter 4 --stage parse
```

Remember that the test suite only checks whether your compiler parses a program successfully or throws an error; it doesn't check that it produced the correct AST. In this chapter, it's especially easy to write a parser that appears to succeed but generates the wrong AST, so you might want to write your own tests to validate the output of your parser.



## Extending TACKY Generation

Next, we need to update the stage that converts the AST to TACKY. We'll start by updating TACKY itself to include binary operations:

```
program = Program(function_definition)
function_definition = Function(identifier, instruction*
body)
instruction = Return(val)
              | Unary(unary_operator, val src, val dst)
              | Binary(binary_operator, val src1, val src2,
val dst)
val = Constant(int) | Var(identifier)
unary_operator = Complement | Negate
binary_operator = Add | Subtract | Multiply | Divide | Mod
```

Listing 4-8

Adding binary expressions to TACKY

The changes here are pretty straightforward: we've just added one new type of instruction to represent binary operations, and defined all the possible operators. Like unary operations, binary operations in TACKY can only be applied to constants and variables, not to nested sub-expressions. We can turn a binary expression into a sequence of TACKY instructions in almost exactly the same way that we handled unary expressions:

```
emit_tacky(e, instructions):
    match e with
    | --snip--
    | Binary(op, e1, e2) ->
        v1 = emit_tacky(e1, instructions)
        v2 = emit_tacky(e2, instructions)
        dst_name = make_temporary()
        dst = Var(dst_name)
        tacky_op = convert_binop(op)
        instructions.append(Binary(tacky_op, v1, v2, dst))
        return dst
```

Listing 4-9

Converting a binary operation to TACKY

We need to emit the TACKY instructions to calculate each operand, then emit the binary expression that uses those source values. The only difference from how we handled unary expression is that we're processing two operands instead of one.

### NO, *YOU'RE* OUT OF ORDER!

In Listing 4-9, we generate code that evaluates the first operand, then the second operand, then the whole operation. Surprisingly, it would be just as correct to evaluate the second operand before the first. According to the C standard, sub-expressions of the same operation are *unsequenced*—they can be evaluated in any order. In the programs we can compile so far, it doesn't matter which

operand we evaluate first; you'll get the same visible behavior either way. That's not the case in the following program:

```
#include <stdio.h>

int main() {
    return printf("Hello, ") + printf("World!");
}
```

You could compile this program with a C standard-compliant compiler, run it, and get either of the following outputs:

```
Hello, World!
World!Hello,
```

There are a few exceptions where the first operand must be evaluated first: the logical `&&` and `||` operators, which we'll cover next chapter; the conditional `?:` operator, which we'll cover a few chapters later; and the comma operator, which we won't implement.

If you're curious, the relevant part of the C18 standard is section 6.5, paragraphs 1-3. There's also a more readable explanation at [https://en.cppreference.com/w/c/language/eval\\_order](https://en.cppreference.com/w/c/language/eval_order).

## Testing TACKY Generation

The TACKY generator should be able to process every valid test case we've seen so far. You can test it with:

```
$ ./test_compiler /path/to/your_compiler --chapter 4 --stage
codegen
```

## Extending Assembly Generation

The next step is converting TACKY into assembly. We'll need several new assembly instructions to handle addition, subtraction, multiplication, division, and the modulo operation. Let's talk through these new instructions.

## Doing Arithmetic in Assembly

The instructions for addition, subtraction, and multiplication all take the form `op src, dst`, where:

`op` is an instruction,

`src` is an immediate value, register, or memory address, and

`dst` is a register or memory address.

Each of these instructions applies `op` to `dst` and `src`, storing the result in `dst`. The instructions for addition, subtraction, and multiplication are

Instruction	Meaning
addl    \$2, %eax	eax = eax + 2
subl    \$2, %eax	eax = eax - 2
imull   \$2, %eax	eax = eax * 2

These instructions are pretty easy to use and understand! In a perfect world, we could perform division in exactly the same way. But we don't live in a perfect world, which is why we're stuck with the `idiv` instruction.

In order to use `idiv`, we need to turn a 32-bit dividend into a 64-bit value spanning both EDX and EAX. Whenever we need to convert a signed integer to a wider format, we'll use an operation called *sign extension*. This operation fills the upper 32 bits of the new, 64-bit value with the sign bit of the original 32-bit value.

[illegible]

into

```
0000000000000000000000000000000000000000000000000000000000000000
0011
```

Both representations have the value 3; the second one just has more leading zeros. To sign extend a negative number, we fill the upper four bytes with ones, which converts -3 from

```
1111111111111111111111111111111111111111111111111111111111111111
1101
```

into

```
1111111111111111111111111111111111111111111111111111111111111111
1101
```

Thanks to the magic of two's complement, the value of both of these binary numbers is -3. (If you're not clear on how this works, you can check out the further reading on two's complement from Chapter 3.)

The `cdq` instruction does exactly what we need here: it sign extends the value from EAX into EDX. If the number in EAX is positive, this instruction will set EDX to all zeros. If EAX is negative, this instruction will set EDX to all ones. Putting it all together, here's how you'd compute `9 / 2`, or `9 % 2`, in assembly:

```
movl    $2, %ebx
movl    $9, %eax
cdq
idiv    %ebx
```

The result of `9 / 2`, the quotient, will be stored in EAX. The result of `9 % 2`, the remainder, will be stored in EDX.

Now we've covered all the new instructions we'll need in this chapter: `add`, `sub`, `imul`, `idiv`, and `cdq`. Next, let's add these new instructions to the assembly AST and update the conversion from TACKY to assembly.

## Converting TACKY to Assembly

Here's the updated assembly AST, with additions bolded:

```
program = Program(function_definition)
function_definition = Function(identifier name, instruction*
instructions)
instruction = Mov(operand src, operand dst)
              | Unary(unary_operator, operand)
              | Binary(binary_operator, operand, operand)
              | Idiv(operand)
              | Cdq
              | AllocateStack(int)
              | Ret
```

```

unary_operator = Neg | Not
binary_operator = Add | Sub | Mult
operand = Imm(int) | Reg(reg) | Pseudo(identifier) |
Stack(int)
reg = AX | DX | R10 | R11

```

Listing 4-10 Adding new instructions to the assembly AST

Since the addition, subtraction, and multiplication instructions all take the same form, we'll represent them all using the `Binary` instruction node. We'll also add instruction nodes for the new `idiv` and `cdq` instructions. We'll add the EDX register to the AST definition, since the `idiv` instruction uses it. We'll also add the R11 register to use along with R10 during the instruction fix-up pass.

Now we need to convert our new binary operations from TACKY to assembly. For addition, subtraction, and multiplication, we'll convert a single TACKY instruction into two assembly instructions:

```
Binary(op, src1, src2, dst)
```

becomes

```

Mov(src1, dst)
Binary(op, src2, dst)

```

Division is a little more complicated; we need to move the first operand into EAX, sign-extend it with `cdq`, issue the `idiv` instruction, and then move the result from EAX to the destination. So

```
Binary(Divide, src1, src2, dst)
```

becomes

```

Mov(src1, Reg(AX))
Cdq
Idiv(src2)
Mov(Reg(AX), dst)

```

The modulo operation looks exactly the same, except that we ultimately want to retrieve the remainder from EDX instead of retrieving the quotient from EAX. So

```
Binary(Mod, src1, src2, dst)
```

becomes

```

Mov(src1, Reg(AX))
Cdq
Idiv(src2)
Mov(Reg(DX), dst)

```

The `idiv` instruction can't operate on immediate values, so the assembly instructions for division and modulo won't be valid if `src2` is a

constant. That’s okay; we’ll fix it during the instruction-rewriting pass.  
 Table 4-3 summarizes the conversion from TACKY to assembly.

Table 4-3 Conversion from TACKY to Assembly

<b>TACKY</b>	<b>Assembly</b>
<b>Top-level constructs</b>	
<code>Program(function_definition)</code>	<code>Program(function_definition)</code>
<code>Function(name, instructions)</code>	<code>Function(name, instructions)</code>
<b>Instructions</b>	
<code>Return(val)</code>	<code>Mov(val, Reg(AX))</code> <code>Ret</code>
<code>Unary(unary_operator, src, dst)</code>	<code>Mov(src, dst)</code> <code>Unary(unary_operator, dst)</code>
<code>Binary(Divide, src1, src2, dst)</code>	<code>Mov(src1, Reg(AX))</code> <code>Cdq</code> <code>Idiv(src2)</code> <code>Mov(Reg(AX), dst)</code>
<code>Binary(Mod, src1, src2, dst)</code>	<code>Mov(src1, Reg(AX))</code> <code>Cdq</code> <code>Idiv(src2)</code> <code>Mov(Reg(DX), dst)</code>
<code>Binary(binary_operator, src1, src2, dst)</code>	<code>Mov(src1, dst)</code> <code>Binary(binary_operator, src2, dst)</code>
<b>Operators</b>	
<code>Complement</code>	<code>Not</code>
<code>Negate</code>	<code>Neg</code>
<code>Add</code>	<code>Add</code>
<code>Subtract</code>	<code>Sub</code>
<code>Multiply</code>	<code>Mult</code>
<b>Operands</b>	

<code>Constant(int)</code>	<code>Imm(int)</code>
<code>Var(identifier)</code>	<code>Pseudo(identifier)</code>

Note that the table above includes three rows for the `Binary` TACKY instruction—one for division, one for modulo, and one for everything else.

## Replacing Pseudoregisters

You'll need to update this pass to handle the new `Binary` and `Idiv` instructions. You can handle them exactly like the existing `Mov` and `Unary` instructions. When you see a pseudoregister in a `Mov`, `Unary`, or `Binary` instruction, replace it with the corresponding stack address. If the register hasn't been assigned a stack address yet, assign it to the next available 4-byte address.

## Fixing Up Instructions

In the last compiler pass before emitting the final program, we rewrite invalid instructions that we produced in earlier stages. Now we need to add a couple more rewrite rules. First, we need to fix `idiv` instructions that take constant operands. Whenever `idiv` needs to operate on a constant, we can just copy that constant into our scratch register first. So

```
| idiv $3
```

is rewritten as

```
| movl $3, %r10d
| idiv %r10d
```

The `add` and `sub` instructions, like `mov`, can't use memory addresses as both source and destination operands. We can rewrite them in the same way as `mov`, so that

```
| addl -4(%rbp), -8(%rbp)
```

becomes

```
| movl -4(%rbp), %r10d
| addl %r10d, -8(%rbp)
```

The `imul` instruction can't use a memory address as a destination, regardless of its source operand. When we need to fix an instruction's destination operand, we'll use the `R11` register instead of `R10`. To fix `imul`, we'll load the destination into `R11`, multiply it by the source operand, and then store the result back to the destination address, so

```
| imull $3, -4(%rbp)
```

becomes

```
movl -4(%rbp), %r11d
imull $3, %r11d
movl %r11d, -4(%rbp)
```

Using different registers to fix source and destination operands will become helpful in Part II, when we'll sometimes need to rewrite the source and destination for the same instruction. We'll need two registers so that the fix-up instructions for the different operands don't clobber each other.

Once you've updated the assembly-generating, pseudoregister-replacing, and instruction-fixing compiler passes, your compiler should be able to generate AST representations of complete, correct assembly programs that perform basic arithmetic. All that's left is emitting those assembly programs in the right format.

## Testing Assembly Generation

To test the assembly generation stages, run:

```
$ ./test_compiler /path/to/your_compiler --chapter 4 --stage
codegen
```

## Extending the Code Emitter

The last step is extending the code emission stage to handle the new assembly instructions we added in this chapter. Here's how to print out each construct, with new and changed constructs bolded:

Table 4-3: Formatting assembly

Assembly Construct	Output
<b>Top-level constructs</b>	
Program(function_definition)	(just print out the function definition)
Function(name, instructions)	<pre> .global &lt;name&gt;  &lt;name&gt;:     pushq    %rbp     movq     %rsp, %rbp     &lt;instructions&gt; </pre>



Instructions		
Mov(src, dst)		movl <src>, <dst>
Ret		<div> <div>movq      %rbp, %rsp</div> <div>popq      %rbp</div> <div>ret</div> </div>
Unary(unary_operator, operand)		<unary_operator> <operand>
Binary(binary_operator, src, dst)		<binary_operator> <src>, <dst>
Idiv(operand)		idivl      <operand>
Cdq		cdq
AllocateStack(int)		subq      \$<int>, %rsp
Operators		
Neg		negl
Not		notl
Add		addl
Sub		subl
Mult		imull
Operands		
Reg(AX)		%eax
Reg(DX)		%edx
Reg(R10)		%r10d
Reg(R11)		%r11d
Stack(int)		<int>(%rbp)
Imm(int)		\$<int>

All the new instructions operate on 32-bit values, so they get **l** suffixes. Note that the `subl` instruction we use to subtract integers and the `subq` instruction that we use to allocate space on the stack are just 32-bit and 64-bit versions of the same instruction.

## Testing the Whole Compiler

Now you're ready to try compiling programs all the way through. To check if you're compiling every test program correctly, run:

```
$ ./test_compiler /path/to/your_compiler --chapter 4
```

## Extra Credit: Bitwise Operators

Now that you've learned how to compile a few binary operators, you know enough to implement the bitwise binary operators on your own. These include bitwise AND (&), OR (|), XOR (^), left shift (<<), and right shift (>>). Your compiler can handle these much like the operators you just added. You'll need to look up the relative precedence of these operators. You'll also need to check the documentation for the x64 instruction set to see how to use the relevant assembly instructions.

Bitwise operations are optional; later test cases don't rely on them. If you do want to implement bitwise operations, you can use the `--bitwise` flag to include test cases that use them, like this:

```
$ ./test_compiler /path/to/your_compiler --chapter 4 --  
bitwise
```

You'll want to include this flag when you run test cases in later chapters too, so that those test cases also include bitwise operators.

## Summary

In this chapter, you implemented several binary arithmetic operations in your compiler. You learned how to use a new technique, precedence climbing, to parse expressions that recursive descent parsers don't handle well. In the next chapter, you'll implement even more unary and binary operations: logical operators like `!`, `&&` and `||`, and relational operators like `==`, `<`, and `>`. Some of these operators don't correspond closely to assembly instructions, so we'll break them down into lower-level instructions in TACKY. We'll also introduce conditional assembly instructions, which will be particularly important when we implement control-flow statements like `if` statements and loops.

# 5

## LOGICAL AND RELATIONAL OPERATORS

Now that you know how to compile binary operators, we're going to add a whole mess of them (plus one more unary operator). We'll cover the logical NOT (!), AND (&&) and OR (||) operators, plus all the relational operators: <, >, ==, and so on. Each of these operators tests some condition, returning 1 if that condition is true, and 0 if it's false.

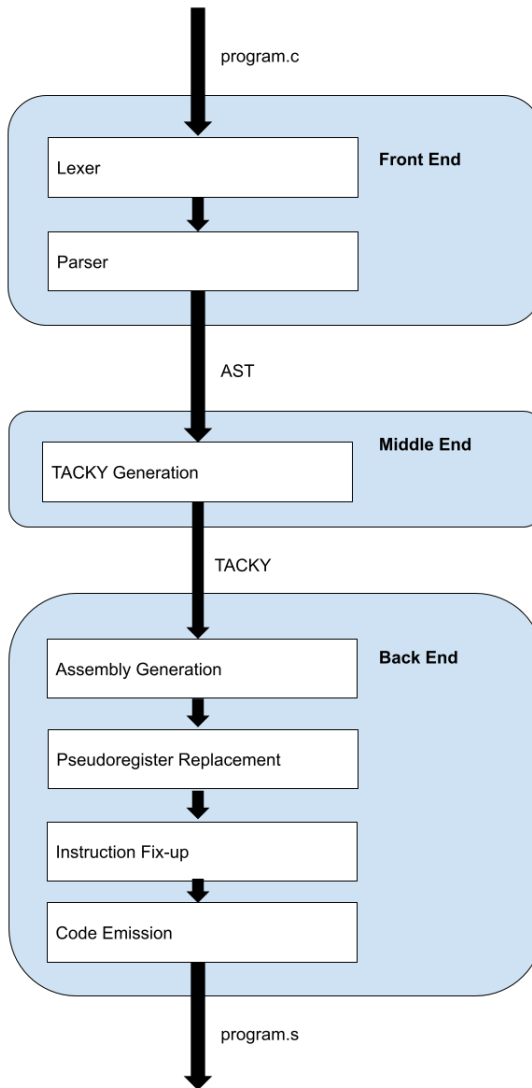


Figure 5-1 Stages of the compiler

The `&&` and `||` operators differ from the binary operators we've seen so far because they *short-circuit*: if you know the result after the first operand, you don't evaluate the second operand. To support short-circuiting logic, we'll add new instructions to TACKY that let us skip over blocks of code. Then, during the code generation pass, we'll introduce several new instructions, including conditional assembly instructions that let us take specific actions only if some condition is met. We'll see how the CPU relies on special-purpose hardware registers to implement these instructions.

Let’s talk about short-circuiting operators before moving on to the compiler passes.

## Short-Circuiting Operators

The C standard guarantees that `&&` and `||` short-circuit when you don’t need the second operand. For example, consider the expression `0 && foo()`. Because the first clause is zero, we know the whole expression will evaluate to zero regardless of what `foo` returns, so we won’t call `foo` at all. Likewise, if the first operand of `||` is non-zero, we don’t evaluate the second operand.

This isn’t just a performance optimization; the second operand might not change the result of the expression, but evaluating it can have visible side effects. For example, the `foo` function might perform I/O or update global variables. If your compiler doesn’t implement `&&` and `||` as short-circuiting operators, some compiled programs will behave incorrectly. (For the record, the standard defines this behavior in section 6.5.13, paragraph 4 for the `&&` operator, and 6.5.14, paragraph 4 for the `||` operator.)

Now that we’ve clarified how these operators work, you’re ready to continue coding.

## Extending the Lexer

In this chapter, you’ll need to add nine new tokens:

`!` : an exclamation point, the logical NOT operator

`&&` : two ampersands, the logical AND operator

`||` : two vertical bars, the logical OR operator

`==` : two equals signs, the “equal to” operator

`!=` : an exclamation point followed by an equal sign, the “not equal to” operator

`<` : the “less than” operator

`>` : the “greater than” operator

`<=` : the “less than or equal” operator

`>=` : the “greater than or equal” operator

Your lexer can handle these the same way as the other operators you’d

added so far. Remember that the lexer should always choose the longest possible match for the next token. If your input is `<=something`, the next token the lexer emits should be `<=`, not `<`.

## Testing the Lexer

You can test the lexer with the usual command:

```
$ ./test_compiler /path/to/your_compiler --chapter 5 --stage
lex
```

Your lexer should succeed on all of this chapter's test cases.

## Extending the Parser

We'll need to add all our new operations to the AST definition. Listing 5-1 gives the updated definition, with new additions bolded:

```
program = Program(function_definition)
function_definition = Function(identifier name, statement
body)
statement = Return(exp)
exp = Constant(int)
    | Unary(unary_operator, exp)
    | Binary(binary_operator, exp, exp)
unary_operator = Complement | Negate | Not
binary_operator = Add | Subtract | Multiply | Divide | Mod |
And | Or
                    | Equal | NotEqual | LessThan | LessOrEqual
                    | GreaterThan | GreaterOrEqual
```

Listing 5-1 Abstract syntax tree with new operations

We'll also need to make corresponding changes to the grammar, which is given in Listing 5-2:

```
<program> ::= <function>
<function> ::= "int" <identifier> "(" ")" "{" <statement>
"}"
<statement> ::= "return" <exp> ";"
<exp> ::= <factor> | <exp> <binop> <exp>
<factor> ::= <int> | <unop> <factor> | "(" <exp> ")"
<unop> ::= "-" | "~" | "!"
<binop> ::= "-" | "+" | "*" | "/" | "%" | "&&" | "||"
    | "==" | "!=" | "<" | "<=" | ">" | ">="
<identifier> ::= ? An identifier token ?
<int> ::= ? A constant token ?
```

Listing 5-2 Grammar with new operations

In Listings 5-1 and 5-2, we've added some new operators, but haven't made any other changes. Now we're ready to update the parsing code. To handle the new `!` operator, you'll need to change `parse_factor`. Your parsing code will treat `!` exactly like the unary `~` and `-` operators you've already implemented.

Next, you'll need to update `parse_exp` to handle all the new binary operators. Remember that in the last chapter, we associated every binary operator with a numeric precedence value. Now, we need to give the new operators precedence values. All the new operators have lower precedence than the ones we've already implemented, and they're all left-associative. Among the new operators, `<`, `<=`, `>`, and `>=` have highest precedence, followed by the equality operators `==` and `!=`. The `&&` operator has lower precedence than any of the relational operators, and `||` has the lowest precedence of all. The precedence values I've chosen are listed in Table 5-1.

Table 5-1      Precedence Values of Old and New Binary Operators

Operator	Precedence
*	50
/	50
%	50
+	45
-	45
<	35
<=	35
>	35
>=	35
==	30
!=	30
&&	10
	5

I chose these values to accommodate the relative precedence of all these operators, plus the optional bitwise operators from the previous chapter. You don't need to use the exact values in this table as long as operators have the same precedence relative to each other.



You'll also need to extend the code that converts tokens into `unary_operator` and `binary_operator` AST nodes. For example, whatever function converts a `+` token into an `Add` node should also be able to convert a `==` token into an `Equal` node. (The pseudocode in the last two chapters called separate functions, `parse_unop` and `parse_binop`, to handle that conversion.)

Once you've updated your parser's table of precedence values, `parse_binop`, and `parse_unop`, you're done! The precedence climbing algorithm we implemented in the last chapter will be able to handle all the new operators without any further changes.

## Testing the Parser

The parser should be able to handle every valid test case in *tests/chapter\_5/valid*, and raise an error on every invalid test case in *tests/chapter\_5/invalid\_parse*. To test your parser, run:

```
$ ./test_compiler /path/to/your_compiler --chapter 5 --stage
parse
```

Once your lexer and parser are working properly, we can venture into less familiar territory: handling the new operators in TACKY.

## Extending TACKY Generation

You can convert relational operators to TACKY in exactly the same way as the binary operators we've already implemented. For example, given the expression `e1 < e2`, the resulting TACKY looks something like this:

```
<instructions for e1>
tmp1 = <result of e1>
<instructions for e2>
tmp2 = <result of e2>
Binary(LessThan, tmp1, tmp2, dst)
```

Listing 5-3

Structure of TACKY for a binary expression

But you can't generate the `&&` and `||` operators this way, because they need a way to short-circuit. The code in Listing 5-3 always evaluates both `e1` and `e2`, but we need to generate code that sometimes skips `e2`. To support short-circuiting operators, we'll add a *jump* instruction, which lets us jump to a different point in the program. We'll also add two *conditional jump* instructions, which only jump if a particular condition is met.

Listing 5-4 shows these new jump instructions, along with the other

## additions to TACKY:

```

program = Program(function_definition)
function_definition = Function(identifier, instruction*
body)
instruction = Return(val)
              | Unary(unary_operator, val src, val dst)
              | Binary(binary_operator, val src1, val src2,
val dst)
              | Copy(val src, val dst)
              | Jump(identifier target)
              | JumpIfZero(val condition, identifier target)
              | JumpIfNotZero(val condition, identifier
target)
              | Label(identifier)
val = Constant(int) | Var(identifier)
unary_operator = Complement | Negate | Not
binary_operator = Add | Subtract | Multiply | Divide | Mod |
Equal | NotEqual
                  | LessThan | LessOrEqual | GreaterThan |
GreaterOrEqual

```

Listing 5-4

Adding conditionals jumps and labels to TACKY

The `Jump` instruction works just like `goto` in C—it causes the program to jump to the point labeled with some identifier, `target`. The `Label` instruction associates an identifier with a location in the program. The following snippet of TACKY shows how `Jump` and `Label` instructions work together:

```

1 Unary(Negate, Constant(1), Var(tmp))
2 Jump("there")
3 Unary(Negate, Constant(2), Var(tmp))
4 Label("there")
5 Return(Var(tmp))

```

Listing 5-5

Example TACKY with a Jump instruction

First, this program will store `-1` in `tmp` **1**. The `Jump` instruction **2** will make it jump to **4** and then execute the return statement at **5**, which will return `-1`. Instruction **3** won't execute at all, because we jumped over it.

The first conditional jump, `JumpIfZero`, says: if the value `condition` is zero, jump to the instruction indicated by `target`. If `condition` is anything other than zero, don't jump to `target`—instead, just execute the next instruction as usual. The second conditional jump, `JumpIfNotZero`, does the opposite: we jump to `target` only if `condition` isn't zero. We don't really need both of these instructions, since any behavior you can express with one can also be expressed with the other plus a `Not` instruction. But adding both of them will let us generate

simpler TACKY for both the `&&` and `||` operations, which will ultimately translate into simpler, shorter assembly.

We'll need one more instruction: `Copy`. Since `&&` and `||` ultimately return 1 or 0, we need this instruction to copy a 1 or 0 into the temporary variable that holds the result of the expression.

Besides these five new instructions, the latest TACKY definition includes all the new binary operators and the unary `Not` operator. Note that TACKY doesn't include binary `And` or `Or` operators, because we won't implement them as binary operations. Instead, we'll implement them using the jump instructions. The TACKY for the expression `e1 && e2` should look something like this:

```

<instructions for e1>
v1 = <result of e1>
1 JumpIfZero(v1, false_label)
<instructions for e2>
v2 = <result of e2>
2 JumpIfZero(v2, false_label)
3 result = 1
4 Jump(end)
5 Label(false_label)
6 result = 0
  Label(end)

```

Listing 5-6

TACKY for `&&` operation

We start by evaluating `e1`. If it's zero, we want to short-circuit and set `result` to 0, without evaluating `e2`. We can accomplish this with the `JumpIfZero` instruction 1; if `v1` is zero, we jump straight to `false_label` 5, and then set `result` to 0 with the `Copy` instruction 6. (I've written this out as `result = 0` instead of `Copy(0, result)` to make it a bit more readable.) If `v1` isn't zero we still need to evaluate `e2`. We can handle the case where `v2` equals zero exactly like the case where `v1` equals zero—by jumping to `false_label`. Once again, we can do this with `JumpIfZero` 2. We'll only reach instruction 3 if we didn't take either conditional jump. That means both `e1` and `e2` must be non-zero, so we set `result` to 1 (using the `Copy` instruction again). Then we need to jump over 6 to the `end` label to avoid overwriting `result`.

Note that Listing 5-6 includes a couple of labels. Labels, like temporary variables, must be globally unique—an instruction like `Jump("foo")` is useless if the label `foo` shows up in multiple places.

Labels differ from temporary variable names in one important way. They'll appear in the final assembly program, so they must be identifiers

that the assembler considers syntactically valid. You can make sure your labels are syntactically valid by including only letters, digits, and underscores, and starting each label with a letter.

You can translate the `||` operation to TACKY in a very similar way, using the `JumpIfNotZero` instruction. I'll leave it to you to implement this on your own. That leaves `!` and all the relational operations; like I mentioned earlier in this section, you can convert these to TACKY in exactly the same way as the unary and binary operations from earlier chapters.

## Testing TACKY Generation

You can test the TACKY generation pass with:

```
$ ./test_compiler /path/to/your_compiler --chapter 5 --stage
codegen
```

## Implementation Tips

**Generate descriptive labels.** Because labels appear in the final assembly program, informative labels can make the program easier to read and debug. For example, when generating instructions for `&&` like the ones in Listing 5-6, the label at `5` could be something like `and_falseN`, where `N` is some globally unique counter.

## Extending Assembly Generation

Before starting on the assembly generation pass, let's talk through the new assembly instructions we'll need. First, we'll discuss the `cmp` instruction, which compares two values, and the *conditional set* instructions, which can set a byte to 1 or 0 based on the result of a comparison. With these instructions, we can implement relational operators like `<`. Next, we'll talk about conditional and unconditional jump instructions.

## Comparisons and Status Flags

The “condition” that all conditional instructions depend on is the state of the RFLAGS register. Unlike EAX, RSP, and the other registers we've encountered, you usually can't set RFLAGS directly. Instead, the CPU updates RFLAGS automatically every time it issues an instruction. Like the name suggests, each bit in the register is a flag that reports some fact about the last instruction or the status of the CPU. Different instructions update different flags—the `add`, `sub`, and `cmp` instructions update all the flags

we'll talk about below, and the `mov` instruction doesn't update any of them. Other instructions have other effects that we can ignore for now. Whenever I refer to the "last instruction" or "last result" while discussing RFLAGS, I mean the last instruction that affects the particular flag I'm talking about.

Right now, we only care about three of these flags:

### The Zero Flag (ZF)

ZF is set to 1 if the result of the last instruction was zero. It's set to 0 if the result of the last instruction was non-zero.

### The Sign Flag (SF)

SF is set to 1 if the most significant bit of the last result was 1. It's set to 0 if the most significant bit of that result is 0. Remember that in two's complement, the most significant bit of a negative number is always 1, and the most significant bit of a positive number is always 0. That means the sign flag can tell us whether the result of the last instruction was positive or negative. (If the last result should be interpreted as an unsigned integer, it can't be negative, and the sign flag is meaningless.)

### The Overflow Flag (OF)

OF is set to 1 if the last instruction resulted in a signed integer overflow, and 0 otherwise. An *integer overflow* occurs when the result of a signed integer operation can't be represented in the number of bits available. A positive result will overflow if it's larger than the maximum value the type can hold. For example, suppose we're operating on 4-bit integers. The largest signed number we can represent is `0111`, or 7. If we add 1 to it with the `add` instruction, the result will be `1000`. This is 8 if we interpret it as an unsigned integer, but -8 if we interpret it as a signed integer. The result of the computation should be positive, but since it overflowed, it appears negative. This computation will set the overflow flag to 1.

We'll also encounter integer overflow in the opposite situation: when the result should be negative, but it's below the smallest possible value. For example, in ordinary math,  $-8 - 1 = -9$ . But if we use the `sub` instruction to subtract 1 from the 4-bit binary representation of -8, which is `1000`, we'll end up with `0111`, or 7. The overflow flag will be set to 1 in this case too.

The result of an unsigned operation can also be too large or small for its type to represent, but we won't refer to this as integer overflow. Instead, we'll just say the result *wrapped around*; that's more consistent with the terminology for unsigned operations in the C standard and in

most discussions of x64 assembly. We draw this distinction because unsigned wrap-around follows different rules from signed integer overflow in the C standard, and the CPU detects it differently. We'll learn exactly how to handle it in Part II.

Tables 5-2 and 5-3 summarize the cases where each kind of integer overflow can occur. Note that Table 5-3 is just Table 5-2 with the columns swapped, since  $A - B$  and  $A + (-B)$  are equivalent. Like SF, OF is meaningless if the result is unsigned.

Table 5-2 Integer overflow and underflow from addition

$A + B$	$B > 0$	$B < 0$
$A > 0$	Overflow from positive to negative	Neither
$A < 0$	Neither	Overflow from negative to positive

Table 5-3 Integer overflow and underflow from subtraction

$A - B$	$B > 0$	$B < 0$
$A > 0$	Neither	Overflow from positive to negative
$A < 0$	Overflow from negative to positive	Neither

The instruction `cmp b, a` computes  $a - b$ , exactly like the `sub` instruction, and has exactly the same impact on RFLAGS, but discards the result instead of storing it in `a`. This is a bit more convenient when you only want to subtract two numbers in order to compare them, but don't necessarily want to overwrite `a`. Let's think about how the instruction `cmp b, a` would impact ZF and SF.

- If  $a == b$ , then  $a - b$  will be 0, so ZF will be 1 and SF will be 0.

- If  $a > b$ , then  $a - b$  will be a positive number, so both SF and ZF will be 0.
- If  $a < b$ , then  $a - b$  will be a negative number, so SF will be 1 and ZF will be 0.

By issuing a `cmp` instruction and then referring to ZF and SF, you can handle every comparison we're implementing in this chapter. But wait! That's not quite true, because  $a - b$  could overflow, which will flip SF. Let's consider how that impacts each of the three cases above:

- If  $a == b$ , then  $a - b$  can't overflow because it's 0.
- If  $a > b$ , then  $a - b$  could overflow when  $a$  is positive and  $b$  is negative. The correct result in this case is positive, but if it overflows, the actual result will be negative. In that case, SF will be 1, and OF will be too.
- If  $a < b$ , then  $a - b$  could overflow when  $a$  is negative and  $b$  is positive. In this case, the correct result is negative, but the actual result will be positive. That means SF will be 0, but OF will be 1.

Table 5-4 gives the values of these flags in every case we've considered.

Table 5-4      Impact of `cmp` instruction on status flags

	ZF	OF	SF
<b>A == B</b>	1	0	0
<b>A &gt; B, no overflow</b>	0	0	0
<b>A &gt; B, overflow</b>	0	1	1
<b>A &lt; B, no overflow</b>	0	0	1
<b>A &lt; B, overflow</b>	0	1	0

Note that you can tell whether  $a$  or  $b$  is larger by checking whether SF and OF are the same. If they are, we know that  $a \geq b$ . Either both are 0, because we got a positive (or zero) result with no overflow, or both are 1, because we got a large positive result that overflowed until it became

negative. If SF and OF are different, we know that  $a < b$ . Either we got a negative result with no overflow, or a negative result that overflowed and became positive.

## UNDEFINED BEHAVIOR ALERT!

If the `add` and `sub` instructions can overflow, why didn't we account for that in Chapter 4? We didn't need to because integer overflow in C is *undefined behavior*, where the standard doesn't tell you what should happen. Compilers are permitted to handle undefined behavior however they want, or ignore it completely.

When an expression in C overflows, for example, the result will *usually* wrap around like the examples we saw earlier. However, it's equally correct for the program to generate a result at random, or raise a signal, or erase your hard drive. That last option may seem unlikely, but production compilers really do handle some undefined behavior in surprising (and arguably undesirable) ways. Take the following program:

```
#include <stdio.h>

int main() {
    for (int i = 2147483646; i > 0; i = i + 1)
        printf("The number is %d\n", i);
    return 0;
}
```

*A program with integer overflow*

The largest value an `int` can hold is 2,147,483,647, so the expression `i + 1` will overflow the second time we execute it. We know that when the `add` assembly instruction overflows, it produces a negative result, so we might expect this loop to execute twice, then stop because the condition `i > 0` no longer holds. And, in fact, that's exactly what happens if you compile Listing 5-7 without optimizations, at least with the versions of Clang and GCC that I tried:

```
$ clang overflow.c
$ ./a.out
The number is 2147483646
The number is 2147483647
```

But if you enable optimizations, the behavior might change completely:

```
$ clang -O overflow.c
$ ./a.out
The number is 2147483646
The number is 2147483647
The number is -2147483648
The number is -2147483647
The number is -2147483646
The number is -2147483645
--snip--
```

What happened? The compiler tried to optimize the program by removing conditional checks that always succeed. Because we initialized `i` to a positive number, and then only incremented it, the compiler concluded that `i > 0` would always be true—which is correct as long as `i` doesn't overflow. It's incorrect if `i` does overflow, of course, but the compiler isn't required to account for that



case. So it removed that condition entirely, resulting in a loop that never terminates.

I used Clang for this example because GCC produced a completely different, even less intuitive behavior. You may well see different results if you compile Listing 5-7 on your own machine. Try it out with a few different optimization levels, and see what happens.

Note that setting the overflow flag in assembly doesn't necessarily indicate overflow in the source program. For example, when we implement a comparison like `a < 10` with `cmp`, that `cmp` instruction may set the overflow flag. But the result of the comparison is 0 or 1, which we can obviously represent as an `int`. Therefore, a comparison in a C program won't produce undefined behavior, regardless of how exactly we implement it in assembly.

These blog posts go into more detail about undefined behavior and the trail of chaos and destruction it leaves in its wake:

- “A Guide to Undefined Behavior in C and C++, Part 1,” by John Regehr, is a good overview of what undefined behavior means in the C standard and how it impacts compiler design. (<https://blog.regehr.org/archives/213>).
- “With Undefined Behavior, Anything is Possible,” by Raph Levien, explores some sources of undefined behavior in C and the history of how it got into the standard to begin with. (<https://raphlinus.github.io/programming/rust/2018/08/17/undefined-behavior.html>).

Now that we understand how to set ZF, OF, and SF, let's take a look at a few instructions that depend on those flags.

## Conditional Set Instructions

To implement a relational operator, we first set some flags using the `cmp` instruction, and then set the result of the expression based on those flags. We'll perform that second step with a *conditional set* instruction. Each conditional set instruction takes a single register or memory address as an operand, which it sets to 0 or 1 based on the state of RFLAGS. The conditional set instructions are all identical, except they test for different conditions. Table 5-5 lists the conditional set instructions we need in this chapter:

One annoying thing about conditional set instructions is that they only set a single byte. If you want to conditionally set EAX to 0 or 1, for example, the instruction must refer to the AL register, which is the least significant byte of EAX. You also need to zero out EAX first, because the conditional set instruction won't clear its upper bytes. For example, if EAX is

and you run

then the new value in EAX will be:

which is, obviously, not 0. The `sete` instruction zeroed out the last byte of EAX, but not the rest of it.

The `jmp` assembly instruction takes a label as an argument, and performs an unconditional jump to that label. Jump assembly instructions manipulate another special-purpose register, RIP. The RIP register always

holds the address of the next instruction to execute (IP stands for “instruction pointer”). To execute a sequence of instructions, the CPU carries out the *fetch-execute cycle*:

1. Fetch an instruction from the memory address in RIP, and store it in a special-purpose *instruction register*. (This register doesn’t have a name because you can’t refer to it at all in assembly.)
2. Increment RIP. Instructions in x64 aren’t all the same length, so the CPU has to inspect the instruction it just fetched, figure out how many bytes long it is, and increment RIP by that many bytes.
3. Run the instruction in the instruction register.
4. Repeat.

Normally, this means that the CPU will execute instructions in the order they appear in memory. But `jmp` puts a new value in RIP, which changes what instruction the CPU executes next. The assembler and linker convert the label in a jump instruction into a *relative offset* that tells you how much to increment or decrement RIP. Consider the following snippet of assembly:

```

    |      addl $1, %eax
    |      jmp foo
1  |      movl $0, %eax
    | foo:
2  |      ret

```

Listing 5-7

Assembly code using the `jmp` instruction

Instruction **1** in Listing 5-8 turns out to be five bytes long. If you want to jump over it and execute instruction **2** instead, you need to increment RIP by an extra five bytes. The assembler and linker will convert `jmp foo` into the machine instruction for `jmp 5`. Then, when the CPU executes this instruction, it will:

1. Fetch the instruction `jmp 5` and store it in the instruction register.
2. Increment RIP to point to the next instruction, `mov $0, %eax`.
3. Execute `jmp 5`. This will add five bytes to RIP, so that it points to `ret`.
4. Fetch the instruction RIP points to, `ret`, and continue the fetch-execute cycle from there.

Note that labels aren't instructions: the CPU doesn't execute them, and they don't appear in the text section of the final executable (the section that contains machine instructions). If you're curious, you can see exactly what's in the text section of an executable with `objdump`; that's how I figured out how long instruction `1` is and verified that label `foo` is resolved to the relative offset `5`. Appendix B explains how to inspect executables with `objdump`.

A *conditional jump* takes a label as an argument, but only jumps to that label if the condition holds. Conditional jumps look a lot like conditional set instructions; they depend on exactly the same conditions, using exactly the same flags in RFLAGS. For example, suppose you wanted to return `3` if two registers were equal, and `0` if they weren't equal. You could write:

```

    cmp %eax, %edx
    je return3
1:   mov $0, %eax
2:   ret
return3:
    mov $3, %eax
    ret

```

Listing 5-8

Assembly code using conditional jumps

If the values in EAX and EDX are equal, `cmp` will set ZF to 1, so at `je` we'll jump to `return3`. Then we'll execute the two instructions following `return3`, causing the function to return `3`. If EAX and EDX aren't equal, we won't jump at `je`. Instead, we'll execute instruction `1` and `2`, causing the function to return `0`. Similarly, `jne` jumps only if ZF is 0. There are also jump instructions that check other conditions, but we don't need them in this chapter.

Now that we've covered jumps, comparisons, and conditional instructions, we're ready to extend the assembly AST and update the assembly generation pass.

## Converting TACKY to Assembly

Here's the latest assembly AST, with additions bolded:

```

program = Program(function_definition)
function_definition = Function(identifier name, instruction*
instructions)
instruction = Mov(operand src, operand dst)
             | Unary(unary_operator, operand)
             | Binary(binary_operator, operand, operand)
             | Cmp(operand, operand)
             | Idiv(operand)

```

```

| Cdq
| Jmp(identifier)
| JmpCC(cond_code, identifier)
| SetCC(cond_code, operand)
| Label(identifier)
| AllocateStack(int)
| Ret

unary_operator = Neg | Not
binary_operator = Add | Sub | Mult
operand = Imm(int) | Reg(reg) | Pseudo(identifier) |
Stack(int)
cond_code = E | NE | G | GE | L | LE
reg = AX | DX | R10 | R11

```

Listing 5-9

Assembly AST with comparisons and conditional instructions

Since all conditional jump instructions have the same form, we'll represent them all with a single `JmpCC` instruction, and just distinguish between them using different condition codes. We'll do the same with conditional set instructions. It's easiest to treat labels like instructions at this stage; however, `Label` isn't really an instruction, since labels aren't executed by the CPU.

We'll implement the `JumpIfZero` and `JumpIfNotZero` instructions from TACKY with the new `JmpCC` instruction. For example, we'll convert

```
JumpIfZero(val, target)
```

to

```
Cmp(Imm(0), val)
JmpCC(E, target)
```

We can implement `JumpIfNotZero` exactly the same way, just by changing the condition code from `E` to `NE`.

Similarly, we can implement all the relational operators using conditional set instructions. For example, the following TACKY instruction:

```
Binary(GreaterThan, src1, src2, dst)
```

becomes

```
Cmp(src2, src1)
Mov(Imm(0), dst)
SetCC(G, dst)
```

For all the other relational operators, just replace `G` with the appropriate condition code. Remember that we have to zero out the destination before

the conditional set instruction, since it only sets the lowest byte. It's safe to perform a `mov` right after the `cmp` instruction because `mov` doesn't change RFLAGS. The one remaining wrinkle is that `SetCC` needs a one-byte operand, but `dst` is four bytes; luckily, we can account for this during code emission. If `dst` is a location in memory, `SetCC` will just operate on the first byte at that location, which is the behavior we want. (Because x64 processors are *little-endian*, the first byte is the least significant, so setting that byte to 1 sets the whole 32-bit value to 1.)

If `dst` is a register, it's a little more complicated: the `Mov` instruction will refer to a 32-bit register, like EAX, and the `SetCC` instruction will refer to the corresponding 8-bit register, like AL. We'll make sure to print out the right register name during code emission. At this stage, we refer to registers without specifying their size, so 8-bit registers don't require any special handling.

Because `!x` is equivalent to `x == 0`, we can also implement the unary `!` operator with a conditional set instruction. We'll convert this TACKY instruction:

```
| Unary(Not, src, dst)
```

into this list of assembly instructions:

```
| Cmp(Imm(0), src)
| Mov(Imm(0), dst)
| SetCC(E, dst)
```

The remaining TACKY instructions, `Jump`, `Label`, and `Copy`, are easy. A TACKY `Jump` becomes an assembly `Jump`, `Label` becomes `Label`, and `Copy` becomes `Mov`. Table 5-6 summarizes how to convert each new TACKY construct to assembly.

Table 5-6 Conversion from TACKY to Assembly

TACKY	Assembly
<b>Top-level constructs</b>	
<code>Program(function_definition)</code>	<code>Program(function_definition)</code>
<code>Function(name, instructions)</code>	<code>Function(name, instructions)</code>
<b>Instructions</b>	
<code>Return(val)</code>	<code>Mov(val, Reg(AX))</code>  <code>Ret</code>
<code>Unary(Not, src, dst)</code>	<code>Cmp(Imm(0), src)</code>  <code>Mov(Imm(0), dst)</code>  <code>SetCC(E, dst)</code>
<code>Unary(unary_operator, src, dst)</code>	<code>Mov(src, dst)</code>  <code>Unary(unary_operator, dst)</code>
<code>Binary(Divide, src1, src2, dst)</code>	<code>Mov(src1, Reg(AX))</code>  <code>Cdq</code>  <code>Idiv(src2)</code>  <code>Mov(Reg(AX), dst)</code>

<code>Binary(Mod, src1, src2, dst)</code>	<code>Mov (src1, Reg (AX))</code>  <code>Cdq</code>  <code>Idiv(src2)</code>  <code>Mov (Reg (DX), dst)</code>
<code>Binary(arithmetic_operator, src1, src2, dst)</code>	<code>Mov(src1, dst)</code>  <code>Binary(arithmetic_operator, src2, dst)</code>
<code>Binary(relational_operator, src1, src2, dst)</code>	<code>Cmp(src2, src1)</code>  <code>Mov(Imm(0), dst)</code>  <code>SetCC(relational_operator, dst)</code>
<code>Jump(target)</code>	<code>Jmp(target)</code>
<code>JumpIfZero(condition, target)</code>	<code>Cmp(Imm(0), condition)</code>  <code>JmpCC(E, target)</code>
<code>JumpIfNotZero(condition, target)</code>	<code>Cmp(Imm(0), condition)</code>  <code>JmpCC(NE, target)</code>
<code>Copy(src, dst)</code>	<code>Mov(src, dst)</code>
<code>Label(identifier)</code>	<code>Label(identifier)</code>



<b>Operators</b>	
Complement	Not
Negate	Neg
Add	Add
Subtract	Sub
Multiply	Mult
<b>Operands</b>	
Constant(int)	Imm(int)
Var(identifier)	Pseudo(identifier)
<b>Conditions</b>	
Equal	E
NotEqual	NE
LessThan	L
LessOrEqual	LE
GreaterThan	G
GreaterOrEqual	GE

Once your compiler can handle the conversion from TACKY to assembly, you're ready to move on to the rest of the code generation pass.

## Replacing Pseudoregisters

You should replace any pseudoregisters used by the new `Cmp` and `SetCC` instructions with stack addresses, just like you're doing for all the other instructions.

## Fixing Up Instructions

The `cmp` instruction, much like the arithmetic instructions we've already implemented, can't use memory addresses for both operands. We'll rewrite it in the usual way, so that

```
|  cmpl -4(%rbp), -8(%rbp)
```

will become

```
|  movl -4(%rbp), %r10d  
|  cmpl %r10d, -8(%rbp)
```

The second operand of a `cmp` instruction can't be a constant. This sort of makes sense if you remember that `cmp` follows the same form as `sub`—the second operand of a `sub`, `add`, or `imul` instruction can't be a constant either, since that operand holds the result. Even though `cmp` doesn't produce a result, the same rules apply. So

```
|  cmpl %eax, $5
```

will become

```
|  movl $5, %r11d  
|  cmpl %eax, %r11d
```

Following the convention that we established in the previous chapter, we use R10 when we need to fix a `cmp` instruction's first operand, and R11 when we need to fix its second operand.

## Testing Assembly Generation

To test the assembly generation stage, run:

```
|  $ ./test_compiler /path/to/your_compiler --chapter 5 --stage  
|  codegen
```

## Extending the Code Emitter

We've now generated a valid assembly program, and we're ready to emit it. Code emission is slightly more complicated in this chapter, for two reasons. The first reason is that we're now dealing with both 8-bit and 32-bit registers. You'll print out a different name for a register depending on

whether it appears in a conditional set instruction, which takes 8-bit operands, or any of the other instructions we’ve encountered so far, which take 32-bit operands.

The second issue is emitting labels. Right now, assembly labels come from two places: some are autogenerated inside the compiler, and some—function names—come from user-defined identifiers. Right now, the only function name is `main`, but eventually we’ll compile programs with arbitrary function names. Because labels must be unique, our auto-generated labels can’t conflict with any function names that could possibly appear in a program. We can avoid conflicts by using labels that aren’t syntactically valid function names. An easy solution is to just add a period at the beginning of each label, since periods can’t appear in the names of C functions. For example, if you autogenerated label `foo123` during the TACKY generation stage, you can emit it as `.foo123`, so it won’t cause problems if the function name `foo123` appears elsewhere in the program.

If you’re compiling on macOS, you don’t have to mangle your labels at all. Remember that your code emission stage already adds underscores to user-defined labels (so that `main` becomes `_main`, for example). As long as your autogenerated labels don’t start with underscores, you don’t have to worry that they’ll conflict with user-defined identifiers.

Otherwise, code emission is pretty straightforward; Table 5-7 summarizes how to print out each construct, with this chapter’s additions bolded:

Table 5-7      Formatting assembly

Assembly Construct	Output	
Top-level constructs		
Program(function_definition)	(just print out the function definition)	
Function(name, instructions)		<pre>        .globl &lt;name&gt; &lt;name&gt;:         pushq    %rbp         movq     %rsp, %rbp         &lt;instructions&gt;</pre>
Instructions		

Mov(src, dst)		movl <src>, <dst>
Ret		movq %rbp, %rsp popq %rbp ret
Unary(unary_operator, operand)		<unary_operator> <operand>
Binary(binary_operator, src, dst)		<binary_operator> <src>, <dst>
Idiv(operand)		idivl <operand>
Cdq		cdq
AllocateStack(int)		subq \$<int>, %rsp
Cmp(operand, operand)		cmpl <operand>, <operand>
Jmp(label)		jmp .<label>
JmpCC(cond_code, label)		j<cond_code> .<label>
SetCC(cond_code, operand)		set<cond_code> <operand>
Label(label)		.<label>:
Operators		

Neg		negl
Not		notl
Add		addl
Sub		subl
Mult		imull
Condition Codes		
E		e
NE		ne
L		l
LE		le
G		g
GE		ge
Operands		

Reg(AX)	4-byte		%eax
	1-byte		%al
Reg(DX)	4-byte		%edx
	1-byte		%dl
Reg(R10)	4-byte		%r10d
	1-byte		%r10b
Reg(R11)	4-byte		%r11d
	1-byte		%r11b
Stack(int)			<int>(%rbp)
Imm(int)			\$<int>

On macOS, you don't need the period before the label in `JmpCC` and `Label` instructions, although it doesn't hurt to include it. Make sure to emit the one-byte version of registers when they appear in `SetCC`, and the four-byte version anywhere else. Note that the `cmp` instruction gets a `l` suffix to indicate that it operates on 32-bit values, but `set` instructions don't. That's because conditional set instructions only take one-byte operands; there are no four-byte or eight-byte variants. Because these instructions only support one possible operand size, you don't need a suffix to indicate which size you want.

## Testing the Whole Compiler

Now you should be able to compile and run programs that use our new operators. To check if you're compiling every test program correctly, run:

```
$ ./test_compiler /path/to/your_compiler --chapter 5
```

## Summary

Your compiler can now handle programs with relational and logical operators. You added conditional jumps to TACKY to support short-circuiting operators, and learned about several new assembly instructions, including conditional instructions. You also learned about how the CPU keeps track of the current instruction and how it records the result of comparisons. The new TACKY and assembly instructions we introduced in this chapter will eventually help you implement complex control structures like `if` statements and loops. But first, we'll implement one of the most essential features of C: variables!