

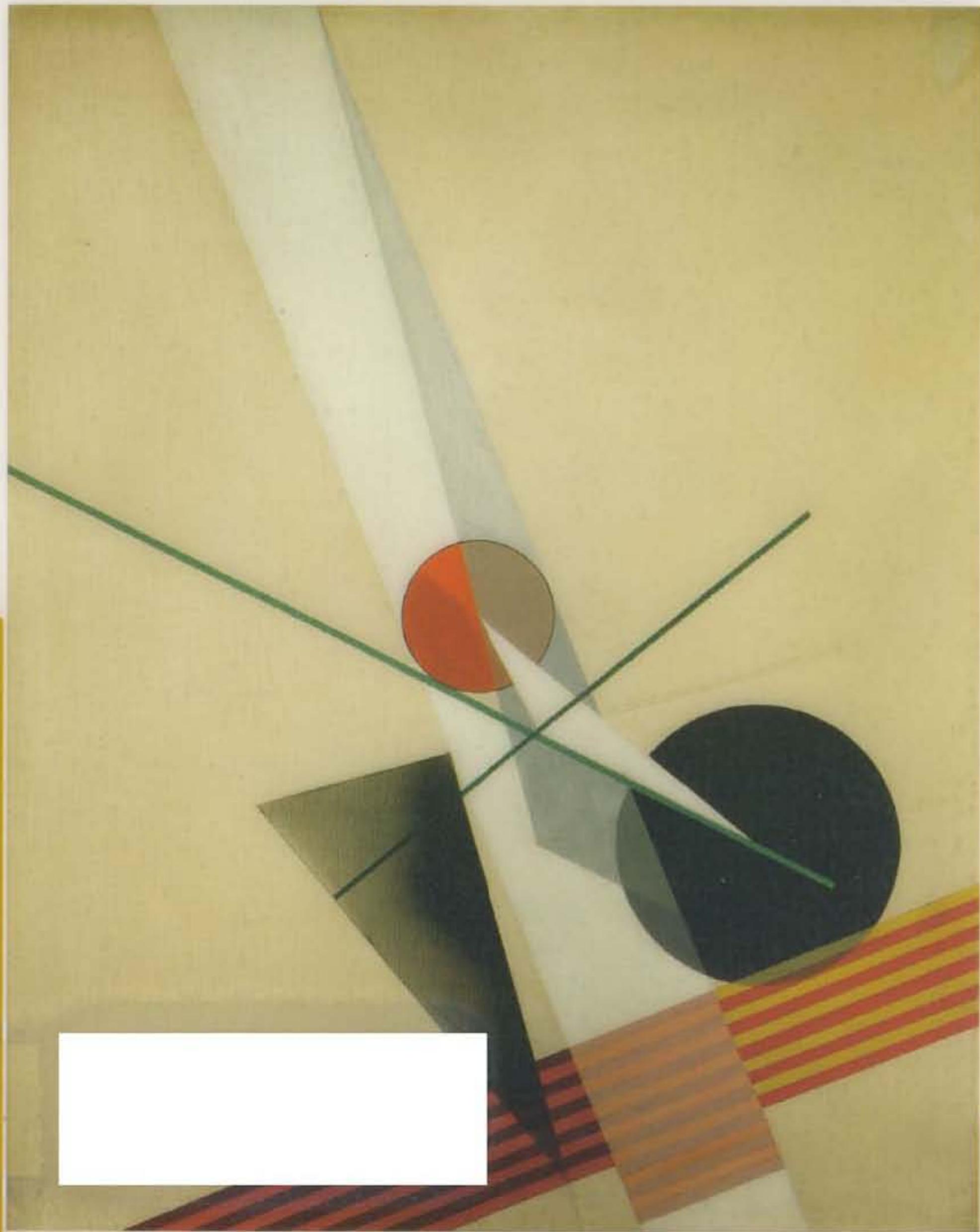
*Covers both C89  
and C99*

K. N. KING

# C PROGRAMMING

*A Modern Approach*

SECOND EDITION



K.N.KING

# C PROGRAMMING

*A Modern Approach*

SECOND EDITION

*A clear, complete, and engaging presentation of the C programming language—now with coverage of both C89 and C99*

The first edition of *C Programming: A Modern Approach* was a hit with students and faculty alike because of its clarity and comprehensiveness as well as its trademark Q&A sections. King's spiral approach made the first edition accessible to a broad range of readers, from beginners to more advanced students. The first edition was used at over 225 colleges, making it one of the leading C textbooks of the last ten years.

#### FEATURES OF THE SECOND EDITION

- Complete coverage of both the C89 standard and the C99 standard, with all C99 changes clearly marked
- Includes a quick reference to all C89 and C99 library functions
- Expanded coverage of GCC
- New coverage of abstract data types
- Updated to reflect today's CPUs and operating systems
- Nearly 500 exercises and programming projects—60% more than in the first edition
- Source code and solutions to selected exercises and programming projects for students, available at the author's website ([knking.com](http://knking.com))
- Password-protected instructor site (also at [knking.com](http://knking.com)) containing solutions to the remaining exercises and projects, plus PowerPoint presentations for most chapters

"I thoroughly enjoyed reading the second edition of *C Programming* and I look forward to using it in future courses."

— Karen Reid, Senior Lecturer, Department of Computer Science, University of Toronto

"The second edition of King's *C Programming* improves on an already impressive base, and is the book I recommend to anyone who wants to learn C."

— Peter Seebach, moderator, *comp.lang.c.moderated*

"I assign *C Programming* to first-year engineering students. It is concise, clear, accessible to the beginner, and yet also covers all aspects of the language."

— Professor Markus Bussmann, Department of Mechanical and Industrial Engineering, University of Toronto

K. N. KING (Ph.D., University of California, Berkeley) is an associate professor of computer science at Georgia State University. He is also the author of *Modula-2: A Complete Guide* and *Java Programming: From the Beginning*.

ISBN 978-0-393-97950-3



9 0000 >  
9 780393 979503

EAN

Cover design by Joan Greenfield

Cover art: László Moholy-Nagy. AXI. Westfälisches  
Landesmuseum, Münster, Westphalia, Germany  
Photo: Erich Lessing / Art Resource, NY  
Copyright: © ARS, NY



W. W. NORTON  
NEW YORK • LONDON

[www.wwnorton.com](http://www.wwnorton.com)

# PREFACE

*In computing, turning the obvious into the useful  
is a living definition of the word “frustration.”*

In the years since the first edition of *C Programming: A Modern Approach* was published, a host of new C-based languages have sprung up—Java and C# foremost among them—and related languages such as C++ and Perl have achieved greater prominence. Still, C remains as popular as ever, plugging away in the background, quietly powering much of the world’s software. It remains the *lingua franca* of the computer universe, as it was in 1996.

But even C must change with the times. The need for a new edition of *C Programming: A Modern Approach* became apparent when the C99 standard was published. Moreover, the first edition, with its references to DOS and 16-bit processors, was becoming dated. The second edition is fully up-to-date and has been improved in many other ways as well.

## What’s New in the Second Edition

Here’s a list of new features and improvements in the second edition:

- **Complete coverage of both the C89 standard and the C99 standard.** The biggest difference between the first and second editions is coverage of the C99 standard. My goal was to cover every significant difference between C89 and C99, including all the language features and library functions added in C99. Each C99 change is clearly marked, either with “C99” in the heading of a section or—in the case of shorter discussions—with a special icon in the left margin. I did this partly to draw attention to the changes and partly so that readers who aren’t interested in C99 or don’t have access to a C99 compiler will know what to skip. Many of the C99 additions are of interest only to a specialized audience, but some of the new features will be of use to nearly all C programmers.

C99

- ***Includes a quick reference to all C89 and C99 library functions.*** Appendix D in the first edition described all C89 standard library functions. In this edition, the appendix covers all C89 and C99 library functions.
- ***Expanded coverage of GCC.*** In the years since the first edition, use of GCC (originally the GNU C Compiler, now the GNU Compiler Collection) has spread. GCC has some significant advantages, including high quality, low (i.e., no) cost, and portability across a variety of hardware and software platforms. In recognition of its growing importance, I've included more information about GCC in this edition, including discussions of how to use it as well as common GCC error messages and warnings.
- ***New coverage of abstract data types.*** In the first edition, a significant portion of Chapter 19 was devoted to C++. This material seems less relevant today, since students may already have learned C++, Java, or C# before reading this book. In this edition, coverage of C++ has been replaced by a discussion of how to set up abstract data types in C.
- ***Expanded coverage of international features.*** Chapter 25, which is devoted to C's international features, is now much longer and more detailed. Information about the Unicode/UCS character set and its encodings is a highlight of the expanded coverage.
- ***Updated to reflect today's CPUs and operating systems.*** When I wrote the first edition, 16-bit architectures and the DOS operating system were still relevant to many readers, but such is not the case today. I've updated the discussion to focus more on 32-bit and 64-bit architectures. The rise of Linux and other versions of UNIX has dictated a stronger focus on that family of operating systems, although aspects of Windows and the Mac OS operating system that affect C programmers are mentioned as well.
- ***More exercises and programming projects.*** The first edition of this book contained 311 exercises. This edition has nearly 500 (498, to be exact), divided into two groups: exercises and programming projects.
- ***Solutions to selected exercises and programming projects.*** The most frequent request I received from readers of the first edition was to provide answers to the exercises. In response to this request, I've put the answers to roughly one-third of the exercises and programming projects on the web at [knking.com/books/c2](http://knking.com/books/c2). This feature is particularly useful for readers who aren't enrolled in a college course and need a way to check their work. Exercises and projects for which answers are provided are marked with a  icon (the "W" stands for "answer available on the Web").
- ***Password-protected instructor website.*** For this edition, I've built a new instructor resource site (accessible through [knking.com/books/c2](http://knking.com/books/c2)) containing solutions to the remaining exercises and projects, plus PowerPoint presentations for most chapters. Faculty may contact me at [cbook@knking.com](mailto:cbook@knking.com) for a password. Please use your campus email address and include a link to your department's website so that I can verify your identity.

I've also taken the opportunity to improve wording and explanations throughout the book. The changes are extensive and painstaking: every sentence has been checked and—if necessary—rewritten.

Although much has changed in this edition, I've tried to retain the original chapter and section numbering as much as possible. Only one chapter (the last one) is entirely new, but many chapters have additional sections. In a few cases, existing sections have been renumbered. One appendix (C syntax) has been dropped, but a new appendix that compares C99 with C89 has been added.

## Goals

The goals of this edition remain the same as those of the first edition:

- ***Be clear, readable, and possibly even entertaining.*** Many C books are too concise for the average reader. Others are badly written or just plain dull. I've tried to give clear, thorough explanations, leavened with enough humor to hold the reader's interest.
- ***Be accessible to a broad range of readers.*** I assume that the reader has at least a little previous programming experience, but I don't assume knowledge of a particular language. I've tried to keep jargon to a minimum and to define the terms that I use. I've also attempted to separate advanced material from more elementary topics, so that the beginner won't get discouraged.
- ***Be authoritative without being pedantic.*** To avoid arbitrarily deciding what to include and what not to include, I've tried to cover all the features of the C language and library. At the same time, I've tried to avoid burdening the reader with unnecessary detail.
- ***Be organized for easy learning.*** My experience in teaching C underscores the importance of presenting the features of C gradually. I use a spiral approach, in which difficult topics are introduced briefly, then revisited one or more times later in the book with details added each time. Pacing is deliberate, with each chapter building gradually on what has come before. For most students, this is probably the best approach: it avoids the extremes of boredom on the one hand, or “information overload” on the other.
- ***Motivate language features.*** Instead of just describing each feature of the language and giving a few simple examples of how the feature is used, I've tried to motivate each feature and discuss how it's used in practical situations.
- ***Emphasize style.*** It's important for every C programmer to develop a consistent style. Rather than dictating what this style should be, though, I usually describe a few possibilities and let the reader choose the one that's most appealing. Knowing alternative styles is a big help when reading other people's programs (which programmers often spend a great deal of time doing).
- ***Avoid dependence on a particular machine, compiler, or operating system.*** Since C is available on such a wide variety of platforms, I've tried to avoid

dependence on any particular machine, compiler, or operating system. All programs are designed to be portable to a wide variety of platforms.

- ***Use illustrations to clarify key concepts.*** I've tried to put in as many figures as I could, since I think these are crucial for understanding many aspects of C. In particular, I've tried to "animate" algorithms whenever possible by showing snapshots of data at different points in the computation.

## What's So Modern about *A Modern Approach*?

One of my most important goals has been to take a "modern approach" to C. Here are some of the ways I've tried to achieve this goal:

- ***Put C in perspective.*** Instead of treating C as the only programming language worth knowing, I treat it as one of many useful languages. I discuss what kind of applications C is best suited for; I also show how to capitalize on C's strengths while minimizing its weaknesses.
- ***Emphasize standard versions of C.*** I pay minimal attention to versions of the language prior to the C89 standard. There are just a few scattered references to K&R C (the 1978 version of the language described in the first edition of Brian Kernighan and Dennis Ritchie's book, *The C Programming Language*). Appendix C lists the major differences between C89 and K&R C.
- ***Debunk myths.*** Today's compilers are often at odds with commonly held assumptions about C. I don't hesitate to debunk some of the myths about C or challenge beliefs that have long been part of the C folklore (for example, the belief that pointer arithmetic is always faster than array subscripting). I've re-examined the old conventions of C, keeping the ones that are still helpful.
- ***Emphasize software engineering.*** I treat C as a mature software engineering tool, emphasizing how to use it to cope with issues that arise during programming-in-the-large. I stress making programs readable, maintainable, reliable, and portable, and I put special emphasis on information hiding.
- ***Postpone C's low-level features.*** These features, although handy for the kind of systems programming originally done in C, are not as relevant now that C is used for a great variety of applications. Instead of introducing them in the early chapters, as many C books do, I postpone them until Chapter 20.
- ***De-emphasize "manual optimization."*** Many books teach the reader to write tricky code in order to gain small savings in program efficiency. With today's abundance of optimizing C compilers, these techniques are often no longer necessary; in fact, they can result in programs that are less efficient.

## Q&A Sections

Each chapter ends with a "Q&A section"—a series of questions and answers related to material covered in the chapter. Topics addressed in these sections include:

- **Frequently asked questions.** I've tried to answer questions that come up frequently in my own courses, in other books, and on newsgroups related to C.
- **Additional discussion and clarification of tricky issues.** Although readers with experience in a variety of languages may be satisfied with a brief explanation and a couple of examples, readers with less experience need more.
- **Side issues that don't belong in the main flow.** Some questions raise technical issues that won't be of interest to all readers.
- **Material too advanced or too esoteric to interest the average reader.** Questions of this nature are marked with an asterisk (\*). Curious readers with a fair bit of programming experience may wish to delve into these questions immediately; others should definitely skip them on a first reading. *Warning:* These questions often refer to topics covered in later chapters.
- **Common differences among C compilers.** I discuss some frequently used (but nonstandard) features provided by particular compilers.

**Q&A**

Some questions in Q&A sections relate directly to specific places in the chapter; these places are marked by a special icon to signal the reader that additional information is available.

## Other Features

In addition to Q&A sections, I've included a number of useful features, many of which are marked with simple but distinctive icons (shown at left).



- **Warnings** alert readers to common pitfalls. C is famous for its traps; documenting them all is a hopeless—if not impossible—task. I've tried to pick out the pitfalls that are most common and/or most important.

- **Cross-references** provide a hypertext-like ability to locate information. Although many of these are pointers to topics covered later in the book, some point to previous topics that the reader may wish to review.

**idiom**

- **Idioms**—code patterns frequently seen in C programs—are marked for quick reference.

**portability tip**

- **Portability tips** give hints for writing programs that are independent of a particular machine, compiler, or operating system.

- **Sidebars** cover topics that aren't strictly part of C but that every knowledgeable C programmer should be aware of. (See “Source Code” on the next page for an example of a sidebar.)

- **Appendices** provide valuable reference information.

## Programs

Choosing illustrative programs isn't an easy job. If programs are too brief and artificial, readers won't get any sense of how the features are used in the real world. On the other hand, if a program is *too* realistic, its point can easily be lost in a forest of

details. I've chosen a middle course, using small, simple examples to make concepts clear when they're first introduced, then gradually building up to complete programs. I haven't included programs of great length; it's been my experience that instructors don't have the time to cover them and students don't have the patience to read them. I don't ignore the issues that arise in the creation of large programs, though—Chapter 15 (Writing Large Programs) and Chapter 19 (Program Design) cover them in detail.

I've resisted the urge to rewrite programs to take advantage of the features of C99, since not every reader may have access to a C99 compiler or wish to use C99. I have, however, used C99's `<stdbool.h>` header in a few programs, because it conveniently defines macros named `bool`, `true`, and `false`. If your compiler doesn't support the `<stdbool.h>` header, you'll need to provide your own definitions for these names.

The programs in this edition have undergone one very minor change. The main function now has the form `int main(void) { ... }` in most cases. This change reflects recommended practice and is compatible with C99, which requires an explicit return type for each function.

---

## *Source Code*

Source code for all programs is available at [knking.com/books/c2](http://knking.com/books/c2). Updates, corrections, and news about the book can also be found at this site.

---

## **Audience**

This book is designed as a primary text for a C course at the undergraduate level. Previous programming experience in a high-level language or assembler is helpful but not necessary for a computer-literate reader (an “adept beginner,” as one of my former editors put it).

Since the book is self-contained and usable for reference as well as learning, it makes an excellent companion text for a course in data structures, compiler design, operating systems, computer graphics, embedded systems, or other courses that use C for project work. Thanks to its Q&A sections and emphasis on practical problems, the book will also appeal to readers who are enrolled in a training class or who are learning C by self-study.

## **Organization**

The book is divided into four parts:

- ***Basic Features of C.*** Chapters 1–10 cover enough of C to allow the reader to write single-file programs using arrays and functions.
- ***Advanced Features of C.*** Chapters 11–20 build on the material in the earlier chapters. The topics become a little harder in these chapters, which provide in-

depth coverage of pointers, strings, the preprocessor, structures, unions, enumerations, and low-level features of C. In addition, two chapters (15 and 19) offer guidance on program design.

- **The Standard C Library.** Chapters 21–27 focus on the C library, a large collection of functions that come with every compiler. These chapters are most likely to be used as reference material, although portions are suitable for lectures.
- **Reference.** Appendix A gives a complete list of C operators. Appendix B describes the major differences between C99 and C89, and Appendix C covers the differences between C89 and K&R C. Appendix D is an alphabetical listing of all functions in the C89 and C99 standard libraries, with a thorough description of each. Appendix E lists the ASCII character set. An annotated bibliography points the reader toward other sources of information.

A full-blown course on C should cover Chapters 1–20 in sequence, with topics from Chapters 21–27 added as needed. (Chapter 22, which includes coverage of file input/output, is the most important chapter of this group.) A shorter course can omit the following topics without losing continuity: Section 8.3 (variable-length arrays), Section 9.6 (recursion), Section 12.4 (pointers and multidimensional arrays), Section 12.5 (pointers and variable-length arrays), Section 14.5 (miscellaneous directives), Section 17.7 (pointers to functions), Section 17.8 (restricted pointers), Section 17.9 (flexible array members), Section 18.6 (inline functions), Chapter 19 (program design), Section 20.2 (bit-fields in structures), and Section 20.3 (other low-level techniques).

## Exercises and Programming Projects

Having a variety of good problems is obviously essential for a textbook. This edition of the book contains both exercises (shorter problems that don't require writing a full program) and programming projects (problems that require writing or modifying an entire program).

A few exercises have nonobvious answers (some individuals uncharitably call these "trick questions"—the nerve!). Since C programs often contain abundant examples of such code, I feel it's necessary to provide some practice. However, I'll play fair by marking these exercises with an asterisk (\*). Be careful with a starred exercise: either pay close attention and think hard or skip it entirely.

## Errors, Lack of (?)

I've taken great pains to ensure the accuracy of this book. Inevitably, however, any book of this size contains a few errors. If you spot one, please contact me at [cbook@knking.com](mailto:cbook@knking.com). I'd also appreciate hearing about which features you found especially helpful, which ones you could do without, and what you'd like to see added.

## Acknowledgments

First, I'd like to thank my editors at Norton, Fred McFarland and Aaron Javicas. Fred got the second edition underway and Aaron stepped in with brisk efficiency to bring it to completion. I'd also like to thank associate managing editor Kim Yi, copy editor Mary Kelly, production manager Roy Tedoff, and editorial assistant Carly Fraser.

I owe a huge debt to the following colleagues, who reviewed some or all of the manuscript for the second edition:

Markus Bussmann, University of Toronto

Jim Clarke, University of Toronto

Karen Reid, University of Toronto

Peter Seebach, moderator of *comp.lang.c.moderated*

Jim and Peter deserve special mention for their detailed reviews, which saved me from a number of embarrassing slips. The reviewers for the first edition, in alphabetical order, were: Susan Anderson-Freed, Manuel E. Bermudez, Lisa J. Brown, Steven C. Cater, Patrick Harrison, Brian Harvey, Henry H. Leitner, Darrell Long, Arthur B. Maccabe, Carolyn Rosner, and Patrick Terry.

I received many useful comments from readers of the first edition; I thank everyone who took the time to write. Students and colleagues at Georgia State University also provided valuable feedback. Ed Bullwinkel and his wife Nancy were kind enough to read much of the manuscript. I'm particularly grateful to my department chair, Yi Pan, who was very supportive of the project.

My wife, Susan Cole, was a pillar of strength as always. Our cats, Dennis, Pounce, and Tex, were also instrumental in the completion of the book. Pounce and Tex were happy to contribute the occasional catfight to help keep me awake while I was working late at night.

Finally, I'd like to acknowledge the late Alan J. Perlis, whose epigrams appear at the beginning of each chapter. I had the privilege of studying briefly under Alan at Yale in the mid-70s. I think he'd be amused at finding his epigrams in a C book.

# BRIEF CONTENTS

## Basic Features of C

1	Introducing C	1
2	C Fundamentals	9
3	Formatted Input/Output	37
4	Expressions	53
5	Selection Statements	73
6	Loops	99
7	Basic Types	125
8	Arrays	161
9	Functions	183
10	Program Organization	219

## The Standard C Library

21	The Standard Library	529
22	Input/Output	539
23	Library Support for Numbers and Character Data	589
24	Error Handling	627
25	International Features	641
26	Miscellaneous Library Functions	677
27	Additional C99 Support for Mathematics	705

## Advanced Features of C

11	Pointers	241
12	Pointers and Arrays	257
13	Strings	277
14	The Preprocessor	315
15	Writing Large Programs	349
16	Structures, Unions, and Enumerations	377
17	Advanced Uses of Pointers	413
18	Declarations	457
19	Program Design	483
20	Low-Level Programming	509

## Reference

A	C Operators	735
B	C99 versus C89	737
C	C89 versus K&R C	743
D	Standard Library Functions	747
E	ASCII Character Set	801
	Bibliography	803
	Index	807

# CONTENTS

<b>Preface</b>	<b>xxi</b>
<b>1 INTRODUCING C</b>	<b>1</b>
<b>1.1 History of C</b>	<b>1</b>
Origins	1
Standardization	2
C-Based Languages	3
<b>1.2 Strengths and Weaknesses of C</b>	<b>4</b>
Strengths	4
Weaknesses	5
Effective Use of C	6
<b>2 C FUNDAMENTALS</b>	<b>9</b>
<b>2.1 Writing a Simple Program</b>	<b>9</b>
Program: Printing a Pun	9
Compiling and Linking	10
Integrated Development Environments	11
<b>2.2 The General Form of a Simple Program</b>	<b>12</b>
Directives	12
Functions	13
Statements	14
Printing Strings	14
<b>2.3 Comments</b>	<b>15</b>
<b>2.4 Variables and Assignment</b>	<b>17</b>
Types	17
Declarations	17
Assignment	18

Printing the Value of a Variable	19
Program: Computing the Dimensional Weight of a Box	20
Initialization	21
Printing Expressions	22
<b>2.5 Reading Input</b>	<b>22</b>
Program: Computing the Dimensional Weight of a Box (Revisited)	22
<b>2.6 Defining Names for Constants</b>	<b>23</b>
Program: Converting from Fahrenheit to Celsius	24
<b>2.7 Identifiers</b>	<b>25</b>
Keywords	26
<b>2.8 Layout of a C Program</b>	<b>27</b>
<b>3 FORMATTED INPUT/OUTPUT</b>	<b>37</b>
<b>3.1 The printf Function</b>	<b>37</b>
Conversion Specifications	38
Program: Using printf to Format Numbers	40
Escape Sequences	41
<b>3.2 The scanf Function</b>	<b>42</b>
How scanf Works	43
Ordinary Characters in Format Strings	45
Confusing printf with scanf	45
Program: Adding Fractions	46
<b>4 EXPRESSIONS</b>	<b>53</b>
<b>4.1 Arithmetic Operators</b>	<b>54</b>
Operator Precedence and Associativity	55
Program: Computing a UPC Check Digit	56
<b>4.2 Assignment Operators</b>	<b>58</b>
Simple Assignment	58
Lvalues	59
Compound Assignment	60
<b>4.3 Increment and Decrement Operators</b>	<b>61</b>
<b>4.4 Expression Evaluation</b>	<b>62</b>
Order of Subexpression Evaluation	64
<b>4.5 Expression Statements</b>	<b>65</b>
<b>5 SELECTION STATEMENTS</b>	<b>73</b>
<b>5.1 Logical Expressions</b>	<b>74</b>
Relational Operators	74
Equality Operators	75
Logical Operators	75
<b>5.2 The if Statement</b>	<b>76</b>
Compound Statements	77

The else Clause	78
Cascaded if Statements	80
Program: Calculating a Broker's Commission	81
The "Dangling else" Problem	82
Conditional Expressions	83
Boolean Values in C89	84
Boolean Values in C99	85
<b>5.3 The switch Statement</b>	<b>86</b>
The Role of the break Statement	88
Program: Printing a Date in Legal Form	89
<b>6 LOOPS</b>	<b>99</b>
<b>6.1 The while Statement</b>	<b>99</b>
Infinite Loops	101
Program: Printing a Table of Squares	102
Program: Summing a Series of Numbers	102
<b>6.2 The do Statement</b>	<b>103</b>
Program: Calculating the Number of Digits in an Integer	104
<b>6.3 The for Statement</b>	<b>105</b>
for Statement Idioms	106
Omitting Expressions in a for Statement	107
for Statements in C99	108
The Comma Operator	109
Program: Printing a Table of Squares (Revisited)	110
<b>6.4 Exiting from a Loop</b>	<b>111</b>
The break Statement	111
The continue Statement	112
The goto Statement	113
Program: Balancing a Checkbook	114
<b>6.5 The Null Statement</b>	<b>116</b>
<b>7 BASIC TYPES</b>	<b>125</b>
<b>7.1 Integer Types</b>	<b>125</b>
Integer Types in C99	128
Integer Constants	128
Integer Constants in C99	129
Integer Overflow	130
Reading and Writing Integers	130
Program: Summing a Series of Numbers (Revisited)	131
<b>7.2 Floating Types</b>	<b>132</b>
Floating Constants	133
Reading and Writing Floating-Point Numbers	134
<b>7.3 Character Types</b>	<b>134</b>
Operations on Characters	135
Signed and Unsigned Characters	136

Arithmetic Types	136
Escape Sequences	137
Character-Handling Functions	138
Reading and Writing Characters using <code>scanf</code> and <code>printf</code>	139
Reading and Writing Characters using <code>getchar</code> and <code>putchar</code>	140
Program: Determining the Length of a Message	141
<b>7.4 Type Conversion</b>	<b>142</b>
The Usual Arithmetic Conversions	143
Conversion During Assignment	145
Implicit Conversions in C99	146
Casting	147
<b>7.5 Type Definitions</b>	<b>149</b>
Advantages of Type Definitions	149
Type Definitions and Portability	150
<b>7.6 The <code>sizeof</code> Operator</b>	<b>151</b>
<b>8 ARRAYS</b>	<b>161</b>
<b>8.1 One-Dimensional Arrays</b>	<b>161</b>
Array Subscripting	162
Program: Reversing a Series of Numbers	164
Array Initialization	164
Designated Initializers	165
Program: Checking a Number for Repeated Digits	166
Using the <code>sizeof</code> Operator with Arrays	167
Program: Computing Interest	168
<b>8.2 Multidimensional Arrays</b>	<b>169</b>
Initializing a Multidimensional Array	171
Constant Arrays	172
Program: Dealing a Hand of Cards	172
<b>8.3 Variable-Length Arrays (C99)</b>	<b>174</b>
<b>9 FUNCTIONS</b>	<b>183</b>
<b>9.1 Defining and Calling Functions</b>	<b>183</b>
Program: Computing Averages	184
Program: Printing a Countdown	185
Program: Printing a Pun (Revisited)	186
Function Definitions	187
Function Calls	189
Program: Testing Whether a Number Is Prime	190
<b>9.2 Function Declarations</b>	<b>191</b>
<b>9.3 Arguments</b>	<b>193</b>
Argument Conversions	194
Array Arguments	195
Variable-Length Array Parameters	198

	Using static in Array Parameter Declarations	200
	Compound Literals	200
9.4	The return Statement	201
9.5	Program Termination	202
	The exit Function	203
9.6	Recursion	204
	The Quicksort Algorithm	205
	Program: Quicksort	207
<b>10</b>	<b>PROGRAM ORGANIZATION</b>	<b>219</b>
10.1	Local Variables	219
	Static Local Variables	220
	Parameters	221
10.2	External Variables	221
	Example: Using External Variables to Implement a Stack	221
	Pros and Cons of External Variables	222
	Program: Guessing a Number	224
10.3	Blocks	227
10.4	Scope	228
10.5	Organizing a C Program	229
	Program: Classifying a Poker Hand	230
<b>11</b>	<b>POINTERS</b>	<b>241</b>
11.1	Pointer Variables	241
	Declaring Pointer Variables	242
11.2	The Address and Indirection Operators	243
	The Address Operator	243
	The Indirection Operator	244
11.3	Pointer Assignment	245
11.4	Pointers as Arguments	247
	Program: Finding the Largest and Smallest Elements in an Array	249
	Using const to Protect Arguments	250
11.5	Pointers as Return Values	251
<b>12</b>	<b>POINTERS AND ARRAYS</b>	<b>257</b>
12.1	Pointer Arithmetic	257
	Adding an Integer to a Pointer	258
	Subtracting an Integer from a Pointer	259
	Subtracting One Pointer from Another	259
	Comparing Pointers	260
	Pointers to Compound Literals	260
12.2	Using Pointers for Array Processing	260
	Combining the * and ++ Operators	262

<b>12.3</b>	<b>Using an Array Name as a Pointer</b>	<b>263</b>
	Program: Reversing a Series of Numbers (Revisited)	264
	Array Arguments (Revisited)	265
	Using a Pointer as an Array Name	266
<b>12.4</b>	<b>Pointers and Multidimensional Arrays</b>	<b>267</b>
	Processing the Elements of a Multidimensional Array	267
	Processing the Rows of a Multidimensional Array	268
	Processing the Columns of a Multidimensional Array	269
	Using the Name of a Multidimensional Array as a Pointer	269
<b>12.5</b>	<b>Pointers and Variable-Length Arrays (C99)</b>	<b>270</b>
<b>13</b>	<b>STRINGS</b>	<b>277</b>
<b>13.1</b>	<b>String Literals</b>	<b>277</b>
	Escape Sequences in String Literals	278
	Continuing a String Literal	278
	How String Literals Are Stored	279
	Operations on String Literals	279
	String Literals versus Character Constants	280
<b>13.2</b>	<b>String Variables</b>	<b>281</b>
	Initializing a String Variable	281
	Character Arrays versus Character Pointers	283
<b>13.3</b>	<b>Reading and Writing Strings</b>	<b>284</b>
	Writing Strings Using <code>printf</code> and <code>puts</code>	284
	Reading Strings Using <code>scanf</code> and <code>gets</code>	285
	Reading Strings Character by Character	286
<b>13.4</b>	<b>Accessing the Characters in a String</b>	<b>287</b>
<b>13.5</b>	<b>Using the C String Library</b>	<b>289</b>
	The <code>strcpy</code> (String Copy) Function	290
	The <code>strlen</code> (String Length) Function	291
	The <code>strcat</code> (String Concatenation) Function	291
	The <code>strcmp</code> (String Comparison) Function	292
	Program: Printing a One-Month Reminder List	293
<b>13.6</b>	<b>String Idioms</b>	<b>296</b>
	Searching for the End of a String	296
	Copying a String	298
<b>13.7</b>	<b>Arrays of Strings</b>	<b>300</b>
	Command-Line Arguments	302
	Program: Checking Planet Names	303
<b>14</b>	<b>THE PREPROCESSOR</b>	<b>315</b>
<b>14.1</b>	<b>How the Preprocessor Works</b>	<b>315</b>
<b>14.2</b>	<b>Preprocessing Directives</b>	<b>318</b>
<b>14.3</b>	<b>Macro Definitions</b>	<b>319</b>
	Simple Macros	319
	Parameterized Macros	321

The # Operator	324
The ## Operator	324
General Properties of Macros	325
Parentheses in Macro Definitions	326
Creating Longer Macros	328
Predefined Macros	329
Additional Predefined Macros in C99	330
Empty Macro Arguments	331
Macros with a Variable Number of Arguments	332
The <code>_func_</code> Identifier	333
<b>14.4 Conditional Compilation</b>	<b>333</b>
The <code>#if</code> and <code>#endif</code> Directives	334
The <code>defined</code> Operator	335
The <code>#ifdef</code> and <code>#ifndef</code> Directives	335
The <code>#elif</code> and <code>#else</code> Directives	336
Uses of Conditional Compilation	337
<b>14.5 Miscellaneous Directives</b>	<b>338</b>
The <code>#error</code> Directive	338
The <code>#line</code> Directive	339
The <code>#pragma</code> Directive	340
The <code>_Pragma</code> Operator	341
<b>15 WRITING LARGE PROGRAMS</b>	<b>349</b>
<b>15.1 Source Files</b>	<b>349</b>
<b>15.2 Header Files</b>	<b>350</b>
The <code>#include</code> Directive	351
Sharing Macro Definitions and Type Definitions	353
Sharing Function Prototypes	354
Sharing Variable Declarations	355
Nested Includes	357
Protecting Header Files	357
<code>#error</code> Directives in Header Files	358
<b>15.3 Dividing a Program into Files</b>	<b>359</b>
Program: Text Formatting	359
<b>15.4 Building a Multiple-File Program</b>	<b>366</b>
Makefiles	366
Errors During Linking	368
Rebuilding a Program	369
Defining Macros Outside a Program	371
<b>16 STRUCTURES, UNIONS, AND ENUMERATIONS</b>	<b>377</b>
<b>16.1 Structure Variables</b>	<b>377</b>
Declaring Structure Variables	378
Initializing Structure Variables	379
Designated Initializers	380
Operations on Structures	381

<b>16.2 Structure Types</b>	<b>382</b>
Declaring a Structure Tag	383
Defining a Structure Type	384
Structures as Arguments and Return Values	384
Compound Literals	386
<b>16.3 Nested Arrays and Structures</b>	<b>386</b>
Nested Structures	387
Arrays of Structures	387
Initializing an Array of Structures	388
Program: Maintaining a Parts Database	389
<b>16.4 Unions</b>	<b>396</b>
Using Unions to Save Space	398
Using Unions to Build Mixed Data Structures	399
Adding a "Tag Field" to a Union	400
<b>16.5 Enumerations</b>	<b>401</b>
Enumeration Tags and Type Names	402
Enumerations as Integers	403
Using Enumerations to Declare "Tag Fields"	404
<b>17 ADVANCED USES OF POINTERS</b>	<b>413</b>
<b>17.1 Dynamic Storage Allocation</b>	<b>414</b>
Memory Allocation Functions	414
Null Pointers	414
<b>17.2 Dynamically Allocated Strings</b>	<b>416</b>
Using <code>malloc</code> to Allocate Memory for a String	416
Using Dynamic Storage Allocation in String Functions	417
Arrays of Dynamically Allocated Strings	418
Program: Printing a One-Month Reminder List (Revisited)	418
<b>17.3 Dynamically Allocated Arrays</b>	<b>420</b>
Using <code>malloc</code> to Allocate Storage for an Array	420
The <code>calloc</code> Function	421
The <code>realloc</code> Function	421
<b>17.4 Deallocating Storage</b>	<b>422</b>
The <code>free</code> Function	423
The "Dangling Pointer" Problem	424
<b>17.5 Linked Lists</b>	<b>424</b>
Declaring a Node Type	425
Creating a Node	425
The <code>-&gt;</code> Operator	426
Inserting a Node at the Beginning of a Linked List	427
Searching a Linked List	429
Deleting a Node from a Linked List	431
Ordered Lists	433
Program: Maintaining a Parts Database (Revisited)	433
<b>17.6 Pointers to Pointers</b>	<b>438</b>

<b>17.7 Pointers to Functions</b>	<b>439</b>
Function Pointers as Arguments	439
The <code>qsort</code> Function	440
Other Uses of Function Pointers	442
Program: Tabulating the Trigonometric Functions	443
<b>17.8 Restricted Pointers (C99)</b>	<b>445</b>
<b>17.9 Flexible Array Members (C99)</b>	<b>447</b>
<b>18 DECLARATIONS</b>	<b>457</b>
<b>18.1 Declaration Syntax</b>	<b>457</b>
<b>18.2 Storage Classes</b>	<b>459</b>
Properties of Variables	459
The <code>auto</code> Storage Class	460
The <code>static</code> Storage Class	461
The <code>extern</code> Storage Class	462
The <code>register</code> Storage Class	463
The Storage Class of a Function	464
Summary	465
<b>18.3 Type Qualifiers</b>	<b>466</b>
<b>18.4 Declarators</b>	<b>467</b>
Deciphering Complex Declarations	468
Using Type Definitions to Simplify Declarations	470
<b>18.5 Initializers</b>	<b>470</b>
Uninitialized Variables	472
<b>18.6 Inline Functions (C99)</b>	<b>472</b>
Inline Definitions	473
Restrictions on Inline Functions	474
Using Inline Functions with GCC	475
<b>19 PROGRAM DESIGN</b>	<b>483</b>
<b>19.1 Modules</b>	<b>484</b>
Cohesion and Coupling	486
Types of Modules	486
<b>19.2 Information Hiding</b>	<b>487</b>
A Stack Module	487
<b>19.3 Abstract Data Types</b>	<b>491</b>
Encapsulation	492
Incomplete Types	492
<b>19.4 A Stack Abstract Data Type</b>	<b>493</b>
Defining the Interface for the Stack ADT	493
Implementing the Stack ADT Using a Fixed-Length Array	495
Changing the Item Type in the Stack ADT	496
Implementing the Stack ADT Using a Dynamic Array	497
Implementing the Stack ADT Using a Linked List	499

19.5	Design Issues for Abstract Data Types	502
	Naming Conventions	502
	Error Handling	502
	Generic ADTs	503
	ADTs in Newer Languages	503
20	LOW-LEVEL PROGRAMMING	509
20.1	Bitwise Operators	509
	Bitwise Shift Operators	510
	Bitwise Complement, <i>And</i> , Exclusive <i>Or</i> , and Inclusive <i>Or</i>	511
	Using the Bitwise Operators to Access Bits	512
	Using the Bitwise Operators to Access Bit-Fields	513
	Program: XOR Encryption	514
20.2	Bit-Fields in Structures	516
	How Bit-Fields Are Stored	517
20.3	Other Low-Level Techniques	518
	Defining Machine-Dependent Types	518
	Using Unions to Provide Multiple Views of Data	519
	Using Pointers as Addresses	520
	Program: Viewing Memory Locations	521
	The <i>volatile</i> Type Qualifier	523
21	THE STANDARD LIBRARY	529
21.1	Using the Library	529
	Restrictions on Names Used in the Library	530
	Functions Hidden by Macros	531
21.2	C89 Library Overview	531
21.3	C99 Library Changes	534
21.4	The <stddef.h> Header: Common Definitions	535
21.5	The <stdbool.h> Header (C99): Boolean Type and Values	536
22	INPUT/OUTPUT	539
22.1	Streams	540
	File Pointers	540
	Standard Streams and Redirection	540
	Text Files versus Binary Files	541
22.2	File Operations	543
	Opening a File	543
	Modes	544
	Closing a File	545
	Attaching a File to an Open Stream	546
	Obtaining File Names from the Command Line	546
	Program: Checking Whether a File Can Be Opened	547

Temporary Files	548
File Buffering	549
Miscellaneous File Operations	551
<b>22.3 Formatted I/O</b>	<b>551</b>
The ...printf Functions	552
...printf Conversion Specifications	552
C99 Changes to ...printf Conversion Specifications	555
Examples of ...printf Conversion Specifications	556
The ...scanf Functions	558
...scanf Format Strings	559
...scanf Conversion Specifications	560
C99 Changes to ...scanf Conversion Specifications	562
scanf Examples	563
Detecting End-of-File and Error Conditions	564
<b>22.4 Character I/O</b>	<b>566</b>
Output Functions	566
Input Functions	567
Program: Copying a File	568
<b>22.5 Line I/O</b>	<b>569</b>
Output Functions	569
Input Functions	570
<b>22.6 Block I/O</b>	<b>571</b>
<b>22.7 File Positioning</b>	<b>572</b>
Program: Modifying a File of Part Records	574
<b>22.8 String I/O</b>	<b>575</b>
Output Functions	576
Input Functions	576
<b>23 LIBRARY SUPPORT FOR NUMBERS AND CHARACTER DATA</b>	<b>589</b>
<b>23.1 The &lt;float.h&gt; Header: Characteristics of Floating Types</b>	<b>589</b>
<b>23.2 The &lt;limits.h&gt; Header: Sizes of Integer Types</b>	<b>591</b>
<b>23.3 The &lt;math.h&gt; Header (C89): Mathematics</b>	<b>593</b>
Errors	593
Trigonometric Functions	594
Hyperbolic Functions	595
Exponential and Logarithmic Functions	595
Power Functions	596
Nearest Integer, Absolute Value, and Remainder Functions	596
<b>23.4 The &lt;math.h&gt; Header (C99): Mathematics</b>	<b>597</b>
IEEE Floating-Point Standard	598
Types	599
Macros	600

Errors	600
Functions	601
Classification Macros	602
Trigonometric Functions	603
Hyperbolic Functions	603
Exponential and Logarithmic Functions	604
Power and Absolute Value Functions	605
Error and Gamma Functions	606
Nearest Integer Functions	606
Remainder Functions	608
Manipulation Functions	608
Maximum, Minimum, and Positive Difference Functions	609
Floating Multiply-Add	610
Comparison Macros	611
<b>23.5 The &lt;ctype.h&gt; Header: Character Handling</b>	<b>612</b>
Character-Classification Functions	612
Program: Testing the Character-Classification Functions	613
Character Case-Mapping Functions	614
Program: Testing the Case-Mapping Functions	614
<b>23.6 The &lt;string.h&gt; Header: String Handling</b>	<b>615</b>
Copying Functions	616
Concatenation Functions	617
Comparison Functions	617
Search Functions	619
Miscellaneous Functions	622
<b>24 ERROR HANDLING</b>	<b>627</b>
<b>24.1 The &lt;assert.h&gt; Header: Diagnostics</b>	<b>628</b>
<b>24.2 The &lt;errno.h&gt; Header: Errors</b>	<b>629</b>
The perror and strerror Functions	630
<b>24.3 The &lt;signal.h&gt; Header: Signal Handling</b>	<b>631</b>
Signal Macros	631
The signal Function	632
Predefined Signal Handlers	633
The raise Function	634
Program: Testing Signals	634
<b>24.4 The &lt;setjmp.h&gt; Header: Nonlocal Jumps</b>	<b>635</b>
Program: Testing setjmp/longjmp	636
<b>25 INTERNATIONAL FEATURES</b>	<b>641</b>
<b>25.1 The &lt;locale.h&gt; Header: Localization</b>	<b>642</b>
Categories	642
The setlocale Function	643
The localeconv Function	644
<b>25.2 Multibyte Characters and Wide Characters</b>	<b>647</b>

Multibyte Characters	648
Wide Characters	649
Unicode and the Universal Character Set	649
Encodings of Unicode	650
Multibyte/Wide-Character Conversion Functions	651
Multibyte/Wide-String Conversion Functions	653
<b>25.3 Digraphs and Trigraphs</b>	<b>654</b>
Trigraphs	654
Digraphs	655
The <code>&lt;iso646.h&gt;</code> Header: Alternative Spellings	656
<b>25.4 Universal Character Names (C99)</b>	<b>656</b>
<b>25.5 The <code>&lt;wchar.h&gt;</code> Header (C99): Extended Multibyte and Wide-Character Utilities</b>	<b>657</b>
Stream Orientation	658
Formatted Wide-Character Input/Output Functions	659
Wide-Character Input/Output Functions	661
General Wide-String Utilities	662
Wide-Character Time-Conversion Functions	667
Extended Multibyte/Wide-Character Conversion Utilities	667
<b>25.6 The <code>&lt;wctype.h&gt;</code> Header (C99): Wide-Character Classification and Mapping Utilities</b>	<b>671</b>
Wide-Character Classification Functions	671
Extensible Wide-Character Classification Functions	672
Wide-Character Case-Mapping Functions	673
Extensible Wide-Character Case-Mapping Functions	673
<b>26 MISCELLANEOUS LIBRARY FUNCTIONS</b>	<b>677</b>
<b>26.1 The <code>&lt;stdarg.h&gt;</code> Header: Variable Arguments</b>	<b>677</b>
Calling a Function with a Variable Argument List	679
The <code>v...printf</code> Functions	680
The <code>v...scanf</code> Functions	681
<b>26.2 The <code>&lt;stdlib.h&gt;</code> Header: General Utilities</b>	<b>682</b>
Numeric Conversion Functions	682
Program: Testing the Numeric Conversion Functions	684
Pseudo-Random Sequence Generation Functions	686
Program: Testing the Pseudo-Random Sequence Generation Functions	687
Communication with the Environment	687
Searching and Sorting Utilities	689
Program: Determining Air Mileage	690
Integer Arithmetic Functions	691
<b>26.3 The <code>&lt;time.h&gt;</code> Header: Date and Time</b>	<b>692</b>
Time Manipulation Functions	693
Time Conversion Functions	695
Program: Displaying the Date and Time	698

<b>27 ADDITIONAL C99 SUPPORT FOR MATHEMATICS</b>	<b>705</b>
<b>27.1 The &lt;stdint.h&gt; Header (C99): Integer Types</b>	<b>705</b>
<stdint.h> Types	706
Limits of Specified-Width Integer Types	707
Limits of Other Integer Types	708
Macros for Integer Constants	708
<b>27.2 The &lt;inttypes.h&gt; Header (C99): Format Conversion of Integer Types</b>	<b>709</b>
Macros for Format Specifiers	710
Functions for Greatest-Width Integer Types	711
<b>27.3 Complex Numbers (C99)</b>	<b>712</b>
Definition of Complex Numbers	713
Complex Arithmetic	714
Complex Types in C99	714
Operations on Complex Numbers	715
Conversion Rules for Complex Types	715
<b>27.4 The &lt;complex.h&gt; Header (C99): Complex Arithmetic</b>	<b>717</b>
<complex.h> Macros	717
The CX_LIMITED_RANGE Pragma	718
<complex.h> Functions	718
Trigonometric Functions	719
Hyperbolic Functions	720
Exponential and Logarithmic Functions	721
Power and Absolute-Value Functions	721
Manipulation Functions	722
Program: Finding the Roots of a Quadratic Equation	722
<b>27.5 The &lt;tgmath.h&gt; Header (C99): Type-Generic Math</b>	<b>723</b>
Type-Generic Macros	724
Invoking a Type-Generic Macro	725
<b>27.6 The &lt;fenv.h&gt; Header (C99): Floating-Point Environment</b>	<b>726</b>
Floating-Point Status Flags and Control Modes	727
<fenv.h> Macros	727
The FENV_ACCESS Pragma	728
Floating-Point Exception Functions	729
Rounding Functions	730
Environment Functions	730
<b>Appendix A C Operators</b>	<b>735</b>
<b>Appendix B C99 versus C89</b>	<b>737</b>
<b>Appendix C C89 versus K&amp;R C</b>	<b>743</b>
<b>Appendix D Standard Library Functions</b>	<b>747</b>
<b>Appendix E ASCII Character Set</b>	<b>801</b>
<b>Bibliography</b>	<b>803</b>
<b>Index</b>	<b>807</b>

# 1 Introducing C

*When someone says "I want a programming language in which I need only say what I wish done," give him a lollipop.\**

What is C? The simple answer—a widely used programming language developed in the early 1970s at Bell Laboratories—conveys little of C’s special flavor. Before we become immersed in the details of the language, let’s take a look at where C came from, what it was designed for, and how it has changed over the years (Section 1.1). We’ll also discuss C’s strengths and weaknesses and see how to get the most out of the language (Section 1.2).

## 1.1 History of C

Let’s take a quick look at C’s history, from its origins, to its coming of age as a standardized language, to its influence on recent languages.

### Origins

C is a by-product of the UNIX operating system, which was developed at Bell Laboratories by Ken Thompson, Dennis Ritchie, and others. Thompson single-handedly wrote the original version of UNIX, which ran on the DEC PDP-7 computer, an early minicomputer with only 8K words of main memory (this was 1969, after all!).

Like other operating systems of the time, UNIX was written in assembly language. Programs written in assembly language are usually painful to debug and hard to enhance; UNIX was no exception. Thompson decided that a higher-level

---

\*The epigrams at the beginning of each chapter are from “Epigrams on Programming” by Alan J. Perlis (*ACM SIGPLAN Notices* (September, 1982): 7–13).

language was needed for the further development of UNIX, so he designed a small language named B. Thompson based B on BCPL, a systems programming language developed in the mid-1960s. BCPL, in turn, traces its ancestry to Algol 60, one of the earliest (and most influential) programming languages.

Ritchie soon joined the UNIX project and began programming in B. In 1970, Bell Labs acquired a PDP-11 for the UNIX project. Once B was up and running on the PDP-11, Thompson rewrote a portion of UNIX in B. By 1971, it became apparent that B was not well-suited to the PDP-11, so Ritchie began to develop an extended version of B. He called his language NB (“New B”) at first, and then, as it began to diverge more from B, he changed the name to C. The language was stable enough by 1973 that UNIX could be rewritten in C. The switch to C provided an important benefit: portability. By writing C compilers for other computers at Bell Labs, the team could get UNIX running on those machines as well.

## Standardization

C continued to evolve during the 1970s, especially between 1977 and 1979. It was during this period that the first book on C appeared. *The C Programming Language*, written by Brian Kernighan and Dennis Ritchie and published in 1978, quickly became the bible of C programmers. In the absence of an official standard for C, this book—known as K&R or the “White Book” to aficionados—served as a de facto standard.

During the 1970s, there were relatively few C programmers, and most of them were UNIX users. By the 1980s, however, C had expanded beyond the narrow confines of the UNIX world. C compilers became available on a variety of machines running under different operating systems. In particular, C began to establish itself on the fast-growing IBM PC platform.

With C’s increasing popularity came problems. Programmers who wrote new C compilers relied on K&R as a reference. Unfortunately, K&R was fuzzy about some language features, so compilers often treated these features differently. Also, K&R failed to make a clear distinction between which features belonged to C and which were part of UNIX. To make matters worse, C continued to change after K&R was published, with new features being added and a few older features removed. The need for a thorough, precise, and up-to-date description of the language soon became apparent. Without such a standard, numerous dialects would have arisen, threatening the portability of C programs, one of the language’s major strengths.

The development of a U.S. standard for C began in 1983 under the auspices of the American National Standards Institute (ANSI). After many revisions, the standard was completed in 1988 and formally approved in December 1989 as ANSI standard X3.159-1989. In 1990, it was approved by the International Organization for Standardization (ISO) as international standard ISO/IEC 9899:1990. This version of the language is usually referred to as C89 or C90, to distinguish it from the

original version of C, often called K&R C. Appendix C summarizes the major differences between C89 and K&R C.

The language underwent a few changes in 1995 (described in a document known as Amendment 1). More significant changes occurred with the publication of a new standard, ISO/IEC 9899:1999, in 1999. The language described in this standard is commonly known as C99. The terms “ANSI C,” “ANSI/ISO C,” and “ISO C”—once used to describe C89—are now ambiguous, thanks to the existence of two standards.

**C99** Because C99 isn’t yet universal, and because of the need to maintain millions (if not billions) of lines of code written in older versions of C, I’ll use a special icon (shown in the left margin) to mark discussions of features that were added in C99. A compiler that doesn’t recognize these features isn’t “C99-compliant.” If history is any guide, it will be some years before all C compilers are C99-compliant, if they ever are. Appendix B lists the major differences between C99 and C89.

## C-Based Languages

C has had a huge influence on modern-day programming languages, many of which borrow heavily from it. Of the many C-based languages, several are especially prominent:

- **C++** includes all the features of C, but adds classes and other features to support object-oriented programming.
- **Java** is based on C++ and therefore inherits many C features.
- **C#** is a more recent language derived from C++ and Java.
- **Perl** was originally a fairly simple scripting language; over time it has grown and adopted many of the features of C.

Considering the popularity of these newer languages, it’s logical to ask whether it’s worth the trouble to learn C. I think it is, for several reasons. First, learning C can give you greater insight into the features of C++, Java, C#, Perl, and the other C-based languages. Programmers who learn one of these languages first often fail to master basic features that were inherited from C. Second, there are a lot of older C programs around; you may find yourself needing to read and maintain this code. Third, C is still widely used for developing new software, especially in situations where memory or processing power is limited or where the simplicity of C is desired.

If you haven’t already used one of the newer C-based languages, you’ll find that this book is excellent preparation for learning these languages. It emphasizes data abstraction, information hiding, and other principles that play a large role in object-oriented programming. C++ includes all the features of C, so you’ll be able to use everything you learn from this book if you later tackle C++. Many of the features of C can be found in the other C-based languages as well.

## 1.2 Strengths and Weaknesses of C

Like any other programming language, C has strengths and weaknesses. Both stem from the language's original use (writing operating systems and other systems software) and its underlying philosophy:

- **C is a low-level language.** To serve as a suitable language for systems programming, C provides access to machine-level concepts (bytes and addresses, for example) that other programming languages try to hide. C also provides operations that correspond closely to a computer's built-in instructions, so that programs can be fast. Since application programs rely on it for input/output, storage management, and numerous other services, an operating system can't afford to be slow.
- **C is a small language.** C provides a more limited set of features than many languages. (The reference manual in the second edition of K&R covers the entire language in 49 pages.) To keep the number of features small, C relies heavily on a "library" of standard functions. (A "function" is similar to what other programming languages might call a "procedure," "subroutine," or "method.")
- **C is a permissive language.** C assumes that you know what you're doing, so it allows you a wider degree of latitude than many languages. Moreover, C doesn't mandate the detailed error-checking found in other languages.

### Strengths

C's strengths help explain why the language has become so popular:

- **Efficiency.** Efficiency has been one of C's advantages from the beginning. Because C was intended for applications where assembly language had traditionally been used, it was crucial that C programs could run quickly and in limited amounts of memory.
- **Portability.** Although program portability wasn't a primary goal of C, it has turned out to be one of the language's strengths. When a program must run on computers ranging from PCs to supercomputers, it is often written in C. One reason for the portability of C programs is that—thanks to C's early association with UNIX and the later ANSI/ISO standards—the language hasn't splintered into incompatible dialects. Another is that C compilers are small and easily written, which has helped make them widely available. Finally, C itself has features that support portability (although there's nothing to prevent programmers from writing nonportable programs).
- **Power.** C's large collection of data types and operators help make it a powerful language. In C, it's often possible to accomplish quite a bit with just a few lines of code.

- **Flexibility.** Although C was originally designed for systems programming, it has no inherent restrictions that limit it to this arena. C is now used for applications of all kinds, from embedded systems to commercial data processing. Moreover, C imposes very few restrictions on the use of its features; operations that would be illegal in other languages are often permitted in C. For example, C allows a character to be added to an integer value (or, for that matter, a floating-point number). This flexibility can make programming easier, although it may allow some bugs to slip through.
- **Standard library.** One of C's great strengths is its standard library, which contains hundreds of functions for input/output, string handling, storage allocation, and other useful operations.
- **Integration with UNIX.** C is particularly powerful in combination with UNIX (including the popular variant known as Linux). In fact, some UNIX tools assume that the user knows C.

## Weaknesses

C's weaknesses arise from the same source as many of its strengths: C's closeness to the machine. Here are a few of C's most notorious problems:

- **C programs can be error-prone.** C's flexibility makes it an error-prone language. Programming mistakes that would be caught in many other languages can't be detected by a C compiler. In this respect, C is a lot like assembly language, where most errors aren't detected until the program is run. To make matters worse, C contains a number of pitfalls for the unwary. In later chapters, we'll see how an extra semicolon can create an infinite loop or a missing & symbol can cause a program crash.
- **C programs can be difficult to understand.** Although C is a small language by most measures, it has a number of features that aren't found in all programming languages (and that consequently are often misunderstood). These features can be combined in a great variety of ways, many of which—although obvious to the original author of a program—can be hard for others to understand. Another problem is the terse nature of C programs. C was designed at a time when interactive communication with computers was tedious at best. As a result, C was purposefully kept terse to minimize the time required to enter and edit programs. C's flexibility can also be a negative factor; programmers who are too clever for their own good can make programs almost impossible to understand.
- **C programs can be difficult to modify.** Large programs written in C can be hard to change if they haven't been designed with maintenance in mind. Modern programming languages usually provide features such as classes and packages that support the division of a large program into more manageable pieces. C, unfortunately, lacks such features.

## Obfuscated C

Even C's most ardent admirers admit that C code can be hard to read. The annual International Obfuscated C Code Contest actually encourages contestants to write the most confusing C programs possible. The winners are truly baffling, as 1990's "Best Small Program" shows:

```
v,i,j,k,l,s,a[99];
main()
{
    for (scanf ("%d", &s); *a-s; v=a[j*=v]-a[i], k=i<s, j+=(v=j<s&&
        (!k&&!printf(2+"\\n\\n%c"-(!l<<!j), "#Q"[l^v?(l^j)&1:2])&&
        ++l||a[i]<s&&v&&v-i+j&&v+i-j))&&! (l%-s), v| (i==j?a[i+=k]=0:
        ++a[i])>=s*k&&++a[--i])
    ;
}
```

This program, written by Doron Osovanski and Baruch Nissenbaum, prints all solutions to the Eight Queens problem (the problem of placing eight queens on a chessboard in such a way that no queen attacks any other queen). In fact, it works for any number of queens between four and 99. For more winning programs, visit [www.ioccc.org](http://www.ioccc.org), the contest's web site.

---

## Effective Use of C

Using C effectively requires taking advantage of C's strengths while avoiding its weaknesses. Here are a few suggestions:

- **Learn how to avoid C pitfalls.** Hints for avoiding pitfalls are scattered throughout this book—just look for the  $\Delta$  symbol. For a more extensive list of pitfalls, see Andrew Koenig's *C Traps and Pitfalls* (Reading, Mass.: Addison-Wesley, 1989). Modern compilers will detect common pitfalls and issue warnings, but no compiler spots them all.
- **Use software tools to make programs more reliable.** C programmers are prolific tool builders (and users). One of the most famous C tools is named `lint`. `lint`, which is traditionally provided with UNIX, can subject a program to a more extensive error analysis than most C compilers. If `lint` (or a similar program) is available, it's a good idea to use it. Another useful tool is a debugger. Because of the nature of C, many bugs can't be detected by a C compiler; these show up instead in the form of run-time errors or incorrect output. Consequently, using a good debugger is practically mandatory for C programmers.
- **Take advantage of existing code libraries.** One of the benefits of using C is that so many other people also use it; it's a good bet that they've written code you can employ in your own programs. C code is often bundled into libraries (collections of functions); obtaining a suitable library is a good way to reduce errors—and save considerable programming effort. Libraries for common

tasks, including user-interface development, graphics, communications, database management, and networking, are readily available. Some libraries are in the public domain, some are open source, and some are sold commercially.

- **Adopt a sensible set of coding conventions.** A coding convention is a style rule that a programmer has decided to adopt even though it's not enforced by the language. Well-chosen conventions help make programs more uniform, easier to read, and easier to modify. Conventions are important when using any programming language, but especially so with C. As noted above, C's highly flexible nature makes it possible for programmers to write code that is all but unreadable. The programming examples in this book follow one set of conventions, but there are other, equally valid, conventions in use. (We'll discuss some of the alternatives from time to time.) Which set you use is less important than adopting *some* conventions and sticking to them.
- **Avoid "tricks" and overly complex code.** C encourages programming tricks. There are usually several ways to accomplish a given task in C; programmers are often tempted to choose the method that's most concise. Don't get carried away; the shortest solution is often the hardest to comprehend. In this book, I'll illustrate a style that's reasonably concise but still understandable.
- **Stick to the standard.** Most C compilers provide language features and library functions that aren't part of the C89 or C99 standards. For portability, it's best to avoid using nonstandard features and libraries unless they're absolutely necessary.

---

## Q & A

### Q: What is this Q&A section anyway?

A: Glad you asked. The Q&A section, which appears at the end of each chapter, serves several purposes.

The primary purpose of Q&A is to tackle questions that are frequently asked by students learning C. Readers can participate in a dialogue (more or less) with the author, much the same as if they were attending one of my C classes.

Another purpose of Q&A is to provide additional information about topics covered in the chapter. Readers of this book will likely have widely varying backgrounds. Some will be experienced in other programming languages, whereas others will be learning to program for the first time. Readers with experience in a variety of languages may be satisfied with a brief explanation and a couple of examples, but readers with less experience may need more. The bottom line: If you find the coverage of a topic to be sketchy, check Q&A for more details.

On occasion, Q&A will discuss common differences among C compilers. For example, we'll cover some frequently used (but nonstandard) features that are provided by particular compilers.

**Q: What does `lint` do? [p. 6]**

A: `lint` checks a C program for a host of potential errors, including—but not limited to—suspicious combinations of types, unused variables, unreachable code, and nonportable code. It produces a list of diagnostic messages, which the programmer must then sift through. The advantage of using `lint` is that it can detect errors that are missed by the compiler. On the other hand, you've got to remember to use `lint`; it's all too easy to forget about it. Worse still, `lint` can produce messages by the hundreds, of which only a fraction refer to actual errors.

**Q: Where did `lint` get its name?**

A: Unlike the names of many other UNIX tools, `lint` isn't an acronym; it got its name from the way it picks up pieces of "fluff" from a program.

**Q: How do I get a copy of `lint`?**

A: `lint` is a standard UNIX utility; if you rely on another operating system, then you probably don't have `lint`. Fortunately, versions of `lint` are available from third parties. An enhanced version of `lint` known as `splint` (Secure Programming Lint) is included in many Linux distributions and can be downloaded for free from [www.splint.org](http://www.splint.org).

**Q: Is there some way to force a compiler to do a more thorough job of error-checking, without having to use `lint`?**

A: Yes. Most compilers will do a more thorough check of a program if asked to. In addition to checking for errors (undisputed violations of the rules of C), most compilers also produce warning messages, indicating potential trouble spots. Some compilers have more than one "warning level"; selecting a higher level causes the compiler to check for more problems than choosing a lower level. If your compiler supports warning levels, it's a good idea to select the highest level, causing the compiler to perform the most thorough job of checking that it's capable of. Error-checking options for the GCC compiler, which is distributed with Linux, are discussed in the Q&A section at the end of Chapter 2.

GCC &gt; 2.1

**\*Q: I'm interested in making my program as reliable as possible. Are there any other tools available besides `lint` and debuggers?**

A: Yes. Other common tools include "bounds-checkers" and "leak-finders." C doesn't require that array subscripts be checked; a bounds-checker adds this capability. A leak-finder helps locate "memory leaks": blocks of memory that are dynamically allocated but never deallocated.

---

\*Starred questions cover material too advanced or too esoteric to interest the average reader, and often refer to topics covered in later chapters. Curious readers with a fair bit of programming experience may wish to delve into these questions immediately; others should definitely skip them on a first reading.

# 2 C Fundamentals

*One man's constant is another man's variable.*

This chapter introduces several basic concepts, including preprocessing directives, functions, variables, and statements, that we'll need in order to write even the simplest programs. Later chapters will cover these topics in much greater detail.

To start off, Section 2.1 presents a small C program and describes how to compile and link it. Section 2.2 then discusses how to generalize the program, and Section 2.3 shows how to add explanatory remarks, known as comments. Section 2.4 introduces variables, which store data that may change during the execution of a program, and Section 2.5 shows how to use the `scanf` function to read data into variables. Constants—data that won't change during program execution—can be given names, as Section 2.6 shows. Finally, Section 2.7 explains C's rules for creating names (identifiers) and Section 2.8 gives the rules for laying out a program.

## 2.1 Writing a Simple Program

In contrast to programs written in some languages, C programs require little “boilerplate”—a complete program can be as short as a few lines.

### PROGRAM Printing a Pun

The first program in Kernighan and Ritchie's classic *The C Programming Language* is extremely short; it does nothing but write the message `hello, world`. Unlike other C authors, I won't use this program as my first example. I will, however, uphold another C tradition: the bad pun. Here's the pun:

To C, or not to C: that is the question.

The following program, which we'll name `pun.c`, displays this message each time it is run.

```
pun.c #include <stdio.h>

int main(void)
{
    printf("To C, or not to C: that is the question.\n");
    return 0;
}
```

Section 2.2 explains the form of this program in some detail. For now, I'll just make a few brief observations. The line

```
#include <stdio.h>
```

is necessary to “include” information about C’s standard I/O (input/output) library. The program’s executable code goes inside `main`, which represents the “main” program. The only line inside `main` is a command to display the desired message. `printf` is a function from the standard I/O library that can produce nicely formatted output. The `\n` code tells `printf` to advance to the next line after printing the message. The line

```
return 0;
```

indicates that the program “returns” the value 0 to the operating system when it terminates.

## Compiling and Linking

Despite its brevity, getting `pun.c` to run is more involved than you might expect. First, we need to create a file named `pun.c` containing the program (any text editor will do). The name of the file doesn’t matter, but the `.c` extension is often required by compilers.

Next, we’ve got to convert the program to a form that the machine can execute. For a C program, that usually involves three steps:

- **Preprocessing.** The program is first given to a *preprocessor*, which obeys commands that begin with # (known as *directives*). A preprocessor is a bit like an editor; it can add things to the program and make modifications.
- **Compiling.** The modified program now goes to a *compiler*, which translates it into machine instructions (*object code*). The program isn’t quite ready to run yet, however.
- **Linking.** In the final step, a *linker* combines the object code produced by the compiler with any additional code needed to yield a complete executable program. This additional code includes library functions (like `printf`) that are used in the program.

Fortunately, this process is often automated, so you won't find it too onerous. In fact, the preprocessor is usually integrated with the compiler, so you probably won't even notice it at work.

The commands necessary to compile and link vary, depending on the compiler and operating system. Under UNIX, the C compiler is usually named `cc`. To compile and link the `pun.c` program, enter the following command in a terminal or command-line window:

```
% cc pun.c
```

(The `%` character is the UNIX prompt, not something that you need to enter.) Linking is automatic when using `cc`; no separate link command is necessary.

After compiling and linking the program, `cc` leaves the executable program in a file named `a.out` by default. `cc` has many options; one of them (the `-o` option) allows us to choose the name of the file containing the executable program. For example, if we want the executable version of `pun.c` to be named `pun`, we would enter the following command:

```
% cc -o pun pun.c
```

---

### The GCC Compiler

One of the most popular C compilers is the GCC compiler, which is supplied with Linux but is available for many other platforms as well. Using this compiler is similar to using the traditional UNIX `cc` compiler. For example, to compile the `pun.c` program, we would use the following command:

```
% gcc -o pun pun.c
```

**Q&A**

The Q&A section at the end of the chapter provides more information about GCC.

---

## Integrated Development Environments

So far, we've assumed the use of a "command-line" compiler that's invoked by entering a command in a special window provided by the operating system. The alternative is to use an *integrated development environment (IDE)*, a software package that allows us to edit, compile, link, execute, and even debug a program without leaving the environment. The components of an IDE are designed to work together. For example, when the compiler detects an error in a program, it can arrange for the editor to highlight the line that contains the error. There's a great deal of variation among IDEs, so I won't discuss them further in this book. However, I would recommend checking to see which IDEs are available for your platform.

## 2.2 The General Form of a Simple Program

Let's take a closer look at `pun.c` and see how we can generalize it a bit. Simple C programs have the form

*directives*

```
int main(void)
{
    statements
}
```

In this template, and in similar templates elsewhere in this book, items printed in Courier would appear in a C program exactly as shown; items in *italics* represent text to be supplied by the programmer.

Notice how the braces show where `main` begins and ends. C uses `{` and `}` in much the same way that some other languages use words like `begin` and `end`. This illustrates a general point about C: it relies heavily on abbreviations and special symbols, one reason that C programs are concise (or—less charitably—cryptic).

**Q&A** Even the simplest C programs rely on three key language features: directives (editing commands that modify the program prior to compilation), functions (named blocks of executable code, of which `main` is an example), and statements (commands to be performed when the program is run). We'll take a closer look at these features now.

### Directives

Before a C program is compiled, it is first edited by a preprocessor. Commands intended for the preprocessor are called directives. Chapters 14 and 15 discuss directives in detail. For now, we're interested only in the `#include` directive.

The `pun.c` program begins with the line

```
#include <stdio.h>
```

This directive states that the information in `<stdio.h>` is to be “included” into the program before it is compiled. `<stdio.h>` contains information about C's standard I/O library. C has a number of *headers* like `<stdio.h>`; each contains information about some part of the standard library. The reason we're including `<stdio.h>` is that C, unlike some programming languages, has no built-in “read” and “write” commands. The ability to perform input and output is provided instead by functions in the standard library.

Directives always begin with a `#` character, which distinguishes them from other items in a C program. By default, directives are one line long; there's no semicolon or other special marker at the end of a directive.

## Functions

**Functions** are like “procedures” or “subroutines” in other programming languages—they’re the building blocks from which programs are constructed. In fact, a C program is little more than a collection of functions. Functions fall into two categories: those written by the programmer and those provided as part of the C implementation. I’ll refer to the latter as *library functions*, since they belong to a “library” of functions that are supplied with the compiler.

The term “function” comes from mathematics, where a function is a rule for computing a value when given one or more arguments:

$$\begin{aligned}f(x) &= x + 1 \\g(y, z) &= y^2 - z^2\end{aligned}$$

C uses the term “function” more loosely. In C, a function is simply a series of statements that have been grouped together and given a name. Some functions compute a value; some don’t. A function that computes a value uses the `return` statement to specify what value it “returns.” For example, a function that adds 1 to its argument might execute the statement

```
return x + 1;
```

while a function that computes the difference of the squares of its arguments might execute the statement

```
return y * y - z * z;
```

Although a C program may consist of many functions, only the main function is mandatory. `main` is special: it gets called automatically when the program is executed. Until Chapter 9, where we’ll learn how to write other functions, `main` will be the only function in our programs.



---

The name `main` is critical; it can’t be `begin` or `start` or even `MAIN`.

---

If `main` is a function, does it return a value? Yes: it returns a status code that is given to the operating system when the program terminates. Let’s take another look at the `pun.c` program:

```
#include <stdio.h>

int main(void)
{
    printf("To C, or not to C: that is the question.\n");
    return 0;
}
```

The word `int` just before `main` indicates that the `main` function returns an integer value. The word `void` in parentheses indicates that `main` has no arguments.

return value of main ➤ 9.5

The statement

```
return 0;
```

has two effects: it causes the `main` function to terminate (thus ending the program) and it indicates that the `main` function returns a value of 0. We'll have more to say about `main`'s return value in a later chapter. For now, we'll always have `main` return the value 0, which indicates normal program termination.

If there's no `return` statement at the end of the `main` function, the program will still terminate. However, many compilers will produce a warning message (because the function was supposed to return an integer but failed to).

## Statements

A ***statement*** is a command to be executed when the program runs. We'll explore statements later in the book, primarily in Chapters 5 and 6. The `pun.c` program uses only two kinds of statements. One is the `return` statement; the other is the ***function call***. Asking a function to perform its assigned task is known as ***calling*** the function. The `pun.c` program, for example, calls the `printf` function to display a string on the screen:

```
printf("To C, or not to C: that is the question.\n");
```

C requires that each statement end with a semicolon. (As with any good rule, there's one exception: the compound statement, which we'll encounter later.) The semicolon shows the compiler where the statement ends; since statements can continue over several lines, it's not always obvious where they end. Directives, on the other hand, are normally one line long, and they *don't* end with a semicolon.

## Printing Strings

`printf` is a powerful function that we'll examine in Chapter 3. So far, we've only used `printf` to display a ***string literal***—a series of characters enclosed in double quotation marks. When `printf` displays a string literal, it doesn't show the quotation marks.

`printf` doesn't automatically advance to the next output line when it finishes printing. To instruct `printf` to advance one line, we must include `\n` (the ***new-line character***) in the string to be printed. Writing a new-line character terminates the current output line; subsequent output goes onto the next line. To illustrate this point, consider the effect of replacing the statement

```
printf("To C, or not to C: that is the question.\n");
```

by two calls of `printf`:

```
printf("To C, or not to C: ");
printf("that is the question.\n");
```

The first call of `printf` writes To C, or not to C: . The second call writes that is the question. and advances to the next line. The net effect is the same as the original `printf`—the user can't tell the difference.

The new-line character can appear more than once in a string literal. To display the message

Brevity is the soul of wit.  
--Shakespeare

we could write

```
printf("Brevity is the soul of wit.\n --Shakespeare\n");
```

## 2.3 Comments

Our `pun.c` program still lacks something important: documentation. Every program should contain identifying information: the program name, the date written, the author, the purpose of the program, and so forth. In C, this information is placed in *comments*. The symbol `/*` marks the beginning of a comment and the symbol `*/` marks the end:

```
/* This is a comment */
```

Comments may appear almost anywhere in a program, either on separate lines or on the same lines as other program text. Here's what `pun.c` might look like with comments added at the beginning:

```
/* Name: pun.c */  
/* Purpose: Prints a bad pun. */  
/* Author: K. N. King */  
  
#include <stdio.h>  
  
int main(void)  
{  
    printf("To C, or not to C: that is the question.\n");  
    return 0;  
}
```

Comments may extend over more than one line; once it has seen the `/*` symbol, the compiler reads (and ignores) whatever follows until it encounters the `*/` symbol. If we like, we can combine a series of short comments into one long comment:

```
/* Name: pun.c  
   Purpose: Prints a bad pun.  
   Author: K. N. King */
```

A comment like this can be hard to read, though, because it's not easy to see where

the comment ends. Putting `*/` on a line by itself helps:

```
/* Name: pun.c
   Purpose: Prints a bad pun.
   Author: K. N. King
*/
```

Even better, we can form a “box” around the comment to make it stand out:

```
***** 
* Name: pun.c
* Purpose: Prints a bad pun.
* Author: K. N. King
***** /
```

Programmers often simplify boxed comments by omitting three of the sides:

```
/*
 * Name: pun.c
 * Purpose: Prints a bad pun.
 * Author: K. N. King
*/
```

A short comment can go on the same line with other program code:

```
int main(void) /* Beginning of main program */
```

A comment like this is sometimes called a “winged comment.”



Forgetting to terminate a comment may cause the compiler to ignore part of your program. Consider the following example:

```
printf("My ");
/* forgot to close this comment...
printf("cat ");
printf("has ");
/* so it ends here */
printf("fleas");
```

Because we’ve neglected to terminate the first comment, the compiler ignores the middle two statements, and the example prints `My fleas`.

**C99**

C99 provides a second kind of comment, which begins with `//` (two adjacent slashes):

```
// This is a comment
```

This style of comment ends automatically at the end of a line. To create a comment that’s more than one line long, we can either use the older comment style (`/* ... */`) or else put `//` at the beginning of each comment line:

```
// Name: pun.c
// Purpose: Prints a bad pun.
// Author: K. N. King
```

The newer comment style has a couple of important advantages. First, because a comment automatically ends at the end of a line, there's no chance that an unterminated comment will accidentally consume part of a program. Second, multiline comments stand out better, thanks to the `//` that's required at the beginning of each line.

## 2.4 Variables and Assignment

Few programs are as simple as the one in Section 2.1. Most programs need to perform a series of calculations before producing output, and thus need a way to store data temporarily during program execution. In C, as in most programming languages, these storage locations are called *variables*.

### Types

Every variable must have a *type*, which specifies what kind of data it will hold. C has a wide variety of types. For now, we'll limit ourselves to just two: `int` and `float`. Choosing the proper type is critical, since the type affects how the variable is stored and what operations can be performed on the variable. The type of a numeric variable determines the largest and smallest numbers that the variable can store; it also determines whether or not digits are allowed after the decimal point.

A variable of type `int` (short for *integer*) can store a whole number such as 0, 1, 392, or -2553. The range of possible values is limited, though. The largest `int` value is typically 2,147,483,647 but can be as small as 32,767.

#### Q&A

A variable of type `float` (short for *floating-point*) can store much larger numbers than an `int` variable. Furthermore, a `float` variable can store numbers with digits after the decimal point, like 379.125. `float` variables have drawbacks, however. Arithmetic on `float` numbers may be slower than arithmetic on `int` numbers. Most significantly, the value of a `float` variable is often just an approximation of the number that was stored in it. If we store 0.1 in a `float` variable, we may later find that the variable has a value such as 0.0999999999999987, thanks to rounding error.

### Declarations

Variables must be *declared*—described for the benefit of the compiler—before they can be used. To declare a variable, we first specify the *type* of the variable, then its *name*. (Variable names are chosen by the programmer, subject to the rules described in Section 2.7.) For example, we might declare variables `height` and `profit` as follows:

```
int height;
float profit;
```

The first declaration states that `height` is a variable of type `int`, meaning that `height` can store an integer value. The second declaration says that `profit` is a variable of type `float`.

If several variables have the same type, their declarations can be combined:

```
int height, length, width, volume;  
float profit, loss;
```

Notice that each complete declaration ends with a semicolon.

Our first template for `main` didn't include declarations. When `main` contains declarations, these must precede statements:

```
int main(void)  
{  
    declarations  
    statements  
}
```

blocks ▶ 10.3 As we'll see in Chapter 9, this is true of functions in general, as well as blocks (statements that contain embedded declarations). As a matter of style, it's a good idea to leave a blank line between the declarations and the statements.

C99

In C99, declarations don't have to come before statements. For example, `main` might contain a declaration, then a statement, and then another declaration. For compatibility with older compilers, the programs in this book don't take advantage of this rule. However, it's common in C++ and Java programs not to declare variables until they're first needed, so this practice can be expected to become popular in C99 programs as well.

## Assignment

A variable can be given a value by means of *assignment*. For example, the statements

```
height = 8;  
length = 12;  
width = 10;
```

assign values to `height`, `length`, and `width`. The numbers 8, 12, and 10 are said to be *constants*.

Before a variable can be assigned a value—or used in any other way, for that matter—it must first be declared. Thus, we could write

```
int height;  
height = 8;
```

but not

```
height = 8;    /*** WRONG ***/  
int height;
```

A constant assigned to a `float` variable usually contains a decimal point. For example, if `profit` is a `float` variable, we might write

```
profit = 2150.48;
```

**Q&A**

It's best to append the letter `f` (for "float") to a constant that contains a decimal point if the number is assigned to a `float` variable:

```
profit = 2150.48f;
```

Failing to include the `f` may cause a warning from the compiler.

An `int` variable is normally assigned a value of type `int`, and a `float` variable is normally assigned a value of type `float`. Mixing types (such as assigning an `int` value to a `float` variable or assigning a `float` value to an `int` variable) is possible but not always safe, as we'll see in Section 4.2.

Once a variable has been assigned a value, it can be used to help compute the value of another variable:

```
height = 8;
length = 12;
width = 10;
volume = height * length * width; /* volume is now 960 */
```

In C, `*` represents the multiplication operator, so this statement multiplies the values stored in `height`, `length`, and `width`, then assigns the result to the variable `volume`. In general, the right side of an assignment can be a formula (or *expression*, in C terminology) involving constants, variables, and operators.

## Printing the Value of a Variable

We can use `printf` to display the current value of a variable. For example, to write the message

```
Height: h
```

where `h` is the current value of the `height` variable, we'd use the following call of `printf`:

```
printf("Height: %d\n", height);
```

`%d` is a placeholder indicating where the value of `height` is to be filled in during printing. Note the placement of `\n` just after `%d`, so that `printf` will advance to the next line after printing the value of `height`.

`%d` works only for `int` variables; to print a `float` variable, we'd use `%f` instead. By default, `%f` displays a number with six digits after the decimal point. To force `%f` to display `p` digits after the decimal point, we can put `.p` between `%` and `f`. For example, to print the line

```
Profit: $2150.48
```

we'd call `printf` as follows:

```
printf("Profit: $%.2f\n", profit);
```

There's no limit to the number of variables that can be printed by a single call of `printf`. To display the values of both the `height` and `length` variables, we could use the following call of `printf`:

```
printf("Height: %d Length: %d\n", height, length);
```

## PROGRAM Computing the Dimensional Weight of a Box

Shipping companies don't especially like boxes that are large but very light, since they take up valuable space in a truck or airplane. In fact, companies often charge extra for such a box, basing the fee on its volume instead of its weight. In the United States, the usual method is to divide the volume by 166 (the allowable number of cubic inches per pound). If this number—the box's "dimensional" or "volumetric" weight—exceeds its actual weight, the shipping fee is based on the dimensional weight. (The 166 divisor is for international shipments; the dimensional weight of a domestic shipment is typically calculated using 194 instead.)

Let's say that you've been hired by a shipping company to write a program that computes the dimensional weight of a box. Since you're new to C, you decide to start off by writing a program that calculates the dimensional weight of a particular box that's  $12'' \times 10'' \times 8''$ . Division is represented by `/` in C, so the obvious way to compute the dimensional weight would be

```
weight = volume / 166;
```

where `weight` and `volume` are integer variables representing the box's weight and volume. Unfortunately, this formula isn't quite what we need. In C, when one integer is divided by another, the answer is "truncated": all digits after the decimal point are lost. The volume of a  $12'' \times 10'' \times 8''$  box will be 960 cubic inches. Dividing by 166 gives the answer 5 instead of 5.783, so we have in effect rounded *down* to the next lowest pound; the shipping company expects us to round *up*. One solution is to add 165 to the volume before dividing by 166:

```
weight = (volume + 165) / 166;
```

A volume of 166 would give a weight of  $331/166$ , or 1, while a volume of 167 would yield  $332/166$ , or 2. Calculating the weight in this fashion gives us the following program.

```
dweight.c /* Computes the dimensional weight of a 12" x 10" x 8" box */
#include <stdio.h>

int main(void)
{
```

```
int height, length, width, volume, weight;  
  
height = 8;  
length = 12;  
width = 10;  
volume = height * length * width;  
weight = (volume + 165) / 166;  
  
printf("Dimensions: %dx%dx%d\n", length, width, height);  
printf("Volume (cubic inches): %d\n", volume);  
printf("Dimensional weight (pounds): %d\n", weight);  
  
return 0;  
}
```

The output of the program is

```
Dimensions: 12x10x8  
Volume (cubic inches): 960  
Dimensional weight (pounds): 6
```

## Initialization

variable initialization ▶ 18.5

Some variables are automatically set to zero when a program begins to execute, but most are not. A variable that doesn't have a default value and hasn't yet been assigned a value by the program is said to be *uninitialized*.



Attempting to access the value of an uninitialized variable (for example, by displaying the variable using `printf` or using it in an expression) may yield an unpredictable result such as 2568, -30891, or some equally strange number. With some compilers, worse behavior—even a program crash—may occur.

We can always give a variable an initial value by using assignment, of course. But there's an easier way: put the initial value of the variable in its declaration. For example, we can declare the `height` variable and initialize it in one step:

```
int height = 8;
```

In C jargon, the value 8 is said to be an *initializer*.

Any number of variables can be initialized in the same declaration:

```
int height = 8, length = 12, width = 10;
```

Notice that each variable requires its own initializer. In the following example, the initializer 10 is good only for the variable `width`, not for `height` or `length` (which remain uninitialized):

```
int height, length, width = 10;
```

## Printing Expressions

`printf` isn't limited to displaying numbers stored in variables; it can display the value of *any* numeric expression. Taking advantage of this property can simplify a program and reduce the number of variables. For instance, the statements

```
volume = height * length * width;
printf("%d\n", volume);
```

could be replaced by

```
printf("%d\n", height * length * width);
```

`printf`'s ability to print expressions illustrates one of C's general principles: *Wherever a value is needed, any expression of the same type will do.*

## 2.5 Reading Input

Because the `dweight.c` program calculates the dimensional weight of just one box, it isn't especially useful. To improve the program, we'll need to allow the user to enter the dimensions.

To obtain input, we'll use the `scanf` function, the C library's counterpart to `printf`. The `f` in `scanf`, like the `f` in `printf`, stands for "formatted"; both `scanf` and `printf` require the use of a *format string* to specify the appearance of the input or output data. `scanf` needs to know what form the input data will take, just as `printf` needs to know how to display output data.

To read an `int` value, we'd use `scanf` as follows:

```
scanf("%d", &i); /* reads an integer; stores into i */
```

& operator ➤ 11.2

The "`%d`" string tells `scanf` to read input that represents an integer; `i` is an `int` variable into which we want `scanf` to store the input. The `&` symbol is hard to explain at this point; for now, I'll just note that it is usually (but not always) required when using `scanf`.

Reading a `float` value requires a slightly different call of `scanf`:

```
scanf("%f", &x); /* reads a float value; stores into x */
```

`%f` works only with variables of type `float`, so I'm assuming that `x` is a `float` variable. The "`%f`" string tells `scanf` to look for an input value in `float` format (the number may contain a decimal point, but doesn't have to).

### PROGRAM Computing the Dimensional Weight of a Box (Revisited)

Here's an improved version of the dimensional weight program in which the user enters the dimensions. Note that each call of `scanf` is immediately preceded by a

call of `printf`. That way, the user will know when to enter input and what input to enter.

```
dweight2.c /* Computes the dimensional weight of a
              box from input provided by the user */

#include <stdio.h>

int main(void)
{
    int height, length, width, volume, weight;

    printf("Enter height of box: ");
    scanf("%d", &height);
    printf("Enter length of box: ");
    scanf("%d", &length);
    printf("Enter width of box: ");
    scanf("%d", &width);
    volume = height * length * width;
    weight = (volume + 165) / 166;

    printf("Volume (cubic inches): %d\n", volume);
    printf("Dimensional weight (pounds): %d\n", weight);

    return 0;
}
```

The output of the program has the following appearance (input entered by the user is underlined):

```
Enter height of box: 8
Enter length of box: 12
Enter width of box: 10
Volume (cubic inches): 960
Dimensional weight (pounds): 6
```

A message that asks the user to enter input (a *prompt*) normally shouldn't end with a new-line character, because we want the user to enter input on the same line as the prompt itself. When the user presses the Enter key, the cursor automatically moves to the next line—the program doesn't need to display a new-line character to terminate the current line.

The `dweight2.c` program suffers from one problem: it doesn't work correctly if the user enters nonnumeric input. Section 3.2 discusses this issue in more detail.

## 2.6 Defining Names for Constants

When a program contains constants, it's often a good idea to give them names. The `dweight.c` and `dweight2.c` programs rely on the constant 166, whose meaning may not be at all clear to someone reading the program later. Using a feature

known as *macro definition*, we can name this constant:

```
#define INCHES_PER_POUND 166
```

`#define` is a preprocessing directive, just as `#include` is, so there's no semicolon at the end of the line.

When a program is compiled, the preprocessor replaces each macro by the value that it represents. For example, the statement

```
weight = (volume + INCHES_PER_POUND - 1) / INCHES_PER_POUND;
```

will become

```
weight = (volume + 166 - 1) / 166;
```

giving the same effect as if we'd written the latter statement in the first place.

The value of a macro can be an expression:

```
#define RECIPROCAL_OF_PI (1.0f / 3.14159f)
```

parentheses in macros ► 14.3

If it contains operators, the expression should be enclosed in parentheses.

Notice that we've used only upper-case letters in macro names. This is a convention that most C programmers follow, not a requirement of the language. (Still, C programmers have been doing this for decades; you wouldn't want to be the first to deviate.)

## PROGRAM

### Converting from Fahrenheit to Celsius

The following program prompts the user to enter a Fahrenheit temperature; it then prints the equivalent Celsius temperature. The output of the program will have the following appearance (as usual, input entered by the user is underlined):

```
Enter Fahrenheit temperature: 212
Celsius equivalent: 100.0
```

The program will allow temperatures that aren't integers; that's why the Celsius temperature is displayed as 100.0 instead of 100. Let's look first at the entire program, then see how it's put together.

```
celsius.c /* Converts a Fahrenheit temperature to Celsius */
#include <stdio.h>

#define FREEZING_PT 32.0f
#define SCALE_FACTOR (5.0f / 9.0f)

int main(void)
{
    float fahrenheit, celsius;
    printf("Enter Fahrenheit temperature: ");
```

```
    scanf("%f", &fahrenheit);

    celsius = (fahrenheit - FREEZING_PT) * SCALE_FACTOR;

    printf("Celsius equivalent: %.1f\n", celsius);

    return 0;
}
```

The statement

```
celsius = (fahrenheit - FREEZING_PT) * SCALE_FACTOR;
```

converts the Fahrenheit temperature to Celsius. Since FREEZING\_PT stands for 32.0f and SCALE\_FACTOR stands for (5.0f / 9.0f), the compiler sees this statement as

```
celsius = (fahrenheit - 32.0f) * (5.0f / 9.0f);
```

Defining SCALE\_FACTOR to be (5.0f / 9.0f) instead of (5 / 9) is important, because C truncates the result when two integers are divided. The value of (5 / 9) would be 0, which definitely isn't what we want.

The call of printf writes the Celsius temperature:

```
printf("Celsius equivalent: %.1f\n", celsius);
```

Notice the use of %.1f to display celsius with just one digit after the decimal point.

## 2.7 Identifiers

As we're writing a program, we'll have to choose names for variables, functions, macros, and other entities. These names are called *identifiers*. In C, an identifier may contain letters, digits, and underscores, but must begin with a letter or underscore. (In C99, identifiers may contain certain "universal character names" as well.)

C99

universal character names ▶ 25.4

Here are some examples of legal identifiers:

```
times10  get_next_char  _done
```

The following are *not* legal identifiers:

```
10times  get-next-char
```

The symbol 10times begins with a digit, not a letter or underscore. get-next-char contains minus signs, not underscores.

C is *case-sensitive*: it distinguishes between upper-case and lower-case letters in identifiers. For example, the following identifiers are all different:

```
job  joB  jOb  jOB  Job  JoB  JOB  JOB
```

These eight identifiers could all be used simultaneously, each for a completely different purpose. (Talk about obfuscation!) Sensible programmers try to make identifiers look different unless they're somehow related.

Since case matters in C, many programmers follow the convention of using only lower-case letters in identifiers (other than macros), with underscores inserted when necessary for legibility:

```
symbol_table current_page name_and_address
```

Other programmers avoid underscores, instead using an upper-case letter to begin each word within an identifier:

```
symbolTable currentPage nameAndAddress
```

(The first letter is sometimes capitalized as well.) Although the former style is common in traditional C, the latter style is becoming more popular thanks to its widespread use in Java and C# (and, to a lesser extent, C++). Other reasonable conventions exist; just be sure to capitalize an identifier the same way each time it appears in a program.

### Q&A

C places no limit on the maximum length of an identifier, so don't be afraid to use long, descriptive names. A name such as `current_page` is a lot easier to understand than a name like `cp`.

## Keywords

The ***keywords*** in Table 2.1 have special significance to C compilers and therefore can't be used as identifiers. Note that five keywords were added in C99.

C99

**Table 2.1**  
Keywords

auto	enum	restrict <sup>†</sup>	unsigned
break	extern	return	void
case	float	short	volatile
char	for	signed	while
const	goto	sizeof	_Bool <sup>†</sup>
continue	if	static	_Complex <sup>†</sup>
default	inline <sup>†</sup>	struct	_Imaginary <sup>†</sup>
do	int	switch	
double	long	typedef	
else	register	union	

<sup>†</sup>C99 only

Because of C's case-sensitivity, keywords must appear in programs exactly as shown in Table 2.1, with all letters in lower case. Names of functions in the standard library (such as `printf`) contain only lower-case letters also. Avoid the plight of the unfortunate programmer who enters an entire program in upper case, only to find that the compiler can't recognize keywords and calls of library functions.



restrictions on identifiers ➤ 21.1

Watch out for other restrictions on identifiers. Some compilers treat certain identifiers (`asm`, for example) as additional keywords. Identifiers that belong to the standard library are restricted as well. Accidentally using one of these names can cause an error during compilation or linking. Identifiers that begin with an underscore are also restricted.

## 2.8 Layout of a C Program

We can think of a C program as a series of *tokens*: groups of characters that can't be split up without changing their meaning. Identifiers and keywords are tokens. So are operators like `+` and `-`, punctuation marks such as the comma and semicolon, and string literals. For example, the statement

```
printf("Height: %d\n", height);
```

consists of seven tokens:

```
printf      (      "Height: %d\n"      ,      height      )      ;  
①          ②          ③          ④          ⑤          ⑥          ⑦
```

Tokens ① and ⑤ are identifiers, token ③ is a string literal, and tokens ②, ④, ⑥, and ⑦ are punctuation.

The amount of space between tokens in a program isn't critical in most cases. At one extreme, tokens can be crammed together with no space between them at all, except where this would cause two tokens to merge into a third token. For example, we could delete most of the space in the `celsius.c` program of Section 2.6, provided that we leave space between tokens such as `int` and `main` and between `float` and `Fahrenheit`:

```
/* Converts a Fahrenheit temperature to Celsius */  
#include <stdio.h>  
#define FREEZING_PT 32.0f  
#define SCALE_FACTOR (5.0f/9.0f)  
int main(void){float fahrenheit,celsius;printf(  
"Enter Fahrenheit temperature: ");scanf("%f", &fahrenheit);  
celsius=(fahrenheit-FREEZING_PT)*SCALE_FACTOR;  
printf("Celsius equivalent: %.1f\n", celsius);return 0;}
```

In fact, if the page were wider, we could put the entire `main` function on a single line. We can't put the whole *program* on one line, though, because each preprocessing directive requires a separate line.

Compressing programs in this fashion isn't a good idea. In fact, adding spaces and blank lines to a program can make it easier to read and understand. Fortunately,

C allows us to insert any amount of space—blanks, tabs, and new-line characters—between tokens. This rule has several important consequences for program layout:

- *Statements can be divided* over any number of lines. The following statement, for example, is so long that it would be hard to squeeze it onto a single line:

```
printf("Dimensional weight (pounds): %d\n",
       (volume + INCHES_PER_POUND - 1) / INCHES_PER_POUND);
```

- *Space between tokens* makes it easier for the eye to separate them. For this reason, I usually put a space before and after each operator:

```
volume = height * length * width;
```

I also put a space after each comma. Some programmers go even further, putting spaces around parentheses and other punctuation.

**Q&A**

- *Indentation* can make nesting easier to spot. For example, we should indent declarations and statements to make it clear that they're nested inside main.
- *Blank lines* can divide a program into logical units, making it easier for the reader to discern the program's structure. A program with no blank lines is as hard to read as a book with no chapters.

The `celsius.c` program of Section 2.6 illustrates several of these guidelines. Let's take a closer look at the `main` function in that program:

```
int main(void)
{
    float fahrenheit, celsius;

    printf("Enter Fahrenheit temperature: ");
    scanf("%f", &fahrenheit);

    celsius = (fahrenheit - FREEZING_PT) * SCALE_FACTOR;

    printf("Celsius equivalent: %.1f\n", celsius);

    return 0;
}
```

First, observe how the space around `=`, `-`, and `*` makes these operators stand out. Second, notice how the indentation of declarations and statements makes it obvious that they all belong to `main`. Finally, note how blank lines divide `main` into five parts: (1) declaring the `fahrenheit` and `celsius` variables; (2) obtaining the Fahrenheit temperature; (3) calculating the value of `celsius`; (4) printing the Celsius temperature; and (5) returning to the operating system.

While we're on the subject of program layout, notice how I've placed the `{` token underneath `main()` and put the matching `}` on a separate line, aligned with `{`. Putting `}` on a separate line lets us insert or delete statements at the end of the function; aligning it with `{` makes it easy to spot the end of `main`.

A final note: Although extra spaces can be added *between* tokens, it's not pos-

sible to add space *within* a token without changing the meaning of the program or causing an error. Writing

```
f1 oat fahrenheit, celsius;    /*** WRONG ***/
```

or

```
f1  
oat fahrenheit, celsius;    /*** WRONG ***/
```

produces an error when the program is compiled. Putting a space inside a string literal is allowed, although it changes the meaning of the string. However, putting a new-line character in a string (in other words, splitting the string over two lines) is illegal:

```
printf("To C, or not to C:  
that is the question.\n");    /*** WRONG ***/
```

Continuing a string from one line to the next requires a special technique that we'll learn in a later chapter.

continuing a string ▶ 13.1

## Q & A

**Q: What does GCC stand for? [p. 11]**

A: GCC originally stood for “GNU C compiler.” It now stands for “GNU Compiler Collection,” because the current version of GCC compiles programs written in a variety of languages, including Ada, C, C++, Fortran, Java, and Objective-C.

**Q: OK, so what does GNU stand for?**

A: GNU stands for “GNU's Not UNIX!” (and is pronounced *guh-NEW*, by the way). GNU is a project of the Free Software Foundation, an organization set up by Richard M. Stallman as a protest against the restrictions of licensed UNIX software. According to its web site, the Free Software Foundation believes that users should be free to “run, copy, distribute, study, change and improve” software. The GNU Project has rewritten much traditional UNIX software from scratch and made it publicly available at no charge.

GCC and other GNU software are crucial to Linux. Linux itself is only the “kernel” of an operating system (the part that handles program scheduling and basic I/O services); the GNU software is necessary to have a fully functional operating system.

For more information on the GNU Project, visit [www.gnu.org](http://www.gnu.org).

**Q: What's the big deal about GCC, anyway?**

A: GCC is significant for many reasons, not least the fact that it's free and capable of compiling a number of languages. It runs under many operating systems and generates code for many different CPUs, including all the widely used ones. GCC is

the primary compiler for many UNIX-based operating systems, including Linux, BSD, and Mac OS X, and it's used extensively for commercial software development. For more information about GCC, visit [gcc.gnu.org](http://gcc.gnu.org).

**Q: How good is GCC at finding errors in programs?**

A: GCC has various command-line options that control how thoroughly it checks programs. When these options are used, GCC is quite good at finding potential trouble spots in a program. Here are some of the more popular options:

-Wall	Causes the compiler to produce warning messages when it detects possible errors. (-W can be followed by codes for specific warnings; -Wall means “all -W options.”) Should be used in conjunction with -O for maximum effect.
-W	Issues additional warning messages beyond those produced by -Wall.
-pedantic	Issues all warnings required by the C standard. Causes programs that use nonstandard features to be rejected.
-ansi	Turns off features of GCC that aren't standard C and enables a few standard features that are normally disabled.
-std=c89	
-std=c99	Specifies which version of C the compiler should use to check the program.

These options are often used in combination:

```
% gcc -O -Wall -W -pedantic -ansi -std=c99 -o pun pun.c
```

**Q: Why is C so terse? It seems as though programs would be more readable if C used begin and end instead of { and }, integer instead of int, and so forth. [p. 12]**

A: Legend has it that the brevity of C programs is due to the environment that existed in Bell Labs at the time the language was developed. The first C compiler ran on a DEC PDP-11 (an early minicomputer); programmers used a teletype—essentially a typewriter connected to a computer—to enter programs and print listings. Because teletypes were very slow (they could print only 10 characters per second), minimizing the number of characters in a program was clearly advantageous.

**Q: In some C books, the main function ends with exit(0) instead of return 0. Are these the same? [p. 14]**

A: When they appear inside main, these statements are indeed equivalent: both terminate the program, returning the value 0 to the operating system. Which one to use is mostly a matter of taste.

**Q: What happens if a program reaches the end of the main function without executing a return statement? [p. 14]**

A: The return statement isn't mandatory; if it's missing, the program will still ter-

**C99** minate. In C89, the value returned to the operating system is undefined. In C99, if `main` is declared to return an `int` (as in our examples), the program returns 0 to the operating system; otherwise, the program returns an unspecified value.

**Q:** Does the compiler remove a comment entirely or replace it with blank space?

**A:** Some old C compilers deleted all the characters in each comment, making it possible to write

```
a/**/b = 0;
```

and have the compiler interpret it as

```
ab = 0;
```

According to the C standard, however, the compiler must replace each comment by a single space character, so this trick doesn't work. Instead, we'd end up with the following (illegal) statement:

```
a b = 0;
```

**Q:** How can I tell if my program has an unterminated comment?

**A:** If you're lucky, the program won't compile because the comment has rendered the program illegal. If the program does compile, there are several techniques that you can use. Stepping through the program line by line with a debugger will reveal if any lines are being skipped. Some IDEs display comments in a distinctive color to distinguish them from surrounding code. If you're using such an environment, you can easily spot unterminated comments, since program text will have a different color if it's accidentally included in a comment. A program such as `lint` can also help.

**Q:** Is it legal to nest one comment inside another?

**A:** Old-style comments (`/* ... */`) can't be nested. For instance, the following code is illegal:

```
/*
    /*** WRONG ***/
*/
```

The `*/` symbol on the second line matches the `/*` symbol on the first line, so the compiler will flag the `*/` symbol on the third line as an error.

**C99**

C's prohibition against nested comments can sometimes be a problem. Suppose we've written a long program containing many short comments. To disable a portion of the program temporarily (during testing, say), our first impulse is to "comment out" the offending lines with `/*` and `*/`. Unfortunately, this method won't work if the lines contain old-style comments. C99 comments (those beginning with `//`) can be nested inside old-style comments, however—another advantage to using this kind of comment.

In any event, there's a better way to disable portions of a program, as we'll see later.

**Q: Where does the `float` type get its name? [p. 17]**

A: `float` is short for "floating-point," a technique for storing numbers in which the decimal point "floats." A `float` value is usually stored in two parts: the fraction (or mantissa) and the exponent. The number 12.0 might be stored as  $1.5 \times 2^3$ , for example, where 1.5 is the fraction and 3 is the exponent. Some programming languages call this type `real` instead of `float`.

**Q: Why do floating-point constants need to end with the letter `f`? [p. 19]**

A: For the full explanation, see Chapter 7. Here's the short answer: a constant that contains a decimal point but doesn't end with `f` has type `double` (short for "double precision"). `double` values are stored more accurately than `float` values. Moreover, `double` values can be larger than `float` values, which is why we need to add the letter `f` when assigning to a `float` variable. Without the `f`, a warning may be generated about the possibility of a number being stored into a `float` variable that exceeds the capacity of the variable.

**\*Q: Is it really true that there's no limit on the length of an identifier? [p. 26]**

A: Yes and no. The C89 standard says that identifiers may be arbitrarily long. However, compilers are only required to remember the first 31 characters (63 characters in C99). Thus, if two names begin with the same 31 characters, a compiler might be unable to distinguish between them.

To make matters even more complicated, there are special rules for identifiers with external linkage; most function names fall into this category. Since these names must be made available to the linker, and since some older linkers can handle only short names, only the first six characters are significant in C89. Moreover, the case of letters may not matter. As a result, ABCDEFG and abcdefh might be treated as the same name. (In C99, the first 31 characters are significant, and the case of letters is taken into account.)

Most compilers and linkers are more generous than the standard, so these rules aren't a problem in practice. Don't worry about making identifiers too long—worry about making them too short.

**Q: How many spaces should I use for indentation? [p. 28]**

A: That's a tough question. Leave too little space, and the eye has trouble detecting indentation. Leave too much, and lines run off the screen (or page). Many C programmers indent nested statements eight spaces (one tab stop), which is probably too much. Studies have shown that the optimum amount of indentation is three spaces, but many programmers feel uncomfortable with numbers that aren't a power of two. Although I normally prefer to indent three or four spaces, I'll use two spaces in this book so that my programs will fit within the margins.

## Exercises

**Section 2.1**

1. Create and run Kernighan and Ritchie's famous "hello, world" program:

```
#include <stdio.h>

int main(void)
{
    printf("hello, world\n");
}
```

Do you get a warning message from the compiler? If so, what's needed to make it go away?

**Section 2.2**

- W 2. Consider the following program:

```
#include <stdio.h>

int main(void)
{
    printf("Parkinson's Law:\nWork expands so as to ");
    printf("fill the time\n");
    printf("available for its completion.\n");
    return 0;
}
```

- (a) Identify the directives and statements in this program.
- (b) What output does the program produce?

**Section 2.4**

- W 3. Condense the `dweight.c` program by (1) replacing the assignments to `height`, `length`, and `width` with initializers and (2) removing the `weight` variable, instead calculating `(volume + 165) / 166` within the last `printf`.
- W 4. Write a program that declares several `int` and `float` variables—without initializing them—and then prints their values. Is there any pattern to the values? (Usually there isn't.)

**Section 2.7**

- W 5. Which of the following are not legal C identifiers?
- (a) `100_bottles`
  - (b) `_100_bottles`
  - (c) `one_hundred_bottles`
  - (d) `bottles_by_the_hundred_`
6. Why is it not a good idea for an identifier to contain more than one adjacent underscore (as in `current__balance`, for example)?
7. Which of the following are keywords in C?
- (a) `for`
  - (b) `If`
  - (c) `main`
  - (d) `printf`
  - (e) `while`

W Answer available on the Web at [knking.com/books/c2](http://knking.com/books/c2).

**Section 2.8**

- W 8. How many tokens are there in the following statement?  
`answer=(3*q-p*p)/3;`
9. Insert spaces between the tokens in Exercise 8 to make the statement easier to read.
10. In the `dweight.c` program (Section 2.4), which spaces are essential?

## Programming Projects

1. Write a program that uses `printf` to display the following picture on the screen:

```
*  
*  
*  
* *  
* *  
*
```

2. Write a program that computes the volume of a sphere with a 10-meter radius, using the formula  $v = \frac{4}{3}\pi r^3$ . Write the fraction  $\frac{4}{3}$  as `4.0f/3.0f`. (Try writing it as `4/3`. What happens?) Hint: C doesn't have an exponentiation operator, so you'll need to multiply  $r$  by itself twice to compute  $r^3$ .
3. Modify the program of Programming Project 2 so that it prompts the user to enter the radius of the sphere.

- W 4. Write a program that asks the user to enter a dollars-and-cents amount, then displays the amount with 5% tax added:

```
Enter an amount: 100.00  
With tax added: $105.00
```

5. Write a program that asks the user to enter a value for  $x$  and then displays the value of the following polynomial:

$$3x^5 + 2x^4 - 5x^3 - x^2 + 7x - 6$$

Hint: C doesn't have an exponentiation operator, so you'll need to multiply  $x$  by itself repeatedly in order to compute the powers of  $x$ . (For example, `x * x * x` is  $x$  cubed.)

6. Modify the program of Programming Project 5 so that the polynomial is evaluated using the following formula:

$$(((3x + 2)x - 5)x - 1)x + 7)x - 6$$

Note that the modified program performs fewer multiplications. This technique for evaluating polynomials is known as **Horner's Rule**.

7. Write a program that asks the user to enter a U.S. dollar amount and then shows how to pay that amount using the smallest number of \$20, \$10, \$5, and \$1 bills:

```
Enter a dollar amount: 93
```

```
$20 bills: 4  
$10 bills: 1  
$5 bills: 0  
$1 bills: 3
```

*Hint:* Divide the amount by 20 to determine the number of \$20 bills needed, and then reduce the amount by the total value of the \$20 bills. Repeat for the other bill sizes. Be sure to use integer values throughout, not floating-point numbers.

8. Write a program that calculates the remaining balance on a loan after the first, second, and third monthly payments:

Enter amount of loan: 20000.00

Enter interest rate: 6.0

Enter monthly payment: 386.66

Balance remaining after first payment: \$19713.34

Balance remaining after second payment: \$19425.25

Balance remaining after third payment: \$19135.71

Display each balance with two digits after the decimal point. *Hint:* Each month, the balance is decreased by the amount of the payment, but increased by the balance times the monthly interest rate. To find the monthly interest rate, convert the interest rate entered by the user to a percentage and divide it by 12.



# 3 Formatted Input/Output

*In seeking the unattainable, simplicity only gets in the way.*

`scanf` and `printf`, which support formatted reading and writing, are two of the most frequently used functions in C. As this chapter shows, both are powerful but tricky to use properly. Section 3.1 describes `printf`, and Section 3.2 covers `scanf`. Neither section gives complete details, which will have to wait until Chapter 22.

## 3.1 The `printf` Function

The `printf` function is designed to display the contents of a string, known as the *format string*, with values possibly inserted at specified points in the string. When it's called, `printf` must be supplied with the format string, followed by any values that are to be inserted into the string during printing:

```
printf(string, expr1, expr2, ...);
```

The values displayed can be constants, variables, or more complicated expressions. There's no limit on the number of values that can be printed by a single call of `printf`.

The format string may contain both ordinary characters and *conversion specifications*, which begin with the % character. A conversion specification is a placeholder representing a value to be filled in during printing. The information that follows the % character *specifies* how the value is *converted* from its internal form (binary) to printed form (characters)—that's where the term “conversion specification” comes from. For example, the conversion specification `%d` specifies that `printf` is to convert an `int` value from binary to a string of decimal digits, while `%f` does the same for a `float` value.

Ordinary characters in a format string are printed exactly as they appear in the string; conversion specifications are replaced by the values to be printed. Consider the following example:

```
int i, j;
float x, y;

i = 10;
j = 20;
x = 43.2892f;
y = 5527.0f;

printf("i = %d, j = %d, x = %f, y = %f\n", i, j, x, y);
```

This call of `printf` produces the following output:

```
i = 10, j = 20, x = 43.289200, y = 5527.000000
```

The ordinary characters in the format string are simply copied to the output line. The four conversion specifications are replaced by the values of the variables `i`, `j`, `x`, and `y`, in that order.



C compilers aren't required to check that the number of conversion specifications in a format string matches the number of output items. The following call of `printf` has more conversion specifications than values to be printed:

```
printf("%d %d\n", i);    /*** WRONG ***/
```

`printf` will print the value of `i` correctly, then print a second (meaningless) integer value. A call with too few conversion specifications has similar problems:

```
printf("%d\n", i, j);    /*** WRONG ***/
```

In this case, `printf` prints the value of `i` but doesn't show the value of `j`.

Furthermore, compilers aren't required to check that a conversion specification is appropriate for the type of item being printed. If the programmer uses an incorrect specification, the program will simply produce meaningless output. Consider the following call of `printf`, in which the `int` variable `i` and the `float` variable `x` are in the wrong order:

```
printf("%f %d\n", i, x);    /*** WRONG ***/
```

Since `printf` must obey the format string, it will dutifully display a `float` value, followed by an `int` value. Unfortunately, both will be meaningless.

## Conversion Specifications

Conversion specifications give the programmer a great deal of control over the appearance of output. On the other hand, they can be complicated and hard to read. In fact, describing conversion specifications in complete detail is too arduous a

task to tackle this early in the book. Instead, we'll just take a brief look at some of their more important capabilities.

In Chapter 2, we saw that a conversion specification can include formatting information. In particular, we used `% . 1f` to display a `float` value with one digit after the decimal point. More generally, a conversion specification can have the form `%m . pX` or `%-m . pX`, where `m` and `p` are integer constants and `X` is a letter. Both `m` and `p` are optional; if `p` is omitted, the period that separates `m` and `p` is also dropped. In the conversion specification `%10 . 2f`, `m` is 10, `p` is 2, and `X` is `f`. In the specification `%10f`, `m` is 10 and `p` (along with the period) is missing, but in the specification `% . 2f`, `p` is 2 and `m` is missing.

The **minimum field width**, `m`, specifies the minimum number of characters to print. If the value to be printed requires fewer than `m` characters, the value is right-justified within the field. (In other words, extra spaces precede the value.) For example, the specification `%4d` would display the number 123 as `•123`. (In this chapter, I'll use `•` to represent the space character.) If the value to be printed requires more than `m` characters, the field width automatically expands to the necessary size. Thus, the specification `%4d` would display the number 12345 as 12345—no digits are lost. Putting a minus sign in front of `m` causes left justification; the specification `%-4d` would display 123 as 123•.

The meaning of the **precision**, `p`, isn't as easily described, since it depends on the choice of `X`, the **conversion specifier**. `X` indicates which conversion should be applied to the value before it's printed. The most common conversion specifiers for numbers are:

**Q&A**

- `d` — Displays an integer in decimal (base 10) form. `p` indicates the minimum number of digits to display (extra zeros are added to the beginning of the number if necessary); if `p` is omitted, it is assumed to have the value 1. (In other words, `%d` is the same as `% . 1d`.)
- `e` — Displays a floating-point number in exponential format (scientific notation). `p` indicates how many digits should appear after the decimal point (the default is 6). If `p` is 0, the decimal point is not displayed.
- `f` — Displays a floating-point number in “fixed decimal” format, without an exponent. `p` has the same meaning as for the `e` specifier.
- `g` — Displays a floating-point number in either exponential format or fixed decimal format, depending on the number's size. `p` indicates the maximum number of significant digits (*not* digits after the decimal point) to be displayed. Unlike the `f` conversion, the `g` conversion won't show trailing zeros. Furthermore, if the value to be printed has no digits after the decimal point, `g` doesn't display the decimal point.

The `g` specifier is especially useful for displaying numbers whose size can't be predicted when the program is written or that tend to vary widely in size. When used to print a moderately large or moderately small number, the `g` specifier uses fixed decimal format. But when used to print a very large or very small number, the `g` specifier switches to exponential format so that the number will require fewer characters.

specifiers for integers ► 7.1  
 specifiers for floats ► 7.2  
 specifiers for characters ► 7.3  
 specifiers for strings ► 13.3

There are many other specifiers besides %d, %e, %f, and %g. I'll gradually introduce many of them in subsequent chapters. For the full list, and for a complete explanation of the other capabilities of conversion specifications, consult Section 22.3.

## PROGRAM Using `printf` to Format Numbers

The following program illustrates the use of `printf` to print integers and floating-point numbers in various formats.

```
tprintf.c /* Prints int and float values in various formats */

#include <stdio.h>

int main(void)
{
    int i;
    float x;

    i = 40;
    x = 839.21f;

    printf("|%d|%5d|%-5d|%5.3d|\n", i, i, i, i);
    printf("|%10.3f|%10.3e|%-10g|\n", x, x, x);

    return 0;
}
```

The | characters in the `printf` format strings are there merely to help show how much space each number occupies when printed; unlike % or \, the | character has no special significance to `printf`. The output of this program is:

40	40	40	040	
	839.210	8.392e+02	839.21	

Let's take a closer look at the conversion specifications used in this program:

- %d — Displays i in decimal form, using a minimum amount of space.
- %5d — Displays i in decimal form, using a minimum of five characters. Since i requires only two characters, three spaces were added.
- %-5d — Displays i in decimal form, using a minimum of five characters; since the value of i doesn't require five characters, the spaces are added afterward (that is, i is left-justified in a field of length five).
- %5.3d — Displays i in decimal form, using a minimum of five characters overall and a minimum of three digits. Since i is only two digits long, an extra zero was added to guarantee three digits. The resulting number is only three characters long, so two spaces were added, for a total of five characters (i is right-justified).
- %10.3f — Displays x in fixed decimal form, using 10 characters overall,

with three digits after the decimal point. Since `x` requires only seven characters (three before the decimal point, three after the decimal point, and one for the decimal point itself), three spaces precede `x`.

- `%10.3e` — Displays `x` in exponential form, using 10 characters overall, with three digits after the decimal point. `x` requires nine characters altogether (including the exponent), so one space precedes `x`.
- `%-10g` — Displays `x` in either fixed decimal form or exponential form, using 10 characters overall. In this case, `printf` chose to display `x` in fixed decimal form. The presence of the minus sign forces left justification, so `x` is followed by four spaces.

## Escape Sequences

escape sequences ► 7.3

The `\n` code that we often use in format strings is called an *escape sequence*. Escape sequences enable strings to contain characters that would otherwise cause problems for the compiler, including nonprinting (control) characters and characters that have a special meaning to the compiler (such as `"`). We'll provide a complete list of escape sequences later; for now, here's a sample:

Alert (bell)	<code>\a</code>
Backspace	<code>\b</code>
New line	<code>\n</code>
Horizontal tab	<code>\t</code>

When they appear in `printf` format strings, these escape sequences represent actions to perform upon printing. Printing `\a` causes an audible beep on most machines. Printing `\b` moves the cursor back one position. Printing `\n` advances the cursor to the beginning of the next line. Printing `\t` moves the cursor to the next tab stop.

### Q&A

A string may contain any number of escape sequences. Consider the following `printf` example, in which the format string contains six escape sequences:

```
printf ("Item\tUnit\tPurchase\n\tPrice\tDate\n");
```

Executing this statement prints a two-line heading:

Item	Unit	Purchase
	Price	Date

Another common escape sequence is `\"`, which represents the `"` character. Since the `"` character marks the beginning and end of a string, it can't appear within a string without the use of this escape sequence. Here's an example:

```
printf ("\\"Hello!\\\"");
```

This statement produces the following output:

```
"Hello!"
```

Incidentally, you can't just put a single \ character in a string; the compiler will assume that it's the beginning of an escape sequence. To print a single \ character, put two \ characters in the string:

```
printf("\\\\"); /* prints one \ character */
```

## 3.2 The `scanf` Function

Just as `printf` prints output in a specified format, `scanf` reads input according to a particular format. A `scanf` format string, like a `printf` format string, may contain both ordinary characters and conversion specifications. The conversions allowed with `scanf` are essentially the same as those used with `printf`.

In many cases, a `scanf` format string will contain only conversion specifications, as in the following example:

```
int i, j;
float x, y;

scanf ("%d%d%f%f", &i, &j, &x, &y);
```

Suppose that the user enters the following input line:

```
1 -20 .3 -4.0e3
```

`scanf` will read the line, converting its characters to the numbers they represent, and then assign 1, -20, 0.3, and -4000.0 to `i`, `j`, `x`, and `y`, respectively. "Tightly packed" format strings like "%d%d%f%f" are common in `scanf` calls. `printf` format strings are less likely to have adjacent conversion specifications.

`scanf`, like `printf`, contains several traps for the unwary. When using `scanf`, the programmer must check that the number of conversion specifications matches the number of input variables and that each conversion is appropriate for the corresponding variable—as with `printf`, the compiler isn't required to check for a possible mismatch. Another trap involves the & symbol, which normally precedes each variable in a `scanf` call. The & is usually (but not always) required, and it's the programmer's responsibility to remember to use it.



Forgetting to put the & symbol in front of a variable in a call of `scanf` will have unpredictable—and possibly disastrous—results. A program crash is a common outcome. At the very least, the value that is read from the input won't be stored in the variable; instead, the variable will retain its old value (which may be meaningless if the variable wasn't given an initial value). Omitting the & is an extremely common error—be careful! Some compilers can spot this error and produce a warning message such as “*format argument is not a pointer*.” (The term *pointer* is defined in Chapter 11; the & symbol is used to create a pointer to a variable.) If you get a warning, check for a missing &.

Calling `scanf` is a powerful but unforgiving way to read data. Many professional C programmers avoid `scanf`, instead reading all data in character form and converting it to numeric form later. We'll use `scanf` quite a bit, especially in the early chapters of this book, because it provides a simple way to read numbers. Be aware, however, that many of our programs won't behave properly if the user enters unexpected input. As we'll see later, it's possible to have a program test whether `scanf` successfully read the requested data (and attempt to recover if it didn't). Such tests are impractical for the programs in this book—they would add too many statements and obscure the point of the examples.

## How `scanf` Works

`scanf` can actually do much more than I've indicated so far. It is essentially a "pattern-matching" function that tries to match up groups of input characters with conversion specifications.

Like the `printf` function, `scanf` is controlled by the format string. When it is called, `scanf` begins processing the information in the string, starting at the left. For each conversion specification in the format string, `scanf` tries to locate an item of the appropriate type in the input data, skipping blank space if necessary. `scanf` then reads the item, stopping when it encounters a character that can't possibly belong to the item. If the item was read successfully, `scanf` continues processing the rest of the format string. If any item is not read successfully, `scanf` returns immediately without looking at the rest of the format string (or the remaining input data).

As it searches for the beginning of a number, `scanf` ignores *white-space characters* (the space, horizontal and vertical tab, form-feed, and new-line characters). As a result, numbers can be put on a single line or spread out over several lines. Consider the following call of `scanf`:

```
scanf ("%d%d%f%f", &i, &j, &x, &y);
```

Suppose that the user enters three lines of input:

```
1  
-20 .3  
-4.0e3
```

`scanf` sees one continuous stream of characters:

```
••1▫-20•••.3▫•••-4.0e3▫
```

(I'm using • to represent the space character and ▫ to represent the new-line character.) Since it skips over white-space characters as it looks for the beginning of each number, `scanf` will be able to read the numbers successfully. In the following diagram, an *s* under a character indicates that it was skipped, and an *x* indicates it was read as part of an input item:

```
••1▫-20•••.3▫•••-4.0e3▫  
ssrsrrrssssrrssssrrrrrr
```

`scanf` “peeks” at the final new-line character without actually reading it. This new-line will be the first character read by the next call of `scanf`.

What rules does `scanf` follow to recognize an integer or a floating-point number? When asked to read an integer, `scanf` first searches for a digit, a plus sign, or a minus sign; it then reads digits until it reaches a nondigit. When asked to read a floating-point number, `scanf` looks for

a plus or minus sign (optional), followed by  
a series of digits (possibly containing a decimal point), followed by  
an exponent (optional). An exponent consists of the letter e (or E), an optional  
sign, and one or more digits.

The `%e`, `%f`, and `%g` conversions are interchangeable when used with `scanf`; all three follow the same rules for recognizing a floating-point number.

When `scanf` encounters a character that can't be part of the current item, the character is “put back” to be read again during the scanning of the next input item or during the next call of `scanf`. Consider the following (admittedly pathological) arrangement of our four numbers:

1-20.3-4.0e3□

Let's use the same call of `scanf` as before:

```
scanf ("%d%d%f%f", &i, &j, &x, &y);
```

Here's how `scanf` would process the new input:

- Conversion specification: `%d`. The first nonblank input character is 1; since integers can begin with 1, `scanf` then reads the next character, -. Recognizing that - can't appear inside an integer, `scanf` stores 1 into `i` and puts the - character back.
- Conversion specification: `%d`. `scanf` then reads the characters -, 2, 0, and . (period). Since an integer can't contain a decimal point, `scanf` stores -20 into `j` and puts the . character back.
- Conversion specification: `%f`. `scanf` reads the characters ., 3, and -. Since a floating-point number can't contain a minus sign after a digit, `scanf` stores 0.3 into `x` and puts the - character back.
- Conversion specification: `%f`. Lastly, `scanf` reads the characters -, 4, ., 0, e, 3, and □ (new-line). Since a floating-point number can't contain a new-line character, `scanf` stores  $-4.0 \times 10^3$  into `y` and puts the new-line character back.

In this example, `scanf` was able to match every conversion specification in the format string with an input item. Since the new-line character wasn't read, it will be left for the next call of `scanf`.

## Ordinary Characters in Format Strings

The concept of pattern-matching can be taken one step further by writing format strings that contain ordinary characters in addition to conversion specifications. The action that `scanf` takes when it processes an ordinary character in a format string depends on whether or not it's a white-space character.

- **White-space characters.** When it encounters one or more consecutive white-space characters in a format string, `scanf` repeatedly reads white-space characters from the input until it reaches a non-white-space character (which is “put back”). The number of white-space characters in the format string is irrelevant; one white-space character in the format string will match any number of white-space characters in the input. (Incidentally, putting a white-space character in a format string doesn't force the input to contain white-space characters. A white-space character in a format string matches *any* number of white-space characters in the input, including none.)
- **Other characters.** When it encounters a non-white-space character in a format string, `scanf` compares it with the next input character. If the two characters match, `scanf` discards the input character and continues processing the format string. If the characters don't match, `scanf` puts the offending character back into the input, then aborts without further processing the format string or reading characters from the input.

For example, suppose that the format string is "%d/%d". If the input is

•5/•96

`scanf` skips the first space while looking for an integer, matches %d with 5, matches / with /, skips a space while looking for another integer, and matches %d with 96. On the other hand, if the input is

•5•/•96

`scanf` skips one space, matches %d with 5, then attempts to match the / in the format string with a space in the input. There's no match, so `scanf` puts the space back; the •/•96 characters remain to be read by the next call of `scanf`. To allow spaces after the first number, we should use the format string "%d /%d" instead.

## Confusing `printf` with `scanf`

Although calls of `scanf` and `printf` may appear similar, there are significant differences between the two functions; ignoring these differences can be hazardous to the health of your program.

One common mistake is to put & in front of variables in a call of `printf`:

```
printf("%d %d\n", &i, &j);    /* *** WRONG ***/
```

Fortunately, this mistake is fairly easy to spot: `printf` will display a couple of odd-looking numbers instead of the values of `i` and `j`.

Since `scanf` normally skips white-space characters when looking for data items, there's often no need for a format string to include characters other than conversion specifications. Incorrectly assuming that `scanf` format strings should resemble `printf` format strings—another common error—may cause `scanf` to behave in unexpected ways. Let's see what happens when the following call of `scanf` is executed:

```
scanf("%d, %d", &i, &j);
```

`scanf` will first look for an integer in the input, which it stores in the variable `i`. `scanf` will then try to match a comma with the next input character. If the next input character is a space, not a comma, `scanf` will terminate without reading a value for `j`.



Although `printf` format strings often end with `\n`, putting a new-line character at the end of a `scanf` format string is usually a bad idea. To `scanf`, a new-line character in a format string is equivalent to a space; both cause `scanf` to advance to the next non-white-space character. For example, if the format string is `"%d\n"`, `scanf` will skip white space, read an integer, then skip to the next non-white-space character. A format string like this can cause an interactive program to “hang” until the user enters a nonblank character.

## PROGRAM

### Adding Fractions

To illustrate `scanf`'s ability to match patterns, consider the problem of reading a fraction entered by the user. Fractions are customarily written in the form *numerator/denominator*. Instead of having the user enter the numerator and denominator of a fraction as separate integers, `scanf` makes it possible to read the entire fraction. The following program, which adds two fractions, illustrates this technique.

```
addfrac.c /* Adds two fractions */

#include <stdio.h>

int main(void)
{
    int num1, denom1, num2, denom2, result_num, result_denom;

    printf("Enter first fraction: ");
    scanf("%d/%d", &num1, &denom1);

    printf("Enter second fraction: ");
    scanf("%d/%d", &num2, &denom2);

    result_num = num1 * denom2 + num2 * denom1;
```

```

result_denom = denom1 * denom2;
printf("The sum is %d/%d\n", result_num, result_denom);

return 0;
}

```

A session with this program might have the following appearance:

```

Enter first fraction: 5/6
Enter second fraction: 3/4
The sum is 38/24

```

Note that the resulting fraction isn't reduced to lowest terms.

## Q & A

**\*Q:** I've seen the `%i` conversion used to read and write integers. What's the difference between `%i` and `%d`? [p. 39]

A: In a `printf` format string, there's no difference between the two. In a `scanf` format string, however, `%d` can only match an integer written in decimal (base 10) form, while `%i` can match an integer expressed in octal (base 8), decimal, or hexadecimal (base 16). If an input number has a `0` prefix (as in `056`), `%i` treats it as an octal number; if it has a `0x` or `0X` prefix (as in `0x56`), `%i` treats it as a hex number. Using `%i` instead of `%d` to read a number can have surprising results if the user should accidentally put `0` at the beginning of the number. Because of this trap, I recommend sticking with `%d`.

**Q:** If `printf` treats `%` as the beginning of a conversion specification, how can I print the `%` character?

A: If `printf` encounters two consecutive `%` characters in a format string, it prints a single `%` character. For example, the statement

```
printf("Net profit: %d%%\n", profit);
```

might print

```
Net profit: 10%
```

**Q:** The `\t` escape is supposed to cause `printf` to advance to the next tab stop. How do I know how far apart tab stops are? [p. 41]

A: You don't. The effect of printing `\t` isn't defined in C; it depends on what your operating system does when asked to print a tab character. Tab stops are typically eight characters apart, but C makes no guarantee.

**Q:** What does `scanf` do if it's asked to read a number but the user enters nonnumeric input?

octal numbers ▶ 7.1

hexadecimal numbers ▶ 7.1

A: Let's look at the following example:

```
printf("Enter a number: ");
scanf("%d", &i);
```

Suppose that the user enters a valid number, followed by nonnumeric characters:

Enter a number: 23foo

In this case, `scanf` reads the 2 and the 3, storing 23 in `i`. The remaining characters (foo) are left to be read by the next call of `scanf` (or some other input function). On the other hand, suppose that the input is invalid from the beginning:

Enter a number: foo

In this case, the value of `i` is undefined and `foo` is left for the next `scanf`.

detecting errors in `scanf` ► 22.3

What can we do about this sad state of affairs? Later, we'll see how to test whether a call of `scanf` has succeeded. If the call fails, we can have the program either terminate or try to recover, perhaps by discarding the offending input and asking the user to try again. (Ways to discard bad input are discussed in the Q&A section at the end of Chapter 22.)

**Q: I don't understand how `scanf` can "put back" characters and read them again later. [p. 44]**

A: As it turns out, programs don't read user input as it is typed. Instead, input is stored in a hidden buffer, to which `scanf` has access. It's easy for `scanf` to put characters back into the buffer for subsequent reading. Chapter 22 discusses input buffering in more detail.

**Q: What does `scanf` do if the user puts punctuation marks (commas, for example) between numbers?**

A: Let's look at a simple example. Suppose that we try to read a pair of integers using `scanf`:

```
printf("Enter two numbers: ");
scanf("%d%d", &i, &j);
```

If the user enters

4,28

`scanf` will read the 4 and store it in `i`. As it searches for the beginning of the second number, `scanf` encounters the comma. Since numbers can't begin with a comma, `scanf` returns immediately. The comma and the second number are left for the next call of `scanf`.

Of course, we can easily solve the problem by adding a comma to the format string if we're sure that the numbers will *always* be separated by a comma:

```
printf("Enter two numbers, separated by a comma: ");
scanf("%d,%d", &i, &j);
```

## Exercises

### Section 3.1

1. What output do the following calls of `printf` produce?
  - (a) `printf("%6d,%4d", 86, 1040);`
  - (b) `printf("%12.5e", 30.253);`
  - (c) `printf("%.4f", 83.162);`
  - (d) `printf("%-6.2g", .0000009979);`
- W 2. Write calls of `printf` that display a `float` variable `x` in the following formats.
  - (a) Exponential notation; left-justified in a field of size 8; one digit after the decimal point.
  - (b) Exponential notation; right-justified in a field of size 10; six digits after the decimal point.
  - (c) Fixed decimal notation; left-justified in a field of size 8; three digits after the decimal point.
  - (d) Fixed decimal notation; right-justified in a field of size 6; no digits after the decimal point.

### Section 3.2

3. For each of the following pairs of `scanf` format strings, indicate whether or not the two strings are equivalent. If they're not, show how they can be distinguished.
  - (a) "%d" versus "%d"
  - (b) "%d-%d-%d" versus "%d -%d -%d"
  - (c) "%f" versus "%f "
  - (d) "%f,%f" versus "%f, %f"

- \*4. Suppose that we call `scanf` as follows:

```
scanf("%d%f%d", &i, &x, &j);
```

If the user enters

10.3 5 6

what will be the values of `i`, `x`, and `j` after the call? (Assume that `i` and `j` are `int` variables and `x` is a `float` variable.)

- W \*5. Suppose that we call `scanf` as follows:

```
scanf("%f%d%f", &x, &i, &y);
```

If the user enters

12.3 45.6 789

what will be the values of `x`, `i`, and `y` after the call? (Assume that `x` and `y` are `float` variables and `i` is an `int` variable.)

6. Show how to modify the `addfrac.c` program of Section 3.2 so that the user is allowed to enter fractions that contain spaces before and after each / character.

\*Starred exercises are tricky—the correct answer is usually not the obvious one. Read the question thoroughly, review the relevant section if necessary, and be careful!

## Programming Projects

- W 1. Write a program that accepts a date from the user in the form *mm/dd/yyyy* and then displays it in the form *yyyymmdd*:

```
Enter a date (mm/dd/yyyy) : 2/17/2011
You entered the date 20110217
```

2. Write a program that formats product information entered by the user. A session with the program should look like this:

```
Enter item number: 583
Enter unit price: 13.5
Enter purchase date (mm/dd/yyyy) : 10/24/2010
```

Item	Unit	Purchase
Price		Date
583	\$ 13.50	10/24/2010

The item number and date should be left justified; the unit price should be right justified. Allow dollar amounts up to \$9999.99. Hint: Use tabs to line up the columns.

- W 3. Books are identified by an International Standard Book Number (ISBN). ISBNs assigned after January 1, 2007 contain 13 digits, arranged in five groups, such as 978-0-393-97950-3. (Older ISBNs use 10 digits.) The first group (the *GS1 prefix*) is currently either 978 or 979. The *group identifier* specifies the language or country of origin (for example, 0 and 1 are used in English-speaking countries). The *publisher code* identifies the publisher (393 is the code for W. W. Norton). The *item number* is assigned by the publisher to identify a specific book (97950 is the code for this book). An ISBN ends with a *check digit* that's used to verify the accuracy of the preceding digits. Write a program that breaks down an ISBN entered by the user:

```
Enter ISBN: 978-0-393-97950-3
GS1 prefix: 978
Group identifier: 0
Publisher code: 393
Item number: 97950
Check digit: 3
```

*Note:* The number of digits in each group may vary; you can't assume that groups have the lengths shown in this example. Test your program with actual ISBN values (usually found on the back cover of a book and on the copyright page).

4. Write a program that prompts the user to enter a telephone number in the form (xxx) xxx-xxxx and then displays the number in the form xxx.xxx.xxx:

```
Enter phone number [(xxx) xxx-xxxx] : (404) 817-6900
You entered 404.817.6900
```

5. Write a program that asks the user to enter the numbers from 1 to 16 (in any order) and then displays the numbers in a 4 by 4 arrangement, followed by the sums of the rows, columns, and diagonals:

```
Enter the numbers from 1 to 16 in any order:
16 3 2 13 5 10 11 8 9 6 7 12 4 15 14 1
```

```
16 3 2 13  
5 10 11 8  
9 6 7 12  
4 15 14 1
```

```
Row sums: 34 34 34 34  
Column sums: 34 34 34 34  
Diagonal sums: 34 34
```

If the row, column, and diagonal sums are all the same (as they are in this example), the numbers are said to form a *magic square*. The magic square shown here appears in a 1514 engraving by artist and mathematician Albrecht Dürer. (Note that the middle numbers in the last row give the date of the engraving.)

6. Modify the addfrac.c program of Section 3.2 so that the user enters both fractions at the same time, separated by a plus sign:

```
Enter two fractions separated by a plus sign: 5/6+3/4  
The sum is 38/24
```



# 4 Expressions

*One does not learn computing by using a hand calculator, but one can forget arithmetic.*

One of C’s distinguishing characteristics is its emphasis on expressions—formulas that show how to compute a value—rather than statements. The simplest expressions are variables and constants. A variable represents a value to be computed as the program runs; a constant represents a value that doesn’t change. More complicated expressions apply operators to operands (which are themselves expressions). In the expression `a + (b * c)`, the `+` operator is applied to the operands `a` and `(b * c)`, both of which are expressions in their own right.

Operators are the basic tools for building expressions, and C has an unusually rich collection of them. To start off, C provides the rudimentary operators that are found in most programming languages:

- Arithmetic operators, including addition, subtraction, multiplication, and division.
- Relational operators to perform comparisons such as “`i` is greater than 0.”
- Logical operators to build conditions such as “`i` is greater than 0 and `i` is less than 10.”

But C doesn’t stop here; it goes on to provide dozens of other operators. There are so many operators, in fact, that we’ll need to introduce them gradually over the first twenty chapters of this book. Mastering so many operators can be a chore, but it’s essential to becoming proficient at C.

In this chapter, we’ll cover some of C’s most fundamental operators: the arithmetic operators (Section 4.1), the assignment operators (Section 4.2), and the increment and decrement operators (Section 4.3). Section 4.1 also explains operator precedence and associativity, which are important for expressions that contain more than one operator. Section 4.4 describes how C expressions are evaluated. Finally, Section 4.5 introduces the expression statement, an unusual feature that allows any expression to serve as a statement.

## 4.1 Arithmetic Operators

The *arithmetic operators*—operators that perform addition, subtraction, multiplication, and division—are the workhorses of many programming languages, including C. Table 4.1 shows C’s arithmetic operators.

**Table 4.1**  
Arithmetic Operators

	Unary	Binary	
	Additive	Multiplicative	
+ unary plus	+ addition	*	multiplication
- unary minus	- subtraction	/	division
		%	remainder

The additive and multiplicative operators are said to be *binary* because they require *two* operands. The *unary* operators require *one* operand:

```
i = +1; /* + used as a unary operator */
j = -i; /* - used as a unary operator */
```

The unary + operator does nothing; in fact, it didn’t even exist in K&R C. It’s used primarily to emphasize that a numeric constant is positive.

The binary operators probably look familiar. The only one that might not is %, the remainder operator. The value of  $i \% j$  is the remainder when  $i$  is divided by  $j$ . For example, the value of  $10 \% 3$  is 1, and the value of  $12 \% 4$  is 0.

### Q&A

The binary operators in Table 4.1—with the exception of %—allow either integer or floating-point operands, with mixing allowed. When `int` and `float` operands are mixed, the result has type `float`. Thus, `9 + 2.5f` has the value 11.5, and `6.7f / 2` has the value 3.35.

The / and % operators require special care:

- The / operator can produce surprising results. When both of its operands are integers, the / operator “truncates” the result by dropping the fractional part. Thus, the value of  $1 / 2$  is 0, not 0.5.
- The % operator requires integer operands; if either operand is not an integer, the program won’t compile.
- Using zero as the right operand of either / or % causes undefined behavior.
- Describing the result when / and % are used with negative operands is tricky. The C89 standard states that if either operand is negative, the result of a division can be rounded either up or down. (For example, the value of  $-9 / 7$  could be either  $-1$  or  $-2$ ). If  $i$  or  $j$  is negative, the sign of  $i \% j$  in C89 depends on the implementation. (For example, the value of  $-9 \% 7$  could be either  $-2$  or  $5$ ). In C99, on the other hand, the result of a division is always truncated toward zero (so  $-9 / 7$  has the value  $-1$ ) and the value of  $i \% j$  has the same sign as  $i$  (hence the value of  $-9 \% 7$  is  $-2$ ).

undefined behavior ▶ 4.4

### Q&A

C99

## Implementation-Defined Behavior

The term **implementation-defined** will arise often enough that it's worth taking a moment to discuss it. The C standard deliberately leaves parts of the language unspecified, with the understanding that an "implementation"—the software needed to compile, link, and execute programs on a particular platform—will fill in the details. As a result, the behavior of the program may vary somewhat from one implementation to another. The behavior of the / and % operators for negative operands in C89 is an example of implementation-defined behavior.

Leaving parts of the language unspecified may seem odd or even dangerous, but it reflects C's philosophy. One of the language's goals is efficiency, which often means matching the way that hardware behaves. Some CPUs yield -1 when -9 is divided by 7, while others produce -2; the C89 standard simply reflects this fact of life.

It's best to avoid writing programs that depend on implementation-defined behavior. If that's not possible, at least check the manual carefully—the C standard requires that implementation-defined behavior be documented.

---

## Operator Precedence and Associativity

When an expression contains more than one operator, its interpretation may not be immediately clear. For example, does  $i + j * k$  mean "add  $i$  and  $j$ , then multiply the result by  $k$ ," or does it mean "multiply  $j$  and  $k$ , then add  $i$ "? One solution to this problem is to add parentheses, writing either  $(i + j) * k$  or  $i + (j * k)$ . As a general rule, C allows the use of parentheses for grouping in all expressions.

What if we don't use parentheses, though? Will the compiler interpret  $i + j * k$  as  $(i + j) * k$  or  $i + (j * k)$ ? Like many other languages, C uses **operator precedence** rules to resolve this potential ambiguity. The arithmetic operators have the following relative precedence:

Highest: + - (unary)

\* / %

Lowest: + - (binary)

Operators listed on the same line (such as + and -) have equal precedence.

When two or more operators appear in the same expression, we can determine how the compiler will interpret the expression by repeatedly putting parentheses around subexpressions, starting with high-precedence operators and working down to low-precedence operators. The following examples illustrate the result:

$i + j * k$  is equivalent to  $i + (j * k)$

$-i * -j$  is equivalent to  $(-i) * (-j)$

$+i + j / k$  is equivalent to  $(+i) + (j / k)$

Operator precedence rules alone aren't enough when an expression contains two or more operators at the same level of precedence. In this situation, the **associativity**

of the operators comes into play. An operator is said to be *left associative* if it groups from left to right. The binary arithmetic operators (\*, /, %, +, and -) are all left associative, so

$$\begin{array}{ll} i - j - k & \text{is equivalent to } (i - j) - k \\ i * j / k & \text{is equivalent to } (i * j) / k \end{array}$$

An operator is *right associative* if it groups from right to left. The unary arithmetic operators (+ and -) are both right associative, so

$$- + i \quad \text{is equivalent to } - (+i)$$

Precedence and associativity rules are important in many languages, but especially so in C. However, C has so many operators (almost fifty!) that few programmers bother to memorize the precedence and associativity rules. Instead, they consult a table of operators when in doubt or just use plenty of parentheses.

table of operators ► Appendix A

## PROGRAM Computing a UPC Check Digit

For a number of years, manufacturers of goods sold in U.S. and Canadian stores have put a bar code on each product. This code, known as a Universal Product Code (UPC), identifies both the manufacturer and the product. Each bar code represents a twelve-digit number, which is usually printed underneath the bars. For example, the following bar code comes from a package of Stouffer's French Bread Pepperoni Pizza:



The digits

0 13800 15173 5

appear underneath the bar code. The first digit identifies the type of item (0 or 7 for most items, 2 for items that must be weighed, 3 for drugs and health-related merchandise, and 5 for coupons). The first group of five digits identifies the manufacturer (13800 is the code for Nestlé USA's Frozen Food Division). The second group of five digits identifies the product (including package size). The final digit is a “check digit,” whose only purpose is to help identify an error in the preceding digits. If the UPC is scanned incorrectly, the first 11 digits probably won’t be consistent with the last digit, and the store’s scanner will reject the entire code.

Here’s one method of computing the check digit:

Add the first, third, fifth, seventh, ninth, and eleventh digits.

Add the second, fourth, sixth, eighth, and tenth digits.

Multiply the first sum by 3 and add it to the second sum.

Subtract 1 from the total.

Compute the remainder when the adjusted total is divided by 10.

Subtract the remainder from 9.

Using the Stouffer's example, we get  $0 + 3 + 0 + 1 + 1 + 3 = 8$  for the first sum and  $1 + 8 + 0 + 5 + 7 = 21$  for the second sum. Multiplying the first sum by 3 and adding the second yields 45. Subtracting 1 gives 44. The remainder upon dividing by 10 is 4. When the remainder is subtracted from 9, the result is 5. Here are a couple of other UPCs, in case you want to try your hand at computing the check digit (raiding the kitchen cabinet for the answer is *not* allowed):

Jif Creamy Peanut Butter (18 oz.):	0	51500	24128	?
Ocean Spray Jellied Cranberry Sauce (8 oz.):	0	31200	01005	?

The answers appear at the bottom of the page.

Let's write a program that calculates the check digit for an arbitrary UPC. We'll ask the user to enter the first 11 digits of the UPC, then we'll display the corresponding check digit. To avoid confusion, we'll ask the user to enter the number in three parts: the single digit at the left, the first group of five digits, and the second group of five digits. Here's what a session with the program will look like:

```
Enter the first (single) digit: 0
Enter first group of five digits: 13800
Enter second group of five digits: 15173
Check digit: 5
```

Instead of reading each digit group as a *five*-digit number, we'll read it as five *one*-digit numbers. Reading the numbers as single digits is more convenient; also, we won't have to worry that one of the five-digit numbers is too large to store in an `int` variable. (Some older compilers limit the maximum value of an `int` variable to 32,767.) To read single digits, we'll use `scanf` with the `%1d` conversion specification, which matches a one-digit integer.

```
upc.c /* Computes a Universal Product Code check digit */

#include <stdio.h>

int main(void)
{
    int d, i1, i2, i3, i4, i5, j1, j2, j3, j4, j5,
        first_sum, second_sum, total;

    printf("Enter the first (single) digit: ");
    scanf("%1d", &d);
    printf("Enter first group of five digits: ");
    scanf("%1d%1d%1d%1d%1d", &i1, &i2, &i3, &i4, &i5);
    printf("Enter second group of five digits: ");
    scanf("%1d%1d%1d%1d%1d", &j1, &j2, &j3, &j4, &j5);
```

```

first_sum = d + i2 + i4 + j1 + j3 + j5;
second_sum = i1 + i3 + i5 + j2 + j4;
total = 3 * first_sum + second_sum;

printf("Check digit: %d\n", 9 - ((total - 1) % 10));

return 0;
}

```

Note that the expression `9 - ((total - 1) % 10)` could have been written as `9 - (total - 1) % 10`, but the extra set of parentheses makes it easier to understand.

## 4.2 Assignment Operators

Once the value of an expression has been computed, we'll often need to store it in a variable for later use. C's = (*simple assignment*) operator is used for that purpose. For updating a value already stored in a variable, C provides an assortment of compound assignment operators.

### Simple Assignment

The effect of the assignment  $v = e$  is to evaluate the expression  $e$  and copy its value into  $v$ . As the following examples show,  $e$  can be a constant, a variable, or a more complicated expression:

```

i = 5;           /* i is now 5 */
j = i;           /* j is now 5 */
k = 10 * i + j; /* k is now 55 */

```

If  $v$  and  $e$  don't have the same type, then the value of  $e$  is converted to the type of  $v$  as the assignment takes place:

```

int i;
float f;

i = 72.99f;    /* i is now 72 */
f = 136;        /* f is now 136.0 */

```

conversion during assignment ▶7.4

We'll return to the topic of type conversion later.

In many programming languages, assignment is a *statement*; in C, however, assignment is an *operator*, just like +. In other words, the act of assignment produces a result, just as adding two numbers produces a result. The value of an assignment  $v = e$  is the value of  $v$  *after* the assignment. Thus, the value of  $i = 72.99f$  is 72 (not 72.99).

## Side Effects

We don't normally expect operators to modify their operands, since operators in mathematics don't. Writing  $i + j$  doesn't modify either  $i$  or  $j$ ; it simply computes the result of adding  $i$  and  $j$ .

Most C operators don't modify their operands, but some do. We say that these operators have **side effects**, since they do more than just compute a value. The simple assignment operator is the first operator we've seen that has side effects; it modifies its left operand. Evaluating the expression  $i = 0$  produces the result 0 and—as a side effect—assigns 0 to  $i$ .

Since assignment is an operator, several assignments can be chained together:

```
i = j = k = 0;
```

The `=` operator is right associative, so this assignment is equivalent to

```
i = (j = (k = 0));
```

The effect is to assign 0 first to  $k$ , then to  $j$ , and finally to  $i$ .



Watch out for unexpected results in chained assignments as a result of type conversion:

```
int i;  
float f;  
  
f = i = 33.3f;
```

$i$  is assigned the value 33, then  $f$  is assigned 33.0 (not 33.3, as you might think).

In general, an assignment of the form  $v = e$  is allowed wherever a value of type  $v$  would be permitted. In the following example, the expression  $j = i$  copies  $i$  to  $j$ ; the new value of  $j$  is then added to 1, producing the new value of  $k$ :

```
i = 1;  
k = 1 + (j = i);  
printf("%d %d %d\n", i, j, k); /* prints "1 1 2" */
```

Using the assignment operator in this fashion usually isn't a good idea. For one thing, “embedded assignments” can make programs hard to read. They can also be a source of subtle bugs, as we'll see in Section 4.4.

## Lvalues

Most C operators allow their operands to be variables, constants, or expressions containing other operators. The assignment operator, however, requires an *lvalue*

as its left operand. An lvalue (pronounced “L-value”) represents an object stored in computer memory, not a constant or the result of a computation. Variables are lvalues; expressions such as 10 or 2 \* i are not. At this point, variables are the only lvalues that we know about; other kinds of lvalues will appear in later chapters.

Since the assignment operator requires an lvalue as its left operand, it’s illegal to put any other kind of expression on the left side of an assignment expression:

```
12 = i;      /*** WRONG ***/
i + j = 0;   /*** WRONG ***/
-i = j;     /*** WRONG **/
```

The compiler will detect errors of this nature, and you’ll get an error message such as “*invalid lvalue in assignment*.”

## Compound Assignment

Assignments that use the old value of a variable to compute its new value are common in C programs. The following statement, for example, adds 2 to the value stored in i:

```
i = i + 2;
```

C’s **compound assignment** operators allow us to shorten this statement and others like it. Using the += operator, we simply write:

```
i += 2; /* same as i = i + 2; */
```

The += operator adds the value of the right operand to the variable on the left.

There are nine other compound assignment operators, including the following:

```
-= *= /= %=
```

other assignment operators ➤ 20.1 (We’ll cover the remaining compound assignment operators in a later chapter.) All compound assignment operators work in much the same way:

*v* += *e* adds *v* to *e*, storing the result in *v*

*v* -= *e* subtracts *e* from *v*, storing the result in *v*

*v* \*= *e* multiplies *v* by *e*, storing the result in *v*

*v* /= *e* divides *v* by *e*, storing the result in *v*

*v* %= *e* computes the remainder when *v* is divided by *e*, storing the result in *v*

Note that I’ve been careful not to say that *v* += *e* is “equivalent” to *v* = *v* + *e*. One problem is operator precedence: *i* \*= *j* + *k* isn’t the same as *i* = *i* \* *j* + *k*. There are also rare cases in which *v* += *e* differs from *v* = *v* + *e* because *v* itself has a side effect. Similar remarks apply to the other compound assignment operators.

### Q&A



When using the compound assignment operators, be careful not to switch the two characters that make up the operator. Switching the characters may yield an expression that is acceptable to the compiler but that doesn’t have the intended meaning. For example, if you meant to write *i* += *j* but typed *i* =+ *j* instead, the

program will still compile. Unfortunately, the latter expression is equivalent to `i = (+j)`, which merely copies the value of `j` into `i`.

The compound assignment operators have the same properties as the `=` operator. In particular, they're right associative, so the statement

```
i += j += k;
```

means

```
i += (j += k);
```

## 4.3 Increment and Decrement Operators

Two of the most common operations on a variable are “incrementing” (adding 1) and “decrementing” (subtracting 1). We can, of course, accomplish these tasks by writing

```
i = i + 1;  
j = j - 1;
```

The compound assignment operators allow us to condense these statements a bit:

```
i += 1;  
j -= 1;
```

But C allows increments and decrements to be shortened even further, using the `++` (**increment**) and `--` (**decrement**) operators.

**Q&A** At first glance, the increment and decrement operators are simplicity itself: `++` adds 1 to its operand, whereas `--` subtracts 1. Unfortunately, this simplicity is misleading—the increment and decrement operators can be tricky to use. One complication is that `++` and `--` can be used as **prefix** operators (`++i` and `--i`, for example) or **postfix** operators (`i++` and `i--`). The correctness of a program may hinge on picking the proper version.

Another complication is that, like the assignment operators, `++` and `--` have side effects: they modify the values of their operands. Evaluating the expression `++i` (a “pre-increment”) yields `i + 1` and—as a side effect—increments `i`:

```
i = 1;  
printf("i is %d\n", ++i); /* prints "i is 2" */  
printf("i is %d\n", i); /* prints "i is 2" */
```

Evaluating the expression `i++` (a “post-increment”) produces the result `i`, but causes `i` to be incremented afterwards:

```
i = 1;  
printf("i is %d\n", i++); /* prints "i is 1" */  
printf("i is %d\n", i); /* prints "i is 2" */
```

The first `printf` shows the original value of `i`, before it is incremented. The second `printf` shows the new value. As these examples illustrate, `++i` means “increment `i` immediately,” while `i++` means “use the old value of `i` for now, but increment `i` later.” How much later? The C standard doesn’t specify a precise time, but it’s safe to assume that `i` will be incremented before the next statement is executed.

**Q&A**

The `--` operator has similar properties:

```
i = 1;
printf("i is %d\n", --i); /* prints "i is 0" */
printf("i is %d\n", i);   /* prints "i is 0" */

i = 1;
printf("i is %d\n", i--); /* prints "i is 1" */
printf("i is %d\n", i);   /* prints "i is 0" */
```

When `++` or `--` is used more than once in the same expression, the result can often be hard to understand. Consider the following statements:

```
i = 1;
j = 2;
k = ++i + j++;
```

What are the values of `i`, `j`, and `k` after these statements are executed? Since `i` is incremented *before* its value is used, but `j` is incremented *after* it is used, the last statement is equivalent to

```
i = i + 1;
k = i + j;
j = j + 1;
```

so the final values of `i`, `j`, and `k` are 2, 3, and 4, respectively. In contrast, executing the statements

```
i = 1;
j = 2;
k = i++ + j++;
```

will give `i`, `j`, and `k` the values 2, 3, and 3, respectively.

For the record, the postfix versions of `++` and `--` have higher precedence than unary plus and minus and are left associative. The prefix versions have the same precedence as unary plus and minus and are right associative.

## 4.4 Expression Evaluation

Table 4.2 summarizes the operators we’ve seen so far. (Appendix A has a similar table that shows *all* operators.) The first column shows the precedence of each

**Table 4.2**  
A Partial List of C Operators

Precedence	Name	Symbol(s)	Associativity
1	increment (postfix) decrement (postfix)	++ --	left
2	increment (prefix) decrement (prefix) unary plus unary minus	++ -- + -	right
3	multiplicative	* / %	left
4	additive	+ -	left
5	assignment	= *= /= %= += -=	right

operator relative to the other operators in the table (the highest precedence is 1; the lowest is 5). The last column shows the associativity of each operator.

Table 4.2 (or its larger cousin in Appendix A) has a variety of uses. Let's look at one of these. Suppose that we run across a complicated expression such as

`a = b += c++ - d + --e / -f`

as we're reading someone's program. This expression would be easier to understand if there were parentheses to show how the expression is constructed from subexpressions. With the help of Table 4.2, adding parentheses to an expression is easy: after examining the expression to find the operator with highest precedence, we put parentheses around the operator and its operands, indicating that it should be treated as a single operand from that point onwards. We then repeat the process until the expression is fully parenthesized.

In our example, the operator with highest precedence is `++`, used here as a postfix operator, so we put parentheses around `++` and its operand:

`a = b += (c++) - d + --e / -f`

We now spot a prefix `--` operator and a unary minus operator (both precedence 2) in the expression:

`a = b += (c++) - d + (--e) / (-f)`

Note that the other minus sign has an operand to its immediate left, so it must be a subtraction operator, not a unary minus operator.

Next, we notice the `/` operator (precedence 3):

`a = b += (c++) - d + ((--e) / (-f))`

The expression contains two operators with precedence 4, subtraction and addition. Whenever two operators with the same precedence are adjacent to an operand, we've got to be careful about associativity. In our example, `-` and `+` are both adjacent to `d`, so associativity rules apply. The `-` and `+` operators group from left to right, so parentheses go around the subtraction first, then the addition:

`a = b += (((c++) - d) + ((--e) / (-f)))`

The only remaining operators are `=` and `+=`. Both operators are adjacent to `b`, so we must take associativity into account. Assignment operators group from right to left, so parentheses go around the `+=` expression first, then the `=` expression:

```
(a = (b += (((c++) - d) + ((--e) / (-f)))) )
```

The expression is now fully parenthesized.

## Order of Subexpression Evaluation

The rules of operator precedence and associativity allow us to break any C expression into subexpressions—to determine uniquely where the parentheses would go if the expression were fully parenthesized. Paradoxically, these rules don't always allow us to determine the value of the expression, which may depend on the order in which its subexpressions are evaluated.

C doesn't define the order in which subexpressions are evaluated (with the exception of subexpressions involving the logical *and*, logical *or*, conditional, and comma operators). Thus, in the expression `(a + b) * (c - d)` we don't know whether `(a + b)` will be evaluated before `(c - d)`.

Most expressions have the same value regardless of the order in which their subexpressions are evaluated. However, this may not be true when a subexpression modifies one of its operands. Consider the following example:

```
a = 5;
c = (b = a + 2) - (a = 1);
```

The effect of executing the second statement is undefined; the C standard doesn't say what will happen. With most compilers, the value of `c` will be either 6 or 2. If the subexpression `(b = a + 2)` is evaluated first, `b` is assigned the value 7 and `c` is assigned 6. But if `(a = 1)` is evaluated first, `b` is assigned 3 and `c` is assigned 2.



Avoid writing expressions that access the value of a variable and also modify the variable elsewhere in the expression. The expression `(b = a + 2) - (a = 1)` accesses the value of `a` (in order to compute `a + 2`) and also modifies the value of `a` (by assigning it 1). Some compilers may produce a warning message such as “*operation on ‘a’ may be undefined*” when they encounter such an expression.

To prevent problems, it's a good idea to avoid using the assignment operators in subexpressions; instead, use a series of separate assignments. For example, the statements above could be rewritten as

```
a = 5;
b = a + 2;
a = 1;
c = b - a;
```

The value of `c` will always be 6 after these statements are executed.

*logical and and or operators* ▶ 5.1

*conditional operator* ▶ 5.2

*comma operator* ▶ 6.3

Besides the assignment operators, the only operators that modify their operands are increment and decrement. When using these operators, be careful that your expressions don't depend on a particular order of evaluation. In the following example, *j* may be assigned one of two values:

```
i = 2;  
j = i * i++;
```

It's natural to assume that *j* is assigned the value 4. However, the effect of executing the statement is undefined, and *j* could just as well be assigned 6 instead. Here's the scenario: (1) The second operand (the original value of *i*) is fetched, then *i* is incremented. (2) The first operand (the new value of *i*) is fetched. (3) The new and old values of *i* are multiplied, yielding 6. "Fetching" a variable means to retrieve the value of the variable from memory. A later change to the variable won't affect the fetched value, which is typically stored in a special location (known as a *register*) inside the CPU.

registers ➤ 18.2

---

### Undefined Behavior

According to the C standard, statements such as `c = (b = a + 2) - (a = 1);` and `j = i * i++;` cause **undefined behavior**, which is different from implementation-defined behavior (see Section 4.1). When a program ventures into the realm of undefined behavior, all bets are off. The program may behave differently when compiled with different compilers. But that's not the only thing that can happen. The program may not compile in the first place, if it compiles it may not run, and if it does run, it may crash, behave erratically, or produce meaningless results. In other words, undefined behavior should be avoided like the plague.

---

## 4.5 Expression Statements

C has the unusual rule that *any* expression can be used as a statement. That is, any expression—regardless of its type or what it computes—can be turned into a statement by appending a semicolon. For example, we could turn the expression `++i` into a statement:

```
++i;
```

When this statement is executed, *i* is first incremented, then the new value of *i* is fetched (as though it were to be used in an enclosing expression). However, since `++i` isn't part of a larger expression, its value is discarded and the next statement executed. (The change to *i* is permanent, of course.)

**Q&A**

Since its value is discarded, there's little point in using an expression as a statement unless the expression has a side effect. Let's look at three examples. In

the first example, 1 is stored into *i*, then the new value of *i* is fetched but not used:

```
i = 1;
```

In the second example, the value of *i* is fetched but not used; however, *i* is decremented afterwards:

```
i--;
```

In the third example, the value of the expression *i \* j - 1* is computed and then discarded:

```
i * j - 1;
```

Since *i* and *j* aren't changed, this statement has no effect and therefore serves no purpose.



A slip of the finger can easily create a “do-nothing” expression statement. For example, instead of entering

```
i = j;
```

we might accidentally type

```
i + j;
```

(This kind of error is more common than you might expect, since the = and + characters usually occupy the same key.) Some compilers can detect meaningless expression statements; you'll get a warning such as “*statement with no effect*.”

## Q & A

**Q: I notice that C has no exponentiation operator. How can I raise a number to a power?**

A: Raising an integer to a small positive integer power is best done by repeated multiplication (*i \* i \* i* is *i* cubed). To raise a number to a noninteger power, call the `pow` function ►23.3

**Q: I want to apply the % operator to a floating-point operand, but my program won't compile. What can I do? [p. 54]**

fmod function ►23.3 A: The % operator requires integer operands. Try the `fmod` function instead.

**Q: Why are the rules for using the / and % operators with negative operands so complicated? [p. 54]**

A: The rules aren't as complicated as they may first appear. In both C89 and C99, the goal is to ensure that the value of  $(a / b) * b + a \% b$  will always be equal to *a*.

(and indeed, both standards guarantee that this is the case, provided that the value of  $a / b$  is “representable”). The problem is that there are two ways for  $a / b$  and  $a \% b$  to satisfy this equality if either  $a$  or  $b$  is negative, as seen in C89, where either  $-9 / 7$  is  $-1$  and  $-9 \% 7$  is  $-2$ , or  $-9 / 7$  is  $-2$  and  $-9 \% 7$  is  $5$ . In the first case,  $(-9 / 7) * 7 + -9 \% 7$  has the value  $-1 \times 7 + -2 = -9$ , and in the second case,  $(-9 / 7) * 7 + -9 \% 7$  has the value  $-2 \times 7 + 5 = -9$ . By the time C99 rolled around, most CPUs were designed to truncate the result of division toward zero, so this was written into the standard as the only allowable outcome.

**Q: If C has lvalues, does it also have rvalues? [p. 59]**

A: Yes, indeed. An *lvalue* is an expression that can appear on the *left* side of an assignment; an *rvalue* is an expression that can appear on the *right* side. Thus, an rvalue could be a variable, constant, or more complex expression. In this book, as in the C standard, we’ll use the term “expression” instead of “rvalue.”

**\*Q: You said that  $v += e$  isn’t equivalent to  $v = v + e$  if  $v$  has a side effect. Can you explain? [p. 60]**

A: Evaluating  $v += e$  causes  $v$  to be evaluated only once; evaluating  $v = v + e$  causes  $v$  to be evaluated twice. Any side effect caused by evaluating  $v$  will occur twice in the latter case. In the following example,  $i$  is incremented once:

```
a[i++] += 2;
```

If we use  $=$  instead of  $+=$ , here’s what the statement will look like:

```
a[i++] = a[i++] + 2;
```

The value of  $i$  is modified as well as used elsewhere in the statement, so the effect of executing the statement is undefined. It’s likely that  $i$  will be incremented twice, but we can’t say with certainty what will happen.

**Q: Why does C provide the `++` and `--` operators? Are they faster than other ways of incrementing and decrementing, or they are just more convenient? [p. 61]**

A: C inherited `++` and `--` from Ken Thompson’s earlier B language. Thompson apparently created these operators because his B compiler could generate a more compact translation for `++i` than for `i = i + 1`. These operators have become a deeply ingrained part of C (in fact, many of C’s most famous idioms rely on them). With modern compilers, using `++` and `--` won’t make a compiled program any smaller or faster; the continued popularity of these operators stems mostly from their brevity and convenience.

**Q: Do `++` and `--` work with float variables?**

A: Yes; the increment and decrement operations can be applied to floating-point numbers as well as integers. In practice, however, it’s fairly rare to increment or decrement a `float` variable.

**\*Q:** When I use the postfix version of `++` or `--`, just when is the increment or decrement performed? [p. 62]

**A:** That's an excellent question. Unfortunately, it's also a difficult one to answer. The C standard introduces the concept of "sequence point" and says that "updating the stored value of the operand shall occur between the previous and the next sequence point." There are various kinds of sequence points in C; the end of an expression statement is one example. By the end of an expression statement, all increments and decrements within the statement must have been performed; the next statement can't begin to execute until this condition has been met.

Certain operators that we'll encounter in later chapters (logical *and*, logical *or*, conditional, and comma) also impose sequence points. So do function calls: the arguments in a function call must be fully evaluated before the call can be performed. If an argument happens to be an expression containing a `++` or `--` operator, the increment or decrement must occur before the call can take place.

**Q:** What do you mean when you say that the value of an expression statement is discarded? [p. 65]

**A:** By definition, an expression represents a value. If `i` has the value 5, for example, then evaluating `i + 1` produces the value 6. Let's turn `i + 1` into a statement by putting a semicolon after it:

```
i + 1;
```

When this statement is executed, the value of `i + 1` is computed. Since we have failed to save this value—or at least use it in some way—it is lost.

**Q:** But what about statements like `i = 1;`? I don't see what is being discarded.

**A:** Don't forget that `=` is an operator in C and produces a value just like any other operator. The assignment

```
i = 1;
```

assigns 1 to `i`. The value of the entire expression is 1, which is discarded. Discarding the expression's value is no great loss, since the reason for writing the statement in the first place was to modify `i`.

## Exercises

### Section 4.1

1. Show the output produced by each of the following program fragments. Assume that `i`, `j`, and `k` are `int` variables.
  - `i = 5; j = 3;  
printf("%d %d", i / j, i % j);`
  - `i = 2; j = 3;  
printf("%d", (i + 10) % j);`
  - `i = 7; j = 8; k = 9;  
printf("%d", (i + 10) % k / j);`

(d) `i = 1; j = 2; k = 3;  
printf("%d", (i + 5) % (j + 2) / k);`

- W \*2. If  $i$  and  $j$  are positive integers, does  $(-i)/j$  always have the same value as  $- (i/j)$ ? Justify your answer.
3. What is the value of each of the following expressions in C89? (Give all possible values if an expression may have more than one value.)
- $8 / 5$
  - $-8 / 5$
  - $8 / -5$
  - $-8 / -5$
4. Repeat Exercise 3 for C99.
5. What is the value of each of the following expressions in C89? (Give all possible values if an expression may have more than one value.)
- $8 \% 5$
  - $-8 \% 5$
  - $8 \% -5$
  - $-8 \% -5$
6. Repeat Exercise 5 for C99.
7. The algorithm for computing the UPC check digit ends with the following steps:  
Subtract 1 from the total.  
Compute the remainder when the adjusted total is divided by 10.  
Subtract the remainder from 9.  
It's tempting to try to simplify the algorithm by using these steps instead:  
Compute the remainder when the total is divided by 10.  
Subtract the remainder from 10.  
Why doesn't this technique work?
8. Would the `upc.c` program still work if the expression `9 - ((total - 1) % 10)` were replaced by `(10 - (total % 10)) % 10`?

#### Section 4.2

- W 9. Show the output produced by each of the following program fragments. Assume that  $i$ ,  $j$ , and  $k$  are `int` variables.
- `i = 7; j = 8;  
i *= j + 1;  
printf("%d %d", i, j);`
  - `i = j = k = 1;  
i += j += k;  
printf("%d %d %d", i, j, k);`
  - `i = 1; j = 2; k = 3;  
i -= j -= k;  
printf("%d %d %d", i, j, k);`
  - `i = 2; j = 1; k = 0;  
i *= j *= k;  
printf("%d %d %d", i, j, k);`

10. Show the output produced by each of the following program fragments. Assume that *i* and *j* are int variables.

```
(a) i = 6;
    j = i += i;
    printf("%d %d", i, j);

(b) i = 5;
    j = (i -= 2) + 1;
    printf("%d %d", i, j);

(c) i = 7;
    j = 6 + (i = 2.5);
    printf("%d %d", i, j);

(d) i = 2; j = 8;
    j = (i = 6) + (j = 3);
    printf("%d %d", i, j);
```

**Section 4.3**

- \*11. Show the output produced by each of the following program fragments. Assume that *i*, *j*, and *k* are int variables.

```
(a) i = 1;
    printf("%d ", i++ - 1);
    printf("%d", i);

(b) i = 10; j = 5;
    printf("%d ", i++ - ++j);
    printf("%d %d", i, j);

(c) i = 7; j = 8;
    printf("%d ", i++ - --j);
    printf("%d %d", i, j);

(d) i = 3; j = 4; k = 5;
    printf("%d ", i++ - j++ + --k);
    printf("%d %d %d", i, j, k);
```

12. Show the output produced by each of the following program fragments. Assume that *i* and *j* are int variables.

```
(a) i = 5;
    j = ++i * 3 - 2;
    printf("%d %d", i, j);

(b) i = 5;
    j = 3 - 2 * i++;
    printf("%d %d", i, j);

(c) i = 7;
    j = 3 * i-- + 2;
    printf("%d %d", i, j);

(d) i = 7;
    j = 3 + --i * 2;
    printf("%d %d", i, j);
```

- W 13. Only one of the expressions  $++i$  and  $i++$  is exactly the same as  $(i += 1)$ ; which is it? Justify your answer.

**Section 4.4**

14. Supply parentheses to show how a C compiler would interpret each of the following expressions.

- (a)  $a * b - c * d + e$
- (b)  $a / b \% c / d$
- (c)  $- a - b + c - + d$
- (d)  $a * - b / c - d$

**Section 4.5**

15. Give the values of *i* and *j* after each of the following expression statements has been executed. (Assume that *i* has the value 1 initially and *j* has the value 2.)
- (a)  $i += j;$
  - (b)  $i--;$
  - (c)  $i * j / i;$
  - (d)  $i \% ++j;$

---

## Programming Projects

1. Write a program that asks the user to enter a two-digit number, then prints the number with its digits reversed. A session with the program should have the following appearance:

Enter a two-digit number: 28  
The reversal is: 82

Read the number using `%d`, then break it into two digits. *Hint:* If *n* is an integer, then *n % 10* is the last digit in *n* and *n / 10* is *n* with the last digit removed.

- W 2. Extend the program in Programming Project 1 to handle *three*-digit numbers.
3. Rewrite the program in Programming Project 2 so that it prints the reversal of a three-digit number without using arithmetic to split the number into digits. *Hint:* See the `upc.c` program of Section 4.1.
4. Write a program that reads an integer entered by the user and displays it in octal (base 8):

Enter a number between 0 and 32767: 1953  
In octal, your number is: 03641

The output should be displayed using five digits, even if fewer digits are sufficient. *Hint:* To convert the number to octal, first divide it by 8; the remainder is the last digit of the octal number (1, in this case). Then divide the original number by 8 and repeat the process to arrive at the next-to-last digit. (`printf` is capable of displaying numbers in base 8, as we'll see in Chapter 7, so there's actually an easier way to write this program.)

5. Rewrite the `upc.c` program of Section 4.1 so that the user enters 11 digits at one time, instead of entering one digit, then five digits, and then another five digits.

Enter the first 11 digits of a UPC: 01380015173  
Check digit: 5

6. European countries use a 13-digit code, known as a European Article Number (EAN) instead of the 12-digit Universal Product Code (UPC) found in North America. Each EAN ends with a check digit, just as a UPC does. The technique for calculating the check digit is also similar:

Add the second, fourth, sixth, eighth, tenth, and twelfth digits.  
Add the first, third, fifth, seventh, ninth, and eleventh digits.  
Multiply the first sum by 3 and add it to the second sum.

Subtract 1 from the total.

Compute the remainder when the adjusted total is divided by 10.

Subtract the remainder from 9.

For example, consider Güllüoglu Turkish Delight Pistachio & Coconut, which has an EAN of 8691484260008. The first sum is  $6 + 1 + 8 + 2 + 0 + 0 = 17$ , and the second sum is  $8 + 9 + 4 + 4 + 6 + 0 = 31$ . Multiplying the first sum by 3 and adding the second yields 82. Subtracting 1 gives 81. The remainder upon dividing by 10 is 1. When the remainder is subtracted from 9, the result is 8, which matches the last digit of the original code. Your job is to modify the `upc.c` program of Section 4.1 so that it calculates the check digit for an EAN. The user will enter the first 12 digits of the EAN as a single number:

```
Enter the first 12 digits of an EAN: 869148426000
Check digit: 8
```

# 5 Selection Statements

*Programmers are not to be measured by their ingenuity and their logic but by the completeness of their case analysis.*

return statement ➤ 2.2  
expression statement ➤ 4.5

Although C has many operators, it has relatively few statements. We've encountered just two so far: the `return` statement and the expression statement. Most of C's remaining statements fall into three categories, depending on how they affect the order in which statements are executed:

- **Selection statements.** The `if` and `switch` statements allow a program to select a particular execution path from a set of alternatives.
- **Iteration statements.** The `while`, `do`, and `for` statements support iteration (looping).
- **Jump statements.** The `break`, `continue`, and `goto` statements cause an unconditional jump to some other place in the program. (The `return` statement belongs in this category, as well.)

The only other statements in C are the compound statement, which groups several statements into a single statement, and the null statement, which performs no action.

This chapter discusses the selection statements and the compound statement. (Chapter 6 covers the iteration statements, the jump statements, and the null statement.) Before we can write `if` statements, we'll need logical expressions: conditions that `if` statements can test. Section 5.1 explains how logical expressions are built from the relational operators (`<`, `<=`, `>`, and `>=`), the equality operators (`==` and `!=`), and the logical operators (`&&`, `||`, and `!`). Section 5.2 covers the `if` statement and compound statement, as well as introducing the conditional operator (`? :`), which can test a condition within an expression. Section 5.3 describes the `switch` statement.

## 5.1 Logical Expressions

Several of C's statements, including the `if` statement, must test the value of an expression to see if it is “true” or “false.” For example, an `if` statement might need to test the expression `i < j`; a true value would indicate that `i` is less than `j`. In many programming languages, an expression such as `i < j` would have a special “Boolean” or “logical” type. Such a type would have only two values, `false` and `true`. In C, however, a comparison such as `i < j` yields an integer: either 0 (false) or 1 (true). With this in mind, let's look at the operators that are used to build logical expressions.

### Relational Operators

C's *relational operators* (Table 5.1) correspond to the `<`, `>`, `≤`, and `≥` operators of mathematics, except that they produce 0 (false) or 1 (true) when used in expressions. For example, the value of `10 < 11` is 1; the value of `11 < 10` is 0.

**Table 5.1**  
Relational Operators

Symbol	Meaning
<code>&lt;</code>	less than
<code>&gt;</code>	greater than
<code>≤</code>	less than or equal to
<code>≥</code>	greater than or equal to

The relational operators can be used to compare integers and floating-point numbers, with operands of mixed types allowed. Thus, `1 < 2.5` has the value 1, while `5.6 < 4` has the value 0.

The precedence of the relational operators is lower than that of the arithmetic operators; for example, `i + j < k - 1` means `(i + j) < (k - 1)`. The relational operators are left associative.



The expression

`i < j < k`

is legal in C, but doesn't have the meaning that you might expect. Since the `<` operator is left associative, this expression is equivalent to

`(i < j) < k`

In other words, the expression first tests whether `i` is less than `j`; the 1 or 0 produced by this comparison is then compared to `k`. The expression does *not* test whether `j` lies between `i` and `k`. (We'll see later in this section that the correct expression would be `i < j && j < k`.)

## Equality Operators

Although the relational operators are denoted by the same symbols as in many other programming languages, the *equality operators* have a unique appearance (Table 5.2). The “equal to” operator is two adjacent = characters, not one, since a single = character represents the assignment operator. The “not equal to” operator is also two characters: ! and =.

**Table 5.2**  
Equality Operators

Symbol	Meaning
<code>==</code>	equal to
<code>!=</code>	not equal to

Like the relational operators, the equality operators are left associative and produce either 0 (false) or 1 (true) as their result. However, the equality operators have *lower* precedence than the relational operators. For example, the expression

`i < j == j < k`

is equivalent to

`(i < j) == (j < k)`

which is true if `i < j` and `j < k` are both true or both false.

Clever programmers sometimes exploit the fact that the relational and equality operators return integer values. For example, the value of the expression `(i >= j) + (i == j)` is either 0, 1, or 2, depending on whether `i` is less than, greater than, or equal to `j`, respectively. Tricky coding like this generally isn’t a good idea, however; it makes programs hard to understand.

## Logical Operators

More complicated logical expressions can be built from simpler ones by using the *logical operators*: *and*, *or*, and *not* (Table 5.3). The `!` operator is unary, while `&&` and `||` are binary.

**Table 5.3**  
Logical Operators

Symbol	Meaning
<code>!</code>	logical negation
<code>&amp;&amp;</code>	logical <i>and</i>
<code>  </code>	logical <i>or</i>

The logical operators produce either 0 or 1 as their result. Often, the operands will have values of 0 or 1, but this isn’t a requirement; the logical operators treat any nonzero operand as a true value and any zero operand as a false value.

The logical operators behave as follows:

- `!expr` has the value 1 if `expr` has the value 0.
- `expr1 && expr2` has the value 1 if the values of `expr1` and `expr2` are both nonzero.

- $expr1 \mid\mid expr2$  has the value 1 if either  $expr1$  or  $expr2$  (or both) has a nonzero value.

In all other cases, these operators produce the value 0.

Both `&&` and `||` perform “short-circuit” evaluation of their operands. That is, these operators first evaluate the left operand, then the right operand. If the value of the expression can be deduced from the value of the left operand alone, then the right operand isn’t evaluated. Consider the following expression:

```
(i != 0) && (j / i > 0)
```

To find the value of this expression, we must first evaluate `(i != 0)`. If `i` isn’t equal to 0, then we’ll need to evaluate `(j / i > 0)` to determine whether the entire expression is true or false. However, if `i` is equal to 0, then the entire expression must be false, so there’s no need to evaluate `(j / i > 0)`. The advantage of short-circuit evaluation is apparent—without it, evaluating the expression would have caused a division by zero.



Be wary of side effects in logical expressions. Thanks to the short-circuit nature of the `&&` and `||` operators, side effects in operands may not always occur. Consider the following expression:

```
i > 0 && ++j > 0
```

Although `j` is apparently incremented as a side effect of evaluating the expression, that isn’t always the case. If `i > 0` is false, then `++j > 0` is not evaluated, so `j` isn’t incremented. The problem can be fixed by changing the condition to `++j > 0 && i > 0` or, even better, by incrementing `j` separately.

The `!` operator has the same precedence as the unary plus and minus operators. The precedence of `&&` and `||` is lower than that of the relational and equality operators; for example, `i < j && k == m` means `(i < j) && (k == m)`. The `!` operator is right associative; `&&` and `||` are left associative.

## 5.2 The if Statement

The `if` statement allows a program to choose between two alternatives by testing the value of an expression. In its simplest form, the `if` statement has the form

**if statement**

`if ( expression ) statement`

Notice that the parentheses around the expression are mandatory; they’re part of the `if` statement, not part of the expression. Also note that the word `then` doesn’t come after the parentheses, as it would in some programming languages.

When an `if` statement is executed, the expression in the parentheses is evaluated; if the value of the expression is nonzero—which C interprets as true—the statement after the parentheses is executed. Here's an example:

```
if (line_num == MAX_LINES)
    line_num = 0;
```

The statement `line_num = 0;` is executed if the condition `line_num == MAX_LINES` is true (has a nonzero value).



Don't confuse `==` (equality) with `=` (assignment). The statement

```
if (i == 0) ...
```

tests whether `i` is equal to 0. However, the statement

```
if (i = 0) ...
```

assigns 0 to `i`, then tests whether the *result* is nonzero. In this case, the test always fails.

Confusing `==` with `=` is perhaps the most common C programming error, probably because `=` means “is equal to” in mathematics (and in certain programming languages). Some compilers issue a warning if they notice `=` where `==` would normally appear.

**Q&A**

Often the expression in an `if` statement will test whether a variable falls within a range of values. To test whether  $0 \leq i < n$ , for example, we'd write

**idiom** `if (0 <= i && i < n) ...`

To test the *opposite* condition (`i` is outside the range), we'd write

**idiom** `if (i < 0 || i >= n) ...`

Note the use of the `||` operator instead of the `&&` operator.

## Compound Statements

In our `if` statement template, notice that *statement* is singular, not plural:

```
if ( expression ) statement
```

What if we want an `if` statement to control *two* or more statements? That's where the **compound statement** comes in. A compound statement has the form

**compound statement**

```
{ statements }
```

By putting braces around a group of statements, we can force the compiler to treat it as a single statement.

Here's an example of a compound statement:

```
{ line_num = 0; page_num++; }
```

For clarity, I'll usually put a compound statement on several lines, with one statement per line:

```
{
    line_num = 0;
    page_num++;
}
```

Notice that each inner statement still ends with a semicolon, but the compound statement itself does not.

Here's what a compound statement would look like when used inside an `if` statement:

```
if (line_num == MAX_LINES) {
    line_num = 0;
    page_num++;
}
```

Compound statements are also common in loops and other places where the syntax of C requires a single statement, but we want more than one.

## The `else` Clause

An `if` statement may have an `else` clause:

### **if statement with else clause**

```
if ( expression ) statement else statement
```

The statement that follows the word `else` is executed if the expression in parentheses has the value 0.

Here's an example of an `if` statement with an `else` clause:

```
if (i > j)
    max = i;
else
    max = j;
```

Notice that both "inner" statements end with a semicolon.

When an `if` statement contains an `else` clause, a layout issue arises: where should the `else` be placed? Many C programmers align it with the `if` at the beginning of the statement, as in the previous example. The inner statements are usually indented, but if they're short they can be put on the same line as the `if` and `else`:

```
if (i > j) max = i;
else max = j;
```

There are no restrictions on what kind of statements can appear inside an `if` statement. In fact, it's not unusual for `if` statements to be nested inside other `if` statements. Consider the following `if` statement, which finds the largest of the numbers stored in `i`, `j`, and `k` and stores that value in `max`:

```
if (i > j)
    if (i > k)
        max = i;
    else
        max = k;
else
    if (j > k)
        max = j;
    else
        max = k;
```

`if` statements can be nested to any depth. Notice how aligning each `else` with the matching `if` makes the nesting easier to see. If you still find the nesting confusing, don't hesitate to add braces:

```
if (i > j) {
    if (i > k)
        max = i;
    else
        max = k;
} else {
    if (j > k)
        max = j;
    else
        max = k;
}
```

Adding braces to statements—even when they're not necessary—is like using parentheses in expressions: both techniques help make a program more readable while at the same time avoiding the possibility that the compiler won't understand the program the way we thought it did.

Some programmers use as many braces as possible inside `if` statements (and iteration statements as well). A programmer who adopts this convention would include a pair of braces for every `if` clause and every `else` clause:

```
if (i > j) {
    if (i > k) {
        max = i;
    } else {
        max = k;
    }
} else {
    if (j > k) {
        max = j;
    } else {
        max = k;
    }
}
```

Using braces even when they're not required has two advantages. First, the program becomes easier to modify, because more statements can easily be added to any `if` or `else` clause. Second, it helps avoid errors that can result from forgetting to use braces when adding statements to an `if` or `else` clause.

## Cascaded if Statements

We'll often need to test a series of conditions, stopping as soon as one of them is true. A "cascaded" `if` statement is often the best way to write such a series of tests. For example, the following cascaded `if` statement tests whether `n` is less than 0, equal to 0, or greater than 0:

```
if (n < 0)
    printf("n is less than 0\n");
else
    if (n == 0)
        printf("n is equal to 0\n");
    else
        printf("n is greater than 0\n");
```

Although the second `if` statement is nested inside the first, C programmers don't usually indent it. Instead, they align each `else` with the original `if`:

```
if (n < 0)
    printf("n is less than 0\n");
else if (n == 0)
    printf("n is equal to 0\n");
else
    printf("n is greater than 0\n");
```

This arrangement gives the cascaded `if` a distinctive appearance:

```
if ( expression )
    statement
else if ( expression )
    statement
...
else if ( expression )
    statement
else
    statement
```

The last two lines (`else statement`) aren't always present, of course. This way of indenting the cascaded `if` statement avoids the problem of excessive indentation when the number of tests is large. Moreover, it assures the reader that the statement is nothing more than a series of tests.

Keep in mind that a cascaded `if` statement isn't some new kind of statement; it's just an ordinary `if` statement that happens to have another `if` statement as its `else` clause (and *that if statement* has another `if` statement as its `else` clause, *ad infinitum*).

**PROGRAM** **Calculating a Broker's Commission**

When stocks are sold or purchased through a broker, the broker's commission is often computed using a sliding scale that depends upon the value of the stocks traded. Let's say that a broker charges the amounts shown in the following table:

<i>Transaction size</i>	<i>Commission rate</i>
Under \$2,500	\$30 + 1.7%
\$2,500–\$6,250	\$56 + 0.66%
\$6,250–\$20,000	\$76 + 0.34%
\$20,000–\$50,000	\$100 + 0.22%
\$50,000–\$500,000	\$155 + 0.11%
Over \$500,000	\$255 + 0.09%

The minimum charge is \$39. Our next program asks the user to enter the amount of the trade, then displays the amount of the commission:

```
Enter value of trade: 30000
Commission: $166.00
```

The heart of the program is a cascaded if statement that determines which range the trade falls into.

```
broker.c /* Calculates a broker's commission */

#include <stdio.h>

int main(void)
{
    float commission, value;

    printf("Enter value of trade: ");
    scanf("%f", &value);

    if (value < 2500.00f)
        commission = 30.00f + .017f * value;
    else if (value < 6250.00f)
        commission = 56.00f + .0066f * value;
    else if (value < 20000.00f)
        commission = 76.00f + .0034f * value;
    else if (value < 50000.00f)
        commission = 100.00f + .0022f * value;
    else if (value < 500000.00f)
        commission = 155.00f + .0011f * value;
    else
        commission = 255.00f + .0009f * value;

    if (commission < 39.00f)
        commission = 39.00f;

    printf("Commission: $%.2f\n", commission);

    return 0;
}
```

The cascaded if statement could have been written this way instead (the changes are indicated in **bold**):

```
if (value < 2500.00f)
    commission = 30.00f + .017f * value;
else if (value >= 2500.00f && value < 6250.00f)
    commission = 56.00f + .0066f * value;
else if (value >= 6250.00f && value < 20000.00f)
    commission = 76.00f + .0034f * value;
...
```

Although the program will still work, the added conditions aren't necessary. For example, the first if clause tests whether value is less than 2500 and, if so, computes the commission. When we reach the second if test (**value** >= 2500.00f && value < 6250.00f), we know that value can't be less than 2500 and therefore must be greater than or equal to 2500. The condition **value** >= 2500.00f will always be true, so there's no point in checking it.

## The “Dangling else” Problem

When if statements are nested, we've got to watch out for the notorious “dangling else” problem. Consider the following example:

```
if (y != 0)
    if (x != 0)
        result = x / y;
else
    printf("Error: y is equal to 0\n");
```

To which if statement does the else clause belong? The indentation suggests that it belongs to the outer if statement. However, C follows the rule that an else clause belongs to the nearest if statement that hasn't already been paired with an else. In this example, the else clause actually belongs to the inner if statement, so a correctly indented version would look like this:

```
if (y != 0)
    if (x != 0)
        result = x / y;
    else
        printf("Error: y is equal to 0\n");
```

To make the else clause part of the outer if statement, we can enclose the inner if statement in braces:

```
if (y != 0) {
    if (x != 0)
        result = x / y;
} else
    printf("Error: y is equal to 0\n");
```

This example illustrates the value of braces; if we'd used them in the original if statement, we wouldn't have gotten into this situation in the first place.

## Conditional Expressions

C's `if` statement allows a program to perform one of two actions depending on the value of a condition. C also provides an *operator* that allows an expression to produce one of two *values* depending on the value of a condition.

The *conditional operator* consists of two symbols (`?` and `:`), which must be used together in the following way:

conditional  
expression

`expr1 ? expr2 : expr3`

`expr1`, `expr2`, and `expr3` can be expressions of any type. The resulting expression is said to be a *conditional expression*. The conditional operator is unique among C operators in that it requires *three* operands instead of one or two. For this reason, it is often referred to as a *ternary* operator.

The conditional expression `expr1 ? expr2 : expr3` should be read "if `expr1` then `expr2` else `expr3`." The expression is evaluated in stages: `expr1` is evaluated first; if its value isn't zero, then `expr2` is evaluated, and its value is the value of the entire conditional expression. If the value of `expr1` is zero, then the value of `expr3` is the value of the conditional.

The following example illustrates the conditional operator:

```
int i, j, k;  
  
i = 1;  
j = 2;  
k = i > j ? i : j;           /* k is now 2 */  
k = (i >= 0 ? i : 0) + j;   /* k is now 3 */
```

The conditional expression `i > j ? i : j` in the first assignment to `k` returns the value of either `i` or `j`, depending on which one is larger. Since `i` has the value 1 and `j` has the value 2, the `i > j` comparison fails, and the value of the conditional is 2, which is assigned to `k`. In the second assignment to `k`, the `i >= 0` comparison succeeds; the conditional expression `(i >= 0 ? i : 0)` has the value 1, which is then added to `j` to produce 3. The parentheses are necessary, by the way; the precedence of the conditional operator is less than that of the other operators we've discussed so far, with the exception of the assignment operators.

Conditional expressions tend to make programs shorter but harder to understand, so it's probably best to avoid them. There are, however, a few places in which they're tempting; one is the `return` statement. Instead of writing

```
if (i > j)  
    return i;  
else  
    return j;
```

many programmers would write

```
return i > j ? i : j;
```

Calls of `printf` can sometimes benefit from condition expressions. Instead of

```
if (i > j)
    printf("%d\n", i);
else
    printf("%d\n", j);
```

we could simply write

```
printf("%d\n", i > j ? i : j);
```

macro definitions ▶ 14.3

Conditional expressions are also common in certain kinds of macro definitions.

## Boolean Values in C89

For many years, the C language lacked a proper Boolean type, and there is none defined in the C89 standard. This omission is a minor annoyance, since many programs need variables that can store either *false* or *true*. One way to work around this limitation of C89 is to declare an `int` variable and then assign it either 0 or 1:

```
int flag;

flag = 0;
...
flag = 1;
```

Although this scheme works, it doesn't contribute much to program readability. It's not obvious that `flag` is to be assigned only Boolean values and that 0 and 1 represent false and true.

To make programs more understandable, C89 programmers often define macros with names such as `TRUE` and `FALSE`:

```
#define TRUE 1
#define FALSE 0
```

Assignments to `flag` now have a more natural appearance:

```
flag = FALSE;
...
flag = TRUE;
```

To test whether `flag` is true, we can write

```
if (flag == TRUE) ...
```

or just

```
if (flag) ...
```

The latter form is better, not only because it's more concise, but also because it will still work correctly if `flag` has a value other than 0 or 1.

To test whether `flag` is false, we can write

```
if (flag == FALSE) ...
```

or

```
if (!flag) ...
```

Carrying this idea one step further, we might even define a macro that can be used as a type:

```
#define BOOL int
```

BOOL can take the place of int when declaring Boolean variables:

```
BOOL flag;
```

It's now clear that flag isn't an ordinary integer variable, but instead represents a Boolean condition. (The compiler still treats flag as an int variable, of course.) In later chapters, we'll discover better ways to set up a Boolean type in C89 by using type definitions and enumerations.

C99

## Boolean Values in C99

Q&A

The longstanding lack of a Boolean type has been remedied in C99, which provides the `_Bool` type. In this version of C, a Boolean variable can be declared by writing

```
_Bool flag;
```

unsigned integer types ► 7.1

`_Bool` is an integer type (more precisely, an *unsigned* integer type), so a `_Bool` variable is really just an integer variable in disguise. Unlike an ordinary integer variable, however, a `_Bool` variable can only be assigned 0 or 1. In general, attempting to store a nonzero value into a `_Bool` variable will cause the variable to be assigned 1:

```
flag = 5; /* flag is assigned 1 */
```

It's legal (although not advisable) to perform arithmetic on `_Bool` variables; it's also legal to print a `_Bool` variable (either 0 or 1 will be displayed). And, of course, a `_Bool` variable can be tested in an if statement:

```
if (flag) /* tests whether flag is 1 */  
...
```

<stdbool.h> header ► 21.5

In addition to defining the `_Bool` type, C99 also provides a new header, `<stdbool.h>`, that makes it easier to work with Boolean values. This header provides a macro, `bool`, that stands for `_Bool`. If `<stdbool.h>` is included, we can write

```
bool flag; /* same as _Bool flag; */
```

The `<stdbool.h>` header also supplies macros named `true` and `false`, which stand for 1 and 0, respectively, making it possible to write

```
flag = false;  
...  
flag = true;
```

Because the `<stdbool.h>` header is so handy, I'll use it in subsequent programs whenever Boolean variables are needed.

## 5.3 The switch Statement

In everyday programming, we'll often need to compare an expression against a series of values to see which one it currently matches. We saw in Section 5.2 that a cascaded `if` statement can be used for this purpose. For example, the following cascaded `if` statement prints the English word that corresponds to a numerical grade:

```
if (grade == 4)  
    printf("Excellent");  
else if (grade == 3)  
    printf("Good");  
else if (grade == 2)  
    printf("Average");  
else if (grade == 1)  
    printf("Poor");  
else if (grade == 0)  
    printf("Failing");  
else  
    printf("Illegal grade");
```

As an alternative to this kind of cascaded `if` statement, C provides the `switch` statement. The following `switch` is equivalent to our cascaded `if`:

```
switch (grade) {  
    case 4: printf("Excellent");  
              break;  
    case 3: printf("Good");  
              break;  
    case 2: printf("Average");  
              break;  
    case 1: printf("Poor");  
              break;  
    case 0: printf("Failing");  
              break;  
    default: printf("Illegal grade");  
              break;  
}
```

When this statement is executed, the value of the variable `grade` is tested against 4, 3, 2, 1, and 0. If it matches 4, for example, the message `Excellent` is printed, then the `break` statement transfers control to the statement following the `switch`. If the value of `grade` doesn't match any of the choices listed, the `default` case applies, and the message `Illegal grade` is printed.

A switch statement is often easier to read than a cascaded if statement. Moreover, switch statements are often faster than if statements, especially when there are more than a handful of cases.

**Q&A**

In its most common form, the switch statement has the form

**switch statement**

```
switch ( expression ) {  
    case constant-expression : statements  
    ...  
    case constant-expression : statements  
    default : statements  
}
```

The switch statement is fairly complex; let's look at its components one by one:

characters ➤ 7.3

- **Controlling expression.** The word switch must be followed by an integer expression in parentheses. Characters are treated as integers in C and thus can be tested in switch statements. Floating-point numbers and strings don't qualify, however.
- **Case labels.** Each case begins with a label of the form

```
case constant-expression :
```

A **constant expression** is much like an ordinary expression except that it can't contain variables or function calls. Thus, 5 is a constant expression, and  $5 + 10$  is a constant expression, but  $n + 10$  isn't a constant expression (unless  $n$  is a macro that represents a constant). The constant expression in a case label must evaluate to an integer (characters are also acceptable).

- **Statements.** After each case label comes any number of statements. No braces are required around the statements. (Enjoy it—this is one of the few places in C where braces aren't required.) The last statement in each group is normally break.

Duplicate case labels aren't allowed. The order of the cases doesn't matter; in particular, the default case doesn't need to come last.

Only one constant expression may follow the word case; however, several case labels may precede the same group of statements:

```
switch (grade) {  
    case 4:  
    case 3:  
    case 2:  
    case 1: printf("Passing");  
              break;  
    case 0: printf("Failing");  
              break;  
    default: printf("Illegal grade");  
              break;  
}
```

To save space, programmers sometimes put several case labels on the same line:

```
switch (grade) {
    case 4: case 3: case 2: case 1:
        printf("Passing");
        break;
    case 0: printf("Failing");
        break;
    default: printf("Illegal grade");
        break;
}
```

Unfortunately, there's no way to write a case label that specifies a range of values, as there is in some programming languages.

A `switch` statement isn't required to have a `default` case. If `default` is missing and the value of the controlling expression doesn't match any of the case labels, control simply passes to the next statement after the `switch`.

## The Role of the `break` Statement

Now, let's take a closer look at the mysterious `break` statement. As we've seen, executing a `break` statement causes the program to "break" out of the `switch` statement; execution continues at the next statement after the `switch`.

The reason that we need `break` has to do with the fact that the `switch` statement is really a form of "computed jump." When the controlling expression is evaluated, control jumps to the case label matching the value of the `switch` expression. A case label is nothing more than a marker indicating a position within the `switch`. When the last statement in the case has been executed, control "falls through" to the first statement in the following case; the case label for the next case is ignored. Without `break` (or some other jump statement), control will flow from one case into the next. Consider the following `switch` statement:

```
switch (grade) {
    case 4: printf("Excellent");
    case 3: printf("Good");
    case 2: printf("Average");
    case 1: printf("Poor");
    case 0: printf("Failing");
    default: printf("Illegal grade");
}
```

If the value of `grade` is 3, the message printed is

GoodAveragePoorFailingIllegal grade



Forgetting to use `break` is a common error. Although omitting `break` is sometimes done intentionally to allow several cases to share code, it's usually just an oversight.

Since deliberately falling through from one case into the next is rare, it's a good idea to point out any deliberate omission of break:

```
switch (grade) {
    case 4: case 3: case 2: case 1:
        num_passing++;
        /* FALL THROUGH */
    case 0: total_grades++;
        break;
}
```

Without the comment, someone might later fix the “error” by adding an unwanted break statement.

Although the last case in a switch statement never needs a break statement, it's common practice to put one there anyway to guard against a “missing break” problem if cases should later be added.

## PROGRAM Printing a Date in Legal Form

Contracts and other legal documents are often dated in the following way:

*Dated this \_\_\_\_\_ day of \_\_\_\_\_, 20\_\_.*

Let's write a program that displays dates in this form. We'll have the user enter the date in month/day/year form, then we'll display the date in “legal” form:

```
Enter date (mm/dd/yy) : 7/19/14
Dated this 19th day of July, 2014.
```

We can get printf to do most of the formatting. However, we're left with two problems: how to add “th” (or “st” or “nd” or “rd”) to the day, and how to print the month as a word instead of a number. Fortunately, the switch statement is ideal for both situations; we'll have one switch print the day suffix and another print the month name.

```
date.c /* Prints a date in legal form */

#include <stdio.h>

int main(void)
{
    int month, day, year;

    printf("Enter date (mm/dd/yy) : ");
    scanf("%d /%d /%d", &month, &day, &year);

    printf("Dated this %d", day);
    switch (day) {
        case 1: case 21: case 31:
            printf("st"); break;
        case 2: case 22:
            printf("nd"); break;
```

```
        case 3: case 23:
            printf("rd"); break;
        default: printf("th"); break;
    }
    printf(" day of ");

    switch (month) {
        case 1: printf("January"); break;
        case 2: printf("February"); break;
        case 3: printf("March"); break;
        case 4: printf("April"); break;
        case 5: printf("May"); break;
        case 6: printf("June"); break;
        case 7: printf("July"); break;
        case 8: printf("August"); break;
        case 9: printf("September"); break;
        case 10: printf("October"); break;
        case 11: printf("November"); break;
        case 12: printf("December"); break;
    }

    printf(", 20%.2d.\n", year);
    return 0;
}
```

Note the use of `% .2d` to display the last two digits of the year. If we had used `%d` instead, single-digit years would be displayed incorrectly (2005 would be printed as 205).

## Q & A

**Q:** My compiler doesn't give a warning when I use `=` instead of `==`. Is there some way to force the compiler to notice the problem? [p. 77]

**A:** Here's a trick that some programmers use: instead of writing

```
if (i == 0) ...
```

they habitually write

```
if (0 == i) ...
```

Now suppose that the `==` operator is accidentally written as `=`:

```
if (0 = i) ...
```

The compiler will produce an error message, since it's not possible to assign a value to 0. I don't use this trick, because I think it makes programs look unnatural. Also, it can be used only when one of the operands in the test condition isn't an lvalue.

Fortunately, many compilers are capable of checking for suspect uses of the `=` operator in `if` conditions. The GCC compiler, for example, will perform this

check if the `-Wparentheses` option is used or if `-Wall` (all warnings) is selected. GCC allows the programmer to suppress the warning in a particular case by enclosing the `if` condition in a second set of parentheses:

```
if ((i == j)) ...
```

**Q: C books seem to use several different styles of indentation and brace placement for compound statements. Which style is best?**

**A:** According to *The New Hacker's Dictionary* (Cambridge, Mass.: MIT Press, 1996), there are four common styles of indentation and brace placement:

- The *K&R style*, used in Kernighan and Ritchie's *The C Programming Language*, is the one I've chosen for the programs in this book. In the K&R style, the left brace appears at the end of a line:

```
if (line_num == MAX_LINES) {  
    line_num = 0;  
    page_num++;  
}
```

The K&R style keeps programs compact by not putting the left brace on a line by itself. A disadvantage: the left brace can be hard to find. (I don't consider this a problem, since the indentation of the inner statements makes it clear where the left brace should be.) The K&R style is the one most often used in Java, by the way.

- The *Allman style*, named after Eric Allman (the author of `sendmail` and other UNIX utilities), puts the left brace on a separate line:

```
if (line_num == MAX_LINES)  
{  
    line_num = 0;  
    page_num++;  
}
```

This style makes it easy to check that braces come in matching pairs.

- The *Whitesmiths style*, popularized by the Whitesmiths C compiler, dictates that braces be indented:

```
if (line_num == MAX_LINES)  
{  
    line_num = 0;  
    page_num++;  
}
```

- The *GNU style*, used in software developed by the GNU Project, indents the braces, then further indents the inner statements:

```
if (line_num == MAX_LINES)  
{  
    line_num = 0;  
    page_num++;  
}
```

Which style you use is mainly a matter of taste; there's no proof that one style is clearly better than the others. In any event, choosing the right style is less important than applying it consistently.

**Q:** If `i` is an `int` variable and `f` is a `float` variable, what is the type of the conditional expression (`i > 0 ? i : f`)?

A: When `int` and `float` values are mixed in a conditional expression, as they are here, the expression has type `float`. If `i > 0` is true, the value of the expression will be the value of `i` after conversion to `float` type.

**Q:** Why doesn't C99 have a better name for its Boolean type? [p. 85]

A: `_Bool` isn't a very elegant name, is it? More common names, such as `bool` or `boolean`, weren't chosen because existing C programs might already define these names, causing older code not to compile.

**Q:** OK, so why wouldn't the name `_Bool` break older programs as well?

A: The C89 standard specifies that names beginning with an underscore followed by an uppercase letter are reserved for future use and should not be used by programmers.

**\*Q:** The template given for the `switch` statement described it as the "most common form." Are there other forms? [p. 87]

A: The `switch` statement is a bit more general than described in this chapter, although the description given here is general enough for virtually all programs. For example, a `switch` statement can contain labels that aren't preceded by the word `case`, which leads to an amusing (?) trap. Suppose that we accidentally misspell the word `default`:

```
switch (...) {  
    ...  
    defualt: ...  
}
```

The compiler may not detect the error, since it assumes that `defualt` is an ordinary label.

**Q:** I've seen several methods of indenting the `switch` statement. Which way is best?

A: There are at least two common methods. One is to put the statements in each case *after* the case label:

```
switch (coin) {  
    case 1: printf("Cent");  
              break;  
    case 5: printf("Nickel");  
              break;  
    case 10: printf("Dime");  
              break;
```

```

    case 25: printf("Quarter");
               break;
}

```

If each case consists of a single action (a call of `printf`, in this example), the `break` statement could even go on the same line as the action:

```

switch (coin) {
    case 1: printf("Cent"); break;
    case 5: printf("Nickel"); break;
    case 10: printf("Dime"); break;
    case 25: printf("Quarter"); break;
}

```

The other method is to put the statements *under* the case label, indenting the statements to make the case label stand out:

```

switch (coin) {
    case 1:
        printf("Cent");
        break;
    case 5:
        printf("Nickel");
        break;
    case 10:
        printf("Dime");
        break;
    case 25:
        printf("Quarter");
        break;
}

```

In one variation of this scheme, each case label is aligned under the word `switch`.

The first method is fine when the statements in each case are short and there are relatively few of them. The second method is better for large `switch` statements in which the statements in each case are complex and/or numerous.

## Exercises

### Section 5.1

1. The following program fragments illustrate the relational and equality operators. Show the output produced by each, assuming that `i`, `j`, and `k` are `int` variables.
  - `i = 2; j = 3;`  
`k = i * j == 6;`  
`printf("%d", k);`
  - `i = 5; j = 10; k = 1;`  
`printf("%d", k > i < j);`
  - `i = 3; j = 2; k = 1;`  
`printf("%d", i < j == j < k);`
  - `i = 3; j = 4; k = 5;`  
`printf("%d", i % j + i < k);`

- W 2. The following program fragments illustrate the logical operators. Show the output produced by each, assuming that *i*, *j*, and *k* are int variables.
- i* = 10; *j* = 5;  
`printf("%d", !i < j);`
  - i* = 2; *j* = 1;  
`printf("%d", !!i + !j);`
  - i* = 5; *j* = 0; *k* = -5;  
`printf("%d", i && j || k);`
  - i* = 1; *j* = 2; *k* = 3;  
`printf("%d", i < j || k);`
- \*3. The following program fragments illustrate the short-circuit behavior of logical expressions. Show the output produced by each, assuming that *i*, *j*, and *k* are int variables.
- i* = 3; *j* = 4; *k* = 5;  
`printf("%d ", i < j || ++j < k);`  
`printf("%d %d %d", i, j, k);`
  - i* = 7; *j* = 8; *k* = 9;  
`printf("%d ", i - 7 && j++ < k);`  
`printf("%d %d %d", i, j, k);`
  - i* = 7; *j* = 8; *k* = 9;  
`printf("%d ", (i = j) || (j = k));`  
`printf("%d %d %d", i, j, k);`
  - i* = 1; *j* = 1; *k* = 1;  
`printf("%d ", ++i || ++j && ++k);`  
`printf("%d %d %d", i, j, k);`
- W \*4. Write a single expression whose value is either -1, 0, or +1, depending on whether *i* is less than, equal to, or greater than *j*, respectively.

### Section 5.2

- \*5. Is the following if statement legal?

```
if (n >= 1 <= 10)
    printf("n is between 1 and 10\n");
```

If so, what does it do when *n* is equal to 0?

- W \*6. Is the following if statement legal?

```
if (n == 1-10)
    printf("n is between 1 and 10\n");
```

If so, what does it do when *n* is equal to 5?

7. What does the following statement print if *i* has the value 17? What does it print if *i* has the value -17?

```
printf("%d\n", i >= 0 ? i : -i);
```

8. The following if statement is unnecessarily complicated. Simplify it as much as possible. (*Hint:* The entire statement can be replaced by a single assignment.)

```
if (age >= 13)
    if (age <= 19)
        teenager = true;
    else
        teenager = false;
else if (age < 13)
    teenager = false;
```

9. Are the following if statements equivalent? If not, why not?

```
if (score >= 90)           if (score < 60)
    printf("A");
else if (score >= 80)     else if (score < 70)
    printf("B");
else if (score >= 70)     else if (score < 80)
    printf("C");
else if (score >= 60)     else if (score < 90)
    printf("D");
else                         else
    printf("F");           printf("A");
```

**Section 5.3**

- W\*10. What output does the following program fragment produce? (Assume that *i* is an integer variable.)

```
i = 1;
switch (i % 3) {
    case 0: printf("zero");
    case 1: printf("one");
    case 2: printf("two");
}
```

11. The following table shows telephone area codes in the state of Georgia along with the largest city in each area:

<i>Area code</i>	<i>Major city</i>
229	Albany
404	Atlanta
470	Atlanta
478	Macon
678	Atlanta
706	Columbus
762	Columbus
770	Atlanta
912	Savannah

Write a switch statement whose controlling expression is the variable *area\_code*. If the value of *area\_code* is in the table, the switch statement will print the corresponding city name. Otherwise, the switch statement will display the message "Area code not recognized". Use the techniques discussed in Section 5.3 to make the switch statement as simple as possible.

## Programming Projects

1. Write a program that calculates how many digits a number contains:

Enter a number: 374

The number 374 has 3 digits

You may assume that the number has no more than four digits. Hint: Use if statements to test the number. For example, if the number is between 0 and 9, it has one digit. If the number is between 10 and 99, it has two digits.

- W 2. Write a program that asks the user for a 24-hour time, then displays the time in 12-hour form:

Enter a 24-hour time: 21:11  
 Equivalent 12-hour time: 9:11 PM  
 Be careful not to display 12:00 as 0:00.

3. Modify the `broker.c` program of Section 5.2 by making both of the following changes:
- Ask the user to enter the number of shares and the price per share, instead of the value of the trade.
  - Add statements that compute the commission charged by a rival broker (\$33 plus 3¢ per share for fewer than 2000 shares; \$33 plus 2¢ per share for 2000 shares or more). Display the rival's commission as well as the commission charged by the original broker.

- W 4. Here's a simplified version of the Beaufort scale, which is used to estimate wind force:

<i>Speed (knots)</i>	<i>Description</i>
Less than 1	Calm
1–3	Light air
4–27	Breeze
28–47	Gale
48–63	Storm
Above 63	Hurricane

Write a program that asks the user to enter a wind speed (in knots), then displays the corresponding description.

5. In one state, single residents are subject to the following income tax:

<i>Income</i>	<i>Amount of tax</i>
Not over \$750	1% of income
\$750–\$2,250	\$7.50 plus 2% of amount over \$750
\$2,250–\$3,750	\$37.50 plus 3% of amount over \$2,250
\$3,750–\$5,250	\$82.50 plus 4% of amount over \$3,750
\$5,250–\$7,000	\$142.50 plus 5% of amount over \$5,250
Over \$7,000	\$230.00 plus 6% of amount over \$7,000

Write a program that asks the user to enter the amount of taxable income, then displays the tax due.

- W 6. Modify the `upc.c` program of Section 4.1 so that it checks whether a UPC is valid. After the user enters a UPC, the program will display either VALID or NOT VALID.

7. Write a program that finds the largest and smallest of four integers entered by the user:

Enter four integers: 21 43 10 35  
 Largest: 43  
 Smallest: 10

Use as few if statements as possible. Hint: Four if statements are sufficient.

8. The following table shows the daily flights from one city to another:

<i>Departure time</i>	<i>Arrival time</i>
8:00 a.m.	10:16 a.m.
9:43 a.m.	11:52 a.m.
11:19 a.m.	1:31 p.m.
12:47 p.m.	3:00 p.m.

2:00 p.m.	4:08 p.m.
3:45 p.m.	5:55 p.m.
7:00 p.m.	9:20 p.m.
9:45 p.m.	11:58 p.m.

Write a program that asks user to enter a time (expressed in hours and minutes, using the 24-hour clock). The program then displays the departure and arrival times for the flight whose departure time is closest to that entered by the user:

Enter a 24-hour time: 13:15

Closest departure time is 12:47 p.m., arriving at 3:00 p.m.

*Hint:* Convert the input into a time expressed in minutes since midnight, and compare it to the departure times, also expressed in minutes since midnight. For example, 13:15 is  $13 \times 60 + 15 = 795$  minutes since midnight, which is closer to 12:47 p.m. (767 minutes since midnight) than to any of the other departure times.

9. Write a program that prompts the user to enter two dates and then indicates which date comes earlier on the calendar:

Enter first date (mm/dd/yy) : 3/6/08

Enter second date (mm/dd/yy) : 5/17/07

5/17/07 is earlier than 3/6/08

- W 10. Using the `switch` statement, write a program that converts a numerical grade into a letter grade:

Enter numerical grade: 84

Letter grade: B

Use the following grading scale: A = 90–100, B = 80–89, C = 70–79, D = 60–69, F = 0–59. Print an error message if the grade is larger than 100 or less than 0. *Hint:* Break the grade into two digits, then use a `switch` statement to test the ten's digit.

11. Write a program that asks the user for a two-digit number, then prints the English word for the number:

Enter a two-digit number: 45

You entered the number forty-five.

*Hint:* Break the number into two digits. Use one `switch` statement to print the word for the first digit ("twenty," "thirty," and so forth). Use a second `switch` statement to print the word for the second digit. Don't forget that the numbers between 11 and 19 require special treatment.



# 6 Loops

*A program without a loop and a structured variable isn't worth writing.*

Chapter 5 covered C's selection statements, `if` and `switch`. This chapter introduces C's iteration statements, which allow us to set up loops.

A *loop* is a statement whose job is to repeatedly execute some other statement (the *loop body*). In C, every loop has a *controlling expression*. Each time the loop body is executed (an *iteration* of the loop), the controlling expression is evaluated; if the expression is true—has a value that's not zero—the loop continues to execute.

C provides three iteration statements: `while`, `do`, and `for`, which are covered in Sections 6.1, 6.2, and 6.3, respectively. The `while` statement is used for loops whose controlling expression is tested *before* the loop body is executed. The `do` statement is used if the expression is tested *after* the loop body is executed. The `for` statement is convenient for loops that increment or decrement a counting variable. Section 6.3 also introduces the comma operator, which is used primarily in `for` statements.

The last two sections of this chapter are devoted to C features that are used in conjunction with loops. Section 6.4 describes the `break`, `continue`, and `goto` statements. `break` jumps out of a loop and transfers control to the next statement after the loop, `continue` skips the rest of a loop iteration, and `goto` jumps to any statement within a function. Section 6.5 covers the `null` statement, which can be used to create loops with empty bodies.

## 6.1 The `while` Statement

Of all the ways to set up loops in C, the `while` statement is the simplest and most fundamental. The `while` statement has the form

**while statement****while ( expression ) statement**

The expression inside the parentheses is the controlling expression; the statement after the parentheses is the loop body. Here's an example:

```
while (i < n) /* controlling expression */
    i = i * 2; /* loop body */
```

Note that the parentheses are mandatory and that nothing goes between the right parenthesis and the loop body. (Some languages require the word *do*.)

When a *while* statement is executed, the controlling expression is evaluated first. If its value is nonzero (true), the loop body is executed and the expression is tested again. The process continues in this fashion—first testing the controlling expression, then executing the loop body—until the controlling expression eventually has the value zero.

The following example uses a *while* statement to compute the smallest power of 2 that is greater than or equal to a number *n*:

```
i = 1;
while (i < n)
    i = i * 2;
```

Suppose that *n* has the value 10. The following trace shows what happens when the *while* statement is executed:

```
i = 1;      i is now 1.
Is i < n?  Yes; continue.
i = i * 2;  i is now 2.
Is i < n?  Yes; continue.
i = i * 2;  i is now 4.
Is i < n?  Yes; continue.
i = i * 2;  i is now 8.
Is i < n?  Yes; continue.
i = i * 2;  i is now 16.
Is i < n?  No; exit from loop.
```

Notice how the loop keeps going as long as the controlling expression (*i < n*) is true. When the expression is false, the loop terminates, and *i* is greater than or equal to *n*, as desired.

Although the loop body must be a single statement, that's merely a technicality. If we want more than one statement, we can just use braces to create a single compound statement:

```
while (i > 0) {
    printf("T minus %d and counting\n", i);
    i--;
}
```

Some programmers always use braces, even when they're not strictly necessary:

```
while (i < n) { /* braces allowed, but not required */
    i = i * 2;
}
```

As a second example, let's trace the execution of the following statements, which display a series of "countdown" messages:

```
i = 10;
while (i > 0) {
    printf("T minus %d and counting\n", i);
    i--;
}
```

Before the `while` statement is executed, the variable `i` is assigned the value 10. Since 10 is greater than 0, the loop body is executed, causing the message `T minus 10 and counting` to be printed and `i` to be decremented. The condition `i > 0` is then tested again. Since 9 is greater than 0, the loop body is executed once more. This process continues until the message `T minus 1 and counting` is printed and `i` becomes 0. The test `i > 0` then fails, causing the loop to terminate.

The countdown example leads us to make several observations about the `while` statement:

- The controlling expression is false when a `while` loop terminates. Thus, when a loop controlled by the expression `i > 0` terminates, `i` must be less than or equal to 0. (Otherwise, we'd still be executing the loop!)
- The body of a `while` loop may not be executed at all. Since the controlling expression is tested *before* the loop body is executed, it's possible that the body isn't executed even once. If `i` has a negative or zero value when the countdown loop is first entered, the loop will do nothing.
- A `while` statement can often be written in a variety of ways. For example, we could make the countdown loop more concise by decrementing `i` inside the call of `printf`:

#### **Q&A**

```
while (i > 0)
    printf("T minus %d and counting\n", i--);
```

## Infinite Loops

A `while` statement won't terminate if the controlling expression always has a nonzero value. In fact, C programmers sometimes deliberately create an *infinite loop* by using a nonzero constant as the controlling expression:

**idiom** `while (1) ...`

A `while` statement of this form will execute forever unless its body contains a statement that transfers control out of the loop (`break`, `goto`, `return`) or calls a function that causes the program to terminate.

**PROGRAM Printing a Table of Squares**

Let's write a program that prints a table of squares. The program will first prompt the user to enter a number  $n$ . It will then print  $n$  lines of output, with each line containing a number between 1 and  $n$  together with its square:

```
This program prints a table of squares.  
Enter number of entries in table: 5
```

1	1
2	4
3	9
4	16
5	25

Let's have the program store the desired number of squares in a variable named  $n$ . We'll need a loop that repeatedly prints a number  $i$  and its square, starting with  $i$  equal to 1. The loop will repeat as long as  $i$  is less than or equal to  $n$ . We'll have to make sure to add 1 to  $i$  each time through the loop.

We'll write the loop as a `while` statement. (Frankly, we haven't got much choice, since the `while` statement is the only kind of loop we've covered so far.) Here's the finished program:

```
square.c /* Prints a table of squares using a while statement */  
  
#include <stdio.h>  
  
int main(void)  
{  
    int i, n;  
  
    printf("This program prints a table of squares.\n");  
    printf("Enter number of entries in table: ");  
    scanf("%d", &n);  
  
    i = 1;  
    while (i <= n) {  
        printf("%10d%10d\n", i, i * i);  
        i++;  
    }  
  
    return 0;  
}
```

Note how `square.c` displays numbers in neatly aligned columns. The trick is to use a conversion specification like `%10d` instead of just `%d`, taking advantage of the fact that `printf` right-justifies numbers when a field width is specified.

**PROGRAM Summing a Series of Numbers**

As a second example of the `while` statement, let's write a program that sums a series of integers entered by the user. Here's what the user will see:

```
This program sums a series of integers.
Enter integers (0 to terminate): 8 23 71 5 0
The sum is: 107
```

Clearly we'll need a loop that uses `scanf` to read a number and then adds the number to a running total.

Letting `n` represent the number just read and `sum` the total of all numbers previously read, we end up with the following program:

```
sum.c /* Sums a series of numbers */

#include <stdio.h>

int main(void)
{
    int n, sum = 0;

    printf("This program sums a series of integers.\n");
    printf("Enter integers (0 to terminate): ");

    scanf("%d", &n);
    while (n != 0) {
        sum += n;
        scanf("%d", &n);
    }
    printf("The sum is: %d\n", sum);

    return 0;
}
```

Notice that the condition `n != 0` is tested just after a number is read, allowing the loop to terminate as soon as possible. Also note that there are two identical calls of `scanf`, which is often hard to avoid when using `while` loops.

## 6.2 The do Statement

The `do` statement is closely related to the `while` statement; in fact, the `do` statement is essentially just a `while` statement whose controlling expression is tested *after* each execution of the loop body. The `do` statement has the form

### do statement

```
do statement while ( expression ) ;
```

As with the `while` statement, the body of a `do` statement must be one statement (possibly compound, of course) and the controlling expression must be enclosed within parentheses.

When a `do` statement is executed, the loop body is executed first, then the controlling expression is evaluated. If the value of the expression is nonzero, the loop

body is executed again and then the expression is evaluated once more. Execution of the `do` statement terminates when the controlling expression has the value 0 *after* the loop body has been executed.

Let's rewrite the countdown example of Section 6.1, using a `do` statement this time:

```
i = 10;
do {
    printf("T minus %d and counting\n", i);
    --i;
} while (i > 0);
```

When the `do` statement is executed, the loop body is first executed, causing the message `T minus 10 and counting` to be printed and `i` to be decremented. The condition `i > 0` is now tested. Since 9 is greater than 0, the loop body is executed a second time. This process continues until the message `T minus 1 and counting` is printed and `i` becomes 0. The test `i > 0` now fails, causing the loop to terminate. As this example shows, the `do` statement is often indistinguishable from the `while` statement. The difference between the two is that the body of a `do` statement is always executed at least once; the body of a `while` statement is skipped entirely if the controlling expression is 0 initially.

Incidentally, it's a good idea to use braces in *all* `do` statements, whether or not they're needed, because a `do` statement without braces can easily be mistaken for a `while` statement:

```
do
    printf("T minus %d and counting\n", i--);
while (i > 0);
```

A careless reader might think that the word `while` was the beginning of a `while` statement.

## PROGRAM Calculating the Number of Digits in an Integer

Although the `while` statement appears in C programs much more often than the `do` statement, the latter is handy for loops that must execute at least once. To illustrate this point, let's write a program that calculates the number of digits in an integer entered by the user:

```
Enter a nonnegative integer: 60
The number has 2 digit(s).
```

Our strategy will be to divide the user's input by 10 repeatedly until it becomes 0; the number of divisions performed is the number of digits. Clearly we'll need some kind of loop, since we don't know how many divisions it will take to reach 0. But should we use a `while` statement or a `do` statement? The `do` statement turns out to be more attractive, because every integer—even 0—has at least *one* digit. Here's the program:

```
numdigits.c /* Calculates the number of digits in an integer */

#include <stdio.h>

int main(void)
{
    int digits = 0, n;

    printf("Enter a nonnegative integer: ");
    scanf("%d", &n);

    do {
        n /= 10;
        digits++;
    } while (n > 0);

    printf("The number has %d digit(s).\n", digits);

    return 0;
}
```

To see why the `do` statement is the right choice, let's see what would happen if we were to replace the `do` loop by a similar `while` loop:

```
while (n > 0) {
    n /= 10;
    digits++;
}
```

If `n` is 0 initially, this loop won't execute at all, and the program would print  
 The number has 0 digit(s).

## 6.3 The for Statement

We now come to the last of C's loops: the `for` statement. Don't be discouraged by the `for` statement's apparent complexity; it's actually the best way to write many loops. The `for` statement is ideal for loops that have a "counting" variable, but it's versatile enough to be used for other kinds of loops as well.

The `for` statement has the form

**for statement**

`for ( expr1 ; expr2 ; expr3 ) statement`

where `expr1`, `expr2`, and `expr3` are expressions. Here's an example:

```
for (i = 10; i > 0; i--)
    printf("T minus %d and counting\n", i);
```

When this `for` statement is executed, the variable `i` is initialized to 10, then `i` is tested to see if it's greater than 0. Since it is, the message T minus 10 and

counting is printed, then *i* is decremented. The condition *i* > 0 is then tested again. The loop body will be executed 10 times in all, with *i* varying from 10 down to 1.

**Q&A** The `for` statement is closely related to the `while` statement. In fact, except in a few rare cases, a `for` loop can always be replaced by an equivalent `while` loop:

```
expr1;
while ( expr2 ) {
    statement
    expr3;
}
```

As this pattern shows, *expr1* is an initialization step that's performed only once, before the loop begins to execute, *expr2* controls loop termination (the loop continues executing as long as the value of *expr2* is nonzero), and *expr3* is an operation to be performed at the end of each loop iteration. Applying this pattern to our previous `for` loop example, we arrive at the following:

```
i = 10;
while (i > 0) {
    printf("T minus %d and counting\n", i);
    i--;
}
```

Studying the equivalent `while` statement can help us understand the fine points of a `for` statement. For example, suppose that we replace `i--` by `--i` in our `for` loop example:

```
for (i = 10; i > 0; --i)
    printf("T minus %d and counting\n", i);
```

How does this change affect the loop? Looking at the equivalent `while` loop, we see that it has no effect:

```
i = 10;
while (i > 0) {
    printf("T minus %d and counting\n", i);
    --i;
}
```

Since the first and third expressions in a `for` statement are executed as statements, their values are irrelevant—they're useful only for their side effects. Consequently, these two expressions are usually assignments or increment/decrement expressions.

## for Statement Idioms

The `for` statement is usually the best choice for loops that “count up” (increment a variable) or “count down” (decrement a variable). A `for` statement that counts up or down a total of *n* times will usually have one of the following forms:

■ ***Counting up from 0 to n-1:***

**idiom**      `for (i = 0; i < n; i++) ...`

■ ***Counting up from 1 to n:***

**idiom**      `for (i = 1; i <= n; i++) ...`

■ ***Counting down from n-1 to 0:***

**idiom**      `for (i = n - 1; i >= 0; i--) ...`

■ ***Counting down from n to 1:***

**idiom**      `for (i = n; i > 0; i--) ...`

Imitating these patterns will help you avoid some of the following errors, which beginning C programmers often make:

- Using `<` instead of `>` (or vice versa) in the controlling expression. Notice that “counting up” loops use the `<` or `<=` operator, while “counting down” loops rely on `>` or `>=`.
- Using `==` in the controlling expression instead of `<`, `<=`, `>`, or `>=`. A controlling expression needs to be true at the beginning of the loop, then later become false so that the loop can terminate. A test such as `i == n` doesn’t make much sense, because it won’t be true initially.
- “Off-by-one” errors such as writing the controlling expression as `i <= n` instead of `i < n`.

## Omitting Expressions in a `for` Statement

The `for` statement is even more flexible than we’ve seen so far. Some `for` loops may not need all three of the expressions that normally control the loop, so C allows us to omit any or all of the expressions.

If the *first* expression is omitted, no initialization is performed before the loop is executed:

```
i = 10;
for (; i > 0; --i)
    printf("T minus %d and counting\n", i);
```

In this example, `i` has been initialized by a separate assignment, so we’ve omitted the first expression in the `for` statement. (Notice that the semicolon between the first and second expressions remains. The two semicolons must always be present, even when we’ve omitted some of the expressions.)

If we omit the *third* expression in a `for` statement, the loop body is responsible for ensuring that the value of the second expression eventually becomes false. Our `for` statement example could be written like this:

```
for (i = 10; i > 0;)
    printf("T minus %d and counting\n", i--);
```

To compensate for omitting the third expression, we've arranged for *i* to be decremented inside the loop body.

When the *first* and *third* expressions are both omitted, the resulting loop is nothing more than a `while` statement in disguise. For example, the loop

```
for (; i > 0;)
    printf("T minus %d and counting\n", i--);
```

is the same as

```
while (i > 0)
    printf("T minus %d and counting\n", i--);
```

The `while` version is clearer and therefore preferable.

If the *second* expression is missing, it defaults to a true value, so the `for` statement doesn't terminate (unless stopped in some other fashion). For example, some programmers use the following `for` statement to establish an infinite loop:

**idiom** `for ( ; ) ...`

### C99 for Statements in C99

In C99, the first expression in a `for` statement can be replaced by a declaration. This feature allows the programmer to declare a variable for use by the loop:

```
for (int i = 0; i < n; i++)
    ...
```

The variable *i* need not have been declared prior to this statement. (In fact, if a declaration of *i* already exists, this statement creates a *new* version of *i* that will be used solely within the loop.)

A variable declared by a `for` statement can't be accessed outside the body of the loop (we say that it's not *visible* outside the loop):

```
for (int i = 0; i < n; i++) {
    ...
    printf("%d", i); /* legal; i is visible inside loop */
    ...
}
printf("%d", i); /* *** WRONG *** /
```

Having a `for` statement declare its own control variable is usually a good idea: it's convenient and it can make programs easier to understand. However, if the program needs to access the variable after loop termination, it's necessary to use the older form of the `for` statement.

Incidentally, a `for` statement may declare more than one variable, provided that all variables have the same type:

```
for (int i = 0, j = 0; i < n; i++)
    ...
```

## The Comma Operator

On occasion, we might like to write a `for` statement with two (or more) initialization expressions or one that increments several variables each time through the loop. We can do this by using a *comma expression* as the first or third expression in the `for` statement.

A comma expression has the form

**comma expression**

*expr1 , expr2*

where *expr1* and *expr2* are any two expressions. A comma expression is evaluated in two steps: First, *expr1* is evaluated and its value discarded. Second, *expr2* is evaluated; its value is the value of the entire expression. Evaluating *expr1* should always have a side effect; if it doesn't, then *expr1* serves no purpose.

For example, suppose that *i* and *j* have the values 1 and 5, respectively. When the comma expression `++i, i + j` is evaluated, *i* is first incremented, then *i + j* is evaluated, so the value of the expression is 7. (And, of course, *i* now has the value 2.) The precedence of the comma operator is less than that of all other operators, by the way, so there's no need to put parentheses around `++i` and `i + j`.

Occasionally, we'll need to chain together a series of comma expressions, just as we sometimes chain assignments together. The comma operator is left associative, so the compiler interprets

`i = 1, j = 2, k = i + j`

as

`((i = 1), (j = 2)), (k = (i + j))`

Since the left operand in a comma expression is evaluated before the right operand, the assignments *i = 1*, *j = 2*, and *k = i + j* will be performed from left to right.

The comma operator is provided for situations where C requires a single expression, but we'd like to have two or more expressions. In other words, the comma operator allows us to "glue" two expressions together to form a single expression. (Note the similarity to the compound statement, which allows us to treat a group of statements as a single statement.)

The need to glue expressions together doesn't arise that often. Certain macro definitions can benefit from the comma operator, as we'll see in a later chapter. The `for` statement is the only other place where the comma operator is likely to be found. For example, suppose that we want to initialize two variables when entering a `for` statement. Instead of writing

```
sum = 0;
for (i = 1; i <= N; i++)
    sum += i;
```

we can write

```
for (sum = 0, i = 1; i <= N; i++)
    sum += i;
```

The expression `sum = 0, i = 1` first assigns 0 to `sum`, then assigns 1 to `i`. With additional commas, the `for` statement could initialize more than two variables.

## PROGRAM Printing a Table of Squares (Revisited)

The `square.c` program (Section 6.1) can be improved by converting its while loop to a for loop:

```
square2.c /* Prints a table of squares using a for statement */

#include <stdio.h>

int main(void)
{
    int i, n;

    printf("This program prints a table of squares.\n");
    printf("Enter number of entries in table: ");
    scanf("%d", &n);

    for (i = 1; i <= n; i++)
        printf("%10d%10d\n", i, i * i);

    return 0;
}
```

We can use this program to illustrate an important point about the `for` statement: C places no restrictions on the three expressions that control its behavior. Although these expressions usually initialize, test, and update the same variable, there's no requirement that they be related in any way. Consider the following version of the same program:

```
square3.c /* Prints a table of squares using an odd method */

#include <stdio.h>

int main(void)
{
    int i, n, odd, square;

    printf("This program prints a table of squares.\n");
    printf("Enter number of entries in table: ");
    scanf("%d", &n);

    i = 1;
    odd = 3;
    for (square = 1; i <= n; odd += 2) {
        printf("%10d%10d\n", i, square);
        ++i;
    }
}
```

```

    square += odd;
}

return 0;
}

```

The `for` statement in this program initializes one variable (`square`), tests another (`i`), and increments a third (`odd`). `i` is the number to be squared, `square` is the square of `i`, and `odd` is the odd number that must be added to the current square to get the next square (allowing the program to compute consecutive squares without performing any multiplications).

linked lists ▶ 17.5

The tremendous flexibility of the `for` statement can sometimes be useful; we'll find it to be a great help when working with linked lists. The `for` statement can easily be misused, though, so don't go overboard. The `for` loop in `square3.c` would be a lot clearer if we rearranged its pieces so that the loop is clearly controlled by `i`.

## 6.4 Exiting from a Loop

We've seen how to write loops that have an exit point before the loop body (using `while` and `for` statements) or after it (using `do` statements). Occasionally, however, we'll need a loop with an exit point in the middle. We may even want a loop to have more than one exit point. The `break` statement makes it possible to write either kind of loop.

After we've examined the `break` statement, we'll look at a couple of related statements: `continue` and `goto`. The `continue` statement makes it possible to skip part of a loop iteration without jumping out of the loop. The `goto` statement allows a program to jump from one statement to another. Thanks to the availability of statements such as `break` and `continue`, the `goto` statement is rarely used.

### The `break` Statement

We've already discussed how a `break` statement can transfer control out of a `switch` statement. The `break` statement can also be used to jump out of a `while`, `do`, or `for` loop.

Suppose that we're writing a program that checks whether a number `n` is prime. Our plan is to write a `for` statement that divides `n` by the numbers between 2 and `n - 1`. We should break out of the loop as soon as any divisor is found; there's no need to try the remaining possibilities. After the loop has terminated, we can use an `if` statement to determine whether termination was premature (hence `n` isn't prime) or normal (`n` is prime):

```

for (d = 2; d < n; d++)
    if (n % d == 0)
        break;

```

```

if (d < n)
    printf("%d is divisible by %d\n", n, d);
else
    printf("%d is prime\n", n);

```

The `break` statement is particularly useful for writing loops in which the exit point is in the middle of the body rather than at the beginning or end. Loops that read user input, terminating when a particular value is entered, often fall into this category:

```

for (;;) {
    printf("Enter a number (enter 0 to stop): ");
    scanf("%d", &n);
    if (n == 0)
        break;
    printf("%d cubed is %d\n", n, n * n * n);
}

```

A `break` statement transfers control out of the *innermost* enclosing `while`, `do`, `for`, or `switch` statement. Thus, when these statements are nested, the `break` statement can escape only one level of nesting. Consider the case of a `switch` statement nested inside a `while` statement:

```

while (...) {
    switch (...) {
        ...
        break;
        ...
    }
}

```

The `break` statement transfers control out of the `switch` statement, but not out of the `while` loop. I'll return to this point later.

## The `continue` Statement

The `continue` statement doesn't really belong here, because it doesn't exit from a loop. It's similar to `break`, though, so its inclusion in this section isn't completely arbitrary. `break` transfers control just *past* the end of a loop, while `continue` transfers control to a point just *before* the end of the loop body. With `break`, control leaves the loop; with `continue`, control remains inside the loop. There's another difference between `break` and `continue`: `break` can be used in `switch` statements and loops (`while`, `do`, and `for`), whereas `continue` is limited to loops.

The following example, which reads a series of numbers and computes their sum, illustrates a simple use of `continue`. The loop terminates when 10 nonzero numbers have been read. Whenever the number 0 is read, the `continue` statement is executed, skipping the rest of the loop body (the statements `sum += i;` and `n++;`) but remaining inside the loop.

```
n = 0;
sum = 0;
while (n < 10) {
    scanf("%d", &i);
    if (i == 0)
        continue;
    sum += i;
    n++;
    /* continue jumps to here */
}
```

If `continue` were not available, we could have written the example as follows:

```
n = 0;
sum = 0;
while (n < 10) {
    scanf("%d", &i);
    if (i != 0) {
        sum += i;
        n++;
    }
}
```

## The `goto` Statement

`break` and `continue` are jump statements that transfer control from one point in the program to another. Both are restricted: the target of a `break` is a point just *beyond* the end of the enclosing loop, while the target of a `continue` is a point just *before* the end of the loop. The `goto` statement, on the other hand, is capable of jumping to *any* statement in a function, provided that the statement has a *label*. (C99 places an additional restriction on the `goto` statement: it can't be used to bypass the declaration of a variable-length array.)

C99

variable-length arrays ▶ 8.3

A label is just an identifier placed at the beginning of a statement:

**labeled statement**

*identifier* : *statement*

A statement may have more than one label. The `goto` statement itself has the form

**goto statement**

`goto` *identifier* ;

Executing the statement `goto L;` transfers control to the statement that follows the label *L*, which must be in the same function as the `goto` statement itself.

If C didn't have a `break` statement, here's how we might use a `goto` statement to exit prematurely from a loop:

```
for (d = 2; d < n; d++)
    if (n % d == 0)
        goto done;
```

```

done:
if (d < n)
    printf("%d is divisible by %d\n", n, d);
else
    printf("%d is prime\n", n);

```

**Q&A**  
exit function ►9.5

The `goto` statement, a staple of older programming languages, is rarely needed in everyday C programming. The `break`, `continue`, and `return` statements—which are essentially restricted `goto` statements—and the `exit` function are sufficient to handle most situations that might require a `goto` in other languages.

Nonetheless, the `goto` statement can be helpful once in a while. Consider the problem of exiting a loop from within a `switch` statement. As we saw earlier, the `break` statement doesn't quite have the desired effect: it exits from the `switch`, but not from the loop. A `goto` statement solves the problem:

```

while (...) {
    switch (...) {
        ...
        goto loop_done; /* break won't work here */
        ...
    }
loop_done: ...

```

The `goto` statement is also useful for exiting from nested loops.

## PROGRAM Balancing a Checkbook

Many simple interactive programs are menu-based: they present the user with a list of commands to choose from. Once the user has selected a command, the program performs the desired action, then prompts the user for another command. This process continues until the user selects an “exit” or “quit” command.

The heart of such a program will obviously be a loop. Inside the loop will be statements that prompt the user for a command, read the command, then decide what action to take:

```

for (;;) {
    prompt user to enter command;
    read command;
    execute command;
}

```

Executing the command will require a `switch` statement (or cascaded `if` statement):

```

for (;;) {
    prompt user to enter command;
    read command;
    switch (command) {
        case command1: perform operation1; break;
}

```

```

    case command2: perform operation2; break;
    :
    :
    case commandn: perform operationn; break;
    default: print error message; break;
}
}

```

To illustrate this arrangement, let's develop a program that maintains a checkbook balance. The program will offer the user a menu of choices: clear the account balance, credit money to the account, debit money from the account, display the current balance, and exit the program. The choices are represented by the integers 0, 1, 2, 3, and 4, respectively. Here's what a session with the program will look like:

```

*** ACME checkbook-balancing program ***
Commands: 0=clear, 1=credit, 2=debit, 3=balance, 4=exit

Enter command: 1
Enter amount of credit: 1042.56
Enter command: 2
Enter amount of debit: 133.79
Enter command: 1
Enter amount of credit: 1754.32
Enter command: 2
Enter amount of debit: 1400
Enter command: 2
Enter amount of debit: 68
Enter command: 2
Enter amount of debit: 50
Enter command: 3
Current balance: $1145.09
Enter command: 4

```

When the user enters the command 4 (exit), the program needs to exit from the switch statement *and* the surrounding loop. The break statement won't help, and we'd prefer not to use a goto statement. Instead, we'll have the program execute a return statement, which will cause the main function to return to the operating system.

```

checking.c /* Balances a checkbook */

#include <stdio.h>

int main(void)
{
    int cmd;
    float balance = 0.0f, credit, debit;

    printf(" *** ACME checkbook-balancing program ***\n");
    printf("Commands: 0=clear, 1=credit, 2=debit, ");
    printf("3=balance, 4=exit\n\n");

```

```

for (;;) {
    printf("Enter command: ");
    scanf("%d", &cmd);
    switch (cmd) {
        case 0:
            balance = 0.0f;
            break;
        case 1:
            printf("Enter amount of credit: ");
            scanf("%f", &credit);
            balance += credit;
            break;
        case 2:
            printf("Enter amount of debit: ");
            scanf("%f", &debit);
            balance -= debit;
            break;
        case 3:
            printf("Current balance: $%.2f\n", balance);
            break;
        case 4:
            return 0;
        default:
            printf("Commands: 0=clear, 1=credit, 2=debit, ");
            printf("3=balance, 4=exit\n\n");
            break;
    }
}
}

```

Note that the `return` statement is not followed by a `break` statement. A `break` immediately following a `return` can never be executed, and many compilers will issue a warning message.

## 6.5 The Null Statement

A statement can be **null**—devoid of symbols except for the semicolon at the end. Here's an example:

```
i = 0; ; j = 1;
```

This line contains three statements: an assignment to `i`, a null statement, and an assignment to `j`.

### Q&A

The null statement is primarily good for one thing: writing loops whose bodies are empty. As an example, recall the prime-finding loop of Section 6.4:

```

for (d = 2; d < n; d++)
    if (n % d == 0)
        break;

```

If we move the `n % d == 0` condition into the loop's controlling expression, the body of the loop becomes empty:

```
for (d = 2; d < n && n % d != 0; d++)
    /* empty loop body */ ;
```

Each time through the loop, the condition `d < n` is tested first; if it's false, the loop terminates. Otherwise, the condition `n % d != 0` is tested, and if that's false, the loop terminates. (In the latter case, `n % d == 0` must be true; in other words, we've found a divisor of `n`.)

Note how we've put the null statement on a line by itself, instead of writing

```
for (d = 2; d < n && n % d != 0; d++) ;
```

### Q&A

C programmers customarily put the null statement on a line by itself. Otherwise, someone reading the program might get confused about whether the statement after the `for` was actually its body:

```
for (d = 2; d < n && n % d != 0; d++);
if (d < n)
    printf("%d is divisible by %d\n", n, d);
```

Converting an ordinary loop into one with an empty body doesn't buy much: the new loop is often more concise but usually no more efficient. In a few cases, though, a loop with an empty body is clearly superior to the alternatives. For example, we'll find these loops to be handy for reading character data.



Accidentally putting a semicolon after the parentheses in an `if`, `while`, or `for` statement creates a null statement, thus ending the `if`, `while`, or `for` prematurely.

- In an `if` statement, putting a semicolon after the parentheses creates an `if` statement that apparently performs the same action regardless of the value of its controlling expression:

```
if (d == 0);                                /*** WRONG ***/
    printf("Error: Division by zero\n");
```

The call of `printf` isn't inside the `if` statement, so it's performed regardless of whether `d` is equal to 0.

- In a `while` statement, putting a semicolon after the parentheses may create an infinite loop:

```
i = 10;
while (i > 0);                                /*** WRONG ***/
{
    printf("T minus %d and counting\n", i);
    --i;
}
```

Another possibility is that the loop terminates, but the statement that should be the loop body is executed only once, after the loop has terminated:

```
i = 11;
while (--i > 0);                                /*** WRONG ***
    printf("T minus %d and counting\n", i);
```

This example prints the message

T minus 0 and counting

- In a `for` statement, putting a semicolon after the parentheses causes the statement that should be the loop body to be executed only once:

```
for (i = 10; i > 0; i--);                      /*** WRONG ***
    printf("T minus %d and counting\n", i);
```

This example also prints the message

T minus 0 and counting

---

## Q & A

**Q: The following loop appears in Section 6.1:**

```
while (i > 0)
    printf("T minus %d and counting\n", i--);
```

**Why not shorten the loop even more by removing the “> 0” test?**

```
while (i)
    printf("T minus %d and counting\n", i--);
```

**This version will stop when `i` reaches 0, so it should be just as good as the original. [p. 101]**

**A:** The new version is certainly more concise, and many C programmers would write the loop in just this way. It does have drawbacks, though.

First, the new loop is not as easy to read as the original. It’s clear that the loop will terminate when `i` reaches 0, but it’s not obvious whether we’re counting up or down. In the original loop, that information can be deduced from the controlling expression, `i > 0`.

Second, the new loop behaves differently than the original if `i` should happen to have a negative value when the loop begins to execute. The original loop terminates immediately, but the new loop doesn’t.

**Q: Section 6.3 says that, except in rare cases, `for` loops can be converted to `while` loops using a standard pattern. Can you give an example of such a case? [p. 106]**

- A: When the body of a `for` loop contains a `continue` statement, the `while` pattern shown in Section 6.3 is no longer valid. Consider the following example from Section 6.4:

```
n = 0;
sum = 0;
while (n < 10) {
    scanf("%d", &i);
    if (i == 0)
        continue;
    sum += i;
    n++;
}
```

At first glance, it looks as though we could convert the `while` loop into a `for` loop:

```
sum = 0;
for (n = 0; n < 10; n++) {
    scanf("%d", &i);
    if (i == 0)
        continue;
    sum += i;
}
```

Unfortunately, this loop isn't equivalent to the original. When `i` is equal to 0, the original loop doesn't increment `n`, but the new loop does.

- Q: Which form of infinite loop is preferable, `while (1)` or `for (;;)`? [p. 108]**

- A: C programmers have traditionally preferred `for (;;)` for reasons of efficiency; older compilers would often force programs to test the `1` condition each time through the `while` loop. With modern compilers, however, there should be no difference in performance.

- Q: I've heard that programmers should never use the `continue` statement. Is this true?**

- A: It's true that `continue` statements are rare. Still, `continue` is handy once in a while. Suppose we're writing a loop that reads some input data, checks that it's valid, and, if so, processes the input in some way. If there are a number of validity tests, or if they're complex, `continue` can be helpful. The loop would look something like this:

```
for (;;) {
    read data;
    if (data fails first test)
        continue;
    if (data fails second test)
        continue;
    .
    .
}
```

```

if (data fails last test)
    continue;
process data;
}

```

**Q: What's so bad about the `goto` statement? [p. 114]**

A: The `goto` statement isn't inherently evil; it's just that we usually have better alternatives. Programs that use more than a few `goto` statements can quickly degenerate into "spaghetti code," with control blithely jumping from here to there. Spaghetti code is hard to understand and hard to modify.

`goto` statements make programs hard to read because they can jump either forward or backward. (In contrast, `break` and `continue` only jump forward.) A program that contains `goto` statements often requires the reader to jump back and forth in an attempt to follow the flow of control.

`goto` statements can make programs hard to modify, since they make it possible for a section of code to serve more than one purpose. For example, a statement that is preceded by a label might be reachable either by "falling through" from the previous statement or by executing one of several `goto` statements.

**Q: Does the null statement have any uses besides indicating that the body of a loop is empty? [p. 116]**

A: Very few. Since the null statement can appear wherever a statement is allowed, there are many *potential* uses for the null statement. In practice, however, there's only one other use of the null statement, and it's rare.

Suppose that we need to put a label at the end of a compound statement. A label can't stand alone; it must always be followed by a statement. Putting a null statement after the label solves the problem:

```

{
    ...
    goto end_of_stmt;
    ...
    end_of_stmt: ;
}

```

**Q: Are there any other ways to make an empty loop body stand out besides putting the null statement on a line by itself? [p. 117]**

A: Some programmers use a dummy `continue` statement:

```

for (d = 2; d < n && n % d != 0; d++)
    continue;

```

Others use an empty compound statement:

```

for (d = 2; d < n && n % d != 0; d++)
    {}

```

## Exercises

**Section 6.1**

- What output does the following program fragment produce?

```
i = 1;
while (i <= 128) {
    printf("%d ", i);
    i *= 2;
}
```

**Section 6.2**

- What output does the following program fragment produce?

```
i = 9384;
do {
    printf("%d ", i);
    i /= 10;
} while (i > 0);
```

**Section 6.3**

- What output does the following `for` statement produce?

```
for (i = 5, j = i - 1; i > 0, j > 0; --i, j = i - 1)
    printf("%d ", i);
```

- W 4. Which one of the following statements is not equivalent to the other two (assuming that the loop bodies are the same)?
- `for (i = 0; i < 10; i++) ...`
  - `for (i = 0; i < 10; ++i) ...`
  - `for (i = 0; i++ < 10; ) ...`
5. Which one of the following statements is not equivalent to the other two (assuming that the loop bodies are the same)?
- `while (i < 10) {...}`
  - `for (; i < 10;) {...}`
  - `do {...} while (i < 10);`
6. Translate the program fragment of Exercise 1 into a single `for` statement.
7. Translate the program fragment of Exercise 2 into a single `for` statement.
- \*8. What output does the following `for` statement produce?
- ```
for (i = 10; i >= 1; i /= 2)
    printf("%d ", i++);
```
9. Translate the `for` statement of Exercise 8 into an equivalent `while` statement. You will need one statement in addition to the `while` loop itself.

**Section 6.4**

- W 10. Show how to replace a `continue` statement by an equivalent `goto` statement.
11. What output does the following program fragment produce?

```

sum = 0;
for (i = 0; i < 10; i++) {
    if (i % 2)
        continue;
    sum += i;
}
printf("%d\n", sum);

```

- W 12. The following “prime-testing” loop appeared in Section 6.4 as an example:

```

for (d = 2; d < n; d++)
    if (n % d == 0)
        break;

```

This loop isn’t very efficient. It’s not necessary to divide  $n$  by all numbers between 2 and  $n - 1$  to determine whether it’s prime. In fact, we need only check divisors up to the square root of  $n$ . Modify the loop to take advantage of this fact. Hint: Don’t try to compute the square root of  $n$ ; instead, compare  $d * d$  with  $n$ .

### Section 6.5

- \*13. Rewrite the following loop so that its body is empty:

```

for (n = 0; m > 0; n++)
    m /= 2;

```

- W\*14. Find the error in the following program fragment and fix it.

```

if (n % 2 == 0);
    printf("n is even\n");

```

## Programming Projects

1. Write a program that finds the largest in a series of numbers entered by the user. The program must prompt the user to enter numbers one by one. When the user enters 0 or a negative number, the program must display the largest nonnegative number entered:

```

Enter a number: 60
Enter a number: 38.3
Enter a number: 4.89
Enter a number: 100.62
Enter a number: 75.2295
Enter a number: 0

```

The largest number entered was 100.62

Notice that the numbers aren’t necessarily integers.

- W 2. Write a program that asks the user to enter two integers, then calculates and displays their greatest common divisor (GCD):

```

Enter two integers: 12 28
Greatest common divisor: 4

```

*Hint:* The classic algorithm for computing the GCD, known as Euclid’s algorithm, goes as follows: Let  $m$  and  $n$  be variables containing the two numbers. If  $n$  is 0, then stop:  $m$  contains the GCD. Otherwise, compute the remainder when  $m$  is divided by  $n$ . Copy  $n$  into  $m$  and copy the remainder into  $n$ . Then repeat the process, starting with testing whether  $n$  is 0.

3. Write a program that asks the user to enter a fraction, then reduces the fraction to lowest terms:

Enter a fraction: 6/12  
In lowest terms: 1/2

*Hint:* To reduce a fraction to lowest terms, first compute the GCD of the numerator and denominator. Then divide both the numerator and denominator by the GCD.

- W 4. Add a loop to the `broker.c` program of Section 5.2 so that the user can enter more than one trade and the program will calculate the commission on each. The program should terminate when the user enters 0 as the trade value:

Enter value of trade: 30000

Commission: \$166.00

Enter value of trade: 20000

Commission: \$144.00

Enter value of trade: 0

5. Programming Project 1 in Chapter 4 asked you to write a program that displays a two-digit number with its digits reversed. Generalize the program so that the number can have one, two, three, or more digits. *Hint:* Use a `do` loop that repeatedly divides the number by 10, stopping when it reaches 0.

- W 6. Write a program that prompts the user to enter a number  $n$ , then prints all even squares between 1 and  $n$ . For example, if the user enters 100, the program should print the following:

4  
16  
36  
64  
100

7. Rearrange the `square3.c` program so that the `for` loop initializes  $i$ , tests  $i$ , and increments  $i$ . Don't rewrite the program; in particular, don't use any multiplications.

- W 8. Write a program that prints a one-month calendar. The user specifies the number of days in the month and the day of the week on which the month begins:

Enter number of days in month: 31  
Enter starting day of the week (1=Sun, 7=Sat) : 3

|    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|
| 1  | 2  | 3  | 4  | 5  |    |    |
| 6  | 7  | 8  | 9  | 10 | 11 | 12 |
| 13 | 14 | 15 | 16 | 17 | 18 | 19 |
| 20 | 21 | 22 | 23 | 24 | 25 | 26 |
| 27 | 28 | 29 | 30 | 31 |    |    |

*Hint:* This program isn't as hard as it looks. The most important part is a `for` statement that uses a variable  $i$  to count from 1 to  $n$ , where  $n$  is the number of days in the month, printing each value of  $i$ . Inside the loop, an `if` statement tests whether  $i$  is the last day in a week; if so, it prints a new-line character.

9. Programming Project 8 in Chapter 2 asked you to write a program that calculates the remaining balance on a loan after the first, second, and third monthly payments. Modify the program so that it also asks the user to enter the number of payments and then displays the balance remaining after each of these payments.

10. Programming Project 9 in Chapter 5 asked you to write a program that determines which of two dates comes earlier on the calendar. Generalize the program so that the user may enter any number of dates. The user will enter 0/0/0 to indicate that no more dates will be entered:

```
Enter a date (mm/dd/yy) : 3/6/08
Enter a date (mm/dd/yy) : 5/17/07
Enter a date (mm/dd/yy) : 6/3/07
Enter a date (mm/dd/yy) : 0/0/0
5/17/07 is the earliest date
```

11. The value of the mathematical constant  $e$  can be expressed as an infinite series:

$$e = 1 + 1/1! + 1/2! + 1/3! + \dots$$

Write a program that approximates  $e$  by computing the value of

$$1 + 1/1! + 1/2! + 1/3! + \dots + 1/n!$$

where  $n$  is an integer entered by the user.

12. Modify Programming Project 11 so that the program continues adding terms until the current term becomes less than  $\epsilon$ , where  $\epsilon$  is a small (floating-point) number entered by the user.

## Programming Projects

# 7 Basic Types

*Make no mistake about it: Computers process numbers—not symbols. We measure our understanding (and control) by the extent to which we can arithmetize an activity.*

So far, we've used only two of C's *basic* (built-in) *types*: `int` and `float`. (We've also seen `_Bool`, which is a basic type in C99.) This chapter describes the rest of the basic types and discusses important issues about types in general. Section 7.1 reveals the full range of integer types, which include long integers, short integers, and unsigned integers. Section 7.2 introduces the `double` and `long double` types, which provide a larger range of values and greater precision than `float`. Section 7.3 covers the `char` type, which we'll need in order to work with character data. Section 7.4 tackles the thorny topic of converting a value of one type to an equivalent value of another. Section 7.5 shows how to use `typedef` to define new type names. Finally, Section 7.6 describes the `sizeof` operator, which measures the amount of storage required for a type.

## 7.1 Integer Types

C supports two fundamentally different kinds of numeric types: integer types and floating types. Values of an *integer type* are whole numbers, while values of a floating type can have a fractional part as well. The integer types, in turn, are divided into two categories: signed and unsigned.

### Signed and Unsigned Integers

The leftmost bit of a *signed* integer (known as the *sign bit*) is 0 if the number is positive or zero, 1 if it's negative. Thus, the largest 16-bit integer has the binary representation

0111111111111111

which has the value 32,767 ( $2^{15} - 1$ ). The largest 32-bit integer is

01111111111111111111111111111111

which has the value 2,147,483,647 ( $2^{31} - 1$ ). An integer with no sign bit (the leftmost bit is considered part of the number's magnitude) is said to be ***unsigned***. The largest 16-bit unsigned integer is 65,535 ( $2^{16} - 1$ ), and the largest 32-bit unsigned integer is 4,294,967,295 ( $2^{32} - 1$ ).

By default, integer variables are signed in C—the leftmost bit is reserved for the sign. To tell the compiler that a variable has no sign bit, we declare it to be ***unsigned***. Unsigned numbers are primarily useful for systems programming and low-level, machine-dependent applications. We'll discuss typical applications for unsigned numbers in Chapter 20; until then, we'll generally avoid them.

---

C's integer types come in different sizes. The ***int*** type is usually 32 bits, but may be 16 bits on older CPUs. Since some programs require numbers that are too large to store in ***int*** form, C also provides ***long*** integers. At times, we may need to conserve memory by instructing the compiler to store a number in less space than normal; such a number is called a ***short*** integer.

To construct an integer type that exactly meets our needs, we can specify that a variable is ***long*** or ***short***, ***signed*** or ***unsigned***. We can even combine specifiers (e.g., ***long unsigned int***). However, only the following six combinations actually produce different types:

```
short int  
unsigned short int  
  
int  
unsigned int  
  
long int  
unsigned long int
```

Other combinations are synonyms for one of these six types. (For example, ***long signed int*** is the same as ***long int***, since integers are always signed unless otherwise specified.) Incidentally, the order of the specifiers doesn't matter; ***unsigned short int*** is the same as ***short unsigned int***.

C allows us to abbreviate the names of integer types by dropping the word ***int***. For example, ***unsigned short int*** may be abbreviated to ***unsigned short***, and ***long int*** may be abbreviated to just ***long***. Omitting ***int*** is a widespread practice among C programmers, and some newer C-based languages (including Java) actually require the programmer to write ***short*** or ***long*** rather than ***short int*** or ***long int***. For these reasons, I'll often omit the word ***int*** when it's not strictly necessary.

The range of values represented by each of the six integer types varies from one machine to another. However, there are a couple of rules that all compilers must obey. First, the C standard requires that `short int`, `int`, and `long int` each cover a certain minimum range of values (see Section 23.2 for details). Second, the standard requires that `int` not be shorter than `short int`, and `long int` not be shorter than `int`. However, it's possible that `short int` represents the same range of values as `int`; also, `int` may have the same range as `long int`.

Table 7.1 shows the usual range of values for the integer types on a 16-bit machine; note that `short int` and `int` have identical ranges.

**Table 7.1**  
Integer Types on a  
16-bit Machine

| Type                            | Smallest Value | Largest Value |
|---------------------------------|----------------|---------------|
| <code>short int</code>          | -32,768        | 32,767        |
| <code>unsigned short int</code> | 0              | 65,535        |
| <code>int</code>                | -32,768        | 32,767        |
| <code>unsigned int</code>       | 0              | 65,535        |
| <code>long int</code>           | -2,147,483,648 | 2,147,483,647 |
| <code>unsigned long int</code>  | 0              | 4,294,967,295 |

Table 7.2 shows the usual ranges on a 32-bit machine; here `int` and `long int` have identical ranges.

**Table 7.2**  
Integer Types on a  
32-bit Machine

| Type                            | Smallest Value | Largest Value |
|---------------------------------|----------------|---------------|
| <code>short int</code>          | -32,768        | 32,767        |
| <code>unsigned short int</code> | 0              | 65,535        |
| <code>int</code>                | -2,147,483,648 | 2,147,483,647 |
| <code>unsigned int</code>       | 0              | 4,294,967,295 |
| <code>long int</code>           | -2,147,483,648 | 2,147,483,647 |
| <code>unsigned long int</code>  | 0              | 4,294,967,295 |

In recent years, 64-bit CPUs have become more common. Table 7.3 shows typical ranges for the integer types on a 64-bit machine (especially under UNIX).

**Table 7.3**  
Integer Types on a  
64-bit Machine

| Type                            | Smallest Value             | Largest Value              |
|---------------------------------|----------------------------|----------------------------|
| <code>short int</code>          | -32,768                    | 32,767                     |
| <code>unsigned short int</code> | 0                          | 65,535                     |
| <code>int</code>                | -2,147,483,648             | 2,147,483,647              |
| <code>unsigned int</code>       | 0                          | 4,294,967,295              |
| <code>long int</code>           | -9,223,372,036,854,775,808 | 9,223,372,036,854,775,807  |
| <code>unsigned long int</code>  | 0                          | 18,446,744,073,709,551,615 |

Once more, let me emphasize that the ranges shown in Tables 7.1, 7.2, and 7.3 aren't mandated by the C standard and may vary from one compiler to another. One way to determine the ranges of the integer types for a particular implementation is to check the `<limits.h>` header, which is part of the standard library. This header defines macros that represent the smallest and largest values of each integer type.

**C99** Integer Types in C99

C99 provides two additional standard integer types, `long long int` and `unsigned long long int`. These types were added because of the growing need for very large integers and the ability of newer processors to support 64-bit arithmetic. Both `long long` types are required to be at least 64 bits wide, so the range of `long long int` values is typically  $-2^{63}$  ( $-9,223,372,036,854,775,808$ ) to  $2^{63} - 1$  ( $9,223,372,036,854,775,807$ ), and range of `unsigned long long int` values is usually 0 to  $2^{64} - 1$  ( $18,446,744,073,709,551,615$ ).

The `short int`, `int`, `long int`, and `long long int` types (along with the `signed char` type) are called *standard signed integer types* in C99. The `unsigned short int`, `unsigned int`, `unsigned long int`, and `unsigned long long int` types (along with the `unsigned char` type and the `_Bool` type) are called *standard unsigned integer types*.

In addition to the standard integer types, the C99 standard allows implementation-defined *extended integer types*, both signed and unsigned. For example, a compiler might provide signed and unsigned 128-bit integer types.

## Integer Constants

Let's turn our attention to *constants*—numbers that appear in the text of a program, not numbers that are read, written, or computed. C allows integer constants to be written in decimal (base 10), octal (base 8), or hexadecimal (base 16).

---

### Octal and Hexadecimal Numbers

An octal number is written using only the digits 0 through 7. Each position in an octal number represents a power of 8 (just as each position in a decimal number represents a power of 10). Thus, the octal number 237 represents the decimal number  $2 \times 8^2 + 3 \times 8^1 + 7 \times 8^0 = 128 + 24 + 7 = 159$ .

A hexadecimal (or hex) number is written using the digits 0 through 9 plus the letters A through F, which stand for 10 through 15, respectively. Each position in a hex number represents a power of 16; the hex number 1AF has the decimal value  $1 \times 16^2 + 10 \times 16^1 + 15 \times 16^0 = 256 + 160 + 15 = 431$ .

- 
- *Decimal* constants contain digits between 0 and 9, but must not begin with a zero:

15 255 32767

- *Octal* constants contain only digits between 0 and 7, and *must* begin with a zero:

017 0377 077777

- **Hexadecimal** constants contain digits between 0 and 9 and letters between a and f, and always begin with 0x:

0xf 0xff 0x7fff

The letters in a hexadecimal constant may be either upper or lower case:

0xff 0xFF 0xFF 0xFF 0Xff 0XFf 0XFF 0XFF

Keep in mind that octal and hexadecimal are nothing more than an alternative way of writing numbers; they have no effect on how the numbers are actually stored. (Integers are always stored in binary, regardless of what notation we've used to express them.) We can switch from one notation to another at any time, and even mix them: 10 + 015 + 0x20 has the value 55 (decimal). Octal and hex are most convenient for writing low-level programs; we won't use these notations much until Chapter 20.

The type of a *decimal* integer constant is normally `int`. However, if the value of the constant is too large to store as an `int`, the constant has type `long int` instead. In the unlikely case that the constant is too large to store as a `long int`, the compiler will try `unsigned long int` as a last resort. The rules for determining the type of an *octal* or *hexadecimal* constant are slightly different: the compiler will go through the types `int`, `unsigned int`, `long int`, and `unsigned long int` until it finds one capable of representing the constant.

To force the compiler to treat a constant as a long integer, just follow it with the letter L (or l):

15L 0377L 0x7ffffL

To indicate that a constant is `unsigned`, put the letter U (or u) after it:

15U 0377U 0x7ffffU

L and U may be used in combination to show that a constant is both `long` *and* `unsigned`: 0xffffffffFUL. (The order of the L and U doesn't matter, nor does their case.)

### C99

## Integer Constants in C99

In C99, integer constants that end with either LL or ll (the case of the two letters must match) have type `long long int`. Adding the letter U (or u) before or after the LL or ll denotes a constant of type `unsigned long long int`.

C99's general rules for determining the type of an integer constant are a bit different from those in C89. The type of a decimal constant with no suffix (U, u, L, l, LL, or ll) is the "smallest" of the types `int`, `long int`, or `long long int` that can represent the value of that constant. For an octal or hexadecimal constant, however, the list of possible types is `int`, `unsigned int`, `long int`, `unsigned long int`, `long long int`, and `unsigned long long int`, in that order. Any suffix at the end of a constant changes the list of possible types. For

example, a constant that ends with `U` (or `u`) must have one of the types `unsigned int`, `unsigned long int`, or `unsigned long long int`. A decimal constant that ends with `L` (or `l`) must have one of the types `long int` or `long long int`. There's also a provision for a constant to have an extended integer type if it's too large to represent using one of the standard integer types.

## Integer Overflow

When arithmetic operations are performed on integers, it's possible that the result will be too large to represent. For example, when an arithmetic operation is performed on two `int` values, the result must be able to be represented as an `int`. If the result can't be represented as an `int` (because it requires too many bits), we say that *overflow* has occurred.

The behavior when integer overflow occurs depends on whether the operands were signed or unsigned. When overflow occurs during an operation on *signed* integers, the program's behavior is undefined. Recall from Section 4.4 that the consequences of undefined behavior may vary. Most likely the result of the operation will simply be wrong, but the program could crash or exhibit other undesirable behavior.

When overflow occurs during an operation on *unsigned* integers, though, the result *is* defined: we get the correct answer modulo  $2^n$ , where  $n$  is the number of bits used to store the result. For example, if we add 1 to the unsigned 16-bit number 65,535, the result is guaranteed to be 0.

## Reading and Writing Integers

Suppose that a program isn't working because one of its `int` variables is overflowing. Our first thought is to change the type of the variable from `int` to `long int`. But we're not done yet; we need to see how the change will affect the rest of the program. In particular, we must check whether the variable is used in a call of `printf` or `scanf`. If so, the format string in the call will need to be changed, since the `%d` conversion works only for the `int` type.

Reading and writing `unsigned`, `short`, and `long` integers requires several new conversion specifiers:

### Q&A

- When reading or writing an *unsigned* integer, use the letter `u`, `o`, or `x` instead of `d` in the conversion specification. If the `u` specifier is present, the number is read (or written) in decimal notation; `o` indicates octal notation, and `x` indicates hexadecimal notation.

```
unsigned int u;

scanf("%u", &u); /* reads u in base 10 */
printf("%u", u); /* writes u in base 10 */
scanf("%o", &u); /* reads u in base 8 */
printf("%o", u); /* writes u in base 8 */
```

```
scanf("%x", &u); /* reads u in base 16 */
printf("%x", u); /* writes u in base 16 */
```

- When reading or writing a *short* integer, put the letter h in front of d, o, u, or x:

```
short s;

scanf("%hd", &s);
printf("%hd", s);
```

- When reading or writing a *long* integer, put the letter l (“ell,” not “one”) in front of d, o, u, or x:

```
long l;

scanf("%ld", &l);
printf("%ld", l);
```

**C99**

- When reading or writing a *long long* integer (C99 only), put the letters ll in front of d, o, u, or x:

```
long long ll;

scanf("%lld", &ll);
printf("%lld", ll);
```

## PROGRAM Summing a Series of Numbers (Revisited)

In Section 6.1, we wrote a program that sums a series of integers entered by the user. One problem with this program is that the sum (or one of the input numbers) might exceed the largest value allowed for an int variable. Here’s what might happen if the program is run on a machine whose integers are 16 bits long:

```
This program sums a series of integers.
Enter integers (0 to terminate): 10000 20000 30000 0
The sum is: -5536
```

The sum was 60,000, which wouldn’t fit in an int variable, so overflow occurred. When overflow occurs with signed numbers, the outcome is undefined. In this case, we got an apparently meaningless number. To improve the program, let’s switch to long variables.

```
sum2.c /* Sums a series of numbers (using long variables) */

#include <stdio.h>

int main(void)
{
    long n, sum = 0;

    printf("This program sums a series of integers.\n");
```

```

printf("Enter integers (0 to terminate): ");

scanf("%ld", &n);
while (n != 0) {
    sum += n;
    scanf("%ld", &n);
}
printf("The sum is: %ld\n", sum);

return 0;
}

```

The change was fairly simple: we declared `n` and `sum` to be `long` variables instead of `int` variables, then we changed the conversion specifications in `scanf` and `printf` to `%ld` instead of `%d`.

## 7.2 Floating Types

The integer types aren't suitable for all applications. Sometimes we'll need variables that can store numbers with digits after the decimal point, or numbers that are exceedingly large or small. Numbers like these are stored in floating-point format (so called because the decimal point "floats"). C provides three *floating types*, corresponding to different floating-point formats:

|                          |                                   |
|--------------------------|-----------------------------------|
| <code>float</code>       | Single-precision floating-point   |
| <code>double</code>      | Double-precision floating-point   |
| <code>long double</code> | Extended-precision floating-point |

`float` is suitable when the amount of precision isn't critical (calculating temperatures to one decimal point, for example). `double` provides greater precision—enough for most programs. `long double`, which supplies the ultimate in precision, is rarely used.

The C standard doesn't state how much precision the `float`, `double`, and `long double` types provide, since different computers may store floating-point numbers in different ways. Most modern computers follow the specifications in IEEE Standard 754 (also known as IEC 60559), so we'll use it as an example.

---

### The IEEE Floating-Point Standard

IEEE Standard 754, developed by the Institute of Electrical and Electronics Engineers, provides two primary formats for floating-point numbers: single precision (32 bits) and double precision (64 bits). Numbers are stored in a form of scientific notation, with each number having three parts: a *sign*, an *exponent*, and a *fraction*. The number of bits reserved for the exponent determines how large (or small) numbers can be, while the number of bits in the fraction determines the precision. In single-precision format, the exponent is 8 bits long, while the fraction occupies 23

bits. As a result, a single-precision number has a maximum value of approximately  $3.40 \times 10^{38}$ , with a precision of about 6 decimal digits.

The IEEE standard also describes two other formats, single extended precision and double extended precision. The standard doesn't specify the number of bits in these formats, although it requires that the single extended type occupy at least 43 bits and the double extended type at least 79 bits. For more information about the IEEE standard and floating-point arithmetic in general, see "What every computer scientist should know about floating-point arithmetic" by David Goldberg (*ACM Computing Surveys*, vol. 23, no. 1 (March 1991): 5–48).

Table 7.4 shows the characteristics of the floating types when implemented according to the IEEE standard. (The table shows the smallest positive *normalized* values. Subnormal numbers can be smaller.) The `long double` type isn't shown in the table, since its length varies from one machine to another, with 80 bits and 128 bits being the most common sizes.

subnormal numbers ► 23.4  
**Table 7.4**  
 Floating Type Characteristics (IEEE Standard)

| Type                | Smallest Positive Value    | Largest Value             | Precision |
|---------------------|----------------------------|---------------------------|-----------|
| <code>float</code>  | $1.17549 \times 10^{-38}$  | $3.40282 \times 10^{38}$  | 6 digits  |
| <code>double</code> | $2.22507 \times 10^{-308}$ | $1.79769 \times 10^{308}$ | 15 digits |

On computers that don't follow the IEEE standard, Table 7.4 won't be valid. In fact, on some machines, `float` may have the same set of values as `double`, or `double` may have the same values as `long double`. Macros that define the characteristics of the floating types can be found in the `<float.h>` header.

In C99, the floating types are divided into two categories. The `float`, `double`, and `long double` types fall into one category, called the *real floating types*. Floating types also include the *complex types* (`float_Complex`, `double_Complex`, and `long double_Complex`), which are new in C99.

## Floating Constants

Floating constants can be written in a variety of ways. The following constants, for example, are all valid ways of writing the number 57.0:

57.0 57. 57.0e0 57E0 5.7e1 5.7e+1 .57e2 570.e-1

A floating constant must contain a decimal point and/or an exponent; the exponent indicates the power of 10 by which the number is to be scaled. If an exponent is present, it must be preceded by the letter E (or e). An optional + or - sign may appear after the E (or e).

By default, floating constants are stored as double-precision numbers. In other words, when a C compiler finds the constant 57.0 in a program, it arranges for the number to be stored in memory in the same format as a `double` variable. This rule generally causes no problems, since `double` values are converted automatically to `float` when necessary.

On occasion, it may be necessary to force the compiler to store a floating constant in `float` or `long double` format. To indicate that only single precision is desired, put the letter `F` (or `f`) at the end of the constant (for example, `57.0F`). To indicate that a constant should be stored in `long double` format, put the letter `L` (or `l`) at the end (`57.0L`).

**C99**

**C99** has a provision for writing floating constants in hexadecimal. Such a constant begins with `0x` or `0X` (like a hexadecimal integer constant). This feature is rarely used.

## Reading and Writing Floating-Point Numbers

As we've discussed, the conversion specifications `%e`, `%f`, and `%g` are used for reading and writing single-precision floating-point numbers. Values of types `double` and `long double` require slightly different conversions:

- When *reading* a value of type `double`, put the letter `l` in front of `e`, `f`, or `g`:

```
double d;
scanf ("%lf", &d);
```

**Q&A**

**Note:** Use `l` only in a `scanf` format string, not a `printf` string. In a `printf` format string, the `e`, `f`, and `g` conversions can be used to write either `float` or `double` values. (C99 legalizes the use of `%le`, `%lf`, and `%lg` in calls of `printf`, although the `l` has no effect.)

- When reading or writing a value of type `long double`, put the letter `L` in front of `e`, `f`, or `g`:

```
long double ld;
scanf ("%Lf", &ld);
printf ("%Lf", ld);
```

**7.3 Character Types**

**Q&A** The only remaining basic type is `char`, the character type. The values of type `char` can vary from one computer to another, because different machines may have different underlying character sets.

### Character Sets

ASCII character set ➤ Appendix E

Today's most popular character set is **ASCII** (American Standard Code for Information Interchange), a 7-bit code capable of representing 128 characters. In ASCII, the digits 0 to 9 are represented by the codes 0110000–0111001, and the uppercase letters A to Z are represented by 1000001–1011010. ASCII is often extended

to a 256-character code known as *Latin-1* that provides the characters necessary for Western European and many African languages.

---

A variable of type `char` can be assigned any single character:

```
char ch;

ch = 'a';      /* lower-case a */
ch = 'A';      /* upper-case A */
ch = '0';      /* zero */
ch = ' ';      /* space */
```

Notice that character constants are enclosed in single quotes, not double quotes.

## Operations on Characters

Working with characters in C is simple, because of one fact: *C treats characters as small integers*. After all, characters are encoded in binary, and it doesn't take much imagination to view these binary codes as integers. In ASCII, for example, character codes range from 0000000 to 1111111, which we can think of as the integers from 0 to 127. The character '`'a'`' has the value 97, '`'A'`' has the value 65, '`'0'`' has the value 48, and '`' '`' has the value 32. The connection between characters and integers in C is so strong that character constants actually have `int` type rather than `char` type (an interesting fact, but not one that will often matter to us).

When a character appears in a computation, C simply uses its integer value. Consider the following examples, which assume the ASCII character set:

```
char ch;
int i;

i = 'a';          /* i is now 97 */
ch = 65;          /* ch is now 'A' */
ch = ch + 1;      /* ch is now 'B' */
ch++;            /* ch is now 'C' */
```

Characters can be compared, just as numbers can. The following `if` statement checks whether `ch` contains a lower-case letter; if so, it converts `ch` to upper case.

```
if ('a' <= ch && ch <= 'z')
    ch = ch - 'a' + 'A';
```

Comparisons such as '`'a' <= ch`' are done using the integer values of the characters involved. These values depend on the character set in use, so programs that use `<`, `<=`, `>`, and `>=` to compare characters may not be portable.

The fact that characters have the same properties as numbers has some advantages. For example, we can easily write a `for` statement whose control variable steps through all the upper-case letters:

```
for (ch = 'A'; ch <= 'Z'; ch++) ...
```

On the other hand, treating characters as numbers can lead to various programming errors that won't be caught by the compiler, and lets us write meaningless expressions such as `'a' * 'b' / 'c'`. It can also hamper portability, since our programs may be based on assumptions about the underlying character set. (Our `for` loop, for example, assumes that the letters from A to Z have consecutive codes.)

## Signed and Unsigned Characters

Since C allows characters to be used as integers, it shouldn't be surprising that the `char` type—like the integer types—exists in both signed and unsigned versions. Signed characters normally have values between -128 and 127, while unsigned characters have values between 0 and 255.

The C standard doesn't specify whether ordinary `char` is a signed or an unsigned type; some compilers treat it as a signed type, while others treat it as an unsigned type. (Some even allow the programmer to select, via a compiler option, whether `char` should be signed or unsigned.)

Most of the time, we don't really care whether `char` is signed or unsigned. Once in a while, though, we do, especially if we're using a character variable to store a small integer. For this reason, C allows the use of the words `signed` and `unsigned` to modify `char`:

```
signed char sch;
unsigned char uch;
```

### portability tip

*Don't assume that `char` is either signed or unsigned by default. If it matters, use `signed char` or `unsigned char` instead of `char`.*

enumerated types ➤ 16.5

In light of the close relationship between characters and integers, C89 uses the term **integral types** to refer to both the integer types and the character types. Enumerated types are also integral types.

C99 doesn't use the term "integral types." Instead, it expands the meaning of "integer types" to include the character types and the enumerated types. C99's `_Bool` type is considered to be an unsigned integer type.

## Arithmetic Types

The integer types and floating types are collectively known as **arithmetic types**. Here's a summary of the arithmetic types in C89, divided into categories and sub-categories:

- Integral types
  - `char`
  - Signed integer types (`signed char`, `short int`, `int`, `long int`)
  - Unsigned integer types (`unsigned char`, `unsigned short int`, `unsigned int`, `unsigned long int`)

C99

`_Bool` type ➤ 5.2

- Enumerated types
- Floating types (`float`, `double`, `long double`)

**C99** C99 has a more complicated hierarchy for its arithmetic types:

- Integer types
  - `char`
  - Signed integer types, both standard (`signed char`, `short int`, `int`, `long int`, `long long int`) and extended
  - Unsigned integer types, both standard (`unsigned char`, `unsigned short int`, `unsigned int`, `unsigned long int`, `unsigned long long int`, `_Bool`) and extended
  - Enumerated types
- Floating types
  - Real floating types (`float`, `double`, `long double`)
  - Complex types (`float _Complex`, `double _Complex`, `long double _Complex`)

## Escape Sequences

A character constant is usually one character enclosed in single quotes, as we've seen in previous examples. However, certain special characters—including the new-line character—can't be written in this way, because they're invisible (non-printing) or because they can't be entered from the keyboard. So that programs can deal with every character in the underlying character set, C provides a special notation, the *escape sequence*.

There are two kinds of escape sequences: *character escapes* and *numeric escapes*. We saw a partial list of character escapes in Section 3.1; Table 7.5 gives the complete set.

**Table 7.5**  
Character Escapes

| Name            | Escape Sequence |
|-----------------|-----------------|
| Alert (bell)    | \a              |
| Backspace       | \b              |
| Form feed       | \f              |
| New line        | \n              |
| Carriage return | \r              |
| Horizontal tab  | \t              |
| Vertical tab    | \v              |
| Backslash       | \\\             |
| Question mark   | \?              |
| Single quote    | \'              |
| Double quote    | \"              |

**Q&A** The \a, \b, \f, \r, \t, and \v escapes represent common ASCII control characters. The \n escape represents the ASCII line-feed character. The \\ escape allows a character constant or string to contain the \ character. The \' escape

**Q&A**

allows a character constant to contain the ' character, while the \" escape allows a string to contain the " character. The \\? escape is rarely used.

Character escapes are handy, but they have a problem: the list of character escapes doesn't include all the nonprinting ASCII characters, just the most common. Character escapes are also useless for representing characters beyond the basic 128 ASCII characters. Numeric escapes, which can represent *any* character, are the solution to this problem.

To write a numeric escape for a particular character, first look up the character's octal or hexadecimal value in a table like the one in Appendix E. For example, the ASCII escape character (decimal value: 27) has the value 33 in octal and 1B in hex. Either of these codes can be used to write an escape sequence:

- An *octal escape sequence* consists of the \ character followed by an octal number with at most three digits. (This number must be representable as an unsigned character, so its maximum value is normally 377 octal.) For example, the escape character could be written \\33 or \\033. Octal numbers in escape sequences—unlike octal constants—don't have to begin with 0.
- A *hexadecimal escape sequence* consists of \\x followed by a hexadecimal number. Although C places no limit on the number of digits in the hexadecimal number, it must be representable as an unsigned character (hence it can't exceed FF if characters are eight bits long). Using this notation, the escape character would be written \\x1b or \\x1B. The x must be in lower case, but the hex digits (such as b) can be upper or lower case.

When used as a character constant, an escape sequence must be enclosed in single quotes. For example, a constant representing the escape character would be written '\\33' (or '\\x1b'). Escape sequences tend to get a bit cryptic, so it's often a good idea to give them names using #define:

```
#define ESC '\\33' /* ASCII escape character */
```

Escape sequences can be embedded in strings as well, as we saw in Section 3.1.

Escape sequences aren't the only special notations for representing characters. Trigraph sequences provide a way to represent the characters #, [, \, ], ^, {, |, }, and ~, which may not be available on keyboards in some countries. C99 adds universal character names, which resemble escape sequences. Unlike escape sequences, however, universal character names are allowed in identifiers.

trigraph sequences ▶ 25.3

**C99**

universal character names ▶ 25.4

## Character-Handling Functions

Earlier in this section, we saw how to write an if statement that converts a lower-case letter to upper-case:

```
if ('a' <= ch && ch <= 'z')
    ch = ch - 'a' + 'A';
```

This isn't the best method, though. A faster—and more portable—way to convert case is to call C's toupper library function:

```
ch = toupper(ch); /* converts ch to upper case */
```

When it's called, `toupper` checks whether its argument (`ch` in this case) is a lower-case letter. If so, it returns the corresponding upper-case letter. Otherwise, `toupper` returns the value of the argument. In our example, we've used the assignment operator to store the return value of `toupper` back into the `ch` variable, although we could just as easily have done something else with it—stored it in another variable, say, or tested it in an `if` statement:

```
if (toupper(ch) == 'A') ...
```

Programs that call `toupper` need to have the following `#include` directive at the top:

```
#include <ctype.h>
```

`toupper` isn't the only useful character-handling function in the C library. Section 23.5 describes them all and gives examples of their use.

## Reading and Writing Characters using `scanf` and `printf`

The `%c` conversion specification allows `scanf` and `printf` to read and write single characters:

```
char ch;
```

```
scanf("%c", &ch); /* reads a single character */
printf("%c", ch); /* writes a single character */
```

`scanf` doesn't skip white-space characters before reading a character. If the next unread character is a space, then the variable `ch` in the previous example will contain a space after `scanf` returns. To force `scanf` to skip white space before reading a character, put a space in its format string just before `%c`:

```
scanf(" %c", &ch); /* skips white space, then reads ch */
```

Recall from Section 3.2 that a blank in a `scanf` format string means “skip zero or more white-space characters.”

Since `scanf` doesn't normally skip white space, it's easy to detect the end of an input line: check to see if the character just read is the new-line character. For example, the following loop will read and ignore all remaining characters in the current input line:

```
do {
    scanf("%c", &ch);
} while (ch != '\n');
```

When `scanf` is called the next time, it will read the first character on the next input line.

## Reading and Writing Characters using `getchar` and `putchar`

### Q&A

C provides other ways to read and write single characters. In particular, we can use the `getchar` and `putchar` functions instead of calling `scanf` and `printf`. `putchar` writes a single character:

```
putchar(ch);
```

Each time `getchar` is called, it reads one character, which it returns. In order to save this character, we must use assignment to store it in a variable:

```
ch = getchar(); /* reads a character and stores it in ch */
```

`getchar` actually returns an `int` value rather than a `char` value (the reason will be discussed in later chapters). As a result, it's not unusual for a variable to have type `int` rather than `char` if it will be used to store a character read by `getchar`. Like `scanf`, `getchar` doesn't skip white-space characters as it reads.

Using `getchar` and `putchar` (rather than `scanf` and `printf`) saves time when the program is executed. `getchar` and `putchar` are fast for two reasons. First, they're much simpler than `scanf` and `printf`, which are designed to read and write many kinds of data in a variety of formats. Second, `getchar` and `putchar` are usually implemented as macros for additional speed.

macros ▶ 14.3

`getchar` has another advantage over `scanf`: because it returns the character that it reads, `getchar` lends itself to various C idioms, including loops that search for a character or skip over all occurrences of a character. Consider the `scanf` loop that we used to skip the rest of an input line:

```
do {
    scanf("%c", &ch);
} while (ch != '\n');
```

Rewriting this loop using `getchar` gives us the following:

```
do {
    ch = getchar();
} while (ch != '\n');
```

Moving the call of `getchar` into the controlling expression allows us to condense the loop:

```
while ((ch = getchar()) != '\n')
;
```

This loop reads a character, stores it into the variable `ch`, then tests if `ch` is not equal to the new-line character. If the test succeeds, the loop body (which is empty) is executed, then the loop test is performed once more, causing a new character to be read. Actually, we don't even need the `ch` variable; we can just compare the return value of `getchar` with the new-line character:

```
idiom while (getchar() != '\n') /* skips rest of line */
;
```

The resulting loop is a well-known C idiom that's cryptic but worth learning.

`getchar` is useful in loops that skip characters as well as loops that search for characters. Consider the following statement, which uses `getchar` to skip an indefinite number of blank characters:

```
idiom while ((ch = getchar()) == ' ') /* skips blanks */
;
```

When the loop terminates, `ch` will contain the first nonblank character that `getchar` encountered.



Be careful if you mix `getchar` and `scanf` in the same program. `scanf` has a tendency to leave behind characters that it has “peeked” at but not read, including the new-line character. Consider what happens if we try to read a number first, then a character:

```
printf("Enter an integer: ");
scanf("%d", &i);
printf("Enter a command: ");
command = getchar();
```

The call of `scanf` will leave behind any characters that weren't consumed during the reading of `i`, including (but not limited to) the new-line character. `getchar` will fetch the first leftover character, which wasn't what we had in mind.

## PROGRAM Determining the Length of a Message

To illustrate how characters are read, let's write a program that calculates the length of a message. After the user enters the message, the program displays the length:

```
Enter a message: Brevity is the soul of wit.
Your message was 27 character(s) long.
```

The length includes spaces and punctuation, but not the new-line character at the end of the message.

We'll need a loop whose body reads a character and increments a counter. The loop will terminate as soon as a new-line character turns up. We could use either `scanf` or `getchar` to read characters; most C programmers would choose `getchar`. Using a straightforward `while` loop, we might end up with the following program.

```
length.c /* Determines the length of a message */

#include <stdio.h>

int main(void)
{
    char ch;
    int len = 0;

    printf("Enter a message: ");
    ch = getchar();
    while (ch != '\n') {
        len++;
        ch = getchar();
    }
    printf("Your message was %d character(s) long.\n", len);

    return 0;
}
```

Recalling our discussion of idioms involving while loops and getchar, we realize that the program can be shortened:

```
length2.c /* Determines the length of a message */

#include <stdio.h>

int main(void)
{
    int len = 0;

    printf("Enter a message: ");
    while (getchar() != '\n')
        len++;
    printf("Your message was %d character(s) long.\n", len);

    return 0;
}
```

## 7.4 Type Conversion

Computers tend to be more restrictive than C when it comes to arithmetic. For a computer to perform an arithmetic operation, the operands must usually be of the same size (the same number of bits) and be stored in the same way. A computer may be able to add two 16-bit integers directly, but not a 16-bit integer and a 32-bit integer or a 32-bit integer and a 32-bit floating-point number.

C, on the other hand, allows the basic types to be mixed in expressions. We can combine integers, floating-point numbers, and even characters in a single expression. The C compiler may then have to generate instructions that convert

some operands to different types so that the hardware will be able to evaluate the expression. If we add a 16-bit `short` and a 32-bit `int`, for example, the compiler will arrange for the `short` value to be converted to 32 bits. If we add an `int` and a `float`, the compiler will arrange for the `int` to be converted to `float` format. This conversion is a little more complicated, since `int` and `float` values are stored in different ways.

Because the compiler handles these conversions automatically, without the programmer's involvement, they're known as *implicit conversions*. C also allows the programmer to perform *explicit conversions*, using the cast operator. I'll discuss implicit conversions first, postponing explicit conversions until later in the section. Unfortunately, the rules for performing implicit conversions are somewhat complex, primarily because C has so many different arithmetic types.

Implicit conversions are performed in the following situations:

- When the operands in an arithmetic or logical expression don't have the same type. (C performs what are known as the *usual arithmetic conversions*.)
- When the type of the expression on the right side of an assignment doesn't match the type of the variable on the left side.
- When the type of an argument in a function call doesn't match the type of the corresponding parameter.
- When the type of the expression in a `return` statement doesn't match the function's return type.

We'll discuss the first two cases now and save the others for Chapter 9.

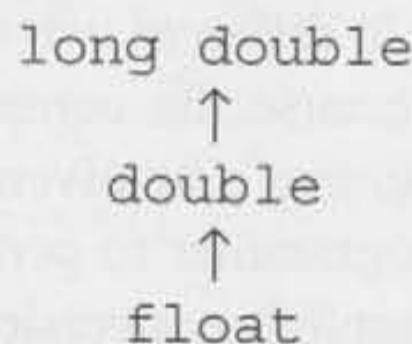
## The Usual Arithmetic Conversions

The usual arithmetic conversions are applied to the operands of most binary operators, including the arithmetic, relational, and equality operators. For example, let's say that `f` has type `float` and `i` has type `int`. The usual arithmetic conversions will be applied to the operands in the expression `f + i`, because their types aren't the same. Clearly it's safer to convert `i` to type `float` (matching `f`'s type) rather than convert `f` to type `int` (matching `i`'s type). An integer can always be converted to `float`; the worst that can happen is a minor loss of precision. Converting a floating-point number to `int`, on the other hand, would cost us the fractional part of the number. Worse still, we'd get a completely meaningless result if the original number were larger than the largest possible integer or smaller than the smallest integer.

The strategy behind the usual arithmetic conversions is to convert operands to the “narrowest” type that will safely accommodate both values. (Roughly speaking, one type is narrower than another if it requires fewer bytes to store.) The types of the operands can often be made to match by converting the operand of the narrower type to the type of the other operand (this act is known as *promotion*). Among the most common promotions are the *integral promotions*, which convert a character or short integer to type `int` (or to `unsigned int` in some cases).

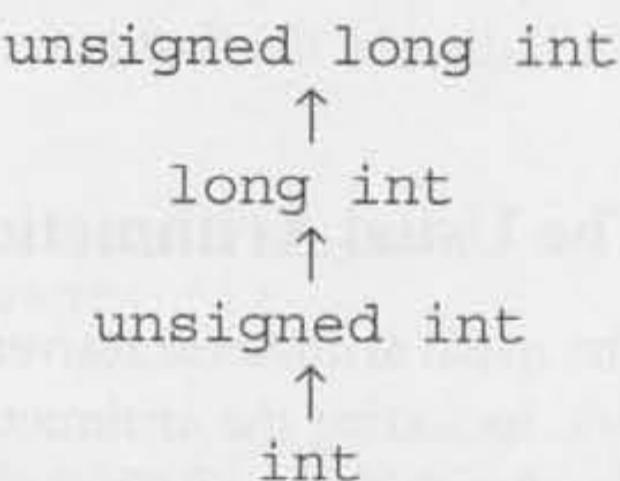
We can divide the rules for performing the usual arithmetic conversions into two cases:

- *The type of either operand is a floating type.* Use the following diagram to promote the operand whose type is narrower:



That is, if one operand has type `long double`, then convert the other operand to type `long double`. Otherwise, if one operand has type `double`, convert the other operand to type `double`. Otherwise, if one operand has type `float`, convert the other operand to type `float`. Note that these rules cover mixtures of integer and floating types: if one operand has type `long int`, for example, and the other has type `double`, the `long int` operand is converted to `double`.

- *Neither operand type is a floating type.* First perform integral promotion on both operands (guaranteeing that neither operand will be a character or short integer). Then use the following diagram to promote the operand whose type is narrower:



There's one special case, but it occurs only when `long int` and `unsigned int` have the same length (32 bits, say). Under these circumstances, if one operand has type `long int` and the other has type `unsigned int`, both are converted to `unsigned long int`.



When a signed operand is combined with an unsigned operand, the signed operand is converted to an unsigned value. The conversion involves adding or subtracting a multiple of  $n + 1$ , where  $n$  is the largest representable value of the unsigned type. This rule can cause obscure programming errors.

Suppose that the `int` variable `i` has the value `-10` and the `unsigned int` variable `u` has the value `10`. If we compare `i` and `u` using the `<` operator, we might expect to get the result `1` (true). Before the comparison, however, `i` is converted to `unsigned int`. Since a negative number can't be represented as an `unsigned int`, the converted value won't be `-10`. Instead, the value `4,294,967,296` is added (assuming that `4,294,967,295` is the largest `unsigned int` value), giving

a converted value of 4,294,967,286. The comparison  $i < u$  will therefore produce 0. Some compilers produce a warning message such as “*comparison between signed and unsigned*” when a program attempts to compare a signed number with an unsigned number.

Because of traps like this one, it’s best to use unsigned integers as little as possible and, especially, never mix them with signed integers.

The following example shows the usual arithmetic conversions in action:

```
char c;
short int s;
int i;
unsigned int u;
long int l;
unsigned long int ul;
float f;
double d;
long double ld;

i = i + c;      /* c is converted to int */ 
i = i + s;      /* s is converted to int */ 
u = u + i;      /* i is converted to unsigned int */ 
l = l + u;      /* u is converted to long int */ 
ul = ul + l;    /* l is converted to unsigned long int */ 
f = f + ul;    /* ul is converted to float */ 
d = d + f;      /* f is converted to double */ 
ld = ld + d;    /* d is converted to long double */
```

## Conversion During Assignment

The usual arithmetic conversions don’t apply to assignment. Instead, C follows the simple rule that the expression on the right side of the assignment is converted to the type of the variable on the left side. If the variable’s type is at least as “wide” as the expression’s, this will work without a snag. For example:

```
char c;
int i;
float f;
double d;

i = c;    /* c is converted to int */ 
f = i;    /* i is converted to float */ 
d = f;    /* f is converted to double */
```

Other cases are problematic. Assigning a floating-point number to an integer variable drops the fractional part of the number:

```
int i;

i = 842.97;    /* i is now 842 */
i = -842.97;   /* i is now -842 */
```

**Q&A**

Moreover, assigning a value to a variable of a narrower type will give a meaningless result (or worse) if the value is outside the range of the variable's type:

```
c = 10000;      /*** WRONG ***/
i = 1.0e20;    /*** WRONG ***/
f = 1.0e100;   /*** WRONG ***/
```

A “narrowing” assignment may elicit a warning from the compiler or from tools such as lint.

It's a good idea to append the `f` suffix to a floating-point constant if it will be assigned to a `float` variable, as we've been doing since Chapter 2:

```
f = 3.14159f;
```

Without the suffix, the constant `3.14159` would have type `double`, possibly causing a warning message.

**C99****Implicit Conversions in C99**

\_Bool type ▶ 5.2

The rules for implicit conversions in C99 are somewhat different from the rules in C89, primarily because C99 has additional types (`_Bool`, long long types, extended integer types, and complex types).

For the purpose of defining conversion rules, C99 gives each integer type an “integer conversion rank.” Here are the ranks from highest to lowest:

1. `long long int`, `unsigned long long int`
2. `long int`, `unsigned long int`
3. `int`, `unsigned int`
4. `short int`, `unsigned short int`
5. `char`, `signed char`, `unsigned char`
6. `_Bool`

For simplicity, I'm ignoring extended integer types and enumerated types.

In place of C89's integral promotions, C99 has “integer promotions,” which involve converting any type whose rank is less than `int` and `unsigned int` to `int` (provided that all values of the type can be represented using `int`) or else to `unsigned int`.

As in C89, the C99 rules for performing the usual arithmetic conversions can be divided into two cases:

- ***The type of either operand is a floating type.*** As long as neither operand has a complex type, the rules are the same as before. (The conversion rules for complex types will be discussed in Section 27.3.)
- ***Neither operand type is a floating type.*** First perform integer promotion on both operands. If the types of the two operands are now the same, the process ends. Otherwise, use the following rules, stopping at the first one that applies:
  - If both operands have signed types or both have unsigned types, convert the

operand whose type has lesser integer conversion rank to the type of the operand with greater rank.

- If the unsigned operand has rank greater or equal to the rank of the type of the signed operand, convert the signed operand to the type of the unsigned operand.
- If the type of the signed operand can represent all of the values of the type of the unsigned operand, convert the unsigned operand to the type of the signed operand.
- Otherwise, convert both operands to the unsigned type corresponding to the type of the signed operand.

Incidentally, all arithmetic types can be converted to `_Bool` type. The result of the conversion is 0 if the original value is 0; otherwise, the result is 1.

## Casting

Although C's implicit conversions are convenient, we sometimes need a greater degree of control over type conversion. For this reason, C provides *casts*. A cast expression has the form

**cast expression**

( *type-name* ) *expression*

*type-name* specifies the type to which the expression should be converted.

The following example shows how to use a cast expression to compute the fractional part of a `float` value:

```
float f, frac_part;
frac_part = f - (int) f;
```

The cast expression `(int) f` represents the result of converting the value of `f` to type `int`. C's usual arithmetic conversions then require that `(int) f` be converted back to type `float` before the subtraction can be performed. The difference between `f` and `(int) f` is the fractional part of `f`, which was dropped during the cast.

Cast expressions enable us to document type conversions that would take place anyway:

```
i = (int) f; /* f is converted to int */
```

They also enable us to overrule the compiler and force it to do conversions that we want. Consider the following example:

```
float quotient;
int dividend, divisor;
quotient = dividend / divisor;
```

As it's now written, the result of the division—an integer—will be converted to `float` form before being stored in `quotient`. We probably want `dividend` and `divisor` converted to `float` *before* the division, though, so that we get a more exact answer. A cast expression will do the trick:

```
quotient = (float) dividend / divisor;
```

`divisor` doesn't need a cast, since casting `dividend` to `float` forces the compiler to convert `divisor` to `float` also.

Incidentally, C regards (*type-name*) as a unary operator. Unary operators have higher precedence than binary operators, so the compiler interprets

```
(float) dividend / divisor
```

as

```
((float) dividend) / divisor
```

If you find this confusing, note that there are other ways to accomplish the same effect:

```
quotient = dividend / (float) divisor;
```

or

```
quotient = (float) dividend / (float) divisor;
```

Casts are sometimes necessary to avoid overflow. Consider the following example:

```
long i;  
int j = 1000;  
  
i = j * j; /* overflow may occur */
```

At first glance, this statement looks fine. The value of `j * j` is 1,000,000, and `i` is a `long`, so it can easily store values of this size, right? The problem is that when two `int` values are multiplied, the result will have `int` type. But `j * j` is too large to represent as an `int` on some machines, causing an overflow. Fortunately, using a cast avoids the problem:

```
i = (long) j * j;
```

Since the cast operator takes precedence over `*`, the first `j` is converted to `long` type, forcing the second `j` to be converted as well. Note that the statement

```
i = (long) (j * j); /* *** WRONG *** /
```

wouldn't work, since the overflow would already have occurred by the time of the cast.

## 7.5 Type Definitions

In Section 5.2, we used the `#define` directive to create a macro that could be used as a Boolean type:

```
#define BOOL int
```

**Q&A** There's a better way to set up a Boolean type, though, using a feature known as a *type definition*:

```
typedef int Bool;
```

Notice that the name of the type being defined comes *last*. Note also that I've capitalized the word `Bool`. Capitalizing the first letter of a type name isn't required; it's just a convention that some C programmers employ.

Using `typedef` to define `Bool` causes the compiler to add `Bool` to the list of type names that it recognizes. `Bool` can now be used in the same way as the built-in type names—in variable declarations, cast expressions, and elsewhere. For example, we might use `Bool` to declare variables:

```
Bool flag; /* same as int flag; */
```

The compiler treats `Bool` as a synonym for `int`; thus, `flag` is really nothing more than an ordinary `int` variable.

### Advantages of Type Definitions

Type definitions can make a program more understandable (assuming that the programmer has been careful to choose meaningful type names). For example, suppose that the variables `cash_in` and `cash_out` will be used to store dollar amounts. Declaring `Dollars` as

```
typedef float Dollars;
```

and then writing

```
Dollars cash_in, cash_out;
```

is more informative than just writing

```
float cash_in, cash_out;
```

Type definitions can also make a program easier to modify. If we later decide that `Dollars` should really be defined as `double`, all we need do is change the type definition:

```
typedef double Dollars;
```

The declarations of `Dollars` variables need not be changed. Without the type definition, we would need to locate all `float` variables that store dollar amounts (not necessarily an easy task) and change their declarations.

## Type Definitions and Portability

Type definitions are an important tool for writing portable programs. One of the problems with moving a program from one computer to another is that types may have different ranges on different machines. If `i` is an `int` variable, an assignment like

```
i = 100000;
```

is fine on a machine with 32-bit integers, but will fail on a machine with 16-bit integers.

**portability tip**

*For greater portability, consider using `typedef` to define new names for integer types.*

Suppose that we're writing a program that needs variables capable of storing product quantities in the range 0–50,000. We could use `long` variables for this purpose (since they're guaranteed to be able to hold numbers up to at least 2,147,483,647), but we'd rather use `int` variables, since arithmetic on `int` values may be faster than operations on `long` values; also, `int` variables may take up less space.

Instead of using the `int` type to declare quantity variables, we can define our own “quantity” type:

```
typedef int Quantity;
```

and use this type to declare variables:

```
Quantity q;
```

When we transport the program to a machine with shorter integers, we'll change the definition of `Quantity`:

```
typedef long Quantity;
```

This technique doesn't solve all our problems, unfortunately, since changing the definition of `Quantity` may affect the way `Quantity` variables are used. At the very least, calls of `printf` and `scanf` that use `Quantity` variables will need to be changed, with `%d` conversion specifications replaced by `%ld`.

The C library itself uses `typedef` to create names for types that can vary from one C implementation to another; these types often have names that end with `_t`, such as `ptrdiff_t`, `size_t`, and `wchar_t`. The exact definitions of these types will vary, but here are some typical examples:

```
typedef long int ptrdiff_t;
typedef unsigned long int size_t;
typedef int wchar_t;
```

**C99**

&lt;stdint.h&gt; header ▶ 27.1

In C99, the `<stdint.h>` header uses `typedef` to define names for integer types with a particular number of bits. For example, `int32_t` is a signed integer type with exactly 32 bits. Using these types is an effective way to make programs more portable.

## 7.6 The `sizeof` Operator

The `sizeof` operator allows a program to determine how much memory is required to store values of a particular type. The value of the expression

**sizeof expression**`sizeof ( type-name )`**Q&A**

is an unsigned integer representing the number of bytes required to store a value belonging to *type-name*. `sizeof(char)` is always 1, but the sizes of the other types may vary. On a 32-bit machine, `sizeof(int)` is normally 4. Note that `sizeof` is a rather unusual operator, since the compiler itself can usually determine the value of a `sizeof` expression.

The `sizeof` operator can also be applied to constants, variables, and expressions in general. If *i* and *j* are `int` variables, then `sizeof(i)` is 4 on a 32-bit machine, as is `sizeof(i + j)`. When applied to an expression—as opposed to a type—`sizeof` doesn’t require parentheses; we could write `sizeof i` instead of `sizeof(i)`. However, parentheses may be needed anyway because of operator precedence. The compiler would interpret `sizeof i + j` as `(sizeof i) + j`, because `sizeof`—a unary operator—takes precedence over the binary `+` operator. To avoid problems, I always use parentheses in `sizeof` expressions.

Printing a `sizeof` value requires care, because the type of a `sizeof` expression is an implementation-defined type named `size_t`. In C89, it’s best to convert the value of the expression to a known type before printing it. `size_t` is guaranteed to be an unsigned integer type, so it’s safest to cast a `sizeof` expression to `unsigned long` (the largest of C89’s unsigned types) and then print it using the `%lu` conversion:

```
printf("Size of int: %lu\n", (unsigned long) sizeof(int));
```

**C99**

In C99, the `size_t` type can be larger than `unsigned long`. However, the `printf` function in C99 is capable of displaying `size_t` values directly, without needing a cast. The trick is to use the letter *z* in the conversion specification, followed by one of the usual integer codes (typically *u*):

```
printf("Size of int: %zu\n", sizeof(int)); /* C99 only */
```

## Q & A

**Q:** Section 7.1 says that `%o` and `%x` are used to write unsigned integers in octal and hex notation. How do I write ordinary (signed) integers in octal or hex? [p. 130]

**A:** You can use `%o` and `%x` to print a signed integer as long as its value isn't negative. These conversions cause `printf` to treat a signed integer as though it were unsigned; in other words, `printf` will assume that the sign bit is part of the number's magnitude. As long as the sign bit is 0, there's no problem. If the sign bit is 1, `printf` will print an unexpectedly large number.

**Q:** But what if the number is negative? How can I write it in octal or hex?

**A:** There's no direct way to print a negative number in octal or hex. Fortunately, the need to do so is pretty rare. You can, of course, test whether the number is negative and print a minus sign yourself:

```
if (i < 0)
    printf("-%x", -i);
else
    printf("%x", i);
```

**Q:** Why are floating constants stored in `double` form rather than `float` form? [p. 133]

**A:** For historical reasons, C gives preference to the `double` type; `float` is treated as a second-class citizen. Consider, for instance, the discussion of `float` in Kernighan and Ritchie's *The C Programming Language*: "The main reason for using `float` is to save storage in large arrays, or, less often, to save time on machines where double-precision arithmetic is particularly expensive." C originally mandated that all floating-point arithmetic be done in double precision. (C89 and C99 have no such requirement.)

**\*Q:** What do hexadecimal floating constants look like, and what are they good for? [p. 134]

**A:** A hexadecimal floating constant begins with `0x` or `0X` and must contain an exponent, which is preceded by the letter `P` (or `p`). The exponent may have a sign, and the constant may end with `f`, `F`, `l`, or `L`. The exponent is expressed in decimal, but represents a power of 2, not a power of 10. For example, `0x1.Bp3` represents the number  $1.6875 \times 2^3 = 13.5$ . The hex digit `B` corresponds to the bit pattern 1011. The `B` occurs to the right of the period, so each 1 bit represents a negative power of 2. Summing these powers of 2 ( $2^{-1} + 2^{-3} + 2^{-4}$ ) yields .6875.

Hexadecimal floating constants are primarily useful for specifying constants that require great precision (including mathematical constants such as  $e$  and  $\pi$ ). Hex numbers have a precise binary representation, whereas a constant written in decimal may be subject to a tiny rounding error when converted to binary. Hexa-

decimal numbers are also useful for defining constants with extreme values, such as the values of the macros in the `<float.h>` header. These constants are easy to write in hex but difficult to write in decimal.

**\*Q: Why do we use `%lf` to read a `double` value but `%f` to print it? [p. 134]**

A: This is a tough question to answer. First, notice that `scanf` and `printf` are unusual functions in that they aren't restricted to a fixed number of arguments. We say that `scanf` and `printf` have variable-length argument lists. When functions with variable-length argument lists are called, the compiler arranges for `float` arguments to be converted automatically to type `double`. As a result, `printf` can't distinguish between `float` and `double` arguments. This explains why `%f` works for both `float` and `double` arguments in calls of `printf`.

`scanf`, on the other hand, is passed a *pointer* to a variable. `%f` tells `scanf` to store a `float` value at the address passed to it, while `%lf` tells `scanf` to store a `double` value at that address. The distinction between `float` and `double` is crucial here. If given the wrong conversion specification, `scanf` will likely store the wrong number of bytes (not to mention the fact that the bit pattern for a `float` isn't the same as that for a `double`).

**Q: What's the proper way to pronounce `char`? [p. 134]**

A: There's no universally accepted pronunciation. Some people pronounce `char` in the same way as the first syllable of "character." Others say "char," as in `char broiled`;

**Q: When does it matter whether a character variable is signed or unsigned? [p. 136]**

A: If we store only 7-bit characters in the variable, it doesn't matter, since the sign bit will be zero. If we plan to store 8-bit characters, however, we'll probably want the variable to have `unsigned char` type. Consider the following example:

```
ch = '\xdb';
```

If `ch` has been declared to have type `char`, the compiler may choose to treat it as a signed character (many compilers do). As long as `ch` is used only as a character, there won't be any problem. But if `ch` is ever used in a context that requires the compiler to convert its value to an integer, we're likely to have trouble: the resulting integer will be negative, since `ch`'s sign bit is 1.

Here's another situation: In some kinds of programs, it's customary to use `char` variables to store one-byte integers. If we're writing such a program, we'll have to decide whether each variable should be `signed char` or `unsigned char`, just as we must decide whether ordinary integer variables should have type `int` or `unsigned int`.

**Q: I don't understand how the new-line character can be the ASCII line-feed character. When a user enters input and presses the Enter key, doesn't the program read this as a carriage-return character or a carriage return plus a line feed? [p. 137]**

A: Nope. As part of C's UNIX heritage, it always regards the end of a line as being marked by a single line-feed character. (In UNIX text files, a single line-feed character—but no carriage return—appears at the end of each line.) The C library takes care of translating the user's keypress into a line-feed character. When a program reads from a file, the I/O library translates the file's end-of-line marker (whatever it may be) into a single line-feed character. The same transformations occur—in reverse—when output is written to the screen or to a file. (See Section 22.1 for details.)

Although these translations may seem confusing, they serve an important purpose: insulating programs from details that may vary from one operating system to another.

**\*Q: What's the purpose of the \? escape sequence? [p. 138]**

A: The \? escape is related to trigraph sequences, which begin with ???. If you should put ??? in a string, there's a possibility that the compiler will mistake it for the beginning of a trigraph. Replacing the second ? by \? fixes the problem.

**Q: If getchar is faster, why would we ever want to use scanf to read individual characters? [p. 140]**

A: Although it's not as fast as getchar, the scanf function is more flexible. As we saw previously, the "%c" format string causes scanf to read the next input character; " %c" causes it to read the next non-white-space character. Also, scanf is good at reading characters that are mixed in with other kinds of data. Let's say that our input data consists of an integer, then a single nonnumeric character, then another integer. By using the format string "%d%c%d", we can get scanf to read all three items.

**\*Q: Under what circumstances do the integral promotions convert a character or short integer to unsigned int? [p. 143]**

A: The integral promotions yield an unsigned int if the int type isn't large enough to include all possible values of the original type. Since characters are usually eight bits long, they are almost always converted to int, which is guaranteed to be at least 16 bits long. Signed short integers can always be converted to int as well. Unsigned short integers are problematic. If short integers have the same length as ordinary integers (as they do on a 16-bit machine), then unsigned short integers will have to be converted to unsigned int, since the largest unsigned short integer (65,535 on a 16-bit machine) is larger than the largest int (32,767).

**Q: Exactly what happens if I assign a value to a variable that's not large enough to hold it? [p. 146]**

A: Roughly speaking, if the value is of an integral type and the variable is of an unsigned type, the extra bits are thrown away; if the variable has a signed type, the result is implementation-defined. Assigning a floating-point number to a variable—integer or floating—that's too small to hold it produces undefined behavior: anything can happen, including program termination.

**\*Q:** Why does C bother to provide type definitions? Isn't defining a `BOOL` macro just as good as defining a `Bool` type using `typedef`? [p. 149]

A: There are two important differences between type definitions and macro definitions. First, type definitions are more powerful than macro definitions. In particular, array and pointer types can't be defined as macros. Suppose that we try to use a macro to define a "pointer to integer" type:

```
#define PTR_TO_INT int *
```

The declaration

```
PTR_TO_INT p, q, r;
```

will become

```
int * p, q, r;
```

after preprocessing. Unfortunately, only `p` is a pointer; `q` and `r` are ordinary integer variables. Type definitions don't have this problem.

Second, `typedef` names are subject to the same scope rules as variables; a `typedef` name defined inside a function body wouldn't be recognized outside the function. Macro names, on the other hand, are replaced by the preprocessor wherever they appear.

**\*Q:** You said that compilers "can usually determine the value of a `sizeof` expression." Can't a compiler *always* determine the value of a `sizeof` expression? [p. 151]

A: In C89, yes. In C99, however, there's one exception. The compiler can't determine the size of a variable-length array, because the number of elements in the array may change during the execution of the program.

## Exercises

### Section 7.1

- Give the decimal value of each of the following integer constants.

- (a) 077
- (b) 0x77
- (c) 0XABC

### Section 7.2

- Which of the following are not legal constants in C? Classify each legal constant as either integer or floating-point.
  - (a) 010E2
  - (b) 32.1E+5
  - (c) 0790
  - (d) 100\_000
  - (e) 3.978e-2

- W 3. Which of the following are not legal types in C?
- short unsigned int
  - short float
  - long double
  - unsigned long
- Section 7.3**
- W 4. If `c` is a variable of type `char`, which one of the following statements is illegal?
- `i += c; /* i has type int */`
  - `c = 2 * c - 1;`
  - `putchar(c);`
  - `printf(c);`
5. Which one of the following is not a legal way to write the number 65? (Assume that the character set is ASCII.)
- 'A'
  - 0b1000001
  - 0101
  - 0x41
6. For each of the following items of data, specify which one of the types `char`, `short`, `int`, or `long` is the smallest one guaranteed to be large enough to store the item.
- Days in a month
  - Days in a year
  - Minutes in a day
  - Seconds in a day
7. For each of the following character escapes, give the equivalent octal escape. (Assume that the character set is ASCII.) You may wish to consult Appendix E, which lists the numerical codes for ASCII characters.
- \b
  - \n
  - \r
  - \t
8. Repeat Exercise 7, but give the equivalent hexadecimal escape.
- Section 7.4**
9. Suppose that `i` and `j` are variables of type `int`. What is the type of the expression `i / j + 'a'`?
- W 10. Suppose that `i` is a variable of type `int`, `j` is a variable of type `long`, and `k` is a variable of type `unsigned int`. What is the type of the expression `i + (int) j * k`?
11. Suppose that `i` is a variable of type `int`, `f` is a variable of type `float`, and `d` is a variable of type `double`. What is the type of the expression `i * f / d`?
- W 12. Suppose that `i` is a variable of type `int`, `f` is a variable of type `float`, and `d` is a variable of type `double`. Explain what conversions take place during the execution of the following statement:
- ```
d = i + f;
```

13. Assume that a program contains the following declarations:

```
char c = '\1';
short s = 2;
int i = -3;
long m = 5;
float f = 6.5f;
double d = 7.5;
```

Give the value and the type of each expression listed below.

- |             |             |                      |
|-------------|-------------|----------------------|
| (a) $c * i$ | (c) $f / c$ | (e) $f - d$          |
| (b) $s + m$ | (d) $d / s$ | (f) $(\text{int}) f$ |

- W 14. Does the following statement always compute the fractional part of  $f$  correctly (assuming that  $f$  and `frac_part` are `float` variables)?

```
frac_part = f - (\text{int}) f;
```

If not, what's the problem?

#### Section 7.5

15. Use `typedef` to create types named `Int8`, `Int16`, and `Int32`. Define the types so that they represent 8-bit, 16-bit, and 32-bit integers on your machine.

## Programming Projects

- W 1. The `square2.c` program of Section 6.3 will fail (usually by printing strange answers) if  $i * i$  exceeds the maximum `int` value. Run the program and determine the smallest value of  $n$  that causes failure. Try changing the type of  $i$  to `short` and running the program again. (Don't forget to update the conversion specifications in the call of `printf`!) Then try `long`. From these experiments, what can you conclude about the number of bits used to store integer types on your machine?
- W 2. Modify the `square2.c` program of Section 6.3 so that it pauses after every 24 squares and displays the following message:

Press Enter to continue...

After displaying the message, the program should use `getchar` to read a character. `getchar` won't allow the program to continue until the user presses the Enter key.

3. Modify the `sum2.c` program of Section 7.1 to sum a series of `double` values.  
4. Write a program that translates an alphabetic phone number into numeric form:

```
Enter phone number: CALLATT
2255288
```

(In case you don't have a telephone nearby, here are the letters on the keys: 2=ABC, 3=DEF, 4=GHI, 5=JKL, 6=MNO, 7=PRS, 8=TUV, 9=WXY.) If the original phone number contains nonalphabetic characters (digits or punctuation, for example), leave them unchanged:

```
Enter phone number: 1-800-COL-LECT
1-800-265-5328
```

You may assume that any letters entered by the user are upper case.

- W 5. In the SCRABBLE Crossword Game, players form words using small tiles, each containing a letter and a face value. The face value varies from one letter to another, based on the letter's rarity. (Here are the face values: 1: AEILNORSTU, 2: DG, 3: BCMP, 4: FHVWY, 5: K, 8: JX, 10: QZ.) Write a program that computes the value of a word by summing the values of its letters:

Enter a word: pitfall  
 Scrabble value: 12

Your program should allow any mixture of lower-case and upper-case letters in the word.  
*Hint:* Use the toupper library function.

- W 6. Write a program that prints the values of `sizeof(int)`, `sizeof(short)`, `sizeof(long)`, `sizeof(float)`, `sizeof(double)` and `sizeof(long double)`.
7. Modify Programming Project 6 from Chapter 3 so that the user may add, subtract, multiply, or divide two fractions (by entering either +, -, \*, or / between the fractions).
8. Modify Programming Project 8 from Chapter 5 so that the user enters a time using the 12-hour clock. The input will have the form *hours : minutes* followed by either A, P, AM, or PM (either lower-case or upper-case). White space is allowed (but not required) between the numerical time and the AM/PM indicator. Examples of valid input:

1:15P  
 1:15PM  
 1:15p  
 1:15pm  
 1:15 P  
 1:15 PM  
 1:15 p  
 1:15 pm

You may assume that the input has one of these forms; there is no need to test for errors.

9. Write a program that asks the user for a 12-hour time, then displays the time in 24-hour form:

Enter a 12-hour time: 9:11 PM  
 Equivalent 24-hour time: 21:11

See Programming Project 8 for a description of the input format.

10. Write a program that counts the number of vowels (*a*, *e*, *i*, *o*, and *u*) in a sentence:

Enter a sentence: And that's the way it is.  
 Your sentence contains 6 vowels.

11. Write a program that takes a first name and last name entered by the user and displays the last name, a comma, and the first initial, followed by a period:

Enter a first and last name: Lloyd Fosdick  
 Fosdick, L.

The user's input may contain extra spaces before the first name, between the first and last names, and after the last name.

12. Write a program that evaluates an expression:

Enter an expression: 1+2.5\*3  
 Value of expression: 10.5

The operands in the expression are floating-point numbers; the operators are `+`, `-`, `*`, and `/`. The expression is evaluated from left to right (no operator takes precedence over any other operator).

13. Write a program that calculates the average word length for a sentence:

Enter a sentence: It was deja vu all over again.

Average word length: 3.4

For simplicity, your program should consider a punctuation mark to be part of the word to which it is attached. Display the average word length to one decimal place.

14. Write a program that uses Newton's method to compute the square root of a positive floating-point number:

Enter a positive number: 3

Square root: 1.73205

Let  $x$  be the number entered by the user. Newton's method requires an initial guess  $y$  for the square root of  $x$  (we'll use  $y = 1$ ). Successive guesses are found by computing the average of  $y$  and  $x/y$ . The following table shows how the square root of 3 would be found:

$x$	$y$	$x/y$	Average of $y$ and $x/y$
3	1	3	2
3	2	1.5	1.75
3	1.75	1.71429	1.73214
3	1.73214	1.73196	1.73205
3	1.73205	1.73205	1.73205

Note that the values of  $y$  get progressively closer to the true square root of  $x$ . For greater accuracy, your program should use variables of type `double` rather than `float`. Have the program terminate when the absolute value of the difference between the old value of  $y$  and the new value of  $y$  is less than the product of `.00001` and  $y$ . *Hint:* Call the `fabs` function to find the absolute value of a `double`. (You'll need to include the `<math.h>` header at the beginning of your program in order to use `fabs`.)

15. Write a program that computes the factorial of a positive integer:

Enter a positive integer: 6

Factorial of 6: 720

- (a) Use a `short` variable to store the value of the factorial. What is the largest value of  $n$  for which the program correctly prints the factorial of  $n$ ?
- (b) Repeat part (a), using an `int` variable instead.
- (c) Repeat part (a), using a `long` variable instead.
- (d) Repeat part (a), using a `long long` variable instead (if your compiler supports the `long long` type).
- (e) Repeat part (a), using a `float` variable instead.
- (f) Repeat part (a), using a `double` variable instead.
- (g) Repeat part (a), using a `long double` variable instead.

In cases (e)–(g), the program will display a close approximation of the factorial, not necessarily the exact value.



# 8 Arrays

*If a program manipulates a large amount of data,  
it does so in a small number of ways.*

So far, the only variables we've seen are *scalar*: capable of holding a single data item. C also supports *aggregate* variables, which can store collections of values. There are two kinds of aggregates in C: arrays and structures. This chapter shows how to declare and use arrays, both one-dimensional (Section 8.1) and multidimensional (Section 8.2). Section 8.3 covers C99's variable-length arrays. The focus of the chapter is on one-dimensional arrays, which play a much bigger role in C than do multidimensional arrays. Later chapters (Chapter 12 in particular) provide additional information about arrays; Chapter 16 covers structures.

## 8.1 One-Dimensional Arrays

An *array* is a data structure containing a number of data values, all of which have the same type. These values, known as *elements*, can be individually selected by their position within the array.

The simplest kind of array has just one dimension. The elements of a one-dimensional array are conceptually arranged one after another in a single row (or column, if you prefer). Here's how we might visualize a one-dimensional array named `a`:



To declare an array, we must specify the *type* of the array's elements and the *number* of elements. For example, to declare that the array `a` has 10 elements of type `int`, we would write

```
int a[10];
```

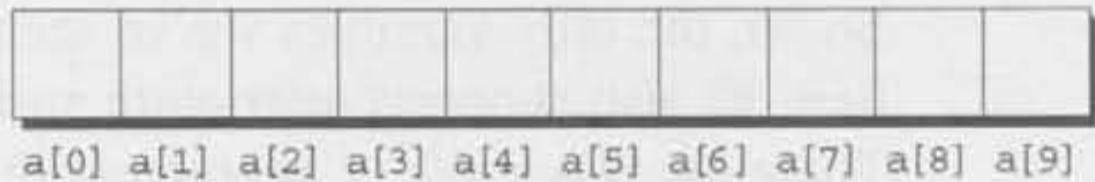
constant expressions ► 5.3

The elements of an array may be of any type; the length of the array can be specified by any (integer) constant expression. Since array lengths may need to be adjusted when the program is later changed, using a macro to define the length of an array is an excellent practice:

```
#define N 10
...
int a[N];
```

## Array Subscripting

**Q&A** To access a particular element of an array, we write the array name followed by an integer value in square brackets (this is referred to as *subscripting* or *indexing* the array). Array elements are always numbered starting from 0, so the elements of an array of length  $n$  are indexed from 0 to  $n - 1$ . For example, if  $a$  is an array with 10 elements, they're designated by  $a[0]$ ,  $a[1]$ , ...,  $a[9]$ , as the following figure shows:



lvalues ► 4.2 Expressions of the form  $a[i]$  are lvalues, so they can be used in the same way as ordinary variables:

```
a[0] = 1;
printf("%d\n", a[5]);
++a[i];
```

In general, if an array contains elements of type  $T$ , then each element of the array is treated as if it were a variable of type  $T$ . In this example, the elements  $a[0]$ ,  $a[5]$ , and  $a[i]$  behave like `int` variables.

Arrays and `for` loops go hand-in-hand. Many programs contain `for` loops whose job is to perform some operation on every element in an array. Here are a few examples of typical operations on an array  $a$  of length  $N$ :

<b>idiom</b>	<code>for (i = 0; i &lt; N; i++)</code>	
	<code>    a[i] = 0;</code>	<code>                /* clears a */</code>
<b>idiom</b>	<code>for (i = 0; i &lt; N; i++)</code>	
	<code>    scanf("%d", &amp;a[i]);</code>	<code>                /* reads data into a */</code>
<b>idiom</b>	<code>for (i = 0; i &lt; N; i++)</code>	
	<code>    sum += a[i];</code>	<code>                /* sums the elements of a */</code>

Notice that we must use the `&` symbol when calling `scanf` to read an array element, just as we would with an ordinary variable.



C doesn't require that subscript bounds be checked; if a subscript goes out of range, the program's behavior is undefined. One cause of a subscript going out of bounds: forgetting that an array with  $n$  elements is indexed from 0 to  $n - 1$ , not 1 to  $n$ . (As one of my professors liked to say, "In this business, you're always off by one." He was right, of course.) The following example illustrates a bizarre effect that can be caused by this common blunder:

```
int a[10], i;

for (i = 1; i <= 10; i++)
    a[i] = 0;
```

With some compilers, this innocent-looking `for` statement causes an infinite loop! When `i` reaches 10, the program stores 0 into `a[10]`. But `a[10]` doesn't exist, so 0 goes into memory immediately after `a[9]`. If the variable `i` happens to follow `a[9]` in memory—as might be the case—then `i` will be reset to 0, causing the loop to start over.

---

An array subscript may be any integer expression:

```
a[i+j*10] = 0;
```

The expression can even have side effects:

```
i = 0;
while (i < N)
    a[i++] = 0;
```

Let's trace this code. After `i` is set to 0, the `while` statement checks whether `i` is less than `N`. If it is, 0 is assigned to `a[0]`, `i` is incremented, and the loop repeats. Note that `a[++i]` wouldn't be right, because 0 would be assigned to `a[1]` during the first loop iteration.




---

Be careful when an array subscript has a side effect. For example, the following loop—which is supposed to copy the elements of the array `b` into the array `a`—may not work properly:

```
i = 0;
while (i < N)
    a[i] = b[i++];
```

The expression `a[i] = b[i++]` accesses the value of `i` and also modifies `i` elsewhere in the expression, which—as we saw in Section 4.4—causes undefined behavior. Of course, we can easily avoid the problem by removing the increment from the subscript:

```
for (i = 0; i < N; i++)
    a[i] = b[i];
```

---

## PROGRAM Reversing a Series of Numbers

Our first array program prompts the user to enter a series of numbers, then writes the numbers in reverse order:

```
Enter 10 numbers: 34 82 49 102 7 94 23 11 50 31
In reverse order: 31 50 11 23 94 7 102 49 82 34
```

Our strategy will be to store the numbers in an array as they're read, then go through the array backwards, printing the elements one by one. In other words, we won't actually reverse the elements in the array, but we'll make the user think we did.

```
reverse.c /* Reverses a series of numbers */

#include <stdio.h>

#define N 10

int main(void)
{
    int a[N], i;

    printf("Enter %d numbers: ", N);
    for (i = 0; i < N; i++)
        scanf("%d", &a[i]);

    printf("In reverse order:");
    for (i = N - 1; i >= 0; i--)
        printf(" %d", a[i]);
    printf("\n");

    return 0;
}
```

This program shows just how useful macros can be in conjunction with arrays. The macro `N` is used four times in the program: in the declaration of `a`, in the `printf` that displays a prompt, and in both `for` loops. Should we later decide to change the size of the array, we need only edit the definition of `N` and recompile the program. Nothing else will need to be altered; even the prompt will still be correct.

## Array Initialization

An array, like any other variable, can be given an initial value at the time it's declared. The rules are somewhat tricky, though, so we'll cover some of them now and save others until later.

The most common form of *array initializer* is a list of constant expressions enclosed in braces and separated by commas:

```
int a[10] = {1, 2, 3, 4, 5, 6, 7, 8, 9, 10};
```

If the initializer is *shorter* than the array, the remaining elements of the array are given the value 0:

```
int a[10] = {1, 2, 3, 4, 5, 6};
/* initial value of a is {1, 2, 3, 4, 5, 6, 0, 0, 0, 0} */
```

Using this feature, we can easily initialize an array to all zeros:

```
int a[10] = {0};
/* initial value of a is {0, 0, 0, 0, 0, 0, 0, 0, 0, 0} */
```

It's illegal for an initializer to be completely empty, so we've put a single 0 inside the braces. It's also illegal for an initializer to be *longer* than the array it initializes.

If an initializer is present, the length of the array may be omitted:

```
int a[] = {1, 2, 3, 4, 5, 6, 7, 8, 9, 10};
```

The compiler uses the length of the initializer to determine how long the array is. The array still has a fixed number of elements (10, in this example), just as if we had specified the length explicitly.

**C99**

## Designated Initializers

It's often the case that relatively few elements of an array need to be initialized explicitly; the other elements can be given default values. Consider the following example:

```
int a[15] = {0, 0, 29, 0, 0, 0, 0, 0, 0, 7, 0, 0, 0, 0, 48};
```

We want element 2 of the array to be 29, element 9 to be 7, and element 14 to be 48, but the other values are just zero. For a large array, writing an initializer in this fashion is tedious and error-prone (what if there were 200 zeros between two of the nonzero values?).

C99's *designated initializers* can be used to solve this problem. Here's how we could redo the previous example using a designated initializer:

```
int a[15] = {[2] = 29, [9] = 7, [14] = 48};
```

Each number in brackets is said to be a *designator*.

Besides being shorter and easier to read (at least for some arrays), designated initializers have another advantage: the order in which the elements are listed no longer matters. Thus, our previous example could also be written in the following way:

```
int a[15] = {[14] = 48, [9] = 7, [2] = 29};
```

Designators must be integer constant expressions. If the array being initialized has length  $n$ , each designator must be between 0 and  $n - 1$ . However, if the length of the array is omitted, a designator can be any nonnegative integer. In the latter case, the compiler will deduce the length of the array from the largest designator.

In the following example, the fact that 23 appears as a designator will force the array to have length 24:

```
int b[] = { [5] = 10, [23] = 13, [11] = 36, [15] = 29 };
```

An initializer may use both the older (element-by-element) technique and the newer (designated) technique:

```
int c[10] = { 5, 1, 9, [4] = 3, 7, 2, [8] = 6 };
```

**Q&A**

This initializer specifies that the array's first three elements will be 5, 1, and 9. Element 4 will have the value 3. The two elements after element 4 will be 7 and 2. Finally, element 8 will have the value 6. All elements for which no value is specified will default to zero.

**PROGRAM****Checking a Number for Repeated Digits**

Our next program checks whether any of the digits in a number appear more than once. After the user enters a number, the program prints either **Repeated digit** or **No repeated digit**:

```
Enter a number: 28212
Repeated digit
```

The number 28212 has a repeated digit (2); a number like 9357 doesn't.

The program uses an array of Boolean values to keep track of which digits appear in a number. The array, named `digit_seen`, is indexed from 0 to 9 to correspond to the 10 possible digits. Initially, every element of the array is false. (The initializer for `digit_seen` is `{false}`, which only initializes the first element of the array. However, the compiler will automatically make the remaining elements zero, which is equivalent to false.)

When given a number `n`, the program examines `n`'s digits one at a time, storing each into the `digit` variable and then using it as an index into `digit_seen`. If `digit_seen[digit]` is true, then `digit` appears at least twice in `n`. On the other hand, if `digit_seen[digit]` is false, then `digit` has not been seen before, so the program sets `digit_seen[digit]` to true and keeps going.

```
repdigit.c /* Checks numbers for repeated digits */

#include <stdbool.h> /* C99 only */
#include <stdio.h>

int main(void)
{
    bool digit_seen[10] = {false};
    int digit;
    long n;

    printf("Enter a number: ");
    scanf("%ld", &n);
```

```

while (n > 0) {
    digit = n % 10;
    if (digit_seen[digit])
        break;
    digit_seen[digit] = true;
    n /= 10;
}

if (n > 0)
    printf("Repeated digit\n");
else
    printf("No repeated digit\n");

return 0;
}

```

**C99**

&lt;stdbool.h&gt; header ▶ 21.5

This program uses the names `bool`, `true`, and `false`, which are defined in C99's `<stdbool.h>` header. If your compiler doesn't support this header, you'll need to define these names yourself. One way to do so is to put the following lines above the `main` function:

```

#define true 1
#define false 0
typedef int bool;

```

Notice that `n` has type `long`, allowing the user to enter numbers up to 2,147,483,647 (or more, on some machines).

## Using the `sizeof` Operator with Arrays

The `sizeof` operator can determine the size of an array (in bytes). If `a` is an array of 10 integers, then `sizeof(a)` is typically 40 (assuming that each integer requires four bytes).

We can also use `sizeof` to measure the size of an array element, such as `a[0]`. Dividing the array size by the element size gives the length of the array:

```
sizeof(a) / sizeof(a[0])
```

Some programmers use this expression when the length of the array is needed. To clear the array `a`, for example, we could write

```
for (i = 0; i < sizeof(a) / sizeof(a[0]); i++)
    a[i] = 0;
```

With this technique, the loop doesn't have to be modified if the array length should change at a later date. Using a macro to represent the array length has the same advantage, of course, but the `sizeof` technique is slightly better, since there's no macro name to remember (and possibly get wrong).

One minor annoyance is that some compilers produce a warning message for the expression `i < sizeof(a) / sizeof(a[0])`. The variable `i` probably has

type `int` (a signed type), whereas `sizeof` produces a value of type `size_t` (an unsigned type). We know from Section 7.4 that comparing a signed integer with an unsigned integer is a dangerous practice, although in this case it's safe because both `i` and `sizeof(a) / sizeof(a[0])` have nonnegative values. To avoid a warning, we can add a cast that converts `sizeof(a) / sizeof(a[0])` to a signed integer:

```
for (i = 0; i < (int) (sizeof(a) / sizeof(a[0])); i++)
    a[i] = 0;
```

Writing `(int) (sizeof(a) / sizeof(a[0]))` is a bit unwieldy; defining a macro that represents it is often helpful:

```
#define SIZE ((int) (sizeof(a) / sizeof(a[0])))

for (i = 0; i < SIZE; i++)
    a[i] = 0;
```

If we're back to using a macro, though, what's the advantage of `sizeof`? We'll answer that question in a later chapter (the trick is to add a parameter to the macro).

## PROGRAM Computing Interest

Our next program prints a table showing the value of \$100 invested at different rates of interest over a period of years. The user will enter an interest rate and the number of years the money will be invested. The table will show the value of the money at one-year intervals—at that interest rate and the next four higher rates—assuming that interest is compounded once a year. Here's what a session with the program will look like:

```
Enter interest rate: 6
Enter number of years: 5
```

Years	6%	7%	8%	9%	10%
1	106.00	107.00	108.00	109.00	110.00
2	112.36	114.49	116.64	118.81	121.00
3	119.10	122.50	125.97	129.50	133.10
4	126.25	131.08	136.05	141.16	146.41
5	133.82	140.26	146.93	153.86	161.05

Clearly, we can use a `for` statement to print the first row. The second row is a little trickier, since its values depend on the numbers in the first row. Our solution is to store the first row in an array as it's computed, then use the values in the array to compute the second row. Of course, this process can be repeated for the third and later rows. We'll end up with two `for` statements, one nested inside the other. The outer loop will count from 1 to the number of years requested by the user. The inner loop will increment the interest rate from its lowest value to its highest value.

```

interest.c /* Prints a table of compound interest */

#include <stdio.h>

#define NUM_RATES ((int) (sizeof(value) / sizeof(value[0])))
#define INITIAL_BALANCE 100.00

int main(void)
{
    int i, low_rate, num_years, year;
    double value[5];

    printf("Enter interest rate: ");
    scanf("%d", &low_rate);
    printf("Enter number of years: ");
    scanf("%d", &num_years);

    printf("\nYears");
    for (i = 0; i < NUM_RATES; i++) {
        printf("%6d%%", low_rate + i);
        value[i] = INITIAL_BALANCE;
    }
    printf("\n");

    for (year = 1; year <= num_years; year++) {
        printf("%3d    ", year);
        for (i = 0; i < NUM_RATES; i++) {
            value[i] += (low_rate + i) / 100.0 * value[i];
            printf("%7.2f", value[i]);
        }
        printf("\n");
    }

    return 0;
}

```

Note the use of `NUM_RATES` to control two of the `for` loops. If we later change the size of the `value` array, the loops will adjust automatically.

## 8.2 Multidimensional Arrays

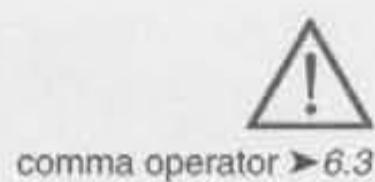
An array may have any number of dimensions. For example, the following declaration creates a two-dimensional array (a *matrix*, in mathematical terminology):

```
int m[5][9];
```

The array `m` has 5 rows and 9 columns. Both rows and columns are indexed from 0, as the following figure shows:

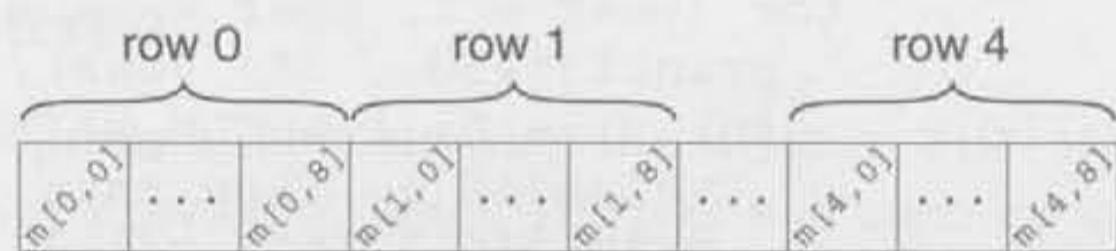
	0	1	2	3	4	5	6	7	8
0									
1									
2									
3									
4									

To access the element of  $m$  in row  $i$ , column  $j$ , we must write  $m[i][j]$ . The expression  $m[i]$  designates row  $i$  of  $m$ , and  $m[i][j]$  then selects element  $j$  in this row.



Resist the temptation to write  $m[i, j]$  instead of  $m[i][j]$ . C treats the comma as an operator in this context, so  $m[i, j]$  is the same as  $m[j]$ .

Although we visualize two-dimensional arrays as tables, that's not the way they're actually stored in computer memory. C stores arrays in **row-major order**, with row 0 first, then row 1, and so forth. For example, here's how the  $m$  array is stored:



We'll usually ignore this detail, but sometimes it will affect our code.

Just as `for` loops go hand-in-hand with one-dimensional arrays, nested `for` loops are ideal for processing multidimensional arrays. Consider, for example, the problem of initializing an array for use as an identity matrix. (In mathematics, an *identity matrix* has 1's on the main diagonal, where the row and column index are the same, and 0's everywhere else.) We'll need to visit each element in the array in some systematic fashion. A pair of nested `for` loops—one that steps through every row index and one that steps through each column index—is perfect for the job:

```
#define N 10

double ident[N][N];
int row, col;

for (row = 0; row < N; row++)
    for (col = 0; col < N; col++)
        if (row == col)
            ident[row][col] = 1.0;
        else
            ident[row][col] = 0.0;
```

Multidimensional arrays play a lesser role in C than in many other programming languages, primarily because C provides a more flexible way to store multidimensional data: arrays of pointers.

## Initializing a Multidimensional Array

We can create an initializer for a two-dimensional array by nesting one-dimensional initializers:

```
int m[5][9] = {{1, 1, 1, 1, 1, 0, 1, 1, 1},
                {0, 1, 0, 1, 0, 1, 0, 1, 0},
                {0, 1, 0, 1, 1, 0, 0, 1, 0},
                {1, 1, 0, 1, 0, 0, 0, 1, 0},
                {1, 1, 0, 1, 0, 0, 1, 1, 1}};
```

Each inner initializer provides values for one row of the matrix. Initializers for higher-dimensional arrays are constructed in a similar fashion.

C provides a variety of ways to abbreviate initializers for multidimensional arrays:

- If an initializer isn't large enough to fill a multidimensional array, the remaining elements are given the value 0. For example, the following initializer fills only the first three rows of *m*; the last two rows will contain zeros:

```
int m[5][9] = {{1, 1, 1, 1, 1, 0, 1, 1, 1},
                {0, 1, 0, 1, 0, 1, 0, 1, 0},
                {0, 1, 0, 1, 1, 0, 0, 1, 0}};
```

- If an inner list isn't long enough to fill a row, the remaining elements in the row are initialized to 0:

```
int m[5][9] = {{1, 1, 1, 1, 1, 0, 1, 1, 1},
                {0, 1, 0, 1, 0, 1, 0, 1, 1},
                {0, 1, 0, 1, 1, 0, 0, 1, 1},
                {1, 1, 0, 1, 0, 0, 0, 1, 1},
                {1, 1, 0, 1, 0, 0, 1, 1, 1}};
```

- We can even omit the inner braces:

```
int m[5][9] = {1, 1, 1, 1, 1, 0, 1, 1, 1,
               0, 1, 0, 1, 0, 1, 0, 1, 0,
               0, 1, 0, 1, 1, 0, 0, 1, 0,
               1, 1, 0, 1, 0, 0, 0, 1, 0,
               1, 1, 0, 1, 0, 0, 1, 1, 1};
```

Once the compiler has seen enough values to fill one row, it begins filling the next.



Omitting the inner braces in a multidimensional array initializer can be risky, since an extra element (or even worse, a missing element) will affect the rest of the initializer. Leaving out the braces causes some compilers to produce a warning message such as “*missing braces around initializer*.”

**C99**

C99’s designated initializers work with multidimensional arrays. For example, we could create a  $2 \times 2$  identity matrix as follows:

```
double ident[2][2] = {[0][0] = 1.0, [1][1] = 1.0};
```

As usual, all elements for which no value is specified will default to zero.

## Constant Arrays

Any array, whether one-dimensional or multidimensional, can be made “constant” by starting its declaration with the word `const`:

```
const char hex_chars[] =
{'0', '1', '2', '3', '4', '5', '6', '7', '8', '9',
'A', 'B', 'C', 'D', 'E', 'F'};
```

An array that’s been declared `const` should not be modified by the program; the compiler will detect direct attempts to modify an element.

Declaring an array to be `const` has a couple of primary advantages. It documents that the program won’t change the array, which can be valuable information for someone reading the code later. It also helps the compiler catch errors, by informing it that we don’t intend to modify the array.

`const` type qualifier ▶ 18.3

`const` isn’t limited to arrays; it works with any variable, as we’ll see later. However, `const` is particularly useful in array declarations, because arrays may contain reference information that won’t change during program execution.

## PROGRAM Dealing a Hand of Cards

Our next program illustrates both two-dimensional arrays and constant arrays. The program deals a random hand from a standard deck of playing cards. (In case you haven’t had time to play games recently, each card in a standard deck has a *suit*—clubs, diamonds, hearts, or spades—and a *rank*—two, three, four, five, six, seven, eight, nine, ten, jack, queen, king, or ace.) We’ll have the user specify how many cards should be in the hand:

```
Enter number of cards in hand: 5
Your hand: 7c 2s 5d as 2h
```

It’s not immediately obvious how we’d write such a program. How do we pick cards randomly from the deck? And how do we avoid picking the same card twice? Let’s tackle these problems separately.

`time` function ▶ 26.3

`srand` function ▶ 26.2

`rand` function ▶ 26.2

To pick cards randomly, we’ll use several C library functions. The `time` function (from `<time.h>`) returns the current time, encoded in a single number. The `srand` function (from `<stdlib.h>`) initializes C’s random number generator. Passing the return value of `time` to `srand` prevents the program from dealing the same cards every time we run it. The `rand` function (also from `<stdlib.h>`) produces an apparently random number each time it’s called. By using the `%` operator, we can scale the return value from `rand` so that it falls between 0 and 3 (for suits) or between 0 and 12 (for ranks).

To avoid picking the same card twice, we’ll need to keep track of which cards have already been chosen. For that purpose, we’ll use an array named `in_hand`

that has four rows (one for each suit) and 13 columns (one for each rank). In other words, each element in the array corresponds to one of the 52 cards in the deck. All elements of the array will be false to start with. Each time we pick a card at random, we'll check whether the element of `in_hand` corresponding to that card is true or false. If it's true, we'll have to pick another card. If it's false, we'll store `true` in that card's array element to remind us later that this card has already been picked.

Once we've verified that a card is "new"—not already selected—we'll need to translate its numerical rank and suit into characters and then display the card. To translate the rank and suit to character form, we'll set up two arrays of characters—one for the rank and one for the suit—and then use the numbers to subscript the arrays. These arrays won't change during program execution, so we may as well declare them to be `const`.

```
deal.c /* Deals a random hand of cards */

#include <stdbool.h>    /* C99 only */
#include <stdio.h>
#include <stdlib.h>
#include <time.h>

#define NUM_SUITS 4
#define NUM_RANKS 13

int main(void)
{
    bool in_hand[NUM_SUITS][NUM_RANKS] = {false};
    int num_cards, rank, suit;
    const char rank_code[] = {'2', '3', '4', '5', '6', '7', '8',
                             '9', 't', 'j', 'q', 'k', 'a'};
    const char suit_code[] = {'c', 'd', 'h', 's'};

    srand((unsigned) time(NULL));

    printf("Enter number of cards in hand: ");
    scanf("%d", &num_cards);

    printf("Your hand:");
    while (num_cards > 0) {
        suit = rand() % NUM_SUITS;      /* picks a random suit */
        rank = rand() % NUM_RANKS;      /* picks a random rank */
        if (!in_hand[suit][rank]) {
            in_hand[suit][rank] = true;
            num_cards--;
            printf(" %c%c", rank_code[rank], suit_code[suit]);
        }
    }
    printf("\n");
}

return 0;
}
```

Notice the initializer for the `in_hand` array:

```
bool in_hand[NUM_SUITS][NUM_RANKS] = {false};
```

Even though `in_hand` is a two-dimensional array, we can use a single pair of braces (at the risk of possibly incurring a warning from the compiler). Also, we've supplied only one value in the initializer, knowing that the compiler will fill in 0 (false) for the other elements.

## 8.3 Variable-Length Arrays (C99)

Section 8.1 stated that the length of an array variable must be specified by a constant expression. In C99, however, it's sometimes possible to use an expression that's *not* constant. The following modification of the `reverse.c` program (Section 8.1) illustrates this ability:

```
reverse2.c /* Reverses a series of numbers using a variable-length
           array - C99 only */

#include <stdio.h>

int main(void)
{
    int i, n;

    printf("How many numbers do you want to reverse? ");
    scanf("%d", &n);

    int a[n]; /* C99 only - length of array depends on n */

    printf("Enter %d numbers: ", n);
    for (i = 0; i < n; i++)
        scanf("%d", &a[i]);

    printf("In reverse order:");
    for (i = n - 1; i >= 0; i--)
        printf(" %d", a[i]);
    printf("\n");

    return 0;
}
```

The array `a` in this program is an example of a *variable-length array* (or *VLA* for short). The length of a VLA is computed when the program is executed, not when the program is compiled. The chief advantage of a VLA is that the programmer doesn't have to pick an arbitrary length when declaring an array; instead, the program itself can calculate exactly how many elements are needed. If the programmer makes the choice, it's likely that the array will be too long (wasting memory) or too short (causing the program to fail). In the `reverse2.c` program, the num-

ber entered by the user determines the length of `a`; the programmer doesn't have to choose a fixed length, unlike in the original version of the program.

The length of a VLA doesn't have to be specified by a single variable. Arbitrary expressions, possibly containing operators, are also legal. For example:

```
int a[3*i+5];
int b[j+k];
```

Like other arrays, VLAs can be multidimensional:

```
int c[m][n];
```

static storage duration ➤ 18.2

The primary restriction on VLAs is that they can't have static storage duration. (We haven't yet seen any arrays with this property.) Another restriction is that a VLA may not have an initializer.

Variable-length arrays are most often seen in functions other than `main`. One big advantage of a VLA that belongs to a function `f` is that it can have a different length each time `f` is called. We'll explore this feature in Section 9.3.

## Q & A

**Q: Why do array subscripts start at 0 instead of 1? [p. 162]**

A: Having subscripts begin at 0 simplifies the compiler a bit. Also, it can make array subscripting marginally faster.

**Q: What if I want an array with subscripts that go from 1 to 10 instead of 0 to 9?**

A: Here's a common trick: declare the array to have 11 elements instead of 10. The subscripts will go from 0 to 10, but you can just ignore element 0.

**Q: Is it possible to use a character as an array subscript?**

A: Yes, because C treats characters as integers. You'll probably need to "scale" the character before you use it as a subscript, though. Let's say that we want the `letter_count` array to keep track of a count for each letter in the alphabet. The array will need 26 elements, so we'd declare it in the following way:

```
int letter_count[26];
```

However, we can't use letters to subscript `letter_count` directly, because their integer values don't fall between 0 and 25. To scale a lower-case letter to the proper range, we can simply subtract '`'a'`'; to scale an upper-case letter, we'll subtract '`'A'`'. For example, if `ch` contains a lower-case letter, we'd write

```
letter_count[ch - 'a'] = 0;
```

to clear the count that corresponds to `ch`. A minor caveat: this technique isn't completely portable, because it assumes that letters have consecutive codes. However, it works with most character sets, including ASCII.

**Q:** It seems like a designated initializer could end up initializing an array element more than once. Consider the following array declaration:

```
int a[] = {4, 9, 1, 8, [0] = 5, 7};
```

**Is this declaration legal, and if so, what is the length of the array? [p. 166]**

**A:** Yes, the declaration is legal. Here's how it works: as it processes an initializer list, the compiler keeps track of which array element is to be initialized next. Normally, the next element is the one following the element that was last initialized. However, when a designator appears in the list, it forces the next element be the one represented by the designator, *even if that element has already been initialized*.

Here's a step-by-step look at how the compiler will process the initializer for the array `a`:

- The 4 initializes element 0; the next element to be initialized is element 1.
- The 9 initializes element 1; the next element to be initialized is element 2.
- The 1 initializes element 2; the next element to be initialized is element 3.
- The 8 initializes element 3; the next element to be initialized is element 4.
- The [0] designator causes the next element to become 0, so the 5 initializes element 0 (replacing the 4 previously stored there). The next element to be initialized is element 1.
- The 7 initializes element 1 (replacing the 9 previously stored there). The next element to be initialized is element 2 (which is irrelevant since we're at the end of the list).

The net effect is the same as if we had written

```
int a[] = {5, 7, 1, 8};
```

Thus, the length of this array is four.

**Q:** The compiler gives me an error message if I try to copy one array into another by using the assignment operator. What's wrong?

**A:** Although it looks quite plausible, the assignment

```
a = b; /* a and b are arrays */
```

is indeed illegal. The reason for its illegality isn't obvious; it has to do with the peculiar relationship between arrays and pointers in C, a topic we'll explore in Chapter 12.

The simplest way to copy one array into another is to use a loop that copies the elements, one by one:

```
for (i = 0; i < N; i++)
    a[i] = b[i];
```

memcpy function ►23.6

Another possibility is to use the `memcpy` ("memory copy") function from the `<string.h>` header. `memcpy` is a low-level function that simply copies bytes from one place to another. To copy the array `b` into the array `a`, use `memcpy` as follows:

```
memcpy(a, b, sizeof(a));
```

Many programmers prefer `memcpy`, especially for large arrays, because it's potentially faster than an ordinary loop.

**\*Q:** Section 6.4 mentioned that C99 doesn't allow a `goto` statement to bypass the declaration of a variable-length array. What's the reason for this restriction?

A: The memory used to store a variable-length array is usually allocated when the declaration of the array is reached during program execution. Bypassing the declaration using a `goto` statement could result in a program accessing the elements of an array that was never allocated.

## Exercises

### Section 8.1

- W 1. We discussed using the expression `sizeof(a) / sizeof(a[0])` to calculate the number of elements in an array. The expression `sizeof(a) / sizeof(t)`, where `t` is the type of `a`'s elements, would also work, but it's considered an inferior technique. Why?
- W 2. The Q&A section shows how to use a *letter* as an array subscript. Describe how to use a *digit* (in character form) as a subscript.
- 3. Write a declaration of an array named `weekend` containing seven `bool` values. Include an initializer that makes the first and last values `true`; all other values should be `false`.
- 4. (C99) Repeat Exercise 3, but this time use a designated initializer. Make the initializer as short as possible.
- 5. The Fibonacci numbers are 0, 1, 1, 2, 3, 5, 8, 13, ..., where each number is the sum of the two preceding numbers. Write a program fragment that declares an array named `fib_numbers` of length 40 and fills the array with the first 40 Fibonacci numbers. Hint: Fill in the first two numbers individually, then use a loop to compute the remaining numbers.

### Section 8.2

- 6. Calculators, watches, and other electronic devices often rely on seven-segment displays for numerical output. To form a digit, such devices "turn on" some of the seven segments while leaving others "off":



Suppose that we want to set up an array that remembers which segments should be "on" for each digit. Let's number the segments as follows:

5	0
4	1
3	2

Here's what the array might look like, with each row representing one digit:

```
const int segments[10][7] = {{1, 1, 1, 1, 1, 1, 0}, ...};
```

I've given you the first row of the initializer; fill in the rest.

- W 7. Using the shortcuts described in Section 8.2, shrink the initializer for the `segments` array (Exercise 6) as much as you can.
8. Write a declaration for a two-dimensional array named `temperature_readings` that stores one month of hourly temperature readings. (For simplicity, assume that a month has 30 days.) The rows of the array should represent days of the month; the columns should represent hours of the day.
9. Using the array of Exercise 8, write a program fragment that computes the average temperature for a month (averaged over all days of the month and all hours of the day).
10. Write a declaration for an  $8 \times 8$  `char` array named `chess_board`. Include an initializer that puts the following data into the array (one character per array element):

```
r n b q k b n r
p p p p p p p p
.
.
.
.
.
.
p p p p p p p p
R N B Q K B N R
```

11. Write a program fragment that declares an  $8 \times 8$  `char` array named `checker_board` and then uses a loop to store the following data into the array (one character per array element):

```
B R B R B R B R
R B R B R B R B
B R B R B R B R
R B R B R B R B
B R B R B R B R
R B R B R B R B
B R B R B R B R
R B R B R B R B
```

*Hint:* The element in row  $i$ , column  $j$ , should be the letter B if  $i + j$  is an even number.

## Programming Projects

1. Modify the `repdigit.c` program of Section 8.1 so that it shows which digits (if any) were repeated:

```
Enter a number: 939577
Repeated digit(s): 7 9
```

- W 2. Modify the `repdigit.c` program of Section 8.1 so that it prints a table showing how many times each digit appears in the number:

```
Enter a number: 41271092
Digit:      0   1   2   3   4   5   6   7   8   9
Occurrences: 1   2   2   0   1   0   0   1   0   1
```

3. Modify the `repdigit.c` program of Section 8.1 so that the user can enter more than one number to be tested for repeated digits. The program should terminate when the user enters a number that's less than or equal to 0.

4. Modify the `reverse.c` program of Section 8.1 to use the expression `(int)(sizeof(a) / sizeof(a[0]))` (or a macro with this value) for the array length.
- W 5. Modify the `interest.c` program of Section 8.1 so that it compounds interest *monthly* instead of *annually*. The form of the output shouldn't change; the balance should still be shown at annual intervals.
6. The prototypical Internet newbie is a fellow named B1FF, who has a unique way of writing messages. Here's a typical B1FF communiqué:

H3Y DUD3, C 15 R1LLY COOL!!!!!!

Write a “B1FF filter” that reads a message entered by the user and translates it into B1FF-speak:

Enter message: Hey dude, C is rilly cool

In B1FF-speak: H3Y DUD3, C 15 R1LLY COOL!!!!!!

Your program should convert the message to upper-case letters, substitute digits for certain letters (A→4, B→8, E→3, I→1, O→0, S→5), and then append 10 or so exclamation marks. *Hint:* Store the original message in an array of characters, then go back through the array, translating and printing characters one by one.

7. Write a program that reads a  $5 \times 5$  array of integers and then prints the row sums and the column sums:

```
Enter row 1: 8 3 9 0 10
Enter row 2: 3 5 17 1 1
Enter row 3: 2 8 6 23 1
Enter row 4: 15 7 3 2 9
Enter row 5: 6 14 2 6 0
```

Row totals: 30 27 40 36 28

Column totals: 34 37 37 32 21

- W 8. Modify Programming Project 7 so that it prompts for five quiz grades for each of five students, then computes the total score and average score for each *student*, and the average score, high score, and low score for each *quiz*.
9. Write a program that generates a “random walk” across a  $10 \times 10$  array. The array will contain characters (all ‘.’ initially). The program must randomly “walk” from element to element, always going up, down, left, or right by one element. The elements visited by the program will be labeled with the letters A through Z, in the order visited. Here's an example of the desired output:

```
A . . . . . .
B C D . . . . .
. F E . . . . .
H G . . . . .
I . . . . .
J . . . . . Z .
K . . R S T U V Y .
L M P Q . . W X .
. N O . . . . .
. . . . . .
```

*Hint:* Use the `srand` and `rand` functions (see `deal.c`) to generate random numbers. After generating a number, look at its remainder when divided by 4. There are four possible values for the remainder—0, 1, 2, and 3—indicating the direction of the next move. Before performing a move, check that (a) it won't go outside the array, and (b) it doesn't take us to

an element that already has a letter assigned. If either condition is violated, try moving in another direction. If all four directions are blocked, the program must terminate. Here's an example of premature termination:

```
A B G H I . . . .
. C F . J K . . .
. D E . M L . . .
. . . N O . . .
. . W X Y P Q . .
. . V U T S R . .
. . . . . . . .
. . . . . . . .
. . . . . . . .
. . . . . . . .
```

Y is blocked on all four sides, so there's no place to put Z.

10. Modify Programming Project 8 from Chapter 5 so that the departure times are stored in an array and the arrival times are stored in a second array. (The times are integers, representing the number of minutes since midnight.) The program will use a loop to search the array of departure times for the one closest to the time entered by the user.
11. Modify Programming Project 4 from Chapter 7 so that the program labels its output:

```
Enter phone number: 1-800-COL-LECT
In numeric form: 1-800-265-5328
```

The program will need to store the phone number (either in its original form or in its numeric form) in an array of characters until it can be printed. You may assume that the phone number is no more than 15 characters long.

12. Modify Programming Project 5 from Chapter 7 so that the SCRABBLE values of the letters are stored in an array. The array will have 26 elements, corresponding to the 26 letters of the alphabet. For example, element 0 of the array will store 1 (because the SCRABBLE value of the letter A is 1), element 1 of the array will store 3 (because the SCRABBLE value of the letter B is 3), and so forth. As each character of the input word is read, the program will use the array to determine the SCRABBLE value of that character. Use an array initializer to set up the array.
13. Modify Programming Project 11 from Chapter 7 so that the program labels its output:

```
Enter a first and last name: Lloyd Fosdick
You entered the name: Fosdick, L.
```

The program will need to store the last name (but not the first name) in an array of characters until it can be printed. You may assume that the last name is no more than 20 characters long.

14. Write a program that reverses the words in a sentence:

```
Enter a sentence: you can cage a swallow can't you?
Reversal of sentence: you can't swallow a cage can you?
```

*Hint:* Use a loop to read the characters one by one and store them in a one-dimensional char array. Have the loop stop at a period, question mark, or exclamation point (the "terminating character"), which is saved in a separate char variable. Then use a second loop to search backward through the array for the beginning of the last word. Print the last word, then search backward for the next-to-last word. Repeat until the beginning of the array is reached. Finally, print the terminating character.

15. One of the oldest known encryption techniques is the Caesar cipher, attributed to Julius Caesar. It involves replacing each letter in a message with another letter that is a fixed number of

positions later in the alphabet. (If the replacement would go past the letter Z, the cipher “wraps around” to the beginning of the alphabet. For example, if each letter is replaced by the letter two positions after it, then Y would be replaced by A, and Z would be replaced by B.) Write a program that encrypts a message using a Caesar cipher. The user will enter the message to be encrypted and the shift amount (the number of positions by which letters should be shifted):

```
Enter message to be encrypted: Go ahead, make my day.  
Enter shift amount (1-25): 3  
Encrypted message: Jr dkhdg, pdnh pb gdb.
```

Notice that the program can decrypt a message if the user enters 26 minus the original key:

```
Enter message to be encrypted: Jr dkhdg, pdnh pb gdb.  
Enter shift amount (1-25): 23  
Encrypted message: Go ahead, make my day.
```

You may assume that the message does not exceed 80 characters. Characters other than letters should be left unchanged. Lower-case letters remain lower-case when encrypted, and upper-case letters remain upper-case. *Hint:* To handle the wrap-around problem, use the expression  $((ch - 'A') + n) \% 26 + 'A'$  to calculate the encrypted version of an upper-case letter, where ch stores the letter and n stores the shift amount. (You'll need a similar expression for lower-case letters.)

16. Write a program that tests whether two words are anagrams (permutations of the same letters):

```
Enter first word: smartest  
Enter second word: mattress  
The words are anagrams.
```

```
Enter first word: dumbest  
Enter second word: stumble  
The words are not anagrams.
```

Write a loop that reads the first word, character by character, using an array of 26 integers to keep track of how many times each letter has been seen. (For example, after the word *smartest* has been read, the array should contain the values 1 0 0 0 1 0 0 0 0 0 0 1 0 0 0 0 1 2 2 0 0 0 0 0, reflecting the fact that *smartest* contains one a, one e, one m, one r, two s's and two t's.) Use another loop to read the second word, except this time decrementing the corresponding array element as each letter is read. Both loops should ignore any characters that aren't letters, and both should treat upper-case letters in the same way as lower-case letters. After the second word has been read, use a third loop to check whether all the elements in the array are zero. If so, the words are anagrams. *Hint:* You may wish to use functions from <ctype.h>, such as isalpha and tolower.

17. Write a program that prints an  $n \times n$  magic square (a square arrangement of the numbers 1, 2, ...,  $n^2$  in which the sums of the rows, columns, and diagonals are all the same). The user will specify the value of n:

This program creates a magic square of a specified size.  
The size must be an odd number between 1 and 99.

Enter size of magic square: 5

17	24	1	8	15
23	5	7	14	16
4	6	13	20	22
10	12	19	21	3
11	18	25	2	9

Store the magic square in a two-dimensional array. Start by placing the number 1 in the middle of row 0. Place each of the remaining numbers 2, 3, ...,  $n^2$  by moving up one row and over one column. Any attempt to go outside the bounds of the array should “wrap around” to the opposite side of the array. For example, instead of storing the next number in row  $-1$ , we would store it in row  $n - 1$  (the last row). Instead of storing the next number in column  $n$ , we would store it in column 0. If a particular array element is already occupied, put the number directly below the previously stored number. If your compiler supports variable-length arrays, declare the array to have  $n$  rows and  $n$  columns. If not, declare the array to have 99 rows and 99 columns.

# 9 Functions

*If you have a procedure with ten parameters, you probably missed some.*

We saw in Chapter 2 that a function is simply a series of statements that have been grouped together and given a name. Although the term “function” comes from mathematics, C functions don’t always resemble math functions. In C, a function doesn’t necessarily have arguments, nor does it necessarily compute a value. (In some programming languages, a “function” returns a value, whereas a “procedure” doesn’t. C lacks this distinction.)

Functions are the building blocks of C programs. Each function is essentially a small program, with its own declarations and statements. Using functions, we can divide a program into small pieces that are easier for us—and others—to understand and modify. Functions can take some of the tedium out of programming by allowing us to avoid duplicating code that’s used more than once. Moreover, functions are reusable: we can take a function that was originally part of one program and use it in others.

Our programs so far have consisted of just the `main` function. In this chapter, we’ll see how to write functions other than `main`, and we’ll learn more about `main` itself. Section 9.1 shows how to define and call functions. Section 9.2 then discusses function declarations and how they differ from function definitions. Next, Section 9.3 examines how arguments are passed to functions. The remainder of the chapter covers the `return` statement (Section 9.4), the related issue of program termination (Section 9.5), and recursion (Section 9.6).

## 9.1 Defining and Calling Functions

Before we go over the formal rules for defining a function, let’s look at three simple programs that define functions.

## PROGRAM Computing Averages

Suppose we often need to compute the average of two double values. The C library doesn't have an "average" function, but we can easily define our own. Here's what it would look like:

```
double average(double a, double b)
{
    return (a + b) / 2;
}
```

The word `double` at the beginning is `average`'s ***return type***: the type of data that the function returns each time it's called. The identifiers `a` and `b` (the function's ***parameters***) represent the two numbers that will be supplied when `average` is called. Each parameter must have a type (just like every variable has a type); in this example, both `a` and `b` have type `double`. (It may look odd, but the word `double` must appear twice, once for `a` and once for `b`.) A function parameter is essentially a variable whose initial value will be supplied later, when the function is called.

Every function has an executable part, called the ***body***, which is enclosed in braces. The body of `average` consists of a single `return` statement. Executing this statement causes the function to "return" to the place from which it was called; the value of `(a + b) / 2` will be the value returned by the function.

To call a function, we write the function name, followed by a list of ***arguments***. For example, `average(x, y)` is a call of the `average` function. Arguments are used to supply information to a function; in this case, `average` needs to know which two numbers to average. The effect of the call `average(x, y)` is to copy the values of `x` and `y` into the parameters `a` and `b`, and then execute the body of `average`. An argument doesn't have to be a variable; any expression of a compatible type will do, allowing us to write `average(5.1, 8.9)` or `average(x/2, y/3)`.

We'll put the call of `average` in the place where we need to use the return value. For example, we could write

```
printf("Average: %g\n", average(x, y));
```

to compute the average of `x` and `y` and then print it. This statement has the following effect:

1. The `average` function is called with `x` and `y` as arguments.
2. `x` and `y` are copied into `a` and `b`.
3. `average` executes its `return` statement, returning the average of `a` and `b`.
4. `printf` prints the value that `average` returns. (The return value of `average` becomes one of `printf`'s arguments.)

Note that the return value of `average` isn't saved anywhere; the program prints it and then discards it. If we had needed the return value later in the program, we could have captured it in a variable:

### Q&A

```
avg = average(x, y);
```

This statement calls `average`, then saves its return value in the variable `avg`.

Now, let's use the `average` function in a complete program. The following program reads three numbers and computes their averages, one pair at a time:

```
Enter three numbers: 3.5 9.6 10.2
Average of 3.5 and 9.6: 6.55
Average of 9.6 and 10.2: 9.9
Average of 3.5 and 10.2: 6.85
```

Among other things, this program shows that a function can be called as often as we need.

```
average.c /* Computes pairwise averages of three numbers */

#include <stdio.h>

double average(double a, double b)
{
    return (a + b) / 2;
}

int main(void)
{
    double x, y, z;

    printf("Enter three numbers: ");
    scanf("%lf%lf%lf", &x, &y, &z);
    printf("Average of %g and %g: %g\n", x, y, average(x, y));
    printf("Average of %g and %g: %g\n", y, z, average(y, z));
    printf("Average of %g and %g: %g\n", x, z, average(x, z));

    return 0;
}
```

Notice that I've put the definition of `average` before `main`. We'll see in Section 9.2 that putting `average` after `main` causes problems.

## PROGRAM Printing a Countdown

Not every function returns a value. For example, a function whose job is to produce output may not need to return anything. To indicate that a function has no return value, we specify that its return type is `void`. (`void` is a type with no values.) Consider the following function, which prints the message `T minus n` and counting, where `n` is supplied when the function is called:

```
void print_count(int n)
{
    printf("T minus %d and counting\n", n);
```

`print_count` has one parameter, `n`, of type `int`. It returns nothing, so I've specified `void` as the return type and omitted the `return` statement. Since `print_count` doesn't return a value, we can't call it in the same way we call `average`. Instead, a call of `print_count` must appear in a statement by itself:

```
print_count(i);
```

Here's a program that calls `print_count` 10 times inside a loop:

```
countdown.c /* Prints a countdown */

#include <stdio.h>

void print_count(int n)
{
    printf("T minus %d and counting\n", n);
}

int main(void)
{
    int i;

    for (i = 10; i > 0; --i)
        print_count(i);

    return 0;
}
```

Initially, `i` has the value 10. When `print_count` is called for the first time, `i` is copied into `n`, so that `n` takes on the value 10 as well. As a result, the first call of `print_count` will print

T minus 10 and counting

`print_count` then returns to the point at which it was called, which happens to be the body of a `for` statement. The `for` statement resumes where it left off, decrementing `i` to 9 and testing whether it's greater than 0. It is, so `print_count` is called again, this time printing

T minus 9 and counting

Each time `print_count` is called, `i` is different, so `print_count` will print 10 different messages.

## PROGRAM Printing a Pun (Revisited)

Some functions have no parameters at all. Consider `print_pun`, which prints a bad pun each time it's called:

```
void print_pun(void)
{
    printf("To C, or not to C: that is the question.\n");
```

The word `void` in parentheses indicates that `print_pun` has no arguments. (This time, we're using `void` as a placeholder that means “nothing goes here.”)

To call a function with no arguments, we write the function's name, followed by parentheses:

```
print_pun();
```

The parentheses *must* be present, even though there are no arguments.

Here's a tiny program that tests the `print_pun` function:

```
pun2.c /* Prints a bad pun */

#include <stdio.h>

void print_pun(void)
{
    printf("To C, or not to C: that is the question.\n");
}

int main(void)
{
    print_pun();
    return 0;
}
```

The execution of this program begins with the first statement in `main`, which happens to be a call of `print_pun`. When `print_pun` begins to execute, it in turn calls `printf` to display a string. When `printf` returns, `print_pun` returns to `main`.

## Function Definitions

Now that we've seen several examples, let's look at the general form of a *function definition*:

### function definition

```
return-type function-name ( parameters )
{
    declarations
    statements
}
```

The return type of a function is the type of value that the function returns. The following rules govern the return type:

- Functions may not return arrays, but there are no other restrictions on the return type.
- Specifying that the return type is `void` indicates that the function doesn't return a value.

**C99**

- If the return type is omitted in C89, the function is presumed to return a value of type `int`. In C99, it's illegal to omit the return type of a function.

As a matter of style, some programmers put the return type *above* the function name:

```
double
average(double a, double b)
{
    return (a + b) / 2;
}
```

Putting the return type on a separate line is especially useful if the return type is lengthy, like `unsigned long int`.

**Q&A**

After the function name comes a list of parameters. Each parameter is preceded by a specification of its type; parameters are separated by commas. If the function has no parameters, the word `void` should appear between the parentheses. *Note:* A separate type must be specified for each parameter, even when several parameters have the same type:

```
double average(double a, b)      /*** WRONG ***/
{
    return (a + b) / 2;
}
```

The body of a function may include both declarations and statements. For example, the `average` function could be written

```
double average(double a, double b)
{
    double sum;          /* declaration */
    sum = a + b;         /* statement */
    return sum / 2;      /* statement */
}
```

**C99**

Variables declared in the body of a function belong exclusively to that function; they can't be examined or modified by other functions. In C89, variable declarations must come first, before all statements in the body of a function. In C99, variable declarations and statements can be mixed, as long as each variable is declared prior to the first statement that uses the variable. (Some pre-C99 compilers also allow mixing of declarations and statements.)

The body of a function whose return type is `void` (which I'll call a "void function") can be empty:

```
void print_pun(void)
{
}
```

Leaving the body empty may make sense during program development; we can leave room for the function without taking the time to complete it, then come back later and write the body.

## Function Calls

A function call consists of a function name followed by a list of arguments, enclosed in parentheses:

```
average(x, y)
print_count(i)
print_pun()
```



If the parentheses are missing, the function won't get called:

```
print_pun;    /* *** WRONG *** /
```

**Q&A**

The result is a legal (albeit meaningless) expression statement that looks correct, but has no effect. Some compilers issue a warning such as “*statement with no effect.*”

A call of a `void` function is always followed by a semicolon to turn it into a statement:

```
print_count(i);
print_pun();
```

A call of a non-`void` function, on the other hand, produces a value that can be stored in a variable, tested, printed, or used in some other way:

```
avg = average(x, y);
if (average(x, y) > 0)
    printf("Average is positive\n");
printf("The average is %g\n", average(x, y));
```

The value returned by a non-`void` function can always be discarded if it's not needed:

```
average(x, y); /* discards return value */
```

expression statements ▶ 4.5

This call of `average` is an example of an expression statement: a statement that evaluates an expression but then discards the result.

Ignoring the return value of `average` is an odd thing to do, but for some functions it makes sense. The `printf` function, for example, returns the number of characters that it prints. After the following call, `num_chars` will have the value 9:

```
num_chars = printf("Hi, Mom!\n");
```

Since we're probably not interested in the number of characters printed, we'll normally discard `printf`'s return value:

```
printf("Hi, Mom!\n"); /* discards return value */
```

To make it clear that we're deliberately discarding the return value of a function, C allows us to put `(void)` before the call:

```
(void) printf("Hi, Mom!\n");
```

casting ▶ 7.4 What we're doing is casting (converting) the return value of `printf` to type `void`. (In C, "casting to `void`" is a polite way of saying "throwing away.") Using `(void)` makes it clear to others that you deliberately discarded the return value, not just forgot that there was one. Unfortunately, there are a great many functions in the C library whose values are routinely ignored; using `(void)` when calling them all can get tiresome, so I haven't done so in this book.

## PROGRAM Testing Whether a Number Is Prime

To see how functions can make programs easier to understand, let's write a program that tests whether a number is prime. The program will prompt the user to enter a number, then respond with a message indicating whether or not the number is prime:

```
Enter a number: 34
```

```
Not prime
```

Instead of putting the prime-testing details in `main`, we'll define a separate function that returns `true` if its parameter is a prime number and `false` if it isn't. When given a number `n`, the `is_prime` function will divide `n` by each of the numbers between 2 and the square root of `n`; if the remainder is ever 0, we know that `n` isn't prime.

```
prime.c /* Tests whether a number is prime */

#include <stdbool.h> /* C99 only */
#include <stdio.h>

bool is_prime(int n)
{
    int divisor;

    if (n <= 1)
        return false;
    for (divisor = 2; divisor * divisor <= n; divisor++)
        if (n % divisor == 0)
            return false;
    return true;
}

int main(void)
{
    int n;

    printf("Enter a number: ");
    scanf("%d", &n);
    if (is_prime(n))
        printf("Prime\n");
    else
        printf("Not prime\n");
```

```

    return 0;
}

```

Notice that `main` contains a variable named `n` even though `is_prime`'s parameter is also named `n`. In general, a function may declare a variable with the same name as a variable in another function. The two variables represent different locations in memory, so assigning a new value to one variable doesn't change the other. (This property extends to parameters as well.) Section 10.1 discusses this point in more detail.

As `is_prime` demonstrates, a function may have more than one `return` statement. However, we can execute just one of these statements during a given call of the function, because reaching a `return` statement causes the function to return to where it was called. We'll learn more about the `return` statement in Section 9.4.

## 9.2 Function Declarations

In the programs in Section 9.1, the definition of each function was always placed *above* the point at which it was called. In fact, C doesn't require that the definition of a function precede its calls. Suppose that we rearrange the `average.c` program by putting the definition of `average` *after* the definition of `main`:

```

#include <stdio.h>

int main(void)
{
    double x, y, z;

    printf("Enter three numbers: ");
    scanf("%lf%lf%lf", &x, &y, &z);
    printf("Average of %g and %g: %g\n", x, y, average(x, y));
    printf("Average of %g and %g: %g\n", y, z, average(y, z));
    printf("Average of %g and %g: %g\n", x, z, average(x, z));

    return 0;
}

double average(double a, double b)
{
    return (a + b) / 2;
}

```

When the compiler encounters the first call of `average` in `main`, it has no information about `average`: it doesn't know how many parameters `average` has, what the types of these parameters are, or what kind of value `average` returns. Instead of producing an error message, though, the compiler assumes that `average` returns an `int` value (recall from Section 9.1 that the return type of a

default argument promotions ► 9.3

function is `int` by default). We say that the compiler has created an *implicit declaration* of the function. The compiler is unable to check that we're passing average the right number of arguments and that the arguments have the proper type. Instead, it performs the default argument promotions and hopes for the best. When it encounters the definition of `average` later in the program, the compiler notices that the function's return type is actually `double`, not `int`, and so we get an error message.

One way to avoid the problem of call-before-definition is to arrange the program so that the definition of each function precedes all its calls. Unfortunately, such an arrangement doesn't always exist, and even when it does, it may make the program harder to understand by putting its function definitions in an unnatural order.

Fortunately, C offers a better solution: declare each function before calling it. A *function declaration* provides the compiler with a brief glimpse at a function whose full definition will appear later. A function declaration resembles the first line of a function definition with a semicolon added at the end:

**function declaration**`return-type function-name ( parameters ) ;`

Needless to say, the declaration of a function must be consistent with the function's definition.

**Q&A**

Here's how our program would look with a declaration of `average` added:

```
#include <stdio.h>

double average(double a, double b); /* DECLARATION */

int main(void)
{
    double x, y, z;

    printf("Enter three numbers: ");
    scanf("%lf%lf%lf", &x, &y, &z);
    printf("Average of %g and %g: %g\n", x, y, average(x, y));
    printf("Average of %g and %g: %g\n", y, z, average(y, z));
    printf("Average of %g and %g: %g\n", x, z, average(x, z));

    return 0;
}

double average(double a, double b) /* DEFINITION */
{
    return (a + b) / 2;
}
```

Function declarations of the kind we've been discussing are known as *function prototypes* to distinguish them from an older style of function declaration in which the parentheses are left empty. A prototype provides a complete description

**Q&A**

of how to call a function: how many arguments to supply, what their types should be, and what type of result will be returned.

Incidentally, a function prototype doesn't have to specify the *names* of the function's parameters, as long as their *types* are present:

```
double average(double, double);
```

It's usually best not to omit parameter names, since they help document the purpose of each parameter and remind the programmer of the order in which arguments must appear when the function is called. However, there are legitimate reasons for omitting parameter names, and some programmers prefer to do so.

**Q&A**
**C99**

C99 has adopted the rule that either a declaration or a definition of a function must be present prior to any call of the function. Calling a function for which the compiler has not yet seen a declaration or definition is an error.

## 9.3 Arguments

Let's review the difference between a parameter and an argument. *Parameters* appear in function *definitions*; they're dummy names that represent values to be supplied when the function is called. *Arguments* are expressions that appear in function *calls*. When the distinction between *argument* and *parameter* isn't important, I'll sometimes use *argument* to mean either.

In C, arguments are *passed by value*: when a function is called, each argument is evaluated and its value assigned to the corresponding parameter. Since the parameter contains a copy of the argument's value, any changes made to the parameter during the execution of the function don't affect the argument. In effect, each parameter behaves like a variable that's been initialized to the value of the matching argument.

The fact that arguments are passed by value has both advantages and disadvantages. Since a parameter can be modified without affecting the corresponding argument, we can use parameters as variables within the function, thereby reducing the number of genuine variables needed. Consider the following function, which raises a number *x* to a power *n*:

```
int power(int x, int n)
{
    int i, result = 1;

    for (i = 1; i <= n; i++)
        result = result * x;

    return result;
}
```

Since *n* is a *copy* of the original exponent, we can modify it inside the function, thus removing the need for *i*:

```
int power(int x, int n)
{
    int result = 1;

    while (n-- > 0)
        result = result * x;

    return result;
}
```

Unfortunately, C's requirement that arguments be passed by value makes it difficult to write certain kinds of functions. For example, suppose that we need a function that will decompose a `double` value into an integer part and a fractional part. Since a function can't *return* two numbers, we might try passing a pair of variables to the function and having it modify them:

```
void decompose(double x, long int_part, double frac_part)
{
    int_part = (long) x; /* drops the fractional part of x */
    frac_part = x - int_part;
}
```

Suppose that we call the function in the following way:

```
decompose(3.14159, i, d);
```

At the beginning of the call, 3.14159 is copied into `x`, `i`'s value is copied into `int_part`, and `d`'s value is copied into `frac_part`. The statements inside `decompose` then assign 3 to `int_part` and .14159 to `frac_part`, and the function returns. Unfortunately, `i` and `d` weren't affected by the assignments to `int_part` and `frac_part`, so they have the same values after the call as they did before the call. With a little extra effort, `decompose` can be made to work, as we'll see in Section 11.4. However, we'll need to cover more of C's features first.

## Argument Conversions

C allows function calls in which the types of the arguments don't match the types of the parameters. The rules governing how the arguments are converted depend on whether or not the compiler has seen a prototype for the function (or the function's full definition) prior to the call:

- ***The compiler has encountered a prototype prior to the call.*** The value of each argument is implicitly converted to the type of the corresponding parameter as if by assignment. For example, if an `int` argument is passed to a function that was expecting a `double`, the argument is converted to `double` automatically.
- ***The compiler has not encountered a prototype prior to the call.*** The compiler performs the **default argument promotions**: (1) `float` arguments are converted to `double`. (2) The integral promotions are performed, causing `char`

**C99**

and short arguments to be converted to int. (In C99, the integer promotions are performed.)



Relying on the default argument promotions is dangerous. Consider the following program:

```
#include <stdio.h>

int main(void)
{
    double x = 3.0;
    printf("Square: %d\n", square(x));

    return 0;
}

int square(int n)
{
    return n * n;
}
```

At the time `square` is called, the compiler hasn't seen a prototype yet, so it doesn't know that `square` expects an argument of type `int`. Instead, the compiler performs the default argument promotions on `x`, with no effect. Since it's expecting an argument of type `int` but has been given a `double` value instead, the effect of calling `square` is undefined. The problem can be fixed by casting `square`'s argument to the proper type:

```
printf("Square: %d\n", square((int) x));
```

**C99**

Of course, a much better solution is to provide a prototype for `square` before calling it. In C99, calling `square` without first providing a declaration or definition of the function is an error.

## Array Arguments

Arrays are often used as arguments. When a function parameter is a one-dimensional array, the length of the array can be (and is normally) left unspecified:

```
int f(int a[]) /* no length specified */
{
    ...
}
```

The argument can be any one-dimensional array whose elements are of the proper type. There's just one problem: how will `f` know how long the array is? Unfortunately, C doesn't provide any easy way for a function to determine the length of an array passed to it. Instead, we'll have to supply the length—if the function needs it—as an additional argument.

**Q&A**



Although we can use the `sizeof` operator to help determine the length of an array *variable*, it doesn't give the correct answer for an array *parameter*:

```
int f(int a[])
{
    int len = sizeof(a) / sizeof(a[0]);
    /*** WRONG: not the number of elements in a ***/
    ...
}
```

Section 12.3 explains why.

The following function illustrates the use of one-dimensional array arguments. When given an array `a` of `int` values, `sum_array` returns the sum of the elements in `a`. Since `sum_array` needs to know the length of `a`, we must supply it as a second argument.

```
int sum_array(int a[], int n)
{
    int i, sum = 0;

    for (i = 0; i < n; i++)
        sum += a[i];

    return sum;
}
```

The prototype for `sum_array` has the following appearance:

```
int sum_array(int a[], int n);
```

As usual, we can omit the parameter names if we wish:

```
int sum_array(int [], int);
```

When `sum_array` is called, the first argument will be the name of an array, and the second will be its length. For example:

```
#define LEN 100

int main(void)
{
    int b[LEN], total;
    ...
    total = sum_array(b, LEN);
    ...
}
```

Notice that we don't put brackets after an array name when passing it to a function:

```
total = sum_array(b[], LEN);    /*** WRONG ***/

```

An important point about array arguments: A function has no way to check that we've passed it the correct array length. We can exploit this fact by telling the function that the array is smaller than it really is. Suppose that we've only stored 50 numbers in the `b` array, even though it can hold 100. We can sum just the first 50 elements by writing

```
total = sum_array(b, 50); /* sums first 50 elements */
```

`sum_array` will ignore the other 50 elements. (Indeed, it won't know that they even exist!)



Be careful not to tell a function that an array argument is *larger* than it really is:

```
total = sum_array(b, 150); /**** WRONG ***/
```

In this example, `sum_array` will go past the end of the array, causing undefined behavior.

Another important thing to know is that a function is allowed to change the elements of an array parameter, and the change is reflected in the corresponding argument. For example, the following function modifies an array by storing zero into each of its elements:

```
void store_zeros(int a[], int n)
{
    int i;

    for (i = 0; i < n; i++)
        a[i] = 0;
}
```

The call

```
store_zeros(b, 100);
```

will store zero into the first 100 elements of the array `b`. This ability to modify the elements of an array argument may seem to contradict the fact that C passes arguments by value. In fact, there's no contradiction, but I won't be able to explain why until Section 12.3.

If a parameter is a multidimensional array, only the length of the first dimension may be omitted when the parameter is declared. For example, if we revise the `sum_array` function so that `a` is a two-dimensional array, we must specify the number of columns in `a`, although we don't have to indicate the number of rows:

```
#define LEN 10

int sum_two_dimensional_array(int a[] [LEN], int n)
{
    int i, j, sum = 0;
```

**Q&A**

arrays of pointers ▶ 13.7

variable-length arrays ▶ 8.3

```

    for (i = 0; i < n; i++)
        for (j = 0; j < LEN; j++)
            sum += a[i][j];

    return sum;
}

```

Not being able to pass multidimensional arrays with an arbitrary number of columns can be a nuisance. Fortunately, we can often work around this difficulty by using arrays of pointers. C99's variable-length array parameters provide an even better solution to the problem.

**C99**

## Variable-Length Array Parameters

C99 adds several new twists to array arguments. The first has to do with variable-length arrays (VLAs), a feature of C99 that allows the length of an array to be specified using a non-constant expression. Variable-length arrays can also be parameters, as it turns out.

Consider the `sum_array` function discussed earlier in this section. Here's the definition of `sum_array`, with the body omitted:

```

int sum_array(int a[], int n)
{
    ...
}

```

As it stands now, there's no direct link between `n` and the length of the array `a`. Although the function body treats `n` as `a`'s length, the actual length of the array could in fact be larger than `n` (or smaller, in which case the function won't work correctly).

Using a variable-length array parameter, we can explicitly state that `a`'s length is `n`:

```

int sum_array(int n, int a[n])
{
    ...
}

```

The value of the first parameter (`n`) specifies the length of the second parameter (`a`). Note that the order of the parameters has been switched; order is important when variable-length array parameters are used.



The following version of `sum_array` is illegal:

```

int sum_array(int a[n], int n)    /*** WRONG ***/
{
    ...
}

```

The compiler will issue an error message at `int a [n]`, because it hasn't yet seen `n`.

There are several ways to write the prototype for our new version of `sum_array`. One possibility is to make it look exactly like the function definition:

```
int sum_array(int n, int a[n]); /* Version 1 */
```

Another possibility is to replace the array length by an asterisk (\*):

```
int sum_array(int n, int a[*]); /* Version 2a */
```

The reason for using the \* notation is that parameter names are optional in function declarations. If the name of the first parameter is omitted, it wouldn't be possible to specify that the length of the array is `n`, but the \* provides a clue that the length of the array is related to parameters that come earlier in the list:

```
int sum_array(int, int [*]); /* Version 2b */
```

It's also legal to leave the brackets empty, as we normally do when declaring an array parameter:

```
int sum_array(int n, int a[]); /* Version 3a */
int sum_array(int, int []); /* Version 3b */
```

Leaving the brackets empty isn't a good choice, because it doesn't expose the relationship between `n` and `a`.

In general, the length of a variable-length array parameter can be any expression. For example, suppose that we were to write a function that concatenates two arrays `a` and `b` by copying the elements of `a`, followed by the elements of `b`, into a third array named `c`:

```
int concatenate(int m, int n, int a[m], int b[n], int c[m+n])
{
    ...
}
```

The length of `c` is the sum of the lengths of `a` and `b`. The expression used to specify the length of `c` involves two other parameters, but in general it could refer to variables outside the function or even call other functions.

Variable-length array parameters with a single dimension—as in all our examples so far—have limited usefulness. They make a function declaration or definition more descriptive by stating the desired length of an array argument. However, no additional error-checking is performed; it's still possible for an array argument to be too long or too short.

It turns out that variable-length array parameters are most useful for multidimensional arrays. Earlier in this section, we tried to write a function that sums the elements in a two-dimensional array. Our original function was limited to arrays with a fixed number of columns. If we use a variable-length array parameter, we can generalize the function to any number of columns:

```

int sum_two_dimensional_array(int n, int m, int a[n][m])
{
    int i, j, sum = 0;

    for (i = 0; i < n; i++)
        for (j = 0; j < m; j++)
            sum += a[i][j];

    return sum;
}

```

Prototypes for this function include the following:

```

int sum_two_dimensional_array(int n, int m, int a[n][m]);
int sum_two_dimensional_array(int n, int m, int a[*][*]);
int sum_two_dimensional_array(int n, int m, int a[] [m]);
int sum_two_dimensional_array(int n, int m, int a[] [*]);

```

### C99 Using static in Array Parameter Declarations

C99 allows the use of the keyword `static` in the declaration of array parameters. (The keyword itself existed before C99. Section 18.2 discusses its traditional uses.)

In the following example, putting `static` in front of the number 3 indicates that the length of `a` is guaranteed to be at least 3:

```

int sum_array(int a[static 3], int n)
{
    ...
}

```

Using `static` in this way has no effect on the behavior of the program. The presence of `static` is merely a “hint” that may allow a C compiler to generate faster instructions for accessing the array. (If the compiler knows that an array will always have a certain minimum length, it can arrange to “prefetch” these elements from memory when the function is called, before the elements are actually needed by statements within the function.)

One last note about `static`: If an array parameter has more than one dimension, `static` can be used only in the first dimension (for example, when specifying the number of rows in a two-dimensional array).

### C99 Compound Literals

Let’s return to the original `sum_array` function one last time. When `sum_array` is called, the first argument is usually the name of an array (the one whose elements are to be summed). For example, we might call `sum_array` in the following way:

```

int b[] = {3, 0, 3, 4, 1};
total = sum_array(b, 5);

```

The only problem with this arrangement is that `b` must be declared as a variable and then initialized prior to the call. If `b` isn't needed for any other purpose, it can be mildly annoying to create it solely for the purpose of calling `sum_array`.

In C99, we can avoid this annoyance by using a ***compound literal***: an unnamed array that's created "on the fly" by simply specifying which elements it contains. The following call of `sum_array` has a compound literal (shown in **bold**) as its first argument:

```
total = sum_array((int []){3, 0, 3, 4, 1}, 5);
```

In this example, the compound literal creates an array containing the five integers 3, 0, 3, 4, and 1. We didn't specify the length of the array, so it's determined by the number of elements in the literal. We also have the option of specifying a length explicitly: `(int [4]) {1, 9, 2, 1}` is equivalent to `(int []) {1, 9, 2, 1}`.

In general, a compound literal consists of a type name within parentheses, followed by a set of values enclosed by braces. A compound literal resembles a cast applied to an initializer. In fact, compound literals and initializers obey the same rules. A compound literal may contain designators, just like a designated initializer, and it may fail to provide full initialization (in which case any uninitialized elements default to zero). For example, the literal `(int [10]) {8, 6}` has 10 elements; the first two have the values 8 and 6, and the remaining elements have the value 0.

Compound literals created inside a function may contain arbitrary expressions, not just constants. For example, we could write

```
total = sum_array((int []){2 * i, i + j, j * k}, 3);
```

where `i`, `j`, and `k` are variables. This aspect of compound literals greatly enhances their usefulness.

A compound literal is an lvalue, so the values of its elements can be changed. If desired, a compound literal can be made "read-only" by adding the word `const` to its type, as in `(const int []) {5, 4}`.

## 9.4 The return Statement

A non-void function must use the `return` statement to specify what value it will return. The `return` statement has the form

**return statement**

```
return expression ;
```

The expression is often just a constant or variable:

```
return 0;
return status;
```

conditional operator ➤5.2

More complex expressions are possible. For example, it's not unusual to see the conditional operator used in a return expression:

```
return n >= 0 ? n : 0;
```

When this statement is executed, the expression `n >= 0 ? n : 0` is evaluated first. The statement returns the value of `n` if it's not negative; otherwise, it returns 0.

If the type of the expression in a `return` statement doesn't match the function's return type, the expression will be implicitly converted to the return type. For example, if a function is declared to return an `int`, but the `return` statement contains a `double` expression, the value of the expression is converted to `int`.

`return` statements may appear in functions whose return type is `void`, provided that no expression is given:

```
return; /* return in a void function */
```

Putting an expression in such a `return` statement will get you a compile-time error. In the following example, the `return` statement causes the function to return immediately when given a negative argument:

```
void print_int(int i)
{
    if (i < 0)
        return;
    printf("%d", i);
}
```

If `i` is less than 0, `print_int` will return without calling `printf`.

A `return` statement may appear at the end of a `void` function:

```
void print_pun(void)
{
    printf("To C, or not to C: that is the question.\n");
    return; /* OK, but not needed */
}
```

Using `return` is unnecessary, though, since the function will return automatically after its last statement has been executed.

If a non-`void` function reaches the end of its body—that is, it fails to execute a `return` statement—the behavior of the program is undefined if it attempts to use the value returned by the function. Some compilers will issue a warning such as “*control reaches end of non-void function*” if they detect the possibility of a non-`void` function “falling off” the end of its body.

## 9.5 Program Termination

Since `main` is a function, it must have a return type. Normally, the return type of `main` is `int`, which is why the programs we've seen so far have defined `main` in the following way:

```
int main(void)
{
    ...
}
```

Older C programs often omit `main`'s return type, taking advantage of the fact that it traditionally defaults to `int`:

```
main()
{
    ...
}
```

**C99**

Omitting the return type of a function isn't legal in C99, so it's best to avoid this practice. Omitting the word `void` in `main`'s parameter list remains legal, but—as a matter of style—it's best to be explicit about the fact that `main` has no parameters. (We'll see later that `main` sometimes *does* have two parameters, usually named `argc` and `argv`.)

**Q&A**

The value returned by `main` is a status code that—in some operating systems—can be tested when the program terminates. `main` should return 0 if the program terminates normally; to indicate abnormal termination, `main` should return a value other than 0. (Actually, there's no rule to prevent us from using the return value for other purposes.) It's good practice to make sure that every C program returns a status code, even if there are no plans to use it, since someone running the program later may decide to test it.

## The `exit` Function

`<stdlib.h>` header ➤ 26.2 Executing a `return` statement in `main` is one way to terminate a program. Another is calling the `exit` function, which belongs to `<stdlib.h>`. The argument passed to `exit` has the same meaning as `main`'s return value: both indicate the program's status at termination. To indicate normal termination, we'd pass 0:

```
exit(0);           /* normal termination */
```

Since 0 is a bit cryptic, C allows us to pass `EXIT_SUCCESS` instead (the effect is the same):

```
exit(EXIT_SUCCESS); /* normal termination */
```

Passing `EXIT_FAILURE` indicates abnormal termination:

```
exit(EXIT_FAILURE); /* abnormal termination */
```

`EXIT_SUCCESS` and `EXIT_FAILURE` are macros defined in `<stdlib.h>`. The values of `EXIT_SUCCESS` and `EXIT_FAILURE` are implementation-defined; typical values are 0 and 1, respectively.

As methods of terminating a program, `return` and `exit` are closely related. In fact, the statement

```
return expression;
```

in main is equivalent to

```
exit (expression) ;
```

The difference between return and exit is that exit causes program termination regardless of which function calls it. The return statement causes program termination only when it appears in the main function. Some programmers use exit exclusively to make it easier to locate all exit points in a program.

## 9.6 Recursion

A function is **recursive** if it calls itself. For example, the following function computes  $n!$  recursively, using the formula  $n! = n \times (n - 1)!$ :

```
int fact(int n)
{
    if (n <= 1)
        return 1;
    else
        return n * fact(n - 1);
}
```

Some programming languages rely heavily on recursion, while others don't even allow it. C falls somewhere in the middle: it allows recursion, but most C programmers don't use it that often.

To see how recursion works, let's trace the execution of the statement

```
i = fact(3);
```

Here's what happens:

```
fact(3) finds that 3 is not less than or equal to 1, so it calls
fact(2), which finds that 2 is not less than or equal to 1, so it calls
fact(1), which finds that 1 is less than or equal to 1, so it returns 1, causing
fact(2) to return  $2 \times 1 = 2$ , causing
fact(3) to return  $3 \times 2 = 6$ .
```

Notice how the unfinished calls of fact "pile up" until fact is finally passed 1. At that point, the old calls of fact begin to "unwind" one by one, until the original call—fact(3)—finally returns with the answer, 6.

Here's another example of recursion: a function that computes  $x^n$ , using the formula  $x^n = x \times x^{n-1}$ .

```
int power(int x, int n)
{
    if (n == 0)
        return 1;
    else
        return x * power(x, n - 1);
}
```

The call `power(5, 3)` would be executed as follows:

```
power(5, 3) finds that 3 is not equal to 0, so it calls
    power(5, 2), which finds that 2 is not equal to 0, so it calls
        power(5, 1), which finds that 1 is not equal to 0, so it calls
            power(5, 0), which finds that 0 is equal to 0, so it returns 1, causing
                power(5, 1) to return  $5 \times 1 = 5$ , causing
                    power(5, 2) to return  $5 \times 5 = 25$ , causing
                        power(5, 3) to return  $5 \times 25 = 125$ .
```

Incidentally, we can condense the `power` function a bit by putting a conditional expression in the `return` statement:

```
int power(int x, int n)
{
    return n == 0 ? 1 : x * power(x, n - 1);
}
```

Both `fact` and `power` are careful to test a “termination condition” as soon as they’re called. When `fact` is called, it immediately checks whether its parameter is less than or equal to 1. When `power` is called, it first checks whether its second parameter is equal to 0. All recursive functions need some kind of termination condition in order to prevent infinite recursion.

## The Quicksort Algorithm

At this point, you may wonder why we’re bothering with recursion; after all, neither `fact` nor `power` really needs it. Well, you’ve got a point. Neither function makes much of a case for recursion, because each calls itself just once. Recursion is much more helpful for sophisticated algorithms that require a function to call itself two or more times.

In practice, recursion often arises naturally as a result of an algorithm design technique known as *divide-and-conquer*, in which a large problem is divided into smaller pieces that are then tackled by the same algorithm. A classic example of the divide-and-conquer strategy can be found in the popular sorting algorithm known as *Quicksort*. The Quicksort algorithm goes as follows (for simplicity, we’ll assume that the array being sorted is indexed from 1 to  $n$ ):

1. Choose an array element  $e$  (the “partitioning element”), then rearrange the array so that elements  $1, \dots, i - 1$  are less than or equal to  $e$ , element  $i$  contains  $e$ , and elements  $i + 1, \dots, n$  are greater than or equal to  $e$ .
2. Sort elements  $1, \dots, i - 1$  by using Quicksort recursively.
3. Sort elements  $i + 1, \dots, n$  by using Quicksort recursively.

After step 1, the element  $e$  is in its proper location. Since the elements to the left of  $e$  are all less than or equal to it, they’ll be in their proper places once they’ve been sorted in step 2; similar reasoning applies to the elements to the right of  $e$ .

Step 1 of the Quicksort algorithm is obviously critical. There are various methods to partition an array, some much better than others. We’ll use a technique

that's easy to understand but not particularly efficient. I'll first describe the partitioning algorithm informally; later, we'll translate it into C code.

The algorithm relies on two "markers" named *low* and *high*, which keep track of positions within the array. Initially, *low* points to the first element of the array and *high* points to the last element. We start by copying the first element (the partitioning element) into a temporary location elsewhere, leaving a "hole" in the array. Next, we move *high* across the array from right to left until it points to an element that's smaller than the partitioning element. We then copy the element into the hole that *low* points to, which creates a new hole (pointed to by *high*). We now move *low* from left to right, looking for an element that's larger than the partitioning element. When we find one, we copy it into the hole that *high* points to. The process repeats, with *low* and *high* taking turns, until they meet somewhere in the middle of the array. At that time, both will point to a hole; all we need do is copy the partitioning element into the hole. The following diagrams illustrate how Quicksort would sort an array of integers:

Let's start with an array containing seven elements. *low* points to the first element; *high* points to the last one.

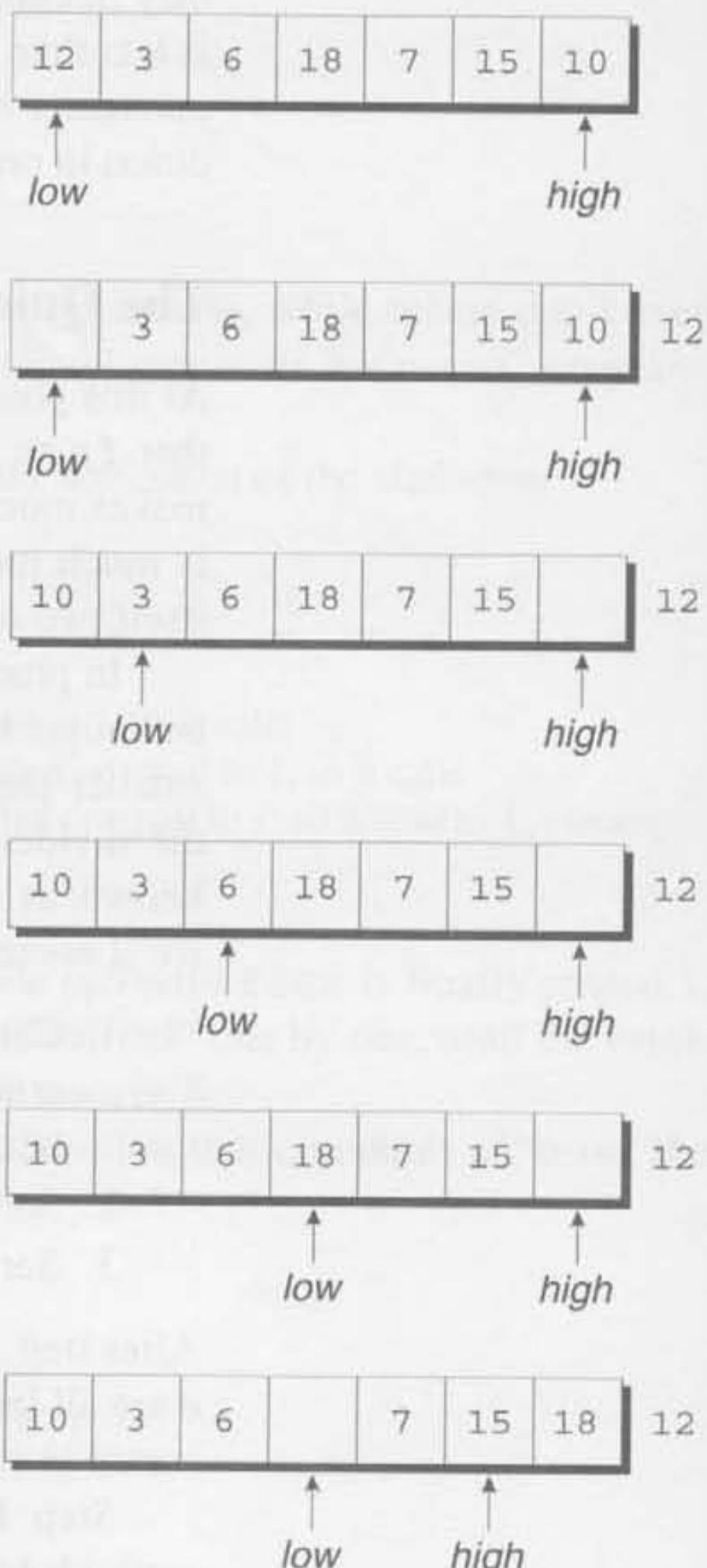
The first element, 12, is the partitioning element. Copying it somewhere else leaves a hole at the beginning of the array.

We now compare the element pointed to by *high* with 12. Since 10 is smaller than 12, it's on the wrong side of the array, so we move it to the hole and shift *low* to the right.

*low* points to the number 3, which is less than 12 and therefore doesn't need to be moved. We shift *low* to the right instead.

Since 6 is also less than 12, we shift *low* again.

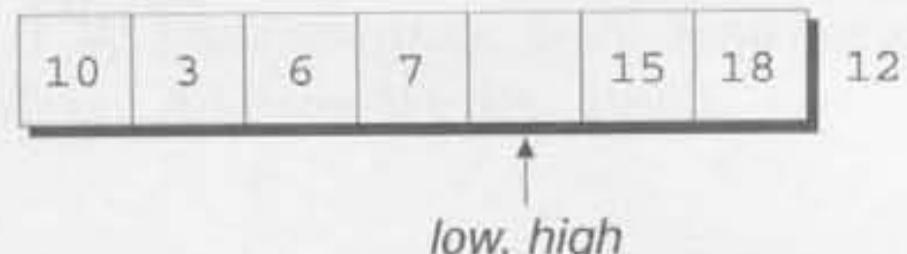
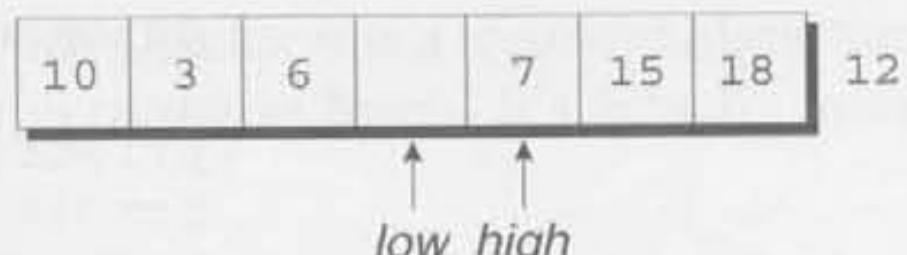
*low* now points to 18, which is larger than 12 and therefore out of position. After moving 18 to the hole, we shift *high* to the left.



*high* points to 15, which is greater than 12 and thus doesn't need to be moved. We shift *high* to the left and continue.

*high* points to 7, which is out of position. After moving 7 to the hole, we shift *low* to the right.

*low* and *high* are now equal, so we move the partitioning element to the hole.



At this point, we've accomplished our objective: all elements to the left of the partitioning element are less than or equal to 12, and all elements to the right are greater than or equal to 12. Now that the array has been partitioned, we can use Quicksort recursively to sort the first four elements of the array (10, 3, 6, and 7) and the last two (15 and 18).

## PROGRAM Quicksort

Let's develop a recursive function named `quicksort` that uses the Quicksort algorithm to sort an array of integers. To test the function, we'll have `main` read 10 numbers into an array, call `quicksort` to sort the array, then print the elements in the array:

```
Enter 10 numbers to be sorted: 9 16 47 82 4 66 12 3 25 51
In sorted order: 3 4 9 12 16 25 47 51 66 82
```

Since the code for partitioning the array is a bit lengthy, I'll put it in a separate function named `split`.

```
qsort.c /* Sorts an array of integers using Quicksort algorithm */

#include <stdio.h>

#define N 10

void quicksort(int a[], int low, int high);
int split(int a[], int low, int high);

int main(void)
{
    int a[N], i;

    printf("Enter %d numbers to be sorted: ", N);
    for (i = 0; i < N; i++)
        scanf("%d", &a[i]);
```

```

        quicksort(a, 0, N - 1);

        printf("In sorted order: ");
        for (i = 0; i < N; i++)
            printf("%d ", a[i]);
        printf("\n");

        return 0;
    }

void quicksort(int a[], int low, int high)
{
    int middle;

    if (low >= high) return;
    middle = split(a, low, high);
    quicksort(a, low, middle - 1);
    quicksort(a, middle + 1, high);
}

int split(int a[], int low, int high)
{
    int part_element = a[low];

    for (;;) {
        while (low < high && part_element <= a[high])
            high--;
        if (low >= high) break;
        a[low++] = a[high];

        while (low < high && a[low] <= part_element)
            low++;
        if (low >= high) break;
        a[high--] = a[low];
    }

    a[high] = part_element;
    return high;
}

```

Although this version of Quicksort works, it's not the best. There are numerous ways to improve the program's performance, including:

- ***Improving the partitioning algorithm.*** Our method isn't the most efficient. Instead of choosing the first element in the array as the partitioning element, it's better to take the median of the first element, the middle element, and the last element. The partitioning process itself can also be sped up. In particular, it's possible to avoid the `low < high` tests in the two `while` loops.
- ***Using a different method to sort small arrays.*** Instead of using Quicksort recursively all the way down to arrays with one element, it's better to use a simpler method for small arrays (those with fewer than, say, 25 elements).

- **Making Quicksort nonrecursive.** Although Quicksort is a recursive algorithm by nature—and is easiest to understand in recursive form—it's actually more efficient if the recursion is removed.

For details about improving Quicksort, consult a book on algorithm design, such as Robert Sedgewick's *Algorithms in C, Parts 1–4: Fundamentals, Data Structures, Sorting, Searching*, Third Edition (Boston, Mass.: Addison-Wesley, 1998).

## Q & A

- Q:** Some C books appear to use terms other than *parameter* and *argument*. Is there any standard terminology? [p. 184]

- A:** As with many other aspects of C, there's no general agreement on terminology, although the C89 and C99 standards use *parameter* and *argument*. The following table should help you translate:

<i>This book:</i>	<i>Other books:</i>
parameter	formal argument, formal parameter
argument	actual argument, actual parameter

Keep in mind that—when no confusion would result—I sometimes deliberately blur the distinction between the two terms, using *argument* to mean either.

- Q:** I've seen programs in which parameter types are specified in separate declarations after the parameter list, as in the following example:

```
double average(a, b)
double a, b;
{
    return (a + b) / 2;
}
```

- Is this practice legal? [p. 188]**

- A:** This method of defining functions comes from K&R C, so you may encounter it in older books and programs. C89 and C99 support this style so that older programs will still compile. I'd avoid using it in new programs, however, for a couple of reasons.

First, functions that are defined in the older way aren't subject to the same degree of error-checking. When a function is defined in the older way—and no prototype is present—the compiler won't check that the function is called with the right number of arguments, nor will it check that the arguments have the proper types. Instead, it will perform the default argument promotions.

Second, the C standard says that the older style is "obsolescent," meaning that its use is discouraged and that it may be dropped from C eventually.

**Q:** Some programming languages allow procedures and functions to be nested within each other. Does C allow function definitions to be nested?

**A:** No. C does not permit the definition of one function to appear in the body of another. Among other things, this restriction simplifies the compiler.

**\*Q:** Why does the compiler allow the use of function names that aren't followed by parentheses? [p. 189]

**A:** We'll see in a later chapter that the compiler treats a function name not followed by parentheses as a *pointer* to the function. Pointers to functions have legitimate uses, so the compiler can't automatically assume that a function name without parentheses is an error. The statement

```
print_pun;
```

is legal because the compiler treats `print_pun` as a pointer and therefore an expression, making this a valid (although pointless) expression statement.

**\*Q:** In the function call `f(a, b)`, how does the compiler know whether the comma is punctuation or whether it's an operator?

**A:** It turns out that the arguments in a function call can't be arbitrary expressions. Instead, they must be "assignment expressions," which can't contain commas used as operators unless they're enclosed in parentheses. In other words, in the call `f(a, b)` the comma is punctuation; in the call `f((a, b))` it's an operator.

**Q:** Do the names of parameters in a function prototype have to match the names given later in the function's definition? [p. 192]

**A:** No. Some programmers take advantage of this fact by giving long names to parameters in the prototype, then using shorter names in the actual definition. Or a French-speaking programmer might use English names in prototypes, then switch to more familiar French names in function definitions.

**Q:** I still don't understand why we bother with function prototypes. If we just put definitions of all the functions before `main`, we're covered, right?

**A:** Wrong. First, you're assuming that only `main` calls the other functions, which is unrealistic. In practice, some of the functions will call each other. If we put all function definitions above `main`, we'll have to watch their order carefully. Calling a function that hasn't been defined yet can lead to big problems.

But that's not all. Suppose that two functions call each other (which isn't as far-fetched as it may sound). No matter which function we define first, it will end up calling a function that hasn't been defined yet.

But there's still more! Once programs reach a certain size, it won't be feasible to put all the functions in one file anymore. When we reach that point, we'll need prototypes to tell the compiler about functions in other files.

**Q:** I've seen function declarations that omit all information about parameters:

```
double average();
```

**Is this practice legal? [p. 192]**

- A: Yes. This declaration informs the compiler that `average` returns a `double` value but provides no information about the number and types of its parameters. (Leaving the parentheses empty doesn't necessarily mean that `average` has no parameters.)

In K&R C, this form of function declaration is the only one allowed; the form that we've been using—the function prototype, in which parameter information *is* included—was introduced in C89. The older kind of function declaration is now obsolescent, although still allowed.

**Q: Why would a programmer deliberately omit parameter names in a function prototype? Isn't it easier to just leave the names? [p. 193]**

- A: Omitting parameter names in prototypes is typically done for defensive purposes. If a macro happens to have the same name as a parameter, the parameter name will be replaced during preprocessing, thereby damaging the prototype in which it appears. This isn't likely to be a problem in a small program written by one person but can occur in large applications written by many people.

**Q: Is it legal to put a function declaration inside the body of another function?**

- A: Yes. Here's an example:

```
int main(void)
{
    double average(double a, double b);
    ...
}
```

This declaration of `average` is valid only for the body of `main`; if other functions need to call `average`, they'll each have to declare it.

The advantage of this practice is that it's clearer to the reader which functions call which other functions. (In this example, we see that `main` will be calling `average`.) On the other hand, it can be a nuisance if several functions need to call the same function. Even worse, trying to add and remove declarations during program maintenance can be a real pain. For these reasons, I'll always put function declarations outside function bodies.

**Q: If several functions have the same return type, can their declarations be combined? For example, since both `print_pun` and `print_count` have `void` as their return type, is the following declaration legal?**

```
void print_pun(void), print_count(int n);
```

- A: Yes. In fact, C even allows us to combine function declarations with variable declarations:

```
double x, y, average(double a, double b);
```

Combining declarations in this way usually isn't a good idea, though; it can easily cause confusion.

- Q:** What happens if I specify a length for a one-dimensional array parameter? [p. 195]

- A:** The compiler ignores it. Consider the following example:

```
double inner_product(double v[3], double w[3]);
```

Other than documenting that `inner_product`'s arguments are supposed to be arrays of length 3, specifying a length doesn't buy us much. The compiler won't check that the arguments actually have length 3, so there's no added security. In fact, the practice is misleading in that it suggests that `inner_product` can only be passed arrays of length 3, when in fact we can pass arrays of arbitrary length.

- \*Q:** Why can the first dimension in an array parameter be left unspecified, but not the other dimensions? [p. 197]

- A:** First, we need to discuss how arrays are passed in C. As Section 12.3 explains, when an array is passed to a function, the function is given a *pointer* to the first element in the array.

Next, we need to know how the subscripting operator works. Suppose that `a` is a one-dimensional array passed to a function. When we write

```
a[i] = 0;
```

the compiler generates instructions that compute the address of `a[i]` by multiplying `i` by the size of an array element and adding the result to the address that `a` represents (the pointer passed to the function). This calculation doesn't depend on the length of `a`, which explains why we can omit it when defining the function.

What about multidimensional arrays? Recall that C stores arrays in row-major order, with the elements in row 0 stored first, then the elements in row 1, and so forth. Suppose that `a` is a two-dimensional array parameter and we write

```
a[i][j] = 0;
```

The compiler generates instructions to do the following: (1) multiply `i` by the size of a single row of `a`; (2) add this result to the address that `a` represents; (3) multiply `j` by the size of an array element; and (4) add this result to the address computed in step 2. To generate these instructions, the compiler must know the size of a row in the array, which is determined by the number of columns. The bottom line: the programmer must declare the number of columns in `a`.

- Q:** Why do some programmers put parentheses around the expression in a `return` statement?

- A:** The examples in the first edition of Kernighan and Ritchie's *The C Programming Language* always have parentheses in `return` statements, even though they aren't required. Programmers (and authors of subsequent books) picked up the habit from K&R. I don't use these parentheses, since they're unnecessary and

contribute nothing to readability. (Kernighan and Ritchie apparently agree: the return statements in the second edition of *The C Programming Language* lack parentheses.)

**Q:** What happens if a non-void function attempts to execute a return statement that has no expression? [p. 202]

A: That depends on the version of C. In C89, executing a return statement without an expression in a non-void function causes undefined behavior (but only if the program attempts to use the value returned by the function). In C99, such a statement is illegal and should be detected as an error by the compiler.

**C99 Q:** How can I test main's return value to see if a program has terminated normally? [p. 203]

A: That depends on your operating system. Many operating systems allow this value to be tested within a “batch file” or “shell script” that contains commands to run several programs. For example, the line

```
if errorlevel 1 command
```

in a Windows batch file will execute *command* if the last program terminated with a status code greater than or equal to 1.

In UNIX, each shell has its own method for testing the status code. In the Bourne shell, the variable \$? contains the status of the last program run. The C shell has a similar variable, but its name is \$status.

**Q:** Why does my compiler produce a “control reaches end of non-void function” warning when it compiles main?

A: The compiler has noticed that main, despite having int as its return type, doesn't have a return statement. Putting the statement

```
return 0;
```

at the end of main will keep the compiler happy. Incidentally, this is good practice even if your compiler doesn't object to the lack of a return statement.

**C99** When a program is compiled using a C99 compiler, this warning shouldn't occur. In C99, it's OK to “fall off” the end of main without returning a value; the standard states that main automatically returns 0 in this situation.

**Q:** With regard to the previous question: Why not just define main's return type to be void?

A: Although this practice is fairly common, it's illegal according to the C89 standard. Even if it weren't illegal, it wouldn't be a good idea, since it presumes that no one will ever test the program's status upon termination.

**C99** C99 opens the door to legalizing this practice, by allowing main to be declared “in some other implementation-defined manner” (with a return type other than int or parameters other than those specified by the standard). However, any such usage isn't portable, so it's best to declare main's return type to be int.

**Q:** Is it legal for a function `f1` to call a function `f2`, which then calls `f1`?

**A:** Yes. This is just an indirect form of recursion in which one call of `f1` leads to another. (But make sure that either `f1` or `f2` eventually terminates!)

## Exercises

### Section 9.1

1. The following function, which computes the area of a triangle, contains two errors. Locate the errors and show how to fix them. (*Hint:* There are no errors in the formula.)

```
double triangle_area(double base, height)
double product;
{
    product = base * height;
    return product / 2;
}
```

- W 2. Write a function `check(x, y, n)` that returns 1 if both `x` and `y` fall between 0 and  $n - 1$ , inclusive. The function should return 0 otherwise. Assume that `x`, `y`, and `n` are all of type `int`.
3. Write a function `gcd(m, n)` that calculates the greatest common divisor of the integers `m` and `n`. (Programming Project 2 in Chapter 6 describes Euclid's algorithm for computing the GCD.)
- W 4. Write a function `day_of_year(month, day, year)` that returns the day of the year (an integer between 1 and 366) specified by the three arguments.
5. Write a function `num_digits(n)` that returns the number of digits in `n` (a positive integer). *Hint:* To determine the number of digits in a number `n`, divide it by 10 repeatedly. When `n` reaches 0, the number of divisions indicates how many digits `n` originally had.
- W 6. Write a function `digit(n, k)` that returns the  $k^{\text{th}}$  digit (from the right) in `n` (a positive integer). For example, `digit(829, 1)` returns 9, `digit(829, 2)` returns 2, and `digit(829, 3)` returns 8. If `k` is greater than the number of digits in `n`, have the function return 0.
7. Suppose that the function `f` has the following definition:

```
int f(int a, int b) { ... }
```

Which of the following statements are legal? (Assume that `i` has type `int` and `x` has type `double`.)

- (a) `i = f(83, 12);`
- (b) `x = f(83, 12);`
- (c) `i = f(3.15, 9.28);`
- (d) `x = f(3.15, 9.28);`
- (e) `f(83, 12);`

### Section 9.2

- W 8. Which of the following would be valid prototypes for a function that returns nothing and has one `double` parameter?
- (a) `void f(double x);`

- (b) void f(double);
- (c) void f(x);
- (d) f(double x);

**Section 9.3**

- \*9. What will be the output of the following program?

```
#include <stdio.h>

void swap(int a, int b);

int main(void)
{
    int i = 1, j = 2;

    swap(i, j);
    printf("i = %d, j = %d\n", i, j);
    return 0;
}

void swap(int a, int b)
{
    int temp = a;
    a = b;
    b = temp;
}
```

- W 10. Write functions that return the following values. (Assume that *a* and *n* are parameters, where *a* is an array of int values and *n* is the length of the array.)
- (a) The largest element in *a*.
  - (b) The average of all elements in *a*.
  - (c) The number of positive elements in *a*.

11. Write the following function:

```
float compute_GPA(char grades[], int n);
```

The *grades* array will contain letter grades (A, B, C, D, or F, either upper-case or lower-case); *n* is the length of the array. The function should return the average of the grades (assume that A = 4, B = 3, C = 2, D = 1, and F = 0).

12. Write the following function:

```
double inner_product(double a[], double b[], int n);
```

The function should return  $a[0] * b[0] + a[1] * b[1] + \dots + a[n-1] * b[n-1]$ .

13. Write the following function, which evaluates a chess position:

```
int evaluate_position(char board[8][8]);
```

*board* represents a configuration of pieces on a chessboard, where the letters K, Q, R, B, N, P represent White pieces, and the letters k, q, r, b, n, and p represent Black pieces. *evaluate\_position* should sum the values of the White pieces (Q = 9, R = 5, B = 3, N = 3, P = 1). It should also sum the values of the Black pieces (done in a similar way). The function will return the difference between the two numbers. This value will be positive if White has an advantage in material and negative if Black has an advantage.

**Section 9.4**

14. The following function is supposed to return *true* if any element of the array *a* has the value 0 and *false* if all elements are nonzero. Sadly, it contains an error. Find the error and show how to fix it:

```
bool has_zero(int a[], int n)
{
    int i;

    for (i = 0; i < n; i++)
        if (a[i] == 0)
            return true;
        else
            return false;
}
```

- W 15. The following (rather confusing) function finds the median of three numbers. Rewrite the function so that it has just one `return` statement.

```
double median(double x, double y, double z)
{
    if (x <= y)
        if (y <= z) return y;
        else if (x <= z) return z;
        else return x;
    if (z <= y) return y;
    if (x <= z) return x;
    return z;
}
```

### Section 9.6

16. Condense the `fact` function in the same way we condensed `power`.
- W 17. Rewrite the `fact` function so that it's no longer recursive.
18. Write a recursive version of the `gcd` function (see Exercise 3). Here's the strategy to use for computing `gcd(m, n)`: If `n` is 0, return `m`. Otherwise, call `gcd` recursively, passing `n` as the first argument and `m % n` as the second.

- W\*19. Consider the following "mystery" function:

```
void pb(int n)
{
    if (n != 0) {
        pb(n / 2);
        putchar('0' + n % 2);
    }
}
```

Trace the execution of the function by hand. Then write a program that calls the function, passing it a number entered by the user. What does the function do?

## Programming Projects

1. Write a program that asks the user to enter a series of integers (which it stores in an array), then sorts the integers by calling the function `selection_sort`. When given an array with  $n$  elements, `selection_sort` must do the following:
  1. Search the array to find the largest element, then move it to the last position in the array.
  2. Call itself recursively to sort the first  $n - 1$  elements of the array.

2. Modify Programming Project 5 from Chapter 5 so that it uses a function to compute the amount of income tax. When passed an amount of taxable income, the function will return the tax due.

3. Modify Programming Project 9 from Chapter 8 so that it includes the following functions:

```
void generate_random_walk(char walk[10][10]);
void print_array(char walk[10][10]);
```

`main` first calls `generate_random_walk`, which initializes the array to contain '.' characters and then replaces some of these characters by the letters A through Z, as described in the original project. `main` then calls `print_array` to display the array on the screen.

4. Modify Programming Project 16 from Chapter 8 so that it includes the following functions:

```
void read_word(int counts[26]);
bool equal_array(int counts1[26], int counts2[26]);
```

`main` will call `read_word` twice, once for each of the two words entered by the user. As it reads a word, `read_word` will use the letters in the word to update the `counts` array, as described in the original project. (`main` will declare two arrays, one for each word. These arrays are used to track how many times each letter occurs in the words.) `main` will then call `equal_array`, passing it the two arrays. `equal_array` will return `true` if the elements in the two arrays are identical (indicating that the words are anagrams) and `false` otherwise.

5. Modify Programming Project 17 from Chapter 8 so that it includes the following functions:

```
void create_magic_square(int n, char magic_square[n][n]);
void print_magic_square(int n, char magic_square[n][n]);
```

After obtaining the number  $n$  from the user, `main` will call `create_magic_square`, passing it an  $n \times n$  array that is declared inside `main`. `create_magic_square` will fill the array with the numbers 1, 2, ...,  $n^2$  as described in the original project. `main` will then call `print_magic_square`, which will display the array in the format described in the original project. *Note:* If your compiler doesn't support variable-length arrays, declare the array in `main` to be  $99 \times 99$  instead of  $n \times n$  and use the following prototypes instead:

```
void create_magic_square(int n, char magic_square[99][99]);
void print_magic_square(int n, char magic_square[99][99]);
```

6. Write a function that computes the value of the following polynomial:

$$3x^5 + 2x^4 - 5x^3 - x^2 + 7x - 6$$

Write a program that asks the user to enter a value for  $x$ , calls the function to compute the value of the polynomial, and then displays the value returned by the function.

7. The `power` function of Section 9.6 can be made faster by having it calculate  $x^n$  in a different way. We first notice that if  $n$  is a power of 2, then  $x^n$  can be computed by squaring. For example,  $x^4$  is the square of  $x^2$ , so  $x^4$  can be computed using only two multiplications instead of three. As it happens, this technique can be used even when  $n$  is not a power of 2. If  $n$  is even, we use the formula  $x^n = (x^{n/2})^2$ . If  $n$  is odd, then  $x^n = x \times x^{n-1}$ . Write a recursive function that computes  $x^n$ . (The recursion ends when  $n = 0$ , in which case the function returns 1.) To test your function, write a program that asks the user to enter values for  $x$  and  $n$ , calls `power` to compute  $x^n$ , and then displays the value returned by the function.

8. Write a program that simulates the game of craps, which is played with two dice. On the first roll, the player wins if the sum of the dice is 7 or 11. The player loses if the sum is 2, 3,

or 12. Any other roll is called the “point” and the game continues. On each subsequent roll, the player wins if he or she rolls the point again. The player loses by rolling 7. Any other roll is ignored and the game continues. At the end of each game, the program will ask the user whether or not to play again. When the user enters a response other than `y` or `Y`, the program will display the number of wins and losses and then terminate.

```
You rolled: 8
Your point is 8
You rolled: 3
You rolled: 10
You rolled: 8
You win!
```

```
Play again? Y
```

```
You rolled: 6
Your point is 6
You rolled: 5
You rolled: 12
You rolled: 3
You rolled: 7
You lose!
```

```
Play again? Y
```

```
You rolled: 11
You win!
```

```
Play again? n
```

```
Wins: 2 Losses: 1
```

Write your program as three functions: `main`, `roll_dice`, and `play_game`. Here are the prototypes for the latter two functions:

```
int roll_dice(void);
bool play_game(void);
```

`roll_dice` should generate two random numbers, each between 1 and 6, and return their sum. `play_game` should play one craps game (calling `roll_dice` to determine the outcome of each dice roll); it will return `true` if the player wins and `false` if the player loses. `play_game` is also responsible for displaying messages showing the results of the player’s dice rolls. `main` will call `play_game` repeatedly, keeping track of the number of wins and losses and displaying the “you win” and “you lose” messages. *Hint:* Use the `rand` function to generate random numbers. See the `deal.c` program in Section 8.2 for an example of how to call `rand` and the related `srand` function.

# 10 Program Organization

*As Will Rogers would have said, "There is no such thing as a free variable."*

Having covered functions in Chapter 9, we're ready to confront several issues that arise when a program contains more than one function. The chapter begins with a discussion of the differences between local variables (Section 10.1) and external variables (Section 10.2). Section 10.3 then considers blocks (compound statements containing declarations). Section 10.4 tackles the scope rules that apply to local names, external names, and names declared in blocks. Finally, Section 10.5 suggests a way to organize function prototypes, function definitions, variable declarations, and the other parts of a C program.

## 10.1 Local Variables

A variable declared in the body of a function is said to be *local* to the function. In the following function, `sum` is a local variable:

```
int sum_digits(int n)
{
    int sum = 0; /* local variable */

    while (n > 0) {
        sum += n % 10;
        n /= 10;
    }

    return sum;
}
```

By default, local variables have the following properties:

- **Automatic storage duration.** The *storage duration* (or *extent*) of a variable is the portion of program execution during which storage for the variable exists. Storage for a local variable is “automatically” allocated when the enclosing function is called and deallocated when the function returns, so the variable is said to have *automatic storage duration*. A local variable doesn’t retain its value when its enclosing function returns. When the function is called again, there’s no guarantee that the variable will still have its old value.
- **Block scope.** The *scope* of a variable is the portion of the program text in which the variable can be referenced. A local variable has *block scope*: it is visible from its point of declaration to the end of the enclosing function body. Since the scope of a local variable doesn’t extend beyond the function to which it belongs, other functions can use the same name for other purposes.

Section 18.2 covers these and other related concepts in more detail.

**C99**

Since C99 doesn’t require variable declarations to come at the beginning of a function, it’s possible for a local variable to have a very small scope. In the following example, the scope of *i* doesn’t begin until the line on which it’s declared, which could be near the end of the function body:

```
void f(void)
{
    ...
    int i;      — scope of i
    ...
}
```

## Static Local Variables

Putting the word *static* in the declaration of a local variable causes it to have *static storage duration* instead of automatic storage duration. A variable with static storage duration has a permanent storage location, so it retains its value throughout the execution of the program. Consider the following function:

```
void f(void)
{
    static int i; /* static local variable */
    ...
}
```

Since the local variable *i* has been declared *static*, it occupies the same memory location throughout the execution of the program. When *f* returns, *i* won’t lose its value.

**Q&A**

A static local variable still has block scope, so it’s not visible to other functions. In a nutshell, a static variable is a place to hide data from other functions but retain it for future calls of the same function.

## Parameters

Parameters have the same properties—automatic storage duration and block scope—as local variables. In fact, the only real difference between parameters and local variables is that each parameter is initialized automatically when a function is called (by being assigned the value of the corresponding argument).

## 10.2 External Variables

Passing arguments is one way to transmit information to a function. Functions can also communicate through *external variables*—variables that are declared outside the body of any function.

The properties of external variables (or *global variables*, as they’re sometimes called) are different from those of local variables:

- **Static storage duration.** External variables have static storage duration, just like local variables that have been declared `static`. A value stored in an external variable will stay there indefinitely.
- **File scope.** An external variable has *file scope*: it is visible from its point of declaration to the end of the enclosing file. As a result, an external variable can be accessed (and potentially modified) by all functions that follow its declaration.

### Example: Using External Variables to Implement a Stack

To illustrate how external variables might be used, let’s look at a data structure known as a *stack*. (Stacks are an abstract concept, not a C feature; they can be implemented in most programming languages.) A stack, like an array, can store multiple data items of the same type. However, the operations on a stack are limited: we can either *push* an item onto the stack (add it to one end—the “stack top”) or *pop* it from the stack (remove it from the same end). Examining or modifying an item that’s not at the top of the stack is forbidden.

One way to implement a stack in C is to store its items in an array, which we’ll call `contents`. A separate integer variable named `top` marks the position of the stack top. When the stack is empty, `top` has the value 0. To push an item on the stack, we simply store the item in `contents` at the position indicated by `top`, then increment `top`. Popping an item requires decrementing `top`, then using it as an index into `contents` to fetch the item that’s being popped.

Based on this outline, here’s a program fragment (not a complete program) that declares the `contents` and `top` variables for a stack and provides a set of functions that represent operations on the stack. All five functions need access to the `top` variable, and two functions need access to `contents`, so we’ll make `contents` and `top` external.

```
#include <stdbool.h> /* C99 only */

#define STACK_SIZE 100

/* external variables */
int contents[STACK_SIZE];
int top = 0;

void make_empty(void)
{
    top = 0;
}

bool is_empty(void)
{
    return top == 0;
}

bool is_full(void)
{
    return top == STACK_SIZE;
}

void push(int i)
{
    if (is_full())
        stack_overflow();
    else
        contents[top++] = i;
}

int pop(void)
{
    if (is_empty())
        stack_underflow();
    else
        return contents[--top];
}
```

## Pros and Cons of External Variables

External variables are convenient when many functions must share a variable or when a few functions share a large number of variables. In most cases, however, it's better for functions to communicate through parameters rather than by sharing variables. Here's why:

- If we change an external variable during program maintenance (by altering its type, say), we'll need to check every function in the same file to see how the change affects it.

- If an external variable is assigned an incorrect value, it may be difficult to identify the guilty function. It's like trying to solve a murder committed at a crowded party—there's no easy way to narrow the list of suspects.
- Functions that rely on external variables are hard to reuse in other programs. A function that depends on external variables isn't self-contained; to reuse the function, we'll have to drag along any external variables that it needs.

Many C programmers rely far too much on external variables. One common abuse: using the same external variable for different purposes in different functions. Suppose that several functions need a variable named `i` to control a `for` statement. Instead of declaring `i` in each function that uses it, some programmers declare it at the top of the program, thereby making the variable visible to all functions. This practice is poor not only for the reasons listed earlier, but also because it's misleading; someone reading the program later may think that the uses of the variable are related, when in fact they're not.

When you use external variables, make sure they have meaningful names. (Local variables don't always need meaningful names: it's often hard to think of a better name than `i` for the control variable in a `for` loop.) If you find yourself using names like `i` and `temp` for external variables, that's a clue that perhaps they should really be local variables.



Making variables external when they should be local can lead to some rather frustrating bugs. Consider the following example, which is supposed to display a  $10 \times 10$  arrangement of asterisks:

```
int i;

void print_one_row(void)
{
    for (i = 1; i <= 10; i++)
        printf("*");
}

void print_all_rows(void)
{
    for (i = 1; i <= 10; i++) {
        print_one_row();
        printf("\n");
    }
}
```

Instead of printing 10 rows, `print_all_rows` prints only one row. When `print_one_row` returns after being called the first time, `i` will have the value 11. The `for` statement in `print_all_rows` then increments `i` and tests whether it's less than or equal to 10. It's not, so the loop terminates and the function returns.

## PROGRAM Guessing a Number

To get more experience with external variables, we'll write a simple game-playing program. The program generates a random number between 1 and 100, which the user attempts to guess in as few tries as possible. Here's what the user will see when the program is run:

```
Guess the secret number between 1 and 100.
```

```
A new number has been chosen.  
Enter guess: 55  
Too low; try again.  
Enter guess: 65  
Too high; try again.  
Enter guess: 60  
Too high; try again.  
Enter guess: 58  
You won in 4 guesses!
```

```
Play again? (Y/N) Y
```

```
A new number has been chosen.  
Enter guess: 78  
Too high; try again.  
Enter guess: 34  
You won in 2 guesses!
```

```
Play again? (Y/N) n
```

This program will need to carry out several different tasks: initializing the random number generator, choosing a secret number, and interacting with the user until the correct number is picked. If we write a separate function to handle each task, we might end up with the following program.

```
guess.c /* Asks user to guess a hidden number */  
  
#include <stdio.h>  
#include <stdlib.h>  
#include <time.h>  
  
#define MAX_NUMBER 100  
  
/* external variable */  
int secret_number;  
  
/* prototypes */  
void initialize_number_generator(void);  
void choose_new_secret_number(void);  
void read_guesses(void);  
  
int main(void)  
{  
    char command;
```

```
printf("Guess the secret number between 1 and %d.\n\n",
      MAX_NUMBER);
initialize_number_generator();
do {
    choose_new_secret_number();
    printf("A new number has been chosen.\n");
    read_guesses();
    printf("Play again? (Y/N) ");
    scanf(" %c", &command);
    printf("\n");
} while (command == 'y' || command == 'Y');

return 0;
}

/*****************
 * initialize_number_generator: Initializes the random *
 * number generator using *
 * the time of day. *
 *****************/
void initialize_number_generator(void)
{
    srand((unsigned) time(NULL));
}

/*****************
 * choose_new_secret_number: Randomly selects a number *
 * between 1 and MAX_NUMBER and *
 * stores it in secret_number. *
 *****************/
void choose_new_secret_number(void)
{
    secret_number = rand() % MAX_NUMBER + 1;
}

/*****************
 * read_guesses: Repeatedly reads user guesses and tells *
 * the user whether each guess is too low, *
 * too high, or correct. When the guess is *
 * correct, prints the total number of *
 * guesses and returns. *
 *****************/
void read_guesses(void)
{
    int guess, num_guesses = 0;

    for (;;) {
        num_guesses++;
        printf("Enter guess: ");
        scanf("%d", &guess);
        if (guess == secret_number) {
            printf("You won in %d guesses!\n\n", num_guesses);
            return;
        } else if (guess < secret_number)
```

```

        printf("Too low; try again.\n");
    else
        printf("Too high; try again.\n");
    }
}

```

time function ▶26.3  
 srand function ▶26.2  
 rand function ▶26.2

For random number generation, the `guess.c` program relies on the `time`, `srand`, and `rand` functions, which we first used in `deal.c` (Section 8.2). This time, we're scaling the return value of `rand` so that it falls between 1 and `MAX_NUMBER`.

Although `guess.c` works fine, it relies on an external variable. We made `secret_number` external so that both `choose_new_secret_number` and `read_guesses` could access it. If we alter `choose_new_secret_number` and `read_guesses` just a little, we should be able to move `secret_number` into the `main` function. We'll modify `choose_new_secret_number` so that it returns the new number, and we'll rewrite `read_guesses` so that `secret_number` can be passed to it as an argument.

Here's our new program, with changes in **bold**:

```

guess2.c /* Asks user to guess a hidden number */

#include <stdio.h>
#include <stdlib.h>
#include <time.h>

#define MAX_NUMBER 100

/* prototypes */
void initialize_number_generator(void);
int new_secret_number(void);
void read_guesses(int secret_number);

int main(void)
{
    char command;
    int secret_number;

    printf("Guess the secret number between 1 and %d.\n\n",
           MAX_NUMBER);
    initialize_number_generator();
    do {
        secret_number = new_secret_number();
        printf("A new number has been chosen.\n");
        read_guesses(secret_number);
        printf("Play again? (Y/N) ");
        scanf(" %c", &command);
        printf("\n");
    } while (command == 'y' || command == 'Y');

    return 0;
}

```

```

*****
 * initialize_number_generator: Initializes the random      *
 *                               number generator using      *
 *                               the time of day.      *
 *****
void initialize_number_generator(void)
{
    srand((unsigned) time(NULL));
}

*****
 * new_secret_number: Returns a randomly chosen number      *
 *                      between 1 and MAX_NUMBER.      *
 *****
int new_secret_number(void)
{
    return rand() % MAX_NUMBER + 1;
}

*****
 * read_guesses: Repeatedly reads user guesses and tells   *
 *                  the user whether each guess is too low,   *
 *                  too high, or correct. When the guess is   *
 *                  correct, prints the total number of   *
 *                  guesses and returns.      *
 *****
void read_guesses(int secret_number)
{
    int guess, num_guesses = 0;

    for (;;) {
        num_guesses++;
        printf("Enter guess: ");
        scanf("%d", &guess);
        if (guess == secret_number) {
            printf("You won in %d guesses!\n\n", num_guesses);
            return;
        } else if (guess < secret_number)
            printf("Too low; try again.\n");
        else
            printf("Too high; try again.\n");
    }
}

```

## 10.3 Blocks

In Section 5.2, we encountered compound statements of the form

{ *statements* }

It turns out that C allows compound statements to contain declarations as well:

<b>block</b>	{ declarations statements }
--------------	-----------------------------

I'll use the term **block** to describe such a compound statement. Here's an example of a block:

```
if (i > j) {
    /* swap values of i and j */
    int temp = i;
    i = j;
    j = temp;
}
```

By default, the storage duration of a variable declared in a block is automatic: storage for the variable is allocated when the block is entered and deallocated when the block is exited. The variable has block scope; it can't be referenced outside the block. A variable that belongs to a block can be declared `static` to give it static storage duration.

The body of a function is a block. Blocks are also useful inside a function body when we need variables for temporary use. In our last example, we needed a variable temporarily so that we could swap the values of `i` and `j`. Putting temporary variables in blocks has two advantages: (1) It avoids cluttering the declarations at the beginning of the function body with variables that are used only briefly. (2) It reduces name conflicts. In our example, the name `temp` can be used elsewhere in the same function for different purposes—the `temp` variable is strictly local to the block in which it's declared.

**C99** C99 allows variables to be declared anywhere within a block, just as it allows variables to be declared anywhere within a function.

## 10.4 Scope

In a C program, the same identifier may have several different meanings. C's scope rules enable the programmer (and the compiler) to determine which meaning is relevant at a given point in the program.

Here's the most important scope rule: When a declaration inside a block names an identifier that's already visible (because it has file scope or because it's declared in an enclosing block), the new declaration temporarily "hides" the old one, and the identifier takes on a new meaning. At the end of the block, the identifier regains its old meaning.

Consider the (somewhat extreme) example at the top of the next page, in which the identifier `i` has four different meanings:

- In Declaration 1, `i` is a variable with static storage duration and file scope.

```

int i; /* Declaration 1 */

void f(int i) /* Declaration 2 */
{
    i = 1;
}

void g(void)
{
    int i = 2; /* Declaration 3 */

    if (i > 0) {
        int i; /* Declaration 4 */

        i = 3;
    }

    i = 4;
}

void h(void)
{
    i = 5;
}

```

- In Declaration 2, *i* is a parameter with block scope.
- In Declaration 3, *i* is an automatic variable with block scope.
- In Declaration 4, *i* is also automatic and has block scope.

*i* is used five times. C's scope rules allow us to determine the meaning of *i* in each case:

- The *i* = 1 assignment refers to the parameter in Declaration 2, not the variable in Declaration 1, since Declaration 2 hides Declaration 1.
- The *i* > 0 test refers to the variable in Declaration 3, since Declaration 3 hides Declaration 1 and Declaration 2 is out of scope.
- The *i* = 3 assignment refers to the variable in Declaration 4, which hides Declaration 3.
- The *i* = 4 assignment refers to the variable in Declaration 3. It can't refer to Declaration 4, which is out of scope.
- The *i* = 5 assignment refers to the variable in Declaration 1.

## 10.5 Organizing a C Program

Now that we've seen the major elements that make up a C program, it's time to develop a strategy for their arrangement. For now, we'll assume that a program

always fits into a single file. Chapter 15 shows how to organize a program that's split over several files.

So far, we've seen that a program may contain the following:

- Preprocessing directives such as `#include` and `#define`
- Type definitions
- Declarations of external variables
- Function prototypes
- Function definitions

C imposes only a few rules on the order of these items: A preprocessing directive doesn't take effect until the line on which it appears. A type name can't be used until it's been defined. A variable can't be used until it's declared. Although C isn't as picky about functions, I strongly recommend that every function be defined or declared prior to its first call. (C99 makes this a requirement anyway.)

**C99** There are several ways to organize a program so that these rules are obeyed. Here's one possible ordering:

```
#include directives  
#define directives  
Type definitions  
Declarations of external variables  
Prototypes for functions other than main  
Definition of main  
Definitions of other functions
```

It makes sense to put `#include` directives first, since they bring in information that will likely be needed in several places within the program. `#define` directives create macros, which are generally used throughout the program. Putting type definitions above the declarations of external variables is logical, since the declarations of these variables may refer to the type names just defined. Declaring external variables next makes them available to all the functions that follow. Declaring all functions except for `main` avoids the problems that arise when a function is called before the compiler has seen its prototype. This practice also makes it possible to arrange the function definitions in any order whatsoever: alphabetically by function name or with related functions grouped together, for example. Defining `main` before the other functions makes it easier for a reader to locate the program's starting point.

A final suggestion: Precede each function definition by a boxed comment that gives the name of the function, explains its purpose, discusses the meaning of each parameter, describes its return value (if any), and lists any side effects it has (such as modifying external variables).

## PROGRAM

### Classifying a Poker Hand

To show how a C program might be organized, let's attempt a program that's a little more complex than our previous examples. The program will read and classify

a poker hand. Each card in the hand will have both a *suit* (clubs, diamonds, hearts, or spades) and a *rank* (two, three, four, five, six, seven, eight, nine, ten, jack, queen, king, or ace). We won't allow the use of jokers, and we'll assume that aces are high. The program will read a hand of five cards, then classify the hand into one of the following categories (listed in order from best to worst):

- straight flush (both a straight and a flush)
- four-of-a-kind (four cards of the same rank)
- full house (a three-of-a-kind and a pair)
- flush (five cards of the same suit)
- straight (five cards with consecutive ranks)
- three-of-a-kind (three cards of the same rank)
- two pairs
- pair (two cards of the same rank)
- high card (any other hand)

If a hand falls into two or more categories, the program will choose the best one.

For input purposes, we'll abbreviate ranks and suits as follows (letters may be either upper- or lower-case):

Ranks: 2 3 4 5 6 7 8 9 t j q k a  
Suits: c d h s

If the user enters an illegal card or tries to enter the same card twice, the program will ignore the card, issue an error message, and then request another card. Entering the number 0 instead of a card will cause the program to terminate.

A session with the program will have the following appearance:

```
Enter a card: 2s
Enter a card: 5s
Enter a card: 4s
Enter a card: 3s
Enter a card: 6s
Straight flush

Enter a card: 8c
Enter a card: as
Enter a card: 8c
Duplicate card; ignored.
Enter a card: 7c
Enter a card: ad
Enter a card: 3h
Pair

Enter a card: 6s
Enter a card: d2
Bad card; ignored.
Enter a card: 2d
Enter a card: 9c
Enter a card: 4h
Enter a card: ts
```

High card

Enter a card:   

From this description of the program, we see that it has three tasks:

Read a hand of five cards.

Analyze the hand for pairs, straights, and so forth.

Print the classification of the hand.

We'll divide the program into three functions—`read_cards`, `analyze_hand`, and `print_result`—that perform these three tasks. `main` does nothing but call these functions inside an endless loop. The functions will need to share a fairly large amount of information, so we'll have them communicate through external variables. `read_cards` will store information about the hand into several external variables. `analyze_hand` will then examine these variables, storing its findings into other external variables for the benefit of `print_result`.

Based on this preliminary design, we can begin to sketch an outline of the program:

```
/* #include directives go here */

/* #define directives go here */

/* declarations of external variables go here */

/* prototypes */
void read_cards(void);
void analyze_hand(void);
void print_result(void);

***** * main: Calls read_cards, analyze_hand, and print_result *
*      repeatedly.
*****
int main(void)
{
    for (;;) {
        read_cards();
        analyze_hand();
        print_result();
    }
}

***** * read_cards: Reads the cards into external variables;
*                  checks for bad cards and duplicate cards.
*****
void read_cards(void)
{
    ...
}
```

```

*****  

* analyze_hand: Determines whether the hand contains a *  

* straight, a flush, four-of-a-kind, *  

* and/or three-of-a-kind; determines the *  

* number of pairs; stores the results into *  

* external variables. *  

*****  

void analyze_hand(void)  

{  

    ...  

}  

*****  

* print_result: Notifies the user of the result, using *  

* the external variables set by *  

* analyze_hand. *  

*****  

void print_result(void)  

{  

    ...  

}

```

The most pressing question that remains is how to represent the hand of cards. Let's see what operations `read_cards` and `analyze_hand` will perform on the hand. During the analysis of the hand, `analyze_hand` will need to know how many cards are in each rank and each suit. This suggests that we use two arrays, `num_in_rank` and `num_in_suit`. The value of `num_in_rank[r]` will be the number of cards with rank `r`, and the value of `num_in_suit[s]` will be the number of cards with suit `s`. (We'll encode ranks as numbers between 0 and 12, and suits as numbers between 0 and 3.) We'll also need a third array, `card_exists`, so that `read_cards` can detect duplicate cards. Each time `read_cards` reads a card with rank `r` and suit `s`, it checks whether the value of `card_exists[r][s]` is true. If so, the card was previously entered; if not, `read_cards` assigns true to `card_exists[r][s]`.

Both the `read_cards` function and the `analyze_hand` function will need access to the `num_in_rank` and `num_in_suit` arrays, so I'll make them external variables. The `card_exists` array is used only by `read_cards`, so it can be local to that function. As a rule, variables should be made external only if necessary.

Having decided on the major data structures, we can now finish the program:

```

poker.c /* Classifies a poker hand */  

#include <stdbool.h>    /* C99 only */  

#include <stdio.h>  

#include <stdlib.h>  

#define NUM_RANKS 13  

#define NUM_SUITS 4  

#define NUM_CARDS 5

```

```
/* external variables */
int num_in_rank[NUM_RANKS];
int num_in_suit[NUM_SUITS];
bool straight, flush, four, three;
int pairs; /* can be 0, 1, or 2 */

/* prototypes */
void read_cards(void);
void analyze_hand(void);
void print_result(void);

/*****
 * main: Calls read_cards, analyze_hand, and print_result
 *        repeatedly.
 *****/
int main(void)
{
    for (;;) {
        read_cards();
        analyze_hand();
        print_result();
    }
}

/*****
 * read_cards: Reads the cards into the external
 *             variables num_in_rank and num_in_suit;
 *             checks for bad cards and duplicate cards.
 *****/
void read_cards(void)
{
    bool card_exists[NUM_RANKS][NUM_SUITS];
    char ch, rank_ch, suit_ch;
    int rank, suit;
    bool bad_card;
    int cards_read = 0;

    for (rank = 0; rank < NUM_RANKS; rank++) {
        num_in_rank[rank] = 0;
        for (suit = 0; suit < NUM_SUITS; suit++)
            card_exists[rank][suit] = false;
    }

    for (suit = 0; suit < NUM_SUITS; suit++)
        num_in_suit[suit] = 0;

    while (cards_read < NUM_CARDS) {
        bad_card = false;

        printf("Enter a card: ");

        rank_ch = getchar();
        switch (rank_ch) {
```

```

        case '0': exit(EXIT_SUCCESS);
        case '2': rank = 0; break;
        case '3': rank = 1; break;
        case '4': rank = 2; break;
        case '5': rank = 3; break;
        case '6': rank = 4; break;
        case '7': rank = 5; break;
        case '8': rank = 6; break;
        case '9': rank = 7; break;
        case 't': case 'T': rank = 8; break;
        case 'j': case 'J': rank = 9; break;
        case 'q': case 'Q': rank = 10; break;
        case 'k': case 'K': rank = 11; break;
        case 'a': case 'A': rank = 12; break;
        default:    bad_card = true;
    }

    suit_ch = getchar();
    switch (suit_ch) {
        case 'c': case 'C': suit = 0; break;
        case 'd': case 'D': suit = 1; break;
        case 'h': case 'H': suit = 2; break;
        case 's': case 'S': suit = 3; break;
        default:    bad_card = true;
    }

    while ((ch = getchar()) != '\n')
        if (ch != ' ') bad_card = true;

    if (bad_card)
        printf("Bad card; ignored.\n");
    else if (card_exists[rank][suit])
        printf("Duplicate card; ignored.\n");
    else {
        num_in_rank[rank]++;
        num_in_suit[suit]++;
        card_exists[rank][suit] = true;
        cards_read++;
    }
}

}

/****************************************
 * analyze_hand: Determines whether the hand contains a *
 * straight, a flush, four-of-a-kind, *
 * and/or three-of-a-kind; determines the *
 * number of pairs; stores the results into *
 * the external variables straight, flush, *
 * four, three, and pairs. *
 ****************************************/
void analyze_hand(void)
{
    int num_consec = 0;
    int rank, suit;

```

```

        straight = false;
        flush = false;
        four = false;
        three = false;
        pairs = 0;

        /* check for flush */
        for (suit = 0; suit < NUM_SUITS; suit++)
            if (num_in_suit[suit] == NUM_CARDS)
                flush = true;

        /* check for straight */
        rank = 0;
        while (num_in_rank[rank] == 0) rank++;
        for (; rank < NUM_RANKS && num_in_rank[rank] > 0; rank++)
            num_consec++;
        if (num_consec == NUM_CARDS) {
            straight = true;
            return;
        }

        /* check for 4-of-a-kind, 3-of-a-kind, and pairs */
        for (rank = 0; rank < NUM_RANKS; rank++) {
            if (num_in_rank[rank] == 4) four = true;
            if (num_in_rank[rank] == 3) three = true;
            if (num_in_rank[rank] == 2) pairs++;
        }
    }

/***** * print_result: Prints the classification of the hand, *
* based on the values of the external *
* variables straight, flush, four, three, *
* and pairs. *
*****/
void print_result(void)
{
    if (straight && flush) printf("Straight flush");
    else if (four)         printf("Four of a kind");
    else if (three &&
             pairs == 1)   printf("Full house");
    else if (flush)        printf("Flush");
    else if (straight)     printf("Straight");
    else if (three)        printf("Three of a kind");
    else if (pairs == 2)   printf("Two pairs");
    else if (pairs == 1)   printf("Pair");
    else                   printf("High card");

    printf("\n\n");
}

```

Notice the use of the `exit` function in `read_cards` (in case '0' of the first switch statement). `exit` is convenient for this program because of its ability to terminate execution from anywhere in the program.

## Q & A

**Q:** What impact do local variables with static storage duration have on recursive functions? [p. 220]

**A:** When a function is called recursively, fresh copies are made of its automatic variables for each call. This doesn't occur for static variables, though. Instead, all calls of the function share the *same* static variables.

**Q:** In the following example, *j* is initialized to the same value as *i*, but there are two variables named *i*:

```
int i = 1;

void f(void)
{
    int j = i;
    int i = 2;
    ...
}
```

Is this code legal? If so, what is *j*'s initial value, 1 or 2?

**A:** The code is indeed legal. The scope of a local variable doesn't begin until its declaration. Therefore, the declaration of *j* refers to the external variable named *i*. The initial value of *j* will be 1.

## Exercises

### Section 10.4

- W 1. The following program outline shows only function definitions and variable declarations.

```
int a;

void f(int b)
{
    int c;
}

void g(void)
{
    int d;
    {
        int e;
    }
}

int main(void)
{
    int f;
}
```

For each of the following scopes, list all variable and parameter names visible in that scope:

- (a) The `f` function
- (b) The `g` function
- (c) The block in which `e` is declared
- (d) The `main` function

2. The following program outline shows only function definitions and variable declarations.

```
int b, c;

void f(void)
{
    int b, d;
}

void g(int a)
{
    int c;
    {
        int a, d;
    }
}

int main(void)
{
    int c, d;
}
```

For each of the following scopes, list all variable and parameter names visible in that scope. If there's more than one variable or parameter with the same name, indicate which one is visible.

- (a) The `f` function
- (b) The `g` function
- (c) The block in which `a` and `d` are declared
- (d) The `main` function

- \*3. Suppose that a program has only one function (`main`). How many different variables named `i` could this program contain?

## Programming Projects

1. Modify the stack example of Section 10.2 so that it stores characters instead of integers. Next, add a `main` function that asks the user to enter a series of parentheses and/or braces, then indicates whether or not they're properly nested:

Enter parentheses and/or braces: `(({}{})()`)  
Parentheses/braces are nested properly

*Hint:* As the program reads characters, have it push each left parenthesis or left brace. When it reads a right parenthesis or brace, have it pop the stack and check that the item popped is a matching parenthesis or brace. (If not, the parentheses/braces aren't nested properly.) When the program reads the new-line character, have it check whether the stack is empty; if so, the parentheses/braces are matched. If the stack *isn't* empty (or if `stack_underflow` is ever

called), the parentheses/braces aren't matched. If `stack_overflow` is called, have the program print the message `Stack overflow` and terminate immediately.

2. Modify the `poker.c` program of Section 10.5 by moving the `num_in_rank` and `num_in_suit` arrays into `main`, which will pass them as arguments to `read_cards` and `analyze_hand`.
- W 3. Remove the `num_in_rank`, `num_in_suit`, and `card_exists` arrays from the `poker.c` program of Section 10.5. Have the program store the cards in a  $5 \times 2$  array instead. Each row of the array will represent a card. For example, if the array is named `hand`, then `hand[0][0]` will store the rank of the first card and `hand[0][1]` will store the suit of the first card.
4. Modify the `poker.c` program of Section 10.5 by having it recognize an additional category, "royal flush" (ace, king, queen, jack, ten of the same suit). A royal flush ranks higher than all other hands.
- W 5. Modify the `poker.c` program of Section 10.5 by allowing "ace-low" straights (ace, two, three, four, five).
6. Some calculators (notably those from Hewlett-Packard) use a system of writing mathematical expressions known as Reverse Polish Notation (RPN). In this notation, operators are placed *after* their operands instead of *between* their operands. For example,  $1 + 2$  would be written  $1\ 2\ +$  in RPN, and  $1 + 2 * 3$  would be written  $1\ 2\ 3\ * +$ . RPN expressions can easily be evaluated using a stack. The algorithm involves reading the operators and operands in an expression from left to right, performing the following actions:

When an operand is encountered, push it onto the stack.

When an operator is encountered, pop its operands from the stack, perform the operation on those operands, and then push the result onto the stack.

Write a program that evaluates RPN expressions. The operands will be single-digit integers. The operators are `+`, `-`, `*`, `/`, and `=`. The `=` operator causes the top stack item to be displayed; afterwards, the stack is cleared and the user is prompted to enter another expression. The process continues until the user enters a character that is not an operator or operand:

```
Enter an RPN expression: 1 2 3 * + =
Value of expression: 7
Enter an RPN expression: 5 8 * 4 9 - / =
Value of expression: -8
Enter an RPN expression: q
```

If the stack overflows, the program will display the message `Expression is too complex` and terminate. If the stack underflows (because of an expression such as `1 2 ++`), the program will display the message `Not enough operands in expression` and terminate. *Hints:* Incorporate the stack code from Section 10.2 into your program. Use `scanf(" %c", &ch)` to read the operators and operands.

7. Write a program that prompts the user for a number and then displays the number, using characters to simulate the effect of a seven-segment display:

```
Enter a number: 491-9014
```



Characters other than digits should be ignored. Write the program so that the maximum number of digits is controlled by a macro named `MAX_DIGITS`, which has the value 10. If

the number contains more than this number of digits, the extra digits are ignored. *Hints:* Use two external arrays. One is the `segments` array (see Exercise 6 in Chapter 8), which stores data representing the correspondence between digits and segments. The other array, `digits`, will be an array of characters with 4 rows (since each segmented digit is four characters high) and `MAX_DIGITS * 4` columns (digits are three characters wide, but a space is needed between `digits` for readability). Write your program as four functions: `main`, `clear_digits_array`, `process_digit`, and `print_digits_array`. Here are the prototypes for the latter three functions:

```
void clear_digits_array(void);  
void process_digit(int digit, int position);  
void print_digits_array(void);
```

`clear_digits_array` will store blank characters into all elements of the `digits` array. `process_digit` will store the seven-segment representation of `digit` into a specified position in the `digits` array (positions range from 0 to `MAX_DIGITS - 1`). `print_digits_array` will display the rows of the `digits` array, each on a single line, producing output such as that shown in the example.

# 11 Pointers

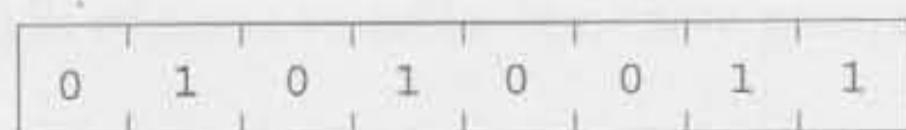
*The 11th commandment was "Thou Shalt Compute" or "Thou Shalt Not Compute"—I forget which.*

Pointers are one of C's most important—and most often misunderstood—features. Because of their importance, we'll devote three chapters to pointers. In this chapter, we'll concentrate on the basics; Chapters 12 and 17 cover more advanced uses of pointers.

We'll start with a discussion of memory addresses and their relationship to pointer variables (Section 11.1). Section 11.2 then introduces the address and indirection operators. Section 11.3 covers pointer assignment. Section 11.4 explains how to pass pointers to functions, while Section 11.5 discusses returning pointers from functions.

## 11.1 Pointer Variables

The first step in understanding pointers is visualizing what they represent at the machine level. In most modern computers, main memory is divided into *bytes*, with each byte capable of storing eight bits of information:



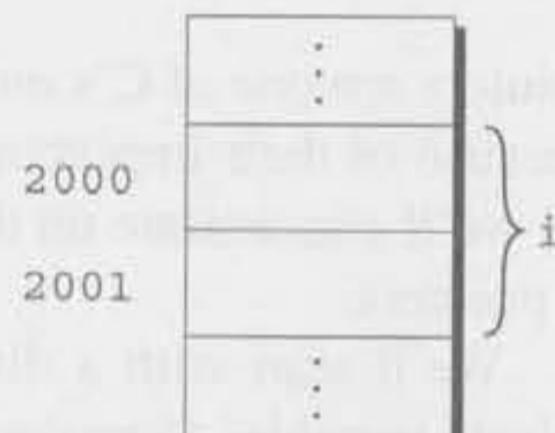
Each byte has a unique *address* to distinguish it from the other bytes in memory. If there are  $n$  bytes in memory, we can think of addresses as numbers that range from 0 to  $n - 1$  (see the figure at the top of the next page).

An executable program consists of both code (machine instructions corresponding to statements in the original C program) and data (variables in the original program). Each variable in the program occupies one or more bytes of memory;

## Address    Contents

0	01010011
1	01110101
2	01110011
3	01100001
4	01101110
:	:
n-1	01000011

the address of the first byte is said to be the address of the variable. In the following figure, the variable *i* occupies the bytes at addresses 2000 and 2001, so *i*'s address is 2000:



Here's where pointers come in. Although addresses are represented by numbers, their range of values may differ from that of integers, so we can't necessarily store them in ordinary integer variables. We can, however, store them in special **pointer variables**. When we store the address of a variable *i* in the pointer variable *p*, we say that *p* "points to" *i*. In other words, a pointer is nothing more than an address, and a pointer variable is just a variable that can store an address.

Instead of showing addresses as numbers in our examples, I'll use a simpler notation. To indicate that a pointer variable *p* stores the address of a variable *i*, I'll show the contents of *p* as an arrow directed toward *i*:



## Declaring Pointer Variables

A pointer variable is declared in much the same way as an ordinary variable. The only difference is that the name of a pointer variable must be preceded by an asterisk:

```
int *p;
```

abstract objects ▶ 19.1

This declaration states that `p` is a pointer variable capable of pointing to *objects* of type `int`. I'm using the term *object* instead of *variable* since—as we'll see in Chapter 17—`p` might point to an area of memory that doesn't belong to a variable. (Be aware that “object” will have a different meaning when we discuss program design in Chapter 19.)

Pointer variables can appear in declarations along with other variables:

```
int i, j, a[10], b[20], *p, *q;
```

In this example, `i` and `j` are ordinary integer variables, `a` and `b` are arrays of integers, and `p` and `q` are pointers to integer objects.

C requires that every pointer variable point only to objects of a particular type (the *referenced type*):

```
int *p;      /* points only to integers */
double *q;    /* points only to doubles */
char *r;     /* points only to characters */
```

pointers to pointers ▶ 17.6

There are no restrictions on what the referenced type may be. In fact, a pointer variable can even point to another pointer.

## 11.2 The Address and Indirection Operators

C provides a pair of operators designed specifically for use with pointers. To find the address of a variable, we use the `&` (address) operator. If `x` is a variable, then `&x` is the address of `x` in memory. To gain access to the object that a pointer points to, we use the `*` (*indirection*) operator. If `p` is a pointer, then `*p` represents the object to which `p` currently points.

### The Address Operator

Declaring a pointer variable sets aside space for a pointer but doesn't make it point to an object:

```
int *p; /* points nowhere in particular */
```

Ivalues ▶ 4.2 It's crucial to initialize `p` before we use it. One way to initialize a pointer variable is to assign it the address of some variable—or, more generally, lvalue—using the `&` operator:

```
int i, *p;
...
p = &i;
```

By assigning the address of `i` to the variable `p`, this statement makes `p` point to `i`:



It's also possible to initialize a pointer variable at the time we declare it:

**Q&A**

```
int i;
int *p = &i;
```

We can even combine the declaration of *i* with the declaration of *p*, provided that *i* is declared first:

```
int i, *p = &i;
```

## The Indirection Operator

Once a pointer variable points to an object, we can use the `*` (indirection) operator to access what's stored in the object. If *p* points to *i*, for example, we can print the value of *i* as follows:

```
printf("%d\n", *p);
```

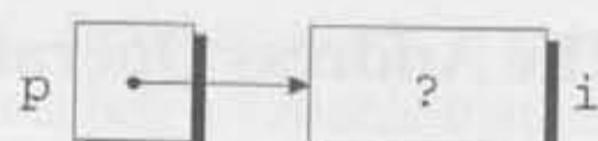
**Q&A** `printf` will display the *value* of *i*, not the *address* of *i*.

The mathematically inclined reader may wish to think of `*` as the inverse of `&`. Applying `&` to a variable produces a pointer to the variable; applying `*` to the pointer takes us back to the original variable:

```
j = *&i; /* same as j = i; */
```

As long as *p* points to *i*, `*p` is an *alias* for *i*. Not only does `*p` have the same value as *i*, but changing the value of `*p` also changes the value of *i*. (`*p` is an lvalue, so assignment to it is legal.) The following example illustrates the equivalence of `*p` and *i*; diagrams show the values of *p* and *i* at various points in the computation.

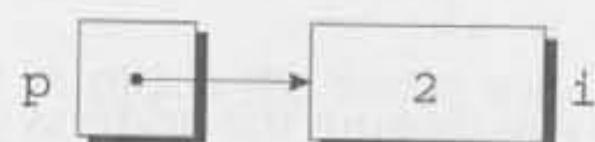
```
p = &i;
```



```
i = 1;
```



```
printf("%d\n", i); /* prints 1 */
printf("%d\n", *p); /* prints 1 */
*p = 2;
```



```
printf("%d\n", i); /* prints 2 */
printf("%d\n", *p); /* prints 2 */
```



Never apply the indirection operator to an uninitialized pointer variable. If a pointer variable *p* hasn't been initialized, attempting to use the value of *p* in any way causes undefined behavior. In the following example, the call of `printf` may print garbage, cause the program to crash, or have some other effect:

```
int *p;
printf("%d", *p);    /*** WRONG ***/

```

Assigning a value to *\*p* is particularly dangerous. If *p* happens to contain a valid memory address, the following assignment will attempt to modify the data stored at that address:

```
int *p;
*p = 1;    /*** WRONG ***/

```

If the location modified by this assignment belongs to the program, it may behave erratically; if it belongs to the operating system, the program will most likely crash. Your compiler may issue a warning that *p* is uninitialized, so pay close attention to any warning messages you get.

## 11.3 Pointer Assignment

C allows the use of the assignment operator to copy pointers, provided that they have the same type. Suppose that *i*, *j*, *p*, and *q* have been declared as follows:

```
int i, j, *p, *q;
```

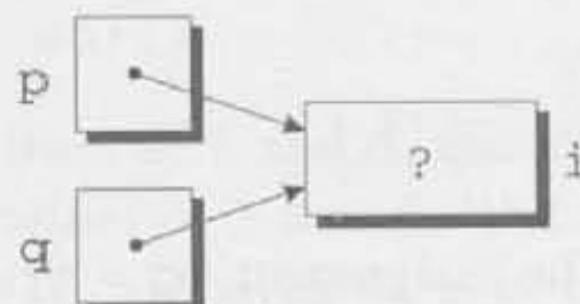
The statement

```
p = &i;
```

is an example of pointer assignment; the address of *i* is copied into *p*. Here's another example of pointer assignment:

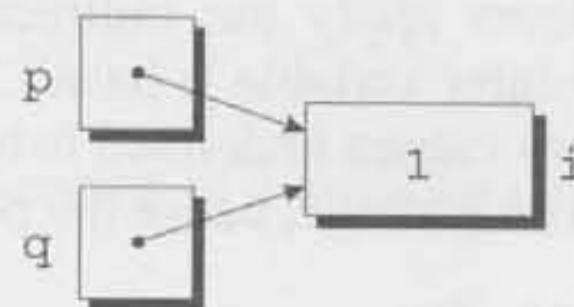
```
q = p;
```

This statement copies the contents of *p* (the address of *i*) into *q*, in effect making *q* point to the same place as *p*:

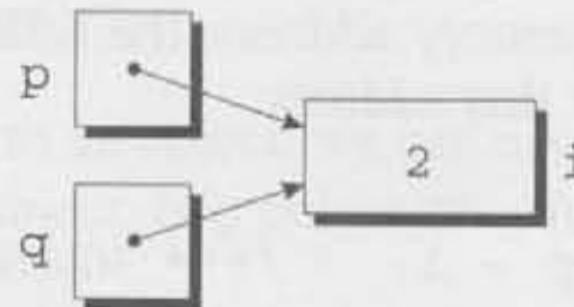


Both *p* and *q* now point to *i*, so we can change *i* by assigning a new value to either *\*p* or *\*q*:

`*p = 1;`



`*q = 2;`



Any number of pointer variables may point to the same object.

Be careful not to confuse

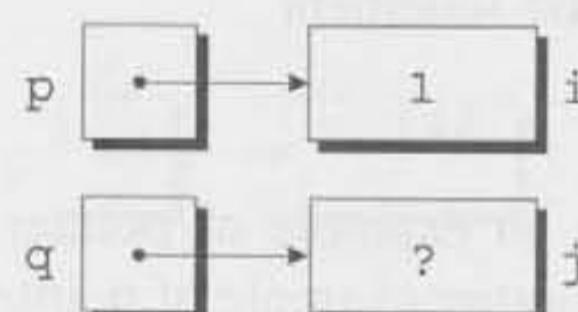
`q = p;`

with

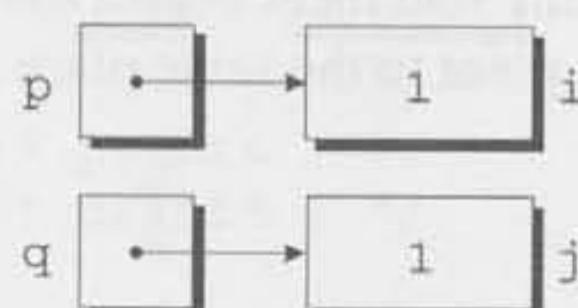
`*q = *p;`

The first statement is a pointer assignment; the second isn't, as the following example shows:

```
p = &i;
q = &j;
i = 1;
```



`*q = *p;`



The assignment `*q = *p` copies the value that p points to (the value of i) into the object that q points to (the variable j).

## 11.4 Pointers as Arguments

So far, we've managed to avoid a rather important question: What are pointers good for? There's no single answer to that question, since pointers have several distinct uses in C. In this section, we'll see how a pointer to a variable can be useful as a function argument. We'll discover other uses for pointers in Section 11.5 and in Chapters 12 and 17.

We saw in Section 9.3 that a variable supplied as an argument in a function call is protected against change, because C passes arguments by value. This property of C can be a nuisance if we want the function to be able to modify the variable. In Section 9.3, we tried—and failed—to write a `decompose` function that could modify two of its arguments.

Pointers offer a solution to this problem: instead of passing a variable `x` as the argument to a function, we'll supply `&x`, a pointer to `x`. We'll declare the corresponding parameter `p` to be a pointer. When the function is called, `p` will have the value `&x`, hence `*p` (the object that `p` points to) will be an alias for `x`. Each appearance of `*p` in the body of the function will be an indirect reference to `x`, allowing the function both to read `x` and to modify it.

To see this technique in action, let's modify the `decompose` function by declaring the parameters `int_part` and `frac_part` to be pointers. The definition of `decompose` will now look like this:

```
void decompose(double x, long *int_part, double *frac_part)
{
    *int_part = (long) x;
    *frac_part = x - *int_part;
}
```

The prototype for `decompose` could be either

```
void decompose(double x, long *int_part, double *frac_part);
```

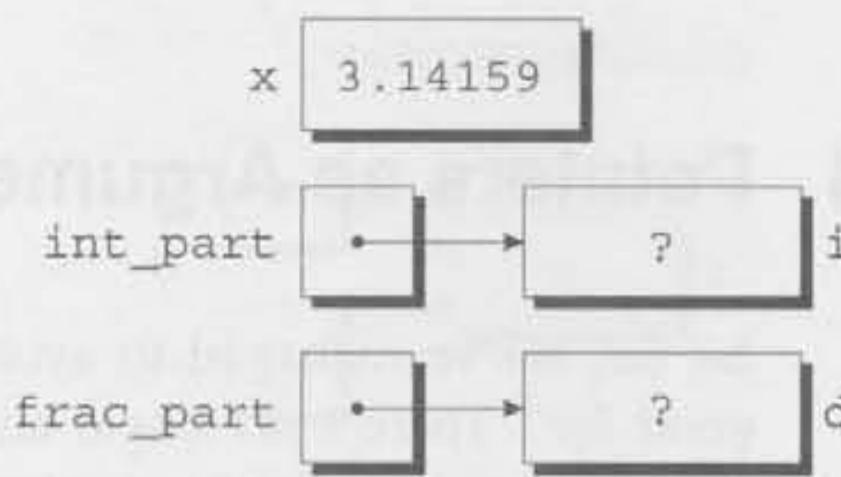
or

```
void decompose(double, long *, double *);
```

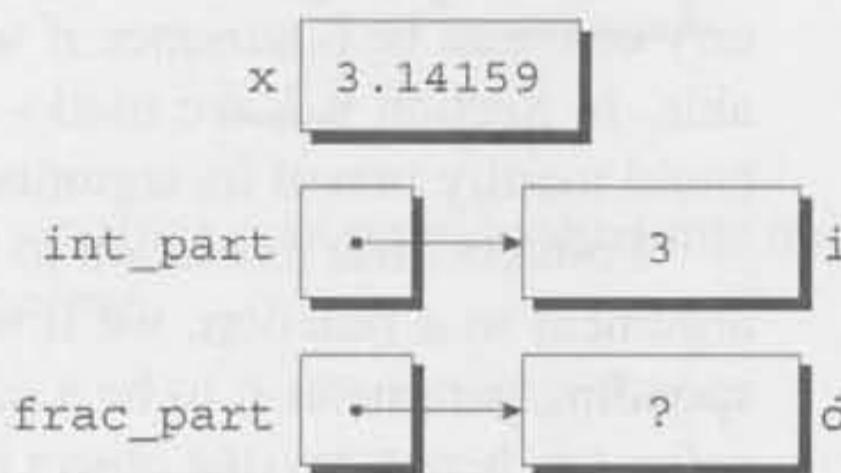
We'll call `decompose` in the following way:

```
decompose(3.14159, &i, &d);
```

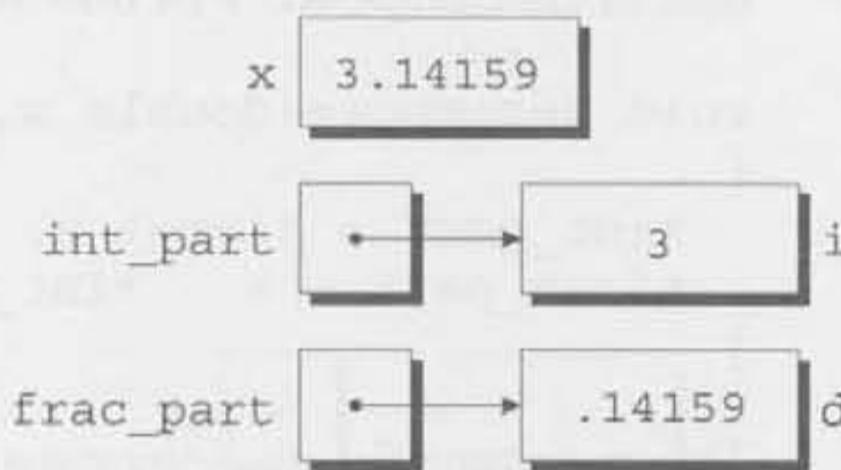
Because of the `&` operator in front of `i` and `d`, the arguments to `decompose` are *pointers* to `i` and `d`, not the *values* of `i` and `d`. When `decompose` is called, the value 3.14159 is copied into `x`, a pointer to `i` is stored in `int_part`, and a pointer to `d` is stored in `frac_part`:



The first assignment in the body of `decompose` converts the value of `x` to type `long` and stores it in the object pointed to by `int_part`. Since `int_part` points to `i`, the assignment puts the value 3 in `i`:



The second assignment fetches the value that `int_part` points to (the value of `i`), which is 3. This value is converted to type `double` and subtracted from `x`, giving `.14159`, which is then stored in the object that `frac_part` points to:



When `decompose` returns, `i` and `d` will have the values 3 and `.14159`, just as we originally wanted.

Using pointers as arguments to functions is actually nothing new; we've been doing it in calls of `scanf` since Chapter 2. Consider the following example:

```
int i;
...
scanf("%d", &i);
```

We must put the `&` operator in front of `i` so that `scanf` is given a *pointer* to `i`; that pointer tells `scanf` where to put the value that it reads. Without the `&`, `scanf` would be supplied with the *value* of `i`.

Although `scanf`'s arguments must be pointers, it's not always true that every argument needs the `&` operator. In the following example, `scanf` is passed a pointer variable:

```

int i, *p;
...
p = &i;
scanf("%d", p);

```

Since `p` contains the address of `i`, `scanf` will read an integer and store it in `i`. Using the `&` operator in the call would be wrong:

```
scanf("%d", &p);    /*** WRONG ***/

```

`scanf` would read an integer and store it in `p` instead of in `i`.



Failing to pass a pointer to a function when one is expected can have disastrous results. Suppose that we call `decompose` without the `&` operator in front of `i` and `d`:

```
decompose(3.14159, i, d);
```

`decompose` is expecting pointers as its second and third arguments, but it's been given the *values* of `i` and `d` instead. `decompose` has no way to tell the difference, so it will use the values of `i` and `d` as though they were pointers. When `decompose` stores values in `*int_part` and `*frac_part`, it will attempt to change unknown memory locations instead of modifying `i` and `d`.

If we've provided a prototype for `decompose` (as we should always do, of course), the compiler will let us know that we're attempting to pass arguments of the wrong type. In the case of `scanf`, however, failing to pass pointers often goes undetected by the compiler, making `scanf` an especially error-prone function.

## PROGRAM Finding the Largest and Smallest Elements in an Array

To illustrate how pointers are passed to functions, let's look at a function named `max_min` that finds the largest and smallest elements in an array. When we call `max_min`, we'll pass it pointers to two variables; `max_min` will then store its answers in these variables. `max_min` has the following prototype:

```
void max_min(int a[], int n, int *max, int *min);
```

A call of `max_min` might have the following appearance:

```
max_min(b, N, &big, &small);
```

`b` is an array of integers; `N` is the number of elements in `b`. `big` and `small` are ordinary integer variables. When `max_min` finds the largest element in `b`, it stores the value in `big` by assigning it to `*max`. (Since `max` points to `big`, an assignment to `*max` will modify the value of `big`.) `max_min` stores the smallest element of `b` in `small` by assigning it to `*min`.

To test `max_min`, we'll write a program that reads 10 numbers into an array, passes the array to `max_min`, and prints the results:

```
Enter 10 numbers: 34 82 49 102 7 94 23 11 50 31
Largest: 102
Smallest: 7
```

Here's the complete program:

```
maxmin.c /* Finds the largest and smallest elements in an array */

#include <stdio.h>

#define N 10

void max_min(int a[], int n, int *max, int *min);

int main(void)
{
    int b[N], i, big, small;

    printf("Enter %d numbers: ", N);
    for (i = 0; i < N; i++)
        scanf("%d", &b[i]);

    max_min(b, N, &big, &small);

    printf("Largest: %d\n", big);
    printf("Smallest: %d\n", small);

    return 0;
}

void max_min(int a[], int n, int *max, int *min)
{
    int i;

    *max = *min = a[0];
    for (i = 1; i < n; i++) {
        if (a[i] > *max)
            *max = a[i];
        else if (a[i] < *min)
            *min = a[i];
    }
}
```

## Using `const` to Protect Arguments

When we call a function and pass it a pointer to a variable, we normally assume that the function will modify the variable (otherwise, why would the function require a pointer?). For example, if we see a statement like

```
f (&x);
```

in a program, we'd probably expect `f` to change the value of `x`. It's possible, though, that `f` merely needs to examine the value of `x`, not change it. The reason for the pointer might be efficiency: passing the value of a variable can waste time and space if the variable requires a large amount of storage. (Section 12.3 covers this point in more detail.)

We can use the word `const` to document that a function won't change an object whose address is passed to the function. `const` goes in the parameter's declaration, just before the specification of its type:

```
void f(const int *p)
{
    *p = 0;    /*** WRONG ***/
}
```

This use of `const` indicates that `p` is a pointer to a “constant integer.” Attempting to modify `*p` is an error that the compiler will detect.

## 11.5 Pointers as Return Values

We can not only pass pointers to functions but also write functions that *return* pointers. Such functions are relatively common; we'll encounter several in Chapter 13.

The following function, when given pointers to two integers, returns a pointer to whichever integer is larger:

```
int *max(int *a, int *b)
{
    if (*a > *b)
        return a;
    else
        return b;
}
```

When we call `max`, we'll pass pointers to two `int` variables and store the result in a pointer variable:

```
int *p, i, j;
...
p = max(&i, &j);
```

During the call of `max`, `*a` is an alias for `i`, while `*b` is an alias for `j`. If `i` has a larger value than `j`, `max` returns the address of `i`; otherwise, it returns the address of `j`. After the call, `p` points to either `i` or `j`.

Although the `max` function returns one of the pointers passed to it as an argument, that's not the only possibility. A function could also return a pointer to an external variable or to a local variable that's been declared `static`.

### Q&A



Never return a pointer to an *automatic local variable*:

```
int *f(void)
{
    int i;
    ...
    return &i;
}
```

The variable *i* doesn't exist once *f* returns, so the pointer to it will be invalid. Some compilers issue a warning such as "*function returns address of local variable*" in this situation.

Pointers can point to array elements, not just ordinary variables. If *a* is an array, then *&a[i]* is a pointer to element *i* of *a*. When a function has an array argument, it's sometimes useful for the function to return a pointer to one of the elements in the array. For example, the following function returns a pointer to the middle element of the array *a*, assuming that *a* has *n* elements:

```
int *find_middle(int a[], int n) {
    return &a[n/2];
}
```

Chapter 12 explores the relationship between pointers and arrays in considerable detail.

## Q & A

**\*Q: Is a pointer always the same as an address? [p. 242]**

A: Usually, but not always. Consider a computer whose main memory is divided into *words* rather than bytes. A word might contain 36 bits, 60 bits, or some other number of bits. If we assume 36-bit words, memory will have the following appearance:

Address	Contents
0	001010011001010011001010011001010011
1	001110101001110101001110101001110101
2	001110011001110011001110011001110011
3	001100001001100001001100001001100001
4	001101110001101110001101110001101110
	⋮
<i>n</i> -1	001000011001000011001000011001000011

When memory is divided into words, each word has an address. An integer usually occupies one word, so a pointer to an integer can just be an address. However, a word can store more than one character. For example, a 36-bit word might store six 6-bit characters:

010011	110101	110011	100001	101110	000011
--------	--------	--------	--------	--------	--------

or four 9-bit characters:

001010011	001110101	001110011	001100001
-----------	-----------	-----------	-----------

For this reason, a pointer to a character may need to be stored in a different form than other pointers. A pointer to a character might consist of an address (the word in which the character is stored) plus a small integer (the position of the character within the word).

On some computers, pointers may be “offsets” rather than complete addresses. For example, CPUs in the Intel x86 family (used in many personal computers) can execute programs in several modes. The oldest of these, which dates back to the 8086 processor of 1978, is called *real mode*. In this mode, addresses are sometimes represented by a single 16-bit number (an *offset*) and sometimes by two 16-bit numbers (a *segment:offset pair*). An offset isn’t a true memory address; the CPU must combine it with a segment value stored in a special register. To support real mode, older C compilers often provide two kinds of pointers: *near pointers* (16-bit offsets) and *far pointers* (32-bit segment:offset pairs). These compilers usually reserve the words *near* and *far* as nonstandard keywords that can be used to declare pointer variables.

**\*Q:** If a pointer can point to *data* in a program, is it possible to have a pointer to *program code*?

A: Yes. We’ll cover pointers to functions in Section 17.7.

**Q:** It seems to me that there’s an inconsistency between the declaration

`int *p = &i;`

and the statement

`p = &i;`

Why isn’t `p` preceded by a `*` symbol in the statement, as it is in the declaration? [p. 244]

A: The source of the confusion is the fact that the `*` symbol can have different meanings in C, depending on the context in which it’s used. In the declaration

`int *p = &i;`

the `*` symbol is *not* the indirection operator. Instead, it helps specify the type of `p`, informing the compiler that `p` is a *pointer* to an `int`. When it appears in a statement,

however, the `*` symbol performs indirection (when used as a unary operator). The statement

```
*p = &i;    /*** WRONG ***/

```

would be wrong, because it assigns the address of `i` to the object that `p` points to, not to `p` itself.

**Q:** Is there some way to print the address of a variable? [p. 244]

**A:** Any pointer, including the address of a variable, can be displayed by calling the `printf` function and using `%p` as the conversion specification. See Section 22.3 for details.

**Q:** The following declaration is confusing:

```
void f(const int *p);
```

**Does this say that `f` can't modify `p`? [p. 251]**

**A:** No. It says that `f` can't change the integer that `p` *points to*; it doesn't prevent `f` from changing `p` itself.

```
void f(const int *p)
{
    int j;

    *p = 0;    /*** WRONG ***/
    p = &j;    /* legal */
}
```

Since arguments are passed by value, assigning `p` a new value—by making it point somewhere else—won't have any effect outside the function.

**\*Q:** When declaring a parameter of a pointer type, is it legal to put the word `const` in front of the parameter's name, as in the following example?

```
void f(int * const p);
```

**A:** Yes, although the effect isn't the same as if `const` precedes `p`'s type. We saw in Section 11.4 that putting `const` *before* `p`'s type protects the object that `p` points to. Putting `const` *after* `p`'s type protects `p` itself:

```
void f(int * const p)
{
    int j;

    *p = 0;    /* legal */
    p = &j;    /*** WRONG ***/
}
```

This feature isn't used very often. Since `p` is merely a copy of another pointer (the argument when the function is called), there's rarely any reason to protect it.

An even greater rarity is the need to protect both `p` *and* the object it points to, which can be done by putting `const` both before and after `p`'s type:

```

void f(const int * const p)
{
    int j;

    *p = 0;      /*** WRONG ***/
    p = &j;      /*** WRONG ***/
}

```

## Exercises

### Section 11.2

1. If *i* is a variable and *p* points to *i*, which of the following expressions are aliases for *i*?
- (a) *\*p*      (c) *\*&p*      (e) *\*i*      (g) *\*&i*  
 (b) *&p*      (d) *&\*p*      (f) *&i*      (h) *&\*i*

### Section 11.3

- W 2. If *i* is an *int* variable and *p* and *q* are pointers to *int*, which of the following assignments are legal?
- (a) *p = i;*      (d) *p = &q;*      (g) *p = \*q;*  
 (b) *\*p = &i;*      (e) *p = \*&q;*      (h) *\*p = q;*  
 (c) *&p = q;*      (f) *p = q;*      (i) *\*p = \*q;*

### Section 11.4

3. The following function supposedly computes the sum and average of the numbers in the array *a*, which has length *n*. *avg* and *sum* point to variables that the function should modify. Unfortunately, the function contains several errors; find and correct them.

```

void avg_sum(double a[], int n, double *avg, double *sum)
{
    int i;

    sum = 0.0;
    for (i = 0; i < n; i++)
        sum += a[i];
    avg = sum / n;
}

```

- W 4. Write the following function:

```
void swap(int *p, int *q);
```

When passed the addresses of two variables, *swap* should exchange the values of the variables:

```
swap(&i, &j); /* exchanges values of i and j */
```

5. Write the following function:

```
void split_time(long total_sec, int *hr, int *min, int *sec);
```

*total\_sec* is a time represented as the number of seconds since midnight. *hr*, *min*, and *sec* are pointers to variables in which the function will store the equivalent time in hours (0–23), minutes (0–59), and seconds (0–59), respectively.

- W 6. Write the following function:

```
void find_two_largest(int a[], int n, int *largest,
                      int *second_largest);
```

When passed an array *a* of length *n*, the function will search *a* for its largest and second-largest elements, storing them in the variables pointed to by *largest* and *second\_largest*, respectively.

7. Write the following function:

```
void split_date(int day_of_year, int year,
                int *month, int *day);
```

*day\_of\_year* is an integer between 1 and 366, specifying a particular day within the year designated by *year*. *month* and *day* point to variables in which the function will store the equivalent month (1–12) and day within that month (1–31).

#### Section 11.5

8. Write the following function:

```
int *find_largest(int a[], int n);
```

When passed an array *a* of length *n*, the function will return a pointer to the array's largest element.

## Programming Projects

1. Modify Programming Project 7 from Chapter 2 so that it includes the following function:

```
void pay_amount(int dollars, int *twenties, int *tens,
                int *fives, int *ones);
```

The function determines the smallest number of \$20, \$10, \$5, and \$1 bills necessary to pay the amount represented by the *dollars* parameter. The *twenties* parameter points to a variable in which the function will store the number of \$20 bills required. The *tens*, *fives*, and *ones* parameters are similar.

2. Modify Programming Project 8 from Chapter 5 so that it includes the following function:

```
void find_closest_flight(int desired_time,
                        int *departure_time,
                        int *arrival_time);
```

This function will find the flight whose departure time is closest to *desired\_time* (expressed in minutes since midnight). It will store the departure and arrival times of this flight (also expressed in minutes since midnight) in the variables pointed to by *departure\_time* and *arrival\_time*, respectively.

3. Modify Programming Project 3 from Chapter 6 so that it includes the following function:

```
void reduce(int numerator, int denominator,
            int *reduced_numerator,
            int *reduced_denominator);
```

*numerator* and *denominator* are the numerator and denominator of a fraction. *reduced\_numerator* and *reduced\_denominator* are pointers to variables in which the function will store the numerator and denominator of the fraction once it has been reduced to lowest terms.

4. Modify the *poker.c* program of Section 10.5 by moving all external variables into *main* and modifying functions so that they communicate by passing arguments. The *analyze\_hand* function needs to change the *straight*, *flush*, *four*, *three*, and *pairs* variables, so it will have to be passed pointers to those variables.

# 12 Pointers and Arrays

*Optimization hinders evolution.*

Chapter 11 introduced pointers and showed how they're used as function arguments and as values returned by functions. This chapter covers another application for pointers. When pointers point to array elements, C allows us to perform arithmetic—addition and subtraction—on the pointers, which leads to an alternative way of processing arrays in which pointers take the place of array subscripts.

The relationship between pointers and arrays in C is a close one, as we'll soon see. We'll exploit this relationship in subsequent chapters, including Chapter 13 (Strings) and Chapter 17 (Advanced Uses of Pointers). Understanding the connection between pointers and arrays is critical for mastering C: it will give you insight into how C was designed and help you understand existing programs. Be aware, however, that one of the primary reasons for using pointers to process arrays—efficiency—is no longer as important as it once was, thanks to improved compilers.

Section 12.1 discusses pointer arithmetic and shows how pointers can be compared using the relational and equality operators. Section 12.2 then demonstrates how we can use pointer arithmetic for processing array elements. Section 12.3 reveals a key fact about arrays—an array name can serve as a pointer to the array's first element—and uses it to show how array arguments really work. Section 12.4 shows how the topics of the first three sections apply to multidimensional arrays. Section 12.5 wraps up the chapter by exploring the relationship between pointers and variable-length arrays, a C99 feature.

## 12.1 Pointer Arithmetic

We saw in Section 11.5 that pointers can point to array elements. For example, suppose that `a` and `p` have been declared as follows:

```
int a[10], *p;
```

We can make p point to a[0] by writing

```
p = &a[0];
```

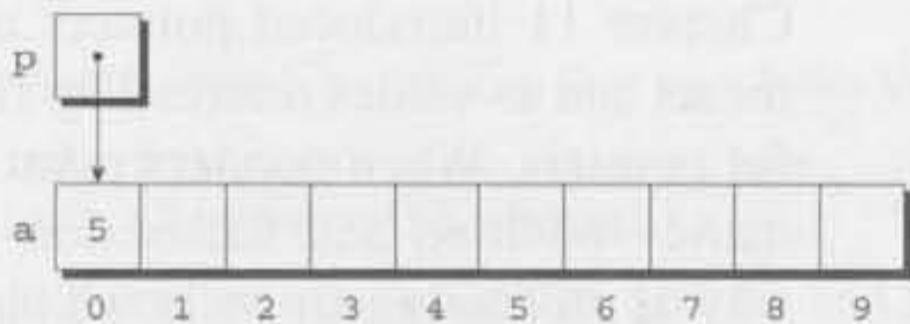
Graphically, here's what we've just done:



We can now access `a[0]` through `p`; for example, we can store the value 5 in `a[0]` by writing

```
*p = 5;
```

Here's our picture now:



Making a pointer `p` point to an element of an array `a` isn't particularly exciting. However, by performing *pointer arithmetic* (or *address arithmetic*) on `p`, we can access the other elements of `a`. C supports three (and only three) forms of pointer arithmetic:

- Adding an integer to a pointer
- Subtracting an integer from a pointer
- Subtracting one pointer from another

Let's take a close look at each of these operations. Our examples assume that the following declarations are in effect:

```
int a[10], *p, *q, i;
```

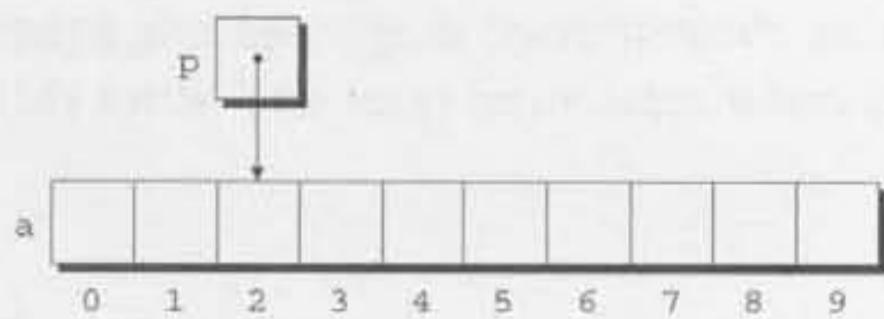
## Adding an Integer to a Pointer

Adding an integer `j` to a pointer `p` yields a pointer to the element `j` places after the one that `p` points to. More precisely, if `p` points to the array element `a[i]`, then `p + j` points to `a[i+j]` (provided, of course, that `a[i+j]` exists).

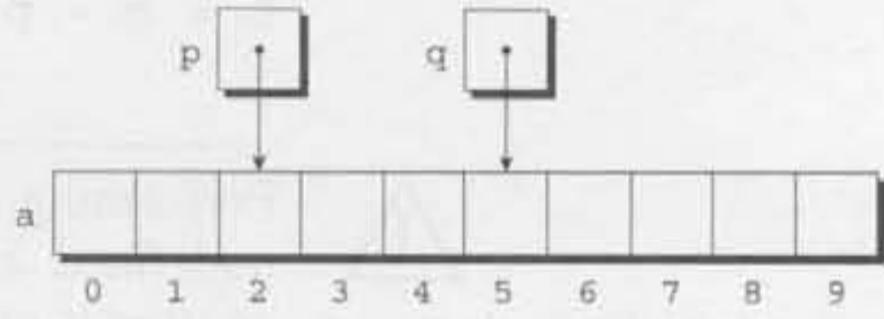
The following example illustrates pointer addition; diagrams show the values of `p` and `q` at various points in the computation.

**Q&A**

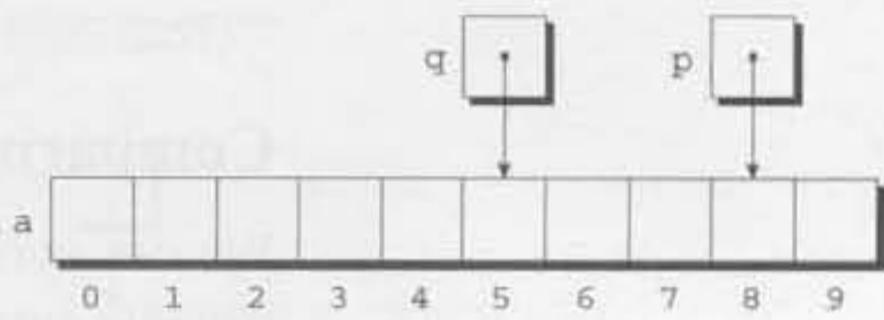
`p = &a[2];`



`q = p + 3;`



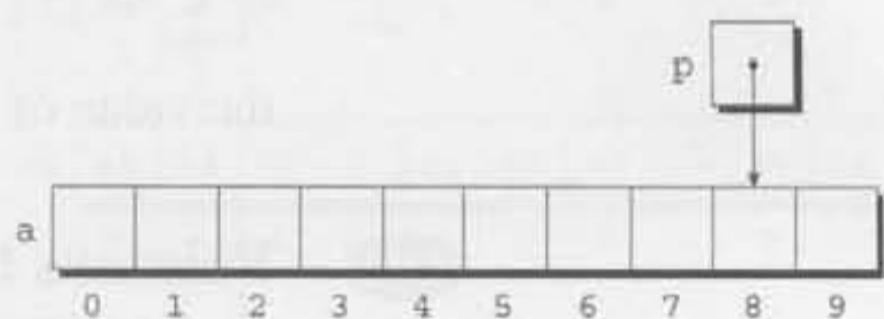
`p += 6;`



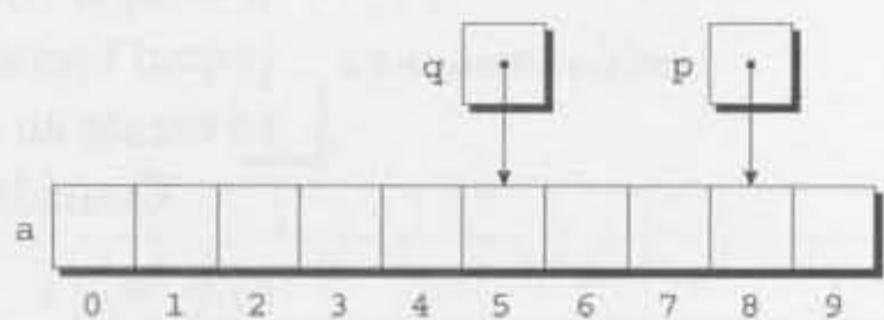
## Subtracting an Integer from a Pointer

If `p` points to the array element `a[i]`, then `p - j` points to `a[i-j]`. For example:

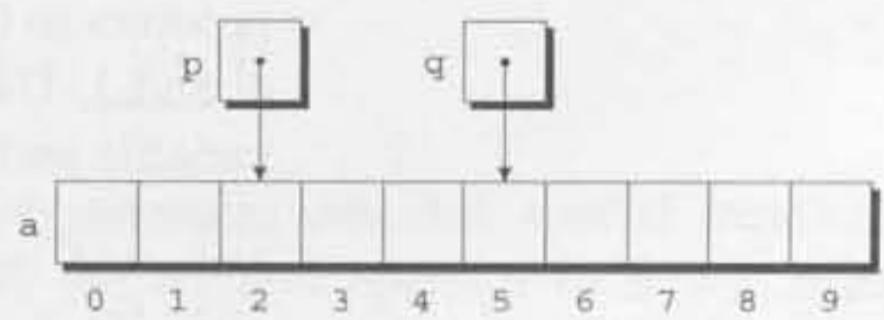
`p = &a[8];`



`q = p - 3;`



`p -= 6;`



## Subtracting One Pointer from Another

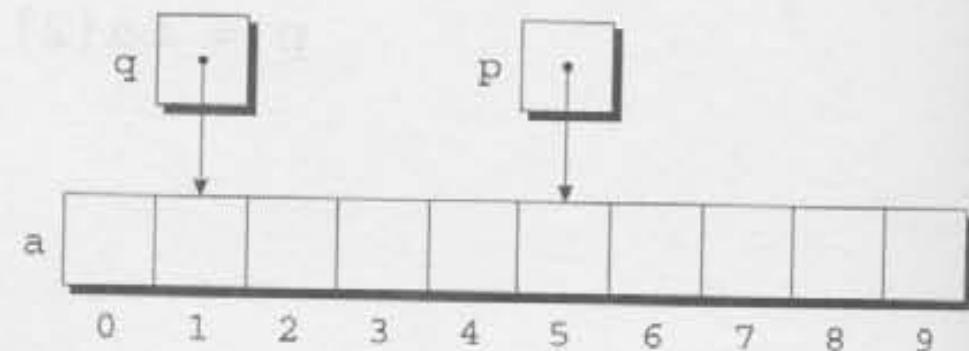
When one pointer is subtracted from another, the result is the distance (measured in array elements) between the pointers. Thus, if `p` points to `a[i]` and `q` points to `a[j]`, then `p - q` is equal to `i - j`. For example:

```

p = &a[5];
q = &a[1];

i = p - q; /* i is 4 */
i = q - p; /* i is -4 */

```



Performing arithmetic on a pointer that doesn't point to an array element causes undefined behavior. Furthermore, the effect of subtracting one pointer from another is undefined unless both point to elements of the *same* array.

## Comparing Pointers

We can compare pointers using the relational operators (`<`, `<=`, `>`, `>=`) and the equality operators (`==` and `!=`). Using the relational operators to compare two pointers is meaningful only when both point to elements of the same array. The outcome of the comparison depends on the relative positions of the two elements in the array. For example, after the assignments

```

p = &a[5];
q = &a[1];

```

the value of `p <= q` is 0 and the value of `p >= q` is 1.

**C99**

## Pointers to Compound Literals

compound literals ▶ 9.3

It's legal for a pointer to point to an element within an array created by a compound literal. A compound literal, you may recall, is a C99 feature that can be used to create an array with no name.

Consider the following example:

```
int *p = (int []){3, 0, 3, 4, 1};
```

`p` points to the first element of a five-element array containing the integers 3, 0, 3, 4, and 1. Using a compound literal saves us the trouble of first declaring an array variable and then making `p` point to the first element of that array:

```
int a[] = {3, 0, 3, 4, 1};
int *p = &a[0];
```

## 12.2 Using Pointers for Array Processing

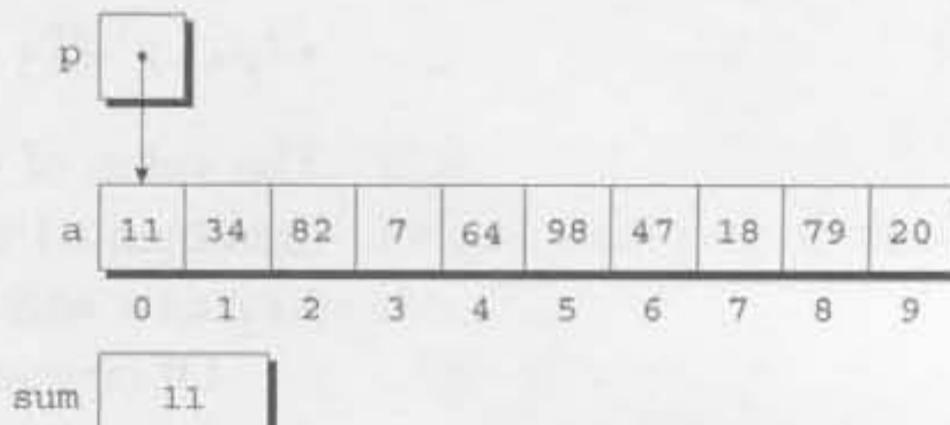
Pointer arithmetic allows us to visit the elements of an array by repeatedly incrementing a pointer variable. The following program fragment, which sums the elements of an array `a`, illustrates the technique. In this example, the pointer variable

`p` initially points to `a[0]`. Each time through the loop, `p` is incremented; as a result, it points to `a[1]`, then `a[2]`, and so forth. The loop terminates when `p` steps past the last element of `a`.

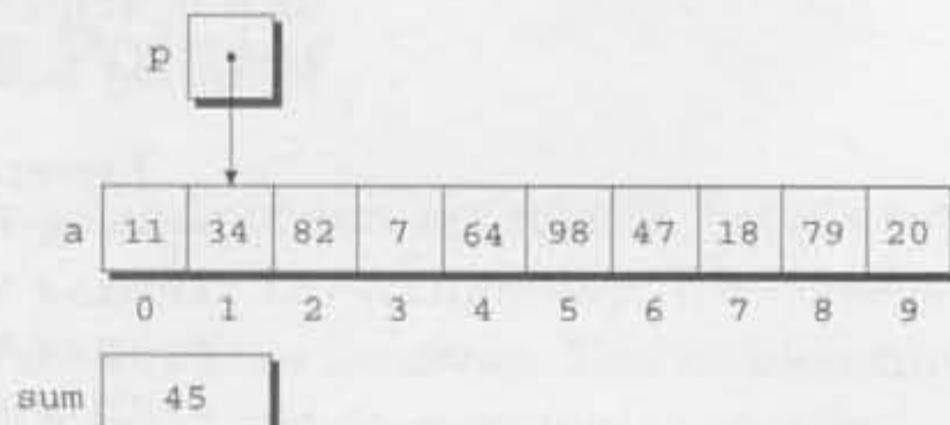
```
#define N 10
...
int a[N], sum, *p;
...
sum = 0;
for (p = &a[0]; p < &a[N]; p++)
    sum += *p;
```

The following figures show the contents of `a`, `sum`, and `p` at the end of the first three loop iterations (before `p` has been incremented).

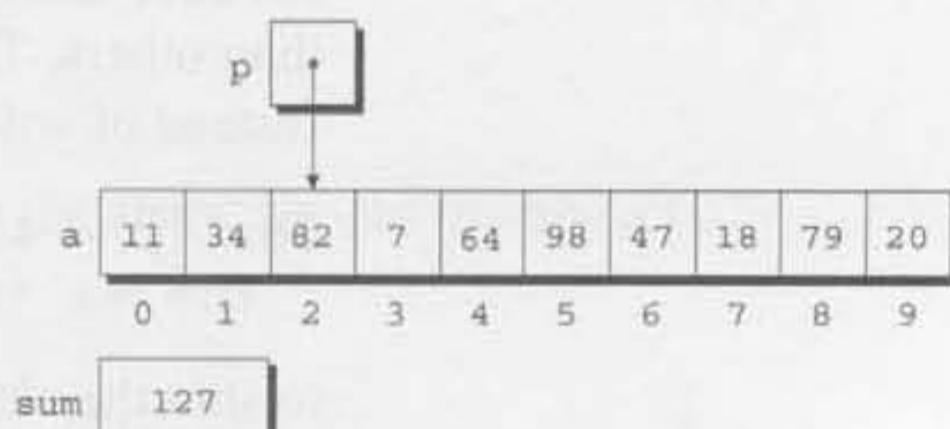
At the end of the first iteration:



At the end of the second iteration:



At the end of the third iteration:



The condition `p < &a[N]` in the `for` statement deserves special mention. Strange as it may seem, it's legal to apply the address operator to `a[N]`, even though this element doesn't exist (`a` is indexed from 0 to  $N - 1$ ). Using `a[N]` in this fashion is perfectly safe, since the loop doesn't attempt to examine its value. The body of the loop will be executed with `p` equal to `&a[0]`, `&a[1]`, ..., `&a[N-1]`, but when `p` is equal to `&a[N]`, the loop terminates.

We could just as easily have written the loop without pointers, of course, using subscripting instead. The argument most often cited in support of pointer arithmetic is that it can save execution time. However, that depends on the implementation—some C compilers actually produce better code for loops that rely on subscripting.

## Combining the \* and ++ Operators

C programmers often combine the \* (indirection) and ++ operators in statements that process array elements. Consider the simple case of storing a value into an array element and then advancing to the next element. Using array subscripting, we might write

```
a[i++] = j;
```

If p is pointing to an array element, the corresponding statement would be

```
*p++ = j;
```

Because the postfix version of ++ takes precedence over \*, the compiler sees this as

```
* (p++) = j;
```

The value of p++ is p. (Since we're using the postfix version of ++, p won't be incremented until after the expression has been evaluated.) Thus, the value of \* (p++) will be \*p—the object to which p is pointing.

Of course, \*p++ isn't the only legal combination of \* and ++. We could write (\*p)++, for example, which returns the value of the object that p points to, and then increments that object (p itself is unchanged). If you find this confusing, the following table may help:

<i>Expression</i>	<i>Meaning</i>
*p++ or * (p++)	Value of expression is *p before increment; increment p later
(*p)++	Value of expression is *p before increment; increment *p later
++p or * (++p)	Increment p first; value of expression is *p after increment
++*p or ++ (*p)	Increment *p first; value of expression is *p after increment

All four combinations appear in programs, although some are far more common than others. The one we'll see most frequently is \*p++, which is handy in loops. Instead of writing

```
for (p = &a[0]; p < &a[N]; p++)
    sum += *p;
```

to sum the elements of the array a, we could write

```
p = &a[0];
while (p < &a[N])
    sum += *p++;
```

The \* and -- operators mix in the same way as \* and ++. For an application that combines \* and --, let's return to the stack example of Section 10.2. The original version of the stack relied on an integer variable named top to keep track of the "top-of-stack" position in the contents array. Let's replace top by a pointer variable that points initially to element 0 of the contents array:

```
int *top_ptr = &contents[0];
```

Here are the new push and pop functions (updating the other stack functions is left as an exercise):

```
void push(int i)
{
    if (is_full())
        stack_overflow();
    else
        *top_ptr++ = i;
}

int pop(void)
{
    if (is_empty())
        stack_underflow();
    else
        return *--top_ptr;
}
```

Note that I've written `*--top_ptr`, not `*top_ptr--`, since I want `pop` to decrement `top_ptr` *before* fetching the value to which it points.

## 12.3 Using an Array Name as a Pointer

Pointer arithmetic is one way in which arrays and pointers are related, but it's not the only connection between the two. Here's another key relationship: *The name of an array can be used as a pointer to the first element in the array*. This relationship simplifies pointer arithmetic and makes both arrays and pointers more versatile.

For example, suppose that `a` is declared as follows:

```
int a[10];
```

Using `a` as a pointer to the first element in the array, we can modify `a[0]`:

```
*a = 7; /* stores 7 in a[0] */
```

We can modify `a[1]` through the pointer `a + 1`:

```
*(a+1) = 12; /* stores 12 in a[1] */
```

In general, `a + i` is the same as `&a[i]` (both represent a pointer to element `i` of `a`) and `*(a+i)` is equivalent to `a[i]` (both represent element `i` itself). In other words, array subscripting can be viewed as a form of pointer arithmetic.

The fact that an array name can serve as a pointer makes it easier to write loops that step through an array. Consider the following loop from Section 12.2:

```
for (p = &a[0]; p < &a[N]; p++)
    sum += *p;
```

To simplify the loop, we can replace `&a[0]` by `a` and `&a[N]` by `a + N`:

**idiom**    for (p = a; p < a + N; p++)  
              sum += \*p;



Although an array name can be used as a pointer, it's not possible to assign it a new value. Attempting to make it point elsewhere is an error:

```
while (*a != 0)
    a++;                                /**** WRONG ****/
```

This is no great loss; we can always copy `a` into a pointer variable, then change the pointer variable:

```
p = a;  
while (*p != 0)  
    p++;
```

## PROGRAM Reversing a Series of Numbers (Revisited)

The `reverse.c` program of Section 8.1 reads 10 numbers, then writes the numbers in reverse order. As the program reads the numbers, it stores them in an array. Once all the numbers are read, the program steps through the array backwards as it prints the numbers.

The original program used subscripting to access elements of the array. Here's a new version in which I've replaced subscripting with pointer arithmetic.

```
reverse3.c /* Reverses a series of numbers (pointer version) */

#include <stdio.h>

#define N 10

int main(void)
{
    int a[N], *p;

    printf("Enter %d numbers: ", N);
    for (p = a; p < a + N; p++)
        scanf("%d", p);

    printf("In reverse order:");
    for (p = a + N - 1; p >= a; p--)
        printf(" %d", *p);
    printf("\n");

    return 0;
}
```

In the original program, an integer variable *i* kept track of the current position within the array. The new version replaces *i* with *p*, a pointer variable. The num-

bers are still stored in an array; we're simply using a different technique to keep track of where we are in the array.

Note that the second argument to `scanf` is `p`, not `&p`. Since `p` points to an array element, it's a satisfactory argument for `scanf`; `&p`, on the other hand, would be a pointer to a pointer to an array element.

## Array Arguments (Revisited)

When passed to a function, an array name is always treated as a pointer. Consider the following function, which returns the largest element in an array of integers:

```
int find_largest(int a[], int n)
{
    int i, max;

    max = a[0];
    for (i = 1; i < n; i++)
        if (a[i] > max)
            max = a[i];
    return max;
}
```

Suppose that we call `find_largest` as follows:

```
largest = find_largest(b, N);
```

This call causes a pointer to the first element of `b` to be assigned to `a`; the array itself isn't copied.

The fact that an array argument is treated as a pointer has some important consequences:

- When an ordinary variable is passed to a function, its value is copied; any changes to the corresponding parameter don't affect the variable. In contrast, an array used as an argument isn't protected against change, since no copy is made of the array itself. For example, the following function (which we first saw in Section 9.3) modifies an array by storing zero into each of its elements:

```
void store_zeros(int a[], int n)
{
    int i;

    for (i = 0; i < n; i++)
        a[i] = 0;
}
```

To indicate that an array parameter won't be changed, we can include the word `const` in its declaration:

```
int find_largest(const int a[], int n)
{
    ...
}
```

If `const` is present, the compiler will check that no assignment to an element of `a` appears in the body of `find_largest`.

- The time required to pass an array to a function doesn't depend on the size of the array. There's no penalty for passing a large array, since no copy of the array is made.
- An array parameter can be declared as a pointer if desired. For example, `find_largest` could be defined as follows:

```
int find_largest(int *a, int n)
{
    ...
}
```

Declaring `a` to be a pointer is equivalent to declaring it to be an array; the compiler treats the declarations as though they were identical.

### Q&A



Although declaring a *parameter* to be an array is the same as declaring it to be a pointer, the same isn't true for a *variable*. The declaration

```
int a[10];
```

causes the compiler to set aside space for 10 integers. In contrast, the declaration

```
int *a;
```

causes the compiler to allocate space for a pointer variable. In the latter case, `a` is not an array; attempting to use it as an array can have disastrous results. For example, the assignment

```
*a = 0;      /*** WRONG ***/
```

will store 0 where `a` is pointing. Since we don't know where `a` is pointing, the effect on the program is undefined.

- A function with an array parameter can be passed an array "slice"—a sequence of consecutive elements. Suppose that we want `find_largest` to locate the largest element in some portion of an array `b`, say elements `b[5], ..., b[14]`. When we call `find_largest`, we'll pass it the address of `b[5]` and the number 10, indicating that we want `find_largest` to examine 10 array elements, starting at `b[5]`:

```
largest = find_largest(&b[5], 10);
```

## Using a Pointer as an Array Name

If we can use an array name as a pointer, will C allow us to subscript a pointer as though it were an array name? By now, you'd probably expect the answer to be yes, and you'd be right. Here's an example:

```
#define N 10
...
int a[N], i, sum = 0, *p = a;
...
for (i = 0; i < N; i++)
    sum += p[i];
```

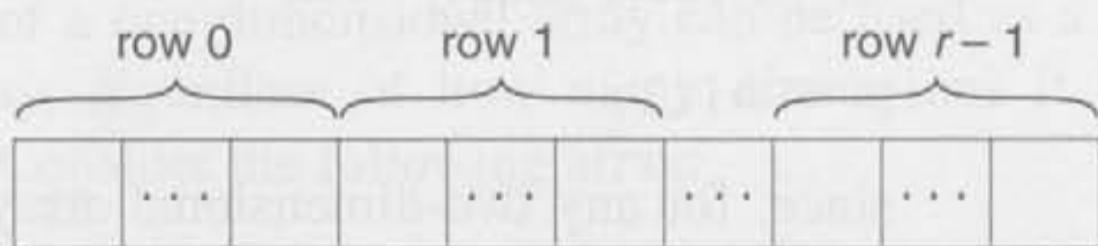
The compiler treats  $p[i]$  as  $*(\text{p} + i)$ , which is a perfectly legal use of pointer arithmetic. Although the ability to subscript a pointer may seem to be little more than a curiosity, we'll see in Section 17.3 that it's actually quite useful.

## 12.4 Pointers and Multidimensional Arrays

Just as pointers can point to elements of one-dimensional arrays, they can also point to elements of multidimensional arrays. In this section, we'll explore common techniques for using pointers to process the elements of multidimensional arrays. For simplicity, I'll stick to two-dimensional arrays, but everything we'll do applies equally to higher-dimensional arrays.

### Processing the Elements of a Multidimensional Array

We saw in Section 8.2 that C stores two-dimensional arrays in row-major order; in other words, the elements of row 0 come first, followed by the elements of row 1, and so forth. An array with  $r$  rows would have the following appearance:



We can take advantage of this layout when working with pointers. If we make a pointer  $p$  point to the first element in a two-dimensional array (the element in row 0, column 0), we can visit every element in the array by incrementing  $p$  repeatedly.

As an example, let's look at the problem of initializing all elements of a two-dimensional array to zero. Suppose that the array has been declared as follows:

```
int a[NUM_ROWS][NUM_COLS];
```

The obvious technique would be to use nested `for` loops:

```
int row, col;
...
for (row = 0; row < NUM_ROWS; row++)
    for (col = 0; col < NUM_COLS; col++)
        a[row][col] = 0;
```

But if we view *a* as a one-dimensional array of integers (which is how it's stored), we can replace the pair of loops by a single loop:

```
int *p;
...
for (p = &a[0][0]; p <= &a[NUM_ROWS-1][NUM_COLS-1]; p++)
    *p = 0;
```

The loop begins with *p* pointing to *a*[0][0]. Successive increments of *p* make it point to *a*[0][1], *a*[0][2], *a*[0][3], and so on. When *p* reaches *a*[0][NUM\_COLS-1] (the last element in row 0), incrementing it again makes *p* point to *a*[1][0], the first element in row 1. The process continues until *p* goes past *a*[NUM\_ROWS-1][NUM\_COLS-1], the last element in the array.

Although treating a two-dimensional array as one-dimensional may seem like cheating, it works with most C compilers. Whether it's a good idea to do so is another matter. Techniques like this one definitely hurt program readability, but—at least with some older compilers—produce a compensating increase in efficiency. With many modern compilers, though, there's often little or no speed advantage.

### Q&A

## Processing the Rows of a Multidimensional Array

What about processing the elements in just one *row* of a two-dimensional array? Again, we have the option of using a pointer variable *p*. To visit the elements of row *i*, we'd initialize *p* to point to element 0 in row *i* in the array *a*:

```
p = &a[i][0];
```

Or we could simply write

```
p = a[i];
```

since, for any two-dimensional array *a*, the expression *a*[*i*] is a pointer to the first element in row *i*. To see why this works, recall the magic formula that relates array subscripting to pointer arithmetic: for any array *a*, the expression *a*[*i*] is equivalent to *\*(a + i)*. Thus, *&a[i][0]* is the same as *&(\*(a[i] + 0))*, which is equivalent to *&a[i]*, which is the same as *a[i]*, since the *&* and *\** operators cancel. We'll use this simplification in the following loop, which clears row *i* of the array *a*:

```
int a[NUM_ROWS][NUM_COLS], *p, i;
...
for (p = a[i]; p < a[i] + NUM_COLS; p++)
    *p = 0;
```

Since *a[i]* is a pointer to row *i* of the array *a*, we can pass *a[i]* to a function that's expecting a one-dimensional array as its argument. In other words, a function that's designed to work with one-dimensional arrays will also work with a row belonging to a two-dimensional array. As a result, functions such as

`find_largest` and `store_zeros` are more versatile than you might expect. Consider `find_largest`, which we originally designed to find the largest element of a one-dimensional array. We can just as easily use `find_largest` to determine the largest element in row  $i$  of the two-dimensional array  $a$ :

```
largest = find_largest(a[i], NUM_COLS);
```

## Processing the Columns of a Multidimensional Array

Processing the elements in a *column* of a two-dimensional array isn't as easy, because arrays are stored by row, not by column. Here's a loop that clears column  $i$  of the array  $a$ :

```
int a[NUM_ROWS][NUM_COLS], (*p)[NUM_COLS], i;
...
for (p = &a[0]; p < &a[NUM_ROWS]; p++)
    (*p)[i] = 0;
```

I've declared  $p$  to be a pointer to an array of length `NUM_COLS` whose elements are integers. The parentheses around  $*p$  in  $(*p)[NUM_COLS]$  are required; without them, the compiler would treat  $p$  as an array of pointers instead of a pointer to an array. The expression  $p++$  advances  $p$  to the beginning of the next row. In the expression  $(*p)[i]$ ,  $*p$  represents an entire row of  $a$ , so  $(*p)[i]$  selects the element in column  $i$  of that row. The parentheses in  $(*p)[i]$  are essential, because the compiler would interpret  $*p[i]$  as  $*(p[i])$ .

## Using the Name of a Multidimensional Array as a Pointer

Just as the name of a one-dimensional array can be used as a pointer, so can the name of *any* array, regardless of how many dimensions it has. Some care is required, though. Consider the following array:

```
int a[NUM_ROWS][NUM_COLS];
```

$a$  is *not* a pointer to  $a[0][0]$ ; instead, it's a pointer to  $a[0]$ . This makes more sense if we look at it from the standpoint of C, which regards  $a$  not as a two-dimensional array but as a one-dimensional array whose elements are one-dimensional arrays. When used as a pointer,  $a$  has type `int (*)[NUM_COLS]` (pointer to an integer array of length `NUM_COLS`).

Knowing that  $a$  points to  $a[0]$  is useful for simplifying loops that process the elements of a two-dimensional array. For example, instead of writing

```
for (p = &a[0]; p < &a[NUM_ROWS]; p++)
    (*p)[i] = 0;
```

to clear column  $i$  of the array  $a$ , we can write

```
for (p = a; p < a + NUM_ROWS; p++)
    (*p)[i] = 0;
```

Another situation in which this knowledge comes in handy is when we want to “trick” a function into thinking that a multidimensional array is really one-dimensional. For example, consider how we might use `find_largest` to find the largest element in `a`. As the first argument to `find_largest`, let’s try passing `a` (the address of the array); as the second, we’ll pass `NUM_ROWS * NUM_COLS` (the total number of elements in `a`):

```
largest = find_largest(a, NUM_ROWS * NUM_COLS); /* WRONG */
```

Unfortunately, the compiler will object to this statement, because the type of `a` is `int (*) [NUM_COLS]` but `find_largest` is expecting an argument of type `int *`. The correct call is

```
largest = find_largest(a[0], NUM_ROWS * NUM_COLS);
```

`a[0]` points to element 0 in row 0, and it has type `int *` (after conversion by the compiler), so the latter call will work correctly.

### Q&A

## 12.5 Pointers and Variable-Length Arrays (C99)

variable-length arrays ➤ 8.3

Pointers are allowed to point to elements of variable-length arrays (VLAs), a feature of C99. An ordinary pointer variable would be used to point to an element of a one-dimensional VLA:

```
void f(int n)
{
    int a[n], *p;
    p = a;
    ...
}
```

When the VLA has more than one dimension, the type of the pointer depends on the length of each dimension except for the first. Let’s look at the two-dimensional case:

```
void f(int m, int n)
{
    int a[m][n], (*p)[n];
    p = a;
    ...
}
```

Since the type of `p` depends on `n`, which isn’t constant, `p` is said to have a **variably modified type**. Note that the validity of an assignment such as `p = a` can’t always be determined by the compiler. For example, the following code will compile but is correct only if `m` and `n` are equal:

```
int a[m][n], (*p)[m];
p = a;
```

If  $m \neq n$ , any subsequent use of  $p$  will cause undefined behavior.

Variably modified types are subject to certain restrictions, just as variable-length arrays are. The most important restriction is that the declaration of a variably modified type must be inside the body of a function or in a function prototype.

Pointer arithmetic works with VLAs just as it does for ordinary arrays. Returning to the example of Section 12.4 that clears a single column of a two-dimensional array  $a$ , let's declare  $a$  as a VLA this time:

```
int a[m][n];
```

A pointer capable of pointing to a row of  $a$  would be declared as follows:

```
int (*p)[n];
```

The loop that clears column  $i$  is almost identical to the one we used in Section 12.4:

```
for (p = a; p < a + m; p++)
    (*p)[i] = 0;
```

## Q & A

**Q: I don't understand pointer arithmetic. If a pointer is an address, does that mean that an expression like  $p + j$  adds  $j$  to the address stored in  $p$ ? [p. 258]**

**A:** No. Integers used in pointer arithmetic are scaled depending on the type of the pointer. If  $p$  is of type `int *`, for example, then  $p + j$  typically adds  $4 \times j$  to  $p$ , assuming that `int` values are stored using 4 bytes. But if  $p$  has type `double *`, then  $p + j$  will probably add  $8 \times j$  to  $p$ , since `double` values are usually 8 bytes long.

**Q: When writing a loop to process an array, is it better to use array subscripting or pointer arithmetic? [p. 261]**

**A:** There's no easy answer to this question, since it depends on the machine you're using and the compiler itself. In the early days of C on the PDP-11, pointer arithmetic yielded a faster program. On today's machines, using today's compilers, array subscripting is often just as good, and sometimes even better. The bottom line: Learn both ways and then use whichever is more natural for the kind of program you're writing.

**\*Q: I read somewhere that  $i[a]$  is the same as  $a[i]$ . Is this true?**

**A:** Yes, it is, oddly enough. The compiler treats  $i[a]$  as  $* (i + a)$ , which is the same as  $* (a + i)$ . (Pointer addition, like ordinary addition, is commutative.) But  $* (a + i)$  is equivalent to  $a[i]$ . Q.E.D. But please don't use  $i[a]$  in programs unless you're planning to enter the next Obfuscated C contest.

**Q:** Why is `*a` the same as `a []` in a parameter declaration? [p. 266]

**A:** Both indicate that the argument is expected to be a pointer. The same operations on `a` are possible in both cases (pointer arithmetic and array subscripting, in particular). And, in both cases, `a` itself can be assigned a new value within the function. (Although C allows us to use the name of an array *variable* only as a “constant pointer,” there’s no such restriction on the name of an array *parameter*.)

**Q:** Is it better style to declare an array parameter as `*a` or `a []`?

**A:** That’s a tough one. From one standpoint, `a []` is the obvious choice, since `*a` is ambiguous (does the function want an array of objects or a pointer to a single object?). On the other hand, many programmers argue that declaring the parameter as `*a` is more accurate, since it reminds us that only a pointer is passed, not a copy of the array. Others switch between `*a` and `a []`, depending on whether the function uses pointer arithmetic or subscripting to access the elements of the array. (That’s the approach I’ll use.) In practice, `*a` is more common than `a []`, so you’d better get used to it. For what it’s worth, Dennis Ritchie now refers to the `a []` notation as “a living fossil” that “serves as much to confuse the learner as to alert the reader.”

**Q:** We’ve seen that arrays and pointers are closely related in C. Would it be accurate to say that they’re interchangeable?

**A:** No. It’s true that array *parameters* are interchangeable with pointer parameters, but array *variables* aren’t the same as pointer variables. Technically, the name of an array isn’t a pointer; rather, the C compiler *converts* it to a pointer when necessary. To see this difference more clearly, consider what happens when we apply the `sizeof` operator to an array `a`. The value of `sizeof(a)` is the total number of bytes in the array—the size of each element multiplied by the number of elements. But if `p` is a pointer variable, `sizeof(p)` is the number of bytes required to store a pointer value.

**Q:** You said that treating a two-dimensional array as one-dimensional works with “most” C compilers. Doesn’t it work with all compilers? [p. 268]

**A:** No. Some modern “bounds-checking” compilers track not only the type of a pointer, but—when it points to an array—also the length of the array. For example, suppose that `p` is assigned a pointer to `a[0][0]`. Technically, `p` points to the first element of `a[0]`, a one-dimensional array. If we increment `p` repeatedly in an effort to visit all the elements of `a`, we’ll go out of bounds once `p` goes past the last element of `a[0]`. A compiler that performs bounds-checking may insert code to check that `p` is used only to access elements in the array pointed to by `a[0]`; an attempt to increment `p` past the end of this array would be detected as an error.

**Q:** If `a` is a two-dimensional array, why can we pass `a[0]`—but not `a` itself—to `find_largest`? Don’t both `a` and `a[0]` point to the same place (the beginning of the array)? [p. 270]

**A:** They do, as a matter of fact—both point to element `a[0][0]`. The problem is that

a has the wrong type. When used as an argument, it's a pointer to an array, but `find_largest` is expecting a pointer to an integer. However, `a[0]` has type `int *`, so it's an acceptable argument for `find_largest`. This concern about types is actually good; if C weren't so picky, we could make all kinds of horrible pointer mistakes without the compiler noticing.

## Exercises

### Section 12.1

- Suppose that the following declarations are in effect:

```
int a[] = {5, 15, 34, 54, 14, 2, 52, 72};
int *p = &a[1], *q = &a[5];
```

- (a) What is the value of `* (p+3)`?
- (b) What is the value of `* (q-3)`?
- (c) What is the value of `q - p`?
- (d) Is the condition `p < q` true or false?
- (e) Is the condition `*p < *q` true or false?

- W 2. Suppose that `high`, `low`, and `middle` are all pointer variables of the same type, and that `low` and `high` point to elements of an array. Why is the following statement illegal, and how could it be fixed?

```
middle = (low + high) / 2;
```

### Section 12.2

- What will be the contents of the `a` array after the following statements are executed?

```
#define N 10
int a[N] = {1, 2, 3, 4, 5, 6, 7, 8, 9, 10};
int *p = &a[0], *q = &a[N-1], temp;
while (p < q) {
    temp = *p;
    *p++ = *q;
    *q-- = temp;
}
```

- W 4. Rewrite the `make_empty`, `is_empty`, and `is_full` functions of Section 10.2 to use the pointer variable `top_ptr` instead of the integer variable `top`.

### Section 12.3

- Suppose that `a` is a one-dimensional array and `p` is a pointer variable. Assuming that the assignment `p = a` has just been performed, which of the following expressions are illegal because of mismatched types? Of the remaining expressions, which are true (have a nonzero value)?

- (a) `p == a[0]`
- (b) `p == &a[0]`
- (c) `*p == a[0]`
- (d) `p[0] == a[0]`

- W 6. Rewrite the following function to use pointer arithmetic instead of array subscripting. (In other words, eliminate the variable `i` and all uses of the `[]` operator.) Make as few changes as possible.

```
int sum_array(const int a[], int n)
{
    int i, sum;
    sum = 0;
    for (i = 0; i < n; i++)
        sum += a[i];
    return sum;
}
```

7. Write the following function:

```
bool search(const int a[], int n, int key);
```

*a* is an array to be searched, *n* is the number of elements in the array, and *key* is the search key. *search* should return *true* if *key* matches some element of *a*, and *false* if it doesn't. Use pointer arithmetic—not subscripting—to visit array elements.

8. Rewrite the following function to use pointer arithmetic instead of array subscripting. (In other words, eliminate the variable *i* and all uses of the [] operator.) Make as few changes as possible.

```
void store_zeros(int a[], int n)
{
    int i;
    for (i = 0; i < n; i++)
        a[i] = 0;
}
```

9. Write the following function:

```
double inner_product(const double *a, const double *b,
                     int n);
```

*a* and *b* both point to arrays of length *n*. The function should return *a*[0] \* *b*[0] + *a*[1] \* *b*[1] + ... + *a*[*n*-1] \* *b*[*n*-1]. Use pointer arithmetic—not subscripting—to visit array elements.

10. Modify the *find\_middle* function of Section 11.5 so that it uses pointer arithmetic to calculate the return value.

11. Modify the *find\_largest* function so that it uses pointer arithmetic—not subscripting—to visit array elements.

12. Write the following function:

```
void find_two_largest(const int *a, int n, int *largest,
                      int *second_largest);
```

*a* points to an array of length *n*. The function searches the array for its largest and second-largest elements, storing them in the variables pointed to by *largest* and *second\_largest*, respectively. Use pointer arithmetic—not subscripting—to visit array elements.

#### Section 12.4

13. Section 8.2 had a program fragment in which two nested *for* loops initialized the array *ident* for use as an identity matrix. Rewrite this code, using a single pointer to step through the array one element at a time. *Hint:* Since we won't be using *row* and *col* index variables, it won't be easy to tell where to store 1. Instead, we can use the fact that the first element of the array should be 1, the next *N* elements should be 0, the next element should

be 1, and so forth. Use a variable to keep track of how many consecutive 0s have been stored; when the count reaches N, it's time to store 1.

14. Assume that the following array contains a week's worth of hourly temperature readings, with each row containing the readings for one day:

```
int temperatures[7][24];
```

Write a statement that uses the `search` function (see Exercise 7) to search the entire `temperatures` array for the value 32.

- W 15. Write a loop that prints all temperature readings stored in row `i` of the `temperatures` array (see Exercise 14). Use a pointer to visit each element of the row.

16. Write a loop that prints the highest temperature in the `temperatures` array (see Exercise 14) for each day of the week. The loop body should call the `find_largest` function, passing it one row of the array at a time.

17. Rewrite the following function to use pointer arithmetic instead of array subscripting. (In other words, eliminate the variables `i` and `j` and all uses of the `[]` operator.) Use a single loop instead of nested loops.

```
int sum_two_dimensional_array(const int a[] [LEN], int n)
{
    int i, j, sum = 0;

    for (i = 0; i < n; i++)
        for (j = 0; j < LEN; j++)
            sum += a[i][j];

    return sum;
}
```

18. Write the `evaluate_position` function described in Exercise 13 of Chapter 9. Use pointer arithmetic—not subscripting—to visit array elements. Use a single loop instead of nested loops.

## Programming Projects

- W 1. (a) Write a program that reads a message, then prints the reversal of the message:

Enter a message: Don't get mad, get even.  
Reversal is: .neve teg ,dam teg t'noD

*Hint:* Read the message one character at a time (using `getchar`) and store the characters in an array. Stop reading when the array is full or the character read is '`\n`'.

(b) Revise the program to use a pointer instead of an integer to keep track of the current position in the array.

2. (a) Write a program that reads a message, then checks whether it's a palindrome (the letters in the message are the same from left to right as from right to left):

Enter a message: He lived as a devil, eh?  
Palindrome

Enter a message: Madam, I am Adam.  
Not a palindrome

Ignore all characters that aren't letters. Use integer variables to keep track of positions in the array.

(b) Revise the program to use pointers instead of integers to keep track of positions in the array.

- W 3. Simplify Programming Project 1(b) by taking advantage of the fact that an array name can be used as a pointer.
- 4. Simplify Programming Project 2(b) by taking advantage of the fact that an array name can be used as a pointer.
- 5. Modify Programming Project 14 from Chapter 8 so that it uses a pointer instead of an integer to keep track of the current position in the array that contains the sentence.
- 6. Modify the `qsort.c` program of Section 9.6 so that `low`, `high`, and `middle` are pointers to array elements rather than integers. The `split` function will need to return a pointer, not an integer.
- 7. Modify the `maxmin.c` program of Section 11.4 so that the `max_min` function uses a pointer instead of an integer to keep track of the current position in the array.

# 13 Strings

*It's difficult to extract sense from strings, but they're the only communication coin we can count on.*

Although we've used `char` variables and arrays of `char` values in previous chapters, we still lack any convenient way to process a series of characters (a *string*, in C terminology). We'll remedy that defect in this chapter, which covers both string *constants* (or *literals*, as they're called in the C standard) and string *variables*, which can change during the execution of a program.

Section 13.1 explains the rules that govern string literals, including the rules for embedding escape sequences in string literals and for breaking long string literals. Section 13.2 then shows how to declare string variables, which are simply arrays of characters in which a special character—the null character—marks the end of a string. Section 13.3 describes ways to read and write strings. Section 13.4 shows how to write functions that process strings, and Section 13.5 covers some of the string-handling functions in the C library. Section 13.6 presents idioms that are often used when working with strings. Finally, Section 13.7 describes how to set up arrays whose elements are pointers to strings of different lengths. This section also explains how C uses such an array to supply command-line information to programs.

## 13.1 String Literals

A *string literal* is a sequence of characters enclosed within double quotes:

"When you come to a fork in the road, take it."

We first encountered string literals in Chapter 2; they often appear as format strings in calls of `printf` and `scanf`.

## Escape Sequences in String Literals

escape sequences ➤ 7.3

String literals may contain the same escape sequences as character constants. We've used character escapes in `printf` and `scanf` format strings for some time. For example, we've seen that each `\n` character in the string

```
"Candy\nIs dandy\nBut liquor\nIs quicker.\n --Ogden Nash\n"
```

causes the cursor to advance to the next line:

```
Candy
Is dandy
But liquor
Is quicker.
--Ogden Nash
```

Although octal and hexadecimal escapes are also legal in string literals, they're not as common as character escapes.



**Q&A**

Be careful when using octal and hexadecimal escape sequences in string literals. An octal escape ends after three digits or with the first non-octal character. For example, the string "`\1234`" contains two characters (`\123` and `4`), and the string "`\189`" contains three characters (`\1`, `8`, and `9`). A hexadecimal escape, on the other hand, isn't limited to three digits; it doesn't end until the first non-hex character. Consider what happens if a string contains the escape `\xfc`, which represents the character *ü* in the Latin1 character set, a common extension of ASCII. The string "`Z\xfcrich`" ("Zürich") has six characters (`Z`, `\xfc`, `r`, `ü`, `i`, and `ch`), but the string "`\xfcber`" (a failed attempt at "über") has only two (`\xfc` and `ber`). Most compilers will object to the latter string, since hex escapes are usually limited to the range `\x00`–`\xFF`.

## Continuing a String Literal

If we find that a string literal is too long to fit conveniently on a single line, C allows us to continue it on the next line, provided that we end the first line with a backslash character (`\`). No other characters may follow `\` on the same line, other than the (invisible) new-line character at the end:

```
printf("When you come to a fork in the road, take it. \
--Yogi Berra");
```

In general, the `\` character can be used to join two or more lines of a program into a single line (a process that the C standard refers to as "splicing"). We'll see more examples of splicing in Section 14.3.

The `\` technique has one drawback: the string must continue at the beginning of the next line, thereby wrecking the program's indented structure. There's a better way to deal with long string literals, thanks to the following rule: when two or more string literals are adjacent (separated only by white space), the compiler will

join them into a single string. This rule allows us to split a string literal over two or more lines:

```
printf("When you come to a fork in the road, take it.\n"
      "--Yogi Berra");
```

## How String Literals Are Stored

We've used string literals often in calls of `printf` and `scanf`. But when we call `printf` and supply a string literal as an argument, what are we actually passing? To answer this question, we need to know how string literals are stored.

In essence, C treats string literals as character arrays. When a C compiler encounters a string literal of length  $n$  in a program, it sets aside  $n + 1$  bytes of memory for the string. This area of memory will contain the characters in the string, plus one extra character—the **null character**—to mark the end of the string. The null character is a byte whose bits are all zero, so it's represented by the `\0` escape sequence.



Don't confuse the null character ('`\0`') with the zero character ('`0`'). The null character has the code 0; the zero character has a different code (48 in ASCII).

For example, the string literal "abc" is stored as an array of four characters (a, b, c, and `\0`):

a	b	c	<code>\0</code>
---	---	---	-----------------

String literals may be empty; the string "" is stored as a single null character:

<code>\0</code>
-----------------

Since a string literal is stored as an array, the compiler treats it as a pointer of type `char *`. Both `printf` and `scanf`, for example, expect a value of type `char *` as their first argument. Consider the following example:

```
printf("abc");
```

When `printf` is called, it's passed the address of "abc" (a pointer to where the letter a is stored in memory).

## Operations on String Literals

In general, we can use a string literal wherever C allows a `char *` pointer. For example, a string literal can appear on the right side of an assignment:

```
char *p;
p = "abc";
```

This assignment doesn't copy the characters in "abc"; it merely makes p point to the first character of the string.

C allows pointers to be subscripted, so we can subscript string literals:

```
char ch;
ch = "abc"[1];
```

The new value of ch will be the letter b. The other possible subscripts are 0 (which would select the letter a), 2 (the letter c), and 3 (the null character). This property of string literals isn't used that much, but occasionally it's handy. Consider the following function, which converts a number between 0 and 15 into a character that represents the equivalent hex digit:

```
char digit_to_hex_char(int digit)
{
    return "0123456789ABCDEF"[digit];
```



Attempting to modify a string literal causes undefined behavior:

```
char *p = "abc";
*p = 'd';    /*** WRONG ***/
```

#### Q&A

A program that tries to change a string literal may crash or behave erratically.

## String Literals versus Character Constants

A string literal containing a single character isn't the same as a character constant. The string literal "a" is represented by a *pointer* to a memory location that contains the character a (followed by a null character). The character constant 'a' is represented by an *integer* (the numerical code for the character).



Don't ever use a character when a string is required (or vice versa). The call

```
printf("\n");
```

is legal, because printf expects a pointer as its first argument. The following call isn't legal, however:

```
printf('\n');    /*** WRONG ***/
```

## 13.2 String Variables

Some programming languages provide a special `string` type for declaring string variables. C takes a different tack: any one-dimensional array of characters can be used to store a string, with the understanding that the string is terminated by a null character. This approach is simple, but has significant difficulties. It's sometimes hard to tell whether an array of characters is being used as a string. If we write our own string-handling functions, we've got to be careful that they deal properly with the null character. Also, there's no faster way to determine the length of a string than a character-by-character search for the null character.

Let's say that we need a variable capable of storing a string of up to 80 characters. Since the string will need a null character at the end, we'll declare the variable to be an array of 81 characters:

```
idiom #define STR_LEN 80
...
char str[STR_LEN+1];
```

We defined `STR_LEN` to be 80 rather than 81, thus emphasizing the fact that `str` can store strings of no more than 80 characters, and then added 1 to `STR_LEN` in the declaration of `str`. This a common practice among C programmers.



When declaring an array of characters that will be used to hold a string, always make the array one character longer than the string, because of the C convention that every string is terminated by a null character. Failing to leave room for the null character may cause unpredictable results when the program is executed, since functions in the C library assume that strings are null-terminated.

Declaring a character array to have length `STR_LEN + 1` doesn't mean that it will always contain a string of `STR_LEN` characters. The length of a string depends on the position of the terminating null character, not on the length of the array in which the string is stored. An array of `STR_LEN + 1` characters can hold strings of various lengths, ranging from the empty string to strings of length `STR_LEN`.

### Initializing a String Variable

A string variable can be initialized at the same time it's declared:

```
char date1[8] = "June 14";
```

The compiler will put the characters from "June 14" in the date1 array, then add a null character so that date1 can be used as a string. Here's what date1 will look like:

date1	J	u	n	e		1	4	\0
-------	---	---	---	---	--	---	---	----

Although "June 14" appears to be a string literal, it's not. Instead, C views it as an abbreviation for an array initializer. In fact, we could have written

```
char date1[8] = {'J', 'u', 'n', 'e', ' ', '1', '4', '\0'};
```

I think you'll agree that the original version is easier to read.

What if the initializer is too short to fill the string variable? In that case, the compiler adds extra null characters. Thus, after the declaration

```
char date2[9] = "June 14";
```

date2 will have the following appearance:

date2	J	u	n	e		1	4	\0	\0
-------	---	---	---	---	--	---	---	----	----

array initializers ▶ 8.1

This behavior is consistent with C's treatment of array initializers in general. When an array initializer is shorter than the array itself, the remaining elements are initialized to zero. By initializing the leftover elements of a character array to \0, the compiler is following the same rule.

What if the initializer is longer than the string variable? That's illegal for strings, just as it's illegal for other arrays. However, C does allow the initializer (not counting the null character) to have exactly the same length as the variable:

```
char date3[7] = "June 14";
```

There's no room for the null character, so the compiler makes no attempt to store one:

date3	J	u	n	e		1	4
-------	---	---	---	---	--	---	---



If you're planning to initialize a character array to contain a string, be sure that the length of the array is longer than the length of the initializer. Otherwise, the compiler will quietly omit the null character, making the array unusable as a string.

The declaration of a string variable may omit its length, in which case the compiler computes it:

```
char date4[] = "June 14";
```

The compiler sets aside eight characters for `date4`, enough to store the characters in "June 14" plus a null character. (The fact that the length of `date4` isn't specified doesn't mean that the array's length can be changed later. Once the program is compiled, the length of `date4` is fixed at eight.) Omitting the length of a string variable is especially useful if the initializer is long, since computing the length by hand is error-prone.

## Character Arrays versus Character Pointers

Let's compare the declaration

```
char date[] = "June 14";
```

which declares `date` to be an *array*, with the similar-looking

```
char *date = "June 14";
```

which declares `date` to be a *pointer*. Thanks to the close relationship between arrays and pointers, we can use either version of `date` as a string. In particular, any function expecting to be passed a character array or character pointer will accept either version of `date` as an argument.

However, we must be careful not to make the mistake of thinking that the two versions of `date` are interchangeable. There are significant differences between the two:

- In the array version, the characters stored in `date` can be modified, like the elements of any array. In the pointer version, `date` points to a string literal, and we saw in Section 13.1 that string literals shouldn't be modified.
- In the array version, `date` is an array name. In the pointer version, `date` is a variable that can be made to point to other strings during program execution.

If we need a string that can be modified, it's our responsibility to set up an array of characters in which to store the string; declaring a pointer variable isn't enough. The declaration

```
char *p;
```

causes the compiler to set aside enough memory for a pointer variable; unfortunately, it doesn't allocate space for a string. (And how could it? We haven't indicated how long the string would be.) Before we can use `p` as a string, it must point to an array of characters. One possibility is to make `p` point to a string variable:

```
char str[STR_LEN+1], *p;  
p = str;
```

`p` now points to the first character of `str`, so we can use `p` as a string. Another possibility is to make `p` point to a dynamically allocated string.



Using an uninitialized pointer variable as a string is a serious error. Consider the following example, which attempts to build the string "abc":

```
char *p;

p[0] = 'a';    /* WRONG */
p[1] = 'b';    /* WRONG */
p[2] = 'c';    /* WRONG */
p[3] = '\0';   /* WRONG */
```

Since `p` hasn't been initialized, we don't know where it's pointing. Using the pointer to write the characters `a`, `b`, `c`, and `\0` into memory causes undefined behavior.

### 13.3 Reading and Writing Strings

Writing a string is easy using either the `printf` or `puts` functions. Reading a string is a bit harder, primarily because of the possibility that the input string may be longer than the string variable into which it's being stored. To read a string in a single step, we can use either `scanf` or `gets`. As an alternative, we can read strings one character at a time.

#### Writing Strings Using `printf` and `puts`

The `%s` conversion specification allows `printf` to write a string. Consider the following example:

```
char str[] = "Are we having fun yet?";

printf("%s\n", str);
```

The output will be

Are we having fun yet?

`printf` writes the characters in a string one by one until it encounters a null character. (If the null character is missing, `printf` continues past the end of the string until—eventually—it finds a null character somewhere in memory.)

To print just part of a string, we can use the conversion specification `%.ps`, where `p` is the number of characters to be displayed. The statement

```
printf("%.6s\n", str);
```

will print

Are we

A string, like a number, can be printed within a field. The `%ms` conversion will display a string in a field of size  $m$ . (A string with more than  $m$  characters will be printed in full, not truncated.) If the string has fewer than  $m$  characters, it will be right-justified within the field. To force left justification instead, we can put a minus sign in front of  $m$ . The  $m$  and  $p$  values can be used in combination: a conversion specification of the form `%m.ps` causes the first  $p$  characters of a string to be displayed in a field of size  $m$ .

`printf` isn't the only function that can write strings. The C library also provides `puts`, which is used in the following way:

```
puts(str);
```

`puts` has only one argument (the string to be printed). After writing the string, `puts` always writes an additional new-line character, thus advancing to the beginning of the next output line.

## Reading Strings Using `scanf` and `gets`

The `%s` conversion specification allows `scanf` to read a string into a character array:

```
scanf("%s", str);
```

There's no need to put the `&` operator in front of `str` in the call of `scanf`; like any array name, `str` is treated as a pointer when passed to a function.

When `scanf` is called, it skips white space, then reads characters and stores them in `str` until it encounters a white-space character. `scanf` always stores a null character at the end of the string.

A string read using `scanf` will never contain white space. Consequently, `scanf` won't usually read a full line of input; a new-line character will cause `scanf` to stop reading, but so will a space or tab character. To read an entire line of input at a time, we can use `gets`. Like `scanf`, the `gets` function reads input characters into an array, then stores a null character. In other respects, however, `gets` is somewhat different from `scanf`:

- `gets` doesn't skip white space before starting to read the string (`scanf` does).
- `gets` reads until it finds a new-line character (`scanf` stops at any white-space character). Incidentally, `gets` discards the new-line character instead of storing it in the array; the null character takes its place.

To see the difference between `scanf` and `gets`, consider the following program fragment:

```
char sentence[SENT_LEN+1];

printf("Enter a sentence:\n");
scanf("%s", sentence);
```

Suppose that after the prompt

Enter a sentence:

the user enters the line

To C, or not to C: that is the question.

`scanf` will store the string "To" in `sentence`. The next call of `scanf` will resume reading the line at the space after the word To.

Now suppose that we replace `scanf` by `gets`:

```
gets(sentence);
```

When the user enters the same input as before, `gets` will store the string

" To C, or not to C: that is the question."

in `sentence`.



As they read characters into an array, `scanf` and `gets` have no way to detect when it's full. Consequently, they may store characters past the end of the array, causing undefined behavior. `scanf` can be made safer by using the conversion specification `%ns` instead of `%s`, where `n` is an integer indicating the maximum number of characters to be stored. `gets`, unfortunately, is inherently unsafe; `fgets` is a much better alternative.

fgets function ▶ 22.5

## Reading Strings Character by Character

Since both `scanf` and `gets` are risky and insufficiently flexible for many applications, C programmers often write their own input functions. By reading strings one character at a time, these functions provide a greater degree of control than the standard input functions.

If we decide to design our own input function, we'll need to consider the following issues:

- Should the function skip white space before beginning to store the string?
- What character causes the function to stop reading: a new-line character, any white-space character, or some other character? Is this character stored in the string or discarded?
- What should the function do if the input string is too long to store: discard the extra characters or leave them for the next input operation?

Suppose we need a function that doesn't skip white-space characters, stops reading at the first new-line character (which isn't stored in the string), and discards extra characters. The function might have the following prototype:

```
int read_line(char str[], int n);
```

`str` represents the array into which we'll store the input, and `n` is the maximum number of characters to be read. If the input line contains more than `n` characters, `read_line` will discard the additional characters. We'll have `read_line` return the number of characters it actually stores in `str` (a number anywhere from 0 to `n`). We may not always need `read_line`'s return value, but it doesn't hurt to have it available.

getchar function ▶ 7.3

### Q&A

`read_line` consists primarily of a loop that calls `getchar` to read a character and then stores the character in `str`, provided that there's room left. The loop terminates when the new-line character is read. (Strictly speaking, we should also have the loop terminate if `getchar` should fail to read a character, but we'll ignore that complication for now.) Here's the complete definition of `read_line`:

```
int read_line(char str[], int n)
{
    int ch, i = 0;

    while ((ch = getchar()) != '\n')
        if (i < n)
            str[i++] = ch;
    str[i] = '\0';           /* terminates string */
    return i;                /* number of characters stored */
}
```

Note that `ch` has `int` type rather than `char` type, because `getchar` returns the character that it reads as an `int` value.

Before returning, `read_line` puts a null character at the end of the string. Standard functions such as `scanf` and `gets` automatically put a null character at the end of an input string; if we're writing our own input function, however, we must take on that responsibility.

## 13.4 Accessing the Characters in a String

Since strings are stored as arrays, we can use subscripting to access the characters in a string. To process every character in a string `s`, for example, we can set up a loop that increments a counter `i` and selects characters via the expression `s[i]`.

Suppose that we need a function that counts the number of spaces in a string. Using array subscripting, we might write the function in the following way:

```
int count_spaces(const char s[])
{
    int count = 0, i;

    for (i = 0; s[i] != '\0'; i++)
        if (s[i] == ' ')
            count++;
    return count;
}
```

I've included `const` in the declaration of `s` to indicate that `count_spaces` doesn't change the array that `s` represents. If `s` were not a string, the function would need a second argument specifying the length of the array. Since `s` is a string, however, `count_spaces` can determine where it ends by testing for the null character.

Many C programmers wouldn't write `count_spaces` as we have. Instead, they'd use a pointer to keep track of the current position within the string. As we saw in Section 12.2, this technique is always available for processing arrays, but it proves to be especially convenient for working with strings.

Let's rewrite the `count_spaces` function using pointer arithmetic instead of array subscripting. We'll eliminate the variable `i` and use `s` itself to keep track of our position in the string. By incrementing `s` repeatedly, `count_spaces` can step through each character in the string. Here's our new version of the function:

```
int count_spaces(const char *s)
{
    int count = 0;

    for (; *s != '\0'; s++)
        if (*s == ' ')
            count++;
    return count;
}
```

Note that `const` doesn't prevent `count_spaces` from modifying `s`; it's there to prevent the function from modifying what `s` points to. And since `s` is a copy of the pointer that's passed to `count_spaces`, incrementing `s` doesn't affect the original pointer.

The `count_spaces` example raises some questions about how to write string functions:

- ***Is it better to use array operations or pointer operations to access the characters in a string?*** We're free to use whichever is more convenient; we can even mix the two. In the second version of `count_spaces`, treating `s` as a pointer simplifies the function slightly by removing the need for the variable `i`. Traditionally, C programmers lean toward using pointer operations for processing strings.
- ***Should a string parameter be declared as an array or as a pointer?*** The two versions of `count_spaces` illustrate the options: the first version declares `s` to be an array; the second declares `s` to be a pointer. Actually, there's no difference between the two declarations—recall from Section 12.3 that the compiler treats an array parameter as though it had been declared as a pointer.
- ***Does the form of the parameter (`s[]` or `*s`) affect what can be supplied as an argument?*** No. When `count_spaces` is called, the argument could be an array name, a pointer variable, or a string literal—`count_spaces` can't tell the difference.

## 13.5 Using the C String Library

Some programming languages provide operators that can copy strings, compare strings, concatenate strings, select substrings, and the like. C's operators, in contrast, are essentially useless for working with strings. Strings are treated as arrays in C, so they're restricted in the same ways as arrays—in particular, they can't be copied or compared using operators.



Direct attempts to copy or compare strings will fail. For example, suppose that `str1` and `str2` have been declared as follows:

```
char str1[10], str2[10];
```

Copying a string into a character array using the `=` operator is not possible:

```
str1 = "abc";    /*** WRONG ***/
str2 = str1;    /*** WRONG ***/
```

We saw in Section 12.3 that using an array name as the left operand of `=` is illegal. *Initializing* a character array using `=` is legal, though:

```
char str1[10] = "abc";
```

In the context of a declaration, `=` is not the assignment operator.

Attempting to compare strings using a relational or equality operator is legal but won't produce the desired result:

```
if (str1 == str2) ...    /*** WRONG ***/
```

This statement compares `str1` and `str2` as *pointers*; it doesn't compare the contents of the two arrays. Since `str1` and `str2` have different addresses, the expression `str1 == str2` must have the value 0.

`<string.h>` header ▶ 23.6

Fortunately, all is not lost: the C library provides a rich set of functions for performing operations on strings. Prototypes for these functions reside in the `<string.h>` header, so programs that need string operations should contain the following line:

```
#include <string.h>
```

Most of the functions declared in `<string.h>` require at least one string as an argument. String parameters are declared to have type `char *`, allowing the argument to be a character array, a variable of type `char *`, or a string literal—all are suitable as strings. Watch out for string parameters that aren't declared `const`, however. Such a parameter may be modified when the function is called, so the corresponding argument shouldn't be a string literal.

There are many functions in `<string.h>`; I'll cover a few of the most basic. In subsequent examples, assume that `str1` and `str2` are character arrays used as strings.

## The `strcpy` (String Copy) Function

The `strcpy` function has the following prototype in `<string.h>`:

```
char *strcpy(char *s1, const char *s2);
```

`strcpy` copies the string `s2` into the string `s1`. (To be precise, we should say “`strcpy` copies the string pointed to by `s2` into the array pointed to by `s1`.”) That is, `strcpy` copies characters from `s2` to `s1` up to (and including) the first null character in `s2`. `strcpy` returns `s1` (a pointer to the destination string). The string pointed to by `s2` isn't modified, so it's declared `const`.

The existence of `strcpy` compensates for the fact that we can't use the assignment operator to copy strings. For example, suppose that we want to store the string "abcd" in `str2`. We can't use the assignment

```
str2 = "abcd";           /* *** WRONG *** /
```

because `str2` is an array name and can't appear on the left side of an assignment. Instead, we can call `strcpy`:

```
strcpy(str2, "abcd");    /* str2 now contains "abcd" */
```

Similarly, we can't assign `str2` to `str1` directly, but we can call `strcpy`:

```
strcpy(str1, str2);      /* str1 now contains "abcd" */
```

Most of the time, we'll discard the value that `strcpy` returns. On occasion, though, it can be useful to call `strcpy` as part of a larger expression in order to use its return value. For example, we could chain together a series of `strcpy` calls:

```
strcpy(str1, strcpy(str2, "abcd"));
/* both str1 and str2 now contain "abcd" */
```



In the call `strcpy(str1, strcpy(str2, "abcd"))`, `strcpy` has no way to check that the string pointed to by `str2` will actually fit in the array pointed to by `str1`. Suppose that `str1` points to an array of length  $n$ . If the string that `str2` points to has no more than  $n - 1$  characters, then the copy will succeed. But if `str2` points to a longer string, undefined behavior occurs. (Since `strcpy` always copies up to the first null character, it will continue copying past the end of the array that `str1` points to.)

Calling the `strncpy` function is a safer, albeit slower, way to copy a string. `strncpy` is similar to `strcpy` but has a third argument that limits the number of characters that will be copied. To copy `str2` into `str1`, we could use the following call of `strncpy`:

```
strncpy(str1, str2, sizeof(str1));
```

As long as `str1` is large enough to hold the string stored in `str2` (including the null character), the copy will be done correctly. `strncpy` itself isn't without danger, though. For one thing, it will leave the string in `str1` without a terminating null character if the length of the string stored in `str2` is greater than or equal to the size of the `str1` array. Here's a safer way to use `strncpy`:

```
strncpy(str1, str2, sizeof(str1) - 1);
str1[sizeof(str1)-1] = '\0';
```

The second statement guarantees that `str1` is always null-terminated, even if `strncpy` fails to copy a null character from `str2`.

## The `strlen` (String Length) Function

The `strlen` function has the following prototype:

```
size_t strlen(const char *s);
```

`size_t` type ▶ 7.6 `size_t`, which is defined in the C library, is a `typedef` name that represents one of C's unsigned integer types. Unless we're dealing with extremely long strings, this technicality need not concern us—we can simply treat the return value of `strlen` as an integer.

`strlen` returns the length of a string `s`: the number of characters in `s` up to, but not including, the first null character. Here are a few examples:

```
int len;

len = strlen("abc"); /* len is now 3 */
len = strlen(""); /* len is now 0 */
strcpy(str1, "abc");
len = strlen(str1); /* len is now 3 */
```

The last example illustrates an important point. When given an array as its argument, `strlen` doesn't measure the length of the array itself; instead, it returns the length of the string stored in the array.

## The `strcat` (String Concatenation) Function

The `strcat` function has the following prototype:

```
char *strcat(char *s1, const char *s2);
```

`strcat` appends the contents of the string `s2` to the end of the string `s1`; it returns `s1` (a pointer to the resulting string).

Here are some examples of `strcat` in action:

```
strcpy(str1, "abc");
strcat(str1, "def"); /* str1 now contains "abcdef" */
```

```
strcpy(str1, "abc");
strcpy(str2, "def");
strcat(str1, str2); /* str1 now contains "abcdef" */
```

As with `strcpy`, the value returned by `strcat` is normally discarded. The following example shows how the return value might be used:

```
strcpy(str1, "abc");
strcpy(str2, "def");
strcat(str1, strcat(str2, "ghi"));
/* str1 now contains "abcdefghi"; str2 contains "defghi" */
```



The effect of the call `strcat(str1, str2)` is undefined if the array pointed to by `str1` isn't long enough to accommodate the additional characters from `str2`. Consider the following example:

```
char str1[6] = "abc";
strcat(str1, "def"); /* *** WRONG ***/
```

`strcat` will attempt to add the characters d, e, f, and \0 to the end of the string already stored in `str1`. Unfortunately, `str1` is limited to six characters, causing `strcat` to write past the end of the array.

strncat function ▶ 23.6

The `strncat` function is a safer but slower version of `strcat`. Like `strncpy`, it has a third argument that limits the number of characters it will copy. Here's what a call might look like:

```
strncat(str1, str2, sizeof(str1) - strlen(str1) - 1);
```

`strncat` will terminate `str1` with a null character, which isn't included in the third argument (the number of characters to be copied). In the example, the third argument calculates the amount of space remaining in `str1` (given by the expression `sizeof(str1) - strlen(str1)`) and then subtracts 1 to ensure that there will be room for the null character.

## The `strcmp` (String Comparison) Function

The `strcmp` function has the following prototype:

```
int strcmp(const char *s1, const char *s2);
```

`strcmp` compares the strings `s1` and `s2`, returning a value less than, equal to, or greater than 0, depending on whether `s1` is less than, equal to, or greater than `s2`. For example, to see if `str1` is less than `str2`, we'd write

```
if (strcmp(str1, str2) < 0) /* is str1 < str2? */
...
```

**Q&A**

To test whether `str1` is less than or equal to `str2`, we'd write

```
if (strcmp(str1, str2) <= 0) /* is str1 <= str2? */  
...
```

By choosing the proper relational operator (`<`, `<=`, `>`, `>=`) or equality operator (`==`, `!=`), we can test any possible relationship between `str1` and `str2`.

`strcmp` compares strings based on their lexicographic ordering, which resembles the way words are arranged in a dictionary. More precisely, `strcmp` considers `s1` to be less than `s2` if either one of the following conditions is satisfied:

- The first  $i$  characters of `s1` and `s2` match, but the  $(i+1)$ st character of `s1` is less than the  $(i+1)$ st character of `s2`. For example, "abc" is less than "bcd", and "abd" is less than "abe".
- All characters of `s1` match `s2`, but `s1` is shorter than `s2`. For example, "abc" is less than "abcd".

As it compares characters from two strings, `strcmp` looks at the numerical codes that represent the characters. Some knowledge of the underlying character set is helpful in order to predict what `strcmp` will do. For example, here are a few important properties of the ASCII character set:

- The characters in each of the sequences A–Z, a–z, and 0–9 have consecutive codes.
- All upper-case letters are less than all lower-case letters. (In ASCII, codes between 65 and 90 represent upper-case letters; codes between 97 and 122 represent lower-case letters.)
- Digits are less than letters. (Codes between 48 and 57 represent digits.)
- Spaces are less than all printing characters. (The space character has the value 32 in ASCII.)

## PROGRAM Printing a One-Month Reminder List

To illustrate the use of the C string library, we'll now develop a program that prints a one-month list of daily reminders. The user will enter a series of reminders, with each prefixed by a day of the month. When the user enters 0 instead of a valid day, the program will print a list of all reminders entered, sorted by day. Here's what a session with the program will look like:

```
Enter day and reminder: 24 Susan's birthday
Enter day and reminder: 5 6:00 - Dinner with Marge and Russ
Enter day and reminder: 26 Movie - "Chinatown"
Enter day and reminder: 7 10:30 - Dental appointment
Enter day and reminder: 12 Movie - "Dazed and Confused"
Enter day and reminder: 5 Saturday class
Enter day and reminder: 12 Saturday class
Enter day and reminder: 0
```

```

Day Reminder
 5 Saturday class
 5 6:00 - Dinner with Marge and Russ
 7 10:30 - Dental appointment
12 Saturday class
12 Movie - "Dazed and Confused"
24 Susan's birthday
26 Movie - "Chinatown"

```

The overall strategy isn't very complicated: we'll have the program read a series of day-and-reminder combinations, storing them in order (sorted by day), and then display them. To read the days, we'll use `scanf`; to read the reminders, we'll use the `read_line` function of Section 13.3.

We'll store the strings in a two-dimensional array of characters, with each row of the array containing one string. After the program reads a day and its associated reminder, it will search the array to determine where the day belongs, using `strcmp` to do comparisons. It will then use `strcpy` to move all strings *below* that point down one position. Finally, the program will copy the day into the array and call `strcat` to append the reminder to the day. (The day and the reminder have been kept separate up to this point.).

Of course, there are always a few minor complications. For example, we want the days to be right-justified in a two-character field, so that their ones digits will line up. There are many ways to handle the problem. I've chosen to have the program use `scanf` to read the day into an integer variable, then call `sprintf` to convert the day back into string form. `sprintf` is a library function that's similar to `printf`, except that it writes output into a string. The call

```
sprintf(day_str, "%2d", day);
```

writes the value of `day` into `day_str`. Since `sprintf` automatically adds a null character when it's through writing, `day_str` will contain a properly null-terminated string.

Another complication is making sure that the user doesn't enter more than two digits. We'll use the following call of `scanf` for this purpose:

```
scanf("%2d", &day);
```

The number 2 between % and d tells `scanf` to stop reading after two digits, even if the input has more digits.

With those details out of the way, here's the program:

```

remind.c /* Prints a one-month reminder list */

#include <stdio.h>
#include <string.h>

#define MAX_REMIND 50    /* maximum number of reminders */
#define MSG_LEN 60        /* max length of reminder message */

```

```

int read_line(char str[], int n);

int main(void)
{
    char reminders[MAX_REMIND] [MSG_LEN+3];
    char day_str[3], msg_str[MSG_LEN+1];
    int day, i, j, num_remind = 0;

    for (;;) {
        if (num_remind == MAX_REMIND) {
            printf("-- No space left --\n");
            break;
        }

        printf("Enter day and reminder: ");
        scanf("%2d", &day);
        if (day == 0)
            break;
        sprintf(day_str, "%2d", day);
        read_line(msg_str, MSG_LEN);

        for (i = 0; i < num_remind; i++)
            if (strcmp(day_str, reminders[i]) < 0)
                break;
        for (j = num_remind; j > i; j--)
            strcpy(reminders[j], reminders[j-1]);

        strcpy(reminders[i], day_str);
        strcat(reminders[i], msg_str);

        num_remind++;
    }

    printf("\nDay Reminder\n");
    for (i = 0; i < num_remind; i++)
        printf(" %s\n", reminders[i]);
}

int read_line(char str[], int n)
{
    int ch, i = 0;

    while ((ch = getchar()) != '\n')
        if (i < n)
            str[i++] = ch;
        str[i] = '\0';
    return i;
}

```

Although *remind.c* is useful for demonstrating the *strcpy*, *strcat*, and *strcmp* functions, it lacks something as a practical reminder program. There are

obviously a number of improvements needed, ranging from minor tweaks to major enhancements (such as saving the reminders in a file when the program terminates). We'll discuss several improvements in the programming projects at the end of this chapter and in later chapters.

## 13.6 String Idioms

Functions that manipulate strings are a particularly rich source of idioms. In this section, we'll explore some of the most famous idioms by using them to write the `strlen` and `strcat` functions. You'll never have to write these functions, of course, since they're part of the standard library, but you may have to write functions that are similar.

The concise style I'll use in this section is popular with many C programmers. You should master this style even if you don't plan to use it in your own programs, since you're likely to encounter it in code written by others.

One last note before we get started. If you want to try out any of the versions of `strlen` and `strcat` in this section, be sure to alter the name of the function (changing `strlen` to `my_strlen`, for example). As Section 21.1 explains, we're not allowed to write a function that has the same name as a standard library function, even when we don't include the header to which the function belongs. In fact, all names that begin with `str` and a lower-case letter are reserved (to allow functions to be added to the `<string.h>` header in future versions of the C standard).

### Searching for the End of a String

Many string operations require searching for the end of a string. The `strlen` function is a prime example. The following version of `strlen` searches its string argument to find the end, using a variable to keep track of the string's length:

```
size_t strlen(const char *s)
{
    size_t n;

    for (n = 0; *s != '\0'; s++)
        n++;
    return n;
}
```

As the pointer `s` moves across the string from left to right, the variable `n` keeps track of how many characters have been seen so far. When `s` finally points to a null character, `n` contains the length of the string.

Let's see if we can condense the function. First, we'll move the initialization of `n` to its declaration:

```
size_t strlen(const char *s)
{
    size_t n = 0;

    for (; *s != '\0'; s++)
        n++;
    return n;
}
```

Next, we notice that the condition `*s != '\0'` is the same as `*s != 0`, because the integer value of the null character is 0. But testing `*s != 0` is the same as testing `*s`; both are true if `*s` isn't equal to 0. These observations lead to our next version of `strlen`:

```
size_t strlen(const char *s)
{
    size_t n = 0;

    for (; *s; s++)
        n++;
    return n;
}
```

But, as we saw in Section 12.2, it's possible to increment `s` and test `*s` in the same expression:

```
size_t strlen(const char *s)
{
    size_t n = 0;

    for (; *s++;)
        n++;
    return n;
}
```

Replacing the `for` statement with a `while` statement, we arrive at the following version of `strlen`:

```
size_t strlen(const char *s)
{
    size_t n = 0;

    while (*s++)
        n++;
    return n;
}
```

Although we've condensed `strlen` quite a bit, it's likely that we haven't increased its speed. Here's a version that *does* run faster, at least with some compilers:

```
size_t strlen(const char *s)
{
    const char *p = s;
```

```

while (*s)
    s++;
return s - p;
}

```

This version of `strlen` computes the length of the string by locating the position of the null character, then subtracting from it the position of the first character in the string. The improvement in speed comes from not having to increment `n` inside the `while` loop. Note the appearance of the word `const` in the declaration of `p`, by the way; without it, the compiler would notice that assigning `s` to `p` places the string that `s` points to at risk.

The statement

**idiom** `while (*s)`  
    `s++;`

and the related

**idiom** `while (*s++)`  
    `;`

are idioms meaning “search for the null character at the end of a string.” The first version leaves `s` pointing to the null character. The second version is more concise, but leaves `s` pointing just past the null character.

## Copying a String

Copying a string is another common operation. To introduce C’s “string copy” idiom, we’ll develop two versions of the `strcat` function. Let’s start with a straightforward but somewhat lengthy version:

```

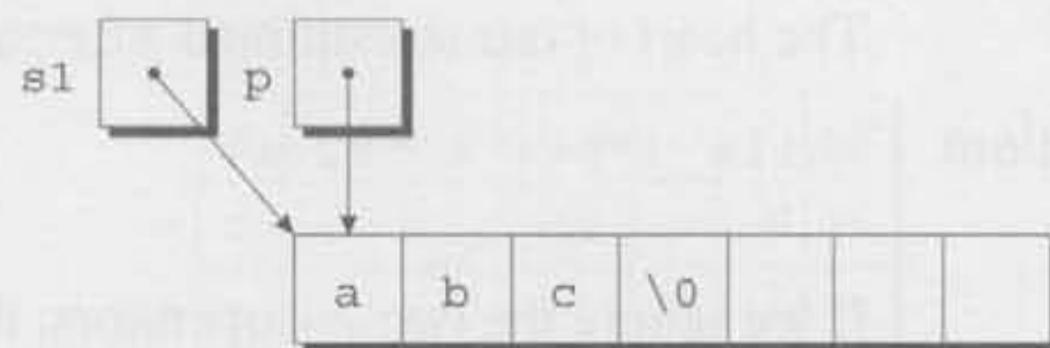
char *strcat(char *s1, const char *s2)
{
    char *p = s1;

    while (*p != '\0')
        p++;
    while (*s2 != '\0') {
        *p = *s2;
        p++;
        s2++;
    }
    *p = '\0';
    return s1;
}

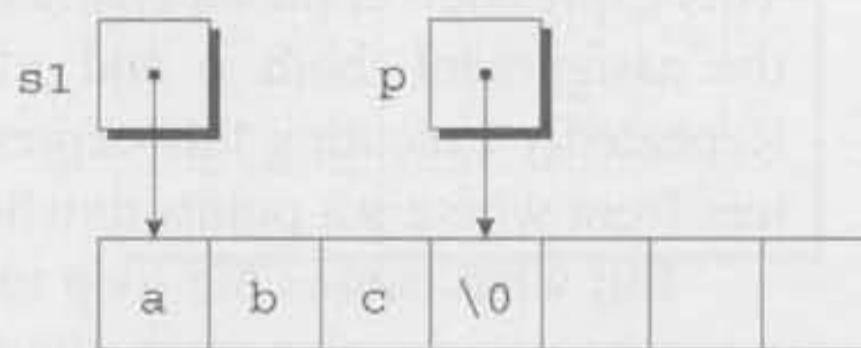
```

This version of `strcat` uses a two-step algorithm: (1) Locate the null character at the end of the string `s1` and make `p` point to it. (2) Copy characters one by one from `s2` to where `p` is pointing.

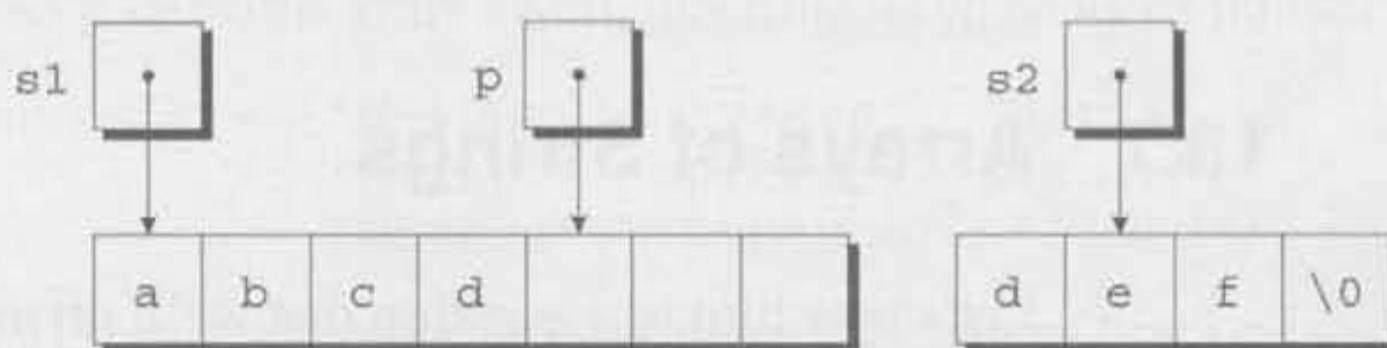
The first `while` statement in the function implements step (1). `p` is set to point to the first character in the `s1` string. Assuming that `s1` points to the string "abc", we have the following picture:



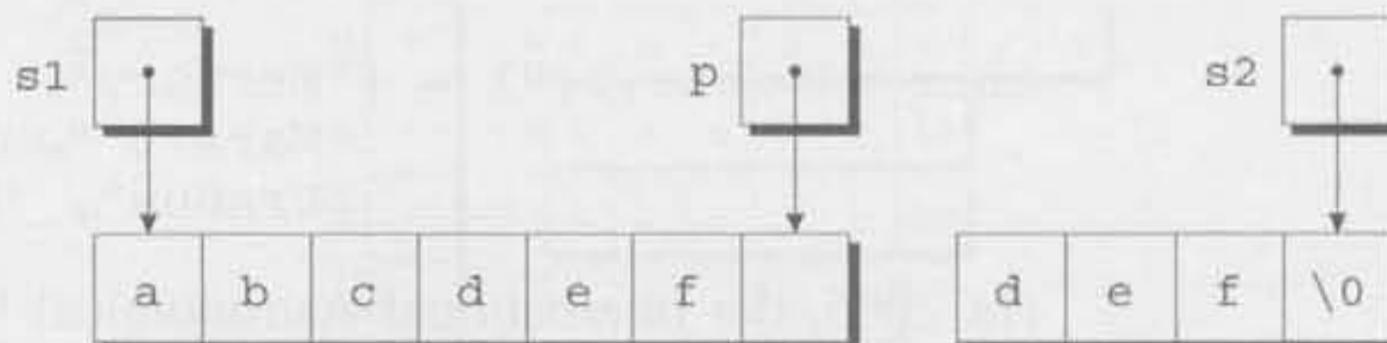
`p` is then incremented as long as it doesn't point to a null character. When the loop terminates, `p` must be pointing to the null character:



The second `while` statement implements step (2). The loop body copies one character from where `s2` points to where `p` points, then increments both `p` and `s2`. If `s2` originally points to the string "def", here's what the strings will look like after the first loop iteration:



The loop terminates when `s2` points to the null character:



After putting a null character where `p` is pointing, `strcat` returns.

By a process similar to the one we used for `strlen`, we can condense the definition of `strcat`, arriving at the following version:

```
char *strcat(char *s1, const char *s2)
{
    char *p = s1;

    while (*p)
        p++;
    while (*p++ = *s2++)
        ;
    return s1;
}
```

The heart of our streamlined `strcat` function is the “string copy” idiom:

```
idiom while (*p++ = *s2++)  
;
```

If we ignore the two `++` operators, the expression inside the parentheses simplifies to an ordinary assignment:

```
*p = *s2
```

This expression copies a character from where `s2` points to where `p` points. After the assignment, both `p` and `s2` are incremented, thanks to the `++` operators. Repeatedly executing this expression has the effect of copying a series of characters from where `s2` points to where `p` points.

But what causes the loop to terminate? Since the primary operator inside the parentheses is assignment, the `while` statement tests the value of the assignment—the character that was copied. All characters except the null character test true, so the loop won’t terminate until the null character has been copied. And since the loop terminates *after* the assignment, we don’t need a separate statement to put a null character at the end of the new string.

## 13.7 Arrays of Strings

Let’s now turn to a question that we’ll often encounter: what’s the best way to store an array of strings? The obvious solution is to create a two-dimensional array of characters, then store the strings in the array, one per row. Consider the following example:

```
char planets[] [8] = { "Mercury", "Venus", "Earth",  
                      "Mars", "Jupiter", "Saturn",  
                      "Uranus", "Neptune", "Pluto" };
```

(In 2006, the International Astronomical Union demoted Pluto from “planet” to “dwarf planet,” but I’ve left it in the `planets` array for old times’ sake.) Note that we’re allowed to omit the number of rows in the `planets` array—since that’s obvious from the number of elements in the initializer—but C requires that we specify the number of columns.

The figure at the top of the next page shows what the `planets` array will look like. Not all our strings were long enough to fill an entire row of the array, so C padded them with null characters. There’s a bit of wasted space in this array, since only three planets have names long enough to require eight characters (including the terminating null character). The `remind.c` program (Section 13.5) is a glaring example of this kind of waste. It stores reminders in rows of a two-dimensional character array, with 60 characters set aside for each reminder. In our example, the reminders ranged from 18 to 37 characters in length, so the amount of wasted space was considerable.

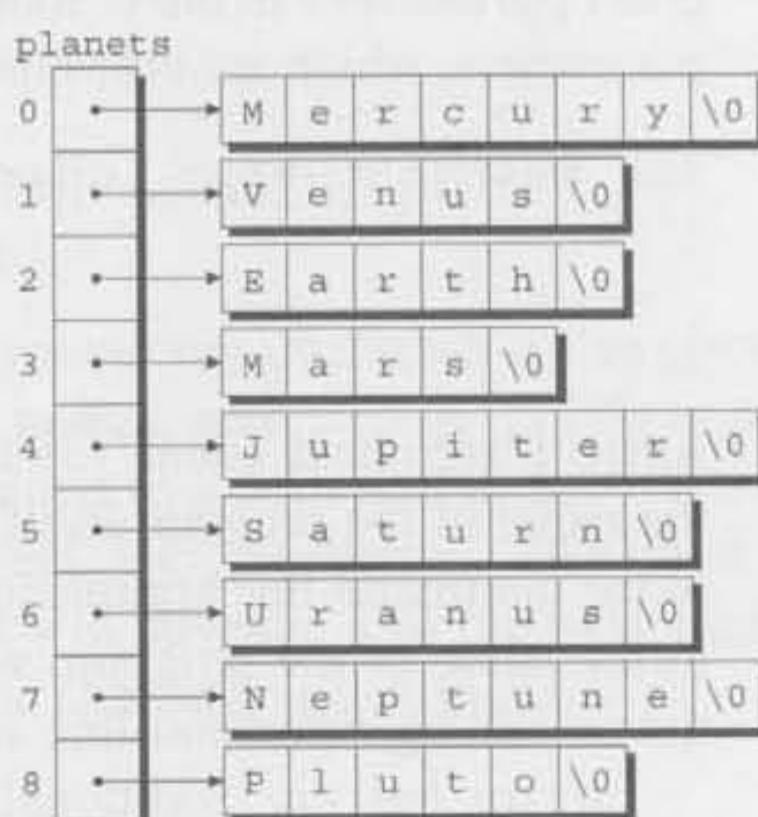
	0	1	2	3	4	5	6	7
0	M	e	r	c	u	r	y	\0
1	V	e	n	u	s	\0	\0	\0
2	E	a	r	t	h	\0	\0	\0
3	M	a	r	s	\0	\0	\0	\0
4	J	u	p	i	t	e	r	\0
5	S	a	t	u	r	n	\0	\0
6	U	r	a	n	u	s	\0	\0
7	N	e	p	t	u	n	e	\0
8	P	l	u	t	o	\0	\0	\0

The inefficiency that's apparent in these examples is common when working with strings, since most collections of strings will have a mixture of long strings and short strings. What we need is a ***ragged array***: a two-dimensional array whose rows can have different lengths. C doesn't provide a "ragged array type," but it does give us the tools to simulate one. The secret is to create an array whose elements are *pointers* to strings.

Here's the `planets` array again, this time as an array of pointers to strings:

```
char *planets[] = {"Mercury", "Venus", "Earth",
                   "Mars", "Jupiter", "Saturn",
                   "Uranus", "Neptune", "Pluto"};
```

Not much of a change, eh? We simply removed one pair of brackets and put an asterisk in front of `planets`. The effect on how `planets` is stored is dramatic, though:



Each element of `planets` is a pointer to a null-terminated string. There are no longer any wasted characters in the strings, although we've had to allocate space for the pointers in the `planets` array.

To access one of the planet names, all we need do is subscript the `planets` array. Because of the relationship between pointers and arrays, accessing a character in a planet name is done in the same way as accessing an element of a two-

dimensional array. To search the `planets` array for strings beginning with the letter M, for example, we could use the following loop:

```
for (i = 0; i < 9; i++)
    if (planets[i][0] == 'M')
        printf("%s begins with M\n", planets[i]);
```

## Command-Line Arguments

When we run a program, we'll often need to supply it with information—a file name, perhaps, or a switch that modifies the program's behavior. Consider the UNIX `ls` command. If we run `ls` by typing

`ls`

at the command line, it will display the names of the files in the current directory. But if we instead type

`ls -l`

then `ls` will display a “long” (detailed) listing of files, showing the size of each file, the file's owner, the date and time the file was last modified, and so forth. To modify the behavior of `ls` further, we can specify that it show details for just one file:

`ls -l remind.c`

`ls` will display detailed information about the file named `remind.c`.

Command-line information is available to all programs, not just operating system commands. To obtain access to these *command-line arguments* (called *program parameters* in the C standard), we must define `main` as a function with two parameters, which are customarily named `argc` and `argv`:

```
int main(int argc, char *argv[])
{
    ...
}
```

`argc` (“argument count”) is the number of command-line arguments (including the name of the program itself). `argv` (“argument vector”) is an array of pointers to the command-line arguments, which are stored in string form. `argv[0]` points to the name of the program, while `argv[1]` through `argv[argc-1]` point to the remaining command-line arguments.

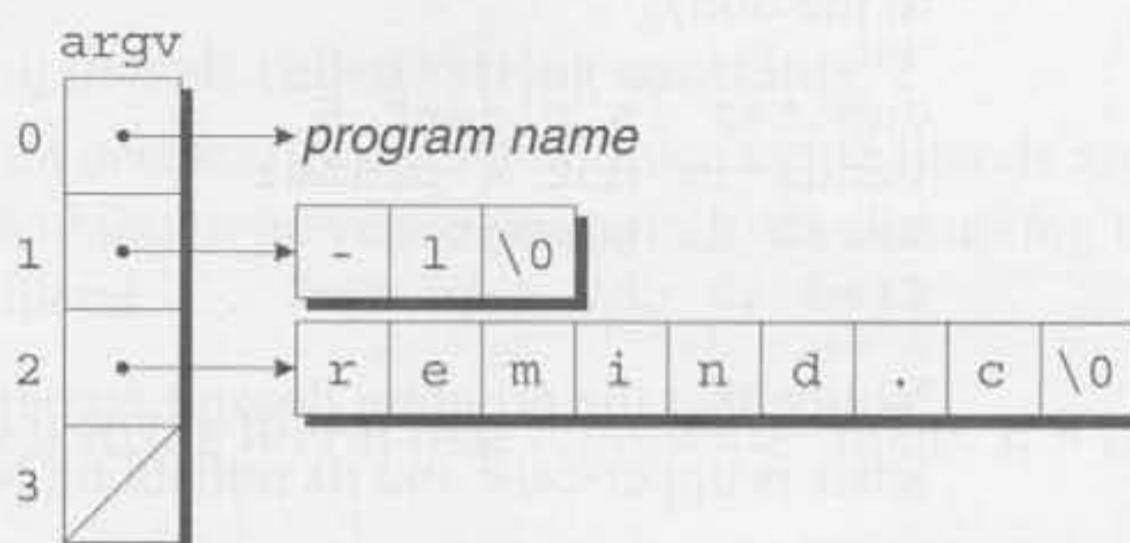
`argv` has one additional element, `argv[argc]`, which is always a *null pointer*—a special pointer that points to nothing. We'll discuss null pointers in a later chapter; for now, all we need to know is that the macro `NULL` represents a null pointer.

If the user enters the command line

`ls -l remind.c`

then `argc` will be 3, `argv[0]` will point to a string containing the program

name, `argv[1]` will point to the string `"-1"`, `argv[2]` will point to the string `"remind.c"`, and `argv[3]` will be a null pointer:



This figure doesn't show the program name in detail, since it may include a path or other information that depends on the operating system. If the program name isn't available, `argv[0]` points to an empty string.

Since `argv` is an array of pointers, accessing command-line arguments is easy. Typically, a program that expects command-line arguments will set up a loop that examines each argument in turn. One way to write such a loop is to use an integer variable as an index into the `argv` array. For example, the following loop prints the command-line arguments, one per line:

```
int i;

for (i = 1; i < argc; i++)
    printf("%s\n", argv[i]);
```

Another technique is to set up a pointer to `argv[1]`, then increment the pointer repeatedly to step through the rest of the array. Since the last element of `argv` is always a null pointer, the loop can terminate when it finds a null pointer in the array:

```
char **p;

for (p = &argv[1]; *p != NULL; p++)
    printf("%s\n", *p);
```

Since `p` is a *pointer to a pointer* to a character, we've got to use it carefully. Setting `p` equal to `&argv[1]` makes sense; `argv[1]` is a pointer to a character, so `&argv[1]` will be a pointer to a pointer. The test `*p != NULL` is OK, since `*p` and `NULL` are both pointers. Incrementing `p` looks good; `p` points to an array element, so incrementing it will advance it to the next element. Printing `*p` is fine, since `*p` points to the first character in a string.

## PROGRAM Checking Planet Names

Our next program, `planet.c`, illustrates how to access command-line arguments. The program is designed to check a series of strings to see which ones are names of planets. When the program is run, the user will put the strings to be tested on the command line:

```
planet Jupiter venus Earth fred
```

The program will indicate whether or not each string is a planet name; if it is, the program will also display the planet's number (with planet 1 being the one closest to the Sun):

```
Jupiter is planet 5
venus is not a planet
Earth is planet 3
fred is not a planet
```

Notice that the program doesn't recognize a string as a planet name unless its first letter is upper-case and its remaining letters are lower-case.

```
planet.c /* Checks planet names */

#include <stdio.h>
#include <string.h>

#define NUM_PLANETS 9

int main(int argc, char *argv[])
{
    char *planets[] = {"Mercury", "Venus", "Earth",
                       "Mars", "Jupiter", "Saturn",
                       "Uranus", "Neptune", "Pluto"};
    int i, j;

    for (i = 1; i < argc; i++) {
        for (j = 0; j < NUM_PLANETS; j++)
            if (strcmp(argv[i], planets[j]) == 0) {
                printf("%s is planet %d\n", argv[i], j + 1);
                break;
            }
        if (j == NUM_PLANETS)
            printf("%s is not a planet\n", argv[i]);
    }

    return 0;
}
```

The program visits each command-line argument in turn, comparing it with the strings in the planets array until it finds a match or reaches the end of the array. The most interesting part of the program is the call of `strcmp`, in which the arguments are `argv[i]` (a pointer to a command-line argument) and `planets[j]` (a pointer to a planet name).

## Q & A

**Q:** How long can a string literal be?

**A:** According to the C89 standard, compilers must allow string literals to be at least

**C99**

509 characters long. (Yes, you read that right—509. Don't ask.) C99 increases the minimum to 4095 characters.

**Q: Why aren't string literals called "string constants"?**

A: Because they're not necessarily constant. Since string literals are accessed through pointers, there's nothing to prevent a program from attempting to modify the characters in a string literal.

**Q: How do we write a string literal that represents "über" if "\xfcber" doesn't work? [p. 278]**

A: The secret is to write two adjacent string literals and let the compiler join them into one. In this example, writing "\xfc" "ber" will give us a string literal that represents the word "über."

**Q: Modifying a string literal seems harmless enough. Why does it cause undefined behavior? [p. 280]**

A: Some compilers try to reduce memory requirements by storing single copies of identical string literals. Consider the following example:

```
char *p = "abc", *q = "abc";
```

A compiler might choose to store "abc" just once, making both p and q point to it. If we were to change "abc" through the pointer p, the string that q points to would also be affected. Needless to say, this could lead to some annoying bugs. Another potential problem is that string literals might be stored in a "read-only" area of memory; a program that attempts to modify such a literal will simply crash.

**Q: Should every array of characters include room for a null character?**

A: Not necessarily, since not every array of characters is used as a string. Including room for the null character (and actually putting one into the array) is necessary only if you're planning to pass it to a function that requires a null-terminated string.

You do *not* need a null character if you'll only be performing operations on individual characters. For example, a program might have an array of characters that it will use to translate from one character set to another:

```
char translation_table[128];
```

The only operation that the program will perform on this array is subscripting. (The value of translation\_table[ch] will be the translated version of the character ch.) We would not consider translation\_table to be a string: it need not contain a null character, and no string operations will be performed on it.

**Q: If printf and scanf expect their first argument to have type char \*, does that mean that the argument can be a string variable instead of a string literal?**

A: Yes, as the following example shows:

```
char fmt[] = "%d\n";
int i;
...
printf(fmt, i);
```

This ability opens the door to some intriguing possibilities—reading a format string as input, for example.

**Q: If I want `printf` to write a string `str`, can't I just supply `str` as the format string, as in the following example?**

```
printf(str);
```

A: Yes, but it's risky. If `str` contains the % character, you won't get the desired result, since `printf` will assume it's the beginning of a conversion specification.

**\*Q: How can `read_line` detect whether `getchar` has failed to read a character? [p. 287]**

A: If it can't read a character, either because of an error or because of end-of-file, `getchar` returns the value EOF, which has type `int`. Here's a revised version of `read_line` that tests whether the return value of `getchar` is EOF. Changes are marked in **bold**:

```
int read_line(char str[], int n)
{
    int ch, i = 0;

    while ((ch = getchar()) != '\n' && ch != EOF)
        if (i < n)
            str[i++] = ch;
        str[i] = '\0';
    return i;
}
```

**Q: Why does `strcmp` return a number that's less than, equal to, or greater than zero? Also, does the exact return value have any significance? [p. 292]**

A: `strcmp`'s return value probably stems from the way the function is traditionally written. Consider the version in Kernighan and Ritchie's *The C Programming Language*:

```
int strcmp(char *s, char *t)
{
    int i;

    for (i = 0; s[i] == t[i]; i++)
        if (s[i] == '\0')
            return 0;
    return s[i] - t[i];
}
```

The return value is the difference between the first “mismatched” characters in the `s` and `t` strings, which will be negative if `s` points to a “smaller” string than `t` and positive if `s` points to a “larger” string. There’s no guarantee that `strcmp` is actually written this way, though, so it’s best not to assume that the magnitude of its return value has any particular meaning.

- Q: My compiler issues a warning when I try to compile the `while` statement in the `strcat` function:**

```
while (*p++ = *s2++)  
;
```

### What am I doing wrong?

- A:** Nothing. Many compilers—but not all, by any means—issue a warning if you use `=` where `==` is normally expected. This warning is valid at least 95% of the time, and it will save you a lot of debugging if you heed it. Unfortunately, the warning isn’t relevant in this particular example; we actually *do* mean to use `=`, not `==`. To get rid of the warning, rewrite the `while` loop as follows:

```
while ((*p++ = *s2++) != 0)  
;
```

Since the `while` statement normally tests whether `*p++ = *s2++` is not 0, we haven’t changed the meaning of the statement. The warning goes away, however, because the statement now tests a condition, not an assignment. With the GCC compiler, putting a pair of parentheses around the assignment is another way to avoid a warning:

```
while ((*p++ = *s2++))  
;
```

- Q: Are the `strlen` and `strcat` functions actually written as shown in Section 13.6?**

- A:** Possibly, although it’s common practice for compiler vendors to write these functions—and many other string functions—in assembly language instead of C. The string functions need to be as fast as possible, since they’re used often and have to deal with strings of arbitrary length. Writing these functions in assembly language makes it possible to achieve great efficiency by taking advantage of any special string-handling instructions that the CPU may provide.

- Q: Why does the C standard use the term “program parameters” instead of “command-line arguments”? [p. 302]**

- A:** Programs aren’t always run from a command line. In a typical graphical user interface, for example, programs are launched with a mouse click. In such an environment, there’s no traditional command line, although there may be other ways of passing information to a program; the term “program parameters” leaves the door open for these alternatives.

**Q:** Do I have to use the names `argc` and `argv` for `main`'s parameters? [p. 302]

**A:** No. Using the names `argc` and `argv` is merely a convention, not a language requirement.

**Q:** I've seen `argv` declared as `**argv` instead of `*argv[]`. Is this legal?

**A:** Certainly. When declaring a parameter, writing `*a` is always the same as writing `a[]`, regardless of the type of `a`'s elements.

**Q:** We've seen how to set up an array whose elements are pointers to string literals. Are there any other applications for arrays of pointers?

**A:** Yes. Although we've focused on arrays of pointers to character strings, that's not the only application of arrays of pointers. We could just as easily have an array whose elements point to any type of data, whether in array form or not. Arrays of pointers are particularly useful in conjunction with dynamic storage allocation.

dynamic storage allocation ▶ 17.1

## Exercises

### Section 13.3

1. The following function calls supposedly write a single new-line character, but some are incorrect. Identify which calls don't work and explain why.

- (a) `printf("%c", '\n');`
- (b) `printf("%c", "\n");`
- (c) `printf("%s", '\n');`
- (d) `printf("%s", "\n");`
- (e) `printf('\n');`
- (f) `printf("\n");`
- (g) `putchar('\n');`
- (h) `putchar("\n");`
- (i) `puts('\n');`
- (j) `puts("\n");`
- (k) `puts("");`

- W 2. Suppose that `p` has been declared as follows:

```
char *p = "abc";
```

Which of the following function calls are legal? Show the output produced by each legal call, and explain why the others are illegal.

- (a) `putchar(p);`
- (b) `putchar(*p);`
- (c) `puts(p);`
- (d) `puts(*p);`

- \*3. Suppose that we call `scanf` as follows:

```
scanf("%d%s%d", &i, s, &j);
```

If the user enters 12abc34 56def78, what will be the values of `i`, `s`, and `j` after the call? (Assume that `i` and `j` are `int` variables and `s` is an array of characters.)

- W 4. Modify the `read_line` function in each of the following ways:

- (a) Have it skip white space before beginning to store input characters.
- (b) Have it stop reading at the first white-space character. Hint: To determine whether or not a character is white space, call the `isspace` function.

isspace function ▶ 23.5

- Section 13.4**
- (c) Have it stop reading at the first new-line character, then store the new-line character in the string.
  - (d) Have it leave behind characters that it doesn't have room to store.

- toupper function ▶ 23.5**
5. (a) Write a function named `capitalize` that capitalizes all letters in its argument. The argument will be a null-terminated string containing arbitrary characters, not just letters. Use array subscripting to access the characters in the string. *Hint:* Use the `toupper` function to convert each character to upper-case.
  - (b) Rewrite the `capitalize` function, this time using pointer arithmetic to access the characters in the string.

- W 6. Write a function named `censor` that modifies a string by replacing every occurrence of `foo` by `xxx`. For example, the string "food fool" would become "xxxd xxxl". Make the function as short as possible without sacrificing clarity.

- Section 13.5**
7. Suppose that `str` is an array of characters. Which one of the following statements is not equivalent to the other three?
    - (a) `*str = 0;`
    - (b) `str[0] = '\0';`
    - (c) `strcpy(str, "");`
    - (d) `strcat(str, "");`  - W \*8. What will be the value of the string `str` after the following statements have been executed?  
`strcpy(str, "tire-bouchon");`  
`strcpy(&str[4], "d-or-wi");`  
`strcat(str, "red?");`
  9. What will be the value of the string `s1` after the following statements have been executed?  
`strcpy(s1, "computer");`  
`strcpy(s2, "science");`  
`if (strcmp(s1, s2) < 0)`  
`strcat(s1, s2);`  
`else`  
`strcat(s2, s1);`  
`s1[strlen(s1)-6] = '\0';`
  - W 10. The following function supposedly creates an identical copy of a string. What's wrong with the function?  

```
char *duplicate(const char *p)
{
    char *q;
    strcpy(q, p);
    return q;
}
```
  11. The Q&A section at the end of this chapter shows how the `strcmp` function might be written using array subscripting. Modify the function to use pointer arithmetic instead.
  12. Write the following function:  
`void get_extension(const char *file_name, char *extension);`

`file_name` points to a string containing a file name. The function should store the extension on the file name in the string pointed to by `extension`. For example, if the file name is "memo.txt", the function will store "txt" in the string pointed to by `extension`. If the file name doesn't have an extension, the function should store an empty string (a single null character) in the string pointed to by `extension`. Keep the function as simple as possible by having it use the `strlen` and `strcpy` functions.

13. Write the following function:

```
void build_index_url(const char *domain, char *index_url);
```

`domain` points to a string containing an Internet domain, such as "knking.com". The function should add "http://www." to the beginning of this string and "/index.html" to the end of the string, storing the result in the string pointed to by `index_url`. (In this example, the result will be "http://www.knking.com/index.html".) You may assume that `index_url` points to a variable that is long enough to hold the resulting string. Keep the function as simple as possible by having it use the `strcat` and `strcpy` functions.

### Section 13.6

- \*14. What does the following program print?

```
#include <stdio.h>

int main(void)
{
    char s[] = "Hsjodi", *p;

    for (p = s; *p; p++)
        --*p;
    puts(s);
    return 0;
}
```

- W\*15. Let `f` be the following function:

```
int f(char *s, char *t)
{
    char *p1, *p2;

    for (p1 = s; *p1; p1++) {
        for (p2 = t; *p2; p2++)
            if (*p1 == *p2) break;
            if (*p2 == '\0') break;
    }
    return p1 - s;
}
```

- (a) What is the value of `f ("abcd", "babc")`?
- (b) What is the value of `f ("abcd", "bcd")`?
- (c) In general, what value does `f` return when passed two strings `s` and `t`?

- W 16. Use the techniques of Section 13.6 to condense the `count_spaces` function of Section 13.4. In particular, replace the `for` statement by a `while` loop.

17. Write the following function:

```
bool test_extension(const char *file_name,
                    const char *extension);
```

**toupper function ▶ 23.5**

`file_name` points to a string containing a file name. The function should return `true` if the file's extension matches the string pointed to by `extension`, ignoring the case of letters. For example, the call `test_extension("memo.txt", "TXT")` would return `true`. Incorporate the “search for the end of a string” idiom into your function. *Hint:* Use the `toupper` function to convert characters to upper-case before comparing them.

18. Write the following function:

```
void remove_filename(char *url);
```

`url` points to a string containing a URL (Uniform Resource Locator) that ends with a file name (such as `"http://www.knking.com/index.html"`). The function should modify the string by removing the file name and the preceding slash. (In this example, the result will be `"http://www.knking.com"`.) Incorporate the “search for the end of a string” idiom into your function. *Hint:* Have the function replace the last slash in the string by a null character.

## Programming Projects

- W 1. Write a program that finds the “smallest” and “largest” in a series of words. After the user enters the words, the program will determine which words would come first and last if the words were listed in dictionary order. The program must stop accepting input when the user enters a four-letter word. Assume that no word is more than 20 letters long. An interactive session with the program might look like this:

```
Enter word: dog
Enter word: zebra
Enter word: rabbit
Enter word: catfish
Enter word: walrus
Enter word: cat
Enter word: fish
```

```
Smallest word: cat
Largest word: zebra
```

*Hint:* Use two strings named `smallest_word` and `largest_word` to keep track of the “smallest” and “largest” words entered so far. Each time the user enters a new word, use `strcmp` to compare it with `smallest_word`; if the new word is “smaller,” use `strcpy` to save it in `smallest_word`. Do a similar comparison with `largest_word`. Use `strlen` to determine when the user has entered a four-letter word.

2. Improve the `remind.c` program of Section 13.5 in the following ways:
  - (a) Have the program print an error message and ignore a reminder if the corresponding day is negative or larger than 31. *Hint:* Use the `continue` statement.
  - (b) Allow the user to enter a day, a 24-hour time, and a reminder. The printed reminder list should be sorted first by day, then by time. (The original program allows the user to enter a time, but it's treated as part of the reminder.)
  - (c) Have the program print a one-year reminder list. Require the user to enter days in the form *month/day*.
3. Modify the `deal.c` program of Section 8.2 so that it prints the full names of the cards it deals:

```

Enter number of cards in hand: 5
Your hand:
Seven of clubs
Two of spades
Five of diamonds
Ace of spades
Two of hearts

```

*Hint:* Replace rank\_code and suit\_code by arrays containing pointers to strings.

- W 4. Write a program named reverse.c that echoes its command-line arguments in reverse order. Running the program by typing

```
reverse void and null
```

should produce the following output:

```
null and void
```

- 5. Write a program named sum.c that adds up its command-line arguments, which are assumed to be integers. Running the program by typing

```
sum 8 24 62
```

should produce the following output:

```
Total: 94
```

atoi function ► 26.2 *Hint:* Use the atoi function to convert each command-line argument from string form to integer form.

- W 6. Improve the planet.c program of Section 13.7 by having it ignore case when comparing command-line arguments with strings in the planets array.

- 7. Modify Programming Project 11 from Chapter 5 so that it uses arrays containing pointers to strings instead of switch statements. For example, instead of using a switch statement to print the word for the first digit, use the digit as an index into an array that contains the strings "twenty", "thirty", and so forth.

- 8. Modify Programming Project 5 from Chapter 7 so that it includes the following function:

```
int compute_scrabble_value(const char *word);
```

The function returns the SCRABBLE value of the string pointed to by word.

- 9. Modify Programming Project 10 from Chapter 7 so that it includes the following function:

```
int compute_vowel_count(const char *sentence);
```

The function returns the number of vowels in the string pointed to by the sentence parameter.

- 10. Modify Programming Project 11 from Chapter 7 so that it includes the following function:

```
void reverse_name(char *name);
```

The function expects name to point to a string containing a first name followed by a last name. It modifies the string so that the last name comes first, followed by a comma, a space, the first initial, and a period. The original string may contain extra spaces before the first name, between the first and last names, and after the last name.

- 11. Modify Programming Project 13 from Chapter 7 so that it includes the following function:

```
double compute_average_word_length(const char *sentence);
```

The function returns the average length of the words in the string pointed to by sentence.

12. Modify Programming Project 14 from Chapter 8 so that it stores the words in a two-dimensional `char` array as it reads the sentence, with each row of the array storing a single word. Assume that the sentence contains no more than 30 words and no word is more than 20 characters long. Be sure to store a null character at the end of each word so that it can be treated as a string.
13. Modify Programming Project 15 from Chapter 8 so that it includes the following function:

```
void encrypt(char *message, int shift);
```

The function expects `message` to point to a string containing the message to be encrypted; `shift` represents the amount by which each letter in the message is to be shifted.
14. Modify Programming Project 16 from Chapter 8 so that it includes the following function:

```
bool are_anagrams(const char *word1, const char *word2);
```

The function returns `true` if the strings pointed to by `word1` and `word2` are anagrams.
15. Modify Programming Project 6 from Chapter 10 so that it includes the following function:

```
int evaluate_RPN_expression(const char *expression);
```

The function returns the value of the RPN expression pointed to by `expression`.
16. Modify Programming Project 1 from Chapter 12 so that it includes the following function:

```
void reverse(char *message);
```

The function reverses the string pointed to by `message`. *Hint:* Use two pointers, one initially pointing to the first character of the string and the other initially pointing to the last character. Have the function reverse these characters and then move the pointers toward each other, repeating the process until the pointers meet.
17. Modify Programming Project 2 from Chapter 12 so that it includes the following function:

```
bool is_palindrome(const char *message);
```

The function returns `true` if the string pointed to by `message` is a palindrome.
18. Write a program that accepts a date from the user in the form `mm/dd/yyyy` and then displays it in the form `month dd, yyyy`, where `month` is the name of the month:

```
Enter a date (mm/dd/yyyy) : 2/17/2011
You entered the date February 17, 2011
```

Store the month names in an array that contains pointers to strings.



# 14 The Preprocessor

*There will always be things we wish to say in our programs  
that in all known languages can only be said poorly.*

In previous chapters, I've used the `#define` and `#include` directives without going into detail about what they do. These directives—and others that we haven't yet covered—are handled by the *preprocessor*, a piece of software that edits C programs just prior to compilation. Its reliance on a preprocessor makes C (along with C++) unique among major programming languages.

The preprocessor is a powerful tool, but it also can be a source of hard-to-find bugs. Moreover, the preprocessor can easily be misused to create programs that are almost impossible to understand. Although some C programmers depend heavily on the preprocessor, I recommend that it—like so many other things in life—be used in moderation.

This chapter begins by describing how the preprocessor works (Section 14.1) and giving some general rules that affect all preprocessing directives (Section 14.2). Sections 14.3 and 14.4 cover two of the preprocessor's major capabilities: macro definition and conditional compilation. (I'll defer detailed coverage of file inclusion, the other major capability, until Chapter 15.) Section 14.5 discusses the preprocessor's lesser-used directives: `#error`, `#line`, and `#pragma`.

## 14.1 How the Preprocessor Works

The behavior of the preprocessor is controlled by *preprocessing directives*: commands that begin with a `#` character. We've encountered two of these directives, `#define` and `#include`, in previous chapters.

The `#define` directive defines a *macro*—a name that represents something else, such as a constant or frequently used expression. The preprocessor responds to a `#define` directive by storing the name of the macro together with its definition.

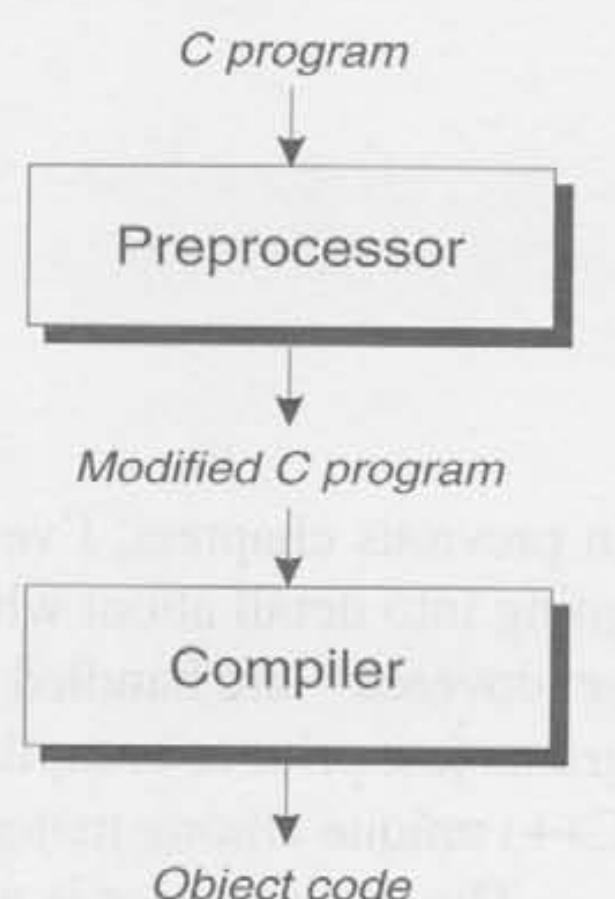
When the macro is used later in the program, the preprocessor “expands” the macro, replacing it by its defined value.

The `#include` directive tells the preprocessor to open a particular file and “include” its contents as part of the file being compiled. For example, the line

```
#include <stdio.h>
```

instructs the preprocessor to open the file named `stdio.h` and bring its contents into the program. (Among other things, `stdio.h` contains prototypes for C’s standard input/output functions.)

The following diagram shows the preprocessor’s role in the compilation process:



The input to the preprocessor is a C program, possibly containing directives. The preprocessor executes these directives, removing them in the process. The output of the preprocessor is another C program: an edited version of the original program, containing no directives. The preprocessor’s output goes directly into the compiler, which checks the program for errors and translates it to object code (machine instructions).

To see what the preprocessor does, let’s apply it to the `celsius.c` program of Section 2.6. Here’s the original program:

```
/* Converts a Fahrenheit temperature to Celsius */

#include <stdio.h>

#define FREEZING_PT 32.0f
#define SCALE_FACTOR (5.0f / 9.0f)

int main(void)
{
    float fahrenheit, celsius;

    printf("Enter Fahrenheit temperature: ");
    scanf("%f", &fahrenheit);

    celsius = (fahrenheit - FREEZING_PT) * SCALE_FACTOR;
```

```

    printf("Celsius equivalent is: %.1f\n", celsius);

    return 0;
}

```

After preprocessing, the program will have the following appearance:

```

Blank line
Blank line
Lines brought in from stdio.h
Blank line
Blank line
Blank line
Blank line
int main(void)
{
    float fahrenheit, celsius;

    printf("Enter Fahrenheit temperature: ");
    scanf("%f", &fahrenheit);

    celsius = (fahrenheit - 32.0f) * (5.0f / 9.0f);

    printf("Celsius equivalent is: %.1f\n", celsius);

    return 0;
}

```

The preprocessor responded to the #include directive by bringing in the contents of stdio.h. The preprocessor also removed the #define directives and replaced FREEZING\_PT and SCALE\_FACTOR wherever they appeared later in the file. Notice that the preprocessor doesn't remove lines containing directives; instead, it simply makes them empty.

As this example shows, the preprocessor does a bit more than just execute directives. In particular, it replaces each comment with a single space character. Some preprocessors go further and remove unnecessary white-space characters, including spaces and tabs at the beginning of indented lines.

In the early days of C, the preprocessor was a separate program that fed its output into the compiler. Nowadays, the preprocessor is often part of the compiler, and some of its output may not necessarily be C code. (For example, including a standard header such as <stdio.h> may have the effect of making its functions available to the program without necessarily copying the contents of the header into the program's source code.) Still, it's useful to think of the preprocessor as separate from the compiler. In fact, most C compilers provide a way to view the output of the preprocessor. Some compilers generate preprocessor output when a certain option is specified (GCC will do so when the -E option is used). Others come with a separate program that behaves like the integrated preprocessor. Check your compiler's documentation for more information.

A word of caution: The preprocessor has only a limited knowledge of C. As a result, it's quite capable of creating illegal programs as it executes directives. Often the original program looks fine, making errors harder to find. In complicated

programs, examining the output of the preprocessor may prove useful for locating this kind of error.

## 14.2 Preprocessing Directives

Most preprocessing directives fall into one of three categories:

- **Macro definition.** The `#define` directive defines a macro; the `#undef` directive removes a macro definition.
- **File inclusion.** The `#include` directive causes the contents of a specified file to be included in a program.
- **Conditional compilation.** The `#if`, `#ifdef`, `#ifndef`, `#elif`, `#else`, and `#endif` directives allow blocks of text to be either included in or excluded from a program, depending on conditions that can be tested by the preprocessor.

The remaining directives—`#error`, `#line`, and `#pragma`—are more specialized and therefore used less often. We'll devote the rest of this chapter to an in-depth examination of preprocessing directives. The only directive we won't discuss in detail is `#include`, since it's covered in Section 15.2.

Before we go further, let's look at a few rules that apply to all directives:

- **Directives always begin with the # symbol.** The `#` symbol need not be at the beginning of a line, as long as only white space precedes it. After the `#` comes the name of the directive, followed by any other information the directive requires.
- **Any number of spaces and horizontal tab characters may separate the tokens in a directive.** For example, the following directive is legal:

```
# define N 100
```

- **Directives always end at the first new-line character, unless explicitly continued.** To continue a directive to the next line, we must end the current line with a `\` character. For example, the following directive defines a macro that represents the capacity of a hard disk, measured in bytes:

```
#define DISK_CAPACITY (SIDES * \
                      TRACKS_PER_SIDE * \
                      SECTORS_PER_TRACK * \
                      BYTES_PER_SECTOR)
```

- **Directives can appear anywhere in a program.** Although we usually put `#define` and `#include` directives at the beginning of a file, other directives are more likely to show up later, even in the middle of function definitions.
- **Comments may appear on the same line as a directive.** In fact, it's good practice to put a comment at the end of a macro definition to explain the meaning of the macro:

```
#define FREEZING_PT 32.0f /* freezing point of water */
```

## 14.3 Macro Definitions

The macros that we've been using since Chapter 2 are known as *simple* macros, because they have no parameters. The preprocessor also supports *parameterized* macros. We'll look first at simple macros, then at parameterized macros. After covering them separately, we'll examine properties shared by both.

### Simple Macros

The definition of a *simple macro* (or *object-like macro*, as it's called in the C standard) has the form

**#define directive  
(simple macro)**

**#define identifier replacement-list**

*replacement-list* is any sequence of *preprocessing tokens*, which are similar to the tokens discussed in Section 2.8. Whenever we use the term "token" in this chapter, it means "preprocessing token."

A macro's replacement list may include identifiers, keywords, numeric constants, character constants, string literals, operators, and punctuation. When it encounters a macro definition, the preprocessor makes a note that *identifier* represents *replacement-list*; wherever *identifier* appears later in the file, the preprocessor substitutes *replacement-list*.



Don't put any extra symbols in a macro definition—they'll become part of the replacement list. Putting the = symbol in a macro definition is a common error:

```
#define N = 100    /*** WRONG ***/
...
int a[N];          /* becomes int a[= 100]; */
```

In this example, we've (incorrectly) defined N to be a pair of tokens (= and 100). Ending a macro definition with a semicolon is another popular mistake:

```
#define N 100;      /*** WRONG ***/
...
int a[N];          /* becomes int a[100;]; */
```

Here N is defined to be the tokens 100 and ;.

The compiler will detect most errors caused by extra symbols in a macro definition. Unfortunately, the compiler will flag each use of the macro as incorrect, rather than identifying the actual culprit—the macro's definition—which will have been removed by the preprocessor.

**Q&A**

Simple macros are primarily used for defining what Kernighan and Ritchie call "manifest constants." Using macros, we can give names to numeric, character, and string values:

```
#define STR_LEN 80
#define TRUE    1
#define FALSE   0
#define PI      3.14159
#define CR      '\r'
#define EOS     '\0'
#define MEM_ERR "Error: not enough memory"
```

Using `#define` to create names for constants has several significant advantages:

- ***It makes programs easier to read.*** The name of the macro—if well-chosen—helps the reader understand the meaning of the constant. The alternative is a program full of “magic numbers” that can easily mystify the reader.
- ***It makes programs easier to modify.*** We can change the value of a constant throughout a program by modifying a single macro definition. “Hard-coded” constants are more difficult to change, especially since they sometimes appear in a slightly altered form. (For example, a program with an array of length 100 may have a loop that goes from 0 to 99. If we merely try to locate occurrences of 100 in the program, we’ll miss the 99.)
- ***It helps avoid inconsistencies and typographical errors.*** If a numerical constant like 3.14159 appears many times in a program, chances are it will occasionally be written 3.1416 or 3.14195 by accident.

Although simple macros are most often used to define names for constants, they do have other applications:

- ***Making minor changes to the syntax of C.*** We can—in effect—alter the syntax of C by defining macros that serve as alternate names for C symbols. For example, programmers who prefer Pascal’s `begin` and `end` to C’s `{` and `}` can define the following macros:

```
#define BEGIN {
#define END }
```

We could go so far as to invent our own language. For example, we might create a `LOOP` “statement” that establishes an infinite loop:

```
#define LOOP for (;;) 
```

Changing the syntax of C usually isn’t a good idea, though, since it can make programs harder for others to understand.

- ***Renaming types.*** In Section 5.2, we created a Boolean type by renaming `int`:

```
#define BOOL int
```

Although some programmers use macros for this purpose, type definitions are a superior way to define type names.

- ***Controlling conditional compilation.*** Macros play an important role in controlling conditional compilation, as we’ll see in Section 14.4. For example, the presence of the following line in a program might indicate that it’s to be com-

piled in “debugging mode,” with extra statements included to produce debugging output:

```
#define DEBUG
```

Incidentally, it’s legal for a macro’s replacement list to be empty, as this example shows.

When macros are used as constants, C programmers customarily capitalize all letters in their names. However, there’s no consensus as to how to capitalize macros used for other purposes. Since macros (especially parameterized macros) can be a source of bugs, some programmers like to draw attention to them by using all upper-case letters in their names. Others prefer lower-case names, following the style of Kernighan and Ritchie’s *The C Programming Language*.

## Parameterized Macros

The definition of a *parameterized macro* (also known as a *function-like macro*) has the form

**#define directive  
(parameterized macro)**

```
#define identifier( x1 , x2 , ... , xn ) replacement-list
```

where  $x_1, x_2, \dots, x_n$  are identifiers (the macro’s *parameters*). The parameters may appear as many times as desired in the replacement list.



There must be *no space* between the macro name and the left parenthesis. If space is left, the preprocessor will assume that we’re defining a simple macro; it will treat  $(x_1, x_2, \dots, x_n)$  as part of the replacement list.

When the preprocessor encounters the definition of a parameterized macro, it stores the definition away for later use. Wherever a macro *invocation* of the form  $identifier(y_1, y_2, \dots, y_n)$  appears later in the program (where  $y_1, y_2, \dots, y_n$  are sequences of tokens), the preprocessor replaces it with *replacement-list*, substituting  $y_1$  for  $x_1$ ,  $y_2$  for  $x_2$ , and so forth.

For example, suppose that we’ve defined the following macros:

```
#define MAX(x,y)    ((x) > (y) ? (x) : (y))
#define IS_EVEN(n)   ((n) %2==0)
```

(The number of parentheses in these macros may seem excessive, but there’s a reason, as we’ll see later in this section.) Now suppose that we invoke the two macros in the following way:

```
i = MAX(j+k, m-n);
if (IS_EVEN(i)) i++;
```

The preprocessor will replace these lines by

```
i = ((j+k) > (m-n) ? (j+k) : (m-n));
if (((i)%2==0)) i++;
```

As this example shows, parameterized macros often serve as simple functions. MAX behaves like a function that computes the larger of two values. IS\_EVEN behaves like a function that returns 1 if its argument is an even number and 0 otherwise.

Here's a more complicated macro that behaves like a function:

```
#define TOUPPER(c) ('a'<=(c)&&(c)<='z'? (c)-'a'+'A':(c))
```

<ctype.h> header ▶ 23.5

This macro tests whether the character c is between 'a' and 'z'. If so, it produces the upper-case version of c by subtracting 'a' and adding 'A'. If not, it leaves c unchanged. (The <ctype.h> header provides a similar function named toupper that's more portable.)

A parameterized macro may have an empty parameter list. Here's an example:

```
#define getchar() getc(stdin)
```

The empty parameter list isn't really needed, but it makes getchar resemble a function. (Yes, this is the same getchar that belongs to <stdio.h>. We'll see in Section 22.4 that getchar is usually implemented as a macro as well as a function.)

Using a parameterized macro instead of a true function has a couple of advantages:

- ***The program may be slightly faster.*** A function call usually requires some overhead during program execution—context information must be saved, arguments copied, and so forth. A macro invocation, on the other hand, requires no run-time overhead. (Note, however, that C99's inline functions provide a way to avoid this overhead without the use of macros.)
- ***Macros are “generic.”*** Macro parameters, unlike function parameters, have no particular type. As a result, a macro can accept arguments of any type, provided that the resulting program—after preprocessing—is valid. For example, we could use the MAX macro to find the larger of two values of type int, long, float, double, and so forth.

But parameterized macros also have disadvantages:

- ***The compiled code will often be larger.*** Each macro invocation causes the insertion of the macro's replacement list, thereby increasing the size of the source program (and hence the compiled code). The more often the macro is used, the more pronounced this effect is. The problem is compounded when macro invocations are nested. Consider what happens when we use MAX to find the largest of three numbers:

```
n = MAX(i, MAX(j, k));
```

Here's the same statement after preprocessing:

```
n = ((i)>(((j)>(k)?(j):(k)))?(i):(((j)>(k)?(j):(k))));
```

C99

inline functions ▶ 18.6

- **Arguments aren't type-checked.** When a function is called, the compiler checks each argument to see if it has the appropriate type. If not, either the argument is converted to the proper type or the compiler produces an error message. Macro arguments aren't checked by the preprocessor, nor are they converted.
- **It's not possible to have a pointer to a macro.** As we'll see in Section 17.7, C allows pointers to functions, a concept that's quite useful in certain programming situations. Macros are removed during preprocessing, so there's no corresponding notion of "pointer to a macro"; as a result, macros can't be used in these situations.
- **A macro may evaluate its arguments more than once.** A function evaluates its arguments only once; a macro may evaluate its arguments two or more times. Evaluating an argument more than once can cause unexpected behavior if the argument has side effects. Consider what happens if one of MAX's arguments has a side effect:

```
n = MAX(i++, j);
```

Here's the same line after preprocessing:

```
n = ((i++) > (j) ? (i++) : (j));
```

If *i* is larger than *j*, then *i* will be (incorrectly) incremented twice and *n* will be assigned an unexpected value.




---

Errors caused by evaluating a macro argument more than once can be difficult to find, because a macro invocation looks the same as a function call. To make matters worse, a macro may work properly most of the time, failing only for certain arguments that have side effects. For self-protection, it's a good idea to avoid side effects in arguments.

---

Parameterized macros are good for more than just simulating functions. In particular, they're often used as patterns for segments of code that we find ourselves repeating. Suppose that we grow tired of writing

```
printf("%d\n", i);
```

every time we need to print an integer *i*. We might define the following macro, which makes it easier to display integers:

```
#define PRINT_INT(n) printf("%d\n", n)
```

Once PRINT\_INT has been defined, the preprocessor will turn the line

```
PRINT_INT(i/j);
```

into

```
printf("%d\n", i/j);
```

## The # Operator

Macro definitions may contain two special operators, # and ##. Neither operator is recognized by the compiler; instead, they're executed during preprocessing.

### Q&A

The # operator converts a macro argument into a string literal; it can appear only in the replacement list of a parameterized macro. (The operation performed by # is known as “stringization,” a term that I’m sure you won’t find in the dictionary.)

There are a number of uses for #; let’s consider just one. Suppose that we decide to use the PRINT\_INT macro during debugging as a convenient way to print the values of integer variables and expressions. The # operator makes it possible for PRINT\_INT to label each value that it prints. Here’s our new version of PRINT\_INT:

```
#define PRINT_INT(n) printf(#n " = %d\n", n)
```

The # operator in front of n instructs the preprocessor to create a string literal from PRINT\_INT’s argument. Thus, the invocation

```
PRINT_INT(i/j);
```

will become

```
printf("i/j" " = %d\n", i/j);
```

We saw in Section 13.1 that the compiler automatically joins adjacent string literals, so this statement is equivalent to

```
printf("i/j = %d\n", i/j);
```

When the program is executed, printf will display both the expression i/j and its value. If i is 11 and j is 2, for example, the output will be

```
i/j = 5
```

## The ## Operator

The ## operator can “paste” two tokens (identifiers, for example) together to form a single token. (Not surprisingly, the ## operation is known as “token-pasting.”) If one of the operands is a macro parameter, pasting occurs after the parameter has been replaced by the corresponding argument. Consider the following macro:

```
#define MK_ID(n) i##n
```

When MK\_ID is invoked (as MK\_ID(1), say), the preprocessor first replaces the parameter n by the argument (1 in this case). Next, the preprocessor joins i and 1 to make a single token (i1). The following declaration uses MK\_ID to create three identifiers:

```
int MK_ID(1), MK_ID(2), MK_ID(3);
```

After preprocessing, this declaration becomes

```
int i1, i2, i3;
```

The ## operator isn't one of the most frequently used features of the preprocessor; in fact, it's hard to think of many situations that require it. To find a realistic application of ##, let's reconsider the MAX macro described earlier in this section. As we observed then, MAX doesn't behave properly if its arguments have side effects. The alternative to using the MAX macro is to write a max function. Unfortunately, one max function usually isn't enough; we may need a max function whose arguments are `int` values, one whose arguments are `float` values, and so on. All these versions of max would be identical except for the types of the arguments and the return type, so it seems a shame to define each one from scratch.

The solution is to write a macro that expands into the definition of a max function. The macro will have a single parameter, `type`, which represents the type of the arguments and the return value. There's just one snag: if we use the macro to create more than one max function, the program won't compile. (C doesn't allow two functions to have the same name if both are defined in the same file.) To solve this problem, we'll use the ## operator to create a different name for each version of max. Here's what the macro will look like:

```
#define GENERIC_MAX(type) \
type type##_max(type x, type y) \
{ \
    return x > y ? x : y; \
}
```

Notice how `type` is joined with `_max` to form the name of the function.

Suppose that we happen to need a max function that works with `float` values. Here's how we'd use `GENERIC_MAX` to define the function:

```
GENERIC_MAX(float)
```

The preprocessor expands this line into the following code:

```
float float_max(float x, float y) { return x > y ? x : y; }
```

## General Properties of Macros

Now that we've discussed both simple and parameterized macros, let's look at some rules that apply to both:

- **A macro's replacement list may contain invocations of other macros.** For example, we could define the macro `TWO_PI` in terms of the macro `PI`:

```
#define PI      3.14159
#define TWO_PI (2*PI)
```

When it encounters `TWO_PI` later in the program, the preprocessor replaces it by `(2*PI)`. The preprocessor then *rescans* the replacement list to see if it

**Q&A**

contains invocations of other macros (PI, in this case). The preprocessor will rescan the replacement list as many times as necessary to eliminate all macro names.

- ***The preprocessor replaces only entire tokens, not portions of tokens.*** As a result, the preprocessor ignores macro names that are embedded in identifiers, character constants, and string literals. For example, suppose that a program contains the following lines:

```
#define SIZE 256
int BUFFER_SIZE;
if (BUFFER_SIZE > SIZE)
    puts("Error: SIZE exceeded");
```

After preprocessing, these lines will have the following appearance:

```
int BUFFER_SIZE;
if (BUFFER_SIZE > 256)
    puts("Error: SIZE exceeded");
```

The identifier `BUFFER_SIZE` and the string "Error: SIZE exceeded" weren't affected by preprocessing, even though both contain the word `SIZE`.

- ***A macro definition normally remains in effect until the end of the file in which it appears.*** Since macros are handled by the preprocessor, they don't obey normal scope rules. A macro defined inside the body of a function isn't local to that function; it remains defined until the end of the file.
- ***A macro may not be defined twice unless the new definition is identical to the old one.*** Differences in spacing are allowed, but the tokens in the macro's replacement list (and the parameters, if any) must be the same.
- ***Macros may be “undefined” by the #undef directive.*** The `#undef` directive has the form

**#undef directive**

**#undef identifier**

where *identifier* is a macro name. For example, the directive

```
#undef N
```

removes the current definition of the macro `N`. (If `N` hasn't been defined as a macro, the `#undef` directive has no effect.) One use of `#undef` is to remove the existing definition of a macro so that it can be given a new definition.

## Parentheses in Macro Definitions

The replacement lists in our macro definitions have been full of parentheses. Is it really necessary to have so many? The answer is an emphatic yes; if we use fewer

parentheses, the macros will sometimes give unexpected—and undesirable—results.

There are two rules to follow when deciding where to put parentheses in a macro definition. First, if the macro's replacement list contains an operator, always enclose the replacement list in parentheses:

```
#define TWO_PI (2*3.14159)
```

Second, if the macro has parameters, put parentheses around each parameter every time it appears in the replacement list:

```
#define SCALE(x) ((x)*10)
```

Without the parentheses, we can't guarantee that the compiler will treat replacement lists and arguments as whole expressions. The compiler may apply the rules of operator precedence and associativity in ways that we didn't anticipate.

To illustrate the importance of putting parentheses around a macro's replacement list, consider the following macro definition, in which the parentheses are missing:

```
#define TWO_PI 2*3.14159
/* needs parentheses around replacement list */
```

During preprocessing, the statement

```
conversion_factor = 360/TWO_PI;
```

becomes

```
conversion_factor = 360/2*3.14159;
```

The division will be performed before the multiplication, yielding a result different from the one intended.

Putting parentheses around the replacement list isn't enough if the macro has parameters—each occurrence of a parameter needs parentheses as well. For example, suppose that SCALE is defined as follows:

```
#define SCALE(x) (x*10) /* needs parentheses around x */
```

During preprocessing, the statement

```
j = SCALE(i+1);
```

becomes

```
j = (i+1*10);
```

Since multiplication takes precedence over addition, this statement is equivalent to

```
j = i+10;
```

Of course, what we wanted was

```
j = (i+1)*10;
```



A shortage of parentheses in a macro definition can cause some of C's most frustrating errors. The program will usually compile and the macro will appear to work, failing only at the least convenient times.

## Creating Longer Macros

The comma operator can be useful for creating more sophisticated macros by allowing us to make the replacement list a series of expressions. For example, the following macro will read a string and then print it:

```
#define ECHO(s) (gets(s), puts(s))
```

Calls of `gets` and `puts` are expressions, so it's perfectly legal to combine them using the comma operator. We can invoke `ECHO` as though it were a function:

```
ECHO(str); /* becomes (gets(str), puts(str)); */
```

Instead of using the comma operator in the definition of `ECHO`, we could have enclosed the calls of `gets` and `puts` in braces to form a compound statement:

```
#define ECHO(s) { gets(s); puts(s); }
```

Unfortunately, this method doesn't work as well. Suppose that we use `ECHO` in an `if` statement:

```
if (echo_flag)
    ECHO(str);
else
    gets(str);
```

Replacing `ECHO` gives the following result:

```
if (echo_flag)
    { gets(str); puts(str); };
else
    gets(str);
```

The compiler treats the first two lines as a complete `if` statement:

```
if (echo_flag)
    { gets(str); puts(str); }
```

It treats the semicolon that follows as a null statement and produces an error message for the `else` clause, since it doesn't belong to any `if`. We could solve the problem by remembering not to put a semicolon after each invocation of `ECHO`, but then the program would look odd.

The comma operator solves this problem for `ECHO`, but not for all macros. Suppose that a macro needs to contain a series of *statements*, not just a series of *expressions*. The comma operator is of no help; it can glue together expressions,

but not statements. The solution is to wrap the statements in a do loop whose condition is false (and which therefore will be executed just once):

```
do { ... } while (0)
```

Notice that the do statement isn't complete—it needs a semicolon at the end. To see this trick (ahem, technique) in action, let's incorporate it into our ECHO macro:

```
#define ECHO(s)      \
    do {           \
        gets(s);   \
        puts(s);   \
    } while (0)
```

When ECHO is used, it must be followed by a semicolon, which completes the do statement:

```
ECHO(str);          \
/* becomes do { gets(str); puts(str); } while (0); */
```

## Predefined Macros

C has several predefined macros. Each macro represents an integer constant or string literal. As Table 14.1 shows, these macros provide information about the current compilation or about the compiler itself.

**Table 14.1**

Predefined Macros

Name	Description
<code>_LINE_</code>	Line number of file being compiled
<code>_FILE_</code>	Name of file being compiled
<code>_DATE_</code>	Date of compilation (in the form "Mmm dd yyyy")
<code>_TIME_</code>	Time of compilation (in the form "hh:mm:ss")
<code>_STDC_</code>	1 if the compiler conforms to the C standard (C89 or C99)

The `_DATE_` and `_TIME_` macros identify when a program was compiled. For example, suppose that a program begins with the following statements:

```
printf("Wacky Windows (c) 2010 Wacky Software, Inc.\n");
printf("Compiled on %s at %s\n", _DATE_, _TIME_);
```

Each time it begins to execute, the program will print two lines of the form

```
Wacky Windows (c) 2010 Wacky Software, Inc.
Compiled on Dec 23 2010 at 22:18:48
```

This information can be helpful for distinguishing among different versions of the same program.

We can use the `_LINE_` and `_FILE_` macros to help locate errors. Consider the problem of detecting the location of a division by zero. When a C program terminates prematurely because it divided by zero, there's usually no indication of which division caused the problem. The following macro can help us pinpoint the source of the error:

```
#define CHECK_ZERO(divisor) \
    if (divisor == 0) \
        printf("**** Attempt to divide by zero on line %d " \
               "of file %s ***\n", __LINE__, __FILE__)
```

The CHECK\_ZERO macro would be invoked prior to a division:

```
CHECK_ZERO(j);  
k = i / j;
```

If j happens to be zero, a message of the following form will be printed:

```
*** Attempt to divide by zero on line 9 of file foo.c ***
```

assert macro ▶ 24.1

Error-detecting macros like this one are quite useful. In fact, the C library has a general-purpose error-detecting macro named assert.

The \_\_STDC\_\_ macro exists and has the value 1 if the compiler conforms to the C standard (either C89 or C99). By having the preprocessor test this macro, a program can adapt to a compiler that predates the C89 standard (see Section 14.4 for an example).

### C99

## Additional Predefined Macros in C99

C99 provides a few additional predefined macros (Table 14.2).

**Table 14.2**  
Additional Predefined  
Macros in C99

Name	Description
__STDC__HOSTED__	1 if this is a hosted implementation; 0 if it is freestanding
__STDC__VERSION__	Version of C standard supported
__STDC_IEC_559__ <sup>†</sup>	1 if IEC 60559 floating-point arithmetic is supported
__STDC_IEC_559_COMPLEX__ <sup>†</sup>	1 if IEC 60559 complex arithmetic is supported
__STDC_ISO_10646__ <sup>†</sup>	yyymmL if wchar_t values match the ISO 10646 standard of the specified year and month

<sup>†</sup>Conditionally defined

complex types ▶ 27.3

### Q&A

To understand the meaning of \_\_STDC\_\_HOSTED\_\_, we need some new vocabulary. An *implementation* of C consists of the compiler plus other software necessary to execute C programs. C99 divides implementations into two categories: hosted and freestanding. A *hosted implementation* must accept any program that conforms to the C99 standard, whereas a *freestanding implementation* doesn't have to compile programs that use complex types or standard headers beyond a few of the most basic. (In particular, a freestanding implementation doesn't have to support the <stdio.h> header.) The \_\_STDC\_\_HOSTED\_\_ macro represents the constant 1 if the compiler is a hosted implementation; otherwise, the macro has the value 0.

The \_\_STDC\_\_VERSION\_\_ macro provides a way to check which version of the C standard is recognized by the compiler. This macro first appeared in Amendment 1 to the C89 standard, where its value was specified to be the long

integer constant 199409L (representing the year and month of the amendment). If a compiler conforms to the C99 standard, the value is 199901L. For each subsequent version of the standard (and each amendment to the standard), this macro will have a different value.

A C99 compiler may (or may not) define three additional macros. Each macro is defined only if the compiler meets a certain requirement:

- `_STDC_IEC_559_` is defined (and has the value 1) if the compiler performs floating-point arithmetic according to the IEC 60559 standard (another name for the IEEE 754 standard).
- `_STDC_IEC_559_COMPLEX_` is defined (and has the value 1) if the compiler performs complex arithmetic according to the IEC 60559 standard.
- `_STDC_ISO_10646_` is defined as an integer constant of the form `yyyymmL` (for example, `199712L`) if values of type `wchar_t` are represented by the codes in the ISO/IEC 10646 standard (with revisions as of the specified year and month).

### C99

## Empty Macro Arguments

C99 allows any or all of the arguments in a macro call to be empty. Such a call will contain the same number of commas as a normal call, however. (That way, it's easy to see which arguments have been omitted.)

In most cases, the effect of an empty argument is clear. Wherever the corresponding parameter name appears in the replacement list, it's replaced by nothing—it simply disappears from the replacement list. Here's an example:

```
#define ADD(x,y) (x+y)
```

After preprocessing, the statement

```
i = ADD(j,k);
```

becomes

```
i = (j+k);
```

whereas the statement

```
i = ADD(,k);
```

becomes

```
i = (+k);
```

When an empty argument is an operand of the # or ## operators, special rules apply. If an empty argument is “stringized” by the # operator, the result is "" (the empty string):

```
#define MK_STR(x) #x
...
char empty_string[] = MK_STR();
```

After preprocessing, the declaration will have the following appearance:

```
char empty_string[] = "";
```

If one of the arguments of the ## operator is empty, it's replaced by an invisible "placeholder" token. Concatenating an ordinary token with a placeholder token yields the original token (the placeholder disappears). If two placeholder tokens are concatenated, the result is a single placeholder. Once macro expansion has been completed, placeholder tokens disappear from the program. Consider the following example:

```
#define JOIN(x,y,z) x##y##z
...
int JOIN(a,b,c), JOIN(a,,c), JOIN(,c);
```

After preprocessing, the declaration will have the following appearance:

```
int abc, ab, ac, c;
```

The missing arguments were replaced by placeholder tokens, which then disappeared when concatenated with any nonempty arguments. All three arguments to the JOIN macro could even be missing, which would yield an empty result.

### C99

## Macros with a Variable Number of Arguments

variable-length argument lists  
► 26.1

In C89, a macro must have a fixed number of arguments, if it has any at all. C99 loosens things up a bit, allowing macros that take an unlimited number of arguments. This feature has long been available for functions, so it's not surprising that macros were finally put on an equal footing.

The primary reason for having a macro with a variable number of arguments is that it can pass these arguments to a function that accepts a variable number of arguments, such as `printf` or `scanf`. Here's an example:

```
#define TEST(condition, ...) ((condition)? \
    printf("Passed test: %s\n", #condition): \
    printf(__VA_ARGS__))
```

The ... token, known as *ellipsis*, goes at the end of a macro's parameter list, preceded by ordinary parameters, if there are any. \_\_VA\_ARGS\_\_ is a special identifier that can appear only in the replacement list of a macro with a variable number of arguments; it represents all the arguments that correspond to the ellipsis. (There must be at least one argument that corresponds to the ellipsis, although that argument may be empty.) The TEST macro requires at least two arguments. The first argument matches the condition parameter; the remaining arguments match the ellipsis.

Here's an example that shows how the TEST macro might be used:

```
TEST(voltage <= max_voltage,
    "Voltage %d exceeds %d\n", voltage, max_voltage);
```

The preprocessor will produce the following output (reformatted for readability):

```
((voltage <= max_voltage) ?
    printf("Passed test: %s\n", "voltage <= max_voltage") :
    printf("Voltage %d exceeds %d\n", voltage, max_voltage));
```

When the program is executed, the program will display the message

```
Passed test: voltage <= max_voltage
```

if `voltage` is no more than `max_voltage`. Otherwise, it will display the values of `voltage` and `max_voltage`:

```
Voltage 125 exceeds 120
```

### C99 The `_func_` Identifier

Another new feature of C99 is the `_func_` identifier. `_func_` has nothing to do with the preprocessor, so it actually doesn't belong in this chapter. However, like many preprocessor features, it's useful for debugging, so I've chosen to discuss it here.

Every function has access to the `_func_` identifier, which behaves like a string variable that stores the name of the currently executing function. The effect is the same as if each function contains the following declaration at the beginning of its body:

```
static const char _func_[] = "function-name";
```

where `function-name` is the name of the function. The existence of this identifier makes it possible to write debugging macros such as the following:

```
#define FUNCTION_CALLED() printf("%s called\n", __func__);
#define FUNCTION_RETURNS() printf("%s returns\n", __func__);
```

Calls of these macros can then be placed inside functions to trace their calls:

```
void f(void)
{
    FUNCTION_CALLED(); /* displays "f called" */
    ...
    FUNCTION_RETURNS(); /* displays "f returns" */
}
```

Another use of `_func_`: it can be passed to a function to let it know the name of the function that called it.

## 14.4 Conditional Compilation

The C preprocessor recognizes a number of directives that support ***conditional compilation***—the inclusion or exclusion of a section of program text depending on the outcome of a test performed by the preprocessor.

## The `#if` and `#endif` Directives

Suppose we're in the process of debugging a program. We'd like the program to print the values of certain variables, so we put calls of `printf` in critical parts of the program. Once we've located the bugs, it's often a good idea to let the `printf` calls remain, just in case we need them later. Conditional compilation allows us to leave the calls in place, but have the compiler ignore them.

Here's how we'll proceed. We'll first define a macro and give it a nonzero value:

```
#define DEBUG 1
```

The name of the macro doesn't matter. Next, we'll surround each group of `printf` calls by an `#if`-`#endif` pair:

```
#if DEBUG
printf("Value of i: %d\n", i);
printf("Value of j: %d\n", j);
#endif
```

During preprocessing, the `#if` directive will test the value of `DEBUG`. Since its value isn't zero, the preprocessor will leave the two calls of `printf` in the program (the `#if` and `#endif` lines will disappear, though). If we change the value of `DEBUG` to zero and recompile the program, the preprocessor will remove all four lines from the program. The compiler won't see the calls of `printf`, so they won't occupy any space in the object code and won't cost any time when the program is run. We can leave the `#if`-`#endif` blocks in the final program, allowing diagnostic information to be produced later (by recompiling with `DEBUG` set to 1) if any problems turn up.

In general, the `#if` directive has the form

**#if directive**

`#if constant-expression`

The `#endif` directive is even simpler:

**#endif directive**

`#endif`

**Q&A** When the preprocessor encounters the `#if` directive, it evaluates the constant expression. If the value of the expression is zero, the lines between `#if` and `#endif` will be removed from the program during preprocessing. Otherwise, the lines between `#if` and `#endif` will remain in the program to be processed by the compiler—the `#if` and `#endif` will have had no effect on the program.

It's worth noting that the `#if` directive treats undefined identifiers as macros that have the value 0. Thus, if we neglect to define `DEBUG`, the test

```
#if DEBUG
```

will fail (but not generate an error message), while the test

```
#if !DEBUG
```

will succeed.

## The `defined` Operator

We encountered the `#` and `##` operators in Section 14.3. There's just one other operator, `defined`, that's specific to the preprocessor. When applied to an identifier, `defined` produces the value 1 if the identifier is a currently defined macro; it produces 0 otherwise. The `defined` operator is normally used in conjunction with the `#if` directive; it allows us to write

```
#if defined(DEBUG)  
...  
#endif
```

The lines between the `#if` and `#endif` directives will be included in the program only if `DEBUG` is defined as a macro. The parentheses around `DEBUG` aren't required; we could simply write

```
#if defined DEBUG
```

Since `defined` tests only whether `DEBUG` is defined or not, it's not necessary to give `DEBUG` a value:

```
#define DEBUG
```

## The `#ifdef` and `#ifndef` Directives

The `#ifdef` directive tests whether an identifier is currently defined as a macro:

**#ifdef directive**

`#ifdef identifier`

Using `#ifdef` is similar to using `#if`:

```
#ifdef identifier
```

*Lines to be included if identifier is defined as a macro*

```
#endif
```

**Q&A**

Strictly speaking, there's no need for `#ifdef`, since we can combine the `#if` directive with the `defined` operator to get the same effect. In other words, the directive

```
#ifdef identifier
```

is equivalent to

```
#if defined(identifier)
```

The `#ifndef` directive is similar to `#ifdef`, but tests whether an identifier is *not* defined as a macro:

**#ifndef directive****#ifndef identifier**

Writing

`#ifndef identifier`

is the same as writing

`#if !defined(identifier)`

## The `#elif` and `#else` Directives

`#if`, `#ifdef`, and `#ifndef` blocks can be nested just like ordinary `if` statements. When nesting occurs, it's a good idea to use an increasing amount of indentation as the level of nesting grows. Some programmers put a comment on each closing `#endif` to indicate what condition the matching `#if` tests:

```
#if DEBUG  
...  
#endif /* DEBUG */
```

This technique makes it easier for the reader to find the beginning of the `#if` block.

For additional convenience, the preprocessor supports the `#elif` and `#else` directives:

**#elif directive****#elif constant-expression****#else directive****#else**

`#elif` and `#else` can be used in conjunction with `#if`, `#ifdef`, or `#ifndef` to test a series of conditions:

```
#if expr1  
Lines to be included if expr1 is nonzero  
#elif expr2  
Lines to be included if expr1 is zero but expr2 is nonzero  
#else  
Lines to be included otherwise  
#endif
```

Although the `#if` directive is shown above, an `#ifdef` or `#ifndef` directive can be used instead. Any number of `#elif` directives—but at most one `#else`—may appear between `#if` and `#endif`.

## Uses of Conditional Compilation

Conditional compilation is certainly handy for debugging, but its uses don't stop there. Here are a few other common applications:

- ***Writing programs that are portable to several machines or operating systems.*** The following example includes one of three groups of lines depending on whether WIN32, MAC\_OS, or LINUX is defined as a macro:

```
#if defined(WIN32)
...
#elif defined(MAC_OS)
...
#elif defined(LINUX)
...
#endif
```

A program might contain many of these `#if` blocks. At the beginning of the program, one (and only one) of the macros will be defined, thereby selecting a particular operating system. For example, defining the `LINUX` macro might indicate that the program is to run under the Linux operating system.

- ***Writing programs that can be compiled with different compilers.*** Different compilers often recognize somewhat different versions of C. Some accept a standard version of C, some don't. Some provide machine-specific language extensions; some don't, or provide a different set of extensions. Conditional compilation can allow a program to adjust to different compilers. Consider the problem of writing a program that might have to be compiled using an older, nonstandard compiler. The `_STDC_` macro allows the preprocessor to detect whether a compiler conforms to the standard (either C89 or C99); if it doesn't, we may need to change certain aspects of the program. In particular, we may have to use old-style function declarations (discussed in the Q&A at the end of Chapter 9) instead of function prototypes. At each point where functions are declared, we can put the following lines:

```
#if __STDC__
Function prototypes
#else
Old-style function declarations
#endif
```

- ***Providing a default definition for a macro.*** Conditional compilation allows us to check whether a macro is currently defined and, if not, give it a default definition. For example, the following lines will define the macro `BUFFER_SIZE` if it wasn't previously defined:

```
#ifndef BUFFER_SIZE
#define BUFFER_SIZE 256
#endif
```

- **Temporarily disabling code that contains comments.** We can't use a `/*...*/` comment to "comment out" code that already contains `/*...*/` comments. Instead, we can use an `#if` directive:

```
#if 0
Lines containing comments
#endif
```

**Q&A**

Disabling code in this way is often called "conditioning out."

Section 15.2 discusses another common use of conditional compilation: protecting header files against multiple inclusion.

## 14.5 Miscellaneous Directives

To end the chapter, we'll take a brief look at the `#error`, `#line`, and `#pragma` directives. These directives are more specialized than the ones we've already examined, and they're used much less frequently.

### The `#error` Directive

The `#error` directive has the form

**#error directive**

**#error message**

where *message* is any sequence of tokens. If the preprocessor encounters an `#error` directive, it prints an error message which must include *message*. The exact form of the error message can vary from one compiler to another; it might be something like

Error directive: *message*

or perhaps just

`#error message`

Encountering an `#error` directive indicates a serious flaw in the program; some compilers immediately terminate compilation without attempting to find other errors.

`#error` directives are frequently used in conjunction with conditional compilation to check for situations that shouldn't arise during a normal compilation. For example, suppose that we want to ensure that a program can't be compiled on a machine whose `int` type isn't capable of storing numbers up to 100,000. The largest possible `int` value is represented by the `INT_MAX` macro, so all we need do is invoke an `#error` directive if `INT_MAX` isn't at least 100,000:

```
#if INT_MAX < 100000
#error int type is too small
#endif
```

Attempting to compile the program on a machine whose integers are stored in 16 bits will produce a message such as

```
Error directive: int type is too small
```

The `#error` directive is often found in the `#else` part of an `#if-#elif-#else` series:

```
#if defined(WIN32)
...
#elif defined(MAC_OS)
...
#elif defined(LINUX)
...
#else
#error No operating system specified
#endif
```

## The `#line` Directive

The `#line` directive is used to alter the way program lines are numbered. (Lines are usually numbered 1, 2, 3, as you'd expect.) We can also use this directive to make the compiler think that it's reading the program from a file with a different name.

The `#line` directive has two forms. In one form, we specify a line number:

**#line directive  
(form 1)**

**#line *n***

**C99**

*n* must be a sequence of digits representing an integer between 1 and 32767 (2147483647 in C99). This directive causes subsequent lines in the program to be numbered *n*, *n* + 1, *n* + 2, and so forth.

In the second form of the `#line` directive, both a line number and a file name are specified:

**#line directive  
(form 2)**

**#line *n* "file"**

The lines that follow this directive are assumed to come from *file*, with line numbers starting at *n*. The values of *n* and/or the *file* string can be specified using macros.

One effect of the `#line` directive is to change the value of the `_LINE_` macro (and possibly the `_FILE_` macro). More importantly, most compilers will use the information from the `#line` directive when generating error messages.

For example, suppose that the following directive appears at the beginning of the file `foo.c`:

```
#line 10 "bar.c"
```

Let's say that the compiler detects an error on line 5 of `foo.c`. The error message will refer to line 13 of file `bar.c`, not line 5 of file `foo.c`. (Why line 13? The directive occupies line 1 of `foo.c`, so the renumbering of `foo.c` begins at line 2, which is treated as line 10 of `bar.c`.)

At first glance, the `#line` directive is mystifying. Why would we want error messages to refer to a different line and possibly a different file? Wouldn't this make programs harder to debug?

In fact, the `#line` directive isn't used very often by programmers. Instead, it's used primarily by programs that generate C code as output. The most famous example of such a program is `yacc` (Yet Another Compiler-Compiler), a UNIX utility that automatically generates part of a compiler. (The GNU version of `yacc` is named `bison`.) Before using `yacc`, the programmer prepares a file that contains information for `yacc` as well as fragments of C code. From this file, `yacc` generates a C program, `y.tab.c`, that incorporates the code supplied by the programmer. The programmer then compiles `y.tab.c` in the usual way. By inserting `#line` directives in `y.tab.c`, `yacc` tricks the compiler into believing that the code comes from the original file—the one written by the programmer. As a result, any error messages produced during the compilation of `y.tab.c` will refer to lines in the original file, not lines in `y.tab.c`. This makes debugging easier, because error messages refer to the file written by the programmer, not the (more complicated) file generated by `yacc`.

## The `#pragma` Directive

The `#pragma` directive provides a way to request special behavior from the compiler. This directive is most useful for programs that are unusually large or that need to take advantage of the capabilities of a particular compiler.

The `#pragma` directive has the form

`#pragma directive`

`#pragma tokens`

where *tokens* are arbitrary tokens. `#pragma` directives can be very simple (a single token) or they can be much more elaborate:

```
#pragma data(heap_size => 1000, stack_size => 2000)
```

Not surprisingly, the set of commands that can appear in `#pragma` directives is different for each compiler; you'll have to consult the documentation for your compiler to see which commands it allows and what those commands do. Incidentally, the preprocessor must ignore any `#pragma` directive that contains an unrecognized command; it's not permitted to give an error message.

- In C89, there are no standard pragmas—they’re all implementation-defined.
- C99** C99 has three standard pragmas, all of which use STDC as the first token following `#pragma`. These pragmas are `FP_CONTRACT` (covered in Section 23.4), `CX_LIMITED_RANGE` (Section 27.4), and `FENV_ACCESS` (Section 27.6).

### **C99** The `_Pragma` Operator

C99 introduces the `_Pragma` operator, which is used in conjunction with the `#pragma` directive. A `_Pragma` expression has the form

`_Pragma expression`

`_Pragma ( string-literal )`

When it encounters such an expression, the preprocessor “destringizes” the string literal (yes, that’s the term used in the C99 standard!) by removing the double quotes around the string and replacing the escape sequences `\"` and `\\` by the characters `"` and `\`, respectively. The result is a series of tokens, which are then treated as though they appear in a `#pragma` directive. For example, writing

```
_Pragma ("data(heap_size => 1000, stack_size => 2000)")
```

is the same as writing

```
#pragma data(heap_size => 1000, stack_size => 2000)
```

The `_Pragma` operator lets us work around a limitation of the preprocessor: the fact that a preprocessing directive can’t generate another directive. `_Pragma`, however, is an operator, not a directive, and can therefore appear in a macro definition. This makes it possible for a macro expansion to leave behind a `#pragma` directive.

Let’s look at an example from the GCC manual. The following macro uses the `_Pragma` operator:

```
#define DO_PRAGMA(x) _Pragma(#x)
```

The macro would be invoked as follows:

```
DO_PRAGMA(GCC dependency "parse.y")
```

After expansion, the result will be

```
#pragma GCC dependency "parse.y"
```

which is one of the pragmas supported by GCC. (It issues a warning if the date of the specified file—`parse.y` in this example—is more recent than the date of the current file—the one being compiled.) Note that the argument to the call of `DO_PRAGMA` is a series of tokens. The `#` operator in the definition of `DO_PRAGMA` causes the tokens to be stringized into `"GCC dependency \"parse.y\"";` this string is then passed to the `_Pragma` operator, which destringizes it, producing a `#pragma` directive containing the original tokens.

## Q & A

**Q:** I've seen programs that contain a # on a line by itself. Is this legal?

A: Yes. This is the *null directive*; it has no effect. Some programmers use null directives for spacing within conditional compilation blocks:

```
#if INT_MAX < 100000
#
#error int type is too small
#
#endif
```

Blank lines would also work, of course, but the # helps the reader see the extent of the block.

**Q:** I'm not sure which constants in a program need to be defined as macros. Are there any guidelines to follow? [p. 319]

A: One rule of thumb says that every numeric constant, other than 0 or 1, should be a macro. Character and string constants are problematic, since replacing a character or string constant by a macro doesn't always improve readability. I recommend using a macro instead of a character constant or string literal provided that (1) the constant is used more than once and (2) the possibility exists that the constant might someday be modified. Because of rule (2), I don't use macros such as

```
#define NUL '\0'
```

although some programmers do.

**Q:** What does the # operator do if the argument that it's supposed to "stringize" contains a " or \ character? [p. 324]

A: It converts " to \" and \ to \\. Consider the following macro:

```
#define STRINGIZE(x) #x
```

The preprocessor will replace STRINGIZE ("foo") by "\\\"foo\\\"".

**\*Q:** I can't get the following macro to work properly:

```
#define CONCAT(x,y) x##y
```

**CONCAT(a,b)** gives ab, as expected, but **CONCAT(a,CONCAT(b,c))** gives an odd result. What's going on?

A: Thanks to rules that Kernighan and Ritchie call "bizarre," macros whose replacement lists depend on ## usually can't be called in a nested fashion. The problem is that **CONCAT(a,CONCAT(b,c))** isn't expanded in a "normal" fashion, with **CONCAT(b,c)** yielding bc, then **CONCAT(a,bc)** giving abc. Macro parameters that are preceded or followed by ## in a replacement list aren't expanded at

the time of substitution. As a result, `CONCAT(a, CONCAT(b, c))` expands to `aCONCAT(b, c)`, which can't be expanded further, since there's no macro named `aCONCAT`.

There's a way to solve the problem, but it's not pretty. The trick is to define a second macro that simply calls the first one:

```
#define CONCAT2(x,y) CONCAT(x,y)
```

Writing `CONCAT2(a, CONCAT2(b, c))` now yields the desired result. As the preprocessor expands the outer call of `CONCAT2`, it will expand `CONCAT2(b, c)` as well; the difference is that `CONCAT2`'s replacement list doesn't contain `##`. If none of this makes any sense, don't worry; it's not a problem that arises often.

The `#` operator has a similar difficulty, by the way. If `#x` appears in a replacement list, where `x` is a macro parameter, the corresponding argument is not expanded. Thus, if `N` is a macro representing `10`, and `STR(x)` has the replacement list `#x`, expanding `STR(N)` yields `"N"`, not `"10"`. The solution is similar to the one we used with `CONCAT`: defining a second macro whose job is to call `STR`.

**\*Q:** Suppose that the preprocessor encounters the original macro name during rescanning, as in the following example:

```
#define N (2*M)
#define M (N+1)

i = N; /* infinite loop? */
```

**The preprocessor will replace `N` by `(2*M)`, then replace `M` by `(N+1)`. Will the preprocessor replace `N` again, thus going into an infinite loop? [p. 326]**

**A:** Some old preprocessors will indeed go into an infinite loop, but newer ones shouldn't. According to the C standard, if the original macro name reappears during the expansion of a macro, the name is not replaced again. Here's how the assignment to `i` will look after preprocessing:

```
i = (2*(N+1));
```

Some enterprising programmers take advantage of this behavior by writing macros whose names match reserved words or functions in the standard library. Consider the `sqrt` library function. `sqrt` computes the square root of its argument, returning an implementation-defined value if the argument is negative. Perhaps we would prefer that `sqrt` return `0` if its argument is negative. Since `sqrt` is part of the standard library, we can't easily change it. We can, however, define a `sqrt` *macro* that evaluates to `0` when given a negative argument:

```
#undef sqrt
#define sqrt(x) ((x)>=0?sqrt(x):0)
```

A later call of `sqrt` will be intercepted by the preprocessor, which expands it into the conditional expression shown here. The call of `sqrt` inside the conditional expression won't be replaced during rescanning, so it will remain for the compiler

to handle. (Note the use of `#undef` to undefine `sqrt` before defining the `sqrt` macro. As we'll see in Section 21.1, the standard library is allowed to have both a macro and a function with the same name. Undefining `sqrt` before defining our own `sqrt` macro is a defensive measure, in case the library has already defined `sqrt` as a macro.)

**Q:** I get an error when I try to use predefined macros such as `_LINE_` and `_FILE_`. Is there a special header that I need to include?

**A:** No. These macros are recognized automatically by the preprocessor. Make sure that you have *two* underscores at the beginning and end of each macro name, not one.

**Q:** What's the purpose of distinguishing between a "hosted implementation" and a "freestanding implementation"? If a freestanding implementation doesn't even support the `<stdio.h>` header, what use is it? [p. 330]

**A:** A hosted implementation is needed for most programs (including the ones in this book), which rely on the underlying operating system for input/output and other essential services. A freestanding implementation of C would be used for programs that require no operating system (or only a minimal operating system). For example, a freestanding implementation would be needed for writing the kernel of an operating system (which requires no traditional input/output and therefore doesn't need `<stdio.h>` anyway). Freestanding implementations are also useful for writing software for embedded systems.

**Q:** I thought the preprocessor was just an editor. How can it evaluate constant expressions? [p. 334]

**A:** The preprocessor is more sophisticated than you might expect; it knows enough about C to be able to evaluate constant expressions, although it doesn't do so in quite the same way as the compiler. (For one thing, the preprocessor treats any undefined name as having the value 0. The other differences are too esoteric to go into here.) In practice, the operands in a preprocessor constant expression are usually constants, macros that represent constants, and applications of the `defined` operator.

**Q:** Why does C provide the `#ifdef` and `#ifndef` directives, since we can get the same effect using the `#if` directive and the `defined` operator? [p. 335]

**A:** The `#ifdef` and `#ifndef` directives have been a part of C since the 1970s. The `defined` operator, on the other hand, was added to C in the 1980s during standardization. So the real question is: Why was `defined` added to the language? The answer is that `defined` adds flexibility. Instead of just being able to test the existence of a single macro using `#ifdef` or `#ifndef`, we can now test any number of macros using `#if` together with `defined`. For example, the following directive checks whether `FOO` and `BAR` are defined but `BAZ` is not defined:

```
#if defined(FOO) && defined(BAR) && !defined(BAZ)
```

- Q:** I wanted to compile a program that I hadn't finished writing, so I "conditioned out" the unfinished part:

```
#if 0
...
#endif
```

When I compiled the program, I got an error message referring to one of the lines between `#if` and `#endif`. Doesn't the preprocessor just ignore these lines? [p. 338]

- A:** No, the lines aren't completely ignored. Comments are processed before preprocessing directives are executed, and the source code is divided into preprocessing tokens. Thus, an unterminated comment between `#if` and `#endif` may cause an error message. Also, an unpaired single quote or double quote character may cause undefined behavior.

## Exercises

### Section 14.3

1. Write parameterized macros that compute the following values.
  - The cube of  $x$ .
  - The remainder when  $n$  is divided by 4.
  - 1 if the product of  $x$  and  $y$  is less than 100, 0 otherwise.

Do your macros always work? If not, describe what arguments would make them fail.
2. Write a macro `NELEMS(a)` that computes the number of elements in a one-dimensional array `a`. *Hint:* See the discussion of the `sizeof` operator in Section 8.1.
3. Let `DOUBLE` be the following macro:
 

```
#define DOUBLE(x) 2*x
```

  - What is the value of `DOUBLE(1+2)`?
  - What is the value of `4/DOMBLE(2)`?
  - Fix the definition of `DOUBLE`.
4. For each of the following macros, give an example that illustrates a problem with the macro and show how to fix it.
  - `#define AVG(x,y) (x+y)/2`
  - `#define AREA(x,y) (x)*(y)`
5. Let `TOUPPER` be the following macro:
 

```
#define TOUPPER(c) ('a' <= (c) && (c) <= 'z' ? (c) - 'a' + 'A' : (c))
```

Let `s` be a string and let `i` be an `int` variable. Show the output produced by each of the following program fragments.

  - `strcpy(s, "abcd");  
i = 0;  
putchar(TOUPPER(s[+i]));`

- (b) `strcpy(s, "0123");  
i = 0;  
putchar(TOUPPER(s[+i]));`
6. (a) Write a macro `DISP(f, x)` that expands into a call of `printf` that displays the value of the function `f` when called with argument `x`. For example,  
`DISP(sqrt, 3.0);`  
should expand into  
`printf("sqrt(%g) = %g\n", 3.0, sqrt(3.0));`
- (b) Write a macro `DISP2(f, x, y)` that's similar to `DISP` but works for functions with two arguments.
- W \*7. Let `GENERIC_MAX` be the following macro:
- ```
#define GENERIC_MAX(type) \
type type##_max(type x, type y) \
{ \
    return x > y ? x : y; \
}
```
- (a) Show the preprocessor's expansion of `GENERIC_MAX(long)`.  
(b) Explain why `GENERIC_MAX` doesn't work for basic types such as `unsigned long`.  
(c) Describe a technique that would allow us to use `GENERIC_MAX` with basic types such as `unsigned long`. Hint: Don't change the definition of `GENERIC_MAX`.
- \*8. Suppose we want a macro that expands into a string containing the current line number and file name. In other words, we'd like to write  
`const char *str = LINE_FILE;`  
and have it expand into  
`const char *str = "Line 10 of file foo.c";`  
where `foo.c` is the file containing the program and 10 is the line on which the invocation of `LINE_FILE` appears. Warning: This exercise is for experts only. Be sure to read the Q&A section carefully before attempting!
9. Write the following parameterized macros.
- (a) `CHECK(x, y, n)` – Has the value 1 if both `x` and `y` fall between 0 and `n - 1`, inclusive.  
(b) `MEDIAN(x, y, z)` – Finds the median of `x`, `y`, and `z`.  
(c) `POLYNOMIAL(x)` – Computes the polynomial  $3x^5 + 2x^4 - 5x^3 - x^2 + 7x - 6$ .
10. Functions can often—but not always—be written as parameterized macros. Discuss what characteristics of a function would make it unsuitable as a macro.
11. (C99) C programmers often use the `fprintf` function to write error messages:  
`fprintf(stderr, "Range error: index = %d\n", index);`  
`stderr` stream ►22.1 `stderr` is C's “standard error” stream; the remaining arguments are the same as those for `printf`, starting with the format string. Write a macro named `ERROR` that generates the call of `fprintf` shown above when given a format string and the items to be displayed:  
`ERROR("Range error: index = %d\n", index);`
- Section 14.4 W 12. Suppose that the macro `M` has been defined as follows:  
`#define M 10`

Which of the following tests will fail?

- (a) #if M
- (b) #ifdef M
- (c) #ifndef M
- (d) #if defined(M)
- (e) #if !defined(M)

13. (a) Show what the following program will look like after preprocessing. You may ignore any lines added to the program as a result of including the `<stdio.h>` header.

```
#include <stdio.h>

#define N 100

void f(void);

int main(void)
{
    f();
#ifndef N
#define N
#endif
    return 0;
}

void f(void)
{
#ifndef defined(N)
    printf("N is %d\n", N);
#else
    printf("N is undefined\n");
#endif
}
```

- (b) What will be the output of this program?

- W\*14. Show what the following program will look like after preprocessing. Some lines of the program may cause compilation errors; find all such errors.

```
#define N = 10
#define INC(x) x+1
#define SUB (x,y) x-y
#define SQR(x) ((x)*(x))
#define CUBE(x) (SQR(x)*(x))
#define M1(x,y) x##y
#define M2(x,y) #x #y

int main(void)
{
    int a[N], i, j, k, m;

#ifndef N
    i = j;
#else
    j = i;
#endif

    i = 10 * INC(j);
```

```
i = SUB(j, k);
i = SQR(SQR(j));
i = CUBE(j);
i = M1(j, k);
puts(M2(i, j));

#undef SQR
i = SQR(j);
#define SQR
i = SQR(j);

return 0;
}
```

15. Suppose that a program needs to display messages in either English, French, or Spanish. Using conditional compilation, write a program fragment that displays one of the following three messages, depending on whether or not the specified macro is defined:

Insert Disk 1 (if ENGLISH is defined)  
Inserez Le Disque 1 (if FRENCH is defined)  
Inserte El Disco 1 (if SPANISH is defined)

#### Section 14.5

- \*16. (C99) Assume that the following macro definitions are in effect:

```
#define IDENT(x) PRAGMA(ident #x)
#define PRAGMA(x) _Pragma(#x)
```

What will the following line look like after macro expansion?

```
IDENT(foo)
```

# 15 Writing Large Programs

*Around computers it is difficult to find the correct unit of time to measure progress. Some cathedrals took a century to complete. Can you imagine the grandeur and scope of a program that would take as long?*

Although some C programs are small enough to be put in a single file, most aren't. Programs that consist of more than one file are the rule rather than the exception. In this chapter, we'll see that a typical program consists of several source files and usually some header files as well. Source files contain definitions of functions and external variables; header files contain information to be shared among source files. Section 15.1 discusses source files, while Section 15.2 covers header files. Section 15.3 describes how to divide a program into source files and header files. Section 15.4 then shows how to "build" (compile and link) a program that consists of more than one file, and how to "rebuild" a program after part of it has been changed.

## 15.1 Source Files

Up to this point, we've assumed that a C program consists of a single file. In fact, a program may be divided among any number of *source files*. By convention, source files have the extension .c. Each source file contains part of the program, primarily definitions of functions and variables. One source file must contain a function named `main`, which serves as the starting point for the program.

For example, suppose that we want to write a simple calculator program that evaluates integer expressions entered in Reverse Polish notation (RPN), in which operators follow operands. If the user enters an expression such as

30 5 - 7 \*

we want the program to print its value (175, in this case). Evaluating an RPN expression is easy if we have the program read the operands and operators, one by one, using a stack to keep track of intermediate results. If the program reads a

number, we'll have it push the number onto the stack. If it reads an operator, we'll have it pop two numbers from the stack, perform the operation, and then push the result back onto the stack. When the program reaches the end of the user's input, the value of the expression will be on the stack. For example, the program will evaluate the expression `30 5 - 7 *` in the following way:

1. Push 30 onto the stack.
2. Push 5 onto the stack.
3. Pop the top two numbers from the stack, subtract 5 from 30, giving 25, and then push the result back onto the stack.
4. Push 7 onto the stack.
5. Pop the top two numbers from the stack, multiply them, and then push the result back onto the stack.

After these steps, the stack will contain the value of the expression (175).

Turning this strategy into a program isn't hard. The program's main function will contain a loop that performs the following actions:

- Read a "token" (a number or an operator).
- If the token is a number, push it onto the stack.
- If the token is an operator, pop its operands from the stack, perform the operation, and then push the result back onto the stack.

When dividing a program like this one into files, it makes sense to put related functions and variables into the same file. The function that reads tokens could go into one source file (`token.c`, say), together with any functions that have to do with tokens. Stack-related functions such as `push`, `pop`, `make_empty`, `is_empty`, and `is_full` could go into a different file, `stack.c`. The variables that represent the stack would also go into `stack.c`. The main function would go into yet another file, `calc.c`.

Splitting a program into multiple source files has significant advantages:

- Grouping related functions and variables into a single file helps clarify the structure of the program.
- Each source file can be compiled separately—a great time-saver if the program is large and must be changed frequently (which is common during program development).
- Functions are more easily reused in other programs when grouped in separate source files. In our example, splitting off `stack.c` and `token.c` from the main function makes it simpler to reuse the stack functions and token functions in the future.

## 15.2 Header Files

When we divide a program into several source files, problems arise: How can a function in one file call a function that's defined in another file? How can a func-

tion access an external variable in another file? How can two files share the same macro definition or type definition? The answer lies with the `#include` directive, which makes it possible to share information—function prototypes, macro definitions, type definitions, and more—among any number of source files.

The `#include` directive tells the preprocessor to open a specified file and insert its contents into the current file. Thus, if we want several source files to have access to the same information, we'll put that information in a file and then use `#include` to bring the file's contents into each of the source files. Files that are included in this fashion are called *header files* (or sometimes *include files*); I'll discuss them in more detail later in this section. By convention, header files have the extension `.h`.

*Note:* The C standard uses the term “source file” to refer to all files written by the programmer, including both `.c` and `.h` files. I'll use “source file” to refer to `.c` files only.

## The `#include` Directive

The `#include` directive has two primary forms. The first form is used for header files that belong to C's own library:

**#include directive  
(form 1)**

`#include <filename>`

The second form is used for all other header files, including any that we write:

**#include directive  
(form 2)**

`#include "filename"`

**Q&A**

The difference between the two is a subtle one having to do with how the compiler locates the header file. Here are the rules that most compilers follow:

- `#include <filename>`: Search the directory (or directories) in which system header files reside. (On UNIX systems, for example, system header files are usually kept in the directory `/usr/include`.)
- `#include "filename"`: Search the current directory, then search the directory (or directories) in which system header files reside.

The places to be searched for header files can usually be altered, often by a command-line option such as `-Ipath`.



Don't use brackets when including header files that you have written:

```
#include <myheader.h>    /*** WRONG ***/
```

The preprocessor will probably look for `myheader.h` where the system header files are kept (and, of course, won't find it).

The file name in an `#include` directive may include information that helps locate the file, such as a directory path or drive specifier:

```
#include "c:\cprogs\utils.h" /* Windows path */
```

```
#include "/cprogs/utils.h" /* UNIX path */
```

Although the quotation marks in the `#include` directive make file names look like string literals, the preprocessor doesn't treat them that way. (That's fortunate, since \c and \u—which appear in the Windows example—would be treated as escape sequences in a string literal.)

#### portability tip

*It's usually best not to include path or drive information in `#include` directives. Such information makes it difficult to compile a program when it's transported to another machine or, worse, another operating system.*

For example, the following Windows `#include` directives specify drive and/or path information that may not always be valid:

```
#include "d:utils.h"
#include "\cprogs\include\utils.h"
#include "d:\cprogs\include\utils.h"
```

The following directives are better; they don't mention specific drives, and paths are relative rather than absolute:

```
#include "utils.h"
#include "..\include\utils.h"
```

The `#include` directive has a third form that's used less often than the other two:

#### #include directive (form 3)

preprocessing tokens ▶ 14.3

#### #include *tokens*

where *tokens* is any sequence of preprocessing tokens. The preprocessor will scan the tokens and replace any macros that it finds. After macro replacement, the resulting directive must match one of the other forms of `#include`. The advantage of the third kind of `#include` is that the file name can be defined by a macro rather than being “hard-coded” into the directive itself, as the following example shows:

```
#if defined(IA32)
#define CPU_FILE "ia32.h"
#elif defined(IA64)
#define CPU_FILE "ia64.h"
#elif defined(AMD64)
#define CPU_FILE "amd64.h"
#endif

#include CPU_FILE
```

## Sharing Macro Definitions and Type Definitions

Most large programs contain macro definitions and type definitions that need to be shared by several source files (or, in the most extreme case, by *all* source files). These definitions should go into header files.

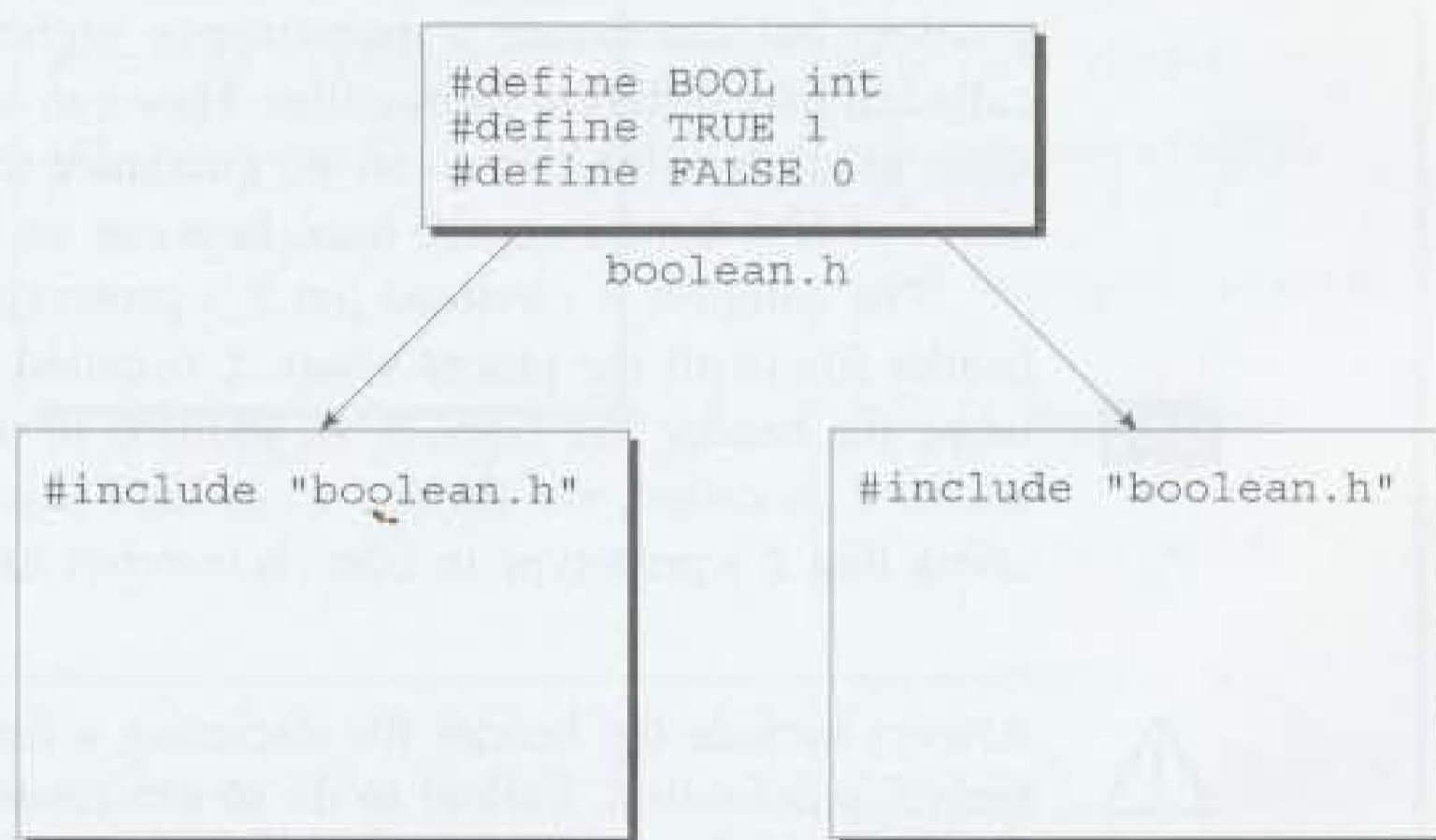
For example, suppose that we're writing a program that uses macros named `BOOL`, `TRUE`, and `FALSE`. (There's no need for these in C99, of course, because the `<stdbool.h>` header defines similar macros.) Instead of repeating the definitions of these macros in each source file that needs them, it makes more sense to put the definitions in a header file with a name like `boolean.h`:

```
#define BOOL int
#define TRUE 1
#define FALSE 0
```

Any source file that requires these macros will simply contain the line

```
#include "boolean.h"
```

In the following figure, two files include `boolean.h`:



Type definitions are also common in header files. For example, instead of defining a `BOOL` macro, we might use `typedef` to create a `Bool` type. If we do, the `boolean.h` file will have the following appearance:

```
#define TRUE 1
#define FALSE 0
typedef int Bool;
```

Putting definitions of macros and types in header files has some clear advantages. First, we save time by not having to copy the definitions into the source files where they're needed. Second, the program becomes easier to modify. Changing the definition of a macro or type requires only that we edit a single header file; we don't have to modify the many source files in which the macro or type is used. Third, we don't have to worry about inconsistencies caused by source files containing different definitions of the same macro or type.

## Sharing Function Prototypes

default argument promotions ➤ 9.3

Suppose that a source file contains a call of a function *f* that's defined in another file, *foo.c*. Calling *f* without declaring it first is risky. Without a prototype to rely on, the compiler is forced to assume that *f*'s return type is *int* and that the number of parameters matches the number of arguments in the call of *f*. The arguments themselves are converted automatically to a kind of "standard form" by the default argument promotions. The compiler's assumptions may well be wrong, but it has no way to check them, since it compiles only one file at a time. If the assumptions are incorrect, the program probably won't work, and there won't be any clues as to why it doesn't. (For this reason, C99 prohibits calling a function for which the compiler has not yet seen a declaration or definition.)



When calling a function *f* that's defined in another file, always make sure that the compiler has seen a prototype for *f* prior to the call.

Our first impulse is to declare *f* in the file where it's called. That solves the problem but can create a maintenance nightmare. Suppose that the function is called in fifty different source files. How can we ensure that *f*'s prototypes are the same in all the files? How can we guarantee that they match the definition of *f* in *foo.c*? If *f* should change later, how can we find all the files where it's used?

The solution is obvious: put *f*'s prototype in a header file, then include the header file in all the places where *f* is called. Since *f* is defined in *foo.c*, let's name the header file *foo.h*. In addition to including *foo.h* in the source files where *f* is called, we'll need to include it in *foo.c*, enabling the compiler to check that *f*'s prototype in *foo.h* matches its definition in *foo.c*.



Always include the header file declaring a function *f* in the source file that contains *f*'s definition. Failure to do so can cause hard-to-find bugs, since calls of *f* elsewhere in the program may not match *f*'s definition.

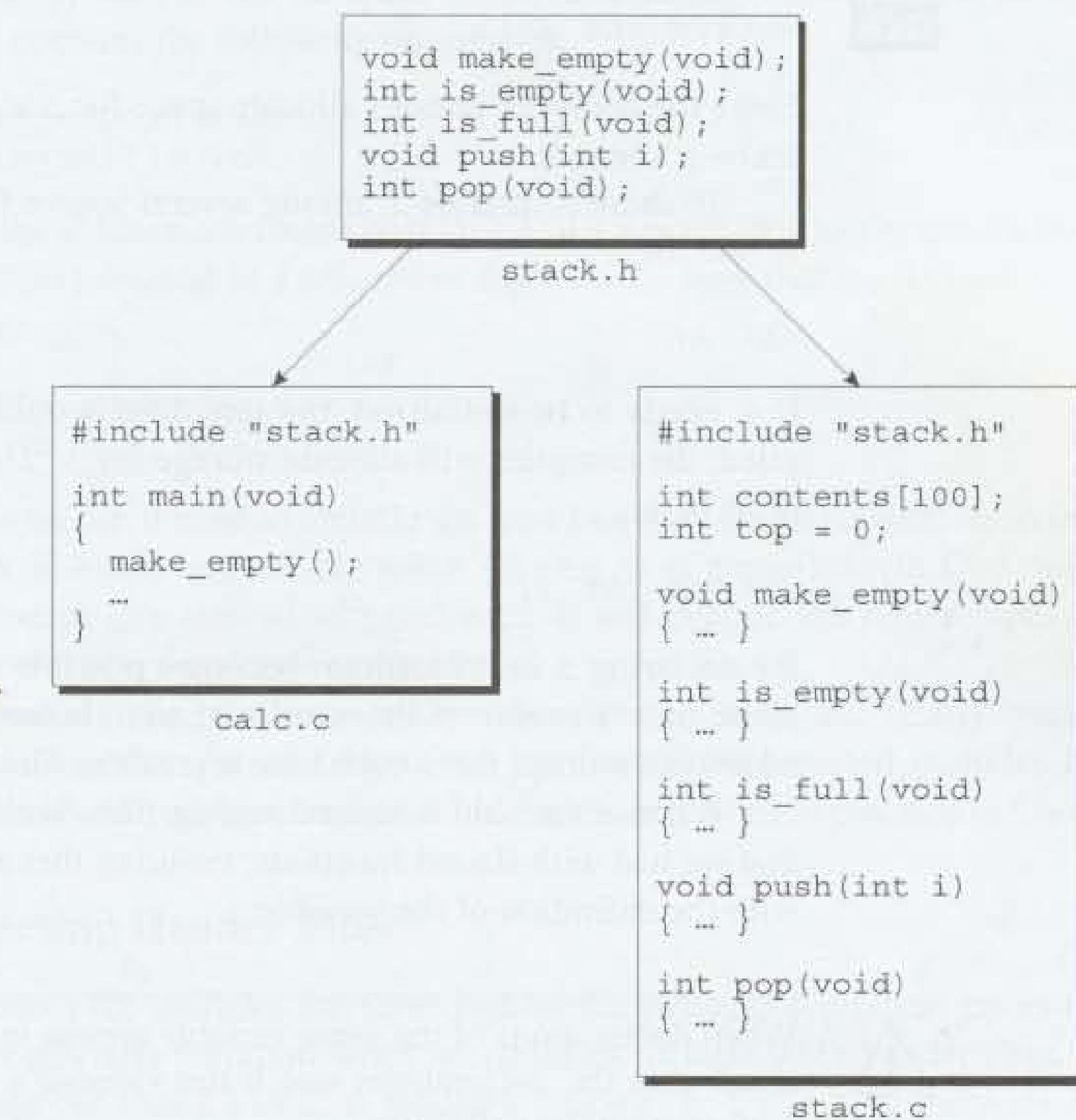
If *foo.c* contains other functions, most of them should be declared in the same header file as *f*. After all, the other functions in *foo.c* are presumably related to *f*; any file that contains a call of *f* probably needs some of the other functions in *foo.c*. Functions that are intended for use only within *foo.c* shouldn't be declared in a header file, however; to do so would be misleading.

To illustrate the use of function prototypes in header files, let's return to the RPN calculator of Section 15.1. The *stack.c* file will contain definitions of the *make\_empty*, *is\_empty*, *is\_full*, *push*, and *pop* functions. The following prototypes for these functions should go in the *stack.h* header file:

```
void make_empty(void);
int is_empty(void);
```

```
int is_full(void);
void push(int i);
int pop(void);
```

(To avoid complicating the example, `is_empty` and `is_full` will return `int` values instead of Boolean values.) We'll include `stack.h` in `calc.c` to allow the compiler to check any calls of stack functions that appear in the latter file. We'll also include `stack.h` in `stack.c` so the compiler can verify that the prototypes in `stack.h` match the definitions in `stack.c`. The following figure shows `stack.h`, `stack.c`, and `calc.c`:



## Sharing Variable Declarations

external variables > 10.2

External variables can be shared among files in much the same way functions are. To share a function, we put its *definition* in one source file, then put *declarations* in other files that need to call the function. Sharing an external variable is done in much the same way.

Up to this point, we haven't needed to distinguish between a variable's declaration and its definition. To declare a variable `i`, we've written

```
int i; /* declares i and defines it as well */
```

extern keyword ▶ 18.2

which not only declares *i* to be a variable of type `int`, but defines *i* as well, by causing the compiler to set aside space for *i*. To declare *i* without defining it, we must put the keyword `extern` at the beginning of its declaration:

```
extern int i; /* declares i without defining it */
```

`extern` informs the compiler that *i* is defined elsewhere in the program (most likely in a different source file), so there's no need to allocate space for it.

`extern` works with variables of all types. When we use it in the declaration of an array, we can omit the length of the array:

**Q&A**

```
extern int a[];
```

Since the compiler doesn't allocate space for *a* at this time, there's no need for it to know *a*'s length.

To share a variable *i* among several source files, we first put a definition of *i* in one file:

```
int i;
```

If *i* needs to be initialized, the initializer would go here. When this file is compiled, the compiler will allocate storage for *i*. The other files will contain declarations of *i*:

```
extern int i;
```

By declaring *i* in each file, it becomes possible to access and/or modify *i* within those files. Because of the word `extern`, however, the compiler doesn't allocate additional storage for *i* each time one of the files is compiled.

When a variable is shared among files, we'll face a challenge similar to one that we had with shared functions: ensuring that all declarations of a variable agree with the definition of the variable.



When declarations of the same variable appear in different files, the compiler can't check that the declarations match the variable's definition. For example, one file may contain the definition

```
int i;
```

while another file contains the declaration

```
extern long i;
```

An error of this kind can cause the program to behave unpredictably.

---

To avoid inconsistency, declarations of shared variables are usually put in header files. A source file that needs access to a particular variable can then include the appropriate header file. In addition, each header file that contains a

variable declaration is included in the source file that contains the variable's definition, enabling the compiler to check that the two match.

Although sharing variables among files is a long-standing practice in the C world, it has significant disadvantages. In Section 19.2, we'll see what the problems are and learn how to design programs that don't need shared variables.

## Nested Includes

A header file may itself contain `#include` directives. Although this practice may seem a bit odd, it can be quite useful in practice. Consider the `stack.h` file, which contains the following prototypes:

```
int is_empty(void);  
int is_full(void);
```

Since these functions return only 0 or 1, it's a good idea to declare their return type to be `Bool` instead of `int`, where `Bool` is the type that we defined earlier in this section:

```
Bool is_empty(void);  
Bool is_full(void);
```

Of course, we'll need to include the `boolean.h` file in `stack.h` so that the definition of `Bool` is available when `stack.h` is compiled. (In C99, we'd include `<stdbool.h>` instead of `boolean.h` and declare the return types of the two functions to be `bool` rather than `Bool`.)

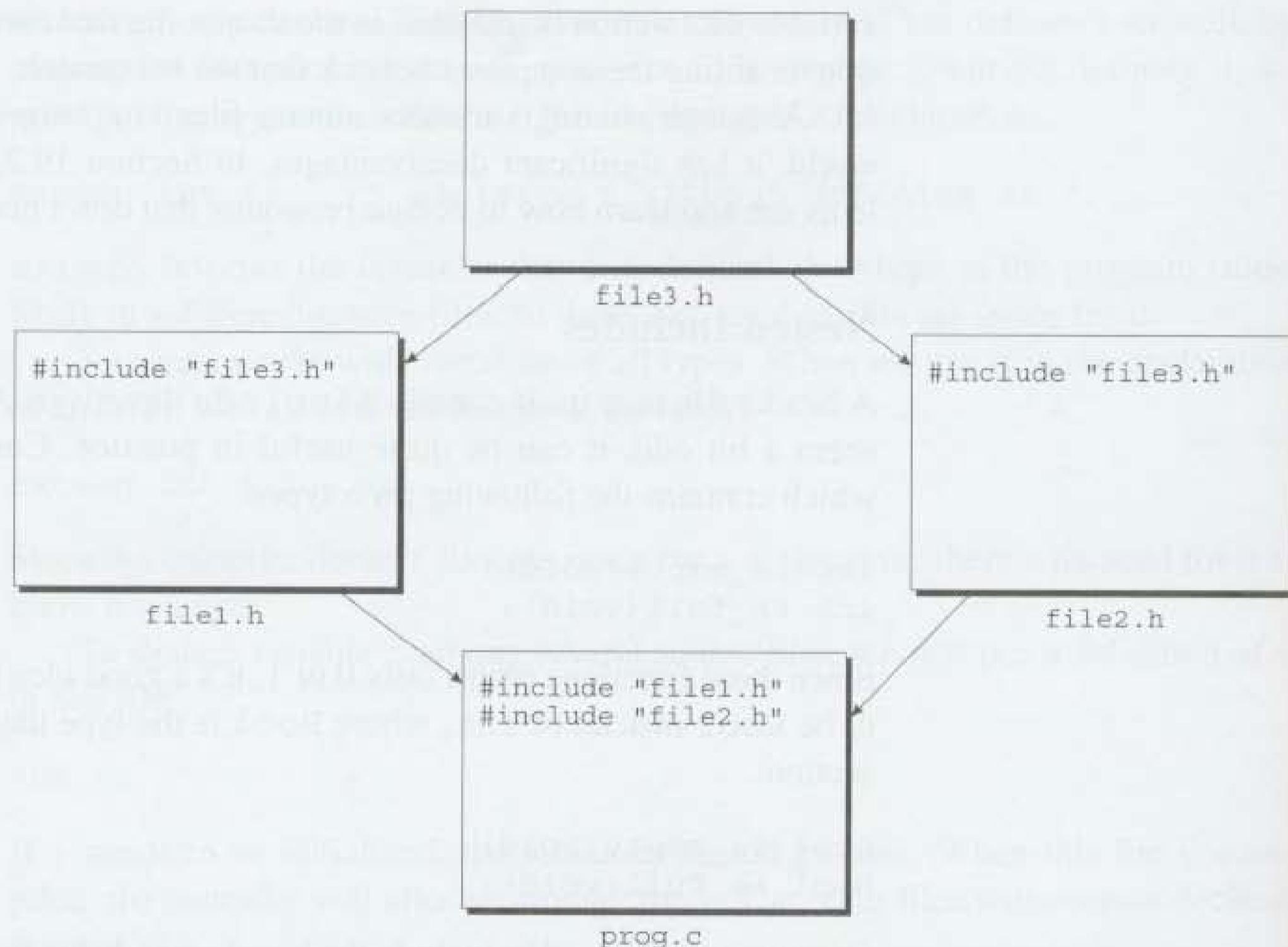
Traditionally, C programmers shun nested includes. (Early versions of C didn't allow them at all.) However, the bias against nested includes has largely faded away, in part because nested includes are common practice in C++.

## Protecting Header Files

If a source file includes the same header file twice, compilation errors may result. This problem is common when header files include other header files. For example, suppose that `file1.h` includes `file3.h`, `file2.h` includes `file3.h`, and `prog.c` includes both `file1.h` and `file2.h` (see the figure at the top of the next page). When `prog.c` is compiled, `file3.h` will be compiled twice.

Including the same header file twice doesn't always cause a compilation error. If the file contains only macro definitions, function prototypes, and/or variable declarations, there won't be any difficulty. If the file contains a type definition, however, we'll get a compilation error.

Just to be safe, it's probably a good idea to protect all header files against multiple inclusion; that way, we can add type definitions to a file later without the risk that we might forget to protect the file. In addition, we might save some time during program development by avoiding unnecessary recompilation of the same header file.



To protect a header file, we'll enclose the contents of the file in an `#ifndef`-`#endif` pair. For example, the `boolean.h` file could be protected in the following way:

```

#ifndef BOOLEAN_H
#define BOOLEAN_H

#define TRUE 1
#define FALSE 0
typedef int Bool;

#endif
  
```

When this file is included the first time, the `BOOLEAN_H` macro won't be defined, so the preprocessor will allow the lines between `#ifndef` and `#endif` to stay. But if the file should be included a second time, the preprocessor will remove the lines between `#ifndef` and `#endif`.

The name of the macro (`BOOLEAN_H`) doesn't really matter. However, making it resemble the name of the header file is a good way to avoid conflicts with other macros. Since we can't name the macro `BOOLEAN.H` (identifiers can't contain periods), a name such as `BOOLEAN_H` is a good alternative.

## #error Directives in Header Files

#error directives ► 14.5

`#error` directives are often put in header files to check for conditions under which the header file shouldn't be included. For example, suppose that a header

file uses a feature that didn't exist prior to the original C89 standard. To prevent the header file from being used with older, nonstandard compilers, it could contain an `#ifndef` directive that tests for the existence of the `_STDC_` macro:

```
#ifndef _STDC_
#error This header requires a Standard C compiler
#endif
```

## 15.3 Dividing a Program into Files

Let's now use what we know about header files and source files to develop a simple technique for dividing a program into files. We'll concentrate on functions, but the same principles apply to external variables as well. We'll assume that the program has already been designed; that is, we've decided what functions the program will need and how to arrange the functions into logically related groups. (We'll discuss program design in Chapter 19.)

Here's how we'll proceed. Each set of functions will go into a separate source file (let's use the name `foo.c` for one such file). In addition, we'll create a header file with the same name as the source file, but with the extension `.h` (`foo.h`, in our case). Into `foo.h`, we'll put prototypes for the functions defined in `foo.c`. (Functions that are designed for use only within `foo.c` need not—and should not—be declared in `foo.h`. The `read_char` function in our next program is an example.) We'll include `foo.h` in each source file that needs to call a function defined in `foo.c`. Moreover, we'll include `foo.h` in `foo.c` so that the compiler can check that the function prototypes in `foo.h` are consistent with the definitions in `foo.c`.

The `main` function will go in a file whose name matches the name of the program—if we want the program to be known as `bar`, then `main` should be in the file `bar.c`. It's possible that there are other functions in the same file as `main`, so long as they're not called from other files in the program.

### PROGRAM Text Formatting

To illustrate the technique that we've just discussed, let's apply it to a small text-formatting program named `justify`. As sample input to `justify`, we'll use a file named `quote` that contains the following (poorly formatted) quotation from “The development of the C programming language” by Dennis M. Ritchie (in *History of Programming Languages II*, edited by T. J. Bergin, Jr., and R. G. Gibson, Jr., Addison-Wesley, Reading, Mass., 1996, pages 671–687):

```
C is quirky, flawed, and an
enormous success. Although accidents of history
surely helped, it evidently satisfied a need
for a system implementation language efficient
```

enough to displace assembly language,  
yet sufficiently abstract and fluent to describe  
algorithms and interactions in a wide variety  
of environments.

-- Dennis M. Ritchie

To run the program from a UNIX or Windows prompt, we'd enter the command

```
justify <quote
```

The < symbol informs the operating system that *justify* will read from the file *quote* instead of accepting input from the keyboard. This feature, supported by input redirection ➤ 22.1 UNIX, Windows, and other operating systems, is called *input redirection*. When given the *quote* file as input, the *justify* program will produce the following output:

C is quirky, flawed, and an enormous success. Although accidents of history surely helped, it evidently satisfied a need for a system implementation language efficient enough to displace assembly language, yet sufficiently abstract and fluent to describe algorithms and interactions in a wide variety of environments. -- Dennis M. Ritchie

The output of *justify* will normally appear on the screen, but we can save it in a

output redirection ➤ 22.1 file by using *output redirection*:

```
justify <quote >newquote
```

The output of *justify* will go into the file *newquote*.

In general, *justify*'s output should be identical to its input, except that extra spaces and blank lines are deleted, and lines are filled and justified. "Filling" a line means adding words until one more word would cause the line to overflow. "Justifying" a line means adding extra spaces between words so that each line has exactly the same length (60 characters). Justification must be done so that the space between words in a line is equal (or as nearly equal as possible). The last line of the output won't be justified.

We'll assume that no word is longer than 20 characters. (A punctuation mark is considered part of the word to which it is adjacent.) That's a bit restrictive, of course, but once the program is written and debugged we can easily increase this limit to the point that it would virtually never be exceeded. If the program encounters a longer word, it must ignore all characters after the first 20, replacing them with a single asterisk. For example, the word

antidisestablishmentarianism

would be printed as

antidisestablishment\*

Now that we understand what the program should do, it's time to think about a design. We'll start by observing that the program can't write the words one by one as they're read. Instead, it will have to store them in a "line buffer" until there are enough to fill a line. After further reflection, we decide that the heart of the program will be a loop that goes something like this:

```
for (;;) {
    read word;
    if (can't read word) {
        write contents of line buffer without justification;
        terminate program;
    }
    if (word doesn't fit in line buffer) {
        write contents of line buffer with justification;
        clear line buffer;
    }
    add word to line buffer;
}
```

Since we'll need functions that deal with words and functions that deal with the line buffer, let's split the program into three source files, putting all functions related to words in one file (`word.c`) and all functions related to the line buffer in another file (`line.c`). A third file (`justify.c`) will contain the `main` function. In addition to these files, we'll need two header files, `word.h` and `line.h`. The `word.h` file will contain prototypes for the functions in `word.c`; `line.h` will play a similar role for `line.c`.

By examining the main loop, we see that the only word-related function that we'll need is a `read_word` function. (If `read_word` can't read a word because it's reached the end of the input file, we'll have it signal the main loop by pretending to read an "empty" word.) Consequently, the `word.h` file is a small one:

```
word.h #ifndef WORD_H
#define WORD_H

/***** * read_word: Reads the next word from the input and *
 * stores it in word. Makes word empty if no *
 * word could be read because of end-of-file. *
 * Truncates the word if its length exceeds *
 * len. */
void read_word(char *word, int len);

#endif
```

Notice how the `WORD_H` macro protects `word.h` from being included more than once. Although `word.h` doesn't really need it, it's good practice to protect all header files in this way.

The `line.h` file won't be as short as `word.h`. Our outline of the main loop reveals the need for functions that perform the following operations:

- Write contents of line buffer without justification
- Determine how many characters are left in line buffer
- Write contents of line buffer with justification
- Clear line buffer
- Add word to line buffer

We'll call these functions `flush_line`, `space_remaining`, `write_line`, `clear_line`, and `add_word`. Here's what the `line.h` header file will look like:

```
line.h #ifndef LINE_H
#define LINE_H

/***** * clear_line: Clears the current line. ****/
void clear_line(void);

/***** * add_word: Adds word to the end of the current line.
* If this is not the first word on the line,
* puts one space before word. ****/
void add_word(const char *word);

/***** * space_remaining: Returns the number of characters left *
* in the current line. ****/
int space_remaining(void);

/***** * write_line: Writes the current line with *
* justification. ****/
void write_line(void);

/***** * flush_line: Writes the current line without *
* justification. If the line is empty, does *
* nothing. ****/
void flush_line(void);

#endif
```

Before we write the `word.c` and `line.c` files, we can use the functions declared in `word.h` and `line.h` to write `justify.c`, the main program. Writing this file is mostly a matter of translating our original loop design into C.

```
justify.c /* Formats a file of text */

#include <string.h>
#include "line.h"
#include "word.h"

#define MAX_WORD_LEN 20

int main(void)
{
    char word[MAX_WORD_LEN+2];
    int word_len;

    clear_line();
    for (;;) {
        read_word(word, MAX_WORD_LEN+1);
        word_len = strlen(word);
        if (word_len == 0) {
            flush_line();
            return 0;
        }
        if (word_len > MAX_WORD_LEN)
            word[MAX_WORD_LEN] = '*';
        if (word_len + 1 > space_remaining()) {
            write_line();
            clear_line();
        }
        add_word(word);
    }
}
```

Including both `line.h` and `word.h` gives the compiler access to the function prototypes in both files as it compiles `justify.c`.

`main` uses a trick to handle words that exceed 20 characters. When it calls `read_word`, `main` tells it to truncate any word that exceeds 21 characters. After `read_word` returns, `main` checks whether `word` contains a string that's longer than 20 characters. If so, the word that was read must have been at least 21 characters long (before truncation), so `main` replaces the word's 21st character by an asterisk.

Now it's time to write `word.c`. Although the `word.h` header file has a prototype for only one function, `read_word`, we can put additional functions in `word.c` if we need to. As it turns out, `read_word` is easier to write if we add a small “helper” function, `read_char`. We'll assign `read_char` the task of reading a single character and, if it's a new-line character or tab, converting it to a space. Having `read_word` call `read_char` instead of `getchar` solves the problem of treating new-line characters and tabs as spaces.

Here's the `word.c` file:

```
word.c #include <stdio.h>
#include "word.h"
```

```

int read_char(void)
{
    int ch = getchar();

    if (ch == '\n' || ch == '\t')
        return ' ';
    return ch;
}

void read_word(char *word, int len)
{
    int ch, pos = 0;

    while ((ch = read_char()) == ' ')
        ;
    while (ch != ' ' && ch != EOF) {
        if (pos < len)
            word[pos++] = ch;
        ch = read_char();
    }
    word[pos] = '\0';
}

```

EOF macro ➤ 22.4

Before we discuss `read_word`, a couple of comments are in order concerning the use of `getchar` in the `read_char` function. First, `getchar` returns an `int` value instead of a `char` value; that's why the variable `ch` in `read_char` is declared to have type `int` and why the return type of `read_char` is `int`. Also, `getchar` returns the value `EOF` when it's unable to continue reading (usually because it has reached the end of the input file).

`read_word` consists of two loops. The first loop skips over spaces, stopping at the first nonblank character. (`EOF` isn't a blank, so the loop stops if it reaches the end of the input file.) The second loop reads characters until encountering a space or `EOF`. The body of the loop stores the characters in `word` until reaching the `len` limit. After that, the loop continues reading characters but doesn't store them. The final statement in `read_word` ends the word with a null character, thereby making it a string. If `read_word` encounters `EOF` before finding a nonblank character, `pos` will be 0 at the end, making `word` an empty string.

The only file left is `line.c`, which supplies definitions of the functions declared in the `line.h` file. `line.c` will also need variables to keep track of the state of the line buffer. One variable, `line`, will store the characters in the current line. Strictly speaking, `line` is the only variable we need. For speed and convenience, however, we'll use two other variables: `line_len` (the number of characters in the current line) and `num_words` (the number of words in the current line).

Here's the `line.c` file:

```

line.c #include <stdio.h>
#include <string.h>
#include "line.h"

```

```
#define MAX_LINE_LEN 60

char line[MAX_LINE_LEN+1];
int line_len = 0;
int num_words = 0;

void clear_line(void)
{
    line[0] = '\0';
    line_len = 0;
    num_words = 0;
}

void add_word(const char *word)
{
    if (num_words > 0) {
        line[line_len] = ' ';
        line[line_len+1] = '\0';
        line_len++;
    }
    strcat(line, word);
    line_len += strlen(word);
    num_words++;
}

int space_remaining(void)
{
    return MAX_LINE_LEN - line_len;
}

void write_line(void)
{
    int extra_spaces, spaces_to_insert, i, j;
    extra_spaces = MAX_LINE_LEN - line_len;
    for (i = 0; i < line_len; i++) {
        if (line[i] != ' ')
            putchar(line[i]);
        else {
            spaces_to_insert = extra_spaces / (num_words - 1);
            for (j = 1; j <= spaces_to_insert + 1; j++)
                putchar(' ');
            extra_spaces -= spaces_to_insert;
            num_words--;
        }
    }
    putchar('\n');
}

void flush_line(void)
{
    if (line_len > 0)
        puts(line);
}
```

Most of the functions in `line.c` are easy to write. The only tricky one is `write_line`, which writes a line with justification. `write_line` writes the characters in `line` one by one, pausing at the space between each pair of words to write additional spaces if needed. The number of additional spaces is stored in `spaces_to_insert`, which has the value `extra_spaces / (num_words - 1)`, where `extra_spaces` is initially the difference between the maximum line length and the actual line length. Since `extra_spaces` and `num_words` change after each word is printed, `spaces_to_insert` will change as well. If `extra_spaces` is 10 initially and `num_words` is 5, then the first word will be followed by 2 extra spaces, the second by 2, the third by 3, and the fourth by 3.

## 15.4 Building a Multiple-File Program

In Section 2.1, we examined the process of compiling and linking a program that fits into a single file. Let's expand that discussion to cover multiple-file programs. Building a large program requires the same basic steps as building a small one:

- **Compiling.** Each source file in the program must be compiled separately. (Header files don't need to be compiled; the contents of a header file are automatically compiled whenever a source file that includes it is compiled.) For each source file, the compiler generates a file containing object code. These files—known as *object files*—have the extension `.o` in UNIX and `.obj` in Windows.
- **Linking.** The linker combines the object files created in the previous step—along with code for library functions—to produce an executable file. Among other duties, the linker is responsible for resolving external references left behind by the compiler. (An external reference occurs when a function in one file calls a function defined in another file or accesses a variable defined in another file.)

Most compilers allow us to build a program in a single step. With the GCC compiler, for example, we'd use the following command to build the `justify` program of Section 15.3:

```
gcc -o justify justify.c line.c word.c
```

The three source files are first compiled into object code. The object files are then automatically passed to the linker, which combines them into a single file. The `-o` option specifies that we want the executable file to be named `justify`.

### Makefiles

Putting the names of all the source files on the command line quickly gets tedious. Worse still, we could waste a lot of time when rebuilding a program if we recompile all source files, not just the ones that were affected by our most recent changes.

To make it easier to build large programs, UNIX originated the concept of the *makefile*, a file containing the information necessary to build a program. A makefile not only lists the files that are part of the program, but also describes *dependencies* among the files. Suppose that the file `foo.c` includes the file `bar.h`. We say that `foo.c` “depends” on `bar.h`, because a change to `bar.h` will require us to recompile `foo.c`.

Here’s a UNIX makefile for the `justify` program. The makefile uses GCC for compilation and linking:

```
justify: justify.o word.o line.o
        gcc -o justify justify.o word.o line.o

justify.o: justify.c word.h line.h
        gcc -c justify.c

word.o: word.c word.h
        gcc -c word.c

line.o: line.c line.h
        gcc -c line.c
```

There are four groups of lines; each group is known as a *rule*. The first line in each rule gives a *target* file, followed by the files on which it depends. The second line is a *command* to be executed if the target should need to be rebuilt because of a change to one of its dependent files. Let’s look at the first two rules; the last two are similar.

In the first rule, `justify` (the executable file) is the target:

```
justify: justify.o word.o line.o
        gcc -o justify justify.o word.o line.o
```

The first line states that `justify` depends on the files `justify.o`, `word.o`, and `line.o`; if any one of these three files has changed since the program was last built, then `justify` needs to be rebuilt. The command on the following line shows how the rebuilding is to be done (by using the `gcc` command to link the three object files).

In the second rule, `justify.o` is the target:

```
justify.o: justify.c word.h line.h
        gcc -c justify.c
```

The first line indicates that `justify.o` needs to be rebuilt if there’s been a change to `justify.c`, `word.h`, or `line.h`. (The reason for mentioning `word.h` and `line.h` is that `justify.c` includes both these files, so it’s potentially affected by a change to either one.) The next line shows how to update `justify.o` (by recompiling `justify.c`). The `-c` option tells the compiler to compile `justify.c` into an object file but not attempt to link it.

### Q&A

Once we’ve created a makefile for a program, we can use the `make` utility to build (or rebuild) the program. By checking the time and date associated with each

file in the program, make can determine which files are out of date. It then invokes the commands necessary to rebuild the program.

If you want to give make a try, here are a few details you'll need to know:

- Each command in a makefile must be preceded by a tab character, not a series of spaces. (In our example, the commands appear to be indented eight spaces, but it's actually a single tab character.)
- A makefile is normally stored in a file named `Makefile` (or `makefile`). When the `make` utility is used, it automatically checks the current directory for a file with one of these names.
- To invoke `make`, use the command

```
make target
```

where *target* is one of the targets listed in the makefile. To build the `justify` executable using our makefile, we would use the command

```
make justify
```

- If no target is specified when `make` is invoked, it will build the target of the first rule. For example, the command

```
make
```

will build the `justify` executable, since `justify` is the first target in our makefile. Except for this special property of the first rule, the order of rules in a makefile is arbitrary.

`make` is complicated enough that entire books have been written about it, so we won't attempt to delve further into its intricacies. Let's just say that real makefiles aren't usually as easy to understand as our example. There are numerous techniques that reduce the amount of redundancy in makefiles and make them easier to modify; at the same time, though, these techniques greatly reduce their readability.

Not everyone uses makefiles, by the way. Other program maintenance tools are also popular, including the "project files" supported by some integrated development environments.

## Errors During Linking

Some errors that can't be detected during compilation will be found during linking. In particular, if the definition of a function or variable is missing from a program, the linker will be unable to resolve external references to it, causing a message such as "*undefined symbol*" or "*undefined reference*".

Errors detected by the linker are usually easy to fix. Here are some of the most common causes:

- *Misspellings.* If the name of a variable or function is misspelled, the linker will report it as missing. For example, if the function `read_char` is defined

in the program but called as `read_cahr`, the linker will report that `read_cahr` is missing.

- **Missing files.** If the linker can't find the functions that are in file `foo.c`, it may not know about the file. Check the makefile or project file to make sure that `foo.c` is listed there.
- **Missing libraries.** The linker may not be able to find all library functions used in the program. A classic example occurs in UNIX programs that use the `<math.h>` header. Simply including the header in a program may not be enough; many versions of UNIX require that the `-lm` option be specified when the program is linked, causing the linker to search a system file that contains compiled versions of the `<math.h>` functions. Failing to use this option may cause “undefined reference” messages during linking.

## Rebuilding a Program

During the development of a program, it's rare that we'll need to compile all its files. Most of the time, we'll test the program, make a change, then build the program again. To save time, the rebuilding process should recompile only those files that might be affected by the latest change.

Let's assume that we've designed our program in the way outlined in Section 15.3, with a header file for each source file. To see how many files will need to be recompiled after a change, we need to consider two possibilities.

The first possibility is that the change affects a single source file. In that case, only that file must be recompiled. (After that, the entire program will need to be relinked, of course.) Consider the `justify` program. Suppose that we decide to condense the `read_char` function in `word.c` (changes are marked in **bold**):

```
int read_char(void)
{
    int ch = getchar();

    return (ch == '\n' || ch == '\t') ? ' ' : ch;
}
```

This modification doesn't affect `word.h`, so we need only recompile `word.c` and relink the program.

The second possibility is that the change affects a header file. In that case, we should recompile all files that include the header file, since they could potentially be affected by the change. (Some of them might not be, but it pays to be conservative.)

As an example, consider the `read_word` function in the `justify` program. Notice that `main` calls `strlen` immediately after calling `read_word`, in order to determine the length of the word that was just read. Since `read_word` already knows the length of the word (`read_word`'s `pos` variable keeps track of the length), it seems silly to use `strlen`. Modifying `read_word` to return the word's length is easy. First, we change the prototype of `read_word` in `word.h`:

```

*****  

* read_word: Reads the next word from the input and *  

* stores it in word. Makes word empty if no *  

* word could be read because of end-of-file. *  

* Truncates the word if its length exceeds *  

* len. Returns the number of characters *  

* stored.  

*****  

int read_word(char *word, int len);

```

Of course, we're careful to change the comment that accompanies `read_word`. Next, we change the definition of `read_word` in `word.c`:

```

int read_word(char *word, int len)
{
    int ch, pos = 0;

    while ((ch = read_char()) == ' ')
        ;
    while (ch != ' ' && ch != EOF) {
        if (pos < len)
            word[pos++] = ch;
        ch = read_char();
    }
    word[pos] = '\0';
    return pos;
}

```

Finally, we modify `justify.c` by removing the include of `<string.h>` and changing `main` as follows:

```

int main(void)
{
    char word[MAX_WORD_LEN+2];
    int word_len;

    clear_line();
    for (;;) {
        word_len = read_word(word, MAX_WORD_LEN+1);
        if (word_len == 0) {
            flush_line();
            return 0;
        }
        if (word_len > MAX_WORD_LEN)
            word[MAX_WORD_LEN] = '*';
        if (word_len + 1 > space_remaining()) {
            write_line();
            clear_line();
        }
        add_word(word);
    }
}

```

Once we've made these changes, we'll rebuild the `justify` program by recompiling `word.c` and `justify.c` and then relinking. There's no need to recompile `line.c`, which doesn't include `word.h` and therefore won't be affected by changes to it. With the GCC compiler, we could use the following command to rebuild the program:

```
gcc -o justify justify.c word.c line.o
```

Note the mention of `line.o` instead of `line.c`.

One of the advantages of using makefiles is that rebuilding is handled automatically. By examining the date of each file, the `make` utility can determine which files have changed since the program was last built. It then recompiles these files, together with all files that depend on them, either directly or indirectly. For example, if we make the indicated changes to `word.h`, `word.c`, and `justify.c` and then rebuild the `justify` program, `make` will perform the following actions:

1. Build `justify.o` by compiling `justify.c` (because `justify.c` and `word.h` were changed).
2. Build `word.o` by compiling `word.c` (because `word.c` and `word.h` were changed).
3. Build `justify` by linking `justify.o`, `word.o`, and `line.o` (because `justify.o` and `word.o` were changed).

## Defining Macros Outside a Program

C compilers usually provide some method of specifying the value of a macro at the time a program is compiled. This ability makes it easy to change the value of a macro without editing any of the program's files. It's especially valuable when programs are built automatically using makefiles.

Most compilers (including GCC) support the `-D` option, which allows the value of a macro to be specified on the command line:

```
gcc -DDEBUG=1 foo.c
```

In this example, the `DEBUG` macro is defined to have the value `1` in the program `foo.c`, just as if the line

```
#define DEBUG 1
```

appeared at the beginning of `foo.c`. If the `-D` option names a macro without specifying its value, the value is taken to be `1`.

predefined macros ► 14.3

Many compilers also support the `-U` option, which "undefines" a macro as if by using `#undef`. We can use `-U` to undefine a predefined macro or one that was defined earlier in the command line using `-D`.

## Q & A

**Q:** You don't have any examples that use the `#include` directive to include a source file. What would happen if we were to do this?

**A:** That's not a good practice, although it's not illegal. Here's an example of the kind of trouble you can get into. Suppose that `foo.c` defines a function `f` that we'll need in `bar.c` and `baz.c`, so we put the directive

```
#include "foo.c"
```

in both `bar.c` and `baz.c`. Each of these files will compile nicely. The problem comes later, when the linker discovers two copies of the object code for `f`. Of course, we would have gotten away with including `foo.c` if only `bar.c` had included it, not `baz.c` as well. To avoid problems, it's best to use `#include` only with header files, not source files.

**Q:** What are the exact search rules for the `#include` directive? [p. 351]

**A:** That depends on your compiler. The C standard is deliberately vague in its description of `#include`. If the file name is enclosed in *brackets*, the preprocessor looks in a “sequence of implementation-defined places,” as the standard obliquely puts it. If the file name is enclosed in *quotation marks*, the file “is searched for in an implementation-defined manner” and, if not found, then searched as if its name had been enclosed in brackets. The reason for this waffling is simple: not all operating systems have hierarchical (tree-like) file systems.

To make matters even more interesting, the standard doesn't require that names enclosed in brackets be file names at all, leaving open the possibility that `#include` directives using `<>` are handled entirely within the compiler.

**Q:** I don't understand why each source file needs its own header file. Why not have one big header file containing macro definitions, type definitions, and function prototypes? By including this file, each source file would have access to all the shared information it needs. [p. 354]

**A:** The “one big header file” approach certainly works; a number of programmers use it. And it does have an advantage: with only one header file, there are fewer files to manage. For large programs, however, the disadvantages of this approach tend to outweigh its advantages.

Using a single header file provides no useful information to someone reading the program later. With multiple header files, the reader can quickly see what other parts of the program are used by a particular source file.

But that's not all. Since each source file depends on the big header file, changing it will cause all source files to be recompiled—a significant drawback in a large program. To make matters worse, the header file will probably change frequently because of the large amount of information it contains.

**Q:** The chapter says that a shared array should be declared as follows:

```
extern int a[];
```

Since arrays and pointers are closely related, would it be legal to write

```
extern int *a;
```

instead? [p. 356]

- A:** No. When used in expressions, arrays “decay” into pointers. (We’ve noticed this behavior when an array name is used as an argument in a function call.) In variable declarations, however, arrays and pointers are distinct types.

**Q:** Does it hurt if a source file includes headers that it doesn’t really need?

- A:** Not unless the header has a declaration or definition that conflicts with one in the source file. Otherwise, the worst that can happen is a minor increase in the time it takes to compile the source file.

**Q:** I needed to call a function in the file `foo.c`, so I included the matching header file, `foo.h`. My program compiled, but it won’t link. Why?

- A:** Compilation and linking are completely separate in C. Header files exist to provide information to the compiler, not the linker. If you want to call a function in `foo.c`, then you have to make sure that `foo.c` is compiled and that the linker is aware that it must search the object file for `foo.c` to find the function. Usually this means naming `foo.c` in the program’s makefile or project file.

**Q:** If my program calls a function in `<stdio.h>`, does that mean that all functions in `<stdio.h>` will be linked with the program?

- A:** No. Including `<stdio.h>` (or any other header) has no effect on linking. In any event, most linkers will link only functions that your program actually needs.

**Q:** Where can I get the `make` utility? [p. 367]

- A:** `make` is a standard UNIX utility. The GNU version, known as GNU Make, is included in most Linux distributions. It’s also available directly from the Free Software Foundation ([www.gnu.org/software/make/](http://www.gnu.org/software/make/)).

## Exercises

### Section 15.1

1. Section 15.1 listed several advantages of dividing a program into multiple source files.
  - (a) Describe several other advantages.
  - (b) Describe some disadvantages.

### Section 15.2

- W 2. Which of the following should *not* be put in a header file? Why not?
  - (a) Function prototypes
  - (b) Function definitions

- (c) Macro definitions  
 (d) Type definitions
3. We saw that writing `#include <file>` instead of `#include "file"` may not work if `file` is one that we've written. Would there be any problem with writing `#include "file"` instead of `#include <file>` if `file` is a system header?
  4. Assume that `debug.h` is a header file with the following contents:

```
#ifdef DEBUG
#define PRINT_DEBUG(n) printf("Value of " #n ": %d\n", n)
#else
#define PRINT_DEBUG(n)
#endif
```

Let `testdebug.c` be the following source file:

```
#include <stdio.h>

#define DEBUG
#include "debug.h"

int main(void)
{
    int i = 1, j = 2, k = 3;

#ifdef DEBUG
    printf("Output if DEBUG is defined:\n");
#else
    printf("Output if DEBUG is not defined:\n");
#endif

    PRINT_DEBUG(i);
    PRINT_DEBUG(j);
    PRINT_DEBUG(k);
    PRINT_DEBUG(i + j);
    PRINT_DEBUG(2 * i + j - k);

    return 0;
}
```

- (a) What is the output when the program is executed?  
 (b) What is the output if the `#define` directive is removed from `testdebug.c`?  
 (c) Explain why the output is different in parts (a) and (b).  
 (d) Is it necessary for the `DEBUG` macro to be defined *before* `debug.h` is included in order for `PRINT_DEBUG` to have the desired effect? Justify your answer.

#### Section 15.4

5. Suppose that a program consists of three source files—`main.c`, `f1.c`, and `f2.c`—plus two header files, `f1.h` and `f2.h`. All three source files include `f1.h`, but only `f1.c` and `f2.c` include `f2.h`. Write a makefile for this program, assuming that the compiler is `gcc` and that the executable file is to be named `demo`.
6. The following questions refer to the program described in Exercise 5.
  - (a) Which files need to be compiled when the program is built for the first time?  
 (b) If `f1.c` is changed after the program has been built, which files need to be recompiled?  
 (c) If `f1.h` is changed after the program has been built, which files need to be recompiled?  
 (d) If `f2.h` is changed after the program has been built, which files need to be recompiled?

## Programming Projects

1. The `justify` program of Section 15.3 justifies lines by inserting extra spaces between words. The way the `write_line` function currently works, the words closer to the end of a line tend to have slightly wider gaps between them than the words at the beginning. (For example, the words closer to the end might have three spaces between them, while the words closer to the beginning might be separated by only two spaces.) Improve the program by having `write_line` alternate between putting the larger gaps at the end of the line and putting them at the beginning of the line.
2. Modify the `justify` program of Section 15.3 by having the `read_word` function (instead of `main`) store the `*` character at the end of a word that's been truncated.
3. Modify the `qsort.c` program of Section 9.6 so that the `quicksort` and `split` functions are in a separate file named `quicksort.c`. Create a header file named `quicksort.h` that contains prototypes for the two functions and have both `qsort.c` and `quicksort.c` include this file.
4. Modify the `remind.c` program of Section 13.5 so that the `read_line` function is in a separate file named `readline.c`. Create a header file named `readline.h` that contains a prototype for the function and have both `remind.c` and `readline.c` include this file.
5. Modify Programming Project 6 from Chapter 10 so that it has separate `stack.h` and `stack.c` files, as described in Section 15.2.



# 16 Structures, Unions, and Enumerations

*Functions delay binding; data structures induce binding.  
Moral: Structure data late in the programming process.*

This chapter introduces three new types: structures, unions, and enumerations. A structure is a collection of values (members), possibly of different types. A union is similar to a structure, except that its members share the same storage; as a result, a union can store one member at a time, but not all members simultaneously. An enumeration is an integer type whose values are named by the programmer.

Of these three types, structures are by far the most important, so I'll devote most of the chapter to them. Section 16.1 shows how to declare structure variables and perform basic operations on them. Section 16.2 then explains how to define structure types, which—among other things—allow us to write functions that accept structure arguments or return structures. Section 16.3 explores how arrays and structures can be nested. The last two sections are devoted to unions (Section 16.4) and enumerations (Section 16.5).

## 16.1 Structure Variables

The only data structure we've covered so far is the array. Arrays have two important properties. First, all elements of an array have the same type. Second, to select an array element, we specify its position (as an integer subscript).

The properties of a *structure* are quite different from those of an array. The elements of a structure (its *members*, in C parlance) aren't required to have the same type. Furthermore, the members of a structure have names; to select a particular member, we specify its name, not its position.

Structures may sound familiar, since most programming languages provide a similar feature. In some languages, structures are called *records*, and members are known as *fields*.

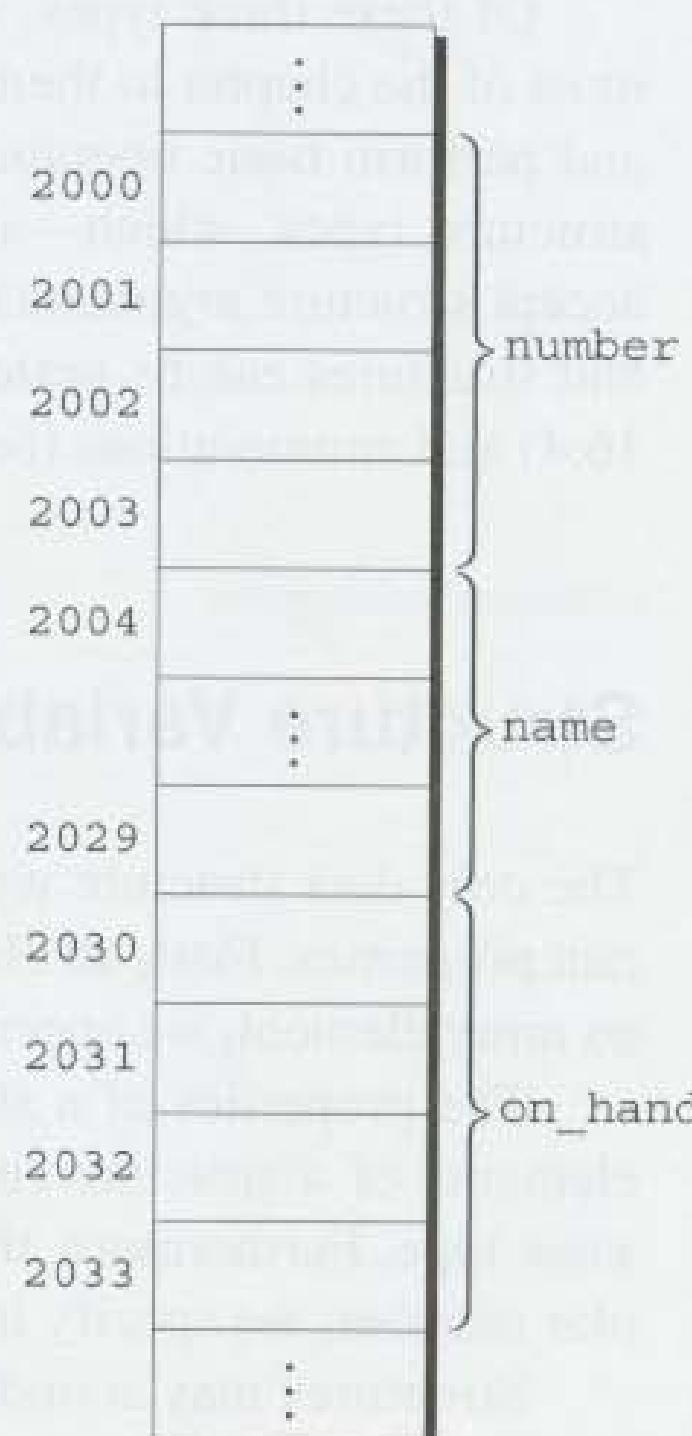
## Declaring Structure Variables

When we need to store a collection of related data items, a structure is a logical choice. For example, suppose that we need to keep track of parts in a warehouse. The information that we'll need to store for each part might include a part number (an integer), a part name (a string of characters), and the number of parts on hand (an integer). To create variables that can store all three items of data, we might use a declaration such as the following:

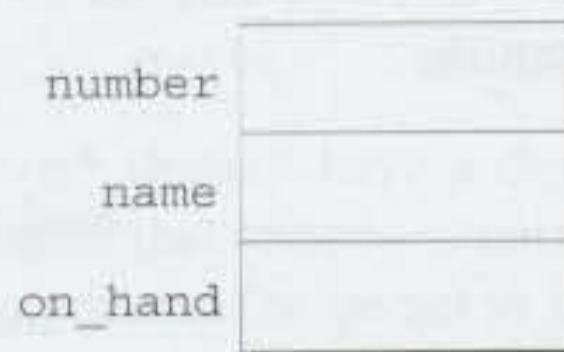
```
struct {
    int number;
    char name [NAME_LEN+1];
    int on_hand;
} part1, part2;
```

Each structure variable has three members: `number` (the part number), `name` (the name of the part), and `on_hand` (the quantity on hand). Notice that this declaration has the same form as other variable declarations in C: `struct { ... }` specifies a type, while `part1` and `part2` are variables of that type.

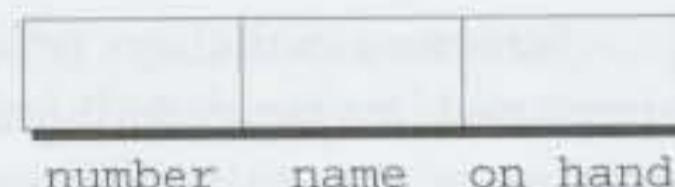
The members of a structure are stored in memory in the order in which they're declared. In order to show what the `part1` variable looks like in memory, let's assume that (1) `part1` is located at address 2000, (2) integers occupy four bytes, (3) `NAME_LEN` has the value 25, and (4) there are no gaps between the members. With these assumptions, `part1` will have the following appearance:



Usually it's not necessary to draw structures in such detail. I'll normally show them more abstractly, as a series of boxes:



I may sometimes draw the boxes horizontally instead of vertically:



Member values will go in the boxes later; for now, I've left them empty.

Each structure represents a new scope; any names declared in that scope won't conflict with other names in a program. (In C terminology, we say that each structure has a separate *name space* for its members.) For example, the following declarations can appear in the same program:

```

struct {
    int number;
    char name[NAME_LEN+1];
    int on_hand;
} part1, part2;

struct {
    char name[NAME_LEN+1];
    int number;
    char sex;
} employee1, employee2;
  
```

The `number` and `name` members in the `part1` and `part2` structures don't conflict with the `number` and `name` members in `employee1` and `employee2`.

## Initializing Structure Variables

Like an array, a structure variable may be initialized at the time it's declared. To initialize a structure, we prepare a list of values to be stored in the structure and enclose it in braces:

```

struct {
    int number;
    char name[NAME_LEN+1];
    int on_hand;
} part1 = {528, "Disk drive", 10},
      part2 = {914, "Printer cable", 5};
  
```

The values in the initializer must appear in the same order as the members of the structure. In our example, the `number` member of `part1` will be 528, the `name` member will be "Disk drive", and so on. Here's how `part1` will look after initialization:

|         |            |
|---------|------------|
| number  | 528        |
| name    | Disk drive |
| on_hand | 10         |

Structure initializers follow rules similar to those for array initializers. Expressions used in a structure initializer must be constant; for example, we couldn't have used a variable to initialize `part1`'s `on_hand` member. (This restriction is relaxed in C99, as we'll see in Section 18.5.) An initializer can have fewer members than the structure it's initializing; as with arrays, any "leftover" members are given 0 as their initial value. In particular, the bytes in a leftover character array will be zero, making it represent the empty string.

**C99**

## Designated Initializers

C99's designated initializers, which were discussed in Section 8.1 in the context of arrays, can also be used with structures. Consider the initializer for `part1` shown in the previous example:

```
{ 528, "Disk drive", 10 }
```

A designated initializer would look similar, but with each value labeled by the name of the member that it initializes:

```
{ .number = 528, .name = "Disk drive", .on_hand = 10 }
```

The combination of the period and the member name is called a *designator*. (Designators for array elements have a different form.)

Designated initializers have several advantages. For one, they're easier to read and check for correctness, because the reader can clearly see the correspondence between the members of the structure and the values listed in the initializer. Another is that the values in the initializer don't have to be placed in the same order that the members are listed in the structure. Our example initializer could be written as follows:

```
{ .on_hand = 10, .name = "Disk drive", .number = 528 }
```

Since the order doesn't matter, the programmer doesn't have to remember the order in which the members were originally declared. Moreover, the order of the members can be changed in the future without affecting designated initializers.

Not all values listed in a designated initializer need be prefixed by a designator. (This is true for arrays as well, as we saw in Section 8.1.) Consider the following example:

```
{ .number = 528, "Disk drive", .on_hand = 10 }
```

The value "Disk drive" doesn't have a designator, so the compiler assumes that it initializes the member that follows `number` in the structure. Any members that the initializer fails to account for are set to zero.

## Operations on Structures

Since the most common array operation is subscripting—selecting an element by position—it's not surprising that the most common operation on a structure is selecting one of its members. Structure members are accessed by name, though, not by position.

To access a member within a structure, we write the name of the structure first, then a period, then the name of the member. For example, the following statements will display the values of `part1`'s members:

```
printf("Part number: %d\n", part1.number);
printf("Part name: %s\n", part1.name);
printf("Quantity on hand: %d\n", part1.on_hand);
```

Ivalues ▶ 4.2

The members of a structure are Ivalues, so they can appear on the left side of an assignment or as the operand in an increment or decrement expression:

```
part1.number = 258;           /* changes part1's part number */
part1.on_hand++;             /* increments part1's quantity on hand */
```

The period that we use to access a structure member is actually a C operator. It has the same precedence as the postfix `++` and `--` operators, so it takes precedence over nearly all other operators. Consider the following example:

```
scanf("%d", &part1.on_hand);
```

The expression `&part1.on_hand` contains two operators (`&` and `.`). The `.` operator takes precedence over the `&` operator, so `&` computes the address of `part1.on_hand`, as we wished.

The other major structure operation is assignment:

```
part2 = part1;
```

The effect of this statement is to copy `part1.number` into `part2.number`, `part1.name` into `part2.name`, and so on.

Since arrays can't be copied using the `=` operator, it comes as something of a surprise to discover that structures can. It's even more surprising when you consider that an array embedded within a structure is copied when the enclosing structure is copied. Some programmers exploit this property by creating "dummy" structures to enclose arrays that will be copied later:

```
struct { int a[10]; } a1, a2;
a1 = a2; /* legal, since a1 and a2 are structures */
```

The = operator can be used only with structures of *compatible* types. Two structures declared at the same time (as part1 and part2 were) are compatible. As we'll see in the next section, structures declared using the same "structure tag" or the same type name are also compatible.

Other than assignment, C provides no operations on entire structures. In particular, we can't use the == and != operators to test whether two structures are equal or not equal.

**Q&A**

## 16.2 Structure Types

Although the previous section showed how to declare structure *variables*, it failed to discuss an important issue: naming structure *types*. Suppose that a program needs to declare several structure variables with identical members. If all the variables can be declared at one time, there's no problem. But if we need to declare the variables at different points in the program, then life becomes more difficult. If we write

```
struct {
    int number;
    char name[NAME_LEN+1];
    int on_hand;
} part1;
```

in one place and

```
struct {
    int number;
    char name[NAME_LEN+1];
    int on_hand;
} part2;
```

in another, we'll quickly run into problems. Repeating the structure information will bloat the program. Changing the program later will be risky, since we can't easily guarantee that the declarations will remain consistent.

But those aren't the biggest problems. According to the rules of C, part1 and part2 don't have compatible types. As a result, part1 can't be assigned to part2, and vice versa. Also, since we don't have a name for the type of part1 or part2, we can't use them as arguments in function calls.

To avoid these difficulties, we need to be able to define a name that represents a *type* of structure, not a particular structure *variable*. As it turns out, C provides two ways to name structures: we can either declare a "structure tag" or use `typedef` to define a type name.

**Q&A**

## Declaring a Structure Tag

A *structure tag* is a name used to identify a particular kind of structure. The following example declares a structure tag named `part`:

```
struct part {
    int number;
    char name[NAME_LEN+1];
    int on_hand;
};
```

Notice the semicolon that follows the right brace—it must be present to terminate the declaration.



Accidentally omitting the semicolon at the end of a structure declaration can cause surprising errors. Consider the following example:

```
struct part {
    int number;
    char name[NAME_LEN+1];
    int on_hand;
}           /*** WRONG: semicolon missing ***/
```

```
f(void)
{
    ...
    return 0; /* error detected at this line */
}
```

The programmer failed to specify the return type of the function `f` (a bit of sloppy programming). Since the preceding structure declaration wasn't terminated properly, the compiler assumes that `f` returns a value of type `struct part`. The error won't be detected until the compiler reaches the first `return` statement in the function. The result: a cryptic error message.

Once we've created the `part` tag, we can use it to declare variables:

```
struct part part1, part2;
```

Unfortunately, we can't abbreviate this declaration by dropping the word `struct`:

```
part part1, part2;    /*** WRONG ***/
```

`part` isn't a type name; without the word `struct`, it is meaningless.

Since structure tags aren't recognized unless preceded by the word `struct`, they don't conflict with other names used in a program. It would be perfectly legal (although more than a little confusing) to have a variable named `part`.

Incidentally, the declaration of a structure *tag* can be combined with the declaration of structure *variables*:

```
struct part {
    int number;
    char name[NAME_LEN+1];
    int on_hand;
} part1, part2;
```

Here, we've declared a structure tag named `part` (making it possible to use `part` later to declare more variables) as well as variables named `part1` and `part2`.

All structures declared to have type `struct part` are compatible with one another:

```
struct part part1 = {528, "Disk drive", 10};
struct part part2;

part2 = part1; /* legal; both parts have the same type */
```

## Defining a Structure Type

As an alternative to declaring a structure tag, we can use `typedef` to define a genuine type name. For example, we could define a type named `Part` in the following way:

```
typedef struct {
    int number;
    char name[NAME_LEN+1];
    int on_hand;
} Part;
```

Note that the name of the type, `Part`, must come at the end, not after the word `struct`.

We can use `Part` in the same way as the built-in types. For example, we might use it to declare variables:

```
Part part1, part2;
```

Since `Part` is a `typedef` name, we're not allowed to write `struct Part`. All `Part` variables, regardless of where they're declared, are compatible.

When it comes time to name a structure, we can usually choose either to declare a structure tag or to use `typedef`. However, as we'll see later, declaring a structure tag is mandatory when the structure is to be used in a linked list. I'll use structure tags rather than `typedef` names in most of my examples.

## Structures as Arguments and Return Values

Functions may have structures as arguments and return values. Let's look at two examples. Our first function, when given a `part` structure as its argument, prints the structure's members:

```
void print_part(struct part p)
{
    printf("Part number: %d\n", p.number);
```

### Q&A

linked lists ► 17.5

```

    printf("Part name: %s\n", p.name);
    printf("Quantity on hand: %d\n", p.on_hand);
}

```

Here's how `print_part` might be called:

```
print_part(part1);
```

Our second function returns a `part` structure that it constructs from its arguments:

```

struct part build_part(int number, const char *name,
                      int on_hand)
{
    struct part p;

    p.number = number;
    strcpy(p.name, name);
    p.on_hand = on_hand;
    return p;
}

```

Notice that it's legal for `build_part`'s parameters to have names that match the members of the `part` structure, since the structure has its own name space. Here's how `build_part` might be called:

```
part1 = build_part(528, "Disk drive", 10);
```

Passing a structure to a function and returning a structure from a function both require making a copy of all members in the structure. As a result, these operations impose a fair amount of overhead on a program, especially if the structure is large. To avoid this overhead, it's sometimes advisable to pass a *pointer* to a structure instead of passing the structure itself. Similarly, we might have a function return a pointer to a structure instead of returning an actual structure. Section 17.5 gives examples of functions that have a pointer to a structure as an argument and/or return a pointer to a structure.

There are other reasons to avoid copying structures besides efficiency. For example, the `<stdio.h>` header defines a type named `FILE`, which is typically a structure. Each `FILE` structure stores information about the state of an open file and therefore must be unique in a program. Every function in `<stdio.h>` that opens a file returns a pointer to a `FILE` structure, and every function that performs an operation on an open file requires a `FILE` pointer as an argument.

On occasion, we may want to initialize a structure variable inside a function to match another structure, possibly supplied as a parameter to the function. In the following example, the initializer for `part2` is the parameter passed to the `f` function:

```

void f(struct part part1)
{
    struct part part2 = part1;
    ...
}

```

automatic storage duration ➤ 10.1

C permits initializers of this kind, provided that the structure we're initializing (part2, in this case) has automatic storage duration (it's local to a function and hasn't been declared `static`). The initializer can be any expression of the proper type, including a function call that returns a structure.

C99

## Compound Literals

Section 9.3 introduced the C99 feature known as the *compound literal*. In that section, compound literals were used to create unnamed arrays, usually for the purpose of passing the array to a function. A compound literal can also be used to create a structure “on the fly,” without first storing it in a variable. The resulting structure can be passed as a parameter, returned by a function, or assigned to a variable. Let's look at a couple of examples.

First, we can use a compound literal to create a structure that will be passed to a function. For example, we could call the `print_part` function as follows:

```
print_part((struct part) {528, "Disk drive", 10});
```

The compound literal (shown in **bold**) creates a `part` structure containing the members 528, "Disk drive", and 10, in that order. This structure is then passed to `print_part`, which displays it.

Here's how a compound literal might be assigned to a variable:

```
part1 = (struct part) {528, "Disk drive", 10};
```

This statement resembles a declaration containing an initializer, but it's not the same—initializers can appear only in declarations, not in statements such as this one.

In general, a compound literal consists of a type name within parentheses, followed by a set of values enclosed by braces. In the case of a compound literal that represents a structure, the type name can be a structure tag preceded by the word `struct`—as in our examples—or a `typedef` name. A compound literal may contain designators, just like a designated initializer:

```
print_part((struct part) { .on_hand = 10,
                           .name = "Disk drive",
                           .number = 528});
```

A compound literal may fail to provide full initialization, in which case any uninitialized members default to zero.

## 16.3 Nested Arrays and Structures

Structures and arrays can be combined without restriction. Arrays may have structures as their elements, and structures may contain arrays and structures as members. We've already seen an example of an array nested inside a structure (the

name member of the part structure). Let's explore the other possibilities: structures whose members are structures and arrays whose elements are structures.

## Nested Structures

Nesting one kind of structure inside another is often useful. For example, suppose that we've declared the following structure, which can store a person's first name, middle initial, and last name:

```
struct person_name {  
    char first[FIRST_NAME_LEN+1];  
    char middle_initial;  
    char last[LAST_NAME_LEN+1];  
};
```

We can use the person\_name structure as part of a larger structure:

```
struct student {  
    struct person_name name;  
    int id, age;  
    char sex;  
} student1, student2;
```

Accessing student1's first name, middle initial, or last name requires two applications of the . operator:

```
strcpy(student1.name.first, "Fred");
```

One advantage of making name a structure (instead of having first, middle\_initial, and last be members of the student structure) is that we can more easily treat names as units of data. For example, if we were to write a function that displays a name, we could pass it just one argument—a person\_name structure—instead of three arguments:

```
display_name(student1.name);
```

Likewise, copying the information from a person\_name structure to the name member of a student structure would take one assignment instead of three:

```
struct person_name new_name;  
...  
student1.name = new_name;
```

## Arrays of Structures

One of the most common combinations of arrays and structures is an array whose elements are structures. An array of this kind can serve as a simple database. For example, the following array of part structures is capable of storing information about 100 parts:

```
struct part inventory[100];
```

To access one of the parts in the array, we'd use subscripting. To print the part stored in position *i*, for example, we could write

```
print_part(inventory[i]);
```

Accessing a member within a *part* structure requires a combination of subscripting and member selection. To assign 883 to the *number* member of *inventory[i]*, we could write

```
inventory[i].number = 883;
```

Accessing a single character in a part name requires subscripting (to select a particular part), followed by selection (to select the *name* member), followed by subscripting (to select a character within the part name). To change the name stored in *inventory[i]* to an empty string, we could write

```
inventory[i].name[0] = '\0';
```

## Initializing an Array of Structures

Initializing an array of structures is done in much the same way as initializing a multidimensional array. Each structure has its own brace-enclosed initializer; the initializer for the array simply wraps another set of braces around the structure initializers.

One reason for initializing an array of structures is that we're planning to treat it as a database of information that won't change during program execution. For example, suppose that we're working on a program that will need access to the country codes used when making international telephone calls. First, we'll set up a structure that can store the name of a country along with its code:

```
struct dialing_code {
    char *country;
    int code;
};
```

Note that *country* is a pointer, not an array of characters. That could be a problem if we were planning to use *dialing\_code* structures as variables, but we're not. When we initialize a *dialing\_code* structure, *country* will end up pointing to a string literal.

Next, we'll declare an array of these structures and initialize it to contain the codes for some of the world's most populous nations:

```
const struct dialing_code country_codes[] =
{{"Argentina", 54}, {"Bangladesh", 880},
 {"Brazil", 55}, {"Burma (Myanmar)", 95},
 {"China", 86}, {"Colombia", 57},
 {"Congo, Dem. Rep. of", 243}, {"Egypt", 20},
 {"Ethiopia", 251}, {"France", 33},
 {"Germany", 49}, {"India", 91},
```

```

    {"Indonesia",           62}, {"Iran",          98},
    {"Italy",               39}, {"Japan",         81},
    {"Mexico",              52}, {"Nigeria",       234},
    {"Pakistan",             92}, {"Philippines",   63},
    {"Poland",                48}, {"Russia",        7},
    {"South Africa",        27}, {"South Korea",  82},
    {"Spain",                  34}, {"Sudan",         249},
    {"Thailand",                 66}, {"Turkey",        90},
    {"Ukraine",                380}, {"United Kingdom", 44},
    {"United States",            1}, {"Vietnam",       84} };

```

The inner braces around each structure value are optional. As a matter of style, however, I prefer not to omit them.

**C99**

Because arrays of structures (and structures containing arrays) are so common, C99's designated initializers allow an item to have more than one designator. Suppose that we want to initialize the `inventory` array to contain a single part. The part number is 528 and the quantity on hand is 10, but the name is to be left empty for now:

```

struct part inventory[100] =
{ [0].number = 528, [0].on_hand = 10, [0].name[0] = '\0' };

```

The first two items in the list use two designators (one to select array element 0—a part structure—and one to select a member within the structure). The last item uses three designators: one to select an array element, one to select the `name` member within that element, and one to select element 0 of `name`.

## PROGRAM Maintaining a Parts Database

To illustrate how nested arrays and structures are used in practice, we'll now develop a fairly long program that maintains a database of information about parts stored in a warehouse. The program is built around an array of structures, with each structure containing information—part number, name, and quantity—about one part. Our program will support the following operations:

- **Add a new part number, part name, and initial quantity on hand.** The program must print an error message if the part is already in the database or if the database is full.
- **Given a part number, print the name of the part and the current quantity on hand.** The program must print an error message if the part number isn't in the database.
- **Given a part number, change the quantity on hand.** The program must print an error message if the part number isn't in the database.
- **Print a table showing all information in the database.** Parts must be displayed in the order in which they were entered.
- **Terminate program execution.**

We'll use the codes `i` (insert), `s` (search), `u` (update), `p` (print), and `q` (quit) to represent these operations. A session with the program might look like this:

```

Enter operation code: i
Enter part number: 528
Enter part name: Disk drive
Enter quantity on hand: 10

Enter operation code: s
Enter part number: 528
Part name: Disk drive
Quantity on hand: 10

Enter operation code: s
Enter part number: 914
Part not found.

Enter operation code: i
Enter part number: 914
Enter part name: Printer cable
Enter quantity on hand: 5

Enter operation code: u
Enter part number: 528
Enter change in quantity on hand: -2

Enter operation code: s
Enter part number: 528
Part name: Disk drive
Quantity on hand: 8

Enter operation code: p
Part Number      Part Name          Quantity on Hand
      528           Disk drive            8
      914           Printer cable        5

Enter operation code: q

```

The program will store information about each part in a structure. We'll limit the size of the database to 100 parts, making it possible to store the structures in an array, which I'll call `inventory`. (If this limit proves to be too small, we can always change it later.) To keep track of the number of parts currently stored in the array, we'll use a variable named `num_parts`.

Since this program is menu-driven, it's fairly easy to sketch the main loop:

```

for (;;) {
    prompt user to enter operation code;
    read code;
    switch (code) {
        case 'i': perform insert operation; break;
        case 's': perform search operation; break;
        case 'u': perform update operation; break;
        case 'p': perform print operation; break;
    }
}

```

```

        case 'q': terminate program;
        default: print error message;
    }
}

```

It will be convenient to have separate functions perform the insert, search, update, and print operations. Since these functions will all need access to `inventory` and `num_parts`, we might want to make these variables external. As an alternative, we could declare the variables inside `main`, and then pass them to the functions as arguments. From a design standpoint, it's usually better to make variables local to a function rather than making them external (see Section 10.2 if you've forgotten why). In this program, however, putting `inventory` and `num_parts` inside `main` would merely complicate matters.

For reasons that I'll explain later, I've decided to split the program into three files: `inventory.c`, which contains the bulk of the program; `readline.h`, which contains the prototype for the `read_line` function; and `readline.c`, which contains the definition of `read_line`. We'll discuss the latter two files later in this section. For now, let's concentrate on `inventory.c`.

```

inventory.c /* Maintains a parts database (array version) */

#include <stdio.h>
#include "readline.h"

#define NAME_LEN 25
#define MAX_PARTS 100

struct part {
    int number;
    char name[NAME_LEN+1];
    int on_hand;
} inventory[MAX_PARTS];

int num_parts = 0; /* number of parts currently stored */

int find_part(int number);
void insert(void);
void search(void);
void update(void);
void print(void);

*****  

* main: Prompts the user to enter an operation code, *
*       then calls a function to perform the requested   *
*       action. Repeats until the user enters the      *
*       command 'q'. Prints an error message if the user *
*       enters an illegal code.                         *
*****/  

int main(void)
{
    char code;

```

```
        for (;;) {
            printf("Enter operation code: ");
            scanf(" %c", &code);
            while (getchar() != '\n') /* skips to end of line */
                ;
            switch (code) {
                case 'i': insert();
                            break;
                case 's': search();
                            break;
                case 'u': update();
                            break;
                case 'p': print();
                            break;
                case 'q': return 0;
                default: printf("Illegal code\n");
            }
            printf("\n");
        }

/***
 * find_part: Looks up a part number in the inventory *
 *             array. Returns the array index if the part *
 *             number is found; otherwise, returns -1. *
 *****/
int find_part(int number)
{
    int i;

    for (i = 0; i < num_parts; i++)
        if (inventory[i].number == number)
            return i;
    return -1;
}

/***
 * insert: Prompts the user for information about a new *
 *         part and then inserts the part into the *
 *         database. Prints an error message and returns *
 *         prematurely if the part already exists or the *
 *         database is full.
 *****/
void insert(void)
{
    int part_number;

    if (num_parts == MAX_PARTS) {
        printf("Database is full; can't add more parts.\n");
        return;
    }

    printf("Enter part number: ");
    scanf("%d", &part_number);
```

```
if (find_part(part_number) >= 0) {
    printf("Part already exists.\n");
    return;
}

inventory[num_parts].number = part_number;
printf("Enter part name: ");
read_line(inventory[num_parts].name, NAME_LEN);
printf("Enter quantity on hand: ");
scanf("%d", &inventory[num_parts].on_hand);
num_parts++;
}

/*****************
 * search: Prompts the user to enter a part number, then *
 *          looks up the part in the database. If the part *
 *          exists, prints the name and quantity on hand;   *
 *          if not, prints an error message.                 *
 *****************/
void search(void)
{
    int i, number;

    printf("Enter part number: ");
    scanf("%d", &number);
    i = find_part(number);
    if (i >= 0) {
        printf("Part name: %s\n", inventory[i].name);
        printf("Quantity on hand: %d\n", inventory[i].on_hand);
    } else
        printf("Part not found.\n");
}

/*****************
 * update: Prompts the user to enter a part number.      *
 *          Prints an error message if the part doesn't   *
 *          exist; otherwise, prompts the user to enter    *
 *          change in quantity on hand and updates the   *
 *          database.                                     *
 *****************/
void update(void)
{
    int i, number, change;

    printf("Enter part number: ");
    scanf("%d", &number);
    i = find_part(number);
    if (i >= 0) {
        printf("Enter change in quantity on hand: ");
        scanf("%d", &change);
        inventory[i].on_hand += change;
    } else
        printf("Part not found.\n");
}
```

```
*****
 * print: Prints a listing of all parts in the database,
 *         showing the part number, part name, and
 *         quantity on hand. Parts are printed in the
 *         order in which they were entered into the
 *         database.
*****
void print(void)
{
    int i;

    printf("Part Number      Part Name
           "Quantity on Hand\n");
    for (i = 0; i < num_parts; i++)
        printf("%7d      %-25s%11d\n", inventory[i].number,
               inventory[i].name, inventory[i].on_hand);
}
```

In the main function, the format string " %c" allows `scanf` to skip over white space before reading the operation code. The space in the format string is crucial; without it, `scanf` would sometimes read the new-line character that terminated a previous line of input.

The program contains one function, `find_part`, that isn't called from `main`. This "helper" function helps us avoid redundant code and simplify the more important functions. By calling `find_part`, the `insert`, `search`, and `update` functions can locate a part in the database (or simply determine if the part exists).

There's just one detail left: the `read_line` function, which the program uses to read the part name. Section 13.3 discussed the issues that are involved in writing such a function. Unfortunately, the version of `read_line` in that section won't work properly in the current program. Consider what happens when the user inserts a part:

```
Enter part number: 528
Enter part name: Disk drive
```

The user presses the Enter key after entering the part number and again after entering the part name, each time leaving an invisible new-line character that the program must read. For the sake of discussion, let's pretend that these characters are visible:

```
Enter part number: 528 
Enter part name: Disk drive 
```

When we call `scanf` to read the part number, it consumes the 5, 2, and 8, but leaves the   character unread. If we try to read the part name using our original `read_line` function, it will encounter the   character immediately and stop reading. This problem is common when numerical input is followed by character input. Our solution will be to write a version of `read_line` that skips white-

space characters before it begins storing characters. Not only will this solve the new-line problem, but it also allows us to avoid storing any blanks that precede the part name.

Since `read_line` is unrelated to the other functions in `inventory.c`, and since it's potentially reusable in other programs, I've decided to separate it from `inventory.c`. The prototype for `read_line` will go in the `readline.h` header file:

```
readline.h #ifndef READLINE_H
#define READLINE_H

/*********************  

 * read_line: Skips leading white-space characters, then *  

 *           reads the remainder of the input line and *  

 *           stores it in str. Truncates the line if its *  

 *           length exceeds n. Returns the number of *  

 *           characters stored.  

 *******************/  

int read_line(char str[], int n);

#endif
```

We'll put the definition of `read_line` in the `readline.c` file:

```
readline.c #include <ctype.h>
#include <stdio.h>
#include "readline.h"

int read_line(char str[], int n)
{
    int ch, i = 0;

    while (isspace(ch = getchar()))
        ;
    while (ch != '\n' && ch != EOF) {
        if (i < n)
            str[i++] = ch;
        ch = getchar();
    }
    str[i] = '\0';
    return i;
}
```

The expression

```
isspace(ch = getchar())
```

isspace function ▶ 23.5

controls the first `while` statement. This expression calls `getchar` to read a character, stores the character into `ch`, and then uses the `isspace` function to test whether `ch` is a white-space character. If not, the loop terminates with `ch` containing a character that's not white space. Section 15.3 explains why `ch` has type `int` instead of `char` and why it's good to test for EOF.

## 16.4 Unions

A *union*, like a structure, consists of one or more members, possibly of different types. However, the compiler allocates only enough space for the largest of the members, which overlay each other within this space. As a result, assigning a new value to one member alters the values of the other members as well.

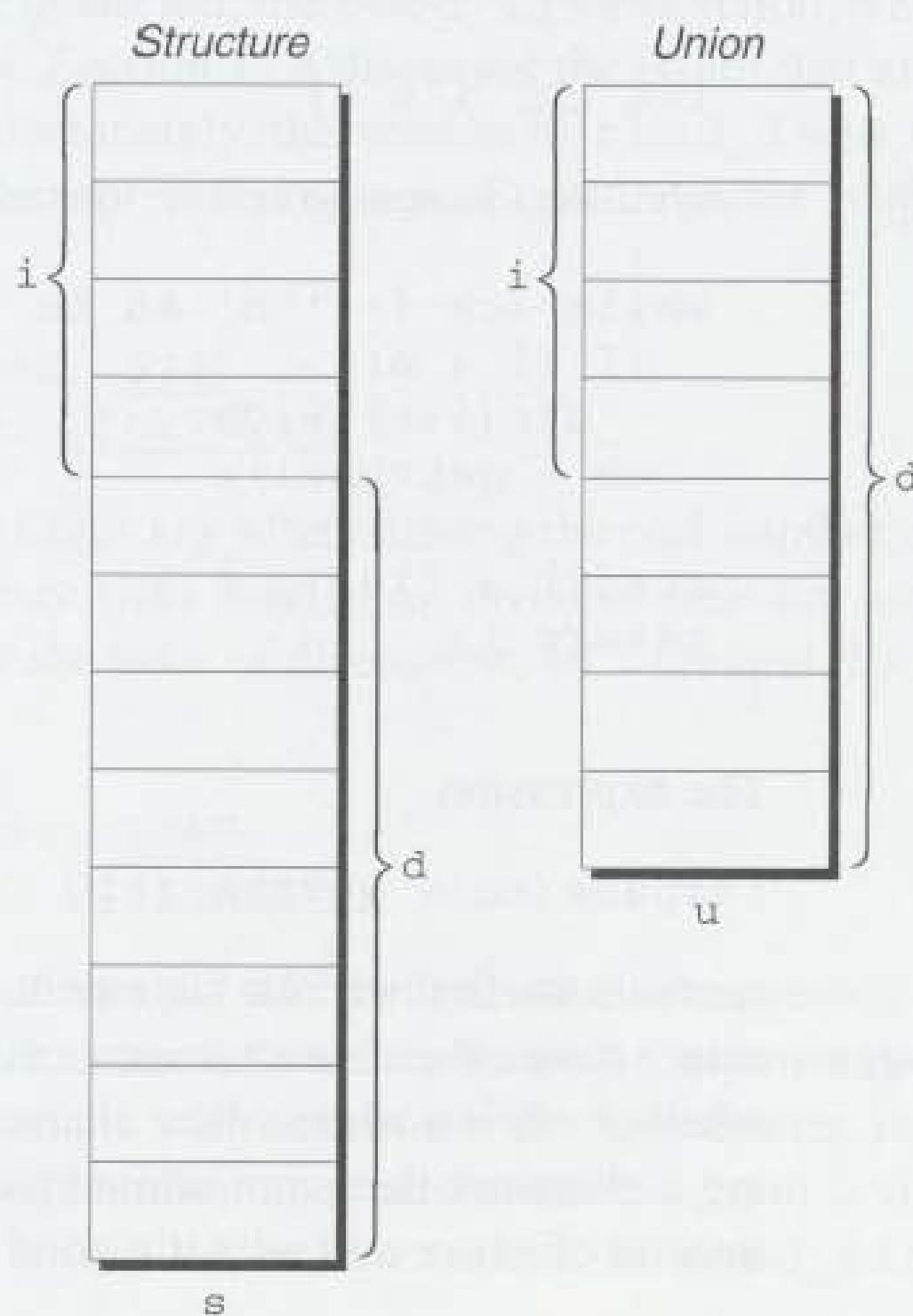
To illustrate the basic properties of unions, let's declare a union variable, *u*, with two members:

```
union {
    int i;
    double d;
} u;
```

Notice how the declaration of a union closely resembles a structure declaration:

```
struct {
    int i;
    double d;
} s;
```

In fact, the structure *s* and the union *u* differ in just one way: the members of *s* are stored at *different* addresses in memory, while the members of *u* are stored at the *same* address. Here's what *s* and *u* will look like in memory (assuming that *int* values require four bytes and *double* values take eight bytes):



In the `s` structure, `i` and `d` occupy different memory locations; the total size of `s` is 12 bytes. In the `u` union, `i` and `d` overlap (`i` is really the first four bytes of `d`), so `u` occupies only eight bytes. Also, `i` and `d` have the same address.

Members of a union are accessed in the same way as members of a structure. To store the number 82 in the `i` member of `u`, we would write

```
u.i = 82;
```

To store the value 74.8 in the `d` member, we would write

```
u.d = 74.8;
```

Since the compiler overlays storage for the members of a union, changing one member alters any value previously stored in any of the other members. Thus, if we store a value in `u.d`, any value previously stored in `u.i` will be lost. (If we examine the value of `u.i`, it will appear to be meaningless.) Similarly, changing `u.i` corrupts `u.d`. Because of this property, we can think of `u` as a place to store either `i` or `d`, not both. (The structure `s` allows us to store `i` and `d`.)

The properties of unions are almost identical to the properties of structures. We can declare union tags and union types in the same way we declare structure tags and types. Like structures, unions can be copied using the `=` operator, passed to functions, and returned by functions.

Unions can even be initialized in a manner similar to structures. However, only the first member of a union can be given an initial value. For example, we can initialize the `i` member of `u` to 0 in the following way:

```
union {
    int i;
    double d;
} u = {0};
```

Notice the presence of the braces, which are required. The expression inside the braces must be constant. (The rules are slightly different in C99, as we'll see in Section 18.5.)

**C99**

Designated initializers, a C99 feature that we've previously discussed in the context of arrays and structures, can also be used with unions. A designated initializer allows us to specify which member of a union should be initialized. For example, we can initialize the `d` member of `u` as follows:

```
union {
    int i;
    double d;
} u = {.d = 10.0};
```

Only one member can be initialized, but it doesn't have to be the first one.

There are several applications for unions. We'll discuss two of these now. Another application—viewing storage in different ways—is highly machine-dependent, so I'll postpone it until Section 20.3.

## Using Unions to Save Space

We'll often use unions as a way to save space in structures. Suppose that we're designing a structure that will contain information about an item that's sold through a gift catalog. The catalog carries only three kinds of merchandise: books, mugs, and shirts. Each item has a stock number and a price, as well as other information that depends on the type of the item:

*Books*: Title, author, number of pages

*Mugs*: Design

*Shirts*: Design, colors available, sizes available

Our first design attempt might result in the following structure:

```
struct catalog_item {
    int stock_number;
    double price;
    int item_type;
    char title[TITLE_LEN+1];
    char author[AUTHOR_LEN+1];
    int num_pages;
    char design[DESIGN_LEN+1];
    int colors;
    int sizes;
};
```

The `item_type` member would have one of the values `BOOK`, `MUG`, or `SHIRT`. The `colors` and `sizes` members would store encoded combinations of colors and sizes.

Although this structure is perfectly usable, it wastes space, since only part of the information in the structure is common to all items in the catalog. If an item is a book, for example, there's no need to store design, colors, and sizes. By putting a union inside the `catalog_item` structure, we can reduce the space required by the structure. The members of the union will be structures, each containing the data that's needed for a particular kind of catalog item:

```
struct catalog_item {
    int stock_number;
    double price;
    int item_type;
    union {
        struct {
            char title[TITLE_LEN+1];
            char author[AUTHOR_LEN+1];
            int num_pages;
        } book;
        struct {
            char design[DESIGN_LEN+1];
        } mug;
```

```

struct {
    char design[DESIGN_LEN+1];
    int colors;
    int sizes;
} shirt;
} item;
};

```

Notice that the union (named `item`) is a member of the `catalog_item` structure, and the `book`, `mug`, and `shirt` structures are members of `item`. If `c` is a `catalog_item` structure that represents a book, we can print the book's title in the following way:

```
printf("%s", c.item.book.title);
```

As this example shows, accessing a union that's nested inside a structure can be awkward: to locate a book title, we had to specify the name of a structure (`c`), the name of the union member of the structure (`item`), the name of a structure member of the union (`book`), and then the name of a member of that structure (`title`).

We can use the `catalog_item` structure to illustrate an interesting aspect of unions. Normally, it's not a good idea to store a value into one member of a union and then access the data through a different member, because assigning to one member of a union causes the values of the other members to be undefined. However, the C standard mentions a special case: two or more of the members of the union are structures, and the structures begin with one or more matching members. (These members need to be in the same order and have compatible types, but need not have the same name.) If one of the structures is currently valid, then the matching members in the other structures will also be valid.

Consider the union embedded in the `catalog_item` structure. It contains three structures as members, two of which (`mug` and `shirt`) begin with a matching member (`design`). Now, suppose that we assign a value to one of the `design` members:

```
strcpy(c.item.mug.design, "Cats");
```

The `design` member in the other structure will be defined and have the same value:

```
printf("%s", c.item.shirt.design); /* prints "Cats" */
```

## Using Unions to Build Mixed Data Structures

Unions have another important application: creating data structures that contain a mixture of data of different types. Let's say that we need an array whose elements are a mixture of `int` and `double` values. Since the elements of an array must be of the same type, it seems impossible to create such an array. Using unions, though, it's relatively easy. First, we define a union type whose members represent the different kinds of data to be stored in the array:

```
typedef union {
    int i;
    double d;
} Number;
```

Next, we create an array whose elements are `Number` values:

```
Number number_array[1000];
```

Each element of `number_array` is a `Number` union. A `Number` union can store either an `int` value or a `double` value, making it possible to store a mixture of `int` and `double` values in `number_array`. For example, suppose that we want element 0 of `number_array` to store 5, while element 1 stores 8.395. The following assignments will have the desired effect:

```
number_array[0].i = 5;
number_array[1].d = 8.395;
```

## Adding a “Tag Field” to a Union

Unions suffer from a major problem: there’s no easy way to tell which member of a union was last changed and therefore contains a meaningful value. Consider the problem of writing a function that displays the value currently stored in a `Number` union. This function might have the following outline:

```
void print_number(Number n)
{
    if (n contains an integer)
        printf("%d", n.i);
    else
        printf("%g", n.d);
}
```

Unfortunately, there’s no way for `print_number` to determine whether `n` contains an integer or a floating-point number.

In order to keep track of this information, we can embed the union within a structure that has one other member: a “tag field” or “discriminant,” whose purpose is to remind us what’s currently stored in the union. In the `catalog_item` structure discussed earlier in this section, `item_type` served this purpose.

Let’s convert the `Number` type into a structure with an embedded union:

```
#define INT_KIND 0
#define DOUBLE_KIND 1

typedef struct {
    int kind; /* tag field */
    union {
        int i;
        double d;
    } u;
} Number;
```

`Number` has two members, `kind` and `u`. The value of `kind` will be either `INT_KIND` or `DOUBLE_KIND`.

Each time we assign a value to a member of `u`, we'll also change `kind` to remind us which member of `u` we modified. For example, if `n` is a `Number` variable, an assignment to the `i` member of `u` would have the following appearance:

```
n.kind = INT_KIND;
n.u.i = 82;
```

Notice that assigning to `i` requires that we first select the `u` member of `n`, then the `i` member of `u`.

When we need to retrieve the number stored in a `Number` variable, `kind` will tell us which member of the union was the last to be assigned a value. The `print_number` function can take advantage of this capability:

```
void print_number(Number n)
{
    if (n.kind == INT_KIND)
        printf("%d", n.u.i);
    else
        printf("%g", n.u.d);
}
```



It's the program's responsibility to change the tag field each time an assignment is made to a member of the union.

## 16.5 Enumerations

In many programs, we'll need variables that have only a small set of meaningful values. A Boolean variable, for example, should have only two possible values: "true" and "false." A variable that stores the suit of a playing card should have only four potential values: "clubs," "diamonds," "hearts," and "spades." The obvious way to deal with such a variable is to declare it as an integer and have a set of codes that represent the possible values of the variable:

```
int s; /* s will store a suit */
...
s = 2; /* 2 represents "hearts" */
```

Although this technique works, it leaves much to be desired. Someone reading the program can't tell that `s` has only four possible values, and the significance of 2 isn't immediately apparent.

Using macros to define a suit "type" and names for the various suits is a step in the right direction:

```
#define SUIT      int
#define CLUBS     0
#define DIAMONDS  1
#define HEARTS    2
#define SPADES    3
```

Our previous example now becomes easier to read:

```
SUIT s;
...
s = HEARTS;
```

This technique is an improvement, but it's still not the best solution. There's no indication to someone reading the program that the macros represent values of the same "type." If the number of possible values is more than a few, defining a separate macro for each will be tedious. Moreover, the names we've defined—CLUBS, DIAMONDS, HEARTS, and SPADES—will be removed by the preprocessor, so they won't be available during debugging.

C provides a special kind of type designed specifically for variables that have a small number of possible values. An *enumerated type* is a type whose values are listed ("enumerated") by the programmer, who must create a name (an *enumeration constant*) for each of the values. The following example enumerates the values (CLUBS, DIAMONDS, HEARTS, and SPADES) that can be assigned to the variables `s1` and `s2`:

```
enum {CLUBS, DIAMONDS, HEARTS, SPADES} s1, s2;
```

Although enumerations have little in common with structures and unions, they're declared in a similar way. Unlike the members of a structure or union, however, the names of enumeration constants must be different from other identifiers declared in the enclosing scope.

Enumeration constants are similar to constants created with the `#define` directive, but they're not equivalent. For one thing, enumeration constants are subject to C's scope rules: if an enumeration is declared inside a function, its constants won't be visible outside the function.

## Enumeration Tags and Type Names

We'll often need to create names for enumerations, for the same reasons that we name structures and unions. As with structures and unions, there are two ways to name an enumeration: by declaring a tag or by using `typedef` to create a genuine type name.

Enumeration tags resemble structure and union tags. To define the tag `suit`, for example, we could write

```
enum suit {CLUBS, DIAMONDS, HEARTS, SPADES};
```

`suit` variables would be declared in the following way:

```
enum suit s1, s2;
```

As an alternative, we could use `typedef` to make `Suit` a type name:

```
typedef enum {CLUBS, DIAMONDS, HEARTS, SPADES} Suit;
Suit s1, s2;
```

In C89, using `typedef` to name an enumeration is an excellent way to create a Boolean type:

```
typedef enum {FALSE, TRUE} Bool;
```

C99 has a built-in Boolean type, of course, so there's no need for a C99 programmer to define a `Bool` type in this way.

## Enumerations as Integers

Behind the scenes, C treats enumeration variables and constants as integers. By default, the compiler assigns the integers 0, 1, 2, ... to the constants in a particular enumeration. In our `suit` enumeration, for example, CLUBS, DIAMONDS, HEARTS, and SPADES represent 0, 1, 2, and 3, respectively.

We're free to choose different values for enumeration constants if we like. Let's say that we want CLUBS, DIAMONDS, HEARTS, and SPADES to stand for 1, 2, 3, and 4. We can specify these numbers when declaring the enumeration:

```
enum suit {CLUBS = 1, DIAMONDS = 2, HEARTS = 3, SPADES = 4};
```

The values of enumeration constants may be arbitrary integers, listed in no particular order:

```
enum dept {RESEARCH = 20, PRODUCTION = 10, SALES = 25};
```

It's even legal for two or more enumeration constants to have the same value.

When no value is specified for an enumeration constant, its value is one greater than the value of the previous constant. (The first enumeration constant has the value 0 by default.) In the following enumeration, BLACK has the value 0, LT\_GRAY is 7, DK\_GRAY is 8, and WHITE is 15:

```
enum EGA_colors {BLACK, LT_GRAY = 7, DK_GRAY, WHITE = 15};
```

Since enumeration values are nothing but thinly disguised integers, C allows us to mix them with ordinary integers:

```
int i;
enum {CLUBS, DIAMONDS, HEARTS, SPADES} s;

i = DIAMONDS;      /* i is now 1          */
s = 0;              /* s is now 0 (CLUBS)   */
s++;                /* s is now 1 (DIAMONDS) */
i = s + 2;          /* i is now 3          */
```

The compiler treats `s` as a variable of some integer type; CLUBS, DIAMONDS, HEARTS, and SPADES are just names for the integers 0, 1, 2, and 3.



Although it's convenient to be able to use an enumeration value as an integer, it's dangerous to use an integer as an enumeration value. For example, we might accidentally store the number 4—which doesn't correspond to any suit—into `s`.

## Using Enumerations to Declare “Tag Fields”

Enumerations are perfect for solving a problem that we encountered in Section 16.4: determining which member of a union was the last to be assigned a value. In the `Number` structure, for example, we can make the `kind` member an enumeration instead of an `int`:

```
typedef struct {
    enum {INT_KIND, DOUBLE_KIND} kind;
    union {
        int i;
        double d;
    } u;
} Number;
```

The new structure is used in exactly the same way as the old one. The advantages are that we've done away with the `INT_KIND` and `DOUBLE_KIND` macros (they're now enumeration constants), and we've clarified the meaning of `kind`—it's now obvious that `kind` has only two possible values: `INT_KIND` and `DOUBLE_KIND`.

## Q & A

**Q:** When I tried using the `sizeof` operator to determine the number of bytes in a structure, I got a number that was larger than the sizes of the members added together. How can this be?

**A:** Let's look at an example:

```
struct {
    char a;
    int b;
} s;
```

If `char` values occupy one byte and `int` values occupy four bytes, how large is `s`? The obvious answer—five bytes—may not be the correct one. Some computers require that the address of certain data items be a multiple of some number of bytes (typically two, four, or eight, depending on the item's type). To satisfy this requirement, a compiler will “align” the members of a structure by leaving “holes” (unused bytes) between adjacent members. If we assume that data items must

begin on a multiple of four bytes, the `a` member of the `s` structure will be followed by a three-byte hole. As a result, `sizeof(s)` will be 8.

By the way, a structure can have a hole at the end, as well as holes between members. For example, the structure

```
struct {
    int a;
    char b;
} s;
```

might have a three-byte hole after the `b` member.

**Q: Can there be a “hole” at the beginning of a structure?**

A: No. The C standard specifies that holes are allowed only *between* members or *after* the last member. One consequence is that a pointer to the first member of a structure is guaranteed to be the same as a pointer to the entire structure. (Note, however, that the two pointers won’t have the same type.)

**Q: Why isn’t it legal to use the == operator to test whether two structures are equal? [p. 382]**

A: This operation was left out of C because there’s no way to implement it that would be consistent with the language’s philosophy. Comparing structure members one by one would be too inefficient. Comparing all bytes in the structures would be better (many computers have special instructions that can perform such a comparison rapidly). If the structures contain holes, however, comparing bytes could yield an incorrect answer; even if corresponding members have identical values, leftover data stored in the holes might be different. The problem could be solved by having the compiler ensure that holes always contain the same value (zero, say). Initializing holes would impose a performance penalty on all programs that use structures, however, so it’s not feasible.

**Q: Why does C provide two ways to name structure types (tags and `typedef` names)? [p. 382]**

A: C originally lacked `typedef`, so tags were the only technique available for naming structure types. When `typedef` was added, it was too late to remove tags. Besides, a tag is still necessary when a member of a structure points to a structure of the same type (see the `node` structure of Section 17.5).

**Q: Can a structure have both a tag and a `typedef` name? [p. 384]**

A: Yes. In fact, the tag and the `typedef` name can even be the same, although that’s not required:

```
typedef struct part {
    int number;
    char name[NAME_LEN+1];
    int on_hand;
} part;
```

**Q: How can I share a structure type among several files in a program?**

A: Put a declaration of the structure tag (or a `typedef`, if you prefer) in a header file, then include the header file where the structure is needed. To share the `part` structure, for example, we'd put the following lines in a header file:

```
struct part {
    int number;
    char name[NAME_LEN+1];
    int on_hand;
};
```

Notice that we're declaring only the structure *tag*, not variables of this type.

Incidentally, a header file that contains a declaration of a structure tag or structure type may need protection against multiple inclusion. Declaring a tag or `typedef` name twice in the same file is an error. Similar remarks apply to unions and enumerations.

**Q: If I include the declaration of the part structure into two different files, will part variables in one file be of the same type as part variables in the other file?**

A: Technically, no. However, the C standard says that the `part` variables in one file have a type that's compatible with the type of the `part` variables in the other file. Variables with compatible types can be assigned to each other, so there's little practical difference between types being "compatible" and being "the same."

**C99** The rules for structure compatibility in C89 and C99 are slightly different. In C89, structures defined in different files are compatible if their members have the same names and appear in the same order, with corresponding members having compatible types. C99 goes one step further: it requires that either both structures have the same tag or neither has a tag.

Similar compatibility rules apply to unions and enumerations (with the same difference between C89 and C99).

**Q: Is it legal to have a pointer to a compound literal?**

A: Yes. Consider the `print_part` function of Section 16.2. Currently, the parameter to this function is a `part` structure. The function would be more efficient if it were modified to accept a *pointer* to a `part` structure instead. Using the function to print a compound literal would then be done by prefixing the argument with the & (address) operator:

```
print_part(&(struct part) {528, "Disk drive", 10});
```

**Q: Allowing a pointer to a compound literal would seem to make it possible to modify the literal. Is that the case?**

**C99** A: Yes. Compound literals are lvalues that can be modified, although doing so is rare.

**Q: I saw a program in which the last constant in an enumeration was followed by a comma, like this:**

```
enum gray_values {
    BLACK = 0,
    DARK_GRAY = 64,
    GRAY = 128,
    LIGHT_GRAY = 192,
};
```

### Is this practice legal?

- A: This practice is indeed legal in C99 (and is supported by some pre-C99 compilers as well). Allowing a “trailing comma” makes enumerations easier to modify, because we can add a constant to the end of an enumeration without changing existing lines of code. For example, we might want to add WHITE to our enumeration:

```
enum gray_values {
    BLACK = 0,
    DARK_GRAY = 64,
    GRAY = 128,
    LIGHT_GRAY = 192,
    WHITE = 255,
};
```

The comma after the definition of LIGHT\_GRAY makes it easy to add WHITE to the end of the list.

- C99 One reason for this change is that C89 allows trailing commas in initializers, so it seemed inconsistent not to allow the same flexibility in enumerations. Incidentally, C99 also allows trailing commas in compound literals.

### Q: Can the values of an enumerated type be used as subscripts?

- A: Yes, indeed. They are integers and have—by default—values that start at 0 and count upward, so they make great subscripts. In C99, moreover, enumeration constants can be used as subscripts in designated initializers. Here’s an example:

```
enum weekdays {MONDAY, TUESDAY, WEDNESDAY, THURSDAY, FRIDAY};
const char *daily_specials[] = {
    [MONDAY] = "Beef ravioli",
    [TUESDAY] = "BLTs",
    [WEDNESDAY] = "Pizza",
    [THURSDAY] = "Chicken fajitas",
    [FRIDAY] = "Macaroni and cheese"
};
```

## Exercises

### Section 16.1

- In the following declarations, the x and y structures have members named x and y:

```
struct { int x, y; } x;
struct { int x, y; } y;
```

Are these declarations legal on an individual basis? Could both declarations appear as shown in a program? Justify your answer.

- W 2. (a) Declare structure variables named `c1`, `c2`, and `c3`, each having members `real` and `imaginary` of type `double`.  
 (b) Modify the declaration in part (a) so that `c1`'s members initially have the values 0.0 and 1.0, while `c2`'s members are 1.0 and 0.0 initially. (`c3` is not initialized.)  
 (c) Write statements that copy the members of `c2` into `c1`. Can this be done in one statement, or does it require two?  
 (d) Write statements that add the corresponding members of `c1` and `c2`, storing the result in `c3`.

**Section 16.2**

3. (a) Show how to declare a tag named `complex` for a structure with two members, `real` and `imaginary`, of type `double`.  
 (b) Use the `complex` tag to declare variables named `c1`, `c2`, and `c3`.  
 (c) Write a function named `make_complex` that stores its two arguments (both of type `double`) in a `complex` structure, then returns the structure.  
 (d) Write a function named `add_complex` that adds the corresponding members of its arguments (both `complex` structures), then returns the result (another `complex` structure).  
 W 4. Repeat Exercise 3, but this time using a *type* named `Complex`.  
 5. Write the following functions, assuming that the `date` structure contains three members: `month`, `day`, and `year` (all of type `int`).  
 (a) `int day_of_year(struct date d);`  
 Returns the day of the year (an integer between 1 and 366) that corresponds to the date `d`.  
 (b) `int compare_dates(struct date d1, struct date d2);`  
 Returns -1 if `d1` is an earlier date than `d2`, +1 if `d1` is a later date than `d2`, and 0 if `d1` and `d2` are the same.  
 6. Write the following function, assuming that the `time` structure contains three members: `hours`, `minutes`, and `seconds` (all of type `int`).  
`struct time split_time(long total_seconds);`  
`total_seconds` is a time represented as the number of seconds since midnight. The function returns a structure containing the equivalent time in hours (0–23), minutes (0–59), and seconds (0–59).  
 7. Assume that the `fraction` structure contains two members: `numerator` and `denominator` (both of type `int`). Write functions that perform the following operations on fractions:  
 (a) Reduce the fraction `f` to lowest terms. *Hint:* To reduce a fraction to lowest terms, first compute the greatest common divisor (GCD) of the numerator and denominator. Then divide both the numerator and denominator by the GCD.  
 (b) Add the fractions `f1` and `f2`.  
 (c) Subtract the fraction `f2` from the fraction `f1`.  
 (d) Multiply the fractions `f1` and `f2`.  
 (e) Divide the fraction `f1` by the fraction `f2`.

The fractions `f`, `f1`, and `f2` will be arguments of type `struct fraction`; each function will return a value of type `struct fraction`. The fractions returned by the functions in parts (b)–(e) should be reduced to lowest terms. *Hint:* You may use the function from part (a) to help write the functions in parts (b)–(e).

8. Let `color` be the following structure:

```
struct color {
    int red;
    int green;
    int blue;
};
```

(a) Write a declaration for a `const` variable named `MAGENTA` of type `struct color` whose members have the values 255, 0, and 255, respectively.

(b) (C99) Repeat part (a), but use a designated initializer that doesn't specify the value of `green`, allowing it to default to 0.

9. Write the following functions. (The `color` structure is defined in Exercise 8.)

(a) `struct color make_color(int red, int green, int blue);`

Returns a `color` structure containing the specified red, green, and blue values. If any argument is less than zero, the corresponding member of the structure will contain zero instead. If any argument is greater than 255, the corresponding member of the structure will contain 255.

(b) `int getRed(struct color c);`

Returns the value of `c`'s `red` member.

(c) `bool equal_color(struct color color1, struct color color2);`

Returns `true` if the corresponding members of `color1` and `color2` are equal.

(d) `struct color brighter(struct color c);`

Returns a `color` structure that represents a brighter version of the color `c`. The structure is identical to `c`, except that each member has been divided by 0.7 (with the result truncated to an integer). However, there are three special cases: (1) If all members of `c` are zero, the function returns a color whose members all have the value 3. (2) If any member of `c` is greater than 0 but less than 3, it is replaced by 3 before the division by 0.7. (3) If dividing by 0.7 causes a member to exceed 255, it is reduced to 255.

(e) `struct color darker(struct color c);`

Returns a `color` structure that represents a darker version of the color `c`. The structure is identical to `c`, except that each member has been multiplied by 0.7 (with the result truncated to an integer).

### Section 16.3

10. The following structures are designed to store information about objects on a graphics screen:

```
struct point { int x, y; };
struct rectangle { struct point upper_left, lower_right; };
```

A `point` structure stores the `x` and `y` coordinates of a point on the screen. A `rectangle` structure stores the coordinates of the upper left and lower right corners of a rectangle. Write functions that perform the following operations on a `rectangle` structure `r` passed as an argument:

(a) Compute the area of `r`.

(b) Compute the center of `r`, returning it as a `point` value. If either the `x` or `y` coordinate of the center isn't an integer, store its truncated value in the `point` structure.

(c) Move `r` by `x` units in the `x` direction and `y` units in the `y` direction, returning the modified version of `r`. (`x` and `y` are additional arguments to the function.)

(d) Determine whether a point `p` lies within `r`, returning `true` or `false`. (`p` is an additional argument of type `struct point`.)

**Section 16.4** **W** 11. Suppose that `s` is the following structure:

```
struct {
    double a;
    union {
        char b[4];
        double c;
        int d;
    } e;
    char f[4];
} s;
```

If `char` values occupy one byte, `int` values occupy four bytes, and `double` values occupy eight bytes, how much space will a C compiler allocate for `s`? (Assume that the compiler leaves no “holes” between members.)

12. Suppose that `u` is the following union:

```
union {
    double a;
    struct {
        char b[4];
        double c;
        int d;
    } e;
    char f[4];
} u;
```

If `char` values occupy one byte, `int` values occupy four bytes, and `double` values occupy eight bytes, how much space will a C compiler allocate for `u`? (Assume that the compiler leaves no “holes” between members.)

13. Suppose that `s` is the following structure (`point` is a structure tag declared in Exercise 10):

```
struct shape {
    int shape_kind;          /* RECTANGLE or CIRCLE */
    struct point center;    /* coordinates of center */
    union {
        struct {
            int height, width;
        } rectangle;
        struct {
            int radius;
        } circle;
    } u;
} s;
```

If the value of `shape_kind` is `RECTANGLE`, the `height` and `width` members store the dimensions of a rectangle. If the value of `shape_kind` is `CIRCLE`, the `radius` member stores the radius of a circle. Indicate which of the following statements are legal, and show how to repair the ones that aren’t:

- `s.shape_kind = RECTANGLE;`
- `s.center.x = 10;`
- `s.height = 25;`
- `s.u.rectangle.width = 8;`
- `s.u.circle = 5;`
- `s.u.radius = 5;`

- W 14. Let `shape` be the structure tag declared in Exercise 13. Write functions that perform the following operations on a `shape` structure `s` passed as an argument:
- Compute the area of `s`.
  - Move `s` by `x` units in the `x` direction and `y` units in the `y` direction, returning the modified version of `s`. (`x` and `y` are additional arguments to the function.)
  - Scale `s` by a factor of `c` (a `double` value), returning the modified version of `s`. (`c` is an additional argument to the function.)

**Section 16.5**

- W 15. (a) Declare a tag for an enumeration whose values represent the seven days of the week.  
 (b) Use `typedef` to define a name for the enumeration of part (a).
16. Which of the following statements about enumeration constants are true?
- An enumeration constant may represent any integer specified by the programmer.
  - Enumeration constants have exactly the same properties as constants created using `#define`.
  - Enumeration constants have the values 0, 1, 2, ... by default.
  - All constants in an enumeration must have different values.
  - Enumeration constants may be used as integers in expressions.

- W 17. Suppose that `b` and `i` are declared as follows:

```
enum { FALSE, TRUE } b;
int i;
```

Which of the following statements are legal? Which ones are “safe” (always yield a meaningful result)?

- `b = FALSE;`
  - `b = i;`
  - `b++;`
  - `i = b;`
  - `i = 2 * b + 1;`
18. (a) Each square of a chessboard can hold one piece—a pawn, knight, bishop, rook, queen, or king—or it may be empty. Each piece is either black or white. Define two enumerated types: `Piece`, which has seven possible values (one of which is “empty”), and `Color`, which has two.
- (b) Using the types from part (a), define a structure type named `Square` that can store both the type of a piece and its color.
- (c) Using the `Square` type from part (b), declare an  $8 \times 8$  array named `board` that can store the entire contents of a chessboard.
- (d) Add an initializer to the declaration in part (c) so that `board`’s initial value corresponds to the usual arrangement of pieces at the start of a chess game. A square that’s not occupied by a piece should have an “empty” piece value and the color black.
19. Declare a structure with the following members whose tag is `pinball_machine`:
- `name` – a string of up to 40 characters
  - `year` – an integer (representing the year of manufacture)
  - `type` – an enumeration with the values `EM` (electromechanical) and `SS` (solid state)
  - `players` – an integer (representing the maximum number of players)
20. Suppose that the `direction` variable is declared in the following way:
- ```
enum { NORTH, SOUTH, EAST, WEST } direction;
```

Let `x` and `y` be `int` variables. Write a switch statement that tests the value of `direction`, incrementing `x` if `direction` is EAST, decrementing `x` if `direction` is WEST, incrementing `y` if `direction` is SOUTH, and decrementing `y` if `direction` is NORTH.

21. What are the integer values of the enumeration constants in each of the following declarations?
  - (a) `enum {NUL, SOH, STX, ETX};`
  - (b) `enum {VT = 11, FF, CR};`
  - (c) `enum {SO = 14, SI, DLE, CAN = 24, EM};`
  - (d) `enum {ENQ = 45, ACK, BEL, LF = 37, ETB, ESC};`
22. Let `chess_pieces` be the following enumeration:
 

```
enum chess_pieces {KING, QUEEN, ROOK, BISHOP, KNIGHT, PAWN};
```

  - (a) Write a declaration (including an initializer) for a constant array of integers named `piece_value` that stores the numbers 200, 9, 5, 3, 3, and 1, representing the value of each chess piece, from king to pawn. (The king's value is actually infinite, since "capturing" the king (checkmate) ends the game, but some chess-playing software assigns the king a large value such as 200.)
  - (b) (C99) Repeat part (a), but use a designated initializer to initialize the array. Use the enumeration constants in `chess_pieces` as subscripts in the designators. (*Hint:* See the last question in Q&A for an example.)

## Programming Projects

- W 1. Write a program that asks the user to enter an international dialing code and then looks it up in the `country_codes` array (see Section 16.3). If it finds the code, the program should display the name of the corresponding country; if not, the program should print an error message.
- 2. Modify the `inventory.c` program of Section 16.3 so that the `p` (print) operation displays the parts sorted by part number.
- W 3. Modify the `inventory.c` program of Section 16.3 by making `inventory` and `num_parts` local to the `main` function.
- 4. Modify the `inventory.c` program of Section 16.3 by adding a `price` member to the `part` structure. The `insert` function should ask the user for the price of a new item. The `search` and `print` functions should display the price. Add a new command that allows the user to change the price of a part.
- 5. Modify Programming Project 8 from Chapter 5 so that the times are stored in a single array. The elements of the array will be structures, each containing a departure time and the corresponding arrival time. (Each time will be an integer, representing the number of minutes since midnight.) The program will use a loop to search the array for the departure time closest to the time entered by the user.
- 6. Modify Programming Project 9 from Chapter 5 so that each date entered by the user is stored in a `date` structure (see Exercise 5). Incorporate the `compare_dates` function of Exercise 5 into your program.

# 17 Advanced Uses of Pointers

*One can only display complex information in the mind.  
Like seeing, movement or flow or alteration of view is more  
important than the static picture, no matter how lovely.*

In previous chapters, we've seen two important uses of pointers. Chapter 11 showed how using a pointer to a variable as a function argument allows the function to modify the variable. Chapter 12 showed how to process arrays by performing arithmetic on pointers to array elements. This chapter completes our coverage of pointers by examining two additional applications: dynamic storage allocation and pointers to functions.

Using dynamic storage allocation, a program can obtain blocks of memory as needed during execution. Section 17.1 explains the basics of dynamic storage allocation. Section 17.2 discusses dynamically allocated strings, which provide more flexibility than ordinary character arrays. Section 17.3 covers dynamic storage allocation for arrays in general. Section 17.4 deals with the issue of storage deallocation—releasing blocks of dynamically allocated memory when they're no longer needed.

Dynamically allocated structures play a big role in C programming, since they can be linked together to form lists, trees, and other highly flexible data structures. Section 17.5 focuses on linked lists, the most fundamental linked data structure. One of the issues that arises in this section—the concept of a “pointer to a pointer”—is important enough to warrant a section of its own (Section 17.6).

Section 17.7 introduces pointers to functions, a surprisingly useful concept. Some of C's most powerful library functions expect function pointers as arguments. We'll examine one of these functions, `qsort`, which is capable of sorting any array.

The last two sections discuss pointer-related features that first appeared in C99: restricted pointers (Section 17.8) and flexible array members (Section 17.9). These features are primarily of interest to advanced C programmers, so both sections can be safely be skipped by the beginner.

## 17.1 Dynamic Storage Allocation

variable-length arrays ▶ 8.3

C's data structures are normally fixed in size. For example, the number of elements in an array is fixed once the program has been compiled. (In C99, the length of a variable-length array is determined at run time, but it remains fixed for the rest of the array's lifetime.) Fixed-size data structures can be a problem, since we're forced to choose their sizes when writing a program; we can't change the sizes without modifying the program and compiling it again.

Consider the inventory program of Section 16.3, which allows the user to add parts to a database. The database is stored in an array of length 100. To enlarge the capacity of the database, we can increase the size of the array and recompile the program. But no matter how large we make the array, there's always the possibility that it will fill up. Fortunately, all is not lost. C supports *dynamic storage allocation*: the ability to allocate storage during program execution. Using dynamic storage allocation, we can design data structures that grow (and shrink) as needed.

Although it's available for all types of data, dynamic storage allocation is used most often for strings, arrays, and structures. Dynamically allocated structures are of particular interest, since we can link them together to form lists, trees, and other data structures.

### Memory Allocation Functions

<stdlib.h> header ▶ 26.2

To allocate storage dynamically, we'll need to call one of the three memory allocation functions declared in the `<stdlib.h>` header:

- `malloc`—Allocates a block of memory but doesn't initialize it.
- `calloc`—Allocates a block of memory and clears it.
- `realloc`—Resizes a previously allocated block of memory.

Of the three, `malloc` is the most used. It's more efficient than `calloc`, since it doesn't have to clear the memory block that it allocates.

When we call a memory allocation function to request a block of memory, the function has no idea what type of data we're planning to store in the block, so it can't return a pointer to an ordinary type such as `int` or `char`. Instead, the function returns a value of type `void *`. A `void *` value is a "generic" pointer—essentially, just a memory address.

### Null Pointers

When a memory allocation function is called, there's always a possibility that it won't be able to locate a block of memory large enough to satisfy our request. If

that should happen, the function will return a **null pointer**. A null pointer is a “pointer to nothing”—a special value that can be distinguished from all valid pointers. After we’ve stored the function’s return value in a pointer variable, we must test to see if it’s a null pointer.



It’s the programmer’s responsibility to test the return value of any memory allocation function and take appropriate action if it’s a null pointer. The effect of attempting to access memory through a null pointer is undefined; the program may crash or behave unpredictably.

**Q&A**

The null pointer is represented by a macro named NULL, so we can test malloc’s return value in the following way:

```
p = malloc(10000);
if (p == NULL) {
    /* allocation failed; take appropriate action */
}
```

Some programmers combine the call of malloc with the NULL test:

```
if ((p = malloc(10000)) == NULL) {
    /* allocation failed; take appropriate action */
}
```

**C99**

The NULL macro is defined in six headers: `<locale.h>`, `<stddef.h>`, `<stdio.h>`, `<stdlib.h>`, `<string.h>`, and `<time.h>`. (The C99 header `<wchar.h>` also defines NULL.) As long as one of these headers is included in a program, the compiler will recognize NULL. A program that uses any of the memory allocation functions will include `<stdlib.h>`, of course, making NULL available.

In C, pointers test true or false in the same way as numbers. All non-null pointers test true; only null pointers are false. Thus, instead of writing

```
if (p == NULL) ...
```

we could write

```
if (!p) ...
```

and instead of writing

```
if (p != NULL) ...
```

we could write

```
if (p) ...
```

As a matter of style, I prefer the explicit comparison with NULL.

## 17.2 Dynamically Allocated Strings

Dynamic storage allocation is often useful for working with strings. Strings are stored in character arrays, and it can be hard to anticipate how long these arrays need to be. By allocating strings dynamically, we can postpone the decision until the program is running.

### Using `malloc` to Allocate Memory for a String

The `malloc` function has the following prototype:

```
void *malloc(size_t size);
```

size\_t type ▶ 7.6 `malloc` allocates a block of `size` bytes and returns a pointer to it. Note that `size` has type `size_t`, an unsigned integer type defined in the C library. Unless we're allocating a very large block of memory, we can just think of `size` as an ordinary integer.

Using `malloc` to allocate memory for a string is easy, because C guarantees that a `char` value requires exactly one byte of storage (`sizeof(char)` is 1, in other words). To allocate space for a string of `n` characters, we'd write

```
p = malloc(n + 1);
```

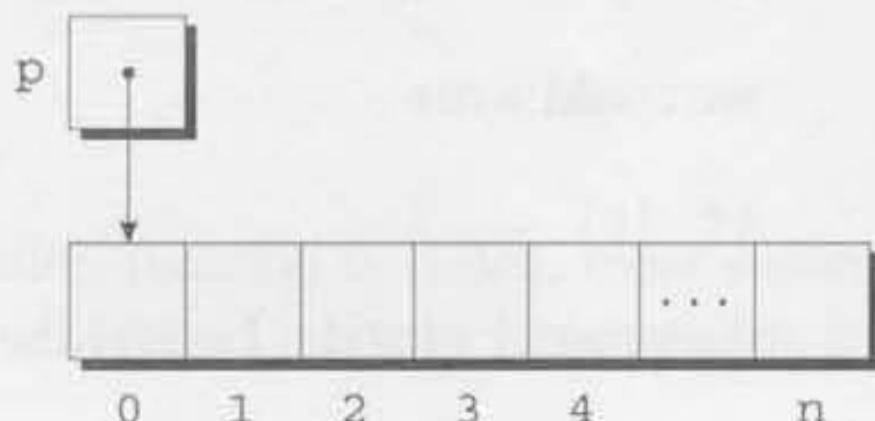
where `p` is a `char *` variable. (The argument is `n + 1` rather than `n` to allow room for the null character.) The generic pointer that `malloc` returns will be converted to `char *` when the assignment is performed; no cast is necessary. (In general, we can assign a `void *` value to a variable of any pointer type and vice versa.) Nevertheless, some programmers prefer to cast `malloc`'s return value:

```
p = (char *) malloc(n + 1);
```



When using `malloc` to allocate space for a string, don't forget to include room for the null character.

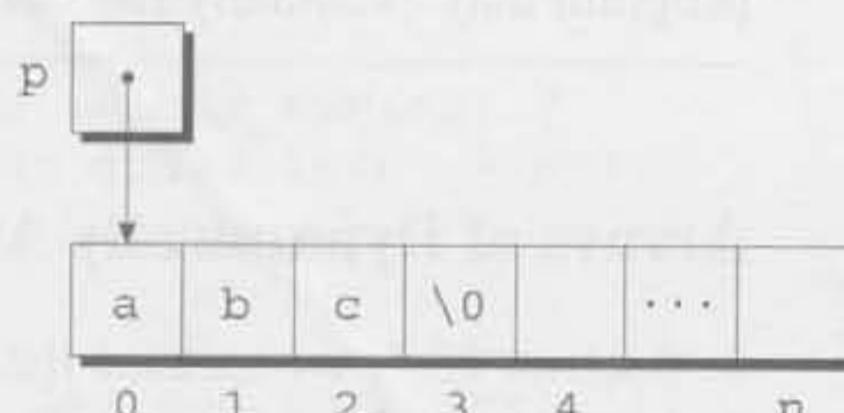
Memory allocated using `malloc` isn't cleared or initialized in any way, so `p` will point to an uninitialized array of `n + 1` characters:



Calling `strcpy` is one way to initialize this array:

```
strcpy(p, "abc");
```

The first four characters in the array will now be a, b, c, and \0:



## Using Dynamic Storage Allocation in String Functions

Dynamic storage allocation makes it possible to write functions that return a pointer to a “new” string—a string that didn’t exist before the function was called. Consider the problem of writing a function that concatenates two strings without changing either one. C’s standard library doesn’t include such a function (`strcat` isn’t quite what we want, since it modifies one of the strings passed to it), but we can easily write our own.

Our function will measure the lengths of the two strings to be concatenated, then call `malloc` to allocate just the right amount of space for the result. The function next copies the first string into the new space and then calls `strcat` to concatenate the second string.

```
char *concat(const char *s1, const char *s2)
{
    char *result;

    result = malloc(strlen(s1) + strlen(s2) + 1);
    if (result == NULL) {
        printf("Error: malloc failed in concat\n");
        exit(EXIT_FAILURE);
    }
    strcpy(result, s1);
    strcat(result, s2);
    return result;
}
```

If `malloc` returns a null pointer, `concat` prints an error message and terminates the program. That’s not always the right action to take; some programs need to recover from memory allocation failures and continue running.

Here’s how the `concat` function might be called:

```
p = concat("abc", "def");
```

After the call, `p` will point to the string "abcdef", which is stored in a dynamically allocated array. The array is seven characters long, including the null character at the end.



free function ▶ 17.4

Functions such as `concat` that dynamically allocate storage must be used with care. When the string that `concat` returns is no longer needed, we'll want to call the `free` function to release the space that the string occupies. If we don't, the program may eventually run out of memory.

## Arrays of Dynamically Allocated Strings

In Section 13.7, we tackled the problem of storing strings in an array. We found that storing strings as rows in a two-dimensional array of characters can waste space, so we tried setting up an array of pointers to string literals. The techniques of Section 13.7 work just as well if the elements of an array are pointers to dynamically allocated strings. To illustrate this point, let's rewrite the `remind.c` program of Section 13.5, which prints a one-month list of daily reminders.

### PROGRAM

## Printing a One-Month Reminder List (Revisited)

The original `remind.c` program stores the reminder strings in a two-dimensional array of characters, with each row of the array containing one string. After the program reads a day and its associated reminder, it searches the array to determine where the day belongs, using `strcmp` to do comparisons. It then uses `strcpy` to move all strings below that point down one position. Finally, the program copies the day into the array and calls `strcat` to append the reminder to the day.

In the new program (`remind2.c`), the array will be one-dimensional; its elements will be pointers to dynamically allocated strings. Switching to dynamically allocated strings in this program will have two primary advantages. First, we can use space more efficiently by allocating the exact number of characters needed to store a reminder, rather than storing the reminder in a fixed number of characters as the original program does. Second, we won't need to call `strcpy` to move existing reminder strings in order to make room for a new reminder. Instead, we'll merely move *pointers* to strings.

Here's the new program, with changes in **bold**. Switching from a two-dimensional array to an array of pointers turns out to be remarkably easy: we'll only need to change eight lines of the program.

```
remind2.c /* Prints a one-month reminder list (dynamic string version) */

#include <stdio.h>
#include <stdlib.h>
#include <string.h>

#define MAX_REMIND 50 /* maximum number of reminders */
#define MSG_LEN 60 /* max length of reminder message */

int read_line(char str[], int n);
```

```
int main(void)
{
    char *reminders[MAX_REMIND];
    char day_str[3], msg_str[MSG_LEN+1];
    int day, i, j, num_remind = 0;

    for (;;) {
        if (num_remind == MAX_REMIND) {
            printf("-- No space left --\n");
            break;
        }

        printf("Enter day and reminder: ");
        scanf("%2d", &day);
        if (day == 0)
            break;
        sprintf(day_str, "%2d", day);
        read_line(msg_str, MSG_LEN);

        for (i = 0; i < num_remind; i++)
            if (strcmp(day_str, reminders[i]) < 0)
                break;
        for (j = num_remind; j > i; j--)
            reminders[j] = reminders[j-1];

        reminders[i] = malloc(2 + strlen(msg_str) + 1);
        if (reminders[i] == NULL) {
            printf("-- No space left --\n");
            break;
        }

        strcpy(reminders[i], day_str);
        strcat(reminders[i], msg_str);

        num_remind++;
    }

    printf("\nDay Reminder\n");
    for (i = 0; i < num_remind; i++)
        printf(" %s\n", reminders[i]);

    return 0;
}

int read_line(char str[], int n)
{
    int ch, i = 0;

    while ((ch = getchar()) != '\n')
        if (i < n)
            str[i++] = ch;
    str[i] = '\0';
    return i;
}
```

## 17.3 Dynamically Allocated Arrays

Dynamically allocated arrays have the same advantages as dynamically allocated strings (not surprisingly, since strings *are* arrays). When we’re writing a program, it’s often difficult to estimate the proper size for an array; it would be more convenient to wait until the program is run to decide how large the array should be. C solves this problem by allowing a program to allocate space for an array during execution, then access the array through a pointer to its first element. The close relationship between arrays and pointers, which we explored in Chapter 12, makes a dynamically allocated array just as easy to use as an ordinary array.

Although `malloc` can allocate space for an array, the `calloc` function is sometimes used instead, since it initializes the memory that it allocates. The `realloc` function allows us to make an array “grow” or “shrink” as needed.

### Using `malloc` to Allocate Storage for an Array

sizeof operator ▶ 7.6

We can use `malloc` to allocate space for an array in much the same way we used it to allocate space for a string. The primary difference is that the elements of an arbitrary array won’t necessarily be one byte long, as they are in a string. As a result, we’ll need to use the `sizeof` operator to calculate the amount of space required for each element.

Suppose we’re writing a program that needs an array of  $n$  integers, where  $n$  is to be computed during the execution of the program. We’ll first declare a pointer variable:

```
int *a;
```

Once the value of  $n$  is known, we’ll have the program call `malloc` to allocate space for the array:

```
a = malloc(n * sizeof(int));
```



Always use `sizeof` when calculating how much space is needed for an array. Failing to allocate enough memory can have severe consequences. Consider the following attempt to allocate space for an array of  $n$  integers:

```
a = malloc(n * 2);
```

If `int` values are larger than two bytes (as they are on most computers), `malloc` won’t allocate a large enough block of memory. When we later try to access elements of the array, the program may crash or behave erratically.

---

Once it points to a dynamically allocated block of memory, we can ignore the fact that `a` is a pointer and use it instead as an array name, thanks to the relation-

ship between arrays and pointers in C. For example, we could use the following loop to initialize the array that `a` points to:

```
for (i = 0; i < n; i++)
    a[i] = 0;
```

We also have the option of using pointer arithmetic instead of subscripting to access the elements of the array.

## The `calloc` Function

Although the `malloc` function can be used to allocate memory for an array, C provides an alternative—the `calloc` function—that's sometimes better. `calloc` has the following prototype in `<stdlib.h>`:

```
void *calloc(size_t nmemb, size_t size);
```

**Q&A** `calloc` allocates space for an array with `nmemb` elements, each of which is `size` bytes long; it returns a null pointer if the requested space isn't available. After allocating the memory, `calloc` initializes it by setting all bits to 0. For example, the following call of `calloc` allocates space for an array of `n` integers, which are all guaranteed to be zero initially:

```
a = calloc(n, sizeof(int));
```

Since `calloc` clears the memory that it allocates but `malloc` doesn't, we may occasionally want to use `calloc` to allocate space for an object other than an array. By calling `calloc` with 1 as its first argument, we can allocate space for a data item of any type:

```
struct point { int x, y; } *p;
p = calloc(1, sizeof(struct point));
```

After this statement has been executed, `p` will point to a structure whose `x` and `y` members have been set to zero.

## The `realloc` Function

Once we've allocated memory for an array, we may later find that it's too large or too small. The `realloc` function can resize the array to better suit our needs. The following prototype for `realloc` appears in `<stdlib.h>`:

```
void *realloc(void *ptr, size_t size);
```

When `realloc` is called, `ptr` must point to a memory block obtained by a previous call of `malloc`, `calloc`, or `realloc`. The `size` parameter represents the new size of the block, which may be larger or smaller than the original size. Although `realloc` doesn't require that `ptr` point to memory that's being used as an array, in practice it usually does.



Be sure that a pointer passed to `realloc` came from a previous call of `malloc`, `calloc`, or `realloc`. If it didn't, calling `realloc` causes undefined behavior.

The C standard spells out a number of rules concerning the behavior of `realloc`:

- When it expands a memory block, `realloc` doesn't initialize the bytes that are added to the block.
- If `realloc` can't enlarge the memory block as requested, it returns a null pointer; the data in the old memory block is unchanged.
- If `realloc` is called with a null pointer as its first argument, it behaves like `malloc`.
- If `realloc` is called with 0 as its second argument, it frees the memory block.

The C standard stops short of specifying exactly how `realloc` works. Still, we expect it to be reasonably efficient. When asked to reduce the size of a memory block, `realloc` should shrink the block "in place," without moving the data stored in the block. By the same token, `realloc` should always attempt to expand a memory block without moving it. If it's unable to enlarge the block (because the bytes following the block are already in use for some other purpose), `realloc` will allocate a new block elsewhere, then copy the contents of the old block into the new one.



Once `realloc` has returned, be sure to update all pointers to the memory block, since it's possible that `realloc` has moved the block elsewhere.

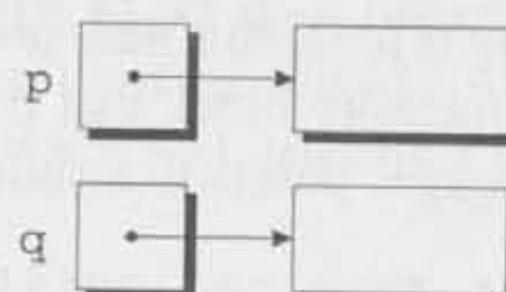
## 17.4 Deallocating Storage

`malloc` and the other memory allocation functions obtain memory blocks from a storage pool known as the *heap*. Calling these functions too often—or asking them for large blocks of memory—can exhaust the heap, causing the functions to return a null pointer.

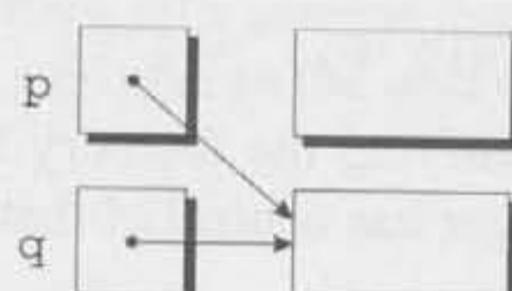
To make matters worse, a program may allocate blocks of memory and then lose track of them, thereby wasting space. Consider the following example:

```
p = malloc(...);  
q = malloc(...);  
p = q;
```

After the first two statements have been executed, p points to one memory block, while q points to another:



After q is assigned to p, both variables now point to the second memory block:



There are no pointers to the first block (shaded), so we'll never be able to use it again.

A block of memory that's no longer accessible to a program is said to be *garbage*. A program that leaves garbage behind has a *memory leak*. Some languages provide a *garbage collector* that automatically locates and recycles garbage, but C doesn't. Instead, each C program is responsible for recycling its own garbage by calling the `free` function to release unneeded memory.

## The `free` Function

The `free` function has the following prototype in `<stdlib.h>`:

```
void free(void *ptr);
```

Using `free` is easy; we simply pass it a pointer to a memory block that we no longer need:

```
p = malloc(...);
q = malloc(...);
free(p);
p = q;
```

Calling `free` releases the block of memory that p points to. This block is now available for reuse in subsequent calls of `malloc` or other memory allocation functions.




---

The argument to `free` must be a pointer that was previously returned by a memory allocation function. (The argument may also be a null pointer, in which case the call of `free` has no effect.) Passing `free` a pointer to any other object (such as a variable or array element) causes undefined behavior.

---

## The “Dangling Pointer” Problem

Although the `free` function allows us to reclaim memory that’s no longer needed, using it leads to a new problem: *dangling pointers*. The call `free(p)` deallocates the memory block that `p` points to, but doesn’t change `p` itself. If we forget that `p` no longer points to a valid memory block, chaos may ensue:

```
char *p = malloc(4);  
...  
free(p);  
...  
strcpy(p, "abc");    /*** WRONG ***/
```

Modifying the memory that `p` points to is a serious error, since our program no longer has control of that memory.



Attempting to access or modify a deallocated memory block causes undefined behavior. Trying to modify a deallocated memory block is likely to have disastrous consequences that may include a program crash.

Dangling pointers can be hard to spot, since several pointers may point to the same block of memory. When the block is freed, all the pointers are left dangling.

## 17.5 Linked Lists

Dynamic storage allocation is especially useful for building lists, trees, graphs, and other linked data structures. We’ll look at linked lists in this section; a discussion of other linked data structures is beyond the scope of this book. For more information, consult a book such as Robert Sedgewick’s *Algorithms in C, Parts 1–4: Fundamentals, Data Structures, Sorting, Searching*, Third Edition (Reading, Mass.: Addison-Wesley, 1998).

A *linked list* consists of a chain of structures (called *nodes*), with each node containing a pointer to the next node in the chain:



The last node in the list contains a null pointer, shown here as a diagonal line.

In previous chapters, we’ve used an array whenever we’ve needed to store a collection of data items; linked lists give us an alternative. A linked list is more flexible than an array: we can easily insert and delete nodes in a linked list, allowing the list to grow and shrink as needed. On the other hand, we lose the “random access” capability of an array. Any element of an array can be accessed in the same

amount of time; accessing a node in a linked list is fast if the node is close to the beginning of the list, slow if it's near the end.

This section describes how to set up a linked list in C. It also shows how to perform several common operations on linked lists: inserting a node at the beginning of a list, searching for a node, and deleting a node.

## Declaring a Node Type

To set up a linked list, the first thing we'll need is a structure that represents a single node in the list. For simplicity, let's assume that a node contains nothing but an integer (the node's data) plus a pointer to the next node in the list. Here's what our node structure will look like:

```
struct node {
    int value;           /* data stored in the node */
    struct node *next;  /* pointer to the next node */
};
```

Notice that the `next` member has type `struct node *`, which means that it can store a pointer to a node structure. There's nothing special about the name `node`, by the way; it's just an ordinary structure tag.

One aspect of the `node` structure deserves special mention. As Section 16.2 explained, we normally have the option of using either a tag or a `typedef` name to define a name for a particular kind of structure. However, when a structure has a member that points to the same kind of structure, as `node` does, we're required to use a structure tag. Without the `node` tag, we'd have no way to declare the type of `next`.

Now that we have the `node` structure declared, we'll need a way to keep track of where the list begins. In other words, we'll need a variable that always points to the first node in the list. Let's name the variable `first`:

```
struct node *first = NULL;
```

Setting `first` to `NULL` indicates that the list is initially empty.

## Creating a Node

As we construct a linked list, we'll want to create nodes one by one, adding each to the list. Creating a node requires three steps:

1. Allocate memory for the node.
2. Store data in the node.
3. Insert the node into the list.

We'll concentrate on the first two steps for now.

When we create a node, we'll need a variable that can point to the node temporarily, until it's been inserted into the list. Let's call this variable `new_node`:

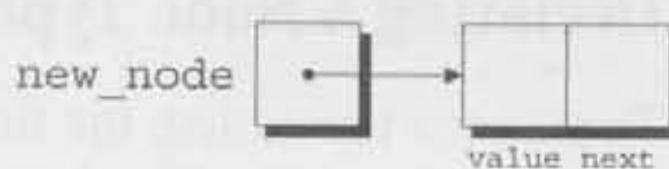
```
struct node *new_node;
```

### Q&A

We'll use `malloc` to allocate memory for the new node, saving the return value in `new_node`:

```
new_node = malloc(sizeof(struct node));
```

`new_node` now points to a block of memory just large enough to hold a node structure:



Be careful to give `sizeof` the name of the *type* to be allocated, not the name of a *pointer* to that type:

```
new_node = malloc(sizeof(new_node));    /*** WRONG ***/
```

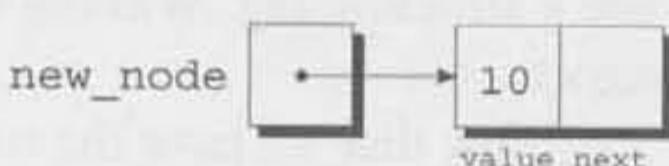
The program will still compile, but `malloc` will allocate only enough memory for a *pointer* to a node structure. The likely result is a crash later, when the program attempts to store data in the node that `new_node` is presumably pointing to.

## Q&A

Next, we'll store data in the `value` member of the new node:

```
(*new_node).value = 10;
```

Here's how the picture will look after this assignment:



To access the `value` member of the node, we've applied the indirection operator `*` (to reference the structure to which `new_node` points), then the selection operator `.` (to select a member of the structure). The parentheses around `*new_node` are mandatory because the `.` operator would otherwise take precedence over the `*` operator.

table of operators ➤ Appendix A

## The `->` Operator

Before we go on to the next step, inserting a new node into a list, let's take a moment to discuss a useful shortcut. Accessing a member of a structure using a pointer is so common that C provides a special operator just for this purpose. This operator, known as **right arrow selection**, is a minus sign followed by `>`. Using the `->` operator, we can write

```
new_node->value = 10;
```

instead of

```
(*new_node).value = 10;
```

The `->` operator is a combination of the `*` and `.` operators; it performs indirection on `new_node` to locate the structure that it points to, then selects the `value` member of the structure.

values ► 4.2

The `->` operator produces an lvalue, so we can use it wherever an ordinary variable would be allowed. We've just seen an example in which `new_node->value` appears on the left side of an assignment. It could just as easily appear in a call of `scanf`:

```
scanf("%d", &new_node->value);
```

Notice that the `&` operator is still required, even though `new_node` is a pointer. Without the `&`, we'd be passing `scanf` the *value* of `new_node->value`, which has type `int`.

## Inserting a Node at the Beginning of a Linked List

One of the advantages of a linked list is that nodes can be added at any point in the list: at the beginning, at the end, or anywhere in the middle. The beginning of a list is the easiest place to insert a node, however, so let's focus on that case.

If `new_node` is pointing to the node to be inserted, and `first` is pointing to the first node in the linked list, then we'll need two statements to insert the node into the list. First, we'll modify the new node's `next` member to point to the node that was previously at the beginning of the list:

```
new_node->next = first;
```

Second, we'll make `first` point to the new node:

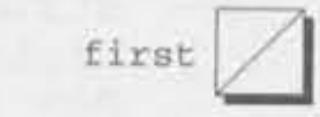
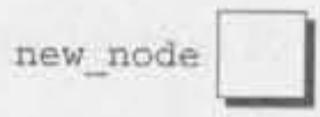
```
first = new_node;
```

Will these statements work if the list is empty when we insert a node? Yes, fortunately. To make sure this is true, let's trace the process of inserting two nodes into an empty list. We'll insert a node containing the number 10 first, followed by a node containing 20. In the figures that follow, null pointers are shown as diagonal lines.

```
first = NULL;
```



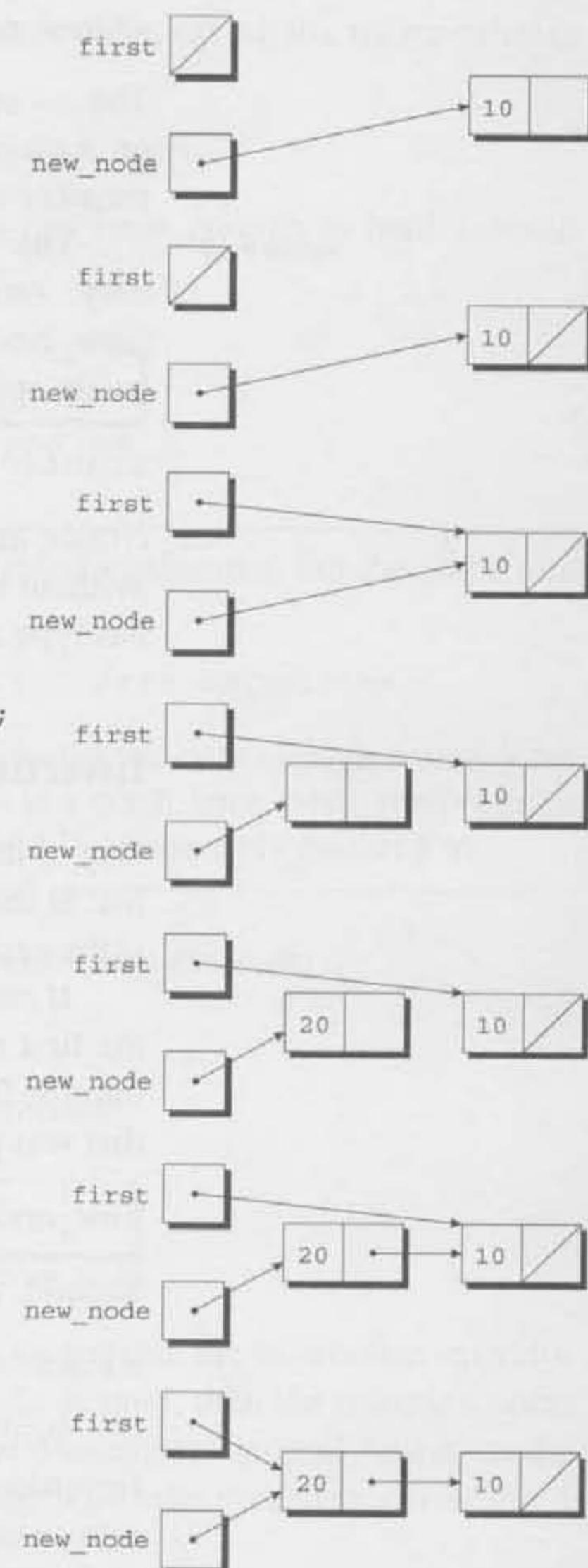
```
new_node = malloc(sizeof(struct node));
```



```

new_node->value = 10;
new_node->next = first;
first = new_node;
new_node = malloc(sizeof(struct node));
new_node->value = 20;

```



Inserting a node into a linked list is such a common operation that we'll probably want to write a function for that purpose. Let's name the function `add_to_list`. It will have two parameters: `list` (a pointer to the first node in the old list) and `n` (the integer to be stored in the new node).

```

struct node *add_to_list(struct node *list, int n)
{
    struct node *new_node;
    new_node = malloc(sizeof(struct node));
    if (new_node == NULL) {
        printf("Error: malloc failed in add_to_list\n");
        exit(EXIT_FAILURE);
    }
}

```

```

    new_node->value = n;
    new_node->next = list;
    return new_node;
}

```

Note that `add_to_list` doesn't modify the `list` pointer. Instead, it returns a pointer to the newly created node (now at the beginning of the list). When we call `add_to_list`, we'll need to store its return value into `first`:

```

first = add_to_list(first, 10);
first = add_to_list(first, 20);

```

These statements add nodes containing 10 and 20 to the list pointed to by `first`. Getting `add_to_list` to update `first` directly, rather than return a new value for `first`, turns out to be tricky. We'll return to this issue in Section 17.6.

The following function uses `add_to_list` to create a linked list containing numbers entered by the user:

```

struct node *read_numbers(void)
{
    struct node *first = NULL;
    int n;

    printf("Enter a series of integers (0 to terminate): ");
    for (;;) {
        scanf("%d", &n);
        if (n == 0)
            return first;
        first = add_to_list(first, n);
    }
}

```

The numbers will be in reverse order within the list, since `first` always points to the node containing the last number entered.

## Searching a Linked List

Once we've created a linked list, we may need to search it for a particular piece of data. Although a `while` loop can be used to search a list, the `for` statement is often superior. We're accustomed to using the `for` statement when writing loops that involve counting, but its flexibility makes the `for` statement suitable for other tasks as well, including operations on linked lists. Here's the customary way to visit the nodes in a linked list, using a pointer variable `p` to keep track of the "current" node:

**idiom** `for (p = first; p != NULL; p = p->next)`

...

The assignment

```
p = p->next
```

advances the `p` pointer from one node to the next. An assignment of this form is invariably used in C when writing a loop that traverses a linked list.

Let's write a function named `search_list` that searches a list (pointed to by the parameter `list`) for an integer `n`. If it finds `n`, `search_list` will return a pointer to the node containing `n`; otherwise, it will return a null pointer. Our first version of `search_list` relies on the "list-traversal" idiom:

```
struct node *search_list(struct node *list, int n)
{
    struct node *p;

    for (p = list; p != NULL; p = p->next)
        if (p->value == n)
            return p;
    return NULL;
}
```

Of course, there are many other ways to write `search_list`. One alternative would be to eliminate the `p` variable, instead using `list` itself to keep track of the current node:

```
struct node *search_list(struct node *list, int n)
{
    for (; list != NULL; list = list->next)
        if (list->value == n)
            return list;
    return NULL;
}
```

Since `list` is a copy of the original `list` pointer, there's no harm in changing it within the function.

Another alternative is to combine the `list->value == n` test with the `list != NULL` test:

```
struct node *search_list(struct node *list, int n)
{
    for (; list != NULL && list->value != n; list = list->next)
        ;
    return list;
}
```

Since `list` is `NULL` if we reach the end of the list, returning `list` is correct even if we don't find `n`. This version of `search_list` might be a bit clearer if we used a `while` statement:

```
struct node *search_list(struct node *list, int n)
{
    while (list != NULL && list->value != n)
        list = list->next;
    return list;
}
```

## Deleting a Node from a Linked List

A big advantage of storing data in a linked list is that we can easily delete nodes that we no longer need. Deleting a node, like creating a node, involves three steps:

1. Locate the node to be deleted.
2. Alter the previous node so that it “bypasses” the deleted node.
3. Call `free` to reclaim the space occupied by the deleted node.

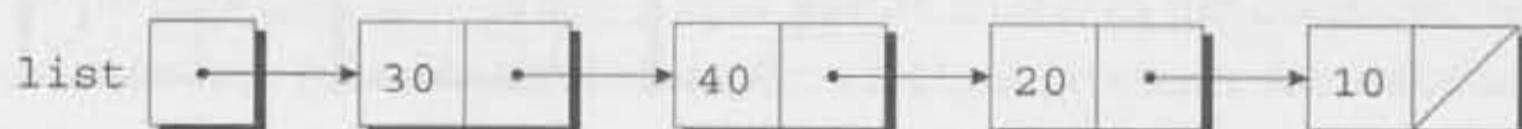
Step 1 is harder than it looks. If we search the list in the obvious way, we’ll end up with a pointer to the node to be deleted. Unfortunately, we won’t be able to perform step 2, which requires changing the *previous* node.

There are various solutions to this problem. We’ll use the “trailing pointer” technique: as we search the list in step 1, we’ll keep a pointer to the previous node (`prev`) as well as a pointer to the current node (`cur`). If `list` points to the list to be searched and `n` is the integer to be deleted, the following loop implements step 1:

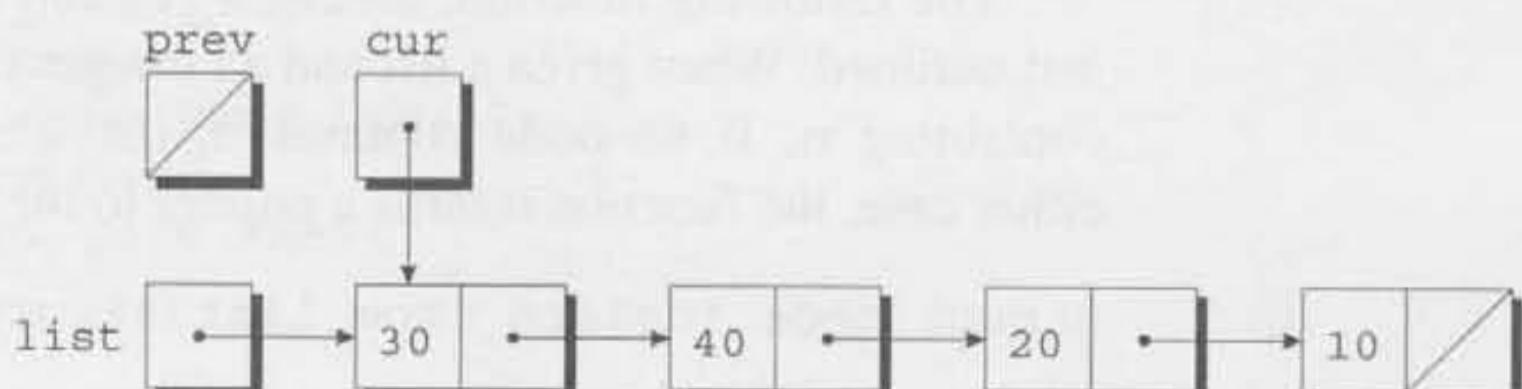
```
for (cur = list, prev = NULL;
     cur != NULL && cur->value != n;
     prev = cur, cur = cur->next)
;
```

Here we see the power of C’s `for` statement. This rather exotic example, with its empty body and liberal use of the comma operator, performs all the actions needed to search for `n`. When the loop terminates, `cur` points to the node to be deleted, while `prev` points to the previous node (if there is one).

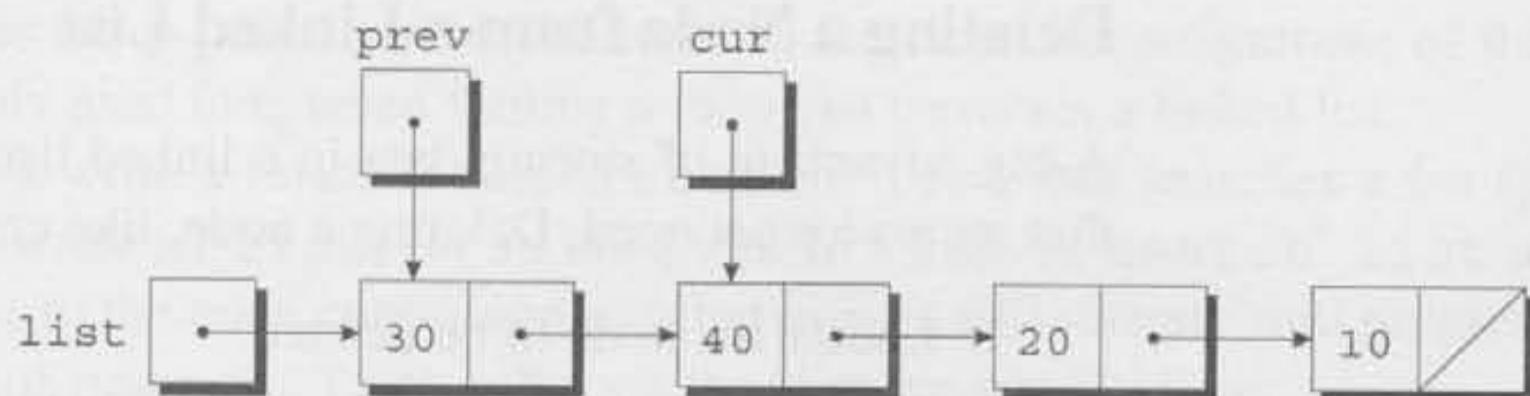
To see how this loop works, let’s assume that `list` points to a list containing 30, 40, 20, and 10, in that order:



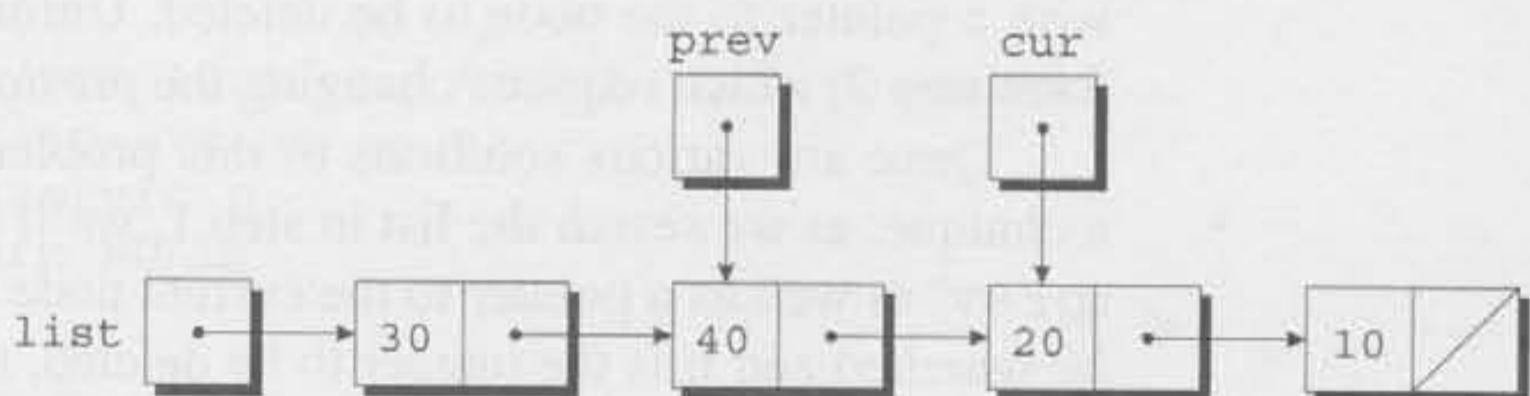
Let’s say that `n` is 20, so our goal is to delete the third node in the list. After `cur = list, prev = NULL` has been executed, `cur` points to the first node in the list:



The test `cur != NULL && cur->value != n` is true, since `cur` is pointing to a node and the node doesn’t contain 20. After `prev = cur, cur = cur->next` has been executed, we begin to see how the `prev` pointer will trail behind `cur`:



Again, the test `cur != NULL && cur->value != n` is true, so `prev = cur`, `cur = cur->next` is executed once more:

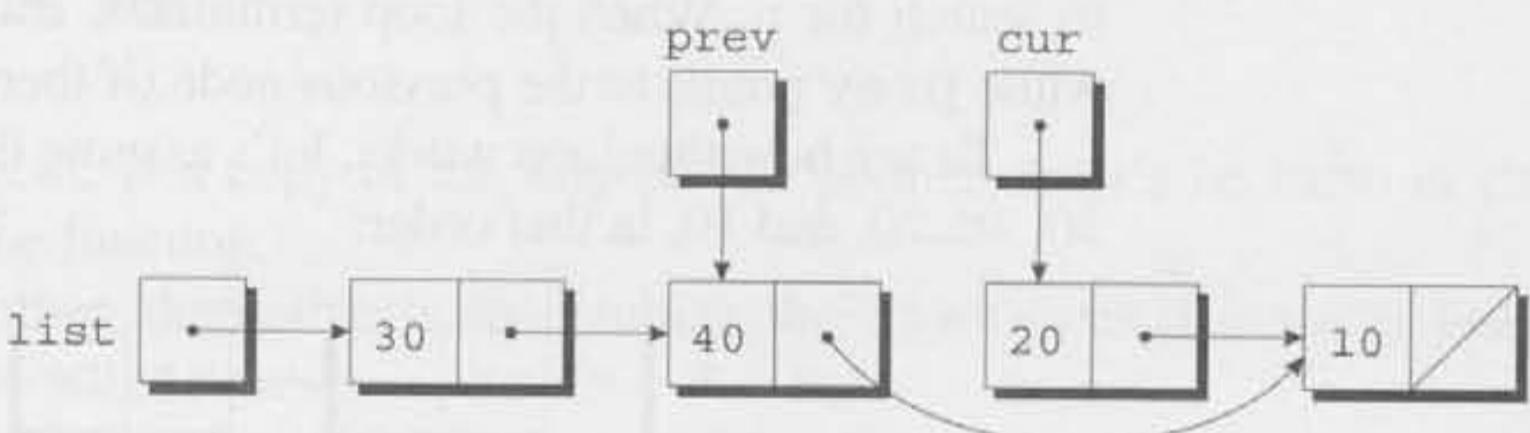


Since `cur` now points to the node containing 20, the condition `cur->value != n` is false and the loop terminates.

Next, we'll perform the bypass required by step 2. The statement

```
prev->next = cur->next;
```

makes the pointer in the previous node point to the node *after* the current node:



We're now ready for step 3, releasing the memory occupied by the current node:

```
free(cur);
```

The following function, `delete_from_list`, uses the strategy that we've just outlined. When given a list and an integer `n`, the function deletes the first node containing `n`. If no node contains `n`, `delete_from_list` does nothing. In either case, the function returns a pointer to the list.

```
struct node *delete_from_list(struct node *list, int n)
{
    struct node *cur, *prev;

    for (cur = list, prev = NULL;
         cur != NULL && cur->value != n;
         prev = cur, cur = cur->next)
    ;
```

```

    if (cur == NULL)
        return list;                      /* n was not found */
    if (prev == NULL)
        list = list->next;                /* n is in the first node */
    else
        prev->next = cur->next;          /* n is in some other node */
        free(cur);
    return list;
}

```

Deleting the first node in the list is a special case. The `prev == NULL` test checks for this case, which requires a different bypass step.

## Ordered Lists

When the nodes of a list are kept in order—sorted by the data stored inside the nodes—we say that the list is *ordered*. Inserting a node into an ordered list is more difficult (the node won’t always be put at the beginning of the list), but searching is faster (we can stop looking after reaching the point at which the desired node would have been located). The following program illustrates both the increased difficulty of inserting a node and the faster search.

### PROGRAM Maintaining a Parts Database (Revisited)

Let’s redo the parts database program of Section 16.3, this time storing the database in a linked list. Using a linked list instead of an array has two major advantages: (1) We don’t need to put a preset limit on the size of the database; it can grow until there’s no more memory to store parts. (2) We can easily keep the database sorted by part number—when a new part is added to the database, we simply insert it in its proper place in the list. In the original program, the database wasn’t sorted.

In the new program, the `part` structure will contain an additional member (a pointer to the next node in the linked list), and the variable `inventory` will be a pointer to the first node in the list:

```

struct part {
    int number;
    char name[NAME_LEN+1];
    int on_hand;
    struct part *next;
};

struct part *inventory = NULL; /* points to first part */

```

Most of the functions in the new program will closely resemble their counterparts in the original program. The `find_part` and `insert` functions will be more complex, however, since we’ll keep the nodes in the `inventory` list sorted by part number.

In the original program, `find_part` returns an index into the `inventory` array. In the new program, `find_part` will return a pointer to the node that contains the desired part number. If it doesn't find the part number, `find_part` will return a null pointer. Since the `inventory` list is sorted by part number, the new version of `find_part` can save time by stopping its search when it finds a node containing a part number that's greater than or equal to the desired part number. `find_part`'s search loop will have the form

```
for (p = inventory;
     p != NULL && number > p->number;
     p = p->next)
;
```

The loop will terminate when `p` becomes `NULL` (indicating that the part number wasn't found) or when `number > p->number` is false (indicating that the part number we're looking for is less than or equal to a number already stored in a node). In the latter case, we still don't know whether or not the desired number is actually in the list, so we'll need another test:

```
if (p != NULL && number == p->number)
    return p;
```

The original version of `insert` stores a new part in the next available array element. The new version must determine where the new part belongs in the list and insert it there. We'll also have `insert` check whether the part number is already present in the list. `insert` can accomplish both tasks by using a loop similar to the one in `find_part`:

```
for (cur = inventory, prev = NULL;
     cur != NULL && new_node->number > cur->number;
     prev = cur, cur = cur->next)
;
```

This loop relies on two pointers: `cur`, which points to the current node, and `prev`, which points to the previous node. Once the loop terminates, `insert` will check whether `cur` isn't `NULL` and `new_node->number` equals `cur->number`; if so, the part number is already in the list. Otherwise `insert` will insert a new node between the nodes pointed to by `prev` and `cur`, using a strategy similar to the one we employed for deleting a node. (This strategy works even if the new part number is larger than any in the list; in that case, `cur` will be `NULL` but `prev` will point to the last node in the list.)

Here's the new program. Like the original program, this version requires the `read_line` function described in Section 16.3; I assume that `readline.h` contains a prototype for this function.

```
inventory2.c /* Maintains a parts database (linked list version) */

#include <stdio.h>
#include <stdlib.h>
#include "readline.h"
```

```

#define NAME_LEN 25

struct part {
    int number;
    char name[NAME_LEN+1];
    int on_hand;
    struct part *next;
};

struct part *inventory = NULL; /* points to first part */

struct part *find_part(int number);
void insert(void);
void search(void);
void update(void);
void print(void);

/******************
 * main: Prompts the user to enter an operation code,
 *       then calls a function to perform the requested
 *       action. Repeats until the user enters the
 *       command 'q'. Prints an error message if the user
 *       enters an illegal code.
*****************/
int main(void)
{
    char code;

    for (;;) {
        printf("Enter operation code: ");
        scanf(" %c", &code);
        while (getchar() != '\n') /* skips to end of line */
        ;
        switch (code) {
            case 'i': insert();
                        break;
            case 's': search();
                        break;
            case 'u': update();
                        break;
            case 'p': print();
                        break;
            case 'q': return 0;
            default: printf("Illegal code\n");
        }
        printf("\n");
    }
}

/******************
 * find_part: Looks up a part number in the inventory
 *             list. Returns a pointer to the node
 *             containing the part number; if the part
 *             number is not found, returns NULL.
*****************/

```

```
struct part *find_part(int number)
{
    struct part *p;

    for (p = inventory;
         p != NULL && number > p->number;
         p = p->next)
    ;
    if (p != NULL && number == p->number)
        return p;
    return NULL;
}

/*****************
 * insert: Prompts the user for information about a new *
 *          part and then inserts the part into the      *
 *          inventory list; the list remains sorted by   *
 *          part number. Prints an error message and    *
 *          returns prematurely if the part already exists *
 *          or space could not be allocated for the part. *
 *****************/
void insert(void)
{
    struct part *cur, *prev, *new_node;

    new_node = malloc(sizeof(struct part));
    if (new_node == NULL) {
        printf("Database is full; can't add more parts.\n");
        return;
    }

    printf("Enter part number: ");
    scanf("%d", &new_node->number);

    for (cur = inventory, prev = NULL;
         cur != NULL && new_node->number > cur->number;
         prev = cur, cur = cur->next)
    ;
    if (cur != NULL && new_node->number == cur->number) {
        printf("Part already exists.\n");
        free(new_node);
        return;
    }

    printf("Enter part name: ");
    read_line(new_node->name, NAME_LEN);
    printf("Enter quantity on hand: ");
    scanf("%d", &new_node->on_hand);

    new_node->next = cur;
    if (prev == NULL)
        inventory = new_node;
    else
        prev->next = new_node;
}
```

```
*****
 * search: Prompts the user to enter a part number, then *
 *          looks up the part in the database. If the part *
 *          exists, prints the name and quantity on hand; *
 *          if not, prints an error message. *
 *****
void search(void)
{
    int number;
    struct part *p;

    printf("Enter part number: ");
    scanf("%d", &number);
    p = find_part(number);
    if (p != NULL) {
        printf("Part name: %s\n", p->name);
        printf("Quantity on hand: %d\n", p->on_hand);
    } else
        printf("Part not found.\n");
}

*****
 * update: Prompts the user to enter a part number. *
 *          Prints an error message if the part doesn't *
 *          exist; otherwise, prompts the user to enter *
 *          change in quantity on hand and updates the *
 *          database. *
 *****
void update(void)
{
    int number, change;
    struct part *p;

    printf("Enter part number: ");
    scanf("%d", &number);
    p = find_part(number);
    if (p != NULL) {
        printf("Enter change in quantity on hand: ");
        scanf("%d", &change);
        p->on_hand += change;
    } else
        printf("Part not found.\n");
}

*****
 * print: Prints a listing of all parts in the database, *
 *         showing the part number, part name, and *
 *         quantity on hand. Part numbers will appear in *
 *         ascending order. *
 *****
void print(void)
{
    struct part *p;
```

```

        printf("Part Number      Part Name
               "Quantity on Hand\n");
    for (p = inventory; p != NULL; p = p->next)
        printf("%7d      %-25s%11d\n", p->number, p->name,
               p->on_hand);
}

```

Notice the use of `free` in the `insert` function. `insert` allocates memory for a part before checking to see if the part already exists. If it does, `insert` releases the space to avoid a memory leak.

---

## 17.6 Pointers to Pointers

In Section 13.7, we came across the notion of a *pointer to a pointer*. In that section, we used an array whose elements were of type `char *`; a pointer to one of the array elements itself had type `char **`. The concept of “pointers to pointers” also pops up frequently in the context of linked data structures. In particular, when an argument to a function is a pointer variable, we’ll sometimes want the function to be able to modify the variable by making it point somewhere else. Doing so requires the use of a pointer to a pointer.

Consider the `add_to_list` function of Section 17.5, which inserts a node at the beginning of a linked list. When we call `add_to_list`, we pass it a pointer to the first node in the original list; it then returns a pointer to the first node in the updated list:

```

struct node *add_to_list(struct node *list, int n)
{
    struct node *new_node;

    new_node = malloc(sizeof(struct node));
    if (new_node == NULL) {
        printf("Error: malloc failed in add_to_list\n");
        exit(EXIT_FAILURE);
    }
    new_node->value = n;
    new_node->next = list;
    return new_node;
}

```

Suppose that we modify the function so that it assigns `new_node` to `list` instead of returning `new_node`. In other words, let’s remove the `return` statement from `add_to_list` and replace it by

```
list = new_node;
```

Unfortunately, this idea doesn’t work. Suppose that we call `add_to_list` in the following way:

```
add_to_list(first, 10);
```

At the point of the call, `first` is copied into `list`. (Pointers, like all arguments, are passed by value.) The last line in the function changes the value of `list`, making it point to the new node. This assignment doesn't affect `first`, however.

Getting `add_to_list` to modify `first` is possible, but it requires passing `add_to_list` a *pointer* to `first`. Here's the correct version of the function:

```
void add_to_list(struct node **list, int n)
{
    struct node *new_node;

    new_node = malloc(sizeof(struct node));
    if (new_node == NULL) {
        printf("Error: malloc failed in add_to_list\n");
        exit(EXIT_FAILURE);
    }
    new_node->value = n;
    new_node->next = *list;
    *list = new_node;
}
```

When we call the new version of `add_to_list`, the first argument will be the address of `first`:

```
add_to_list(&first, 10);
```

Since `list` is assigned the address of `first`, we can use `*list` as an alias for `first`. In particular, assigning `new_node` to `*list` will modify `first`.

## 17.7 Pointers to Functions

We've seen that pointers may point to various kinds of data, including variables, array elements, and dynamically allocated blocks of memory. But C doesn't require that pointers point only to *data*; it's also possible to have pointers to *functions*. Pointers to functions aren't as odd as you might think. After all, functions occupy memory locations, so every function has an address, just as each variable has an address.

### Function Pointers as Arguments

We can use function pointers in much the same way we use pointers to data. In particular, passing a function pointer as an argument is fairly common in C. Suppose that we're writing a function named `integrate` that integrates a mathematical function `f` between points `a` and `b`. We'd like to make `integrate` as general as possible by passing it `f` as an argument. To achieve this effect in C, we'll declare `f` to be a pointer to a function. Assuming that we want to integrate functions that have

a double parameter and return a double result, the prototype for `integrate` will look like this:

```
double integrate(double (*f)(double), double a, double b);
```

The parentheses around `*f` indicate that `f` is a pointer to a function, not a function that returns a pointer. It's also legal to declare `f` as though it were a function:

```
double integrate(double f(double), double a, double b);
```

From the compiler's standpoint, this prototype is identical to the previous one.

sin function ► 23.3 When we call `integrate`, we'll supply a function name as the first argument. For example, the following call will integrate the `sin` (sine) function from 0 to  $\pi/2$ :

```
result = integrate(sin, 0.0, PI / 2);
```

Notice that there are no parentheses after `sin`. When a function name isn't followed by parentheses, the C compiler produces a pointer to the function instead of generating code for a function call. In our example, we're not calling `sin`; instead, we're passing `integrate` a pointer to `sin`. If this seems confusing, think of how C handles arrays. If `a` is the name of an array, then `a[i]` represents one element of the array, while `a` by itself serves as a pointer to the array. In a similar way, if `f` is a function, C treats `f(x)` as a *call* of the function but `f` by itself as a *pointer* to the function.

Within the body of `integrate`, we can call the function that `f` points to:

```
y = (*f)(x);
```

`*f` represents the function that `f` points to; `x` is the argument to the call. Thus, during the execution of `integrate(sin, 0.0, PI / 2)`, each call of `*f` is actually a call of `sin`. As an alternative to `(*f)(x)`, C allows us to write `f(x)` to call the function that `f` points to. Although `f(x)` looks more natural, I'll stick with `(*f)(x)` as a reminder that `f` is a pointer to a function, not a function name.

## The `qsort` Function

Although it might seem that pointers to functions aren't relevant to the average programmer, that couldn't be further from the truth. In fact, some of the most useful functions in the C library require a function pointer as an argument. One of these is `qsort`, which belongs to the `<stdlib.h>` header. `qsort` is a general-purpose sorting function that's capable of sorting any array, based on any criteria that we choose.

**Q&A** Since the elements of the array that it sorts may be of any type—even a structure or union type—`qsort` must be told how to determine which of two array elements is “smaller.” We'll provide this information to `qsort` by writing a **comparison function**. When given two pointers `p` and `q` to array elements, the comparison function must return an integer that is *negative* if `*p` is “less than” `*q`,

*zero* if  $*p$  is “equal to”  $*q$ , and *positive* if  $*p$  is “greater than”  $*q$ . The terms “less than,” “equal to,” and “greater than” are in quotes because it’s our responsibility to determine how  $*p$  and  $*q$  are compared.

`qsort` has the following prototype:

```
void qsort(void *base, size_t nmemb, size_t size,
           int (*compar)(const void *, const void *));
```

`base` must point to the first element in the array. (If only a portion of the array is to be sorted, we’ll make `base` point to the first element in this portion.) In the simplest case, `base` is just the name of the array. `nmemb` is the number of elements to be sorted (not necessarily the number of elements in the array). `size` is the size of each array element, measured in bytes. `compar` is a pointer to the comparison function. When `qsort` is called, it sorts the array into ascending order, calling the comparison function whenever it needs to compare array elements.

To sort the `inventory` array of Section 16.3, we’d use the following call of `qsort`:

```
qsort(inventory, num_parts, sizeof(struct part), compare_parts);
```

Notice that the second argument is `num_parts`, not `MAX_PARTS`; we don’t want to sort the entire `inventory` array, just the portion in which parts are currently stored. The last argument, `compare_parts`, is a function that compares two `part` structures.

Writing the `compare_parts` function isn’t as easy as you might expect. `qsort` requires that its parameters have type `void *`, but we can’t access the members of a `part` structure through a `void *` pointer; we need a pointer of type `struct part *` instead. To solve the problem, we’ll have `compare_parts` assign its parameters, `p` and `q`, to variables of type `struct part *`, thereby converting them to the desired type. `compare_parts` can now use these variables to access the members of the structures that `p` and `q` point to. Assuming that we want to sort the `inventory` array into ascending order by part number, here’s how the `compare_parts` function might look:

```
int compare_parts(const void *p, const void *q)
{
    const struct part *p1 = p;
    const struct part *q1 = q;

    if (p1->number < q1->number)
        return -1;
    else if (p1->number == q1->number)
        return 0;
    else
        return 1;
}
```

The declarations of `p1` and `q1` include the word `const` to avoid getting a warning from the compiler. Since `p` and `q` are `const` pointers (indicating that the objects

## Q&A

to which they point should not be modified), they should be assigned only to pointer variables that are also declared to be `const`.

Although this version of `compare_parts` works, most C programmers would write the function more concisely. First, notice that we can replace `p1` and `q1` by cast expressions:

```
int compare_parts(const void *p, const void *q)
{
    if (((struct part *) p)->number <
        ((struct part *) q)->number)
        return -1;
    else if (((struct part *) p)->number ==
              ((struct part *) q)->number))
        return 0;
    else
        return 1;
}
```

The parentheses around `((struct part *) p)` are necessary; without them, the compiler would try to cast `p->number` to type `struct part *`.

We can make `compare_parts` even shorter by removing the `if` statements:

```
int compare_parts(const void *p, const void *q)
{
    return ((struct part *) p)->number -
           ((struct part *) q)->number;
}
```

Subtracting `q`'s part number from `p`'s part number produces a negative result if `p` has a smaller part number, zero if the part numbers are equal, and a positive result if `p` has a larger part number. (Note that subtracting two integers is potentially risky because of the danger of overflow. I'm assuming that part numbers are positive integers, so that shouldn't happen here.)

To sort the `inventory` array by part name instead of part number, we'd use the following version of `compare_parts`:

```
int compare_parts(const void *p, const void *q)
{
    return strcmp(((struct part *) p)->name,
                  ((struct part *) q)->name);
}
```

All `compare_parts` has to do is call `strcmp`, which conveniently returns a negative, zero, or positive result.

## Other Uses of Function Pointers

Although I've emphasized the usefulness of function pointers as arguments to other functions, that's not all they're good for. C treats pointers to functions just like pointers to data; we can store function pointers in variables or use them as ele-

ments of an array or as members of a structure or union. We can even write functions that return function pointers.

Here's an example of a variable that can store a pointer to a function:

```
void (*pf)(int);
```

`pf` can point to any function with an `int` parameter and a return type of `void`. If `f` is such a function, we can make `pf` point to `f` in the following way:

```
pf = f;
```

Notice that there's no ampersand preceding `f`. Once `pf` points to `f`, we can call `f` by writing either

```
(*pf)(i);
```

or

```
pf(i);
```

Arrays whose elements are function pointers have a surprising number of applications. For example, suppose that we're writing a program that displays a menu of commands for the user to choose from. We can write functions that implement these commands, then store pointers to the functions in an array:

```
void (*file_cmd[]) (void) = { new_cmd,
                             open_cmd,
                             close_cmd,
                             close_all_cmd,
                             save_cmd,
                             save_as_cmd,
                             save_all_cmd,
                             print_cmd,
                             exit_cmd
                           };
```

If the user selects command `n`, where `n` falls between 0 and 8, we can subscript the `file_cmd` array and call the corresponding function:

```
(*file_cmd[n])(); /* or file_cmd[n](); */
```

Of course, we could get a similar effect with a `switch` statement. Using an array of function pointers gives us more flexibility, however, since the elements of the array can be changed as the program is running.

## PROGRAM Tabulating the Trigonometric Functions

The following program prints tables showing the values of the `cos`, `sin`, and `tan` functions (all three belong to `<math.h>`). The program is built around a function named `tabulate` that, when passed a function pointer `f`, prints a table showing the values of `f`.

```
tabulate.c /* Tabulates values of trigonometric functions */

#include <math.h>
#include <stdio.h>

void tabulate(double (*f)(double), double first,
              double last, double incr);

int main(void)
{
    double final, increment, initial;

    printf("Enter initial value: ");
    scanf("%lf", &initial);

    printf("Enter final value: ");
    scanf("%lf", &final);

    printf("Enter increment: ");
    scanf("%lf", &increment);

    printf("\n      x      cos(x)\n"
           "      -----  ----- \n");
    tabulate(cos, initial, final, increment);

    printf("\n      x      sin(x)\n"
           "      -----  ----- \n");
    tabulate(sin, initial, final, increment);

    printf("\n      x      tan(x)\n"
           "      -----  ----- \n");
    tabulate(tan, initial, final, increment);

    return 0;
}

void tabulate(double (*f)(double), double first,
              double last, double incr)
{
    double x;
    int i, num_intervals;

    num_intervals = ceil((last - first) / incr);
    for (i = 0; i <= num_intervals; i++) {
        x = first + i * incr;
        printf("%10.5f %10.5f\n", x, (*f)(x));
    }
}
```

tabulate uses the ceil function, which also in <math.h>. When given an argument x of double type, ceil returns the smallest integer that's greater than or equal to x.

Here's what a session with `tabulate.c` might look like:

```
Enter initial value: 0
Enter final value: .5
Enter increment: .1
```

x	cos(x)
0.00000	1.00000
0.10000	0.99500
0.20000	0.98007
0.30000	0.95534
0.40000	0.92106
0.50000	0.87758

x	sin(x)
0.00000	0.00000
0.10000	0.09983
0.20000	0.19867
0.30000	0.29552
0.40000	0.38942
0.50000	0.47943

x	tan(x)
0.00000	0.00000
0.10000	0.10033
0.20000	0.20271
0.30000	0.30934
0.40000	0.42279
0.50000	0.54630

## 17.8 Restricted Pointers (C99)

This section and the next discuss two of C99's pointer-related features. Both are primarily of interest to advanced C programmers; most readers will want to skip these sections.

In C99, the keyword `restrict` may appear in the declaration of a pointer:

```
int * restrict p;
```

A pointer that's been declared using `restrict` is called a **restricted pointer**. The intent is that if `p` points to an object that is later modified, then that object is not accessed in any way other than through `p`. (Alternative ways to access the object include having another pointer to the same object or having `p` point to a named variable.) Having more than one way to access an object is often called **aliasing**.

Let's look at an example of the kind of behavior that restricted pointers are supposed to discourage. Suppose that `p` and `q` have been declared as follows:

```
int * restrict p;
int * restrict q;
```

Now suppose that `p` is made to point to a dynamically allocated block of memory:

```
p = malloc(sizeof(int));
```

(A similar situation would arise if `p` were assigned the address of a variable or an array element.) Normally it would be legal to copy `p` into `q` and then modify the integer through `q`:

```
q = p;
*q = 0; /* causes undefined behavior */
```

Because `p` is a restricted pointer, however, the effect of executing the statement `*q = 0;` is undefined. By making `p` and `q` point to the same object, we caused `*p` and `*q` to be aliases.

extern storage class ▶ 18.2

blocks ▶ 10.3

file scope ▶ 10.2

<string.h> header ▶ 23.6

If a restricted pointer `p` is declared as a local variable without the `extern` storage class, `restrict` applies only to `p` when the block in which `p` is declared is being executed. (Note that the body of a function is a block.) `restrict` can be used with function parameters of pointer type, in which case it applies only when the function is executing. When `restrict` is applied to a pointer variable with file scope, however, the restriction lasts for the entire execution of the program.

The exact rules for using `restrict` are rather complex; see the C99 standard for details. There are even situations in which an alias created from a restricted pointer is legal. For example, a restricted pointer `p` can be legally copied into another restricted pointer variable `q`, provided that `p` is local to a function and `q` is defined inside a block nested within the function's body.

To illustrate the use of `restrict`, let's look at the `memcpy` and `memmove` functions, which belong to the `<string.h>` header. `memcpy` has the following prototype in C99:

```
void *memcpy(void * restrict s1, const void * restrict s2,
            size_t n);
```

`memcpy` is similar to `strcpy`, except that it copies bytes from one object to another (`strcpy` copies characters from one string into another). `s2` points to the data to be copied, `s1` points to the destination of the copy, and `n` is the number of bytes to be copied. The use of `restrict` with both `s1` and `s2` indicates that the source of the copy and the destination shouldn't overlap. (It doesn't guarantee that they don't overlap, however.)

In contrast, `restrict` doesn't appear in the prototype for `memmove`:

```
void *memmove(void *s1, const void *s2, size_t n);
```

`memmove` does the same thing as `memcpy`: it copies bytes from one place to another. The difference is that `memmove` is guaranteed to work even if the source and destination overlap. For example, we could use `memmove` to shift the elements of an array by one position:

```
int a[100];
...
```

```
memmove(&a[0], &a[1], 99 * sizeof(int));
```

Prior to C99, there was no way to document the difference between `memcpy` and `memmove`. The prototypes for the two functions were nearly identical:

```
void *memcpy(void *s1, const void *s2, size_t n);
void *memmove(void *s1, const void *s2, size_t n);
```

The use of `restrict` in the C99 version of `memcpy`'s prototype lets the programmer know that `s1` and `s2` should point to objects that don't overlap, or else the function isn't guaranteed to work.

Although using `restrict` in function prototypes is useful documentation, that's not the primary reason for its existence. `restrict` provides information to the compiler that may enable it to produce more efficient code—a process known as *optimization*. (The `register` storage class serves the same purpose.) Not every compiler attempts to optimize programs, however, and the ones that do normally allow the programmer to disable optimization. As a result, the C99 standard guarantees that `restrict` has no effect on the behavior of a program that conforms to the standard: if all uses of `restrict` are removed from such a program, it should behave the same.

Most programmers won't use `restrict` unless they're fine-tuning a program to achieve the best possible performance. Still, it's worth knowing about `restrict` because it appears in the C99 prototypes for a number of standard library functions.

## 17.9 Flexible Array Members (C99)

Every once in a while, we'll need to define a structure that contains an array of an unknown size. For example, we might want to store strings in a form that's different from the usual one. Normally, a string is an array of characters, with a null character marking the end. However, there are advantages to storing strings in other ways. One alternative is to store the length of the string along with the string's characters (but with no null character). The length and the characters could be stored in a structure such as this one:

```
struct vstring {
    int len;
    char chars[N];
};
```

Here `N` is a macro that represents the maximum length of a string. Using a fixed-length array such as this is undesirable, however, because it forces us to limit the length of the string, plus it wastes memory (since most strings won't need all `N` characters in the array).

C programmers have traditionally solved this problem by declaring the length of `chars` to be 1 (a dummy value) and then dynamically allocating each string:

```

struct vstring {
    int len;
    char chars[1];
};

...
struct vstring *str = malloc(sizeof(struct vstring) + n - 1);
str->len = n;

```

We're "cheating" by allocating more memory than the structure is declared to have (in this case, an extra  $n - 1$  characters), and then using the memory to store additional elements of the `chars` array. This technique has become so common over the years that it has a name: the "struct hack."

The struct hack isn't limited to character arrays; it has a variety of uses. Over time, it has become popular enough to be supported by many compilers. Some (including GCC) even allow the `chars` array to have zero length, which makes this trick a little more explicit. Unfortunately, the C89 standard doesn't guarantee that the struct hack will work, nor does it allow zero-length arrays.

In recognition of the struct hack's usefulness, C99 has a feature known as the **flexible array member** that serves the same purpose. When the last member of a structure is an array, its length may be omitted:

```

struct vstring {
    int len;
    char chars[]; /* flexible array member - C99 only */
};

```

The length of the `chars` array isn't determined until memory is allocated for a `vstring` structure, normally using a call of `malloc`:

```

struct vstring *str = malloc(sizeof(struct vstring) + n);
str->len = n;

```

In this example, `str` points to a `vstring` structure in which the `chars` array occupies  $n$  characters. The `sizeof` operator ignores the `chars` member when computing the size of the structure. (A flexible array member is unusual in that it takes up no space within a structure.)

A few special rules apply to a structure that contains a flexible array member. The flexible array member must appear last in the structure, and the structure must have at least one other member. Copying a structure that contains a flexible array member will copy the other members but not the flexible array itself.

A structure that contains a flexible array member is an **incomplete type**. An incomplete type is missing part of the information needed to determine how much memory it requires. Incomplete types, which are discussed further in one of the Q&A questions at the end of this chapter and in Section 19.3, are subject to various restrictions. In particular, an incomplete type (and hence a structure that contains a flexible array member) can't be a member of another structure or an element of an array. However, an array may contain pointers to structures that have a flexible array member; Programming Project 7 at the end of this chapter is built around such an array.

## Q & A

**Q: What does the `NULL` macro represent? [p. 415]**

A: `NULL` actually stands for 0. When we use 0 in a context where a pointer would be required, C compilers treat it as a null pointer instead of the integer 0. The `NULL` macro is provided merely to help avoid confusion. The assignment

```
p = 0;
```

could be assigning the value 0 to a numeric variable or assigning a null pointer to a pointer variable; we can't easily tell which. In contrast, the assignment

```
p = NULL;
```

makes it clear that `p` is a pointer.

**\*Q: In the header files that come with my compiler, `NULL` is defined as follows:**

```
#define NULL (void *) 0
```

**What's the advantage of casting 0 to `void *`?**

A: This trick, which is allowed by the C standard, enables compilers to spot incorrect uses of the null pointer. For example, suppose that we try to assign `NULL` to an integer variable:

```
i = NULL;
```

If `NULL` is defined as 0, this assignment is perfectly legal. But if `NULL` is defined as `(void *) 0`, the compiler can warn us that we're assigning a pointer to an integer variable.

Defining `NULL` as `(void *) 0` has a second, more important, advantage. Suppose that we call a function with a variable-length argument list and pass `NULL` as one of the arguments. If `NULL` is defined as 0, the compiler will incorrectly pass a zero integer value. (In an ordinary function call, `NULL` works fine because the compiler knows from the function's prototype that it expects a pointer. When a function has a variable-length argument list, however, the compiler lacks this knowledge.) If `NULL` is defined as `(void *) 0`, the compiler will pass a null pointer.

To make matters even more confusing, some header files define `NULL` to be `0L` (the long version of 0). This definition, like the definition of `NULL` as 0, is a holdover from C's earlier years, when pointers and integers were compatible. For most purposes, though, it really doesn't matter how `NULL` is defined; just think of it as a name for the null pointer.

**Q: Since 0 is used to represent the null pointer, I guess a null pointer is just an address with all zero bits, right?**

- A: Not necessarily. Each C compiler is allowed to represent null pointers in a different way, and not all compilers use a zero address. For example, some compilers use a nonexistent memory address for the null pointer; that way, attempting to access memory through a null pointer can be detected by the hardware.

How the null pointer is stored inside the computer shouldn't concern us; that's a detail for compiler experts to worry about. The important thing is that, when used in a pointer context, 0 is converted to the proper internal form by the compiler.

**Q: Is it acceptable to use `NULL` as a null character?**

- A: Definitely not. `NULL` is a macro that represents the null *pointer*, not the null *character*. Using `NULL` as a null character will work with some compilers, but not with all (since some define `NULL` as `(void *) 0`). In any event, using `NULL` as anything other than a pointer can lead to a great deal of confusion. If you want a name for the null character, define the following macro:

```
#define NUL '\0'
```

**\*Q: When my program terminates, I get the message “*Null pointer assignment*.” What does this mean?**

- A: This message, which is produced by programs compiled with some older DOS-based C compilers, indicates that the program has stored data in memory using a bad pointer (but not necessarily a null pointer). Unfortunately, the message isn't displayed until the program terminates, so there's no clue as to which statement caused the error. The “*Null pointer assignment*” message can be caused by a missing & in `scanf`:

```
scanf("%d", i); /* should have been scanf("%d", &i); */
```

Another possibility is an assignment involving a pointer that's uninitialized or null:

```
*p = i; /* p is uninitialized or null */
```

**\*Q: How does a program know that a “*null pointer assignment*” has occurred?**

- A: The message depends on the fact that, in the small and medium memory models, data is stored in a single segment, with addresses beginning at 0. The compiler leaves a “hole” at the beginning of the data segment—a small block of memory that's initialized to 0 but otherwise isn't used by the program. When the program terminates, it checks to see if any data in the “hole” area is nonzero. If so, it must have been altered through a bad pointer.

**Q: Is there any advantage to casting the return value of `malloc` or the other memory allocation functions? [p. 416]**

- A: Not usually. Casting the `void *` pointer that these functions return is unnecessary, since pointers of type `void *` are automatically converted to any pointer type upon assignment. The habit of casting the return value is a holdover from older versions of C, in which the memory allocation functions returned a `char *` value, making the cast necessary. Programs that are designed to be compiled as C++ code

may benefit from the cast, but that's about the only reason to do it.

In C89, there's actually a small advantage to *not* performing the cast. Suppose that we've forgotten to include the `<stdlib.h>` header in our program. When we call `malloc`, the compiler will assume that its return type is `int` (the default return value for any C function). If we don't cast the return value of `malloc`, a C89 compiler will produce an error (or at least a warning), since we're trying to assign an integer value to a pointer variable. On the other hand, if we cast the return value to a pointer, the program may compile, but likely won't run properly. With C99, this advantage disappears. Forgetting to include the `<stdlib.h>` header will cause an error when `malloc` is called, because C99 requires that a function be declared before it's called.

**Q:** **The `calloc` function initializes a memory block by setting its bits to zero. Does this mean that all data items in the block become zero? [p. 421]**

**A:** Usually, but not always. Setting an integer to zero bits always makes the integer zero. Setting a floating-point number to zero bits usually makes the number zero, but this isn't guaranteed—it depends on how floating-point numbers are stored. The story is the same for pointers; a pointer whose bits are zero isn't necessarily a null pointer.

**\*Q:** **I see how the structure tag mechanism allows a structure to contain a pointer to itself. But what if two structures each have a member that points to the other? [p. 425]**

**A:** Here's how we'd handle that situation:

```
struct s1; /* incomplete declaration of s1 */

struct s2 {
    ...
    struct s1 *p;
    ...
};

struct s1 {
    ...
    struct s2 *q;
    ...
};
```

incomplete types ▶ 19.3

The first declaration of `s1` creates an incomplete structure type, since we haven't specified the members of `s1`. The second declaration of `s1` "completes" the type by describing the members of the structure. Incomplete declarations of a structure type are permitted in C, although their uses are limited. Creating a pointer to such a type (as we did when declaring `p`) is one of these uses.

**Q:** **Calling `malloc` with the wrong argument—causing it to allocate too much memory or too little memory—seems to be a common error. Is there a safer way to use `malloc`? [p. 426]**

- A: Yes, there is. Some programmers use the following idiom when calling `malloc` to allocate memory for a single object:

```
p = malloc(sizeof(*p));
```

Since `sizeof(*p)` is the size of the object to which `p` will point, this statement guarantees that the correct amount of memory will be allocated. At first glance, this idiom looks fishy: it's likely that `p` is uninitialized, making the value of `*p` undefined. However, `sizeof` doesn't evaluate `*p`, it merely computes its size, so the idiom works even if `p` is uninitialized or contains a null pointer.

To allocate memory for an array with `n` elements, we can use a slightly modified version of the idiom:

```
p = malloc(n * sizeof(*p));
```

- Q: Why isn't the `qsort` function simply named `sort`? [p. 440]**

- A: The name `qsort` comes from the Quicksort algorithm published by C. A. R. Hoare in 1962 (and discussed in Section 9.6). Ironically, the C standard doesn't require that `qsort` use the Quicksort algorithm, although many versions of `qsort` do.

- Q: Isn't it necessary to cast `qsort`'s first argument to type `void *`, as in the following example? [p. 441]**

```
qsort((void *) inventory, num_parts, sizeof(struct part),
      compare_parts);
```

- A: No. A pointer of any type can be converted to `void *` automatically.

- \*Q: I want to use `qsort` to sort an array of integers, but I'm having trouble writing a comparison function. What's the secret?**

- A: Here's a version that works:

```
int compare_ints(const void *p, const void *q)
{
    return *(int *)p - *(int *)q;
}
```

Bizarre, eh? The expression `(int *)p` casts `p` to type `int *`, so `*(int *)p` would be the integer that `p` points to. A word of warning, though: Subtracting two integers may cause overflow. If the integers being sorted are completely arbitrary, it's safer to use `if` statements to compare `*(int *)p` with `*(int *)q`.

- \*Q: I needed to sort an array of strings, so I figured I'd just use `strcmp` as the comparison function. When I passed it to `qsort`, however, the compiler gave me a warning. I tried to fix the problem by embedding `strcmp` in a comparison function:**

```
int compare_strings(const void *p, const void *q)
{
    return strcmp(p, q);
}
```

**Now my program compiles, but `qsort` doesn't seem to sort the array. What am I doing wrong?**

- A: First, you can't pass `strcmp` itself to `qsort`, since `qsort` requires a comparison function with two `const void *` parameters. Your `compare_strings` function doesn't work because it incorrectly assumes that `p` and `q` are strings (`char *` pointers). In fact, `p` and `q` point to array elements containing `char *` pointers. To fix `compare_strings`, we'll cast `p` and `q` to type `char **`, then use the `*` operator to remove one level of indirection:

```
int compare_strings(const void *p, const void *q)
{
    return strcmp(*((char **)p), *((char **)q));
}
```

## Exercises

### Section 17.1

- Having to check the return value of `malloc` (or any other memory allocation function) each time we call it can be an annoyance. Write a function named `my_malloc` that serves as a “wrapper” for `malloc`. When we call `my_malloc` and ask it to allocate `n` bytes, it in turn calls `malloc`, tests to make sure that `malloc` doesn't return a null pointer, and then returns the pointer from `malloc`. Have `my_malloc` print an error message and terminate the program if `malloc` returns a null pointer.

### Section 17.2

- W 2. Write a function named `duplicate` that uses dynamic storage allocation to create a copy of a string. For example, the call

```
p = duplicate(str);
```

would allocate space for a string of the same length as `str`, copy the contents of `str` into the new string, and return a pointer to it. Have `duplicate` return a null pointer if the memory allocation fails.

### Section 17.3

- Write the following function:

```
int *create_array(int n, int initial_value);
```

The function should return a pointer to a dynamically allocated `int` array with `n` members, each of which is initialized to `initial_value`. The return value should be `NULL` if the array can't be allocated.

### Section 17.5

- Suppose that the following declarations are in effect:

```
struct point { int x, y; };
struct rectangle { struct point upper_left, lower_right; };
struct rectangle *p;
```

Assume that we want `p` to point to a `rectangle` structure whose upper left corner is at (10, 25) and whose lower right corner is at (20, 15). Write a series of statements that allocate such a structure and initialize it as indicated.

- W 5. Suppose that `f` and `p` are declared as follows:

```
struct {
    union {
        char a, b;
        int c;
    } d;
    int e[5];
} f, *p = &f;
```

Which of the following statements are legal?

- (a) `p->b = ' ';`
- (b) `p->e[3] = 10;`
- (c) `(*p).d.a = '*' ;`
- (d) `p->d->c = 20;`

6. Modify the `delete_from_list` function so that it uses only one pointer variable instead of two (`cur` and `prev`).

- W 7. The following loop is supposed to delete all nodes from a linked list and release the memory that they occupy. Unfortunately, the loop is incorrect. Explain what's wrong with it and show how to fix the bug.

```
for (p = first; p != NULL; p = p->next)
    free(p);
```

- W 8. Section 15.2 describes a file, `stack.c`, that provides functions for storing integers in a stack. In that section, the stack was implemented as an array. Modify `stack.c` so that a stack is now stored as a linked list. Replace the `contents` and `top` variables by a single variable that points to the first node in the list (the “top” of the stack). Write the functions in `stack.c` so that they use this pointer. Remove the `is_full` function, instead having `push` return either `true` (if memory was available to create a node) or `false` (if not).

9. True or false: If `x` is a structure and `a` is a member of that structure, then `(&x) ->a` is the same as `x.a`. Justify your answer.

10. Modify the `print_part` function of Section 16.2 so that its parameter is a *pointer* to a `part` structure. Use the `->` operator in your answer.

11. Write the following function:

```
int count_occurrences(struct node *list, int n);
```

The `list` parameter points to a linked list; the function should return the number of times that `n` appears in this list. Assume that the `node` structure is the one defined in Section 17.5.

12. Write the following function:

```
struct node *find_last(struct node *list, int n);
```

The `list` parameter points to a linked list. The function should return a pointer to the *last* node that contains `n`; it should return `NULL` if `n` doesn't appear in the list. Assume that the `node` structure is the one defined in Section 17.5.

13. The following function is supposed to insert a new node into its proper place in an ordered list, returning a pointer to the first node in the modified list. Unfortunately, the function

doesn't work correctly in all cases. Explain what's wrong with it and show how to fix it. Assume that the node structure is the one defined in Section 17.5.

```
struct node *insert_into_ordered_list(struct node *list,
                                      struct node *new_node)
{
    struct node *cur = list, *prev = NULL;
    while (cur->value <= new_node->value) {
        prev = cur;
        cur = cur->next;
    }
    prev->next = new_node;
    new_node->next = cur;
    return list;
}
```

### Section 17.6

14. Modify the `delete_from_list` function (Section 17.5) so that its first parameter has type `struct node **` (a pointer to a pointer to the first node in a list) and its return type is `void`. `delete_from_list` must modify its first argument to point to the list after the desired node has been deleted.

### Section 17.7

- W 15. Show the output of the following program and explain what it does.

```
#include <stdio.h>

int f1(int (*f)(int));
int f2(int i);

int main(void)
{
    printf("Answer: %d\n", f1(f2));
    return 0;
}

int f1(int (*f)(int))
{
    int n = 0;

    while ((*f)(n)) n++;
    return n;
}

int f2(int i)
{
    return i * i + i - 12;
}
```

16. Write the following function. The call `sum(g, i, j)` should return `g(i) + ... + g(j)`.

```
int sum(int (*f)(int), int start, int end);
```

- W 17. Let `a` be an array of 100 integers. Write a call of `qsort` that sorts only the *last* 50 elements in `a`. (You don't need to write the comparison function).
18. Modify the `compare_parts` function so that parts are sorted with their numbers in *descending* order.
19. Write a function that, when given a string as its argument, searches the following array of structures for a matching command name, then calls the function associated with that name.

```

struct {
    char *cmd_name;
    void (*cmd_pointer)(void);
} file_cmd[] =
{ {"new",      new_cmd},
  {"open",     open_cmd},
  {"close",   close_cmd},
  {"close all", close_all_cmd},
  {"save",     save_cmd},
  {"save as",  save_as_cmd},
  {"save all", save_all_cmd},
  {"print",   print_cmd},
  {"exit",    exit_cmd}
};

```

## Programming Projects

- W 1. Modify the `inventory.c` program of Section 16.3 so that the `inventory` array is allocated dynamically and later reallocated when it fills up. Use `malloc` initially to allocate enough space for an array of 10 part structures. When the array has no more room for new parts, use `realloc` to double its size. Repeat the doubling step each time the array becomes full.
- W 2. Modify the `inventory.c` program of Section 16.3 so that the `p` (print) command calls `qsort` to sort the `inventory` array before it prints the parts.
- 3. Modify the `inventory2.c` program of Section 17.5 by adding an `e` (erase) command that allows the user to remove a part from the database.
- 4. Modify the `justify` program of Section 15.3 by rewriting the `line.c` file so that it stores the current line in a linked list. Each node in the list will store a single word. The `line` array will be replaced by a variable that points to the node containing the first word. This variable will store a null pointer whenever the line is empty.
- 5. Write a program that sorts a series of words entered by the user:

```

Enter word: foo
Enter word: bar
Enter word: baz
Enter word: quux
Enter word:

```

In sorted order: bar baz foo quux

Assume that each word is no more than 20 characters long. Stop reading when the user enters an empty word (i.e., presses Enter without entering a word). Store each word in a dynamically allocated string, using an array of pointers to keep track of the strings, as in the `remind2.c` program (Section 17.2). After all words have been read, sort the array (using any sorting technique) and then use a loop to print the words in sorted order. *Hint:* Use the `read_line` function to read each word, as in `remind2.c`.

- 6. Modify Programming Project 5 so that it uses `qsort` to sort the array of pointers.
- 7. (C99) Modify the `remind2.c` program of Section 17.2 so that each element of the `reminders` array is a pointer to a `vstring` structure (see Section 17.9) rather than a pointer to an ordinary string.

# 18 Declarations

*Making something variable is easy.  
Controlling duration of constancy is the trick.*

Declarations play a central role in C programming. By declaring variables and functions, we furnish vital information that the compiler will need in order to check a program for potential errors and translate it into object code.

Previous chapters have provided examples of declarations without going into full details; this chapter fills in the gaps. It explores the sophisticated options that can be used in declarations and reveals that variable declarations and function declarations have quite a bit in common. It also provides a firm grounding in the important concepts of storage duration, scope, and linkage.

Section 18.1 examines the syntax of declarations in their most general form, a topic that we've avoided up to this point. The next four sections focus on the items that appear in declarations: storage classes (Section 18.2), type qualifiers (Section 18.3), declarators (Section 18.4), and initializers (Section 18.5). Section 18.6 discusses the `inline` keyword, which can appear in C99 function declarations.

## 18.1 Declaration Syntax

Declarations furnish information to the compiler about the meaning of identifiers. When we write

```
int i;
```

we're informing the compiler that, in the current scope, the name `i` represents a variable of type `int`. The declaration

```
float f(float);
```

tells the compiler that `f` is a function that returns a `float` value and has one argument, also of type `float`.

In general, a declaration has the following appearance:

**declaration**

*declaration-specifiers declarators ;*

**Declaration specifiers** describe the properties of the variables or functions being declared. **Declarators** give their names and may provide additional information about their properties.

Declaration specifiers fall into three categories:

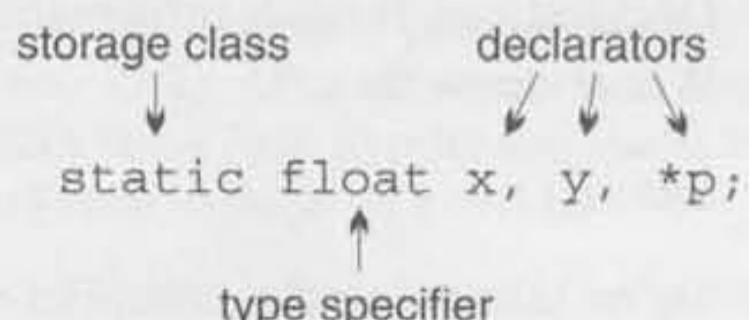
- **Storage classes.** There are four storage classes: `auto`, `static`, `extern`, and `register`. At most one storage class may appear in a declaration; if present, it should come first.
- **Type qualifiers.** In C89, there are only two type qualifiers: `const` and `volatile`. C99 has a third type qualifier, `restrict`. A declaration may contain zero or more type qualifiers.
- **Type specifiers.** The keywords `void`, `char`, `short`, `int`, `long`, `float`, `double`, `signed`, and `unsigned` are all type specifiers. These words may be combined as described in Chapter 7; the order in which they appear doesn't matter (`int unsigned long` is the same as `long unsigned int`). Type specifiers also include specifications of structures, unions, and enumerations (for example, `struct point { int x, y; }`, `struct { int x, y; }`, or `struct point`). Type names created using `typedef` are type specifiers as well.

C99

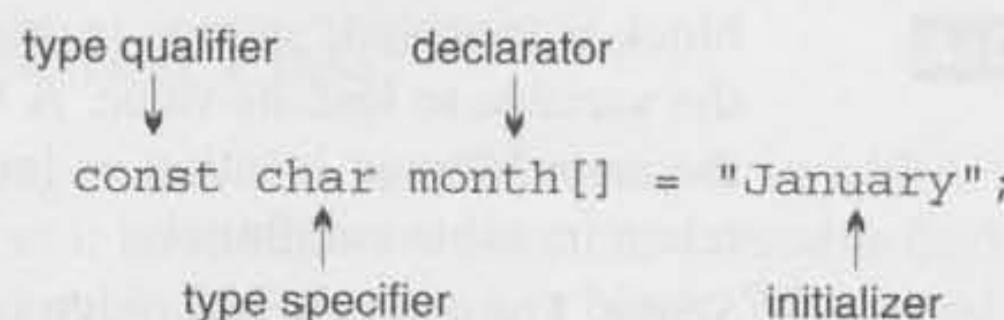
(C99 has a fourth kind of declaration specifier, the **function specifier**, which is used only in function declarations. This category has just one member, the keyword `inline`.) Type qualifiers and type specifiers should follow the storage class, but there are no other restrictions on their order. As a matter of style, I'll put type qualifiers before type specifiers.

Declarators include identifiers (names of simple variables), identifiers followed by `[]` (array names), identifiers preceded by `*` (pointer names), and identifiers followed by `()` (function names). Declarators are separated by commas. A declarator that represents a variable may be followed by an initializer.

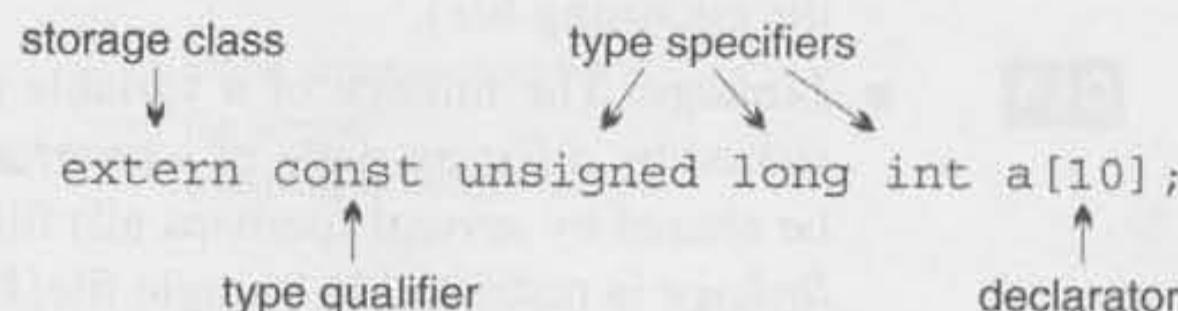
Let's look at a few examples that illustrate these rules. Here's a declaration with a storage class and three declarators:



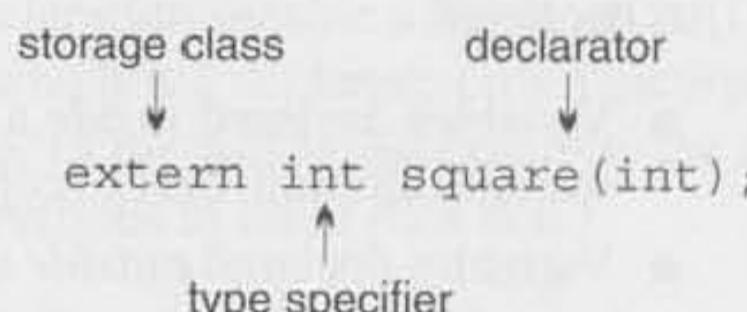
The following declaration has a type qualifier but no storage class. It also has an initializer:



The following declaration has both a storage class and a type qualifier. It also has three type specifiers; their order isn't important:



Function declarations, like variable declarations, may have a storage class, type qualifiers, and type specifiers. The following declaration has a storage class and a type specifier:



The next four sections cover storage classes, type qualifiers, declarators, and initializers in detail.

## 18.2 Storage Classes

Storage classes can be specified for variables and—to a lesser extent—functions and parameters. We'll concentrate on variables for now.

Recall from Section 10.3 that the term *block* refers to the body of a function (the part enclosed in braces) or a compound statement, possibly containing declarations. In C99, selection statements (`if` and `switch`) and iteration statements (`while`, `do`, and `for`)—along with the “inner” statements that they control—are considered to be blocks as well, although this is primarily a technicality.

**C99**

**Q&A**

### Properties of Variables

Every variable in a C program has three properties:

- **Storage duration.** The storage duration of a variable determines when memory is set aside for the variable and when that memory is released. Storage for a variable with *automatic storage duration* is allocated when the surrounding

**Q&A**

block is executed; storage is deallocated when the block terminates, causing the variable to lose its value. A variable with *static storage duration* stays at the same storage location as long as the program is running, allowing it to retain its value indefinitely.

- **Scope.** The scope of a variable is the portion of the program text in which the variable can be referenced. A variable can have either *block scope* (the variable is visible from its point of declaration to the end of the enclosing block) or *file scope* (the variable is visible from its point of declaration to the end of the enclosing file).

**Q&A**

- **Linkage.** The linkage of a variable determines the extent to which it can be shared by different parts of a program. A variable with *external linkage* may be shared by several (perhaps all) files in a program. A variable with *internal linkage* is restricted to a single file, but may be shared by the functions in that file. (If a variable with the same name appears in another file, it's treated as a different variable.) A variable with *no linkage* belongs to a single function and can't be shared at all.

The default storage duration, scope, and linkage of a variable depend on where it's declared:

- Variables declared *inside* a block (including a function body) have *automatic storage duration*, *block scope*, and *no linkage*.
- Variables declared *outside* any block, at the outermost level of a program, have *static storage duration*, *file scope*, and *external linkage*.

The following example shows the default properties of the variables *i* and *j*:

```
int i;           // static storage duration, file scope, external linkage
void f(void)
{
    int j;        // automatic storage duration, block scope, no linkage
}
```

For many variables, the default storage duration, scope, and linkage are satisfactory. When they aren't, we can alter these properties by specifying an explicit storage class: *auto*, *static*, *extern*, or *register*.

## The **auto** Storage Class

The *auto* storage class is legal only for variables that belong to a block. An *auto* variable has automatic storage duration (not surprisingly), block scope, and no linkage. The *auto* storage class is almost never specified explicitly, since it's the default for variables declared inside a block.

## The `static` Storage Class

The `static` storage class can be used with all variables, regardless of where they're declared, but it has a different effect on a variable declared outside a block than it does on a variable declared inside a block. When used *outside* a block, the word `static` specifies that a variable has internal linkage. When used *inside* a block, `static` changes the variable's storage duration from automatic to static. The following figure shows the effect of declaring `i` and `j` to be `static`:

```
static int i;           static storage duration
                      file scope
                      internal linkage

void f(void)
{
    static int j;       static storage duration
                      block scope
                      no linkage
}
```

When used in a declaration outside a block, `static` essentially hides a variable within the file in which it's declared; only functions that appear in the same file can see the variable. In the following example, the functions `f1` and `f2` both have access to `i`, but functions in other files don't:

```
static int i;

void f1(void)
{
    /* has access to i */
}

void f2(void)
{
    /* has access to i */
}
```

This use of `static` can help implement a technique known as information hiding.

A `static` variable declared within a block resides at the same storage location throughout program execution. Unlike automatic variables, which lose their values each time the program leaves the enclosing block, a `static` variable will retain its value indefinitely. `static` variables have some interesting properties:

- A `static` variable in a block is initialized only once, prior to program execution. An `auto` variable is initialized every time it comes into existence (provided, of course, that it has an initializer).
- Each time a function is called recursively, it gets a new set of `auto` variables. If it has a `static` variable, on the other hand, that variable is shared by all calls of the function.

- Although a function shouldn't return a pointer to an `auto` variable, there's nothing wrong with it returning a pointer to a `static` variable.

Declaring one of its variables to be `static` allows a function to retain information between calls in a “hidden” area that the rest of the program can't access. More often, however, we'll use `static` to make programs more efficient. Consider the following function:

```
char digit_to_hex_char(int digit)
{
    const char hex_chars[16] = "0123456789ABCDEF";
    return hex_chars[digit];
}
```

Each time the `digit_to_hex_char` function is called, the characters `0123456789ABCDEF` will be copied into the `hex_chars` array to initialize it. Now, let's make the array `static`:

```
char digit_to_hex_char(int digit)
{
    static const char hex_chars[16] = "0123456789ABCDEF";
    return hex_chars[digit];
}
```

Since `static` variables are initialized only once, we've improved the speed of `digit_to_hex_char`.

## The `extern` Storage Class

The `extern` storage class enables several source files to share the same variable. Section 15.2 covered the essentials of using `extern`, so I won't devote much space to it here. Recall that the declaration

```
extern int i;
```

informs the compiler that `i` is an `int` variable, but doesn't cause it to allocate memory for `i`. In C terminology, this declaration is not a *definition* of `i`; it merely informs the compiler that we need access to a variable that's defined elsewhere (perhaps later in the same file, or—more often—in another file). A variable can have many *declarations* in a program but should have only one *definition*.

There's one exception to the rule that an `extern` declaration of a variable isn't a definition. An `extern` declaration that initializes a variable serves as a definition of the variable. For example, the declaration

```
extern int i = 0;
```

is effectively the same as

```
int i = 0;
```

This rule prevents multiple `extern` declarations from initializing a variable in different ways.

### Q&A

A variable in an `extern` declaration always has static storage duration. The scope of the variable depends on the declaration's placement. If the declaration is inside a block, the variable has block scope; otherwise, it has file scope:

```
static storage duration
extern int i;           file scope
                        ? linkage

void f(void)
{
    static storage duration
    extern int j;           block scope
                            ? linkage
}
```

Determining the linkage of an `extern` variable is a bit harder. If the variable was declared `static` earlier in the file (outside of any function definition), then it has internal linkage. Otherwise (the normal case), the variable has external linkage.

## The `register` Storage Class

Using the `register` storage class in the declaration of a variable asks the compiler to store the variable in a register instead of keeping it in main memory like other variables. (A *register* is a storage area located in a computer's CPU. Data stored in a register can be accessed and updated faster than data stored in ordinary memory.) Specifying the storage class of a variable to be `register` is a request, not a command. The compiler is free to store a `register` variable in memory if it chooses.

The `register` storage class is legal only for variables declared in a block. A `register` variable has the same storage duration, scope, and linkage as an `auto` variable. However, a `register` variable lacks one property that an `auto` variable has: since registers don't have addresses, it's illegal to use the `&` operator to take the address of a `register` variable. This restriction applies even if the compiler has elected to store the variable in memory.

`register` is best used for variables that are accessed and/or updated frequently. For example, the loop control variable in a `for` statement is a good candidate for `register` treatment:

```
int sum_array(int a[], int n)
{
    register int i;
    int sum = 0;

    for (i = 0; i < n; i++)
        sum += a[i];
    return sum;
}
```

`register` isn't nearly as popular among C programmers as it once was. Today's compilers are much more sophisticated than early C compilers; many can determine automatically which variables would benefit the most from being kept in registers. Still, using `register` provides useful information that can help the compiler optimize the performance of a program. In particular, the compiler knows that a `register` variable can't have its address taken, and therefore can't be modified through a pointer. In this respect, the `register` keyword is related to C99's `restrict` keyword.

## The Storage Class of a Function

Function declarations (and definitions), like variable declarations, may include a storage class, but the only options are `extern` and `static`. The word `extern` at the beginning of a function declaration specifies that the function has external linkage, allowing it to be called from other files. `static` indicates internal linkage, limiting use of the function's name to the file in which it's defined. If no storage class is specified, the function is assumed to have external linkage.

Consider the following function declarations:

```
extern int f(int i);
static int g(int i);
int h(int i);
```

`f` has external linkage, `g` has internal linkage, and `h` (by default) has external linkage. Because it has internal linkage, `g` can't be called directly from outside the file in which it's defined. (Declaring `g` to be `static` doesn't completely prevent it from being called in another file; an indirect call via a function pointer is still possible.)

Declaring functions to be `extern` is like declaring variables to be `auto`—it serves no purpose. For that reason, I don't use `extern` in function declarations. Be aware, however, that some programmers use `extern` extensively, which certainly does no harm.

Declaring functions to be `static`, on the other hand, is quite useful. In fact, I recommend using `static` when declaring any function that isn't intended to be called from other files. The benefits of doing so include:

- **Easier maintenance.** Declaring a function `f` to be `static` guarantees that `f` isn't visible outside the file in which its definition appears. As a result, someone modifying the program later knows that changes to `f` won't affect functions in other files. (One exception: a function in another file that's passed a pointer to `f` might be affected by changes to `f`. Fortunately, that situation is easy to spot by examining the file in which `f` is defined, since the function that passes `f` must also be defined there.)
- **Reduced “name space pollution.”** Since functions declared `static` have internal linkage, their names can be reused in other files. Although we proba-

bly wouldn't deliberately reuse a function name for some other purpose, it can be hard to avoid in large programs. An excessive number of names with external linkage can result in what C programmers call "name space pollution": names in different files accidentally conflicting with each other. Using `static` helps prevent this problem.

Function parameters have the same properties as `auto` variables: automatic storage duration, block scope, and no linkage. The only storage class that can be specified for parameters is `register`.

## Summary

Now that we've covered the various storage classes, let's summarize what we know. The following program fragment shows all possible ways to include—or omit—storage classes in declarations of variables and parameters.

```
int a;
extern int b;
static int c;

void f(int d, register int e)
{
    auto int g;
    int h;
    static int i;
    extern int j;
    register int k;
}
```

Table 18.1 shows the properties of each variable and parameter in this example.

**Table 18.1**

Properties of Variables  
and Parameters

Name	Storage Duration	Scope	Linkage
a	static	file	external
b	static	file	†
c	static	file	internal
d	automatic	block	none
e	automatic	block	none
g	automatic	block	none
h	automatic	block	none
i	static	block	none
j	static	block	†
k	automatic	block	none

<sup>†</sup>The definitions of `b` and `j` aren't shown, so it's not possible to determine the linkage of these variables. In most cases, the variables will be defined in another file and will have external linkage.

Of the four storage classes, the most important are `static` and `extern`. `auto` has no effect, and modern compilers have made `register` less important.

## 18.3 Type Qualifiers

**C99**

restricted pointers ▶ 17.8

There are two type qualifiers: `const` and `volatile`. (C99 has a third type qualifier, `restrict`, which is used only with pointers.) Since the use of `volatile` is limited to low-level programming, I'll postpone discussing it until Section 20.3. `const` is used to declare objects that resemble variables but are “read-only”: a program may access the value of a `const` object, but can't change it. For example, the declaration

```
const int n = 10;
```

creates a `const` object named `n` whose value is 10. The declaration

```
const int tax_brackets[] = {750, 2250, 3750, 5250, 7000};
```

creates a `const` array named `tax_brackets`.

Declaring an object to be `const` has several advantages:

- It's a form of documentation: it alerts anyone reading the program to the read-only nature of the object.
- The compiler can check that the program doesn't inadvertently attempt to change the value of the object.
- When programs are written for certain types of applications (embedded systems, in particular), the compiler can use the word `const` to identify data to be stored in ROM (read-only memory).

At first glance, it might appear that `const` serves the same role as the `#define` directive, which we've used in previous chapters to create names for constants. There are significant differences between `#define` and `const`, however:

- We can use `#define` to create a name for a numerical, character, or string constant. `const` can be used to create read-only objects of *any* type, including arrays, pointers, structures, and unions.
- `const` objects are subject to the same scope rules as variables; constants created using `#define` aren't. In particular, we can't use `#define` to create a constant with block scope.
- The value of a `const` object, unlike the value of a macro, can be viewed in a debugger.
- Unlike macros, `const` objects can't be used in constant expressions. For example, we can't write

```
const int n = 10;
int a[n];           /*** WRONG ***/
```

**C99**

since array bounds must be constant expressions. (In C99, this example would

be legal if `a` has automatic storage duration—it would be treated as a variable-length array—but not if it has static storage duration.)

- It's legal to apply the address operator (`&`) to a `const` object, since it has an address. A macro doesn't have an address.

There are no absolute rules that dictate when to use `#define` and when to use `const`. I recommend using `#define` for constants that represent numbers or characters. That way, you'll be able to use the constants as array dimensions, in `switch` statements, and in other places where constant expressions are required.

## 18.4 Declarators

A declarator consists of an identifier (the name of the variable or function being declared), possibly preceded by the `*` symbol or followed by `[]` or `()`. By combining `*`, `[]`, and `()`, we can create declarators of mind-numbing complexity.

Before we look at the more complicated declarators, let's review the declarators that we've seen in previous chapters. In the simplest case, a declarator is just an identifier, like `i` in the following example:

```
int i;
```

Declarators may also contain the symbols `*`, `[]`, and `()`:

- A declarator that begins with `*` represents a pointer:

```
int *p;
```

- A declarator that ends with `[]` represents an array:

```
int a[10];
```

The brackets may be left empty if the array is a parameter, if it has an initializer, or if its storage class is `extern`:

```
extern int a[];
```

Since `a` is defined elsewhere in the program, the compiler doesn't need to know its length here. (In the case of a multidimensional array, only the first set of brackets can be empty.) C99 provides two additional options for what goes between the brackets in the declaration of an array parameter. One option is the keyword `static`, followed by an expression that specifies the array's minimum length. The other is the `*` symbol, which can be used in a function prototype to indicate a variable-length array argument. Section 9.3 discusses both C99 features.

- A declarator that ends with `()` represents a function:

```
int abs(int i);
void swap(int *a, int *b);
int find_largest(int a[], int n);
```

C99

C allows parameter names to be omitted in a function declaration:

```
int abs(int);
void swap(int *, int *);
int find_largest(int [], int);
```

The parentheses can even be left empty:

```
int abs();
void swap();
int find_largest();
```

The declarations in the last group specify the return types of the `abs`, `swap`, and `find_largest` functions, but provide no information about their arguments. Leaving the parentheses empty isn't the same as putting the word `void` between them, which indicates that there are no arguments. The empty-parentheses style of function declaration has largely disappeared. It's inferior to the prototype style introduced in C89, since it doesn't allow the compiler to check whether function calls have the right arguments.

If all declarators were as simple as these, C programming would be a snap. Unfortunately, declarators in actual programs often combine the `*`, `[]`, and `()` notations. We've seen examples of such combinations already. We know that

```
int *ap[10];
```

declares an array of 10 pointers to integers. We know that

```
float *fp(float);
```

declares a function that has a `float` argument and returns a pointer to a `float`. And, in Section 17.7, we learned that

```
void (*pf)(int);
```

declares a pointer to a function with an `int` argument and a `void` return type.

## Deciphering Complex Declarations

So far, we haven't had too much trouble understanding declarators. But what about declarators like the one in the following declaration?

```
int *(*x[10])(void);
```

This declarator combines `*`, `[]`, and `()`, so it's not obvious whether `x` is a pointer, an array, or a function.

Fortunately, there are two simple rules that will allow us to understand any declaration, no matter how convoluted:

- **Always read declarators from the inside out.** In other words, locate the identifier that's being declared, and start deciphering the declaration from there.

- When there's a choice, always favor [] and () over \*. If \* precedes the identifier and [] follows it, the identifier represents an array, not a pointer. Likewise, if \* precedes the identifier and () follows it, the identifier represents a function, not a pointer. (Of course, we can always use parentheses to override the normal priority of [] and () over \*.)

Let's apply these rules to our simple examples first. In the declaration

```
int *ap[10];
```

the identifier is ap. Since \* precedes ap and [] follows it, we give preference to [], so ap is an *array of pointers*. In the declaration

```
float *fp(float);
```

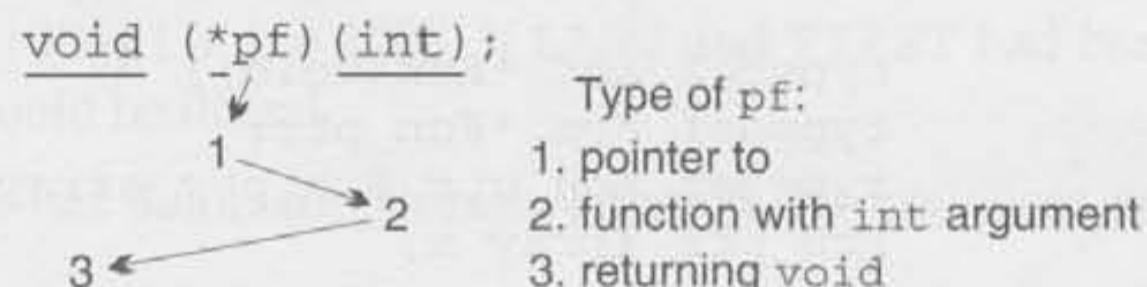
the identifier is fp. Since \* precedes fp and () follows it, we give preference to (), so fp is a *function that returns a pointer*.

The declaration

```
void (*pf)(int);
```

is a little trickier. Since \*pf is enclosed in parentheses, pf must be a pointer. But (\*pf) is followed by (int), so pf must point to a function with an int argument. The word void represents the return type of this function.

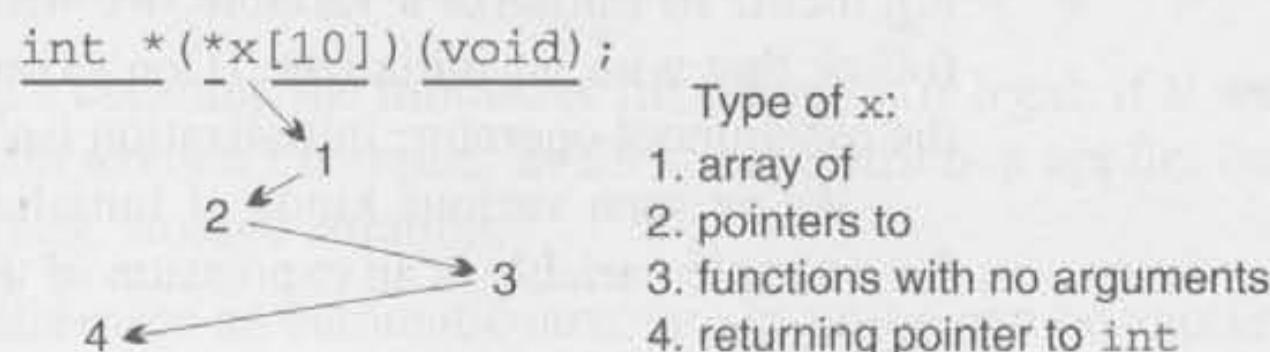
As the last example shows, understanding a complex declarator often involves zigzagging from one side of the identifier to the other:



Let's use this zigzagging technique to decipher the declaration given earlier:

```
int *(*x[10])(void);
```

First, we locate the identifier being declared (x). We see that x is preceded by \* and followed by [] ; since [] have priority over \*, we go right (x is an array). Next, we go left to find out the type of the elements in the array (pointers). Next, we go right to find out what kind of data the pointers point to (functions with no arguments). Finally, we go left to see what each function returns (a pointer to an int). Graphically, here's what the process looks like:



Mastering C declarations takes time and practice. The only good news is that there are certain things that can't be declared in C. Functions can't return arrays:

```
int f(int) [] ;      /*** WRONG ***/
```

Functions can't return functions:

```
int g(int) (int) ;   /*** WRONG ***/
```

Arrays of functions aren't possible, either:

```
int a[10] (int) ;   /*** WRONG ***/
```

In each case, we can use pointers to get the desired effect. A function can't return an array, but it can return a *pointer* to an array. A function can't return a function, but it can return a *pointer* to a function. Arrays of functions aren't allowed, but an array may contain *pointers* to functions. (Section 17.7 has an example of such an array.)

## Using Type Definitions to Simplify Declarations

Some programmers use type definitions to help simplify complex declarations. Consider the declaration of *x* that we examined earlier in this section:

```
int *(*x[10]) (void) ;
```

To make *x*'s type easier to understand, we could use the following series of type definitions:

```
typedef int *Fcn(void) ;
typedef Fcn *Fcn_ptr;
typedef Fcn_ptr Fcn_ptr_array[10];
Fcn_ptr_array x;
```

If we read these lines in reverse order, we see that *x* has type *Fcn\_ptr\_array*, a *Fcn\_ptr\_array* is an array of *Fcn\_ptr* values, a *Fcn\_ptr* is a pointer to type *Fcn*, and a *Fcn* is a function that has no arguments and returns a pointer to an *int* value.

---

## 18.5 Initializers

For convenience, C allows us to specify initial values for variables as we're declaring them. To initialize a variable, we write the = symbol after its declarator, then follow that with an initializer. (Don't confuse the = symbol in a declaration with the assignment operator; initialization isn't the same as assignment.)

We've seen various kinds of initializers in previous chapters. The initializer for a simple variable is an expression of the same type as the variable:

```
int i = 5 / 2; /* i is initially 2 */
```

If the types don't match, C converts the initializer using the same rules as for conversion during assignment ➤ 7.4

```
int j = 5.5;      /* converted to 5 */
```

The initializer for a pointer variable must be a pointer expression of the same type as the variable or of type `void *`:

```
int *p = &i;
```

The initializer for an array, structure, or union is usually a series of values enclosed in braces:

```
int a[5] = {1, 2, 3, 4, 5};
```

**C99**

In C99, brace-enclosed initializers can have other forms, thanks to designated initializers.

To complete our coverage of declarations, let's take a look at some additional rules that govern initializers:

- An initializer for a variable with static storage duration must be constant:

```
#define FIRST 1
#define LAST 100

static int i = LAST - FIRST + 1;
```

Since `LAST` and `FIRST` are macros, the compiler can compute the initial value of `i` ( $100 - 1 + 1 = 100$ ). If `LAST` and `FIRST` had been variables, the initializer would be illegal.

- If a variable has automatic storage duration, its initializer need not be constant:

```
int f(int n)
{
    int last = n - 1;
    ...
}
```

- A brace-enclosed initializer for an array, structure, or union must contain only constant expressions, never variables or function calls:

```
#define N 2

int powers[5] = {1, N, N * N, N * N * N, N * N * N * N};
```

**C99**

Since `N` is a constant, the initializer for `powers` is legal; if `N` were a variable, the program wouldn't compile. In C99, this restriction applies only if the variable has static storage duration.

- The initializer for an automatic structure or union can be another structure or union:

```
void g(struct part part1)
{
    struct part part2 = part1;
    ...
}
```

The initializer doesn't have to be a variable or parameter name, although it does need to be an expression of the proper type. For example, `part2`'s initializer could be `*p`, where `p` is of type `struct part *`, or `f(part1)`, where `f` is a function that returns a `part` structure.

## Uninitialized Variables

In previous chapters, we've implied that uninitialized variables have undefined values. That's not always true; the initial value of a variable depends on its storage duration:

- Variables with *automatic* storage duration have no default initial value. The initial value of an automatic variable can't be predicted and may be different each time the variable comes into existence.
- Variables with *static* storage duration have the value zero by default. Unlike memory allocated by `calloc`, which is simply set to zero bits, a static variable is correctly initialized based on its type: integer variables are initialized to 0, floating variables are initialized to 0.0, and pointer variables contain a null pointer.

calloc function ▶ 17.3

As a matter of style, it's better to provide initializers for static variables rather than rely on the fact that they're guaranteed to be zero. If a program accesses a variable that hasn't been initialized explicitly, someone reading the program later can't easily determine whether the variable is assumed to be zero or whether it's initialized by an assignment somewhere in the program.

---

## 18.6 Inline Functions (C99)

C99 function declarations have an additional option that doesn't exist in C89: they may contain the keyword `inline`. This keyword is a new breed of declaration specifier, distinct from storage classes, type qualifiers, and type specifiers. To understand the effect of `inline`, we'll need to visualize the machine instructions that are generated by a C compiler to handle the process of calling a function and returning from a function.

At the machine level, several instructions may need to be executed to prepare for the call, the call itself requires jumping to the first instruction in the function, and there may be additional instructions executed by the function itself as it begins to execute. If the function has arguments, they'll need to be copied (because C passes its arguments by value). Returning from a function requires a similar

amount of effort on both the part of the function that was called and the one that called it. The cumulative work required to call a function and later return from it is often referred to as “overhead,” since it’s extra work above and beyond what the function is really supposed to accomplish. Although the overhead of a function call slows the program by only a tiny amount, it may add up in certain situations, such as when a function is called millions or billions of times, when an older, slower processor is in use (as might be the case in an embedded system), or when a program has to meet very strict deadlines (as in a real-time system).

parameterized macros ➤ 14.3

In C89, the only way to avoid the overhead of a function call is to use a parameterized macro. Parameterized macros have certain drawbacks, though. C99 offers a better solution to this problem: create an *inline function*. The word “inline” suggests an implementation strategy in which the compiler replaces each call of the function by the machine instructions for the function. This technique avoids the usual overhead of a function call, although it may cause a minor increase in the size of the compiled program.

Declaring a function to be *inline* doesn’t actually force the compiler to “*inline*” the function, however. It merely suggests that the compiler should try to make calls of the function as fast as possible, perhaps by performing an *inline expansion* when the function is called. The compiler is free to ignore this suggestion. In this respect, *inline* is similar to the *register* and *restrict* keywords, which the compiler may use to improve the performance of a program but may also choose to ignore.

## Inline Definitions

An *inline function* has the keyword *inline* as one of its declaration specifiers:

```
inline double average(double a, double b)
{
    return (a + b) / 2;
}
```

Here’s where things get a bit complicated. *average* has external linkage, so other source files may contain calls of *average*. However, the definition of *average* isn’t considered to be an external definition by the compiler (it’s an *inline definition* instead), so attempting to call *average* from another file will be considered an error.

There are two ways to avoid this error. One option is to add the word *static* to the function definition:

```
static inline double average(double a, double b)
{
    return (a + b) / 2;
}
```

*average* now has internal linkage, so it can’t be called from other files. Other files may contain their own definitions of *average*, which might be the same as this definition or might be different.

The other option is to provide an external definition for `average` so that calls are permitted from other files. One way to do this is to write the `average` function a second time (without using `inline`) and put the second definition in a different source file. Doing so is legal, but it's not a good idea to have two versions of the same function, because we can't guarantee that they'll remain consistent when the program is modified.

Here's a better approach. First, we'll put the inline definition of `average` in a header file (let's name it `average.h`):

```
#ifndef AVERAGE_H
#define AVERAGE_H

inline double average(double a, double b)
{
    return (a + b) / 2;
}

#endif
```

Next, we'll create a matching source file, `average.c`:

```
#include "average.h"

extern double average(double a, double b);
```

Now, any file that needs to call the `average` function may simply include `average.h`, which contains the inline definition of `average`. The `average.c` file contains a prototype for `average` that uses the `extern` keyword, which causes the definition of `average` included from `average.h` to be treated as an external definition in `average.c`.

The general rule in C99 is that if all top-level declarations of a function in a particular file include `inline` but not `extern`, then the definition of the function in that file is `inline`. If the function is used anywhere in the program (including the file that contains its `inline` definition), then an external definition of the function will need to be provided by some other file. When the function is called, the compiler may choose to perform an ordinary call (using the function's external definition) or perform inline expansion (using the function's `inline` definition). There's no way to tell which choice the compiler will make, so it's crucial that the two definitions be consistent. The technique that we just discussed (using the `average.h` and `average.c` files) guarantees that the definitions are the same.

## Restrictions on Inline Functions

Since inline functions are implemented in a way that's quite different from ordinary functions, they're subject to different rules and restrictions. Variables with static storage duration are a particular problem for inline functions with external linkage. Consequently, C99 imposes the following restrictions on an inline function with external linkage (but not on one with internal linkage):

- The function may not define a modifiable `static` variable.
- The function may not contain references to variables with internal linkage.

Such a function is allowed to define a variable that is both `static` and `const`, but each inline definition of the function may create its own copy of the variable.

## Using Inline Functions with GCC

Some compilers, including GCC, supported inline functions prior to the C99 standard. As a result, their rules for using inline functions may vary from the standard. In particular, the scheme described earlier (using the `average.h` and `average.c` files) may not work with these compilers. Version 4.3 of GCC (not available at the time this book was written) is expected to support inline functions in the way described in the C99 standard.

Functions that are specified to be both `static` and `inline` should work fine, regardless of the version of GCC. This strategy is legal in C99 as well, so it's the safest bet. A `static inline` function can be used within a single file or placed in a header file and included into any source file that needs to call the function.

There's another way to share an inline function among multiple files that works with older versions of GCC but conflicts with C99. This technique involves putting a definition of the function in a header file, specifying that the function is both `extern` and `inline`, then including the header file into any source file that contains a call of the function. A second copy of the definition—without the words `extern` and `inline`—is placed in one of the source files. (That way, if the compiler is unable to “inline” the function for any reason, it will still have a definition.)

A final note about GCC: Functions are “inlined” only when optimization is requested via the `-O` command-line option.

## Q & A

**\*Q:** Why are selection statements and iteration statements (and their “inner” statements) considered to be blocks in C99? [p. 459]

**C99**

A: This rather surprising rule stems from a problem that can occur when compound literals are used in selection statements and iteration statements. The problem has to do with the storage duration of compound literals, so let's take a moment to discuss that issue first.

The C99 standard states that the object represented by a compound literal has static storage duration if the compound literal occurs outside the body of a function. Otherwise, it has automatic storage duration; as a result, the memory occupied by the object is deallocated at the end of the block in which the compound literal appears. Consider the following function, which returns a `point` structure created using a compound literal:

```
struct point create_point(int x, int y)
{
    return (struct point) {x, y};
}
```

This function works correctly, because the object created by the compound literal will be copied when the function returns. The original object will no longer exist, but the copy will remain. Now suppose that we change the function slightly:

```
struct point *create_point(int x, int y)
{
    return &(struct point) {x, y};
}
```

This version of `create_point` suffers from undefined behavior, because it returns a pointer to an object that has automatic storage duration and won't exist after the function returns.

Now let's return to the question we started with: Why are selection statements and iteration statements considered to be blocks? Consider the following example:

```
/* Example 1 - if statement without braces */

double *coefficients, value;

if (polynomial_selected == 1)
    coefficients = (double[3]) {1.5, -3.0, 6.0};
else
    coefficients = (double[3]) {4.5, 1.0, -3.5};
value = evaluate_polynomial(coefficients);
```

This program fragment apparently behaves in the desired fashion (but read on). `coefficients` will point to one of two objects created by compound literals, and this object will still exist at the time `evaluate_polynomial` is called. Now consider what happens if we put braces around the “inner” statements—the ones controlled by the `if` statement:

```
/* Example 2 - if statement with braces */

double *coefficients, value;

if (polynomial_selected == 1) {
    coefficients = (double[3]) {1.5, -3.0, 6.0};
} else {
    coefficients = (double[3]) {4.5, 1.0, -3.5};
}
value = evaluate_polynomial(coefficients);
```

Now we're in trouble. Each compound literal causes an object to be created, but that object exists only within the block formed by the braces that enclose the statement in which the literal appears. By the time `evaluate_polynomial` is called, `coefficients` points to an object that no longer exists. The result: undefined behavior.

The creators of C99 were unhappy with this state of affairs, because programmers were unlikely to expect that simply adding braces within an `if` statement would cause undefined behavior. To avoid the problem, they decided that the inner statements would always be considered blocks. As a result, Example 1 and Example 2 are equivalent, with both exhibiting undefined behavior.

A similar problem can arise when a compound literal is part of the controlling expression of a selection statement or iteration statement. For this reason, each entire selection statement and iteration statement is considered to be a block as well (as though an invisible set of braces surrounds the entire statement). So, for example, an `if` statement with an `else` clause consists of three blocks: each of the two inner statements is a block, as is the entire `if` statement.

**Q: You said that storage for a variable with automatic storage duration is allocated when the surrounding block is executed. Is this true for C99's variable-length arrays? [p. 460]**

**A:** No. Storage for a variable-length array isn't allocated at the beginning of the surrounding block, because the length of the array isn't yet known. Instead, it's allocated when the declaration of the array is reached during the execution of the block. In this respect, variable-length arrays are different from all other automatic variables.

**Q: What exactly is the difference between "scope" and "linkage"? [p. 460]**

**A:** Scope is for the benefit of the compiler, while linkage is for the benefit of the linker. The compiler uses the scope of an identifier to determine whether or not it's legal to refer to the identifier at a given point in a file. When the compiler translates a source file into object code, it notes which names have external linkage, eventually storing these names in a table inside the object file. Thus, the linker has access to names with external linkage; names with internal linkage or no linkage are invisible to the linker.

**Q: I don't understand how a name could have block scope but external linkage. Could you elaborate? [p. 463]**

**A:** Certainly. Suppose that one source file defines a variable `i`:

```
int i;
```

Let's assume that the definition of `i` lies outside any function, so `i` has external linkage by default. In another file, there's a function `f` that needs to access `i`, so the body of `f` declares `i` as `extern`:

```
void f(void)
{
    extern int i;
    ...
}
```

In the first file, `i` has file scope. Within `f`, however, `i` has block scope. If other functions besides `f` need access to `i`, they'll need to declare it separately. (Or we

can simply move the declaration of *i* outside *f* so that *i* has file scope.) What's confusing about this entire business is that each declaration or definition of *i* establishes a different scope; sometimes it's file scope, and sometimes it's block scope.

**\*Q:** Why can't `const` objects be used in constant expressions? `const` means "constant," right? [p. 466]

A: In C, `const` means "read-only," not "constant." Let's look at a few examples that illustrate why `const` objects can't be used in constant expressions.

To start with, a `const` object might only be constant during its *lifetime*, not throughout the execution of the program. Suppose that a `const` object is declared inside a function:

```
void f(int n)
{
    const int m = n / 2;
    ...
}
```

When *f* is called, *m* will be initialized to the value of *n* / 2. The value of *m* will then remain constant until *f* returns. When *f* is called the next time, *m* will likely be given a different value. That's where the problem arises. Suppose that *m* appears in a `switch` statement:

```
void f(int n)
{
    const int m = n / 2;
    ...
    switch (...) {
        ...
        case m: ... /* *** WRONG ***/
        ...
    }
    ...
}
```

The value of *m* won't be known until *f* is called, which violates C's rule that the values of case labels must be constant expressions.

Next, let's look at `const` objects declared outside blocks. These objects have external linkage and can be shared among files. If C allowed the use of `const` objects in constant expressions, we could easily find ourselves in the following situation:

```
extern const int n;
int a[n]; /* *** WRONG ***/
```

*n* is probably defined in another file, making it impossible for the compiler to determine *a*'s length. (I'm assuming that *a* is an external variable, so it can't be a variable-length array.)

volatile type qualifier ▶ 20.3

If that's not enough to convince you, consider this: If a `const` object is also declared to be `volatile`, its value may change at any time during execution. Here's an example from the C standard:

```
extern const volatile int real_time_clock;
```

The `real_time_clock` variable may not be changed by the program (because it's declared `const`), yet its value may change via some other mechanism (because it's declared `volatile`).

**Q: Why is the syntax of declarators so odd?**

- A: Declarations are intended to mimic use. A pointer declarator has the form `*p`, which matches the way the indirection operator will later be applied to `p`. An array declarator has the form `a [...]`, which matches the way the array will later be subscripted. A function declarator has the form `f (...)`, which matches the syntax of a function call. This reasoning extends to even the most complicated declarators. Consider the `file_cmd` array of Section 17.7, whose elements are pointers to functions. The declarator for `file_cmd` has the form

```
(*file_cmd[]) (void)
```

and a call of one of the functions has the form

```
(*file_cmd[n]) () ;
```

The parentheses, brackets, and `*` are in identical positions.

## Exercises

### Section 18.1

- For each of the following declarations, identify the storage class, type qualifiers, type specifiers, declarators, and initializers.
  - `static char **lookup(int level);`
  - `volatile unsigned long io_flags;`
  - `extern char *file_name[MAX_FILES], path[];`
  - `static const char token_buf[] = "";`

### Section 18.2

- Answer each of the following questions with `auto`, `extern`, `register`, and/or `static`.
  - Which storage class is used primarily to indicate that a variable or function can be shared by several files?
  - Suppose that a variable `x` is to be shared by several functions in one file but hidden from functions in other files. Which storage class should `x` be declared to have?
  - Which storage classes can affect the storage duration of a variable?
- List the storage duration (static or automatic), scope (block or file), and linkage (internal, external, or none) of each variable and parameter in the following file:

```

extern float a;

void f(register double b)
{
    static int c;
    auto char d;
}

```

- W 4. Let *f* be the following function. What will be the value of *f*(10) if *f* has never been called before? What will be the value of *f*(10) if *f* has been called five times previously?

```

int f(int i)
{
    static int j = 0;
    return i * j++;
}

```

5. State whether each of the following statements is true or false. Justify each answer.

- (a) Every variable with static storage duration has file scope.
- (b) Every variable declared inside a function has no linkage.
- (c) Every variable with internal linkage has static storage duration.
- (d) Every parameter has block scope.

6. The following function is supposed to print an error message. Each message is preceded by an integer, indicating the number of times the function has been called. Unfortunately, the function always displays 1 as the number of the error message. Locate the error and show how to fix it without making any changes outside the function.

```

void print_error(const char *message)
{
    int n = 1;
    printf("Error %d: %s\n", n++, message);
}

```

### Section 18.3

7. Suppose that we declare *x* to be a `const` object. Which one of the following statements about *x* is *false*?
- (a) If *x* is of type `int`, it can be used as the value of a case label in a `switch` statement.
  - (b) The compiler will check that no assignment is made to *x*.
  - (c) *x* is subject to the same scope rules as variables.
  - (d) *x* can be of any type.

### Section 18.4

- W 8. Write a complete description of the type of *x* as specified by each of the following declarations.
- (a) `char (*x[10])(int);`
  - (b) `int (*x(int))[5];`
  - (c) `float *(*x(void))(int);`
  - (d) `void (*x(int, void (*y)(int)))(int);`
9. Use a series of type definitions to simplify each of the declarations in Exercise 8.
- W 10. Write declarations for the following variables and functions:
- (a) *p* is a pointer to a function with a character pointer argument that returns a character pointer.









































































































































































































































































































































































































































































































































































































































































