

Lecture Slides for the C++ Programming Language

(Version: 2016-01-18)

Current with the C++14 Standard

Michael D. Adams

Department of Electrical and Computer Engineering
University of Victoria
Victoria, British Columbia, Canada



For additional information and resources related to these lecture slides (including errata and *lecture videos* covering the material on many of these slides), please visit:

<http://www.ece.uvic.ca/~mdadams/cppbook>

If you like these lecture slides, *please show your support* by posting a review of them on Google Play:

<https://play.google.com/store/search?q=Michael%20D%20Adams%20C%2B%2B&c=books>

The author has taken care in the preparation of this document, but makes no expressed or implied warranty of any kind and assumes no responsibility for errors or omissions. No liability is assumed for incidental or consequential damages in connection with or arising out of the use of the information or programs contained herein.

Copyright © 2015, 2016 Michael D. Adams

Published by the University of Victoria, Victoria, British Columbia, Canada

This document is licensed under a Creative Commons Attribution-NonCommercial-NoDerivs 3.0 Unported (CC BY-NC-ND 3.0) License. A copy of this license can be found on page iii of this document. For a simple explanation of the rights granted by this license, see:

<http://creativecommons.org/licenses/by-nc-nd/3.0/>

This document was typeset with L^AT_EX.

ISBN 978-1-55058-582-7 (paperback)

ISBN 978-1-55058-583-4 (PDF)

License I

Creative Commons Legal Code

Attribution-NonCommercial-NoDerivs 3.0 Unported

CREATIVE COMMONS CORPORATION IS NOT A LAW FIRM AND DOES NOT PROVIDE LEGAL SERVICES. DISTRIBUTION OF THIS LICENSE DOES NOT CREATE AN ATTORNEY-CLIENT RELATIONSHIP. CREATIVE COMMONS PROVIDES THIS INFORMATION ON AN "AS-IS" BASIS. CREATIVE COMMONS MAKES NO WARRANTIES REGARDING THE INFORMATION PROVIDED, AND DISCLAIMS LIABILITY FOR DAMAGES RESULTING FROM ITS USE.

License

THE WORK (AS DEFINED BELOW) IS PROVIDED UNDER THE TERMS OF THIS CREATIVE COMMONS PUBLIC LICENSE ("CCPL" OR "LICENSE"). THE WORK IS PROTECTED BY COPYRIGHT AND/OR OTHER APPLICABLE LAW. ANY USE OF THE WORK OTHER THAN AS AUTHORIZED UNDER THIS LICENSE OR COPYRIGHT LAW IS PROHIBITED.

BY EXERCISING ANY RIGHTS TO THE WORK PROVIDED HERE, YOU ACCEPT AND AGREE TO BE BOUND BY THE TERMS OF THIS LICENSE. TO THE EXTENT THIS LICENSE MAY BE CONSIDERED TO BE A CONTRACT, THE LICENSOR GRANTS YOU THE RIGHTS CONTAINED HERE IN CONSIDERATION OF YOUR ACCEPTANCE OF SUCH TERMS AND CONDITIONS.

1. Definitions

- a. "Adaptation" means a work based upon the Work, or upon the Work and other pre-existing works, such as a translation, adaptation, derivative work, arrangement of music or other alterations of a literary or artistic work, or phonogram or performance and includes cinematographic adaptations or any other form in which the Work may be recast, transformed, or adapted including in any form recognizably derived from the original, except that a work that constitutes a Collection will not be considered an Adaptation for the purpose of this License. For the avoidance of doubt, where the Work is a musical work, performance or phonogram, the synchronization of the Work in timed-relation with a moving image ("synching") will be considered an Adaptation for the purpose of this License.
- b. "Collection" means a collection of literary or artistic works, such as encyclopedias and anthologies, or performances, phonograms or

License II

broadcasts, or other works or subject matter other than works listed in Section 1(f) below, which, by reason of the selection and arrangement of their contents, constitute intellectual creations, in which the Work is included in its entirety in unmodified form along with one or more other contributions, each constituting separate and independent works in themselves, which together are assembled into a collective whole. A work that constitutes a Collection will not be considered an Adaptation (as defined above) for the purposes of this License.

- c. "Distribute" means to make available to the public the original and copies of the Work through sale or other transfer of ownership.
- d. "Licensor" means the individual, individuals, entity or entities that offer(s) the Work under the terms of this License.
- e. "Original Author" means, in the case of a literary or artistic work, the individual, individuals, entity or entities who created the Work or if no individual or entity can be identified, the publisher; and in addition (i) in the case of a performance the actors, singers, musicians, dancers, and other persons who act, sing, deliver, declaim, play in, interpret or otherwise perform literary or artistic works or expressions of folklore; (ii) in the case of a phonogram the producer being the person or legal entity who first fixes the sounds of a performance or other sounds; and, (iii) in the case of broadcasts, the organization that transmits the broadcast.
- f. "Work" means the literary and/or artistic work offered under the terms of this License including without limitation any production in the literary, scientific and artistic domain, whatever may be the mode or form of its expression including digital form, such as a book, pamphlet and other writing; a lecture, address, sermon or other work of the same nature; a dramatic or dramatico-musical work; a choreographic work or entertainment in dumb show; a musical composition with or without words; a cinematographic work to which are assimilated works expressed by a process analogous to cinematography; a work of drawing, painting, architecture, sculpture, engraving or lithography; a photographic work to which are assimilated works expressed by a process analogous to photography; a work of applied art; an illustration, map, plan, sketch or three-dimensional work relative to geography, topography, architecture or science; a performance; a broadcast; a phonogram; a compilation of data to the extent it is protected as a copyrightable work; or a work performed by a variety or circus performer to the extent it is not otherwise considered a literary or artistic work.

License III

- g. "You" means an individual or entity exercising rights under this License who has not previously violated the terms of this License with respect to the Work, or who has received express permission from the Licensor to exercise rights under this License despite a previous violation.
- h. "Publicly Perform" means to perform public recitations of the Work and to communicate to the public those public recitations, by any means or process, including by wire or wireless means or public digital performances; to make available to the public Works in such a way that members of the public may access these Works from a place and at a place individually chosen by them; to perform the Work to the public by any means or process and the communication to the public of the performances of the Work, including by public digital performance; to broadcast and rebroadcast the Work by any means including signs, sounds or images.
- i. "Reproduce" means to make copies of the Work by any means including without limitation by sound or visual recordings and the right of fixation and reproducing fixations of the Work, including storage of a protected performance or phonogram in digital form or other electronic medium.

2. Fair Dealing Rights. Nothing in this License is intended to reduce, limit, or restrict any uses free from copyright or rights arising from limitations or exceptions that are provided for in connection with the copyright protection under copyright law or other applicable laws.

3. License Grant. Subject to the terms and conditions of this License, Licensor hereby grants You a worldwide, royalty-free, non-exclusive, perpetual (for the duration of the applicable copyright) license to exercise the rights in the Work as stated below:

- a. to Reproduce the Work, to incorporate the Work into one or more Collections, and to Reproduce the Work as incorporated in the Collections; and,
- b. to Distribute and Publicly Perform the Work including as incorporated in Collections.

The above rights may be exercised in all media and formats whether now known or hereafter devised. The above rights include the right to make such modifications as are technically necessary to exercise the rights in other media and formats, but otherwise you have no rights to make

License IV

Adaptations. Subject to 8(f), all rights not expressly granted by Licensor are hereby reserved, including but not limited to the rights set forth in Section 4(d).

4. Restrictions. The license granted in Section 3 above is expressly made subject to and limited by the following restrictions:

- a. You may Distribute or Publicly Perform the Work only under the terms of this License. You must include a copy of, or the Uniform Resource Identifier (URI) for, this License with every copy of the Work You Distribute or Publicly Perform. You may not offer or impose any terms on the Work that restrict the terms of this License or the ability of the recipient of the Work to exercise the rights granted to that recipient under the terms of the License. You may not sublicense the Work. You must keep intact all notices that refer to this License and to the disclaimer of warranties with every copy of the Work You Distribute or Publicly Perform. When You Distribute or Publicly Perform the Work, You may not impose any effective technological measures on the Work that restrict the ability of a recipient of the Work from You to exercise the rights granted to that recipient under the terms of the License. This Section 4(a) applies to the Work as incorporated in a Collection, but this does not require the Collection apart from the Work itself to be made subject to the terms of this License. If You create a Collection, upon notice from any Licensor You must, to the extent practicable, remove from the Collection any credit as required by Section 4(c), as requested.
- b. You may not exercise any of the rights granted to You in Section 3 above in any manner that is primarily intended for or directed toward commercial advantage or private monetary compensation. The exchange of the Work for other copyrighted works by means of digital file-sharing or otherwise shall not be considered to be intended for or directed toward commercial advantage or private monetary compensation, provided there is no payment of any monetary compensation in connection with the exchange of copyrighted works.
- c. If You Distribute, or Publicly Perform the Work or Collections, You must, unless a request has been made pursuant to Section 4(a), keep intact all copyright notices for the Work and provide, reasonable to the medium or means You are utilizing: (i) the name of the Original Author (or pseudonym, if applicable) if supplied, and/or if the Original Author and/or Licensor designate another party or parties (e.g., a sponsor institute, publishing entity, journal) for

attribution ("Attribution Parties") in Licensor's copyright notice, terms of service or by other reasonable means, the name of such party or parties; (ii) the title of the Work if supplied; (iii) to the extent reasonably practicable, the URI, if any, that Licensor specifies to be associated with the Work, unless such URI does not refer to the copyright notice or licensing information for the Work. The credit required by this Section 4(c) may be implemented in any reasonable manner; provided, however, that in the case of a Collection, at a minimum such credit will appear, if a credit for all contributing authors of Collection appears, then as part of these credits and in a manner at least as prominent as the credits for the other contributing authors. For the avoidance of doubt, You may only use the credit required by this Section for the purpose of attribution in the manner set out above and, by exercising Your rights under this License, You may not implicitly or explicitly assert or imply any connection with, sponsorship or endorsement by the Original Author, Licensor and/or Attribution Parties, as appropriate, of You or Your use of the Work, without the separate, express prior written permission of the Original Author, Licensor and/or Attribution Parties.

d. For the avoidance of doubt:

- i. Non-waivable Compulsory License Schemes. In those jurisdictions in which the right to collect royalties through any statutory or compulsory licensing scheme cannot be waived, the Licensor reserves the exclusive right to collect such royalties for any exercise by You of the rights granted under this License;
- ii. Waivable Compulsory License Schemes. In those jurisdictions in which the right to collect royalties through any statutory or compulsory licensing scheme can be waived, the Licensor reserves the exclusive right to collect such royalties for any exercise by You of the rights granted under this License if Your exercise of such rights is for a purpose or use which is otherwise than noncommercial as permitted under Section 4(b) and otherwise waives the right to collect royalties through any statutory or compulsory licensing scheme; and,
- iii. Voluntary License Schemes. The Licensor reserves the right to collect royalties, whether individually or, in the event that the Licensor is a member of a collecting society that administers voluntary licensing schemes, via that society, from any exercise by You of the rights granted under this License that is for a

License VI

purpose or use which is otherwise than noncommercial as permitted under Section 4(b).

- e. Except as otherwise agreed in writing by the Licensor or as may be otherwise permitted by applicable law, if You Reproduce, Distribute or Publicly Perform the Work either by itself or as part of any Collections, You must not distort, mutilate, modify or take other derogatory action in relation to the Work which would be prejudicial to the Original Author's honor or reputation.

5. Representations, Warranties and Disclaimer

UNLESS OTHERWISE MUTUALLY AGREED BY THE PARTIES IN WRITING, LICENSOR OFFERS THE WORK AS-IS AND MAKES NO REPRESENTATIONS OR WARRANTIES OF ANY KIND CONCERNING THE WORK, EXPRESS, IMPLIED, STATUTORY OR OTHERWISE, INCLUDING, WITHOUT LIMITATION, WARRANTIES OF TITLE, MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE, NONINFRINGEMENT, OR THE ABSENCE OF LATENT OR OTHER DEFECTS, ACCURACY, OR THE PRESENCE OF ABSENCE OF ERRORS, WHETHER OR NOT DISCOVERABLE. SOME JURISDICTIONS DO NOT ALLOW THE EXCLUSION OF IMPLIED WARRANTIES, SO SUCH EXCLUSION MAY NOT APPLY TO YOU.

6. Limitation on Liability. EXCEPT TO THE EXTENT REQUIRED BY APPLICABLE LAW, IN NO EVENT WILL LICENSOR BE LIABLE TO YOU ON ANY LEGAL THEORY FOR ANY SPECIAL, INCIDENTAL, CONSEQUENTIAL, PUNITIVE OR EXEMPLARY DAMAGES ARISING OUT OF THIS LICENSE OR THE USE OF THE WORK, EVEN IF LICENSOR HAS BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

7. Termination

- a. This License and the rights granted hereunder will terminate automatically upon any breach by You of the terms of this License. Individuals or entities who have received Collections from You under this License, however, will not have their licenses terminated provided such individuals or entities remain in full compliance with those licenses. Sections 1, 2, 5, 6, 7, and 8 will survive any termination of this License.
- b. Subject to the above terms and conditions, the license granted here is perpetual (for the duration of the applicable copyright in the Work). Notwithstanding the above, Licensor reserves the right to release the Work under different license terms or to stop distributing the Work at any time; provided, however that any such election will not serve to withdraw this License (or any other license that has been, or is

License VII

required to be, granted under the terms of this License), and this License will continue in full force and effect unless terminated as stated above.

8. Miscellaneous

- a. Each time You Distribute or Publicly Perform the Work or a Collection, the Licensor offers to the recipient a license to the Work on the same terms and conditions as the license granted to You under this License.
- b. If any provision of this License is invalid or unenforceable under applicable law, it shall not affect the validity or enforceability of the remainder of the terms of this License, and without further action by the parties to this agreement, such provision shall be reformed to the minimum extent necessary to make such provision valid and enforceable.
- c. No term or provision of this License shall be deemed waived and no breach consented to unless such waiver or consent shall be in writing and signed by the party to be charged with such waiver or consent.
- d. This License constitutes the entire agreement between the parties with respect to the Work licensed here. There are no understandings, agreements or representations with respect to the Work not specified here. Licensor shall not be bound by any additional provisions that may appear in any communication from You. This License may not be modified without the mutual written agreement of the Licensor and You.
- e. The rights granted under, and the subject matter referenced, in this License were drafted utilizing the terminology of the Berne Convention for the Protection of Literary and Artistic Works (as amended on September 28, 1979), the Rome Convention of 1961, the WIPO Copyright Treaty of 1996, the WIPO Performances and Phonograms Treaty of 1996 and the Universal Copyright Convention (as revised on July 24, 1971). These rights and subject matter take effect in the relevant jurisdiction in which the License terms are sought to be enforced according to the corresponding provisions of the implementation of those treaty provisions in the applicable national law. If the standard suite of rights granted under applicable copyright law includes additional rights not granted under this License, such additional rights are deemed to be included in the License; this License is not intended to restrict the license of any rights under applicable law.

License VIII

Creative Commons Notice

Creative Commons is not a party to this License, and makes no warranty whatsoever in connection with the Work. Creative Commons will not be liable to You or any party on any legal theory for any damages whatsoever, including without limitation any general, special, incidental or consequential damages arising in connection to this license. Notwithstanding the foregoing two (2) sentences, if Creative Commons has expressly identified itself as the Licensor hereunder, it shall have all rights and obligations of Licensor.

Except for the limited purpose of indicating to the public that the Work is licensed under the CCPL, Creative Commons does not authorize the use by either party of the trademark "Creative Commons" or any related trademark or logo of Creative Commons without the prior written consent of Creative Commons. Any permitted use will be in compliance with Creative Commons' then-current trademark usage guidelines, as may be published on its website or otherwise made available upon request from time to time. For the avoidance of doubt, this trademark restriction does not form part of this License.

Creative Commons may be contacted at <http://creativecommons.org/>.

Other Textbooks and Lecture Slides by the Author I

- 1 M. D. Adams, *Multiresolution Signal and Geometry Processing: Filter Banks, Wavelets, and Subdivision (Version 2013-09-26)*, University of Victoria, Victoria, BC, Canada, Sept. 2013, xxxviii + 538 pages, ISBN 978-1-55058-507-0 (print), ISBN 978-1-55058-508-7 (PDF). Available from Google Books, Google Play Books, University of Victoria Bookstore, and author's web site <http://www.ece.uvic.ca/~mdadams/waveletbook>.
- 2 M. D. Adams, *Lecture Slides for Multiresolution Signal and Geometry Processing (Version 2015-02-03)*, University of Victoria, Victoria, BC, Canada, Feb. 2015, xi + 587 slides, ISBN 978-1-55058-535-3 (print), ISBN 978-1-55058-536-0 (PDF). Available from Google Books, Google Play Books, University of Victoria Bookstore, and author's web site <http://www.ece.uvic.ca/~mdadams/waveletbook>.

Other Textbooks and Lecture Slides by the Author II

- ③ M. D. Adams, *Continuous-Time Signals and Systems (Version 2013-09-11)*, University of Victoria, Victoria, BC, Canada, Sept. 2013, xxx + 308 pages, ISBN 978-1-55058-495-0 (print), ISBN 978-1-55058-506-3 (PDF). Available from Google Books, Google Play Books, University of Victoria Bookstore, and author's web site <http://www.ece.uvic.ca/~mdadams/sigsysbook>.
- ④ M. D. Adams, *Lecture Slides for Continuous-Time Signals and Systems (Version 2013-09-11)*, University of Victoria, Victoria, BC, Canada, Dec. 2013, 286 slides, ISBN 978-1-55058-517-9 (print), ISBN 978-1-55058-518-6 (PDF). Available from Google Books, Google Play Books, University of Victoria Bookstore, and author's web site <http://www.ece.uvic.ca/~mdadams/sigsysbook>.

- 5 M. D. Adams, *Lecture Slides for Signals and Systems (Version 2016-01-25)*, University of Victoria, Victoria, BC, Canada, Jan. 2016, xvi + 481 slides, ISBN 978-1-55058-584-1 (print), ISBN 978-1-55058-585-8 (PDF). Available from Google Books, Google Play Books, University of Victoria Bookstore, and author's web site <http://www.ece.uvic.ca/~mdadams/sigsysbook>.

Part 0

Preface

About These Lecture Slides

- This document constitutes a detailed set of lecture slides on the C++ programming language and is current with the *C++14* standard.
- Many aspects of the C++ language are covered from introductory to more advanced.
- Some aspects of the C++ standard library are also introduced.
- In addition, various general programming-related topics are considered.

Acknowledgments

- The author would like to thank Robert Leahy for reviewing various drafts of many of these slides and providing many useful comments that allowed the quality of these materials to be improved significantly.

- Many code examples are included throughout these slides.
- Often, in order to make an example short enough to fit on a slide, compromises had to be made in terms of good programming style.
- These deviations from good style include (but are not limited to) such things as:
 - ① frequently formatting source code in unusual ways to conserve vertical space in listings;
 - ② not fully documenting source code with comments;
 - ③ using short meaningless identifier names; and
 - ④ engaging other evil behavior such as using many global variables and employing constructs like “**using namespace** std;”.

Typesetting Conventions

- In a definition, the term being defined is often typeset in a font **like this**.
- To emphasize particular words, the words are typeset in a font *like this*.

Part 1

Software

Why Is Software Important?

- almost all electronic devices run some software
- automobile engine control system, implantable medical devices, remote controls, office machines (e.g., photocopiers), appliances (e.g., televisions, refrigerators, washers/dryers, dishwashers, air conditioner), power tools, toys, mobile phones, media players, computers, printers, photocopies, disk drives, scanners, webcams, MRI machines

Why Software-Based Solutions?

- more cost effective to implement functionality in software than hardware
- software bugs easy to fix, give customer new software upgrade
- hardware bugs extremely costly to repair, customer sends in old device and manufacturer sends replacement
- systems increasingly complex, bugs unavoidable
- allows new features to be added later
- implement only absolute minimal functionality in hardware, do the rest in software

Software-Related Jobs

- many more software jobs than hardware jobs
- relatively small team of hardware designers produce platform like iPhone
- thousands of companies develop applications for platform
- only implement directly in hardware when absolutely necessary (e.g., for performance reasons)

Which Language to Learn?

- C, C++, Fortran, Java, MATLAB, C#, Objective C
- programming language popularity
- <http://www.tiobe.com/> TIOBE Software Programming Community Index Jan 2011 all in top four: Java, C, C++ MATLAB (23rd) Fortran (27th)
- Programming Language Popularity Normalized Comparison <http://www.langpop.com/> top three languages: C, Java, C++
- international standard
- vendor neutral

- created by Dennis Ritchie, AT&T Bell Labs in 1970s
- international standard ISO/IEC 9899:2011 (informally known as “C11”)
- available on wide range of platforms, from microcontrollers to supercomputers; very few platforms for which C compiler not available
- procedural, provides language constructs that map efficiently to machine instructions
- does not directly support object-oriented or generic programming
- application domains: system software, device drivers, embedded applications, application software
- greatly influenced development of C++
- when something lasts in computer industry for more than 40 years (outliving its creator), must be good

- created by Bjarne Stroustrup, Bell Labs
- originally C with Classes, renamed as C++ in 1983
- most recent specification of language in ISO/IEC 14882:2014 (informally known as “C++14”)
- procedural
- loosely speaking is superset of C
- directly supports object-oriented and generic programming
- maintains efficiency of C
- application domains: systems software, application software, device drivers, embedded software, high-performance server and client applications, entertainment software such as video games, native code for Android applications
- greatly influenced development of C# and Java

- developed in 1990s by James Gosling at Sun Microsystems (later bought by Oracle Corporation)
- de facto standard but not international standard
- usually less efficient than C and C++
- simplified memory management (with garbage collection)
- direct support for object-oriented programming
- application domains: web applications, Android applications

- proprietary language, developed by The MathWorks
- not general-purpose programming language
- application domain: numerical computing
- used to design and simulate systems
- not used to implement real-world systems

- designed by John Backus, IBM, in 1950s
- international standard ISO/IEC 1539-1:2010 (informally known as "Fortran 2008")
- application domain: scientific and engineering applications, intensive supercomputing tasks such as weather and climate modelling, finite element analysis, computational fluid dynamics, computational physics, computational chemistry

- developed by Microsoft, team led by Anders Hejlsberg
- ECMA-334 and ISO/IEC 23270:2006
- most recent language specifications not standardized by ECMA or ISO/IEC
- intellectual property concerns over Microsoft patents
- object oriented

Objective C

- developed by Tom Love and Brad Cox of Stepstone (later bought by NeXT and subsequently Apple)
- used primarily on Apple Mac OS X and iOS
- strict superset of C
- no official standard that describes Objective C
- authoritative manual on Objective-C 2.0 available from Apple

Why Learn C++?

- vendor neutral
- international standard
- general purpose
- powerful yet efficient
- loosely speaking, includes C as subset; so can learn two languages (C++ and C) for price of one
- easy to move from C++ to other languages but often not in other direction
- many other popular languages inspired by C++

Part 2

C++

Section 2.1

History of C++

Motivation

- developed by Bjarne Stroustrup starting in 1979 at Computing Science Research Center of Bell Laboratories, Murray Hill, NJ, USA
- doctoral work in Computing Laboratory of University of Cambridge, Cambridge, UK
- study alternatives for organization of system software for distributed systems
- required development of relatively large and detailed simulator
- dissertation:
 - B. Stroustrup. *Communication and Control in Distributed Computer Systems*.
PhD thesis, University of Cambridge, Cambridge, UK, 1979.
- in 1979, joined Bell Laboratories after having finished doctorate
- work started with attempt to analyze UNIX kernel to determine to what extent it could be distributed over network of computers connected by LAN
- needed way to model module structure of system and pattern of communication between modules
- no suitable tools available

Objectives

- had bad experiences writing simulator during Ph.D. studies; originally used Simula for simulator; later forced to rewrite in BCPL for speed; more low level than C; BCPL was horrible to use
- notion of what properties good tool would have motivated by these experiences
- suitable tool for projects like simulator, operating system, other systems programming tasks should:
 - support for effective program organization (like in Simula) (i.e., classes, some form of class hierarchies, some form of support for concurrency, strong checking of type system based on classes)
 - produce programs that run fast (like with BCPL)
 - be able to easily combine separately compilable units into program (like with BCPL)
 - have simple linkage convention, essential for combining units written in languages such as C, Algol68, Fortran, BCPL, assembler into single program
 - allow highly portable implementations (only very limited ties to operating system)

Timeline for C with Classes (1979–1983) I

- May 1979 work on C with Classes starts
- Oct 1979 initial version of Cpre, preprocessor that added Simula-like classes to C; language accepted by preprocessor later started being referred to as C with Classes
- Mar 1980 Cpre supported one real project and several experiments (used on about 16 systems)
- Apr 1980 first internal Bell Labs paper on C with Classes published (later to appear in ACM SIGPLAN Notices in Jan. 1982)

B. Stroustrup. [Classes: An abstract data type facility for the C language.](#)

Bell Laboratories Computer Science Technical Report CSTR-84, Apr. 1980.

Timeline for C with Classes (1979–1983) II

1980 initial 1980 implementation had following features:

- classes
- derived classes
- public/private access control
- constructors and destructors
- call and return functions (call function implicitly called before every call of every member function; return function implicitly called after every return from every member function; can be used for synchronization)
- friend classes
- type checking and conversion of function arguments

1981 in 1981, added:

- inline functions
- default arguments
- overloading of assignment operator

Jan 1982 first external paper on C with Classes published

Timeline for C with Classes (1979–1983) III

B. Stroustrup. [Classes: An abstract data type facility for the C language.](#)

ACM SIGPLAN Notices, 17(1):42–51, Jan. 1982.

Feb 1983 more detailed paper on C with Classes published

B. Stroustrup. [Adding classes to the C language: An exercise in language evolution.](#)

Software: Practice and Experience, 13(2):139–161, Feb. 1983.

- C with Classes proved very successful; generated considerable interest
- first real application of C with Classes was network simulators

Timeline for C84 to C++98 (1982–1998) I

- started to work on cleaned up and extended successor to C with Classes, initially called C84 and later renamed C++

Spring 1982 started work on Cfront compiler front-end for C84; initially written in C with Classes and then transcribed to C84; traditional compiler front-end performing complete check of syntax and semantics of language, building internal representation of input, analyzing and rearranging representation, and finally producing output for some code generator; generated C code as output; difficult to bootstrap on machine without C84 compiler; Cfront software included special “half-processed” version of C code resulting from compiling Cfront, which could be compiled with native C compiler and resulting executable then used to compile Cfront

Timeline for C84 to C++98 (1982–1998) II

Dec 1983 C84 (C with Classes) renamed C++;
name used in following paper prepared in Dec. 1983

B. Stroustrup. [Data abstraction in C.](#)

Bell Labs Technical Journal, 63(8):1701–1732, Oct. 1984.

(name C++ suggested by Rick Mascitti)

1983 virtual functions added

Note: going from C with Classes to C84 added: virtual functions,
function name and operator overloading, references, constants
(**const**), user-controlled free-store memory control, improved
type checking

Jan 1984 first C++ manual

B. Stroustrup. [The C++ reference manual.](#)

AT&T Bell Labs Computer Science Technical Report No.
108, Jan. 1984.

Sep 1984 paper describing operator overloading published

Timeline for C84 to C++98 (1982–1998) III

B. Stroustrup. [Operator overloading in C++](#).

In *Proc. IFIP WG2.4 Conference on System Implementation Languages: Experience & Assessment*, Sept. 1984.

1984 stream I/O library first implemented and later presented in

B. Stroustrup. [An extensible I/O facility for C++](#).

In *Proc. of Summer 1985 USENIX Conference*, pages 57–70, June 1985.

Feb 1985 Cfront Release E (first external release); “E” for “Educational”; available to universities

Oct 1985 Cfront Release 1.0 (first commercial release)

Oct 1985 first edition of C++PL written

B. Stroustrup. [The C++ Programming Language](#).

Addison Wesley, 1986.

Timeline for C84 to C++98 (1982–1998) IV

(Cfront Release 1.0 corresponded to language as defined in this book)

Oct 1985 tutorial paper on C++

B. Stroustrup. [A C++ tutorial.](#)

In *Proceedings of the ACM annual conference on the range of computing: mid-80's perspective*, pages 56–64, Oct. 1985.

Jun 1986 Cfront Release 1.1; mainly bug fix release

Aug 1986 first exposition of set of techniques for which C++ was aiming to provide support (rather than what features are already implemented and in use)

B. Stroustrup. [What is object-oriented programming?](#)

In *Proc. of 14th Association of Simula Users Conference*, Stockholm, Sweden, Aug. 1986.

Timeline for C84 to C++98 (1982–1998) V

- Sep 1986 first Object-Oriented Programming, Systems, Languages, and Applications (OOPSLA) conference (start of OO hype centered on Smalltalk)
- Nov 1986 first commercial Cfront PC port (Cfront 1.1, Glockenspiel [in Ireland])
- Feb 1987 Cfront Release 1.2; primarily bug fixes but also added:
- pointers to members
 - protected members
- Nov 1987 first conference devoted to C++:
USENIX C++ conference (Santa Fe, NM, USA)
- Dec 1987 first GNU C++ release (1.13)
- Jan 1988 first Oregon Software (a.k.a. TauMetric) C++ release
- Jun 1988 first Zortech C++ release
- Oct 1988 first presented templates at USENIX C++ conference (Denver, CO, USA) in paper:

Timeline for C84 to C++98 (1982–1998) VI

B. Stroustrup. [Parameterized types for C++](#).

In *Proc. of USENIX C++ Conference*, pages 1–18, Denver, CO, USA, Oct. 1988.

Oct 1988 first USENIX C++ implementers workshop (Estes Park, CO, USA)

Jan 1989 first C++ journal “The C++ Report” (from SIGS publications) started publishing

Jun 1989 Cfront Release 2.0 major cleanup; new features included:

- multiple inheritance
- type-safe linkage
- better resolution of overloaded functions
- recursive definition of assignment and initialization
- better facilities for user-defined memory management
- abstract classes
- static member functions
- const member functions

Timeline for C84 to C++98 (1982–1998) VII

- protected member functions (first provided in release 1.2)
- overloading of operator \rightarrow
- pointers to members (first provided in release 1.2)

1989 main features of Cfront 2.0 summarized in

B. Stroustrup. [The evolution of C++: 1985–1989](#).
USENIX Computer Systems, 2(3), Summer 1989.

first presented in

B. Stroustrup. [The evolution of C++: 1985–1987](#).
In *Proc. of USENIX C++ Conference*, pages 1–22, Santa Fe, NM, USA, Nov. 1987.

Nov 1989 paper describing exceptions published

A. Koenig and B. Stroustrup. [Exception handling for C++](#).
In *Proc. of “C++ at Work” Conference*, Nov. 1989.

followed up by

Timeline for C84 to C++98 (1982–1998) VIII

A. Koenig and B. Stroustrup. [Exception handling for C++](#).
In *Proc. of USENIX C++ Conference*, Apr. 1990.

Dec 1989 ANSI X3J16 organizational meeting (Washington, DC, USA)

Mar 1990 first ANSI X3J16 technical meeting (Somerset, NJ, USA)

Apr 1990 Cfront Release 2.1; bug fix release to bring Cfront mostly into line with ARM

May 1990 annotated reference manual (ARM) published

M. A. Ellis and B. Stroustrup. [The Annotated C++ Reference Manual](#).

Addison Wesley, May 1990.

(formed basis for ANSI standardization)

May 1990 first Borland C++ release

Jul 1990 templates accepted (Seattle, WA, USA)

Nov 1990 exceptions accepted (Palo Alto, CA, USA)

Timeline for C84 to C++98 (1982–1998) IX

Jun 1991 second edition of C++PL published

B. Stroustrup. *The C++ Programming Language*.
Addison Wesley, 2nd edition, June 1991.

Jun 1991 first ISO WG21 meeting (Lund, Sweden)

Sep 1991 Cfront Release 3.0; added templates (as specified in ARM)

Oct 1991 estimated number of C++ users 400,000

Feb 1992 first DEC C++ release (including templates and exceptions)

Mar 1992 run-time type identification (RTTI) described in

B. Stroustrup and D. Lenkov. [Run-time type identification for C++](#).

The C++ Report, Mar. 1992.

(RTTI in C++ based on this paper)

Mar 1992 first Microsoft C++ release (did not support templates or exceptions)

Timeline for C84 to C++98 (1982–1998) X

- May 1992 first IBM C++ release (including templates and exceptions)
- Mar 1993 RTTI accepted (Portland, OR, USA)
 - Jul 1993 namespaces accepted (Munich, Germany)
 - 1993 further work on Cfront Release 4.0 abandoned after failed attempt to add exception support
- Aug 1994 ANSI/ISO Committee Draft registered
- Aug 1994 Standard Template Library (STL) accepted (Waterloo, ON, CA); described in
 - A. Stepanov and M. Lee. [The standard template library](#). Technical Report HPL-94-34 (R.1), HP Labs, Aug. 1994.
- Aug 1996 **export** accepted (Stockholm, Sweden)
 - 1997 third edition of C++PL published
 - B. Stroustrup. [The C++ Programming Language](#). Addison Wesley Longman, Reading, MA, USA, 3rd edition, 1997.

Timeline for C84 to C++98 (1982–1998) XI

- Nov 1997 final committee vote on complete standard (Morristown, NJ, USA)
- Jul 1998 Microsoft releases VC++ 6.0, first Microsoft compiler to provide close-to-complete set of ISO C++
- Sep 1998 ISO/IEC 14882:1998 (informally known as C++98) published
ISO/IEC 14882:1998 — programming languages — C++,
Sept. 1998.
- 1998 Beman Dawes starts Boost (provides peer-reviewed portable C++ source libraries)
- Feb 2000 special edition of C++PL published
B. Stroustrup. *The C++ Programming Language*.
Addison Wesley, Reading, MA, USA, special edition, Feb.
2000.

Timeline After C++98 (1998–Present) I

- Apr 2001** motion passed to request new work item: technical report on libraries (Copenhagen, Denmark); later to become ISO/IEC TR 19768:2007
- Oct 2003** ISO/IEC 14882:2003 (informally known as C++03) published; essentially bug fix release; no changes to language from programmer's point of view
 - ISO/IEC 14882:2003 — programming languages — C++, Oct. 2003.
- 2003** work on C++0x (now known as C++11) starts
- Oct 2004** estimated number of C++ users 3,270,000
- Apr 2005** first votes on features for C++0x (Lillehammer, Norway)
 - 2005** **auto**, **static_assert**, and rvalue references accepted in principle
- Apr 2006** first full committee (official) votes on features for C++0x (Berlin, Germany)

Timeline After C++98 (1998–Present) II

Sep 2006 performance technical report (TR 18015) published:
ISO/IEC TR 18015:2006 — information technology — programming languages, their environments and system software interfaces — technical report on C++ performance, Sept. 2006.

work spurred by earlier proposal to standardize subset of C++ for embedded systems called Embedded C++ (or just EC++); EC++ motivated by performance concerns

Apr 2006 decision to move special mathematical functions to separate ISO standard (Berlin, Germany); deemed too specialized for most programmers

Nov 2007 ISO/IEC TR 19768:2007 (informally known as C++TR1) published;

ISO/IEC TR 19768:2007 — information technology — programming languages — technical report on C++ library extensions, Nov. 2007.

Timeline After C++98 (1998–Present) III

specifies series of library extensions to be considered for adoption later in C++

2009 another particularly notable book on C++ published

B. Stroustrup. *Programming: Principles and Practice Using C++*.

Addison Wesley, Upper Saddle River, NJ, USA, 2009.

Aug 2011 ISO/IEC 14882:2011 (informally known as C++11) ratified ISO/IEC 14882:2011 — information technology — programming languages — C++, Sept. 2011.

2013 fourth edition of C++PL published

B. Stroustrup. *The C++ Programming Language*.

Addison Wesley, 4th edition, 2013.

2014 ISO/IEC 14882:2014 (informally known as C++14) ratified ISO/IEC 14882:2014 — information technology — programming languages — C++, 2014.

Additional Comments

- reasons for using C as starting point:
 - flexibility (can be used for most application areas)
 - efficiency
 - availability (C compilers available for most platforms)
 - portability (source code relatively portable from one platform to another)
- main sources for ideas for C++ (aside from C) were Simula, Algol68, BCPL, Ada, Clu, ML; in particular:
 - Simula gave classes
 - Algol68 gave operator overloading, references, ability to declare variables anywhere in block
 - BCPL gave `//` comments
 - exceptions influenced by ML
 - templates influenced by generics in Ada and parameterized modules in Clu

C++ User Population

Time	Estimated Number of Users
Oct 1979	1
Oct 1980	16
Oct 1981	38
Oct 1982	85
Oct 1983	??+2 (no Cpre count)
Oct 1984	??+50 (no Cpre count)
Oct 1985	500
Oct 1986	2,000
Oct 1987	4,000
Oct 1988	15,000
Oct 1989	50,000
Oct 1990	150,000
Oct 1991	400,000
Oct 2004	over 3,270,000

- above numbers are conservative
- 1979 to 1991: C++ user population doubled approximately every 7.5 months
- stable growth thereafter

Success of C++

- C++ very successful programming language
- not luck or solely because based on C
- efficient, provides low-level access to hardware, but also supports abstraction
- non-proprietary: in 1989, all rights to language transferred to standards bodies (first ANSI and later ISO) from AT&T
- multi-paradigm language, supporting procedural, object-oriented, generic, and functional (e.g., lambda functions) programming
- does not force particular programming style
- reasonably portable
- has continued to evolve, incorporating new ideas (e.g., templates, exceptions, STL)
- stable: high degree of compatibility with earlier versions of language
- very strong bias towards providing general-purpose facilities rather than more application-specific ones

Application Areas

- banking and financial (funds transfer, financial modelling, teller machines)
- classical systems programming (compilers, operating systems, device drivers, network layers, editors, database systems)
- small business applications (inventory systems)
- desktop publishing (document viewers/editors, image editing)
- embedded systems (cameras, cell phones, airplanes, medical systems, appliances)
- entertainment (games)
- GUI
- hardware design and verification
- scientific and numeric computation (physics, engineering, simulations, data analysis, geometry processing)
- servers (web servers, billing systems)
- telecommunication systems (phones, networking, monitoring, billing, operations systems)

Section 2.1.1

References

- B. Stroustrup. *A history of C++: 1979–1991*.
In *Proc. of ACM History of Programming Languages Conference*, pages 271–298, Mar. 1993
- B. Stroustrup. *The Design and Evolution of C++*.
Addison Wesley, Mar. 1994.
- B. Stroustrup. *Evolving a language in and for the real world: C++ 1991–2006*.
In *Proc. of the ACM SIGPLAN Conference on History of Programming Languages*, pages 4–1–4–59, 2007.
- Cfront software available from Computer History Museum’s Software Preservation Group <http://www.softwarepreservation.org>.
(See http://www.softwarepreservation.org/projects/c_plus_plus/cfront).
- ISO JTC1/SC22/WG21 web site. <http://www.open-std.org/jtc1/sc22/wg21/>.

- ISO/IEC 14882:1998 — programming languages — C++, Sept. 1998.
- ISO/IEC 14882:2003 — programming languages — C++, Oct. 2003.
- ISO/IEC TR 18015:2006 — information technology — programming languages, their environments and system software interfaces — technical report on C++ performance, Sept. 2006.
- ISO/IEC TR 19768:2007 — information technology — programming languages — technical report on C++ library extensions, Nov. 2007.
- ISO/IEC 14882:2011 — information technology — programming languages — C++, Sept. 2011.
- ISO/IEC 14882:2014 — information technology — programming languages — C++, 2014.
- ISO JTC1/SC22/WG21 web site. <http://www.open-std.org/jtc1/sc22/wg21/>.

Section 2.2

Getting Started

Section 2.2.1

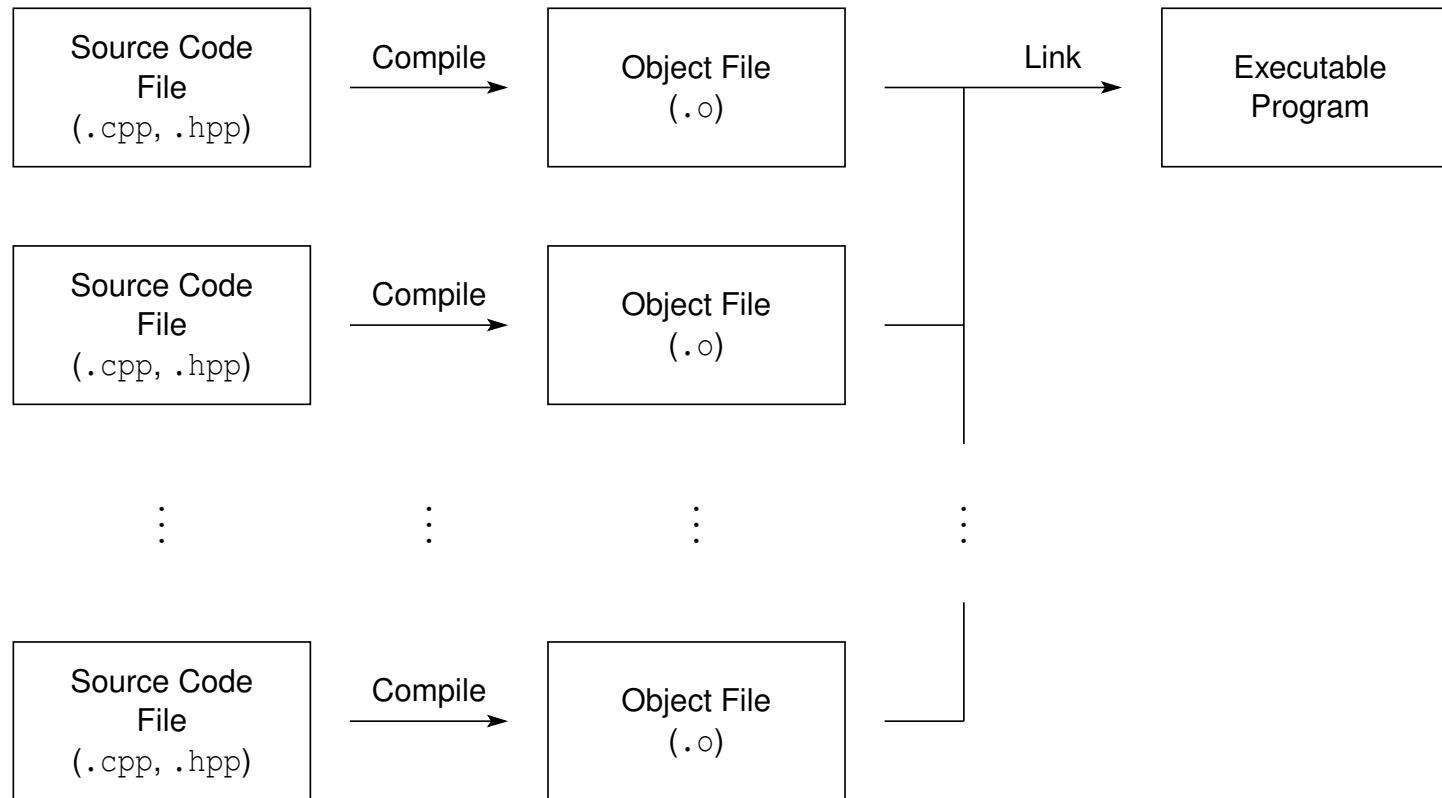
Building Programs: Compiling and Linking

hello Program: hello.cpp

```
1  #include <iostream>
2
3  int main ()
4  {
5      std::cout << "Hello, world!\n";
6  }
```

- program prints message “Hello, world!” and then exits
- starting point for execution of C++ program is function called `main`; every C++ program must define function called `main`
- `#include` preprocessor directive to include complete contents of file
- `iostream` standard header file that defines various types and variables related to I/O
- `std::cout` is standard output stream (defaults to user’s terminal)
- operator `<<` is used for output

Software Build Process



- start with C++ source code files (.cpp, .hpp)
- compile: convert source code to object code
- object code stored in object file (.o)
- link: combine contents of one or more object files (and possibly some libraries) to produce executable program
- executable program can then be run directly

- `g++` command provides both compiling and linking functionality
- command-line usage:
`g++ [options] input_file ...`
- many command-line options are supported
- some particularly useful command-line options listed on next slide
- compile C++ source file *file.cpp* to produce object code file *file.o*:
`g++ -c file.cpp`
- link object files *file_1.o*, *file_2.o*, ... to produce executable file *executable*:
`g++ -o executable file_1.o file_2.o ...`
- web page: <http://www.gnu.org/software/gcc>

Common g++ Command-Line Options

Option	Description
-c	compile only (i.e., do not link)
-o <i>file</i>	use file <i>file</i> for output
-g	include debugging information
-On	set optimization level to <i>n</i> (0 almost none; 3 full)
-std=c++14	conform to C++14 standard
-pthread	enable concurrency support (via pthreads library)
-I <i>dir</i>	specify additional directory <i>dir</i> to search for include files
-L <i>dir</i>	specify additional directory <i>dir</i> to search for libraries
-l <i>lib</i>	link with library <i>lib</i>
-pedantic-errors	strictly enforce compliance with standard
-Wall	enable most warning messages
-Wextra	enable some extra warning messages not enabled by -Wall
-Wpedantic	warn about deviations from strict standard compliance
-Werror	treat all warnings as errors

Manually Building `hello` Program

- numerous ways in which `hello` program could be built
- often advantageous to compile each source file separately
- can compile and link as follows:
 - 1 compile source code file `hello.cpp` to produce object file `hello.o`:

```
g++ -c hello.cpp
```
 - 2 link object file `hello.o` to produce executable program `hello`:

```
g++ -o hello hello.o
```
- generally, manual building of program is quite tedious, especially when program consists of multiple source files and additional compiler options need to be specified
- in practice, we use tools to automate build process (e.g., `make` utility)

Section 2.2.2

Make Utility

Make

- `make` command
- controls generation of executables and/or other non-source files from program's source files
- extremely popular tool for automating build process
- available on many platforms (e.g., Unix, Microsoft Windows, Mac OS X); used extensively on Unix systems
- very flexible
- can handle building multiple programs consisting of hundreds of source files or single program consisting of only one source file
- can be used to build almost anything (i.e., need not be a program)
- for example, all materials for this course typeset using \LaTeX (e.g., coursepack, slides, handouts, exams), and `make` utility used to compile \LaTeX source code into PDF documents
- one of most popular implementations of `make` is GNU Make
- GNU Make web page: <http://www.gnu.org/software/make>

make Command

- target is something that can be built, typically (but not necessarily) file such as executable file or object file
- `make` command driven by data file called `makefile`
- `makefile` usually named `Makefile` or `makefile`
- command-line usage:
 `make [options] [targets]`
- *targets*: zero or more targets to be built
- *options*: zero or more options
- by default, looks for `makefile` called `makefile` and then `Makefile`
- if no targets are specified, will build first target specified in `makefile`
- only builds files that are out of date
- most common command-line options include:
 - n show commands that would be executed but do not actually execute them
 - f *makefile* use `makefile` *makefile*

Makefiles

- comment starts at hash character (i.e., “#”) and continues until end of line; example:

```
# This comment continues until the end of the line.
```

- supports variables
- some important variables used by built-in rules:

Name	Description
CXX	C++ compiler command
CXXFLAGS	C++ compiler options
LDFLAGS	linker options

- to assign value to variable, use equal sign; example:

```
CXX = g++
```

- to substitute value of variable, use dollar sign followed by variable name in parentheses; example:

```
$(CXX)
```

Makefiles (Continued 1)

- makefile specifies targets and rules for building targets
- each rule in makefile has following form:

```
targets : prerequisites  
_____ commands  
_____ ...
```

- indentation shown above must be with tab character and not spaces
- *targets*: list of one or more targets
- *prerequisites*: files on which targets depend (i.e., files used to produce targets)
- *commands*: actions that must be carried out to produce target from its prerequisites

Makefiles (Continued 2)

- normally, each target associated with file of same name (and building target will create this file)
- phony target: target that is not associated with any file
- to identify target as phony make it prerequisite of special target called “**.PHONY**”; example (specify `all` as phony target):

```
.PHONY: all
```

- some special built-in variables that can be used in rules:

Name	Description
<code>\$@</code>	target
<code>\$<</code>	name of first prerequisite
<code>\$^</code>	names of all of prerequisites separated by spaces

Makefile for hello Program

```
1  CXX = g++           # The C++ compiler command.
2  CXXFLAGS = -g -O    # The C++ compiler options.
3  LDFLAGS =           # The linker options (if any).
4
5  # The all target builds all of the programs handled by
6  # the makefile.
7  # This target has the dependency chain:
8  #     all -> hello -> hello.o -> hello.cpp
9  all: hello
10
11 # The clean target removes all of the executable files
12 # and object files produced by the build process.
13 clean:
14  _____rm -f hello *.o
15
16 # The hello target builds the hello executable.
17 hello: hello.o
18  _____$(CXX) $(CXXFLAGS) -o $@ $^ $(LDFLAGS)
19
20 # Indicate that the all and clean targets do not
21 # correspond to actual files.
22 .PHONY: all clean
23
24 # The following rule is effectively built into make and
25 # therefore need not be explicitly specified:
26 # hello.o: hello.cpp
27 # _____$(CXX) $(CXXFLAGS) -c $<
```

Commentary on Makefile for `hello` Program

- `all` target: builds all of the programs handled by the makefile (e.g., `hello`)
- `clean` target: removes all of the executable files and object files produced by build process (e.g., `hello`, `hello.o`)
- although `all` and `clean` have no special meaning to make, very common practice to provide targets with these particular names in all makefiles
- `hello` target: compiles and links the hello program
- chain of dependencies for `all` target:
 $all \rightarrow hello \rightarrow hello.o \rightarrow hello.cpp$
- `all` and `clean` examples of phony targets

Section 2.2.3

Debugging Tools

Source-Level Debuggers

- unfortunately, software does not always work as intended due to errors in code (i.e., bugs)
- how does one go about fixing bugs in time-efficient manner?
- source-level debugger is essential tool
- single stepping: step through execution of code, one source-code line at a time
- breakpoints: pause execution at particular points in code
- watchpoints: pause execution when the value of variable is changed
- print values of variables

GNU Debugger (GDB)

- GNU Debugger (GDB) is powerful source-level debugger
- home page: <http://www.gnu.org/software/gdb>
- available on most platforms (e.g., Unix, Microsoft Windows)
- most popular source-level debugger on Unix systems
- allows one to see what is happening inside program as it executes or what a program was doing at the moment it crashed
- has all of the standard functionality of a source-level debugger (e.g., breakpoints, watchpoints, single-stepping)
- gdb command
- command-line usage:

`gdb [options] executable`

gdb Commands

help

Print help information.

quit

Exit debugger.

run [*arglist*]

Start the program (with *arglist* if specified).

print *expr*

Display the value of the expression *expr*.

bt

Display a stack backtrace.

list

Type the source code lines in the vicinity of where the program is currently stopped.

`break function`

Set a breakpoint at the function *function*.

`watch expr`

Set a watchpoint for the expression *expr*.

`c`

Continue running the program (e.g., after stopping at a breakpoint).

`next`

Execute the next program line, stepping over any function calls in the line.

`step`

Execute the next program line, stepping into any function calls in the line.

GNU Data Display Debugger (DDD)

- graphical front-end to command-line debuggers such as GDB
- has some fancy graphical data display functionality
- all `gdb` commands available in text window, but can use graphical interface to enter commands as well
- home page: <http://www.gnu.org/software/ddd>
- `ddd` command

- can detect many memory management and threading bugs
- can profile programs in detail
- home page: <http://www.valgrind.org>
- `valgrind` command
- `valkyrie` command (GUI for Memcheck and Helgrind tools in Valgrind)

Valgrind References I

- 1 P. Floyd. [Valgrind part 1 — introduction.](#)
Overload, 108:14–15, Apr. 2012.
- 2 P. Floyd. [Valgrind part 2 — basic memcheck.](#)
Overload, 109:24–28, June 2012.
- 3 P. Floyd. [Valgrind part 3 — advanced memcheck.](#)
Overload, 110:4–7, Aug. 2012.
- 4 P. Floyd. [Valgrind part 4 — cachegrind and callgrind.](#)
Overload, 111:4–7, Oct. 2012.
- 5 P. Floyd. [Valgrind part 5 — massif.](#)
Overload, 112:20–24, Dec. 2012.

Section 2.3

C++ Basics

The C++ Programming Language

- created by Bjarne Stroustrup of Bell Labs
- originally known as C with Classes; renamed as C++ in 1983
- most recent specification of language in ISO/IEC 14882:2014 (informally known as “C++14”)
- next version of standard expected in 2017
- procedural
- loosely speaking is superset of C
- directly supports object-oriented and generic programming
- maintains efficiency of C
- application domains: systems software, application software, device drivers, embedded software, high-performance server and client applications, entertainment software such as video games, native code for Android applications
- greatly influenced development of C# and Java

Comments

- two styles of comments provided
- comment starts with `//` and proceeds to end of line
- comment starts with `/*` and proceeds to first `*/`

```
// This is an example of a comment.  
/* This is another example of a comment. */  
/* This is an example of a comment that  
   spans  
   multiple lines. */
```

- comments of `/* ... */` style *do not nest*

```
/*  
/* This sentence is part of a comment. */  
This sentence is not part of any comment and  
will probably cause a compile error.  
*/
```

Identifiers

- identifiers used to name entities such as: types, objects (i.e., variables), and functions
- valid identifier is sequence of one or more letters, digits, and underscore characters that does not begin with a digit
- identifiers that begin with underscore (in many cases) or contain double underscores are reserved for use by C++ implementation and should be avoided
- examples of valid identifiers:
 - `event_counter`
 - `eventCounter`
 - `sqrt_2`
 - `f_o_o_b_a_r_4_2`
- identifiers are case sensitive (e.g., `counter` and `cOuNtEr` are distinct identifiers)
- identifiers cannot be any of reserved keywords (see next slide)
- **scope** of identifier is context in which identifier is valid (e.g., block, function, global)

Reserved Keywords

<code>alignas</code>	<code>default</code>	<code>noexcept</code>	<code>this</code>
<code>alignof</code>	<code>delete</code>	<code>not</code>	<code>thread_local</code>
<code>and</code>	<code>do</code>	<code>not_eq</code>	<code>throw</code>
<code>and_eq</code>	<code>double</code>	<code>nullptr</code>	<code>true</code>
<code>asm</code>	<code>dynamic_cast</code>	<code>operator</code>	<code>try</code>
<code>auto</code>	<code>else</code>	<code>or</code>	<code>typedef</code>
<code>bitand</code>	<code>enum</code>	<code>or_eq</code>	<code>typeid</code>
<code>bitor</code>	<code>explicit</code>	<code>private</code>	<code>typename</code>
<code>bool</code>	<code>export</code>	<code>protected</code>	<code>union</code>
<code>break</code>	<code>extern</code>	<code>public</code>	<code>unsigned</code>
<code>case</code>	<code>false</code>	<code>register</code>	<code>using</code>
<code>catch</code>	<code>float</code>	<code>reinterpret_cast</code>	<code>virtual</code>
<code>char</code>	<code>for</code>	<code>return</code>	<code>void</code>
<code>char16_t</code>	<code>friend</code>	<code>short</code>	<code>volatile</code>
<code>char32_t</code>	<code>goto</code>	<code>signed</code>	<code>wchar_t</code>
<code>class</code>	<code>if</code>	<code>sizeof</code>	<code>while</code>
<code>compl</code>	<code>inline</code>	<code>static</code>	<code>xor</code>
<code>const</code>	<code>int</code>	<code>static_assert</code>	<code>xor_eq</code>
<code>constexpr</code>	<code>long</code>	<code>static_cast</code>	<code>override*</code>
<code>const_cast</code>	<code>mutable</code>	<code>struct</code>	<code>final*</code>
<code>continue</code>	<code>namespace</code>	<code>switch</code>	
<code>decltype</code>	<code>new</code>	<code>template</code>	

*Note: context sensitive

Section 2.3.1

Objects, Types, and Values

Fundamental Types

- boolean type: **bool**
- character types:
 - **char** (may be signed or unsigned)
 - **signed char**
 - **unsigned char**
 - **char16_t**
 - **char32_t**
 - **wchar_t**
- **char** is distinct type from **signed char** and **unsigned char**
- standard signed integer types:
 - **signed char**
 - **signed short int**
 - **signed int**
 - **signed long int**
 - **signed long long int**
- standard unsigned integer types:
 - **unsigned char**
 - **unsigned short int**
 - **unsigned int**
 - **unsigned long int**
 - **unsigned long long int**

Fundamental Types (Continued)

- “**int**” may be omitted from names of (non-character) integer types (e.g., “**unsigned**” equivalent to “**unsigned int**” and “**signed**” equivalent to “**signed int**”)
- “**signed**” may be omitted from names of signed integer types, excluding **signed char** (e.g., “**int**” equivalent to “**signed int**”)
- boolean, character, and (signed and unsigned) integer types collectively called **integral types**
- floating-point types:
 - **float**
 - **double**
 - **long double**
- void (i.e., incomplete/valueless) type: **void**
- null pointer type: `std::nullptr_t` (defined in header file `cstddef`)

- **literal** (a.k.a. literal constant) is value written exactly as it is meant to be interpreted
- examples of literals:

```
"Hello, world"
```

```
"Bjarne"
```

```
'a'
```

```
'A'
```

```
123
```

```
123U
```

```
1'000'000'000
```

```
3.1415
```

```
1.0L
```

```
1.23456789e-10
```

Character Literals

- character literal consists of optional prefix followed by one or more characters enclosed in single quotes
- type of character literal determined by prefix (or lack thereof) as follows:

Prefix	Literal	Type
None	ordinary	normally char (in special cases int)
u ⁸ [since C++17]	UTF-8	char
u	UCS-2	char16_t
U	UCS-4	char32_t
L	wide	wchar_t

- special characters can be represented by escape sequence:

Character	Escape Sequence
newline (LF)	<code>\n</code>
horizontal tab (HT)	<code>\t</code>
vertical tab (VT)	<code>\v</code>
backspace (BS)	<code>\b</code>
carriage return (CR)	<code>\r</code>
form feed (FF)	<code>\f</code>
alert (BEL)	<code>\a</code>

Character	Escape Sequence
backslash (\)	<code>\\</code>
question mark (?)	<code>\?</code>
single quote (')	<code>\'</code>
double quote (")	<code>\"</code>
octal number ooo	<code>\ooo</code>
hex number hhh	<code>\xhhh</code>

- examples of character literals:

`'a'` `'1'` `'!'` `'\n'` `u'a'` `U'a'` `L'a'` `u8'a'`



Character Literals (Continued)

- decimal digit characters guaranteed to be consecutive in value (e.g., '1' must equal '0' + 1)
- in case of ordinary character literals, alphabetic characters are *not* guaranteed to be consecutive in value (e.g., 'b' is not necessarily 'a' + 1)

String Literals

- string literal consists of optional prefix followed by zero or more characters enclosed in double quotes
- string literal has character array type
- type of string literal determined by prefix (or lack thereof) as follows:

Prefix	Literal	Type
None	narrow	<code>const char []</code>
u8	UTF-8	<code>const char []</code>
u	UTF-16	<code>const char16_t []</code>
U	UTF-32	<code>const char32_t []</code>
L	wide	<code>const wchar_t []</code>

- examples of string literals:

```
"Hello, World!\n"
```

```
"123"
```

```
"ABCDEFGH"
```

- adjacent string literals are concatenated (e.g., "Hel" "lo" equivalent to "Hello")
- string literals implicitly terminated by null character (i.e., '\0')
- so, for example, "Hi" means 'H' followed by 'i' followed by '\0'

Integer Literals

- can be specified in decimal, binary, hexadecimal, and octal
- number base indicated by prefix (or lack thereof) as follows:

Prefix	Number Base
None	decimal
Leading 0	octal
0b or 0B	binary
0x or 0X	hexadecimal

- various suffixes can be specified to control type of literal:

- u or U
- l or L
- both u or U and l or L
- ll or LL
- both u or U and ll or LL

- can use single quote as digit separator (e.g., 1'000'000)

- examples of integer literals:

42

1'000'000'000'000ULL

0xdeadU

- integer literal always nonnegative; so, for example, -1 is integer literal 1 with negation operation applied

Integer Literals (Continued)

Suffix	Decimal Literal	Non-Decimal Literal
None	<code>int</code> <code>long int</code> <code>long long int</code>	<code>int</code> <code>unsigned int</code> <code>long int</code> <code>unsigned long int</code> <code>long long int</code> <code>unsigned long long int</code>
<code>u</code> or <code>U</code>	<code>unsigned int</code> <code>unsigned long int</code> <code>unsigned long long int</code>	<code>unsigned int</code> <code>unsigned long int</code> <code>unsigned long long int</code>
<code>l</code> or <code>L</code>	<code>long int</code> <code>long long int</code>	<code>long int</code> <code>unsigned long int</code> <code>long long int</code> <code>unsigned long long int</code>
Both <code>u</code> or <code>U</code> and <code>l</code> or <code>L</code>	<code>unsigned long int</code> <code>unsigned long long int</code>	<code>unsigned long int</code> <code>unsigned long long int</code>
<code>ll</code> or <code>LL</code>	<code>long long int</code>	<code>long long int</code> <code>unsigned long long int</code>
Both <code>u</code> or <code>U</code> and <code>ll</code> or <code>LL</code>	<code>unsigned long long int</code>	<code>unsigned long long int</code>

Floating-Point Literals

- type of literal indicated by suffix (or lack thereof) as follows:

Suffix	Type
None	double
f or F	float
l or L	long double

- examples of **double** literals:

1.414

1.25e-8

- examples of **float** literals:

1.414f

1.25e-8f

- examples of **long double** literals:

1.5L

1.25e-20L

- floating-point literals always nonnegative; so, for example, `-1.0` is literal `1.0` with negation operator applied

Boolean and Pointer Literals

- boolean literals:

true

false

- pointer literal:

nullptr

Declarations and Definitions

- **declaration** introduces identifier for type, object (i.e., variable), or function (without necessarily providing full information about identifier)
 - in case of object, specifies type (of object)
 - in case of function, specifies number of parameters, type of each parameter, and type of return value (if not automatically deduced)
- each identifier must be declared before it can be used (i.e., referenced)
- **definition** provides full information about identifier and causes entity associated with identifier (if any) to be created
 - in case of type, provides full details about type
 - in case of object, causes storage to be allocated for object and object to be created
 - in case of function, provides code for function body
- in case of objects, in most (but not all) contexts, declaring object also defines it
- can declare identifier multiple times but can define only once
- above terminology often abused, with “declaration” and “definition” being used interchangeably

Examples of Declarations and Definitions

```
int count; // declare and define count
extern double alpha; // (only) declare alpha

void func() { // declare and define func
    int n; // declare and define n
    double x = 1.0; // declare and define x
    // ...
}

bool isOdd(int); // declare isOdd
bool isOdd(int x); // declare isOdd (x ignored)

bool isOdd(int x) { // declare and define isOdd
    return x % 2;
}

struct Thing; // declare Thing

struct Vector2 { // declare and define Vector2
    double x;
    double y;
};
```

Variable Declarations and Definitions

- **variable declaration** (a.k.a. object declaration) introduces identifier that names object and specifies type of object
- **variable definition** (a.k.a. object definition) provides all information included in variable declaration and also causes object to be created (e.g., storage allocated for object)
- example:

```
int count;  
    // declare and define count  
double alpha;  
    // declare and define alpha  
extern double gamma;  
    // declare (but do not define) gamma
```

Arrays

- **array** is collection of one or more objects of *same* type that are stored *contiguously* in memory
- each element in array identified by (unique) integer index, with indices starting from *zero*
- array denoted by []
- example:

```
double x[10]; // array of 10 doubles  
int data[512][512]; // 512 by 512 array of ints
```

- elements of array accessed using subscripting operator []
- example:

```
int x[10];  
// elements of arrays are x[0], x[1], ..., x[9]
```

- in C++ rarely ever need to use arrays
- use `std::array` or `std::vector` type instead (as this has many practical advantages over array)
- will revisit `std::array` and `std::vector` types later

Array Example

- code:

```
int a[4] = {1, 2, 3, 4};
```

- assumptions (for some completely fictitious C++ language implementation):
 - **sizeof(int)** is 4
 - array a starts at address 1000
- memory layout:

Address		Name
1000	1	a[0]
1004	2	a[1]
1008	3	a[2]
1012	4	a[3]

Pointers

- **pointer** is object whose value is address in memory where another object is stored
- pointer to object of type T denoted by T^*
- **null pointer** is special pointer value that does not refer to any valid memory location
- null pointer value provided by **nullptr** keyword
- accessing object to which pointer refers called **dereferencing**
- dereferencing pointer performed by *indirection operator* (i.e., “*”)
- if p is pointer, $*p$ is object to which pointer refers
- if x is object of type T , $\&x$ is address of object (which has type T^*)
- example:

```
char c;  
char* cp = nullptr; // cp is pointer to char  
char* cp2 = &c; // cp2 is pointer to char
```


Pointer Example

- code:

```
int i = 42;  
int* p = &i;  
assert(*p == 42);
```

- assumptions (for some completely fictitious C++ language implementation):

- **sizeof(int)** is 4
- **sizeof(int*)** is 4
- **&i** is **((int*)1000)**
- **&p** is **((int*)1004)**

- memory layout:

Address		Name
1000	42	i
1004	1000	p

References

- **reference** is alias (i.e., nickname) for *already existing* object
- two kinds of references:
 - ① lvalue reference
 - ② rvalue reference
- lvalue reference to object of type T denoted by $T\&$
- rvalue reference to object of type T denoted by $T\&\&$
- initializing reference called **reference binding**
- lvalue and rvalue references differ in their binding properties (i.e., to what kinds of objects reference can be bound)
- in most contexts, lvalue references usually needed
- rvalue references used in context of move constructors and move assignment operators (to be discussed later)
- example:

```
int x;  
int& y = x; // y is lvalue reference to int  
int&& tmp = 3; // tmp is rvalue reference to int
```

References Example

- code:

```
int i = 42;  
int& j = i;  
assert(j == 42);
```

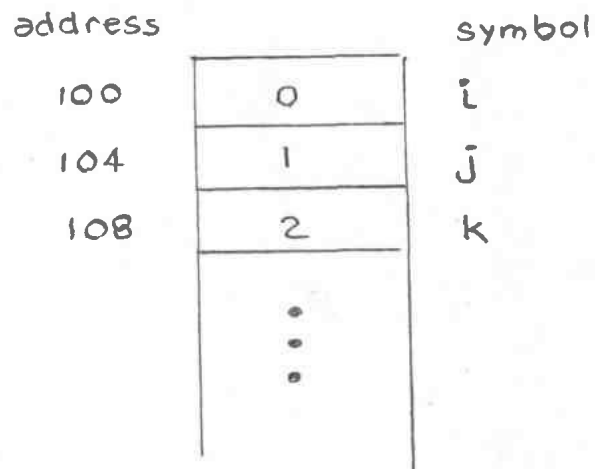
- assumptions (for some completely fictitious C++ language implementation):
 - **sizeof(int)** is 4
 - **&i** is **((int*)1000)**
- memory layout:

Address		Name
1000	<div style="border: 1px solid black; display: inline-block; padding: 5px 20px;">42</div>	i, j

Addresses, Pointers, and References

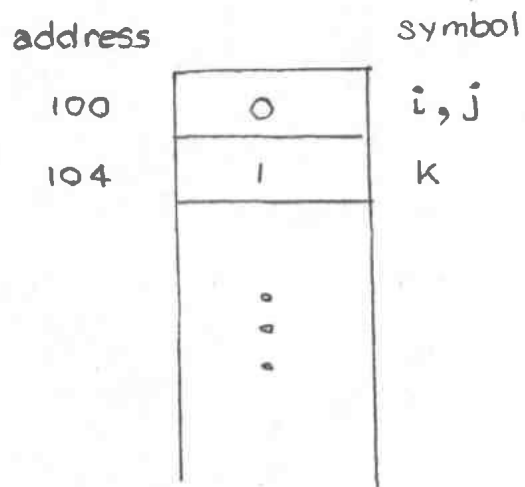
Assume `sizeof(int)` is 4, `sizeof(int*)` is 4

```
int i = 0;  
int j = 1;  
int k = 2;
```



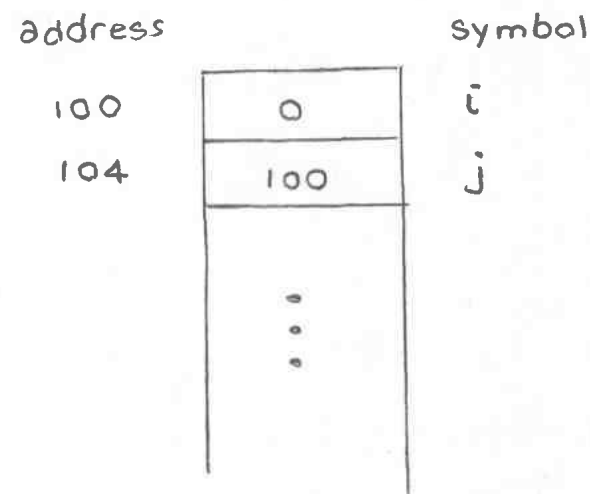
`i = 3;` vs. `j = 3;` ?

```
int i = 0;  
int &j = i;  
int k = 1;
```



`i = 3;` vs. `j = 3;` ?

```
int i = 0;  
int* j = &i;
```



`i = 3;` vs. `*j = 3;` ?

References Versus Pointers

- references and pointers similar in that both can be used to refer to some other entity (e.g., object or function)
- two key differences between references and pointers:
 - ① reference must refer to something, while pointer can have null value (**`nullptr`**)
 - ② references cannot be rebound, while pointers can be changed to point to different entity
- references have cleaner syntax than pointers, since pointers must be dereferenced upon each use (and dereference operations tend to clutter code)
- use of pointers often implies need for memory management (i.e., memory allocation, deallocation, etc.), and memory management can introduce numerous kinds of bugs when done incorrectly
- often faced with decision of using pointer or reference in code
- generally advisable to prefer use of references over use of pointers unless compelling reason to do otherwise, such as:
 - must be able to handle case of referring to nothing
 - must be able to change entity being referred to

Unscoped Enumerations

- **enumerated type** provides way to describe range of values that are represented by named constants called **enumerators**
- object of enumerated type can take any one of enumerators as value
- enumerator values represented by some *integral type*
- enumerator can be assigned specific value (which may be negative)
- if enumerator not assigned specific value, value defaults to zero if first enumerator in enumeration and one greater than value for previous enumerator otherwise
- example:

```
enum Suit {  
    Clubs, Diamonds, Hearts, Spades  
};  
  
Suit suit = Clubs;
```

- example:

```
enum Suit {  
    Clubs = 1, Diamonds = 2, Hearts = 4, Spades = 8  
};
```

Scoped Enumerations

- scoped enumeration similar to unscoped enumeration, except
 - all enumerators are placed in scope of enumeration itself
 - integral type to used to hold enumerator values can be explicitly specified
 - conversions involving scoped enumerations are stricter (i.e., more type safe)
- **class** or **struct** added after **enum** keyword to make enumeration scoped
- scope resolution operator (i.e., “::”) used to access enumerators
- scoped enumerations should probably be preferred to unscoped ones
- example:

```
enum struct Season {
    spring, summer, fall, winter
};

enum struct Suit : unsigned char {
    clubs, diamonds, hearts, spades
};

Season season = Season::summer;
Suit suit = Suit::spades;
```

Type Aliases with `typedef` Keyword

- `typedef` keyword used to create alias for existing type
- example:

```
typedef long long BigInt;  
BigInt i; // i has type long long
```

```
typedef char* CharPtr;  
CharPtr p; // p has type char*
```


Type Aliases with `using` Statement

- `using` statement can be used to create alias for existing type
- probably preferable to use `using` statement over `typedef`
- example:

```
using BigInt = long long;  
BigInt i; // i has type long long
```

```
using CharPtr = char*;  
CharPtr p; // p has type char*
```

The `extern` Keyword

- **translation unit**: basic unit of compilation in C++ (i.e., single source code file plus all of its directly and indirectly included header files)
- **extern** keyword used to declare object/function in separate translation unit
- example:

```
extern int evil_global_variable;  
    // declaration only  
    // actual definition in another file
```

The `const` Qualifier

- `const` qualifier specifies that object has value that is *constant* (i.e., cannot be changed)

- following defines `x` as `int` with value 42 that cannot be modified:

```
const int x = 42;
```

- example:

```
const int x = 42;  
x = 13; // ERROR: x is const  
const int& x1 = x; // OK  
const int* p1 = &x; // OK  
int& x2 = x; // ERROR: x const, x2 not const  
int* p2 = &x; // ERROR: x const, *p2 not const
```

- example:

```
int x = 0;  
const int& y = x;  
x = 42; // OK  
// y also changed to 42 since y refers to x  
// y cannot be used to change x, however  
// i.e., the following would cause compile error:  
// y = 24; // ERROR: y is const
```

The `volatile` Qualifier

- **`volatile`** qualifier used to indicate that object can change due to agent *external to program* (e.g., memory-mapped device, signal handler)
- compiler cannot optimize away read and write operations on **`volatile`** objects (e.g., repeated reads without intervening writes cannot be optimized away)
- **`volatile`** qualifier typically used when object:
 - corresponds to register of memory-mapped device
 - may be modified by signal handler (namely, object of type **`volatile std::sig_atomic_t`**)
- example:

```
volatile int x;  
volatile unsigned char* deviceStatus;
```

The `auto` Keyword

- in various contexts, `auto` keyword can be used as place holder for type
- in such contexts, implication is that compiler must deduce type
- example:

```
auto i = 3; // i has type int
auto j = i; // j has type int
auto& k = i; // k has type int&
const auto& n = i; // n has type const int&
auto x = 3.14; // x has type double
```

- very useful in generic programming (covered later) when types not always easy to determine
- can potentially save typing long type names
- can lead to more readable code (if well used)
- if overused, can lead to bugs (sometimes very subtle ones) and difficult to read code

Section 2.3.2

Operators and Expressions

Arithmetic Operators

Operator Name	Syntax
addition	$a + b$
subtraction	$a - b$
unary plus	$+a$
unary minus	$-a$
multiplication	$a * b$
division	a / b
modulo (i.e., remainder)	$a \% b$
pre-increment	$++a$
post-increment	$a++$
pre-decrement	$--a$
post-decrement	$a--$

Bitwise Operators

Operator Name	Syntax
bitwise NOT	$\sim a$
bitwise AND	$a \& b$
bitwise OR	$a b$
bitwise XOR	$a \wedge b$
arithmetic left shift	$a \ll b$
arithmetic right shift	$a \gg b$

Assignment and Compound-Assignment Operators

Operator Name	Syntax
assignment	<code>a = b</code>
addition assignment	<code>a += b</code>
subtraction assignment	<code>a -= b</code>
multiplication assignment	<code>a *= b</code>
division assignment	<code>a /= b</code>
modulo assignment	<code>a %= b</code>
bitwise AND assignment	<code>a &= b</code>
bitwise OR assignment	<code>a = b</code>
bitwise XOR assignment	<code>a ^= b</code>
arithmetic left shift assignment	<code>a <<= b</code>
arithmetic right shift assignment	<code>a >>= b</code>

Operators (Continued 2)

Logical/Relational Operators

Operator Name	Syntax
equal	<code>a == b</code>
not equal	<code>a != b</code>
greater than	<code>a > b</code>
less than	<code>a < b</code>
greater than or equal	<code>a >= b</code>
less than or equal	<code>a <= b</code>
logical negation	<code>!a</code>
logical AND	<code>a && b</code>
logical OR	<code>a b</code>

Member and Pointer Operators

Operator Name	Syntax
array subscript	<code>a[b]</code>
indirection	<code>*a</code>
address of	<code>&a</code>
member selection	<code>a.b</code>
member selection	<code>a->b</code>
member selection	<code>a.*b</code>
member selection	<code>a->*b</code>

Operators (Continued 3)

Other Operators

Operator Name	Syntax
function call	<code>a (...)</code>
comma	<code>a, b</code>
ternary conditional	<code>a ? b : c</code>
scope resolution	<code>a :: b</code>
sizeof	sizeof (a)
parameter-pack sizeof	sizeof... (a)
alignof	alignof (T)
allocate storage	new T
allocate storage (array)	new T[a]
deallocate storage	delete a
deallocate storage (array)	delete [] a

Other Operators (Continued)

Operator Name	Syntax
type ID	typeid (a)
type cast	(T) a
const cast	const_cast <T> (a)
static cast	static_cast <T> (a)
dynamic cast	dynamic_cast <T> (a)
reinterpret cast	reinterpret_cast <T> (a)
throw	throw a
noexcept	noexcept (e)

Operator Precedence

Precedence	Operator	Name	Associativity
1	::	scope resolution	none
2	. -> [] () ++ --	member selection (object) member selection (pointer) subscripting function call postfix increment postfix decrement	left to right

Operator Precedence (Continued 1)

Precedence	Operator	Name	Associativity
3	sizeof	size of object/type	right to left
	++	prefix increment	
	--	prefix decrement	
	~	bitwise NOT	
	!	logical NOT	
	-	unary minus	
	+	unary plus	
	&	address of	
	*	indirection	
	new	allocate storage	
	new []	allocate storage (array)	
	delete	deallocate storage	
	delete []	deallocate storage (array)	
	()	cast	

Operator Precedence (Continued 2)

Precedence	Operator	Name	Associativity
4	. [*] -> [*]	member selection (objects) member selection (pointers)	left to right
5	* / %	multiplication division modulus	left to right
6	+ -	addition subtraction	left to right
7	<< >>	left shift right shift	left to right
8	< <= > >=	less than less than or equal greater than greater than or equal	left to right
9	== !=	equality inequality	left to right

Operator Precedence (Continued 3)

Precedence	Operator	Name	Associativity
10	&	bitwise AND	left to right
11	^	bitwise XOR	left to right
12		bitwise OR	left to right
13	&&	logical AND	left to right
14		logical OR	left to right
15	? :	ternary conditional	right to left

Operator Precedence (Continued 4)

Precedence	Operator	Name	Associativity
16	= *= /= %= += -= <<= >>= &= = ^=	assignment multiplication assignment division assignment modulus assignment addition assignment subtraction assignment left shift assignment right shift assignment bitwise AND assignment bitwise OR assignment bitwise XOR assignment	right to left
17	throw	throw exception	right to left
18	,	comma	left to right

Alternative Tokens

Alternative	Primary
and	&&
bitor	
or	
xor	^
compl	~
bitand	&
and_eq	&=
or_eq	=
xor_eq	^=
not	!
not_eq	!=

- alternative tokens above probably best avoided as they lead to more verbose code

Expressions

- An **expression** is a sequence of operators and operands that specifies a computation.
- An expression has a type and, if the type is not **void**, a value.
- A **constant expression** is an expression that can be evaluated at compile time (e.g., $1 + 1$).
- Example:

```
int x = 0;  
int y = 0;  
int* p = &x;  
double d = 0.0;  
// Evaluate some  
// expressions here.
```

Expression	Type	Value
x	int	0
y = x	int &	reference to y
x + 1	int	1
x * x + 2 * x	int	0
y = x * x	int &	reference to y
x == 42	bool	false
*p	int &	reference to x
p == &x	bool	true
x > 2 * y	bool	false
std::sin(d)	double	0.0

Operator Precedence/Associativity Example

Expression	Fully-Parentthesized Expression
<code>a + b + c</code>	<code>((a + b) + c)</code>
<code>a = b = c</code>	<code>(a = (b = c))</code>
<code>c = a + b</code>	<code>(c = (a + b))</code>
<code>d = a && !b c</code>	<code>(d = ((a && (!b)) c))</code>
<code>++*p++</code>	<code>((++(*p++)))</code>
<code>a ~b & c ^ d</code>	<code>(a (((~b) & c) ^ d))</code>
<code>a[0]++ + a[1]++</code>	<code>((a[0]++) + (a[1]++))</code>
<code>a + b * c / d % -g</code>	<code>(a + (((b * c) / d) % (-g)))</code>
<code>++p[i]</code>	<code>((++(p[i])))</code>
<code>--*++p</code>	<code>((--(*++p)))</code>
<code>a += b += c += d</code>	<code>(a += (b += (c += d)))</code>
<code>z = a == b ? ++c : --d</code>	<code>(z = ((a == b) ? (++c) : (--d)))</code>

Short-Circuit Evaluation

- logical and operator (i.e., &&):
 - groups left-to-right
 - result true if both operands are true, and false otherwise
 - second operand is *not evaluated* if first operand is false
- logical or operator (i.e., ||):
 - groups left-to-right
 - result is true if either operand is true, and false otherwise
 - second operand is *not evaluated* if first operand is true
- example:

```
int x = 0;  
bool b = (x == 0 || ++x == 1);  
// b equals true; x equals 0  
b = (x != 0 && ++x == 1);  
// b equals false; x equals 0
```

- above behavior referred to as short circuit evaluation

The `sizeof` Operator

- **`sizeof`** operator is used to query size of object or object type (i.e., amount of storage required)
- for object type `T`, **`sizeof (T)`** yields size of `T` in bytes (e.g., **`sizeof (int)`**, **`sizeof (int [10])`**)
- for expression `e`, **`sizeof e`** yields size of object required to hold result of `e` in bytes (e.g., **`sizeof (&x)`** where `x` is some object)
- **`sizeof (char)`**, **`sizeof (signed char)`**, and **`sizeof (unsigned char)`** guaranteed to be 1
- byte is at least 8 bits (usually exactly 8 bits except on more exotic platforms)

The `alignof` Operator

- object type can have restriction on address at which object of type can start called **alignment requirement**
- for given object type T , starting address for objects of type T must be integer multiple of N bytes, where integer N is called **alignment** of type
- alignment of 1 corresponds to no restriction on alignment (since starting address of object can be any address in memory)
- alignment of 2 restricts starting address of object to be even (i.e., integer multiple of 2)
- for efficiency reasons and due to restrictions imposed by hardware, alignment of particular type may be greater than 1
- **alignof** operator is used to query alignment of type
- for object type T , **alignof** (T) yields alignment used for objects of this type
- **alignof** (`char`), **alignof** (`signed char`), and **alignof** (`unsigned char`) guaranteed to be 1
- fundamental types of size greater than 1 often have alignment greater than 1

The `constexpr` Qualifier for Variables

- `constexpr` qualifier indicates object has value that is *constant expression* (i.e., can be evaluated at compile time)
- `constexpr` implies `const` (but converse not necessarily true)
- following defines `x` as constant expression with type `const int` and value 42:

```
constexpr int x = 42;
```

- example:

```
constexpr int x = 42;  
int y = 1;  
x = 0; // ERROR: x is const  
const int& x1 = x; // OK  
const int* p1 = &x; // OK  
int& x2 = x; // ERROR: x const, x2 not const  
int* p2 = &x; // ERROR: x const, *p2 not const  
int a1[x]; // OK: x is constexpr  
int a2[y]; // ERROR: y is not constexpr
```

The `static_assert` Statement

- `static_assert` allows testing of boolean condition at compile time
- used to test sanity of code or test validity of assumptions made by code
- `static_assert` has two arguments:
 - 1 boolean constant expression (condition to test)
 - 2 string literal for error message to print if boolean expression not true
- as of C++17, second argument is optional
- failed static assertion results in compile error
- example:

```
static_assert(sizeof(int) >= 4, "int is too small");  
static_assert(1 + 1 == 2, "compiler is buggy");
```


Section 2.3.3

Control-Flow Constructs: Selection and Looping

The `if` Statement

- allows conditional execution of code
- syntax has form:

```
if (expression)  
    statement1  
else  
    statement2
```

- if expression *expression* is true, execute statement *statement*₁; otherwise, execute statement *statement*₂
- **else** clause can be omitted leading to simpler form:

```
if (expression)  
    statement1
```

- conditional execution based on more than one condition can be achieved using construct like:

```
if (expression1)  
    statement1  
else if (expression2)  
    statement2  
...  
else  
    statementn
```

The `if` Statement (Continued)

- to include multiple statements in branch of `if`, *must group statements* into single statement using brace brackets

```
if (expression) {  
    statement1,1  
    statement1,2  
    statement1,3  
    ...  
} else {  
    statement2,1  
    statement2,2  
    statement2,3  
    ...  
}
```

- advisable to *always include brace brackets* even when not necessary, as this avoids potential bugs caused by forgetting to include brackets later when more statements added to branch of `if`

The `if` Statement: Example

- example with **else** clause:

```
int x = someValue;
if (x % 2 == 0) {
    std::cout << "x is even\n";
} else {
    std::cout << "x is odd\n";
}
```

- example without **else** clause:

```
int x = someValue;
if (x % 2 == 0) {
    std::cout << "x is divisible by 2\n";
}
```

- example that tests for more than one condition:

```
int x = someValue;
if (x > 0) {
    std::cout << "x is positive\n";
} else if (x < 0) {
    std::cout << "x is negative\n";
} else {
    std::cout << "x is zero\n";
}
```

The `switch` Statement

- allows conditional execution of code based on value of integer expression
- syntax has form:

```
switch (expression) {  
  case const_expr1 :  
    statements1  
  case const_expr2 :  
    statements2  
  ...  
  case const_exprn :  
    statementsn  
  default :  
    statements  
}
```

- *expression* is integer expression; *const_expr*_{*i*} is constant integer expression (e.g., 2, 5+3, 3*5-11)
- if expression *expression* equals *const_expr*_{*i*}, jump to beginning of statements *statements*_{*i*};
if expression *expr* does not equal *const_expr*_{*i*} for any *i*, jump to beginning of statements *statements*;
then, continue executing statements until **break** statement is encountered

The `switch` Statement: Example

```
int x = someValue;
switch (x) {
case 0:
    // Note that there is no break here.
case 1:
    std::cout << "x is 0 or 1\n";
    break;
case 2:
    std::cout << "x is 2\n";
    break;
default:
    std::cout << "x is not 0, 1, or 2\n";
    break;
}
```

The `while` Statement

- looping construct
- syntax has form:

```
while (expression)  
    statement
```

- if expression *expression* is true, statement *statement* is executed; this process repeats until expression *expression* becomes false
- to allow multiple statements to be executed in loop body, ***must group multiple statements*** into single statement with brace brackets

```
while (expression) {  
    statement1  
    statement2  
    statement3  
    ...  
}
```

- advisable to ***always use brace brackets***, even when loop body consists of only one statement

The `while` Statement: Example

```
// print hello 10 times
int n = 10;
while (n > 0) {
    std::cout << "hello\n";
    --n;
}

// loop forever, printing hello
while (true) {
    std::cout << "hello\n";
}
```


The `for` Statement

- looping construct
- has following syntax:

```
for (statement1; expression; statement2)  
    statement3
```
- first, execute statement *statement*₁; then, while expression *expression* is true, execute statement *statement*₃ followed by statement *statement*₂
- *statement*₁ and *statement*₂ may be omitted; *expression* treated as **true** if omitted
- to include multiple statements in loop body, **must group multiple statements** into single statement using brace brackets; advisable to **always use brace brackets**, even when loop body consists of only one statement:

```
for (statement1; expression; statement2) {  
    statement3,1  
    statement3,2  
    ...  
}
```

- any objects declared in *statement*₁ go out of scope as soon as **for** loop ends

The `for` Statement (Continued)

- consider **for** loop:

```
for (statement1; expression; statement2)  
    statement3
```

- above **for** loop can be equivalently expressed in terms of **while** loop as follows (except for behavior of **continue** statement, yet to be discussed):

```
{  
    statement1;  
    while (expression) {  
        statement3  
        statement2;  
    }  
}
```

The `for` Statement: Example

- example with single statement in loop body:

```
// Print the integers from 0 to 9 inclusive.  
for (int i = 0; i < 10; ++i)  
    std::cout << i << '\n';
```

- example with multiple statements in loop body:

```
int values[10];  
// ...  
int sum = 0;  
for (int i = 0; i < 10; ++i) {  
    // Stop if value is negative.  
    if (values[i] < 0) {  
        break;  
    }  
    sum += values[i];  
}
```

- example with error in assumption about scoping rules:

```
for (int i = 0; i < 10; ++i) {  
    std::cout << i << '\n';  
}  
++i; // ERROR: i no longer exists
```

Range-Based `for` Statement

- variant of `for` loop for iterating over elements in range
- example:

```
int array[4] = {1, 2, 3, 4};  
// Triple the value of each element in the array.  
for (int& x : array) {  
    x *= 3;  
}
```

- range-based `for` loop nice in that it clearly expresses programmer intent (i.e., iterate over each element of collection)

The do Statement

- looping construct
- has following general syntax:

```
do
    statement
while (expression) ;
```

- statement *statement* executed;
then, expression *expression* evaluated;
if expression *expression* is true, entire process repeats from beginning
- to execute multiple statements in body of loop, must group multiple statements into single statement using brace brackets

```
do {
    statement1
    statement2
    ...
} while (expression) ;
```

- advisable to *always use brace brackets*, even when loop body consists of only one statement

The do Statement: Example

- example with single statement in loop body:

```
// delay by looping 10000 times  
int n = 0;  
do  
    ++n;  
while (n < 10000);
```

- example with multiple statements in loop body:

```
// print integers from 0 to 9 inclusive  
int n = 0;  
do {  
    std::cout << n << '\n';  
    ++n;  
} while (n < 10);
```

The `break` Statement

- **break** statement causes enclosing loop or switch to be terminated immediately
- example:

```
// Read integers from standard input until an  
// error or end-of-file is encountered or a  
// negative integer is read.  
int x;  
while (std::cin >> x) {  
    if (x < 0) {  
        break;  
    }  
    std::cout << x << '\n';  
}
```

The `continue` Statement

- **`continue`** statement causes next iteration of enclosing loop to be started immediately
- example:

```
int values[10];  
...  
// Print the nonzero elements of the array.  
for (int i = 0; i < 10; ++i) {  
    if (values[i] == 0) {  
        // Skip over zero elements.  
        continue;  
    }  
    // Print the (nonzero) element.  
    std::cout << values[i] << '\n';  
}
```


The goto Statement

- **goto** statement transfers control to another statement specified by label
- should generally try to *avoid use of goto statement*
- well written code rarely has legitimate use for **goto** statement
- example:

```
int i = 0;
loop: // label for goto statement
do {
    if (i == 3) {
        ++i;
        goto loop;
    }
    std::cout << i << '\n';
    ++i;
} while (i < 10);
```

- some restrictions on use of **goto** (e.g., cannot jump over initialization in same block as **goto**)

```
goto skip; // ERROR
int i = 0;
skip:
++i;
```

Section 2.3.4

Functions

Function Parameters, Arguments, and Return Values

- **argument** (a.k.a. **actual parameter**): argument is value supplied to function by caller; appears in parentheses of function-call operator
- **parameter** (a.k.a. **formal parameter**): parameter is object/reference declared as part of function that acquires value on entry to function; appears in function definition/declaration
- although abuse of terminology, parameter and argument often used interchangeably
- **return value**: result passed from function back to caller

```
int square(int i) { // i is parameter  
    return i * i; // return value is i * i  
}
```

```
void compute() {  
    int i = 3;  
    int j = square(i); // i is argument  
}
```

Function Declarations and Definitions

- **function declaration** introduces identifier that names function and specifies following properties of function:
 - number of parameters
 - type of each parameter
 - type of return value (if not automatically deduced)

- example:

```
bool isOdd(int); // declare isOdd
bool isOdd(int x); // declare isOdd (x ignored)
```

- **function definition** provides all information included in function declaration as well as code for body of function

- example:

```
bool isOdd(int x) { // declare and define isOdd
    return x % 2;
}
```

Basic Syntax (Leading Return Type)

- most basic syntax for function declarations and definitions places return type at start (i.e., leading return-type syntax)
- basic syntax for function declaration:

```
return_type function_name (parameter_declarations) ;
```

- examples of function declarations:

```
int min(int, int);  
double square(double);
```

- basic syntax for function definition:

```
return_type function_name (parameter_declarations)  
{  
    statements  
}
```

- examples of function definitions:

```
int min(int x, int y) {return x < y ? x : y;}  
double square(double x) {return x * x;}
```

Trailing Return-Type Syntax

- with trailing return-type syntax, return type comes after parameter declarations and **auto** used as placeholder for where return type would normally be placed
- trailing return-type syntax for function declaration:

```
auto function_name (parameter_declarations) -> return_type;
```

- examples of function declarations:

```
auto min(int, int) -> int;  
auto square(double) -> double;
```

- trailing return-type syntax for function definition:

```
auto function_name (parameter_declarations) -> return_type  
{  
    statements  
}
```

- examples of function definitions:

```
auto min(int x, int y) -> int  
    {return x < y ? x : y;}  
auto square(double x) -> double {return x * x;}
```

The `return` Statement

- **return** statement used to exit function, passing specified return value (if any) back to caller
- code in function executes until **return** statement is reached or execution falls off end of function
- if function return type is not **void**, **return** statement takes single parameter indicating value to be returned
- if function return type is **void**, function does not return any value and **return** statement takes no parameter
- falling off end of function equivalent to executing **return** statement with no value
- example:

```
double unit_step(double x) {  
    if (x >= 0.0) {  
        return 1.0; // exit with return value 1.0  
    }  
    return 0.0; // exit with return value 0.0  
}
```

Automatic Return-Type Deduction

- with both leading and trailing return-type syntax, can specify return type as **auto**
- in this case, return type of function will be automatically deduced
- if function definition has no **return** statement, return type deduced to be **void**
- otherwise, return type deduced to match type in expression of **return** statement or, if **return** statement has no expression, as **void**
- if multiple return statements, must use same type for all **return** expressions
- when return-type deduction used, function definition must be visible in order to call function (since return type cannot be determined otherwise)
- example:

```
auto square(double x) {  
    return x * x;  
    // x * x has type double  
    // deduced return type is double  
}
```


The `main` Function

- entry point to program is always function called `main`
- has return type of `int`
- can be declared to take either no arguments or two arguments as follows (although other possibilities may also be supported by implementation):

```
int main();
```

```
int main(int argc, char* argv[]);
```

- two-argument variant allows arbitrary number of C-style strings to be passed to program from environment in which program run
- `argc`: number of C-style strings provided to program
- `argv`: array of pointers to C-style strings
- `argv[0]` is name by which program invoked
- `argv[argc]` is guaranteed to be 0 (i.e., null pointer)
- `argv[1]`, `argv[2]`, ..., `argv[argc - 1]` typically correspond to command line options

The `main` Function (Continued)

- suppose that following command line given to shell:

```
program one two three
```

- `main` function would be invoked as follows:

```
int argc = 4;
char* argv[] = {
    "program", "one", "two", "three", 0
};
main(argc, argv);
```

- return value of `main` typically passed back to operating system
- can also use function `void exit(int)` to terminate program, passing integer return value back to operating system
- return statement in `main` is optional
- if control reaches end of `main` without encountering return statement, effect is that of executing “`return 0;`”

- **lifetime** of object is period of time in which object exists (e.g., block, function, global)

```
int x;
```

```
void wasteTime ()
```

```
{  
    int j = 10000;  
    while (j > 0) {  
        --j;  
    }  
    for (int i = 0; i < 10000; ++i) {  
    }  
}
```

- in above example: *x* global scope and lifetime; *j* function scope and lifetime; *i* block scope and lifetime

Parameter Passing

- function parameter can be passed by value or by reference
- **pass by value**: function given copy of object from caller
- **pass by reference**: function given reference to object from caller
- to pass parameter by reference, use *reference type* for parameter
- example:

```
void increment(int& x)
    // x is passed by reference
{
    ++x;
}
```

```
double square(double x)
    // x is passed by value
{
    return x * x;
}
```

Pass-By-Value Versus Pass-By-Reference

- if object being passed to function is *expensive to copy* (e.g., a very large data type), always faster to pass by reference
- if function needs to *change value of object in caller*, must pass by reference
- example:

```
void increment0(int x) {  
    ++x; // Increment x by one.  
}
```

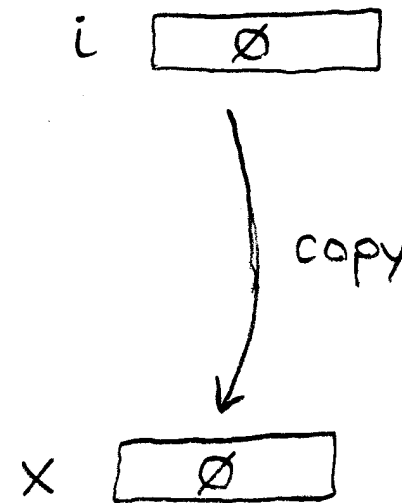
```
void increment(int& x) {  
    ++x; // Increment x by one.  
}
```

```
void func() {  
    int i = 0;  
    increment0(i); // i is passed by value  
    // i still equals 0 (i was not incremented)  
    increment(i); // i is passed by reference  
    // i equals 1 (i was incremented)  
}
```

Pass By Value

```
void func() {  
    int i = 0;  
    increment(i); // i unchanged  
}
```

```
void increment(int x) {  
    ++x;  
}
```



Pass By Reference

```
void func() {
```

```
    int i = 0;
```

```
    increment(i); // i is incremented
```

```
}
```

```
void increment(int& x) {
```

```
    ++x;
```

```
}
```

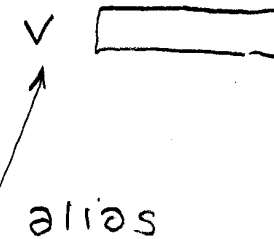


alias

Pass-By-Reference Example

```
void func() {  
    double a;  
    RealVector v = getVector();  
    a = average(v);  
}
```

```
double average(RealVector & x) {  
    double a;  
    // initialize a here  
    return a;  
}
```

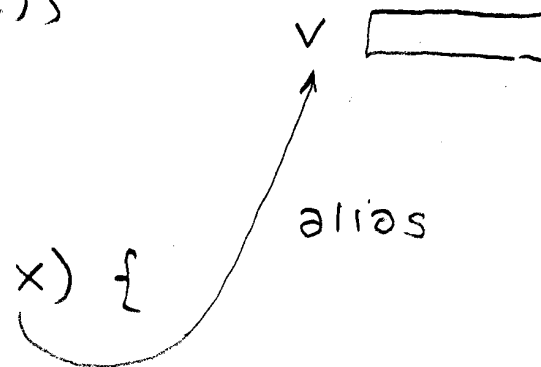


- above code is incorrect

Pass-By-Reference Example (Continued)

```
void func() {  
    double a;  
    const RealVector v = getVector();  
    a = average(v);  
}
```

```
double average(RealVector& x) {  
    double a;  
    // initialize a here  
    return a;  
}
```



- code will not compile

Inline Functions

- **inline function**: function for which compiler copies code from function definition directly into code of calling function rather than creating separate set of instructions in memory
- since code copied directly into calling function, no need to transfer control to separate piece of code and back again to caller, *eliminating performance overhead* of function call
- can request function be made inline by including **inline** qualifier along with function return type
- inline typically used for *very short functions* (where overhead of calling function is large relative to cost of executing code within function itself)
- inline function definition must be visible at point of use
- example:

```
inline bool isEven(int x) {  
    return x % 2 == 0;  
}
```

Inlining of a Function

- inlining of `isEven` function transforms code fragment 1 into code fragment 2
- Code fragment 1:

```
inline bool isEven(int x) {  
    return x % 2 == 0;  
}  
  
void myFunction() {  
    int i = 3;  
    bool result = isEven(i);  
}
```

- Code fragment 2:

```
void myFunction() {  
    int i = 3;  
    bool result = (i % 2 == 0);  
}
```

The `constexpr` Qualifier for Functions

- `constexpr` qualifier indicates return value of function is constant expression (i.e., can be evaluated at compile time) *provided that all arguments to function are constant expressions*
- `constexpr` function required to be evaluated at compile time if all arguments are constant expressions and return value *used in constant expression*
- `constexpr` functions are implicitly inline
- `constexpr` function very restricted in what it can do (e.g., no external state, can only call `constexpr` functions)
- example:

```
constexpr int factorial(int n) {
    return n >= 2 ? (n * factorial(n - 1)) : 1;
}

int u[factorial(5)];
    // OK: factorial(5) is constant expression

int x = 5;
int v[factorial(x)];
    // ERROR: factorial(x) is not constant
    // expression
```

Constexpr Function Example: square

```
1  #include <iostream>
2
3  constexpr double square(double x) {
4      return x * x;
5  }
6
7  int main() {
8      constexpr double a = square(2.0);
9      // must be computed at compile time
10
11     double b = square(0.5);
12     // might be computed at compile time
13
14     double t;
15     if (!(std::cin >> t)) {
16         return 1;
17     }
18     const double c = square(t);
19     // must be computed at run time
20
21     std::cout << a << ' ' << b << ' ' << c << '\n';
22 }
```

constexpr Function Example: power_int (Recursive)

```
1  #include <iostream>
2
3  constexpr double power_int_helper(double x, int n) {
4      return (n > 0) ? x * power_int_helper(x, n - 1) : 1;
5  }
6
7  constexpr double power_int(double x, int n) {
8      return (n < 0) ? power_int_helper(1.0 / x, -n) :
9          power_int_helper(x, n);
10 }
11
12 int main() {
13     constexpr double a = power_int(0.5, 8);
14     // must be computed at compile time
15
16     double b = power_int(0.5, 8);
17     // might be computed at compile time
18
19     double x;
20     if (!(std::cin >> x)) {return 1;}
21     const double c = power_int(x, 2);
22     // must be computed at run time
23
24     std::cout << a << ' ' << b << ' ' << c << '\n';
25 }
```

Constexpr Function Example: `power_int` (Iterative)

```
1  #include <iostream>
2
3  constexpr double power_int(double x, int n) {
4      double result = 1.0;
5      if (n < 0) {
6          x = 1.0 / x;
7          n = -n;
8      }
9      while (--n >= 0) {
10         result *= x;
11     }
12     return result;
13 }
14
15 int main() {
16     constexpr double a = power_int(0.5, 8);
17     // must be computed at compile time
18
19     double b = power_int(0.5, 8);
20     // might be computed at compile time
21
22     double x;
23     if (!(std::cin >> x)) {return 1;}
24     const double c = power_int(x, 2);
25     // must be computed at run time
26
27     std::cout << a << ' ' << b << ' ' << c << '\n';
28 }
```

Compile-Time Versus Run-Time Computation

- constexpr variables and constexpr functions provide mechanism for moving computation from run time to compile time
- benefits of compile-time computation include:
 - ① no execution-time cost at run-time
 - ② can reduce code size since code used only for compile-time computation does not need to be included in executable
 - ③ can find errors at compile-time and link-time instead of at run time
 - ④ no synchronization concerns
- when floating point is involved, compile-time and run-time computations can yield different results, due to differences in such things as
 - rounding mode in effect
 - processor architecture used for computation (when cross compiling)

Function Overloading

- **function overloading**: multiple functions can have same name as long as they differ in number/type of their arguments
- example:

```
void print(int x) {  
    std::cout << "int has value " << x << '\n';  
}
```

```
void print(double x) {  
    std::cout << "double has value " << x << '\n';  
}
```

```
void demo() {  
    int i = 5;  
    double d = 1.414;  
    print(i); // calls print(int)  
    print(d); // calls print(double)  
    print(42); // calls print(int)  
    print(3.14); // calls print(double)  
}
```

Default Arguments

- can specify default values for arguments to functions
- example:

```
// Compute log base b of x.
```

```
double logarithm(double x, double b) {  
    return std::log(x) / std::log(b);  
}
```

```
// Declaration of logarithm with a default argument.
```

```
double logarithm(double, double = 10.0);
```

```
void demo() {
```

```
    double x =
```

```
        logarithm(100.0); // calls logarithm(100.0, 10.0)
```

```
    double y =
```

```
        logarithm(4.0, 2.0); // calls logarithm(4.0, 2.0)
```

```
}
```

Argument Matching

- call of given function name chooses function that best matches actual arguments
- consider all functions in scope for which set of conversions exists so function could possibly be called
- best match is intersection of sets of functions that best match on each argument
- matches attempted in following order:
 - ① exact match with zero or more trivial conversions (e.g., `T` to `T&`, `T&` to `T`, adding **`const`** and/or **`volatile`**); of these, those that do not add **`const`** and/or **`volatile`** to pointer/reference better than those that do
 - ② match with promotions (e.g., **`int`** to **`long`**, **`float`** to **`double`**)
 - ③ match with standard conversions (e.g., **`float`** to **`int`**, **`double`** to **`int`**)
 - ④ match with user-defined conversions
 - ⑤ match with ellipsis
- if set of best matches contains exactly one element, this element chosen as function to call
- if set of best matches is either empty or contains more than one element, function call is invalid (since either no matches found or multiple equally-good matches found)

Argument Matching: Example

```
int max(int, int);  
double max(double, double);  
  
int i, j, k;  
double a, b, c;  
  
// ...  
k = max(i, j);  
    // best match on first argument: max(int, int)  
    // best match on second argument: max(int, int)  
    // best match: max(int, int)  
    // OK: calls max(int, int)  
c = max(a, b);  
    // best match on first argument: max(double, double)  
    // best match on second argument: max(double, double)  
    // best match: max(double, double)  
    // OK: calls max(double, double)  
c = max(i, b);  
    // best match on first argument: max(int, int)  
    // best match on second argument: max(double, double)  
    // best match: empty set  
    // ERROR: ambiguous function call
```

The `assert` Macro

- `assert` macro allows testing of boolean condition at run time
- typically used to test sanity of code (e.g., test preconditions, postconditions, or other invariants) or test validity of assumptions made by code
- defined in header file `cassert`
- macro takes single argument: boolean expression
- if assertion fails, program is terminated by calling `std::abort`
- if `NDEBUG` preprocessor symbol is defined at time `cassert` header file included, all assertions are disabled (i.e., not checked)
- example:

```
#include <cassert>

double sqrt(double x) {
    assert(x >= 0);
    // ...
}
```

Section 2.3.5

Input/Output (I/O)

- relevant declarations and such in header file `iostream`
- `std::istream`: stream from which characters/data can be read (i.e., input stream)
- `std::ostream`: stream to which characters/data can be written (i.e., output stream)
- `std::istream` `std::cin` standard input stream
- `std::ostream` `std::cout` standard output stream
- `std::ostream` `std::cerr` standard error stream
- in most environments, above three streams refer to user's terminal by default
- output operator (inserter) `<<`
- input operator (extractor) `>>`
- stream can be used as **bool** expression; converts to **true** if stream has not encountered any errors and **false** otherwise (e.g., if invalid data read or I/O error occurred)

Basic I/O: Example

- example:

```
std::cout << "Enter an integer: ";  
int x;  
std::cin >> x;  
if (std::cin) {  
    std::cout << "The integer entered was "  
        << x << '\n';  
} else {  
    std::cerr <<  
        "End-of-file reached or I/O error" << '\n';  
}
```


I/O Manipulators

- manipulators provide way to control formatting of data values written to streams as well as parsing of data values read from streams
- declarations related information for manipulators can be found in header files: `ios`, `iomanip`, `istream`, and `ostream`
- most manipulators used to control output formatting
- focus here on manipulators as they pertain to output
- manipulator may have *immediate* effect (e.g., `endl`), only affect *next* data value output (e.g., `setw`), or affect *all* subsequent data values output (e.g., `setprecision`)

I/O Manipulators (Continued)

Name	Description
<code>setw</code>	set field width
<code>setfill</code>	set fill character
<code>endl</code>	insert newline and flush
<code>flush</code>	flush stream
<code>dec</code>	use decimal
<code>hex</code>	use hexadecimal
<code>oct</code>	use octal
<code>showpos</code>	show positive sign
<code>noshowpos</code>	do not show positive sign
<code>left</code>	left align
<code>right</code>	right align
<code>fixed</code>	write floating-point values in fixed-point notation
<code>scientific</code>	write floating-point values in scientific notation
<code>setprecision</code>	for default notation, specify maximum number of meaningful digits to display before and after decimal point; for fixed and scientific notations, specify exactly how many digits to display after decimal point (padding with trailing zeros if necessary)

I/O Manipulators Example

● example:

```
#include <iostream>
#include <iomanip>

int main(int argc, char** argv)
{
    const double pi = 3.1415926535;
    const double big = 123456789.0;

    // default notation
    std::cout << pi << ' ' << big << '\n';
    // fixed-point notation
    std::cout << std::fixed << pi << ' ' << big << '\n';
    // scientific notation
    std::cout << std::scientific << pi << ' ' << big << '\n';
    // fixed-point notation with 7 digits after decimal point
    std::cout << std::fixed << std::setprecision(7) << pi << ' '
        << big << '\n';
    // fixed-point notation with precision and width specified
    std::cout << std::setw(8) << std::fixed << std::setprecision(2)
        << pi << ' ' << std::setw(20) << big << '\n';
    // fixed-point notation with precision, width, and fill specified
    std::cout << std::setw(8) << std::setfill('x') << std::fixed
        << std::setprecision(2) << pi << ' ' << std::setw(20) << big << '\n';

    return 0;
}
```

● output:

```
3.14159 1.23457e+08
3.141593 123456789.000000
3.141593e+00 1.234568e+08
3.1415927 123456789.0000000
    3.14          123456789.00
xxxx3.14 xxxxxxxx123456789.00
```

Section 2.3.6

Miscellany

Namespaces

- mechanism for reducing likelihood of naming conflicts (i.e., attempt to use same identifier to have different meaning in various places in code)
- has general syntax:

```
namespace name {  
    code  
}
```

- all identifiers (e.g., variable names, function names, type names) declared/defined in code *code* (i.e., code contained in namespace body) made to belong to namespace *name*
- identifiers only have to be unique within a single namespace
- same identifier can be re-used in different namespaces
- scope-resolution operator (i.e., `::`) used to specify namespace to which particular identifier belongs
- **using** statement can be used to make identifiers declared in different namespaces appear as if they were in current namespace

Namespaces: Example

```
using std::cout;

namespace mike {
    int someValue;
    void initialize() {
        cout << "mike::initialize called\n";
        someValue = 0;
    }
}

namespace fred {
    double someValue;
    void initialize() {
        cout << "fred::initialize called\n";
        someValue = 1.0;
    }
}

mike::initialize(); // call initialize in namespace mike
fred::initialize(); // call initialize in namespace fred
using mike::initialize;
initialize(); // call initialize in mike namespace
```

Memory Allocation: `new` and `delete`

- to allocate memory, use **new** statement
- to deallocate memory allocated with **new** statement, use **delete** statement
- similar to `malloc` and `free` in C
- two forms of allocation: 1) single object (i.e., nonarray case) and 2) array of objects
- array version of `new/delete` distinguished by `[]`
- example:

```
char* buffer = new char[64]; // allocate
                               // array of 64 chars
delete [] buffer; // deallocate array
double* x = new double; // allocate single double
delete x; // deallocate single object
```

- important to match nonarray and array versions of **new** and **delete**:

```
char* buffer = new char[64]; // allocate
delete buffer; // ERROR: nonarray delete to
                // delete array
                // may compile fine, but crash
```

User-Defined Literals Example

```
1  #include <iostream>
2  #include <complex>
3
4  std::complex<long double> operator "" _i(long double d) {
5      return std::complex<long double>(0.0, d);
6  }
7
8  int main() {
9      auto z = 3.14_i;
10     std::cout << z << '\n';
11 }
12
13 // Program output:
14 // (0,3.14)
```


Section 2.4

Classes

Section 2.4.1

Classes, Members, and Access Specifiers

Classes

- **class** is user-defined type
- class specifies:
 - ① how objects of class are *represented*
 - ② *operations* that can be performed on objects of class
- class consists of *zero or more members*
- members can be of various types: data member, function member, and others (e.g., type member)
- **data members** define representation of object of class
- **function members** (also called member functions) provide operations on such objects
- **type members** specify any types associated with class
- **interface** is part of class that is directly accessible to its users
- **implementation** is part of class that its users access only indirectly through interface

Access Specifiers (Public and Private)

- can control *level of access* that users of class have to its members
- three levels of access: private, protected, and public
- **private**: member can only be accessed by other members of class and friends of class
- **public**: member can be accessed by any code
- **protected**: relates to inheritance (discussion deferred until later)
- public members constitute class interface
- private members constitute class implementation

Class Example

- class typically has form:

```
class MyClass // The class is named MyClass.
{
public:
    // public members
    // (i.e., the interface to users)
    // usually functions and types (but not data)
private:
    // private members
    // (i.e., the implementation details only
    // accessible by members of class)
    // usually functions, types, and data
};
```

Default Member Access

- class members are private by default
- two code examples below are exactly equivalent:

```
class MyClass {  
    // ...  
};
```

```
class MyClass {  
private:  
    // ...  
};
```

The `struct` Keyword

- `struct` is class where members public by default
- two code examples below are exactly equivalent:

```
struct MyClass {  
    // ...  
};
```

```
class MyClass {  
public:  
    // ...  
};
```

Data Members

- class example:

```
class Vector_2 { // Two-dimensional vector class.
public:
    double x; // The x component of the vector.
    double y; // The y component of the vector.
};

void func() {
    Vector_2 v;
    v.x = 1.0; // Set data member x to 1.0
    v.y = 2.0; // Set data member y to 2.0
}
```

- above class has data members `x` and `y`
- members accessed by *member-selection operator* (i.e., “. ”)

Function Members

- class example:

```
class Vector_2 { // Two-dimensional vector class.
public:
    double x; // The x component of the vector.
    double y; // The y component of the vector.
    void initialize(double x_, double y_);
};

void Vector_2::initialize(double x_, double y_) {
    x = x_;
    y = y_;
}

void func() {
    Vector_2 v; // Create Vector_2 called v.
    v.initialize(1.0, 2.0); // Initialize v to (1.0, 2.0).
}
```

- above class has member function `initialize`
- to refer to member of class outside of class body must use *scope-resolution operator* (i.e., `::`)
- for example, in case of `initialize` function, we use `Vector_2::initialize`
- member function always has object of class as *implicit parameter*

The `this` Keyword

- member function always has object of class as *implicit parameter*
- implicit parameter passed in form of pointer using special variable called **this**
- normally, we do not explicitly write “**this**”, however
- example:

```
class MyClass {
public:
    int updateValue(int newValue) {
        int oldValue = value;
        value = newValue; // "value" means "this->value"
        return oldValue;
    }
private:
    int value;
};

void func() {
    MyClass x;
    x.updateValue(5);
    // in MyClass::updateValue, variable this equals &x
}
```

Definition of Function Members in Class Body

- member function whose definition is provided in body of class is automatically **inline**
- two code examples below are exactly equivalent:

```
class MyInteger {  
public:  
    // Set the value of the integer and return the old value.  
    int setValue(int newValue) {  
        int oldValue = value;  
        value = newValue;  
        return oldValue;  
    }  
private:  
    int value;  
};
```

```
class MyInteger {  
public:  
    // Set the value of the integer and return the old value.  
    int setValue(int newValue);  
private:  
    int value;  
};  
  
inline int MyInteger::setValue(int newValue) {  
    int oldValue = value;  
    value = newValue;  
    return oldValue;  
}
```

Type Members

- example:

```
class Point_2 { // Two-dimensional point class.  
public:  
    typedef double Coordinate; // Coordinate type.  
    Coordinate x; // The x coordinate of the point.  
    Coordinate y; // The y coordinate of the point.  
};  
  
void func() {  
    Point_2 p;  
    // ...  
    Point_2::Coordinate x = p.x;  
    // Point_2::Coordinate same as double  
}
```

- above class has type member `Coordinate`
- to refer to type member outside of class body, we must use *scope-resolution operator* (i.e., `::`)

- normally, only class has access to its private members
- sometimes, necessary to allow another class or function to have access to private members of class
- friend of class is function/class that is allowed to access private members of class
- to make function or class friend of another class, use **friend** statement
- example:

```
class SomeClass; // forward declaration of SomeClass

class MyClass {
    // ...
    friend void myFunc(); // function myFunc is
                        // friend of MyClass
    friend class SomeClass; // class SomeClass is
                          // friend of MyClass
    // ...
};
```

Class Example

```
class MyClass {
public:
    int setValue(int newValue) { // member function
        int oldValue = value; // save old value
        value = newValue; // change value to new value
        return oldValue; // return old value
    }
private:
    friend void wasteTime();
    void doNothing() {}
    int value; // data member
};

void wasteTime() {
    MyClass x;
    x.doNothing(); // OK: friend
    x.value = 5; // OK: friend
}

void func() {
    MyClass x; // x is object of type MyClass
    x.setValue(5); // call MyClass's setValue member
                  // (sets x.value to 5)
    x.value = 5; // ERROR: value is private
    x.doNothing(); // ERROR: doNothing is private
}
```

const Member Functions

- need way to indicate if member function can change value of object
- **const** member function cannot change value of object

```
class Counter {  
public:  
    int getCount () const {return count;}  
    void setCount (int newCount) {count = newCount;}  
    void incrementCount () {++count;}  
private:  
    int count;  
};  
  
void func () {  
    Counter ctr;  
    ctr.setCount (0);  
    int count = ctr.getCount ();  
    const Counter& ctr2 = ctr;  
    count = ctr2.getCount (); // getCount better be const  
}
```

Propagating Values: Copying and Moving

- Suppose that we have two objects of the same type and we want to propagate the value of one object (i.e., the source) to the other object (i.e., the destination).
- This can be accomplished in one of two ways: 1) copying or 2) moving.
- **Copying** propagates the value of the source object to the destination object *without modifying the source object*.
- **Moving** propagates the value of the source object to the destination object and is *permitted to modify the source object*.
- Moving is always at least as efficient as copying, and for many types, moving is *more efficient* than copying.
- For some types, *copying does not make sense*, while moving does (e.g., `std::ostream`, `std::istream`).

Section 2.4.2

Constructors and Destructors

Constructors

- when new object created usually desirable to immediately initialize it to some known state
- prevents object from accidentally being used before it is initialized
- **constructor** is member function that is *called automatically* when object created in order to *initialize* its value
- constructor has *same name as class* (i.e., constructor for class `T` is function `T::T`)
- constructor has *no return type* (not even `void`)
- constructor *cannot be called directly* (although placement new provides mechanism for achieving similar effect, in rare cases when needed)
- constructor *can be overloaded*
- before constructor body is entered, all data members of class type are first constructed
- in certain circumstances, constructors may be automatically provided
- sometimes, automatically provided constructors *will not* have correct behavior

Default Constructor

- constructor that can be called with no arguments known as **default constructor**
- if *no constructors* specified, default constructor *automatically provided* that calls default constructor for each data member of class type (does nothing for data member of built-in type)

```
class Vector { // Two-dimensional vector class.  
public:  
    Vector() { // Default constructor.  
        x_ = 0.0; y_ = 0.0;  
    }  
    // ...
```

```
private:  
    double x_; // The x component of the vector.  
    double y_; // The y component of the vector.  
};
```

```
Vector u; // calls Vector(); u set to (0,0)  
Vector x(); // declares function x that returns Vector
```

Copy Constructor

- for class `T`, constructor taking lvalue reference to `T` as first parameter that can be called with one argument known as **copy constructor**
- used to create object by copying from already-existing object
- copy constructor for class `T` typically is of form `T (const T&)`
- if *no copy constructor* specified (and no move constructor or move assignment operator specified), copy constructor is *automatically provided* that copies each data member (using copy constructor for class and bitwise copy for built-in type)

```
class Vector { // Two-dimensional vector class.
public:
    // ... (e.g., default constructor)
    Vector(const Vector& v) { // Copy constructor.
        x_ = v.x_; y_ = v.y_;
    }
    // ...
private:
    double x_; // The x component of the vector.
    double y_; // The y component of the vector.
};

Vector v;
Vector w(v); // calls Vector(const Vector&)
Vector u = v; // calls Vector(const Vector&)
```

Move Constructor

- for class `T`, constructor taking rvalue reference to `T` as first parameter that can be called with one argument known as **move constructor**
- used to create object by moving from already-existing object
- move constructor for class `T` typically is of form `T (T&&)`
- if **no move constructor** specified (and no destructor, copy constructor, or copy/move assignment operator specified), move constructor is **automatically provided** that moves each data member (using move for class and bitwise copy for built-in type)

```
class Vector { // Two-dimensional vector class.
public:
    // ...
    Vector(Vector&& v) { // Move constructor.
        x_ = v.x_; y_ = v.y_;
    }
    // ...
private:
    double x_; // The x component of the vector.
    double y_; // The y component of the vector.
};
```

```
Vector x(); // declares function x that returns Vector
Vector y = x(); // calls Vector(Vector&&) if move not elided
```

Constructor Example

```
class Vector { // Two-dimensional vector class.
public:
    Vector() { // Default constructor.
        x_ = 0.0; y_ = 0.0;
    }
    Vector(const Vector& v) { // Copy constructor.
        x_ = v.x_; y_ = v.y_;
    }
    Vector(Vector&& v) { // Move constructor.
        x_ = v.x_; y_ = v.y_;
    }
    Vector(double x, double y) { // Another constructor.
        x_ = x; y_ = y;
    }
    // ...
private:
    double x_; // The x component of the vector.
    double y_; // The y component of the vector.
};

Vector u; // calls Vector(); u set to (0,0)
Vector v(1.0, 2.0); // calls Vector(double, double)
Vector w(v); // calls Vector(const Vector&)
Vector z = u; // calls Vector(const Vector&)
Vector x(); // declares function x that returns Vector
Vector y = x(); // calls Vector(Vector&&) if move not elided
```

- four constructors provided

Initializer Lists

- in constructor of class, often we want to control which constructor is used to initialize each data member
- since all data members are constructed *before* body of constructor is entered, this cannot be controlled inside body of constructor
- to allow control over which constructors are used to initialize individual data members, mechanism called **initializer lists** provided
- initializer list forces specific constructors to be used to initialize individual data members before body of constructor is entered
- data members always initialized in *order of declaration*, regardless of order in initializer list

Initializer List Example

```
class ArrayDouble { // array of doubles class
public:
    ArrayDouble(); // create empty array
    ArrayDouble(int size); // create array of specified size
    // ...
private:
    // ...
};

class Vector { // n-dimensional real vector class
public:
    Vector(int size) : data_(size) {}
    // force data_ to be constructed with
    // ArrayDouble::ArrayDouble(int)
    // ...
private:
    ArrayDouble data_; // elements of vector
};
```


Destructors

- when object reaches end of lifetime, typically some cleanup required before object passes out of existence
- **destructor** is member function that is *automatically called* when object reaches end of lifetime in order to perform any necessary cleanup
- often object may have allocated resources associated with it (e.g., memory, files, devices, network connections, processes/threads)
- when object destroyed, must ensure that any resources associated with object are released
- destructors often serve to *release resources* associated with object
- destructor for class `T` always has *name* `T::~~T`
- destructor has *no return type* (not even `void`)
- destructor *cannot be overloaded*
- destructor always takes *no parameters*
- if *no destructor* is specified, destructor *automatically provided* that calls destructor for each data member of class type
- sometimes, automatically provided destructor *will not* have correct behavior

Destructor Example

- example:

```
class MyClass {  
public:  
    MyClass(int bufferSize) { // Constructor.  
        // allocate some memory for buffer  
        bufferPtr = new char[bufferSize];  
    }  
    ~MyClass() { // Destructor.  
        // free memory previously allocated  
        delete [] bufferPtr;  
    }  
    // copy constructor, assignment operator, ...  
private:  
    char* bufferPtr; // pointer to start of buffer  
};
```

- without explicitly-provided destructor (i.e., with destructor automatically provided by compiler), memory associated with `bufferPtr` would not be freed

Section 2.4.3

Operator Overloading

Operator Overloading

- can specify the meaning of operator whose operands are one or more user-defined types through process known as **operator overloading**
- operators that can be overloaded:

arithmetic	+ - * / %
bitwise	^ & ~ << >>
logical	! &&
relational	< > <= >= == !=
assignment	=
compound assignment	+= -= *= /= %= ^= &= = <<= >>=
increment/decrement	++ --
subscript	[]
function call	()
address, indirection	& *
others	->* , -> new delete

- not possible to change precedence/associativity or syntax of operators
- meaning of operator specified by operator function, where name of function is **operator** followed by operator itself (e.g., **operator+**)

Operator Overloading (Continued 1)

- binary operator can be defined either by: 1) member function taking one argument, or 2) global function taking two arguments
- for any binary operator @, $a@b$ can be interpreted as `a.operator@(b)` or `operator@(a, b)`
- unary operator can be defined either by: 1) member function taking no arguments, or 2) global function taking one argument
- for any unary operator @, @a can be interpreted as `a.operator@()` or `operator@(a)`
- for any postfix unary operator @, $a@$ can be interpreted as `a.operator@(int)` or `operator@(a, int)` (where second argument only exists to distinguish postfix operators from prefix ones)
- if member and global functions both defined, argument matching rules determine which is called
- assignment, function-call, subscript, and member-selection operators must be overloaded as member functions
- if first operand of overloaded operator not object of class type, must use global function

Operator Overloading (Continued 2)

- for most part, operators can be defined quite arbitrarily for user-defined types
- for example, no requirement that “++x”, “x += 1”, and “x = x + 1” be equivalent
- of course, probably not advisable to define operators in very counterintuitive ways, as will inevitably lead to bugs in code
- some examples showing how expressions translated into function calls are as follows:

Expression	Member Function	Global Function
y = x	y.operator=(x)	—
y += x	y.operator+=(x)	operator+=(y, x)
x + y	x.operator+(y)	operator+(x, y)
++x	x.operator++()	operator++(x)
x++	x.operator++(int)	operator++(x, int)
x == y	x.operator==(y)	operator==(x, y)
x < y	x.operator<(y)	operator<(x, y)

Operator Overloading Example: Vector

```
class Vector { // Two-dimensional vector class
public:
    Vector() : x_(0.0), y_(0.0) {}
    Vector(double x, double y) : x_(x), y_(y) {}
    double x() const { return x_; }
    double y() const { return y_; }
private:
    double x_; // The x component
    double y_; // The y component
};

// Vector addition
Vector operator+(const Vector& u, const Vector& v)
    {return Vector(u.x() + v.x(), u.y() + v.y());}

// Dot product
double operator*(const Vector& u, const Vector& v)
    {return u.x() * v.x() + u.y() * v.y();}

void func() {
    Vector u(1.0, 2.0);
    Vector v(u);
    Vector w;
    w = u + v; // w.operator=(operator+(u, v))
    double c = u * v; // calls operator*(u, v)
    // since c is built-in type, assignment operator
    // does not require function call
}
```

Operator Overloading Example: Array10

```
class Array10 { // Ten-element real array class
public:
    Array10() {
        for (int i = 0; i < 10; ++i) { // Zero array
            data_[i] = 0;
        }
    }
    const double& operator[(int index) const] {
        return data_[index];
    }
    double& operator[(int index)] {
        return data_[index];
    }
private:
    double data_[10]; // array data
};

void func() {
    Array10 v;
    v[1] = 3.5; // calls Array10::operator[(int)]
    double c = v[1]; // calls Array10::operator[(int)]
    const Array10 u;
    u[1] = 2.5; // ERROR: u[1] is const
    double d = u[1]; // calls Array10::operator[(int)] const
}
```


Operator Overloading: Global Versus Member Functions

- some considerations: access to private data; whether first operand has class type

```
class Complex { // Complex number type.
public:
    Complex(double re, double im) : re_(re), im_(im) {}
    double real() const { return re_; }
    double imag() const { return im_; }
    Complex operator+(const double&);
private:
    double re_; // The real part.
    double im_; // The imaginary part.
};

// Overload as global function.
Complex operator+(const Complex& a, const double& b) {
    return Complex(a.real() + b, a.imag());
}

// Overload as member function.
Complex Complex::operator+(const double& b) {
    return Complex(real() + b, imag());
}

// This can only be accomplished with global function.
Complex operator+(const double& b, const Complex& a) {
    return Complex(b + a.real(), a.imag());
}

void myFunc() {
    Complex a(1.0, 2.0);
    Complex b(1.0, -2.0);
    double r = 2.0;
    Complex c = a + r; // could use global or member function
                       // operator+(a, r) or a.operator+(r)
    Complex d = r + a; // must use global function
                       // operator+(r, a)
                       // since r.operator+(a) will not work
}
```

Copy Assignment Operator

- for class `T`, `T::operator=` having exactly one parameter that is lvalue reference to `T` known as **copy assignment operator**
- used to assign, to already-existing object, value of another object by *copying*
- if no copy assignment operator specified (and no move constructor or move assignment operator specified), copy assignment operator *automatically provided* that copy assigns to each data member (using data member's copy assignment operator for class and bitwise copy for built-in type)
- copy assignment operator for class `T` typically is of form `T& operator=(const T&)` (returning reference to `*this`)
- copy assignment operator returns (nonconstant) reference in order to allow for statements like following to be valid (where `x`, `y`, and `z` are of type `T` and `T::modify` is a non-const member function):

```
x = y = z; // x.operator=(y.operator=(z))
(x = y) = z; // (x.operator=(y)).operator=(z)
(x = y).modify(); // (x.operator=(y)).modify()
```
- be careful to correctly consider case of self-assignment

Self-Assignment Example

- in practice, self assignment typically occurs when references (or pointers) are involved
- example:

```
void doSomething(SomeType& x, SomeType& y) {  
    x = y; // self assignment if &x == &y  
    // ...  
}
```

```
void myFunc() {  
    SomeType z;  
    // ...  
    doSomething(z, z); // results in self assignment  
    // ...  
}
```

Move Assignment Operator

- for class T , $T::\mathbf{operator}=\mathbf{}$ having exactly one parameter that is rvalue reference to T known as **move assignment operator**
- used to assign, to already-existing object, value of another object by *moving*
- if no move assignment operator specified (and no destructor, copy/move constructor, or copy assignment operator specified), move assignment operator *automatically provided* that move assigns to each data member (using move for class and bitwise copy for built-in type)
- move assignment operator for class T typically is of form $T\& \mathbf{operator}=(T\&\&)$ (returning reference to $*\mathbf{this}$)
- move assignment operator returns (nonconstant) reference for same reason as in case of copy assignment operator
- self-assignment should probably not occur in move case (but might be prudent to protect against “insane” code with assertion) (library effectively forbids self-assignment for move)

Copy/Move Assignment Operator Example: Complex

```
class Complex {
public:
    Complex(double re = 0.0, double im = 0.0) :
        re_(re), im_(im) {}
    Complex(const Complex& a) : re_(a.re_), im_(a.im_) {}
    Complex(Complex&& a) : re_(a.re_), im_(a.im_) {}
    Complex& operator=(const Complex& a) { // Copy assign
        if (this != &a) {
            re_ = a.re_; im_ = a.im_;
        }
        return *this;
    }
    Complex& operator=(Complex&& a) { // Move assign
        re_ = a.re_; im_ = a.im_;
        return *this;
    }
private:
    double re_; // The real part.
    double im_; // The imaginary part.
};

int main() {
    Complex z(1.0, 2.0);
    Complex v(1.5, 2.5);
    v = z; // v.operator=(z)
    v = Complex(0.0, 1.0); // v.operator=(Complex(0.0, 1.0))
}
```

Section 2.4.4

Miscellany

std::initializer_list Example

```
1  #include <iostream>
2  #include <vector>
3
4  class Sequence {
5  public:
6      Sequence(std::initializer_list<int> list) {
7          for (std::initializer_list<int>::const_iterator i =
8              list.begin(); i != list.end(); ++i)
9              elements_.push_back(*i);
10     }
11     void print() const {
12         for (std::vector<int>::const_iterator i =
13             elements_.begin(); i != elements_.end(); ++i)
14             std::cout << *i << '\n';
15     }
16 private:
17     std::vector<int> elements_;
18 };
19
20 int main() {
21     Sequence seq = {1, 2, 3, 4, 5, 6};
22     seq.print();
23 }
```

Explicit Constructors

- constructor callable with *single* argument can be used in implicit conversions (e.g., when attempting to obtain matching type for function parameter in function call)
- often, desirable to prevent constructor from being used for implicit conversions
- to accommodate this, constructor can be marked as explicit
- **explicit constructor** is constructor that cannot be used to perform implicit conversions
- prefixing constructor declaration with **explicit** keyword makes constructor explicit
- example:

```
class Widget {  
public:  
    explicit Widget(int); // explicit constructor  
    // ...  
};
```


Example Without Explicit Constructor

```
1  #include <cstdlib>
2
3  // one-dimensional integer array class
4  class IntArray {
5  public:
6      // create array of int with size elements
7      IntArray(std::size_t size) { /* ... */ };
8      // ...
9  };
10
11 void processArray(const IntArray& x) {
12     // ...
13 }
14
15 int main() {
16     // following lines of code almost certain to be
17     // incorrect, but valid due to implicit type
18     // conversion provided by
19     // IntArray::IntArray(std::size_t)
20     IntArray a = 42;
21     // probably incorrect
22     // implicit conversion effectively yields code:
23     // IntArray a = IntArray(42);
24     processArray(42);
25     // probably incorrect
26     // implicit conversion effectively yields code:
27     // processArray(IntArray(42));
28 }
```

Example With Explicit Constructor

```
1  #include <cstdlib>
2
3  // one-dimensional integer array class
4  class IntArray {
5  public:
6      // create array of int with size elements
7      explicit IntArray(std::size_t size) { /* ... */ };
8      // ...
9  };
10
11 void processArray(const IntArray& x) {
12     // ...
13 }
14
15 int main() {
16     IntArray a = 42; // ERROR: cannot convert
17     processArray(42); // ERROR: cannot convert
18 }
```

Explicitly Deleted/Defaulted Special Member Functions

- can explicitly default or delete special member functions (i.e., default constructor, copy constructor, move constructor, destructor, copy assignment operator, and move assignment operator)
- can also delete non-special member functions
- example:

```
class Thing {  
public:  
    Thing() = default;  
  
    // Prevent copying.  
    Thing(const Thing&) = delete;  
    Thing& operator=(const Thing&) = delete;  
  
    Thing(Thing&&) = default;  
    Thing& operator=(Thing&&) = default;  
    ~Thing() = default;  
    // ...  
};  
// Thing is movable but not copyable.
```

Assignment Operator Example: Buffer

- example:

```
class Buffer { // Character buffer class.
public:
    Buffer(int bufferSize) { // Constructor.
        bufSize_ = bufferSize;
        bufPtr_ = new char[bufferSize];
    }
    Buffer(const Buffer& buffer) { // Copy constructor.
        bufSize_ = buffer.bufSize_;
        bufPtr_ = new char[bufSize_];
        for (int i = 0; i < bufSize_; ++i)
            bufPtr_[i] = buffer.bufPtr_[i];
    }
    ~Buffer() { // Destructor.
        delete [] bufPtr_;
    }
    Buffer& operator=(const Buffer& buffer) { // Copy assignment operator.
        if (this != &buffer) {
            delete [] bufPtr_;
            bufSize_ = buffer.bufSize_;
            bufPtr_ = new char[bufSize_];
            for (int i = 0; i < bufSize_; ++i)
                bufPtr_[i] = buffer.bufPtr_[i];
        }
        return *this;
    }
    // ...
private:
    int bufSize_; // buffer size
    char* bufPtr_; // pointer to start of buffer
};
```

- without explicitly-provided assignment operator (i.e., with assignment operator automatically provided by compiler), memory leaks and memory corruption would result

Delegating Constructors

- sometimes, one constructor of class needs to performs all work of another constructor followed by some additional work
- rather than duplicate common code in both constructors, one constructor can use its initializer list to invoke other constructor (which must be only one in initializer list)
- constructor that invokes another constructor via initializer list called **delegating constructor**
- example:

```
class Widget {
public:
    Widget(char c, int i) : c_(c), i_(i) {}
    Widget(int i) : Widget('a', i) {}
    // delegating constructor
    // ...
private:
    char c_;
    int i_;
};

int main() {
    Widget w('A', 42);
    Widget v(42);
}
```

Static Data Members

- sometimes want to have object that is shared by all objects of class
- data member that is shared by all objects of class is called **static data member**
- to make data member static, declare using **static** qualifier
- static data member must (in most cases) be defined outside body of class
- example:

```
class Widget {  
public:  
    Widget () {++count_;}  
    Widget (const Widget&) {++count_;}  
    Widget (Widget&&) {++count_;}  
    ~Widget () {--count_;}  
    // ...  
private:  
    static int count_; // total number of Widget  
                       // objects in existence  
};  
  
// Define (and initialize) count member.  
int Widget::count_ = 0;
```

Static Member Functions

- sometimes want to have member function that does not operate on objects of class
- member function of class that does not operate on object of class (i.e., has no **this** variable) called **static member function**
- to make member function static, declare using **static** qualifier
- example:

```
class MyClass {
public:
    // ...
    // convert degrees to radians
    static double degToRad(double deg)
        {return (M_PI / 180.0) * deg;}
private:
    // ...
};

void func() {
    double rad;
    rad = MyClass::degToRad(45.0);
    rad = x.degToRad(45.0); // x is ignored
}
```

constexpr Member Functions

- like non-member functions, member functions can also be qualified as **constexpr** to indicate function can be computed *at compile time* provided that all arguments to function are constant expressions
- some additional restrictions on constexpr member functions relative to nonmember case (e.g., cannot be virtual)
- constexpr member function *implicitly inline*
- constexpr member function *not implicitly const* (as of C++14)

- constructors can also be qualified as **constexpr** to indicate object construction can be performed *at compile time* provided that all arguments to constructor are constant expressions
- constexpr constructor *implicitly inline*

Example: constexpr Constructors and Member Functions

```
// Two-dimensional vector class.
class Vector {
public:
    constexpr Vector() : x_(0), y_(0) {}
    constexpr Vector(double x, double y) : x_(x), y_(y) {}
    constexpr Vector(const Vector& v) : x_(v.x_), y_(v.y_) {}
    constexpr Vector(Vector&& v) : x_(v.x_), y_(v.y_) {}
    Vector& operator=(const Vector& v) {
        if (this != &v) {
            x_ = v.x_; y_ = v.y_;
        }
        return *this;
    }
    constexpr double x() const {return x_;}
    constexpr double y() const {return y_;}
    constexpr double squaredLength() const {
        return x_ * x_ + y_ * y_;
    }
    // ...
private:
    double x_; // The x component of the vector.
    double y_; // The y component of the vector.
};
```

The `mutable` Qualifier

- type for data member can be qualified as **`mutable`** meaning that member does not affect externally visible state of class
- mutable data member can be modified in const member function
- **`mutable`** qualifier often used for mutexes, condition variables, cached values, statistical information for performance analysis or debugging

Example: Mutable Qualifier for Statistical Information

```
#include <iostream>
#include <string>

class Employee {
public:
    Employee(int id, std::string& name, double salary) :
        id_(id), name_(name), salary_(salary), accessCount_(0) {}
    int getId() const {
        ++accessCount_; return id_;
    }
    std::string getName() const {
        ++accessCount_; return name_;
    }
    double getSalary() const {
        ++accessCount_; return salary_;
    }
    // ...
    // for debugging
    void outputDebugInfo(std::ostream& out) const {
        out << accessCount_ << '\n';
    }
private:
    int id_; // employee ID
    std::string name_; // employee name
    double salary_; // employee salary
    mutable unsigned long accessCount_; // for debugging
};
```

Stream Inserters

- stream inserters write data to output stream
- overload **operator**<<
- have general form
std::ostream& **operator**<<(std::ostream&, T) where type T is typically const lvalue reference type

- example:

```
std::ostream& operator<<(std::ostream& outStream,  
    const Complex& a)  
{  
    outStream << a.real() << ' ' << a.imag();  
    return outStream;  
}
```

- inserter and extractor should use *compatible formats* (i.e., what is written by extractor should be readable by inserter)

Stream Extractors

- stream extractors read data from input stream
- overload **operator>>**
- have general form
std::istream& **operator>>**(std::istream&, T) where type T is typically non-const lvalue reference type
- example:

```
std::istream& operator>>(std::istream& inStream,  
    Complex& a)  
{  
    double real = 0.0;  
    double imag = 0.0;  
    inStream >> real >> imag;  
    a = Complex(real, imag);  
    return inStream;  
}
```

Section 2.4.5

Temporary Objects

Temporary Objects

- A **temporary object** is an unnamed object introduced by the compiler.
- Temporary objects are used during:
 - evaluation of expressions
 - argument passing
 - function returns (that return by value)
 - reference initialization
- It is important to understand when temporary objects can be introduced, since the introduction of temporaries impacts performance.

- Evaluation of expression:

```
std::string s1("Hello ");  
std::string s2("World");  
std::string s;  
s = s1 + s2; // must create temporary  
// std::string _tmp(s1 + s2);  
// s = _tmp;
```

- Argument passing:

```
double func(const double& x);  
func(3); // must create temporary  
// double _tmp = 3;  
// func(_tmp);
```


Temporary Objects (Continued)

- Reference initialization:

```
int i = 2;  
const double& d = i; // must create temporary  
    // double _tmp = i;  
    // const double& d = _tmp;
```

- Function return:

```
std::string getMessage();  
std::string s;  
s = getMessage(); // must create temporary  
    // std::string _tmp(getMessage());  
    // s = _tmp;
```

- In most (but not all) circumstances, a temporary object is destroyed as the last step in evaluating the full expression that contains the point where the temporary object was created .

Temporary Objects Example

```
1  class Complex {
2  public:
3      Complex(double re = 0.0, double im = 0.0) : re_(re),
4          im_(im) {}
5      Complex(const Complex& a) = default;
6      Complex(Complex&& a) = default;
7      Complex& operator=(const Complex& a) = default;
8      Complex& operator=(Complex&& a) = default;
9      ~Complex() = default;
10     double real() const {return re_;}
11     double imag() const {return im_;}
12 private:
13     double re_; // The real part.
14     double im_; // The imaginary part.
15 };
16
17 Complex operator+(const Complex& a, const Complex& b) {
18     return Complex(a.real() + b.real(), a.imag() + b.imag());
19 }
20
21 int main() {
22     Complex a(1.0, 2.0);
23     Complex b(a + a);
24     b = a + b;
25 }
```

Temporary Objects Example (Continued)

Original code:

```
int main() {  
    Complex a(1.0, 2.0);  
    Complex b(a + a);  
    b = a + b;  
}
```

Code showing temporaries (assuming no optimization):

```
int main() {  
    Complex a(1.0, 2.0);  
    Complex _tmp1(a + a);  
    Complex b(_tmp1);  
    Complex _tmp2(a + b);  
    b = _tmp2;  
}
```

Original code:

```
Complex operator+(const Complex& a, const Complex& b) {  
    return Complex(a.real() + b.real(), a.imag() + b.imag());  
}
```

Code showing temporaries:

```
Complex operator+(const Complex& a, const Complex& b) {  
    Complex _tmp(a.real() + b.real(), a.imag() + b.imag());  
    return _tmp;  
}
```

Prefix Versus Postfix Increment/Decrement

```
1  class Counter {
2  public:
3      Counter() : count_(0) {}
4      int getCount() const {return count_;}
5      Counter& operator++() { // prefix increment
6          ++count_;
7          return *this;
8      }
9      Counter operator++(int) { // postfix increment
10         Counter old(*this);
11         ++count_;
12         return old;
13     }
14 private:
15     int count_; // counter value
16 };
17
18 int main() {
19     Counter x;
20     Counter y;
21     y = ++x; // no temporaries, int increment, operator=
22     y = x++; // 1 temporary, 1 named, 2 constructors,
23             // 2 destructors, int increment, operator=
24 }
```

Compound Assignment Versus Separate Assignment

```
1  #include <complex>
2  using std::complex;
3
4  int main() {
5      complex<double> a(1.0, 1.0);
6      complex<double> b(1.0, -1.0);
7      complex<double> z(0.0, 0.0);
8
9      // 2 temporary objects
10     // 2 constructors, 2 destructors
11     // 1 operator=, 1 operator+, 1 operator*
12     z = b * (z + a);
13
14     // no temporary objects
15     // only 1 operator+= and 1 operator*=
16     z += a;
17     z *= b;
18 }
```

Lifetime of Temporary Objects

- Normally, a temporary object is destroyed as the last step in evaluating the full expression that contains point where temporary object was created.
- First exception: When a default constructor with one or more default arguments is called to initialize an element of an array.
- Second exception: When a *reference is bound to a temporary* (or a subobject of a temporary), the lifetime of the temporary is extended to *match the lifetime* of the reference, with following *exceptions*:
 - A temporary bound to a reference member in a constructor initializer list persists until the constructor exits.
 - A temporary bound to a reference parameter in a function call persists until the completion of the full expression containing the call.
 - A temporary bound to the return value of a function in a return statement is not extended, and is destroyed at end of the full expression in the return statement.
 - A temporary bound to a reference in an initializer used in a new-expression persists until the end of the full expression containing that new-expression.

Lifetime of Temporary Objects Examples

- Example:

```
void func() {  
    std::string s1("Hello");  
    std::string s2(" ");  
    std::string s3("World!\n");  
    const std::string& s = s1 + s2 + s3;  
    std::cout << s; // OK?  
}
```

- Example:

```
const std::string& getString() {  
    return std::string("Hello");  
}  
void func() {  
    std::cout << getString(); // OK?  
}
```

Return Value Optimization (RVO)

- **return value optimization (RVO)** is compiler optimization technique that eliminates copy of return value from local object in function to object in caller

- example:

```
SomeType function() {  
    return SomeType(); // returns temporary object  
}
```

```
void caller() {  
    SomeType x = function(); // copy construction  
}
```

- without RVO: return value of function (which is local to function) is copied to new temporary object (so return value not lost when function returns); then, value of new temporary object copied to object that is to hold return value
- with RVO: return value of function is placed directly in object (in caller) that is to hold return value
- by avoiding need for temporary object to hold return value, eliminates one copy constructor and destructor call
- any good compiler should support RVO, although RVO cannot always be applied in all circumstances

Named Return Value Optimization (NRVO)

- **named return value optimization (NRVO)** is variation on RVO where return value is named object (i.e., not temporary object)

- example:

```
SomeType function() {  
    SomeType result;  
    // ...  
    return result; // returns named object  
}
```

```
void caller() {  
    SomeType x = function(); // copy construction  
}
```

- compiler optimizes away `result` in `function` and return value constructed directly in `x`
- effectively, `result` becomes reference to `x`
- code with NRVO more efficient (i.e., move/copy constructor and destructor calls eliminated)

Section 2.4.6

Functors

- **function object** (also known as **functor**) is object that can be invoked or called as if it were ordinary function
- class that provides member function that overloads **operator ()** is called **functor class** and object of that class is **functor**
- functors more flexible than functions as functors are objects and can therefore carry arbitrary state information
- functors are extremely useful, especially in generic programming
- as we will see later, standard library makes heavy use of functors

Functor Example: Less Than

```
struct LessThan { // Functor class
    bool operator() (double x, double y) {
        return x < y;
    }
};

void myFunc () {
    double a = 1.0;
    double b = 2.0;
    LessThan lessThan; // Functor
    bool result = lessThan(a, b);
    // calls LessThan::operator() (double, double)
    // lessThan is functor, not function
    // result == true
}
```

Functor Example With State

```
class IsGreater { // Functor class
public:
    IsGreater(int threshold) : threshold_(threshold) {}
    bool operator()(int x) const {
        return x > threshold_;
    }
private:
    // state information for functor
    int threshold_; // threshold for comparison
};

void myFunc() {
    IsGreater isGreater(5); // functor
    int x = 3;
    bool result = isGreater(x);
    // calls IsGreater::operator()(int)
    // result == false
}
```

Section 2.5

Templates

Templates

- **generic programming**: algorithms written in terms of types to be specified later (i.e., algorithms are generic in sense of being applicable to any type that meets only some very basic constraints)
- templates facilitate generic programming
- extremely important language feature
- avoids code duplication
- leads to highly efficient and customizable code
- promotes code reuse
- C++ standard library makes very heavy use of templates (actually, most of standard library consists of templates)
- many other libraries make heavy use of templates (e.g., CGAL, Boost)

Section 2.5.1

Function Templates

Motivation for Function Templates

- consider following functions:

```
int max(int x, int y)
    {return x > y ? x : y;}
```

```
double max(double x, double y)
    {return x > y ? x : y;}
```

```
// more similar-looking max functions...
```

- each of above functions has *same general form*; that is, for some type T , we have:

```
T max(T x, T y)
    {return x > y ? x : y;}
```

- would be nice if we did not have to repeatedly type, debug, test, and maintain nearly identical code
- in effect, would like code to be parameterized on type T

Function Templates

- **function template** is family of functions parameterized by one or parameters
- each template parameter can be: non-type (integral constant), type, template, or parameter pack (in case of variadic template)
- syntax for template function has general form:

```
template <parameter_list> function_declaration
```
- *parameter_list*: parameters on which template function depends
- *function*: function declaration
- type parameter designated by **class** or **typename** keyword
- template parameter designated by **template** keyword
- template template parameter must use **class** keyword
- non-type (integral constant) parameter designated by its type (e.g., **int**)
- example:

```
// declaration of function template
template <class T> T max(T x, T y);

// definition of function template
template <class T> T max(T x, T y)
    { return x > y ? x : y; }
```

Function Templates (Continued)

- to explicitly identify particular instance of template, use syntax:

function<parameters>

- example:

for function template declaration:

```
template <class T> T max (T x, T y);
```

max<**int**> refers to **int** max(**int**, **int**)

max<**double**> refers to **double** max(**double**, **double**)

- compiler only creates code for function template when it is instantiated (i.e., used)
- therefore, definition of function template must be visible in place where it is instantiated
- consequently, function template definitions usually appear in header file
- template code only needs to pass basic syntax checks, unless actually instantiated

Function Template Examples

```
1 // compute minimum of two values
2 template <class T>
3 T min(T x, T y) {
4     return x < y ? x : y;
5 }
6
7 // compute square of value
8 template <typename T>
9 T sqr(T x) {
10     return x * x;
11 }
12
13 // swap two values
14 template <class T>
15 void swap(T& x, T& y) {
16     T tmp = x;
17     x = y;
18     y = tmp;
19 }
20
21 // increment value by constant
22 template <int N = 1, typename T>
23 T& increment_by(T& n) {
24     n += N;
25     return n;
26 }
```

Template Function Overloading Resolution

- overload resolution proceeds (in order) as follows:
 - ① look for an exact match with zero or more trivial conversions on (nontemplate) functions; if found call it
 - ② look for function template from which function that can be called with exact match with zero or more trivial conversions can be generated; if found, call it
 - ③ try ordinary overloading resolution for functions; if function found, call it; otherwise, call is error
- in each step, if more than one match found, call is ambiguous and is error
- template function only used in case of exact match (unless explicitly forced)
- example:

```
template <class T> T max(T x, T y) {  
    return x > y ? x : y;  
}  
  
double x, y, z;  
int i, j, k;  
// ...  
z = max(x, y); // calls max<double>  
k = max(i, j); // calls max<int>  
z = max(i, x); // ERROR: no match  
z = max<double>(i, x); // calls max<double>
```

Section 2.5.2

Class Templates

Motivation for Class Templates

- consider almost identical complex number classes:

```
1  class ComplexDouble {
2      ComplexDouble(double re = 0.0, double im = 0.0) : re_(re), im_(im) {}
3      double real() const { return re_; }
4      double imag() const { return im_; }
5      // ...
6  private:
7      double re_; // real part
8      double im_; // imaginary part
9  };
10
11 class ComplexFloat {
12     ComplexFloat(float re = 0.0, float im = 0.0) : re_(re), im_(im) {}
13     float real() const { return re_; }
14     float imag() const { return im_; }
15     // ...
16 private:
17     float re_; // real part
18     float im_; // imaginary part
19 };
```

- both of above classes are special cases of following class parameterized on type T:

```
1  class Complex {
2      Complex(T re = T(0), T im = T(0)) : re_(re), im_(im) {}
3      T real() const { return re_; }
4      T imag() const { return im_; }
5      // ...
6  private:
7      T re_; // real part
8      T im_; // imaginary part
9  };
```

- again, would be nice if we did not have to repeatedly type, debug, test, and maintain nearly identical code

Class Templates

- **class template** is family of classes parameterized on one or more parameters
- each template parameter can be: non-type (integral constant), type, template, or parameter pack (in case of variadic template)
- syntax has general form:

```
template <parameter_list> class
```

- *parameter_list*: parameter list for class
- *class*: class/struct declaration or definition
- example:

```
// declaration of class template  
template <class T, unsigned int size>  
class MyArray;  
  
// definition of class template  
template <class T, unsigned int size>  
class MyArray {  
    // ...  
    T array_[size];  
};  
  
MyArray<double, 100> x;
```


Class Templates (Continued)

- compiler only generates code for class template when it is instantiated (i.e., used)
- since compiler only generates code for class template when it is instantiated, definition of template must be visible at point where instantiated
- consequently, class template code usually placed in header file
- template code only needs to pass basic syntax checks, unless actually instantiated
- compile errors related to class templates can often be very long and difficult to parse (especially, when template class has parameters that are template classes which, in turn, have parameters that are template classes, and so on)
- be careful when nesting angle brackets, since `<<` and `>>` may be parsed as left shift and right shift operators in some contexts (e.g., prior to C++11 `std::vector<std::complex<double>>` would lead to parsing error)

Class Template Example

```
1  template <class T>
2  class Complex { // complex number class template
3  public:
4      Complex(T re = T(0), T im = T(0)) :
5          re_(re), im_(im) {}
6      T real() const {
7          return re_;
8      }
9      T imag() const {
10         return im_;
11     }
12     // ...
13 private:
14     T re_; // real part
15     T im_; // imaginary part
16 };
17
18 Complex<int> zi;
19 Complex<double> zd;
```

Class-Template Default Parameters

- class template parameters can have *default values*
- example:

```
template <class T = int, unsigned int size = 2>  
struct MyArray {  
    T data[size];  
};
```

```
MyArray<> a; // MyArray<int, 2>  
MyArray<double> b; // MyArray<double, 2>  
MyArray<double, 10> b; // MyArray<double, 10>
```

Qualified Names

- **qualified name** is name that specifies scope
- example:

```
#include <iostream>
int main(int argc, char** argv)
{
    for (int i = 0; i < 10; ++i)
        std::cout << "Hello, world!" << std::endl;
    return 0;
}
```

- in above example, names `std::cout` and `std::endl` are qualified, while names `main`, `argc`, `argv`, and `i`, are not qualified

Dependent Names

- **dependent name** is name that depends on template parameter
- example:

```
template <class T>
class MyClass
{
public:
    struct Thing {
        T array[3];
    };
    Thing x;
    typedef T* Pointer;
    int i;
};
```

- names `Thing` and `Pointer` are dependent

Qualified Dependent Names

- to avoid any potential ambiguities, compiler will automatically assume qualified dependent name does not name type unless **typename** keyword is used
- must precede qualified dependent name that names type by **typename**
- following code is invalid and will cause compile error:

```
template <class T>
class MyClass {
    std::vector<T> vec; // ERROR?
    std::vector<T>::iterator iter; // ERROR
    std::vector<T>::value_type val; // ERROR
    // ...
};
```

- must use code like following instead:

```
template <class T>
class MyClass {
    typename std::vector<T> vec;
    typename std::vector<T>::iterator iter;
    typename std::vector<T>::value_type val;
    // ...
};
```

Why `typename` is Needed

```
1  int x = 42;
2
3  template <class T> void func() {
4      // The compiler must be able to check syntactic
5      // correctness of this template code without
6      // knowing T. Without knowing T, however, the
7      // meaning of following line of code is ambiguous.
8      // Is it a declaration of a variable x or an
9      // expression consisting of a binary operator*
10     // with operands T::foo and x?
11     T::foo* x;
12     // ...
13 }
14
15 struct ContainsType {
16     using foo = int*; // foo is type
17     // ...
18 };
19
20 struct ContainsValue {
21     static int foo; // foo is value
22     // ...
23 };
24
25 int main() {
26     // Which one of the following should be invalid?
27     func<ContainsType>();
28     func<ContainsValue>();
29 }
```

Template Template Parameter Example

```
1  #include <vector>
2  #include <list>
3  #include <deque>
4  #include <memory>
5
6  template <template <class, class> class Container,
7          typename Value>
8  class Stack {
9  public:
10     // ...
11 private:
12     Container<Value, std::allocator<Value>> data_;
13 };
14
15 int main() {
16     Stack<std::vector, int> s1;
17     Stack<std::list, int> s2;
18     Stack<std::deque, int> s3;
19 }
```


Section 2.5.3

Variable Templates

Variable Templates

- **variable template** is family of variables parameterized on one or more parameters
- each template parameter can be: non-type (integral constant), type, template, or parameter pack (in case of variadic templates)
- although less frequently used than function and class templates, variable templates quite useful in some situations
- syntax has general form:

```
template <parameter_list> variable
```

- *parameter_list*: parameter list for variable template
- *variable*: variable declaration
- example:

```
template <class T>  
T meaning_of_life = T(42);  
  
int x = meaning_of_life<int>;
```

Variable Template Example: pi

```
1  #include <limits>
2  #include <iostream>
3
4  template <typename T>
5  constexpr T pi =
6      T(3.14159265358979323846264338327950288419716939937510L);
7
8  int main() {
9      std::cout.precision(
10         std::numeric_limits<long double>::max_digits10);
11     std::cout
12         << pi<int> << '\n'
13         << pi<float> << '\n'
14         << pi<double> << '\n'
15         << pi<long double> << '\n';
16 }
```

Section 2.5.4

Alias Templates

Alias Templates

- **alias template** is family of types parameterized on one or more parameters
- each template parameter can be: non-type (integral constant), type, template, or parameter pack (in case of variadic templates)
- syntax has general form:

```
template <parameter_list> alias
```

- *parameter_list*: parameter list for class
- *alias*: alias declaration (i.e., with **using**)
- example:

```
template <class Value,  
        class Compare = std::greater<Value>,  
        class Alloc = std::allocator<Value>>  
using GreaterSet = typename std::set<Value, Compare,  
        Alloc>;
```

```
GreaterSet<int> x{4, 1, 3, 2};
```

Alias Template Example

```
1  #include <iostream>
2  #include <set>
3
4  template <typename Compare, typename Alloc =
5      std::allocator<int>>
6  using IntSet = typename std::set<int, Compare, Alloc>;
7
8  int main() {
9      IntSet<std::less<int>> x{1, 4, 3, 2};
10     IntSet<std::greater<int>> y{1, 4, 3, 2};
11     for (auto i : x) {
12         std::cout << i << '\n';
13     }
14     std::cout << '\n';
15     for (auto i : y) {
16         std::cout << i << '\n';
17     }
18 }
```

Section 2.5.5

Template Specialization

Template Specialization

- sometimes can be desirable to provide customized version of template for certain choices of template parameters
- customized version of templates can be specified through language feature known as **template specialization**
- two kinds of specialization: explicit and partial
- **explicit specialization** (less formally known as full specialization): customized version of template where all template parameters are fixed
- **partial specialization**: customized version of template where only some of template parameters are fixed
- class templates, function templates, and variable templates can all be specialized
- alias templates cannot be specialized
- class templates and variable templates can be partially or explicitly specialized
- function templates can only be explicitly specialized (not partially)

Explicitly-Specialized Class Template: `is_void`

```
1  template <class T>
2  struct is_void
3      {static constexpr bool value = false;};
4
5  template <>
6  struct is_void<void>
7      {static constexpr bool value = true;};
8
9  template <>
10 struct is_void<const void>
11     {static constexpr bool value = true;};
12
13 template <>
14 struct is_void<volatile void>
15     {static constexpr bool value = true;};
16
17 template <>
18 struct is_void<const volatile void>
19     {static constexpr bool value = true;};
20
21 static_assert(is_void<int>::value == false, "");
22 static_assert(is_void<double*>::value == false, "");
23 static_assert(is_void<void>::value == true, "");
24 static_assert(is_void<const void>::value == true, "");
25 static_assert(is_void<volatile void>::value == true, "");
26 static_assert(is_void<const volatile void>::value == true,
27     "");
28
29 int main() {}
```

Partially-Specialized Class Template

```
1  #include <iostream>
2
3  // unspecialized version
4  template <typename T, typename V>
5  struct Widget {
6      Widget() {std::cout << "unspecialized\n";}
7  };
8
9  // partial specialization
10 template <typename T>
11 struct Widget<int, T> {
12     Widget() {std::cout << "partial\n";}
13 };
14
15 // explicit specialization
16 template <>
17 struct Widget<int, int> {
18     Widget() {std::cout << "explicit\n";}
19 };
20
21 int main() {
22     Widget<double, int> w1; // unspecialized verion
23     Widget<int, double> w2; // partial specialization
24     Widget<int, int> w3; // explicit specialization
25 }
```

Partially-Specialized Class Template: `std::vector`

- `std::vector` class employs specialization
- consider vector of elements of type `T`
- most natural way to store elements is as array of `T`
- if `T` is `bool`, such an approach makes very inefficient use of memory, since each `bool` object requires one byte of storage
- if `T` is `bool`, would be much more memory-efficient to use array of, say, `unsigned char` and pack multiple `bool` objects in each byte
- `std::vector` accomplishes this by providing (partial) specialization for case that `T` is `bool`
- declaration of base template for `std::vector` and its partial specialization for case when `T` is `bool` are as follows:

```
template <class T, class Alloc = allocator<T>>  
class vector; // unspecialized version
```

```
template <class Alloc>  
class vector<bool, Alloc>; // partial specialization
```

Explicitly-Specialized Function Template: printPointee

```
1  #include <iostream>
2
3  // unspecialized version
4  template <class T>
5  typename std::ostream& printPointee(
6      typename std::ostream& out, const T* p)
7  {
8      return out << *p << '\n';
9  }
10
11 // specialization
12 template <>
13 typename std::ostream& printPointee<void>(
14     typename std::ostream& out, const void* p)
15 {
16     return out << *static_cast<const char*>(p) << '\n';
17 }
18
19 int main() {
20     int i = 42;
21     const int* ip = &i;
22     char c = 'A';
23     const void* vp = &c;
24     printPointee(std::cout, ip);
25     printPointee(std::cout, vp);
26 }
```

Explicitly-Specialized Variable Template: `is_void_v`

```
1  template <class T>
2  constexpr bool is_void_v = false;
3
4  template <>
5  constexpr bool is_void_v<void> = true;
6
7  template <>
8  constexpr bool is_void_v<const void> = true;
9
10 template <>
11 constexpr bool is_void_v<volatile void> = true;
12
13 template <>
14 constexpr bool is_void_v<const volatile void> = true;
15
16 static_assert(is_void_v<int> == false, "");
17 static_assert(is_void_v<double*> == false, "");
18 static_assert(is_void_v<void> == true, "");
19 static_assert(is_void_v<const void> == true, "");
20 static_assert(is_void_v<volatile void> == true, "");
21 static_assert(is_void_v<const volatile void> == true, "");
22
23 int main() {}
```

Explicitly-Specialized Variable Template: factorial

```
1  template <unsigned long long N>
2  constexpr unsigned long long
3     factorial = N * factorial<N - 1>;
4
5  template <>
6  constexpr unsigned long long
7     factorial<0> = 1;
8
9  int main() {
10     static_assert(factorial<5> == 120,
11         "factorial<5> failed");
12     static_assert(factorial<12> == 479'001'600,
13         "factorial<12> failed");
14 }
```

Partially-Specialized Variable Template: quotient

```
1  #include <limits>
2
3  // unspecialized version
4  template <int X, int Y>
5  constexpr int quotient = X / Y;
6
7  // partial specialization (which prevents division by zero)
8  template <int X>
9  constexpr int quotient<X, 0> = (X < 0) ?
10     std::numeric_limits<int>::min() :
11     std::numeric_limits<int>::max();
12
13 static_assert(quotient<4, 2> == 2, "");
14 static_assert(quotient<5, 3> == 1, "");
15 static_assert(quotient<4, 0> ==
16     std::numeric_limits<int>::max(), "");
17 static_assert(quotient<-4, 0> ==
18     std::numeric_limits<int>::min(), "");
19
20 int main() {}
```

Section 2.5.6

Variadic Templates

Variadic Templates

- language provides ability to specify template that can take variable number of arguments
- template that can take variable number of arguments called **variadic template**
- alias templates, class templates, function templates, and variable templates may be variadic
- variable number of arguments specified by using what is called parameter pack
- parameter pack is parameter that accepts (i.e., is placeholder for) zero or more arguments (of same kind)
- parameter pack used in parameter list of template to allow to variable number of template parameters
- ellipsis (i.e., “. . .”) is used in various contexts relating to parameter packs
- ellipsis after designator for kind of template argument in template parameter list designates argument is parameter pack
- ellipsis after parameter pack parameter expands parameter pack in context-sensitive manner

Variadic Template Examples

```
1  #include <tuple>
2
3  // variadic alias template
4  template <class... T>
5  using My_tuple = std::tuple<bool, T...>;
6
7  // variadic class template
8  template <int... Values>
9  class Integer_sequence {
10     // ...
11 };
12
13 // variadic function template
14 template <class... Ts>
15 void print(const Ts&... values) {
16     // ...
17 }
18
19 // variadic variable template
20 template <typename T, T... Values>
21 constexpr T array[] = {Values...};
22
23 int main() {
24     Integer_sequence<1, 3, 4, 2> x;
25     auto a = array<int, 1, 2, 4, 8>;
26     My_tuple<int, double> t(true, 42, 42.0);
27     print(1'000'000, 1, 43.2, "Hello");
28 }
```

Parameter Pack Expansion

- parameter pack expansion allowed in following contexts:
 - inside parentheses of function call operator
 - in template argument list
 - in function parameter list
 - in template parameter list
 - base class specifiers in class declaration
 - member initializer lists
 - braced initializer lists
 - lambda captures

The `sizeof...` Operator

- `sizeof...` operator yields number of elements in parameter pack
- example:

```
template <int... Values>
constexpr int num_parms = sizeof...(Values);

static_assert(num_parms < 1, 2, 3> == 3, "");
static_assert(num_parms <> == 0, "");
```

- example:

```
#include <cassert>

template <typename... Ts>
int number_of_arguments(const Ts&... args) {
    return sizeof...(args);
}

int main() {
    assert(number_of_arguments(1, 2, 3) == 3);
    assert(number_of_arguments() == 0);
}
```

Variadic Function Template: maximum

```
1  #include <iostream>
2  #include <type_traits>
3  #include <string>
4
5  using namespace std::literals;
6
7  template <typename T1, typename T2>
8  typename std::common_type_t<const T1&, const T2&>
9  maximum(const T1 &a, const T2 &b) {
10     return a > b ? a : b;
11 }
12
13 template <typename T1, typename T2, typename ... Args>
14 typename std::common_type_t<const T1&, const T2&,
15     const Args& ...>
16 maximum(const T1 &a, const T2 &b, const Args& ... args) {
17     return maximum(maximum(a, b), args...);
18 }
19
20 int main() {
21     std::cout << maximum(1, 2, 3, 4, -1.4) << '\n';
22     std::cout << maximum(-1'000'000L, -42L, 10, 42) << '\n';
23     std::cout << maximum("apple"s, "zebra"s, "c++"s) << '\n';
24 }
```

Variadic Function Template: print

```
1  #include <iostream>
2  #include <iomanip>
3
4  template <typename T>
5  std::ostream& print(std::ostream& out, const T& value) {
6      return out << value;
7  }
8
9  template <typename T, typename... Args>
10 std::ostream& print(std::ostream& out, const T& value,
11     Args... args) {
12     if (!(out << value)) {
13         return out;
14     }
15     return print(out, args...);
16 }
17
18 int main() {
19     print(std::cout, "Hello, World!\n",
20         std::left, std::setfill('x'),
21         std::setw(10), 1234, ' ',
22         std::setw(10), 3.1415, '\n');
23 }
```

Variadic Function Template With Template Template

Parameter: print_container

```
1  #include <iostream>
2  #include <vector>
3  #include <string>
4  #include <set>
5
6  template <template <typename, typename...>
7      class ContainerType, typename ValueType, typename... Args>
8  bool print_container(const ContainerType<ValueType, Args...>&
9      c) {
10     for (auto i = c.begin(); i != c.end(); i) {
11         std::cout << *i;
12         if (++i != c.end()) {std::cout << ' ';}
13     }
14     std::cout << '\n';
15     return bool(std::cout);
16 }
17
18 int main() {
19     std::vector<int> vi{1, 2, 3, 4, 5};
20     std::set<int> si{5, 4, 3, 2, 1};
21     std::set<std::string> ss{"world", "hello"};
22     print_container(vi);
23     print_container(si);
24     print_container(ss);
25 }
```

Variadic Class Template: Integer_sequence

```
1  #include <iostream>
2  #include <cstdlib>
3
4  template <class T, T... Values>
5  class Integer_sequence {
6  public:
7      using value_type = T;
8      using const_iterator = const T*;
9      constexpr std::size_t size() const
10         {return sizeof...(Values);}
11     constexpr T operator[](int i) const {return values_[i];}
12     constexpr const_iterator begin() const
13         {return &values_[0];}
14     constexpr const_iterator end() const
15         {return &values_[size()];}
16 private:
17     static constexpr T values_[sizeof...(Values)] =
18         {Values...};
19 };
20
21 template <class T, T... Values>
22 constexpr T
23 Integer_sequence<T, Values...>::values_[sizeof...(Values)];
24
25 int main() {
26     Integer_sequence<std::size_t, 1, 2, 4, 8> seq;
27     std::cout << seq.size() << '\n' << seq[0] << '\n';
28     for (auto i : seq) {std::cout << i << '\n';}
29 }
```


Variadic Variable Template: `int_array`

```
1  #include <iostream>
2
3  template <int... Args>
4  constexpr int int_array[] = {Args...};
5
6  int main() {
7      for (auto i : int_array<1,2,4,8>) {
8          std::cout << i << '\n';
9      }
10 }
```

Variadic Alias Template: My_tuple

```
1  #include <iostream>
2  #include <string>
3  #include <tuple>
4
5  template <class... Ts>
6  using My_tuple = std::tuple<bool, Ts...>;
7
8  int main() {
9      My_tuple<int, std::string> t(true, 42,
10     "meaning of life");
11     std::cout << std::get<0>(t) << ' '
12     << std::get<1>(t) << ' '
13     << std::get<2>(t) << '\n';
14 }
```

Section 2.5.7

Miscellany

Overload Resolution and Substitution Failure

- when creating candidate set (of functions) for overload resolution, some or all candidates of that set may be result of instantiated templates with template arguments substituted for corresponding template parameters
- process of substituting template arguments for corresponding template parameters can lead to invalid code, known as **substitution failure**
- substitution failure *not treated as error*
- instead, substitution failure simply causes overload to be *removed from candidate set*
- this behavior often referred to by term “*substitution failure is not an error (SFINAE)*”
- if substitution failure were treated as error, this would place heavy burden on programmer to avoid compile errors when using function templates
- SFINAE behavior often exploited in template metaprogramming

Some Kinds of Substitution Failures

- attempting to instantiate pack expansion containing multiple parameter packs of differing lengths
- attempting to create array with element type that is **void**, function type, reference type, or abstract class type
- attempting to create array with size that is zero or negative
- attempting to use type that is not class or enumeration type in qualified name
- attempting to use type in nested name specifier of qualified ID, when type does not contain specified member, or
 - specified member is not type where type is required
 - specified member is not template where template is required
 - specified member is not non-type where non-type is required
- attempting to create pointer to reference type
- attempting to create reference to **void**

Some Kinds of Substitution Failures (Continued)

- attempting to create pointer to member of \mathbb{T} when \mathbb{T} is not class type
- attempting to give invalid type to non-type template parameter
- attempting to perform invalid conversion in either template argument expression, or expression used in function declaration
- attempting to create function type in which parameter has type of **void**, or in which return type is function type or array type
- attempting to create function type in which parameter type or return type is abstract class

std::enable_if and std::enable_if_t

- to make SFINAE more convenient to exploit, class template `std::enable_if` and alias template `std::enable_if_t` are provided
- declaration of class template `enable_if`:

```
template <bool B, class T = void>
struct enable_if;
```
- if B is **true**, class has member type `type` defined as T; otherwise, class has no `type` member
- possible implementation of `enable_if`:

```
1  template <bool B, class T = void>
2  struct enable_if {};
3
4  template <class T>
5  struct enable_if<true, T> {
6      typedef T type;
7  };
```
- declaration of alias template `enable_if_t`:

```
template <bool B, class T = void>
using enable_if_t = typename enable_if<B, T>::type;
```
- if `enable_if_t` is used with its first parameter as **false**, substitution failure will result

SFINAE Example

```
1  #include <iostream>
2  #include <type_traits>
3  #include <cmath>
4
5  // integral version (uses left shift)
6  template <class T>
7  std::enable_if_t<std::is_integral<T>::value, T>
8  mult_by_pow_2(T x, unsigned int n) {
9      std::cerr << "integral version\n";
10     return x << n;
11 }
12
13 // floating-point version (uses pow)
14 template <class T>
15 std::enable_if_t<std::is_floating_point<T>::value, T>
16 mult_by_pow_2(T x, unsigned int n) {
17     std::cerr << "floating-point version\n";
18     return x * std::pow(2.0, n);
19 }
20
21 int main() {
22     std::cout << mult_by_pow_2(2, 2) << '\n';
23     // will call integral version
24     std::cout << mult_by_pow_2(0.5, 4) << '\n';
25     // will call floating-point version
26 }
```


Section 2.5.8

References

- 1 D. Vandevoorde and N. M. Josuttis. *C++ Templates: The Complete Guide*. Addison Wesley, 2002.
- 2 P. Sommerlad. *Variadic and variable templates*. *Overload*, 126:14–17, Apr. 2015.
Available online at <http://accu.org/index.php/journals/2087>.

- 1 Peter Sommerlad. Variadic Templates in C++11/C++14: An Introduction, CppCon, Bellevue, WA, USA, Sept 19–25, 2015.

Section 2.6

Lambda Expressions

Motivation for Lambda Expressions

- functor classes extremely useful, especially for generic programming
- writing definitions of functor classes somewhat tedious, especially if many such classes
- functor classes all have same general structure (i.e., constructor, function-call operator, zero or more data members)
- would be nice if functor could be created without need to explicitly write functor-class definition
- lambda expressions provide compact notation for creating functors
- convenience feature (not fundamentally anything new that can be done with lambda expressions that could not already have been done without them)

Lambda Expressions

- lambda expression consists of:
 - ① introducer: *capture list* in square brackets
 - ② declarator: *parameter list* in parentheses followed by *return type* using trailing return-type syntax
 - ③ *compound statement* in brace brackets
- capture list specifies objects to be captured as data members
- declarator specifies parameter list and return type of function-call operator
- compound statement specifies body of function-call operator
- if no declarator specified, defaults to ()
- if no return type specified, defaults to type of expression in return statement, or void if no return statement
- when evaluated, lambda expression yields object called **closure** (which is essentially a functor)
- examples:

```
[ ] ( double x ) -> int { return floor ( x ); }  
[ ] ( int x , int y ) { return x < y ; }  
[ ] { std :: cout << "Hello , World ! \n " ; }
```

Lambda Expressions (Continued)

- closure object is unnamed (temporary object)
- closure type is unnamed
- **operator** () is always inline
- **operator** () is const member function unless **mutable** keyword used
- if no capture, closure type provides conversion function to pointer to function having same parameter and return types as closure type's function call operator; value returned is address of function that, when invoked, has same effect as invoking closure type's function call operator (function pointer not tied to lifetime of closure object)
- although **operator** () in closure very similar to case of normal functor, not everything same (e.g., **operator** () member in closure type cannot access **this** pointer for closure type)

Hello World Program Revisited

```
1  #include <iostream>
2
3  int main() {
4      []{std::cout << "Hello, World!\n";}();
5  }
```

```
1  #include <iostream>
2
3  struct Hello {
4      void operator()() const {
5          std::cout << "Hello, World!\n";
6      }
7  };
8
9  int main() {
10     Hello hello;
11     hello();
12 }
```


Comparison Functor Example

```
1  #include <iostream>
2  #include <algorithm>
3  #include <cstdlib>
4
5  int main() {
6      std::vector<int> v{-3, 3, 4, 0, -2, -1, 2, 1, -4};
7      std::sort(v.begin(), v.end(),
8              [](int x, int y) {return abs(x) < abs(y);});
9      for (auto x : v) std::cout << x << '\n';
10 }
```

```
1  #include <iostream>
2  #include <algorithm>
3  #include <cstdlib>
4
5  struct abs_less {
6      bool operator()(int x, int y) const
7          {return abs(x) < abs(y);}
8  };
9
10 int main() {
11     std::vector<int> v{-3, 3, 4, 0, -2, -1, 2, 1, -4};
12     std::sort(v.begin(), v.end(), abs_less());
13     for (auto x : v) std::cout << x << '\n';
14 }
```

Capturing Objects

- locals only available if captured; non-locals always available
- can capture by value or by reference
- different locals can be captured differently
- can specify default capture mode
- can explicitly list objects to be captured or not
- might be wise to explicitly list all objects to be captured (when practical) to avoid capturing objects accidentally (e.g., due to typos)
- to capture class members within member function, capture **this**
- capture of **this** must be done by value

- (unary version of) `std::transform` applies given (unary) operator to each element in range specified by pair of iterators and writes result to location specified by another iterator
- definition of `std::transform` would typically resemble:

```
template <class InputIterator, class OutputIterator,  
         class UnaryOperator>  
OutputIterator transform(InputIterator first,  
                        InputIterator last, OutputIterator result,  
                        UnaryOperator op) {  
    while (first != last) {  
        *result = op(*first);  
        ++result;  
        ++first;  
    }  
    return result;  
}
```

Modulus Example

```
1  #include <iostream>
2  #include <vector>
3  #include <algorithm>
4
5  int main() {
6      int m = 2;
7      std::vector<int> v{0, 1, 2, 3};
8      std::transform(v.begin(), v.end(), v.begin(),
9                    [m](int x){return x % m;});
10     for (auto x : v) std::cout << x << '\n';
11 }
```

```
1  #include <iostream>
2  #include <vector>
3  #include <algorithm>
4
5  class mod {
6  public:
7      mod(int m_) : m(m_) {}
8      int operator()(int x) const {return x % m;}
9  private:
10     int m;
11 };
12
13 int main() {
14     int m = 2;
15     std::vector<int> v{0, 1, 2, 3};
16     std::transform(v.begin(), v.end(), v.begin(), mod(m));
17     for (auto x : v) std::cout << x << '\n';
18 }
```

Modulus Example: Without Lambda Expression

```
1  #include <iostream>
2  #include <vector>
3  #include <algorithm>
4
5  class mod {
6  public:
7      mod(int m_) : m(m_) {}
8      int operator()(int x) const {return x % m;}
9  private:
10     int m;
11 };
12
13 int main() {
14     int m = 2;
15     std::vector<int> v{0, 1, 2, 3};
16     std::transform(v.begin(), v.end(), v.begin(), mod(m));
17     for (auto x : v) std::cout << x << '\n';
18 }
```

- approximately 8.5 lines of code to generate functor

Modulus Example: With Lambda Expression

```
1  #include <iostream>
2  #include <vector>
3  #include <algorithm>
4
5  int main() {
6      int m = 2;
7      std::vector<int> v{0, 1, 2, 3};
8      std::transform(v.begin(), v.end(), v.begin(),
9                    [m](int x){return x % m;});
10     for (auto x : v) std::cout << x << '\n';
11 }
```

- m captured by value
- approximately 0.5 lines of code to generate functor

- `std::for_each` applies given function/functor to each element in range specified by pair of iterators
- definition of `std::for_each` would typically resemble:

```
template<class InputIterator, class Function>
Function for_each(InputIterator first,
InputIterator last, Function func) {
    while (first != last) {
        func(*first);
        ++first;
    }
    return move(func);
}
```

Product Example

```
1  #include <iostream>
2  #include <vector>
3  #include <algorithm>
4
5  int main() {
6      std::vector<int> v{2, 3, 4};
7      int prod = 1;
8      std::for_each(v.begin(), v.end(),
9                  [&prod](int x)->void{prod *= x;});
10     std::cout << prod << '\n';
11 }
```

```
1  #include <iostream>
2  #include <vector>
3  #include <algorithm>
4
5  class cum_prod {
6  public:
7      cum_prod(int& prod_) : prod(prod_) {}
8      void operator()(int x) const {prod *= x;}
9  private:
10     int& prod;
11 };
12
13 int main() {
14     std::vector<int> v{2, 3, 4};
15     int prod = 1;
16     std::for_each(v.begin(), v.end(), cum_prod(prod));
17     std::cout << prod << '\n';
18 }
```


Product Example: Without Lambda Expression

```
1  #include <iostream>
2  #include <vector>
3  #include <algorithm>
4
5  class cum_prod {
6  public:
7      cum_prod(int& prod_) : prod(prod_) {}
8      void operator()(int x) const {prod *= x;}
9  private:
10     int& prod;
11 };
12
13 int main() {
14     std::vector<int> v{2, 3, 4};
15     int prod = 1;
16     std::for_each(v.begin(), v.end(), cum_prod(prod));
17     std::cout << prod << '\n';
18 }
```

- approximately 8.5 lines of code to generate functor

Product Example: With Lambda Expression

```
1  #include <iostream>
2  #include <vector>
3  #include <algorithm>
4
5  int main() {
6      std::vector<int> v{2, 3, 4};
7      int prod = 1;
8      std::for_each(v.begin(), v.end(),
9          [&prod](int x)->void{prod *= x;});
10     std::cout << prod << '\n';
11 }
```

- prod captured by reference
- approximately 1 line of code to generate functor

More Variations on Capture

```
double a = 2.14;
```

```
double b = 3.14;
```

```
double c = 42.0;
```

```
// capture all objects by reference (i.e., a, b, and c)  
[&] (double x, double y) {return a * x + b * y + c;}
```

```
// capture all objects by value (i.e., a, b, and c)  
[=] (double x, double y) {return a * x + b * y + c;}
```

```
// capture all objects by value, except a  
// which is captured by reference  
[=, &a] (double x, double y) {return a * x + b * y + c;}
```

```
// capture all objects by reference, except a  
// which is captured by value  
[&, a] (double x, double y) {return a * x + b * y + c;}
```

Generalized Lambda Capture

- can specify name for captured object in closure type

```
int a = 1;  
auto f = [x = a]() {return x;};
```

- can capture result of expression (e.g., to perform move instead of copy or to add arbitrary new state to closure type)

```
std::vector<int> v(1000, 1);  
auto f = [v = std::move(v)]() ->  
    const std::vector<int>& {return v;};
```

Generalized Lambda Capture Example

```
1  #include <iostream>
2
3  int main() {
4      int x = 0;
5      int y = 1;
6      auto f = [&count = x, inc = y + 1]() {
7          return count += inc;
8      };
9      std::cout << f() << ' ';
10     std::cout << f() << '\n';
11 }
12
13 // output: 2 4
```

Generic Lambda Expressions

- can allow compiler to deduce type of lambda function parameters
- generates closure type with templated function-call operator
- one template type parameter for each occurrence of **auto** in lambda expression's parameter declaration clause

Generic Lambda Expression Example [Generic]

```
1  #include <iostream>
2  #include <complex>
3  #include <string>
4
5  int main() {
6      using namespace std::literals;
7      auto add = [](auto x, auto y) {return x + y;};
8      std::cout << add(1, 2) << ' ' << add(1.0, 2.0) << ' '
9          << add(1.0, 2.0i) << ' ' << add("Jell"s, "o"s) << '\n';
10 }
```

```
1  #include <iostream>
2  #include <complex>
3  #include <string>
4
5  struct Add {
6      template <class T, class U>
7      auto operator()(T x, U y) {return x + y;};
8  };
9
10 int main() {
11     using namespace std::literals;
12     Add add;
13     std::cout << add(1, 2) << ' ' << add(1.0, 2.0) << ' '
14         << add(1.0, 2.0i) << ' ' << add("Jell"s, "o"s) << '\n';
15 }
```

Generic Lambda Expression Example [Convenience]

```
1 #include <iostream>
2 #include <vector>
3 #include <algorithm>
4
5 int main() {
6     std::vector<int> v{0, 1, 2, 3, 4, 5, 6, 7};
7     // sort elements of vector in descending order
8     std::sort(v.begin(), v.end(),
9         [](auto i, auto j) {return i > j;});
10    std::for_each(v.begin(), v.end(),
11        [](auto i) {std::cout << i << '\n';});
12 }
```


Dealing With Unnamed Types

- fact that closure types unnamed causes complications when need arises to refer to closure type
- helpful language features: **auto**, **decltype**
- helpful library features: `std::function`
- closures can be stored using **auto** or `std::function`
- closures that do not capture can be “stored” by assigning to function pointer

Using `auto`, `decltype`, and `std::function`

```
1  #include <iostream>
2  #include <functional>
3
4  std::function<double(double)> linear(double a, double b) {
5      return [=](double x){return a * x + b;};
6  }
7
8  int main() {
9      // type of f is std::function<double(double)>
10     auto f = linear(2.0, -1.0);
11     // g has closure type
12     auto g = [](double x){return 2.0 * x - 1.0;};
13     double (*u)(double) = [](double x){return 2.0 * x - 1.0;};
14     // h has same type as g
15     decltype(g) h = g;
16     for (double x = 0.0; x < 10.0; x += 1.0) {
17         std::cout << x << ' ' << f(x) << ' ' << g(x) <<
18             ' ' << h(x) << (*u)(x) << '\n';
19     }
20 }
```

- applying function-call operator to `f` much slower than in case of `g` and `h`
- when `std::function` used, inlining of called function probably not possible
- when functor used directly (via function-call operator) inlining is very likely
- prefer `auto` over `std::function` for storing closures

operator () as Non-const Member

```
1  #include <iostream>
2
3  int main()
4  {
5      int count = 5;
6      // Must use mutable in order to be able to
7      // modify count member.
8      auto get_count = [count]() mutable -> int {
9          return count++;
10     };
11
12     int c;
13     while ((c = get_count()) < 10) {
14         std::cout << c << '\n';
15     }
16 }
```

- **operator ()** is declared as const member function unless **mutable** keyword used
- const member function cannot change (non-static) data members

Comparison Functors for Containers

```
1  #include <iostream>
2  #include <vector>
3  #include <set>
4
5  int main() {
6      // The following two lines are the only important ones:
7      auto cmp = [](int* x, int* y){return *x < *y;};
8      std::set<int*, decltype(cmp)> s(cmp);
9
10     // Just for something to do:
11     // Print the elements of v in sorted order with
12     // duplicates removed.
13     std::vector<int> v = {4, 1, 3, 2, 1, 1, 1, 1};
14     for (auto& x : v) {
15         s.insert(&x);
16     }
17     for (auto x : s) {
18         std::cout << *x << '\n';
19     }
20 }
```

- note that `s` is not default constructed
- since closure types not default constructible, following would fail:
`std::set<int*, decltype(cmp)> s;`
- note use of `decltype` in order to specify type of functor

What Could Possibly Go Wrong?

```
1  #include <iostream>
2  #include <vector>
3  #include <functional>
4
5  std::vector<int> vec{2000, 4000, 6000, 8000, 10000};
6  std::function<int(int)> func;
7
8  void do_stuff()
9  {
10     int modulus = 10000;
11     func = [&](int x){return x % modulus;};
12     for (auto x : vec) {
13         std::cout << func(x) << '\n';
14     }
15 }
16
17 int main()
18 {
19     do_stuff();
20     for (auto x : vec) {
21         std::cout << func(x) << '\n';
22     }
23 }
```

- above code has very serious bug; what is it?

Dangling References

- if some objects captured by reference, closure can hold dangling references
- responsibility of programmer to avoid such problems
- if will not cause performance issues, may be advisable to capture by value (to avoid problem of dangling references)
- dangling-reference example:

```
1  #include <iostream>
2  #include <functional>
3
4  std::function<double(double)> linear(double a, double b) {
5      return [&](double x){return a * x + b;};
6  }
7
8  int main() {
9      auto f = linear(2.0, -1.0);
10     // bad things will happen here
11     std::cout << f(1.0) << '\n';
12 }
```

Triangle Scan Conversion

- in SPLEL software, triangle scan conversion performed by `scan_triangle` template function

- declaration:

```
template <class T, class F>
void scan_triangle(T a_x, T a_y, T b_x, T b_y,
    T c_x, T c_y, unsigned mask, F scan_line);
```

- `scan_line` is functor to handle single horizontal scan within triangle
- `scan_line` has signature:

```
void scan_line(T y, T x_min, T x_max,
    unsigned left_mask, unsigned right_mask,
    unsigned mid_mask);
```

Section 2.6.1

References

- 1 Herb Sutter. Lambdas, Lambdas Everywhere, Professional Developers Conference (PDC), Redmond, WA, USA, October 27–29, 2010.
- 2 Herb Sutter. C++0x Lambda Functions, Northwest C++ Users' Group (NWCPP), Redmond, WA, USA, May 18, 2011.

Section 2.7

Classes and Inheritance

Section 2.7.1

Derived Classes and Class Hierarchies

Derived Classes

- sometimes, want to express commonality between classes
- want to create new class from existing class by adding new members or replacing (i.e., hiding/overriding) existing members
- can be achieved through language feature known as *inheritance*
- generate new class with all members of already existing class, excluding special member functions (i.e., constructors, assignment operators, and destructor)
- new class called **derived class** and original class called **base class**
- derived class said to **inherit** from base class
- can add new members (not in base class) to derived class
- can hide or override member functions from base class with new version
- syntax for specifying derived class:

```
class derived_class : base_class_specifiers
```
- *derived_class* is name of derived class; *base_class_specifiers* provide base-class information

Derived Classes (Continued)

- can more clearly express intent by explicitly identifying relationship between classes
- can facilitate code reuse by leverage existing code
- interface inheritance: allow different derived classes to be used interchangeably through interface provided by common base class
- implementation inheritance: save implementation effort by sharing capabilities provided by base class

Person Class

```
1  #include <string>
2
3  class Person {
4  public:
5      Person(const std::string& family_name,
6            const std::string& given_name) :
7          family_name_(family_name), given_name_(given_name) {}
8      std::string family_name() const {return family_name_;}
9      std::string given_name() const {return given_name_;}
10     std::string full_name() const
11         {return family_name_ + ", " + given_name_;}
12     // ...
13 private:
14     std::string family_name_;
15     std::string given_name_;
16 };
```

Student Class Without Inheritance

```
1  #include <string>
2
3  class Student {
4  public:
5      Student(const std::string& family_name,
6              const std::string& given_name) :
7          family_name_(family_name), given_name_(given_name) {}
8          // NEW
9      std::string family_name() const {return family_name_;}
10     std::string given_name() const {return given_name_;}
11     std::string full_name() const
12         {return family_name_ + ", " + given_name_;}
13     std::string student_id() {return student_id_;} // NEW
14 private:
15     std::string family_name_;
16     std::string given_name_;
17     std::string student_id_; // NEW
18 };
```

Student Class With Inheritance

```
1 // include definition of Person class here
2
3 class Student : public Person {
4 public:
5     Student(const std::string& family_name,
6           const std::string& given_name,
7           const std::string& student_id) :
8         Person(family_name, given_name),
9         student_id_(student_id) {}
10    std::string student_id() {return student_id_;}
11 private:
12    std::string student_id_;
13 };
```


Complete Inheritance Example

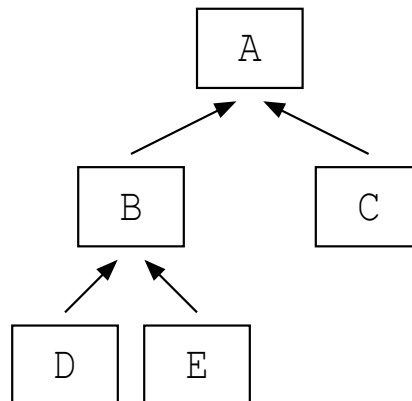
```
1  #include <string>
2
3  class Person {
4  public:
5      Person(const std::string& family_name,
6             const std::string& given_name) :
7          family_name_(family_name), given_name_(given_name) {}
8      std::string family_name() const {return family_name_;}
9      std::string given_name() const {return given_name_;}
10     std::string full_name() const
11         {return family_name_ + ", " + given_name_;}
12     // ... (including virtual destructor)
13 private:
14     std::string family_name_;
15     std::string given_name_;
16 };
17
18 class Student : public Person {
19 public:
20     Student(const std::string& family_name,
21            const std::string& given_name,
22            const std::string& student_id) :
23         Person(family_name, given_name),
24         student_id_(student_id) {}
25     std::string student_id() {return student_id_;}
26 private:
27     std::string student_id_;
28 };
```

Class Hierarchies

- inheritance relationships between classes form what is called **class hierarchy**
- often class hierarchy represented by directed (acyclic) graph, where nodes correspond to classes and edges correspond to inheritance relationships
- class definitions:

```
class A { /* ... */ };  
class B : public A { /* ... */ };  
class C : public A { /* ... */ };  
class D : public B { /* ... */ };  
class E : public B { /* ... */ };
```

- inheritance diagram:

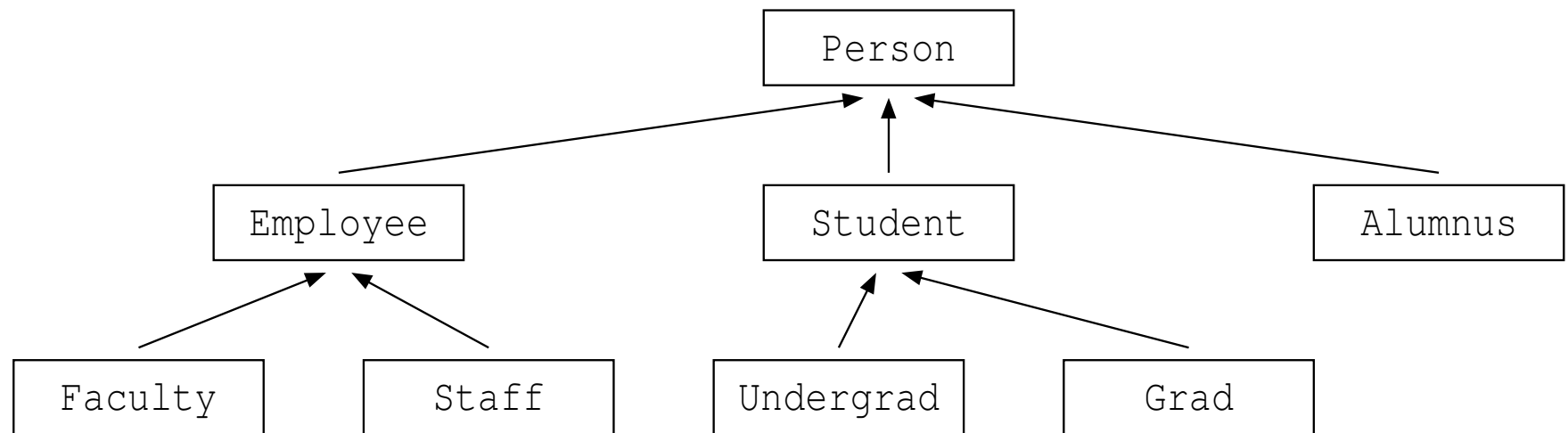


Class Hierarchy Example

- class definitions:

```
class Person { /* ... */ };  
class Employee : public Person { /* ... */ };  
class Student : public Person { /* ... */ };  
class Alumnus : public Person { /* ... */ };  
class Faculty : public Employee { /* ... */ };  
class Staff : public Employee { /* ... */ };  
class Grad : public Student { /* ... */ };  
class Undergrad : public Student { /* ... */ };
```

- inheritance diagram:



- each of Employee, Student, and Alumnus is a Person; each of Faculty and Staff is an Employee; each of Undergrad and Grad is a Student

Member Access Specifiers: **protected**

- earlier, introduced **public** and **private** access specifiers for class members
- in context of inheritance, another access specifier becomes relevant, namely, **protected**
- member declared in protected section of class can only be accessed by
 - member functions and friends of that class; and
 - by member functions and friends of *derived classes*
- protected members used to provide developers of derived classes access to some inner workings of base class without exposing such inner workings to everyone
- usually, bad idea to use protected access for data members (for similar reasons that using public access for data members is usually bad)
- protected access usually employed for function members

Types of Inheritance

- three types of inheritance with respect to access protection: public, protected, and private
- these three types of inheritance differ in terms of accessibility, in derived class, of members inherited from base class
- private parts of base class are always inaccessible in derived class, regardless of whether public, protected, or private inheritance used
- if this were not case, all access protection could simply be bypassed by using inheritance
- access specifiers for members accessible in derived class chosen as follows:

Access Specifier in Base Class	Access Specifier in Derived Class		
	Public Inheritance	Protected Inheritance	Private Inheritance
public	public	protected	private
protected	protected	protected	private

Types of Inheritance (Continued)

- for struct, defaults to public inheritance
- for class, defaults to private inheritance
- public and protected/private inheritance have different use cases, as we will see later

Inheritance and Member Access Example

```
1  class Base {
2  public:
3      void f();
4  protected:
5      void g();
6  private:
7      int x;
8  };
9
10 class Derived_1 : public Base {
11     // f is public
12     // g is protected
13     // x is not accessible from Derived_1
14 };
15
16 class Derived_2 : protected Base {
17     // f is protected
18     // g is protected
19     // x is not accessible from Derived_2
20 };
21
22 class Derived_3 : private Base {
23     // f is private
24     // g is private
25     // x is not accessible from Derived_3
26 };
```

Public Inheritance Example

```
1  class Base {
2  public:
3      void func_1 ();
4  protected:
5      void func_2 ();
6  private:
7      int x_;
8  };
9
10 class Derived : public Base {
11 public:
12     void func_3 () {
13         func_1 (); // OK
14         func_2 (); // OK
15         x_ = 0; // ERROR: inaccessible
16     }
17 };
18
19 struct Widget : public Derived {
20     void func_4 () { func_2 (); } // OK
21 };
22
23 int main () {
24     Derived d;
25     d.func_1 (); // OK
26     d.func_2 (); // ERROR: inaccessible
27     d.x_ = 0; // ERROR: inaccessible
28 }
```


Protected Inheritance Example

```
1  class Base {
2  public:
3      void func_1 ();
4  protected:
5      void func_2 ();
6  private:
7      int x_;
8  };
9
10 class Derived : protected Base {
11 public:
12     void func_3 () {
13         func_1 (); // OK
14         func_2 (); // OK
15         x_ = 0; // ERROR: inaccessible
16     }
17 };
18
19 struct Widget : public Derived {
20     void func_4 () { func_2 (); } // OK
21 };
22
23 int main () {
24     Derived d; // OK: defaulted constructor is public
25     d.func_1 (); // ERROR: inaccessible
26     d.func_2 (); // ERROR: inaccessible
27     d.x_ = 0; // ERROR: inaccessible
28 }
```

Private Inheritance Example

```
1  class Base {
2  public:
3      void func_1 ();
4  protected:
5      void func_2 ();
6  private:
7      int x_;
8  };
9
10 class Derived : private Base {
11 public:
12     void func_3 () {
13         func_1 (); // OK
14         func_2 (); // OK
15         x_ = 0; // ERROR: inaccessible
16     }
17 };
18
19 struct Widget : public Derived {
20     void func_4 () { func_2 (); } // ERROR: inaccessible
21 };
22
23 int main () {
24     Derived d; // OK: defaulted constructor is public
25     d.func_1 (); // ERROR: inaccessible
26     d.func_2 (); // ERROR: inaccessible
27     d.x_ = 0; // ERROR: inaccessible
28 }
```

Public Inheritance

- public inheritance is inheritance in traditional object-oriented programming sense
- public inheritance models an *is-a* relationship (i.e., derived class object is a base class object)
- most common form of inheritance
- inheritance relationship visible to all code

Public Inheritance Example

```
1  #include <string>
2
3  class Person {
4  public:
5      Person(const std::string& family_name, const std::string&
6          given_name) : family_name_(family_name),
7          given_name_(given_name) {}
8      std::string family_name() const
9          {return family_name_;}
10     std::string given_name() const
11         {return given_name_;}
12     std::string full_name() const
13         {return family_name_ + ", " + given_name_;}
14 private:
15     std::string family_name_;
16     std::string given_name_;
17 };
18
19 class Student : public Person {
20 public:
21     Student(const std::string& family_name, const std::string&
22         given_name, const std::string& student_id) :
23         Person(family_name, given_name), student_id_(student_id)
24     std::string student_id()
25         {return student_id_;}
26 private:
27     std::string student_id_;
28 };
```

Protected and Private Inheritance

- protected and private inheritance not inheritance in traditional object-oriented programming sense (i.e., no is-a relationship)
- form of implementation inheritance
- *implemented-in-terms-of* relationship (i.e., derived class object implemented in terms of a base class object)
- in case of protected inheritance, inheritance relationship only seen by derived classes and their friends and class itself and its friends
- in case of private inheritance, inheritance relationship only seen by class itself and its friends (not derived classes and their friends)
- except in special circumstances, normally bad idea to use inheritance for composition
- one good use case for private/protected inheritance is in policy-based design, which exploits empty base optimization (EBO)

Policy-Based Design Example: Inefficient Memory Usage

```
1  #include <mutex>
2
3  class ThreadSafePolicy {
4  public:
5      void lock() {mutex_.lock();}
6      void unlock() {mutex_.unlock();}
7  private:
8      std::mutex mutex_;
9  };
10
11 class ThreadUnsafePolicy {
12 public:
13     void lock() {} // no-op
14     void unlock() {} // no-op
15 };
16
17 template<class ThreadSafetyPolicy>
18 class Widget {
19     ThreadSafetyPolicy policy_;
20     // ...
21 };
22
23 int main() {
24     Widget<ThreadUnsafePolicy> w;
25     // w.policy_ has no data members, but
26     // sizeof(w.policy_) >= 1
27     // inefficient use of memory
28 }
```

Policy-Based Design Example: Private Inheritance and EBO

```
1  #include <mutex>
2
3  class ThreadSafePolicy {
4  public:
5      void lock() {mutex_.lock();}
6      void unlock() {mutex_.unlock();}
7  private:
8      std::mutex mutex_;
9  };
10
11 class ThreadUnsafePolicy {
12 public:
13     void lock() {} // no-op
14     void unlock() {} // no-op
15 };
16
17 template<class ThreadSafetyPolicy>
18 class Widget : ThreadSafetyPolicy {
19     // ...
20 };
21
22 int main() {
23     Widget<ThreadUnsafePolicy> w;
24     // empty-base optimization (EBO) can be applied
25     // no memory overhead for no-op thread-safety policy
26 }
```

Inheritance and Constructors

- by default, constructors not inherited
- often, derived class introduces new data members not in base class
- since base-class constructors cannot initialize derived-class data members, inheriting constructors from base class by default would be bad idea (e.g., could lead to uninitialized data members)
- in some cases, however, base-class constructors may be sufficient to initialize derived-class objects
- in such cases, can inherit all non-special base-class constructors with **using** statement
- special constructors (i.e., default, copy, and move constructors) cannot be inherited
- constructors to be inherited with **using** statement may still be hidden by constructors in derived class

Inheriting Constructors Example

```
1  class Base {
2  public:
3      Base() : i_(0.0), j_(0) {}
4      Base(int i) : i_(i), j_(0) {}
5      Base(int i, int j) : i_(i), j_(j) {}
6      // ... (other non-constructor members)
7  private:
8      int i_, j_;
9  };
10
11 class Derived : public Base {
12 public:
13     // inherit non-special constructors from Base
14     // (default constructor not inherited)
15     using Base::Base;
16     // default constructor is implicitly declared and
17     // not inherited
18 };
19
20 int main() {
21     Derived a;
22     // invokes non-inherited Derived::Derived()
23     Derived b(42, 42);
24     // invokes inherited Base::Base(int, int)
25 }
```

Inheriting Constructors Example

```
1  class Base {
2  public:
3      Base() : i_(0), j_(0), k_(0) {}
4      Base(int i, int j) : i_(i), j_(j), k_(0) {}
5      Base(int i, int j, int k) : i_(i), j_(j), k_(k) {}
6      // ... (other non-constructor members)
7  private:
8      int i_, j_, k_;
9  };
10
11 class Derived : public Base {
12 public:
13     // inherit non-special constructors from Base
14     // (default constructor not inherited)
15     using Base::Base;
16     // following constructor hides inherited constructor
17     Derived(int i, int j, int k) : Base(-i, -j, -k) {}
18     // no implicitly-generated default constructor
19 };
20
21 int main() {
22     Derived b(1, 2);
23     // invokes inherited Base::Base(int, int)
24     Derived c(1, 2, 3);
25     // invokes Derived::Derived(int, int, int)
26     // following would produce compile-time error:
27     // Derived a; // ERROR: no default constructor
28 }
```

Inheritance, Assignment Operators, and Destructors

- by default, assignment operators not inherited (for similar reasons as in case of constructors)
- can inherit all non-special base-class assignment operators with **using** statement
- copy and move assignment operators cannot be inherited
- assignment operators to be inherited with **using** statement may still be hidden by assignment operators in derived class
- cannot inherit destructor

Inheriting Assignment Operators Example

```
1  class Base {
2  public:
3      explicit Base(int i) : i_(i) {}
4      Base& operator=(int i) {
5          i_ = i;
6          return *this;
7      }
8      // ...
9  private:
10     int i_;
11 };
12
13 class Derived : public Base {
14 public:
15     // inherit non-special constructors
16     using Base::Base;
17     // inherit non-special assignment operators
18     using Base::operator==;
19     // ...
20 };
21
22 int main() {
23     Derived d(0);
24     // invokes inherited Base::Base(int)
25     d = 42;
26     // invokes inherited Base::operator=(int)
27 }
```

Construction and Destruction Order

- during construction of object, all of its base class objects constructed first
- order of construction:
 - ① *base class objects* as listed in type definition left to right
 - ② data members as listed in type definition top to bottom
 - ③ constructor body
- order of destruction is exact reverse of order of construction, namely:
 - ① destructor body
 - ② data members as listed in type definition bottom to top
 - ③ *base class objects* as listed in type definition right to left

Order of Construction

```
1  #include <vector>
2  #include <string>
3
4  class Base {
5  public:
6      Base(int n) : v_(n, 0) {}
7      // ...
8  private:
9      std::vector<char> v_;
10 };
11
12 class Derived : public Base {
13 public:
14     Derived(const std::string& s) : Base(1024), s_(s)
15         { i_ = 0; }
16     // ...
17 private:
18     std::string s_;
19     int i_;
20 };
21
22 int main() {
23     Derived d("hello");
24 }
```

- construction order for Derived constructor: 1) Base class object, 2) data member s_, 3) Derived constructor body (initializes data member i_)

Hiding Base-Class Member Functions in Derived Class

- can provide new versions of member functions in derived class to hide original functions in base class

```
1  #include <iostream>
2
3  class Fruit {
4  public:
5      void print() const {std::cout << "fruit\n";}
6  };
7
8  class Apple : public Fruit {
9  public:
10     void print() const {std::cout << "apple\n";}
11 };
12
13 class Banana : public Fruit {
14 public:
15     void print() const {std::cout << "banana\n";}
16 };
17
18 int main() {
19     Fruit f;
20     Apple a;
21     Banana b;
22     f.print(); // calls Fruit::print
23     a.print(); // calls Apple::print
24     b.print(); // calls Banana::print
25 }
```

Upcasting

- derived-class object always has base-class subobject
- given reference or pointer to derived-class object, may want to find reference or pointer to corresponding base-class object
- **upcasting**: converting derived-class pointer or reference to base-class pointer or reference
- upcasting allows us to treat derived-class object as base-class object
- upcasting always safe in sense that cannot result in incorrect type (since every derived-class object is also a base-class object)
- in case of public inheritance, can upcast without explicit type-cast operator
- in case of protected or private inheritance, cannot upcast, except with C-style cast (which also can bypass access protection)
- example:

```
class Base { /* ... */ };  
class Derived : public Base { /* ... */ };  
void func() {  
    Derived d;  
    Base* bp = &d;  
}
```


Downcasting

- **downcasting**: converting base-class pointer or reference to derived-class pointer or reference
- downcasting allows us to force base-class object to be treated as derived-class object
- downcasting is not always safe (since not every base-class object is necessarily also derived-class object)
- must only downcast when known that object actually has derived type
- downcasting always requires explicit cast (**static_cast** for non-polymorphic case, **dynamic_cast** for polymorphic case, C-style cast for either case)
- example:

```
class Base { /* ... (nonpolymorphic) */ };  
class Derived : public Base { /* ... */ };  
void func() {  
    Derived d;  
    Base* bp = &d;  
    Derived* dp = static_cast<Derived*>(bp);  
}
```

Upcasting/Downcasting Example

```
1  class Base { /* ... (nonpolymorphic) */ };
2
3  class Derived : public Base { /* ... */ };
4
5  int main() {
6      Base b;
7      Derived d;
8      Base* bp = nullptr;
9      Derived* dp = nullptr;
10     bp = &d;
11     // OK: upcast does not require explicit cast
12     dp = bp;
13     // ERROR: downcast requires explicit cast
14     dp = static_cast<Derived*>(bp);
15     // OK: downcast with explicit cast and
16     // pointer (bp) refers to Derived object
17     Base& br = d;
18     // OK: upcast does not require explicit cast
19     Derived& dr1 = *bp;
20     // ERROR: downcast requires explicit cast
21     Derived& dr2 = *static_cast<Derived*>(bp);
22     // OK: downcast with explicit cast and
23     // object (*bp) is of Derived type
24     dp = static_cast<Derived*>(&b);
25     // BUG: pointer (&b) does not refer to Derived object
26 }
```

Upcasting Example

```
1  class Base { /* ... */ };
2
3  class Derived : public Base { /* ... */ };
4
5  void func_1 (Base& b) { /* ... */ }
6
7  void func_2 (Base* b) { /* ... */ }
8
9  int main () {
10     Base b;
11     Derived d;
12     func_1 (b);
13     func_1 (d); // OK: Derived& upcast to Base&
14     func_2 (&b);
15     func_2 (&d); // OK: Derived* upcast to Base*
16 }
```

Nonpolymorphic Behavior

```
1  #include <iostream>
2  #include <string>
3
4  class Person {
5  public:
6      Person(const std::string& family, const std::string& given) :
7          family_(family), given_(given) {}
8      void print() const {std::cout << "person: " << family_ << ',' << given_ << '\n';}
9  protected:
10     std::string family_; // family name
11     std::string given_; // given name
12 };
13
14 class Student : public Person {
15 public:
16     Student(const std::string& family, const std::string& given,
17             const std::string& id) : Person(family, given), id_(id) {}
18     void print() const {
19         std::cout << "student: " << family_ << ',' << given_ << ',' << id_ << '\n';
20     }
21 private:
22     std::string id_; // student ID
23 };
24
25 void processPerson(const Person& p) {
26     p.print(); // always calls Person::print
27     // ...
28 }
29
30 int main() {
31     Person p("Ritchie", "Dennis");
32     Student s("Doe", "John", "12345678");
33     processPerson(p); // invokes Person::print
34     processPerson(s); // invokes Person::print
35 }
```

- would be nice if `processPerson` called version of `print` that corresponds to *actual* type of object referenced by function parameter `p`

Slicing

- **slicing**: copying or moving object of derived class to object of base class (e.g., during construction or assignment), losing part of information in so doing
- example:

```
1  class Base {
2      // ...
3      int x_;
4  };
5
6  class Derived : public Base {
7      // ...
8      int y_;
9  };
10
11 int main() {
12     Derived d1, d2;
13     Base b = d1;
14     // slicing occurs
15     Base& r = d1;
16     r = d2;
17     // more treacherous case of slicing
18     // slicing occurs
19     // d1 now contains mixture of d1 and d2
20     // (i.e., base part of d2 and derived part of d1)
21 }
```

Inheritance and Overloading

- functions do not overload across scopes
- can employ **using** statement to bring base members into scope for overloading

Inheritance and Overloading Example

```
1  #include<iostream>
2
3  class Base {
4  public:
5      double f(double d) const {return d;}
6      // ...
7  };
8
9  class Derived : public Base {
10 public:
11     int f(int i) const {return i;}
12     // ...
13 };
14
15 int main()
16 {
17     Derived d;
18     std::cout << d.f(0) << '\n';
19     // calls Derived::f(int) const
20     std::cout << d.f(0.5) << '\n';
21     // calls Derived::f(int) const; probably not intended
22     Derived* dp = &d;
23     std::cout << dp->f(0) << '\n';
24     // calls Derived::f(int) const
25     std::cout << dp->f(0.5) << '\n';
26     // calls Derived::f(int) const; probably not intended
27 }
```

Using Base Members Example

```
1  #include<iostream>
2
3  class Base {
4  public:
5      double f(double d) const {return d;}
6      // ...
7  };
8
9  class Derived : public Base {
10 public:
11     using Base::f; // bring Base::f into scope
12     int f(int i) const {return i;}
13     // ...
14 };
15
16 int main()
17 {
18     Derived d;
19     std::cout << d.f(0) << '\n';
20     // calls Derived::f(int) const
21     std::cout << d.f(0.5) << '\n';
22     // calls Base::f(double) const
23     Derived* dp = &d;
24     std::cout << dp->f(0) << '\n';
25     // calls Derived::f(int) const
26     std::cout << dp->f(0.5) << '\n';
27     // calls Base::f(double) const
28 }
```


Section 2.7.2

Virtual Functions and Run-Time Polymorphism

Run-Time Polymorphism

- **polymorphism** is characteristic of being able to assign different meaning to something in different contexts
- polymorphism that occurs at run time called **run-time polymorphism** (also known as **dynamic polymorphism**)
- in context of inheritance, key type of run-time polymorphism is polymorphic function call (also known as dynamic dispatch)
- when inheritance relationship exists between two classes, type of reference or pointer to object may not correspond to actual dynamic (i.e., run-time) type of object referenced by reference or pointer
- that is, reference or pointer to type T may, in fact, refer to object of type D , where D is either directly or indirectly derived from T
- when calling member function through pointer or reference, may want actual function invoked to be determined by *dynamic* type of object referenced by pointer or reference
- function call with this property said to be **polymorphic**

Virtual Functions

- in context of class hierarchies, polymorphic function calls achieved through use of virtual functions
- **virtual function** is member function with polymorphic behavior
- when call made to virtual function through reference or pointer, actual function invoked will be determined by *dynamic* type of referenced object
- to make member function virtual, add keyword **virtual** to function declaration
- example:

```
class Base {  
public:  
    virtual void func(); // virtual function  
    // ...  
};
```

Virtual Functions (Continued)

- once function made virtual, it will *automatically* be virtual in all derived classes, regardless of whether **virtual** keyword is used in derived classes
- therefore, not necessary to repeat **virtual** qualifier in derived classes (and perhaps preferable not to do so)
- virtual function must be defined in class where first declared unless pure virtual function (to be discussed shortly)
- derived class inherits definition of each virtual function from its base class, but may override each virtual function with new definition
- function in derived class with same name and same set of argument types as virtual function in base class overrides base class version of virtual function

Virtual Function Example

```
1  #include <iostream>
2  #include <string>
3
4  class Person {
5  public:
6      Person(const std::string& family, const std::string& given) :
7          family_(family), given_(given) {}
8      virtual void print() const
9          {std::cout << "person: " << family_ << ',' << given_ << '\n';}
10 protected:
11     std::string family_; // family name
12     std::string given_; // given name
13 };
14
15 class Student : public Person {
16 public:
17     Student(const std::string& family, const std::string& given,
18             const std::string& id) : Person(family, given), id_(id) {}
19     void print() const {
20         std::cout << "student: " << family_ << ',' << given_ << ',' << id_ << '\n';
21     }
22 private:
23     std::string id_; // student ID
24 };
25
26 void processPerson(const Person& p) {
27     p.print(); // polymorphic function call
28     // ...
29 }
30
31 int main() {
32     Person p("Ritchie", "Dennis");
33     Student s("Doe", "John", "12345678");
34     processPerson(p); // invokes Person::print
35     processPerson(s); // invokes Student::print
36 }
```

Override Control: The `override` Qualifier

- when looking at code for derived class, often not possible to determine if member function intended to override virtual function in base class (or one of its base classes)
- can sometimes lead to bugs where programmer expects member function to override virtual function when function not virtual
- **`override`** qualifier used to indicate that member function is expected to override virtual function in parent class; must come at end of function declaration
- example:

```
class Person {
public:
    virtual void print() const;
    // ...
};

class Employee : public Person {
public:
    void print() const override; // must be virtual
    // ...
};
```

Override Control: The `final` Qualifier

- sometimes, may want to prevent any further overriding of virtual function in any subsequent derived classes
- adding `final` qualifier to declaration of virtual function prevents function from being overridden in any subsequent derived classes
- preventing further overriding can sometimes allow for better optimization by compiler (e.g., via devirtualization)
- example:

```
class A {
public:
    virtual void doStuff();
    // ...
};

class B : public A {
public:
    void doStuff() final; // prevent further overriding
    // ...
};

class C : public B {
public:
    void doStuff(); // ERROR: cannot override
    // ...
};
```

final Qualifier Example

```
1  class Worker {
2  public:
3      virtual void prepareEnvelope ();
4      // ...
5  };
6
7  class SpecialWorker : public Worker {
8  public:
9      // prevent overriding function responsible for
10     // overall envelope preparation process
11     // but allow functions for individual steps in
12     // process to be overridden
13     void prepareEnvelope () final {
14         stuffEnvelope (); // step 1
15         lickEnvelope (); // step 2
16         sealEnvelope (); // step 3
17     }
18     virtual void stuffEnvelope ();
19     virtual void lickEnvelope ();
20     virtual void sealEnvelope ();
21     // ...
22 };
```


Constructors, Destructors, and Virtual Functions

- except in very rare cases, destructors in class hierarchy need to be virtual
- otherwise, invoking destructor through base-class pointer/reference would only destroy base-class part of object, leaving remainder of derived-class object untouched
- normally, bad idea to call virtual function inside constructor or destructor
- dynamic type of object changes during construction and changes again during destruction
- final overrider of virtual function will change depending where in hierarchy virtual function call is made
- when constructor/destructor being executed, object is of exactly that type, never type derived from it
- although semantics of virtual function calls during construction and destruction well defined, easy to write code where actual overrider not what expected (and might even be pure virtual)

Problematic Code with Non-Virtual Destructor

```
1  class Base {
2  public:
3      Base() {}
4      ~Base() {} // non-virtual destructor
5      // ...
6  };
7
8  class Derived : public Base {
9  public:
10     Derived() : buffer_(new char[10'000]) {}
11     ~Derived() {delete[] buffer_;}
12     // ...
13 private:
14     char* buffer_;
15 };
16
17 void process(Base* bp) {
18     // ...
19     delete bp; // always invokes only Base::~~Base
20 }
21
22 int main() {
23     process(new Base);
24     process(new Derived); // leaks memory
25 }
```

Corrected Code with Virtual Destructor

```
1  class Base {
2  public:
3      Base() {}
4      virtual ~Base() {} // virtual destructor
5      // ...
6  };
7
8  class Derived : public Base {
9  public:
10     Derived() : buffer_(new char[10'000]) {}
11     ~Derived() {delete[] buffer_;}
12     // ...
13 private:
14     char* buffer_;
15 };
16
17 void process(Base* bp) {
18     // ...
19     delete bp; // invokes destructor polymorphically
20 }
21
22 int main() {
23     process(new Base);
24     process(new Derived);
25 }
```

Preventing Creation of Derived Classes

- in some situations, may want to prevent deriving from class
- language provides means for accomplishing this
- in class/struct declaration, after name of class can add keyword **final** to prevent deriving from class
- example:

```
class Widget final { /* ... */ };  
class Gadget : public Widget { /* ... */ };  
    // ERROR: cannot derive from Widget
```

- might want to prevent deriving from class with destructor that is not virtual
- preventing derivation can sometimes also facilitate better compiler optimization (e.g., via devirtualization)
- might want to prevent derivation so that objects can be copied safely without fear of slicing

Covariant Return Type

- in some special cases, language allows relaxation of rule that type of overriding function f must be same as type of virtual function f overrides
- in particular, requirement that return type be same is relaxed
- return type of derived-class function is permitted to be type derived (directly or indirectly) from return type of base-class function
- this relaxation of return type more formally known as **covariant return type**
- *case of pointer return type*: if original return type B^* then return type of overriding function may be D^* , provided B is public base of D (i.e., may return pointer to more derived type)
- *case of reference return type*: if original return type $B\&$ then return type of overriding function may be $D\&$, provided B is public base of D (i.e., may return reference to more derived type)
- covariant return type can sometimes be exploited in order to avoid need for type casts

Covariant Return Type Example: Cloning

```
1  class Base {
2  public:
3      virtual Base* clone() const {
4          return new Base(*this);
5      }
6      // ...
7  };
8
9  class Derived : public Base {
10 public:
11     // use covariant return type
12     Derived* clone() const override {
13         return new Derived(*this);
14     }
15     // ...
16 };
17
18 int main() {
19     Derived* d = new Derived;
20     Derived* d2 = d->clone();
21     // OK: return type is Derived*
22     // without covariant return type, would need cast:
23     // Derived* d2 = (Derived*)d->clone();
24 }
```

Pure Virtual Functions

- sometimes desirable to require derived class to override virtual function
- **pure virtual function**: virtual function that must be overridden in every derived class
- to declare virtual function as pure, add “= 0” at end of declaration
- example:

```
class Widget {  
public:  
    virtual void doStuff() = 0; // pure virtual  
    // ...  
};
```

- pure virtual function can still be defined, although likely only useful in case of virtual destructor

Abstract Classes

- class with one or more pure virtual functions called **abstract class**
- cannot directly instantiate objects of abstract class (can only use them as base class objects)
- class that derives from abstract class need not override all of its pure virtual methods
- class that does not override all pure virtual methods of abstract base class will also be abstract
- most commonly, abstract classes have no state (i.e., data members) and used to provide interfaces, which can be inherited by other classes
- if class has no pure virtual functions and abstract class is desired, can make destructor pure virtual (but must provide definition of destructor since invoked by derived classes)

Abstract Class Example

```
1  #include <cmath>
2
3  class Shape {
4  public:
5      virtual bool isPolygon() const = 0;
6      virtual float area() const = 0;
7      virtual ~Shape() {};
8  };
9
10 class Rectangle : public Shape {
11 public:
12     Rectangle(float w, float h) : w_(w), h_(h) {}
13     bool isPolygon() const override {return true;}
14     float area() const override {return w_ * h_;}
15 private:
16     float w_; // width of rectangle
17     float h_; // height of rectangle
18 };
19
20 class Circle : public Shape {
21 public:
22     Circle(float r) : r_(r) {}
23     float area() const override {return M_PI * r_ * r_;}
24     bool isPolygon() const override {return false;}
25 private:
26     float r_; // radius of circle
27 };
```

Pure Virtual Destructor Example

```
1 class Abstract {
2 public:
3     virtual ~Abstract() = 0; // pure virtual destructor
4     // ... (no other virtual functions)
5 };
6
7 inline Abstract::~~Abstract()
8     { /* possibly empty */ }
```

The `dynamic_cast` Operator

- often need to upcast and downcast (as well as cast sideways) in inheritance hierarchy
- **`dynamic_cast`** can be used to safely perform type conversions on pointers and references to classes
- syntax: **`dynamic_cast`**<T> (*expr*)
- types involved must be *polymorphic* (i.e., have at least one virtual function)
- inspects run-time information about types to determine whether cast can be safely performed
- if conversion is valid (i.e., *expr* can validly be cast to T), casts *expr* to type T and returns result
- if conversion is not valid, cast fails
- if *expr* is of pointer type, **`nullptr`** is returned upon failure
- if *expr* is of reference type, `std::bad_cast` exception is thrown upon failure

dynamic_cast Example

```
1  #include <cassert>
2
3  class Base {
4  public:
5      virtual void doStuff() { /* ... */ };
6      // ...
7  };
8
9  class Derived1 : public Base { /* ... */ };
10 class Derived2 : public Base { /* ... */ };
11
12 bool isDerived1(Base& b) {
13     return dynamic_cast<Derived1*>(&b) != nullptr;
14 }
15
16 int main() {
17     Base b;
18     Derived1 d1;
19     Derived2 d2;
20     assert(isDerived1(b) == false);
21     assert(isDerived1(d1) == true);
22     assert(isDerived1(d2) == false);
23 }
```

Cost of Run-Time Polymorphism

- typically, run-time polymorphism does not come without run-time cost in terms of both time and memory
- in some contexts, cost can be significant
- typically, virtual functions implemented using virtual function table
- each polymorphic class has virtual function table containing pointers to all virtual functions for class
- each polymorphic class object has pointer to virtual function table
- memory cost to store virtual function table and pointer to table in each polymorphic object
- in most cases, impossible for compiler to inline virtual function calls since function to be called cannot be known until run time
- each virtual function call is made through pointer, which adds overhead

Curiously-Recurring Template Pattern (CRTP)

- when derived type known at compile time, may want behavior similar to virtual functions but without run-time cost (by performing binding at compile time instead of run time)
- can be achieved with technique known as **curiously-recurring template pattern (CRTP)**
- class `Derived` derives from class template instantiation using `Derived` itself as template argument
- example:

```
template <class Derived>
class Base {
    // ...
};

class Derived : public Base<Derived> {
    // ...
};
```

C RTP Example: Static Polymorphism

```
1  #include <iostream>
2
3  template <class Derived>
4  class Base {
5  public:
6      void interface() {
7          std::cout << "Base::interface called\n";
8          static_cast<Derived*>(this)->implementation();
9      }
10     // ...
11 };
12
13 class Derived : public Base<Derived> {
14 public:
15     void implementation() {
16         std::cout << "Derived::implementation called\n";
17     }
18     // ...
19 };
20
21 int main() {
22     Derived d;
23     d.interface();
24     // calls Base::interface which, in turn, calls
25     // Derived::implementation
26     // no virtual function call, however
27
28 }
```

C RTP Example: Static Polymorphism

```
1  class TreeNode {
2  public:
3      enum Kind {RED, BLACK}; // kinds of nodes
4      TreeNode *left(); // get left child node
5      TreeNode *right(); // get right child node
6      Kind kind(); // get kind of node
7      // ...
8  };
9
10 template <class Derived>
11 class GenericVisitor {
12 public:
13     void visit_preorder(TreeNode* node) {
14         if (node) {
15             process_node(node);
16             visit_preorder(node->left());
17             visit_preorder(node->right());
18         }
19     }
20     void visit_inorder(TreeNode* node) { /* ... */ }
21     void visit_postorder(TreeNode* node) { /* ... */ }
22     void process_red_node(TreeNode* node) { /* ... */ };
23     void process_black_node(TreeNode* node) { /* ... */ };
24 private:
25     Derived& derived() {return *static_cast<Derived*>(this);}
26     void process_node(TreeNode* node) {
27         if (node->kind() == TreeNode::RED) {
28             derived().process_red_node(node);
29         } else {
30             derived().process_black_node(node);
31         }
32     }
33 };
34
35 class SpecialVisitor : public GenericVisitor<SpecialVisitor> {
36 public:
37     void process_red_node(TreeNode* node) { /* ... */ }
38 };
39
40 int main() {SpecialVisitor v;}
```


C RTP Example: Comparisons

```
1  #include <cassert>
2
3  template<class Derived>
4  struct Comparisons {
5      friend bool operator==(const Comparisons<Derived>& x,
6          const Comparisons<Derived>& y) {
7          const Derived& xr = static_cast<const Derived&>(x);
8          const Derived& yr = static_cast<const Derived&>(y);
9          return !(xr < yr) && !(yr < xr);
10     }
11     // operator!= and others
12 };
13
14 class Widget : public Comparisons<Widget> {
15 public:
16     Widget(bool b, int i) : b_(b), i_(i) {}
17     friend bool operator<(const Widget& x, const Widget& y)
18         {return x.i_ < y.i_;}
19 private:
20     bool b_;
21     int i_;
22 };
23
24 int main() {
25     Widget w1(true, 1);
26     Widget w2(false, 1);
27     assert(w1 == w2);
28 }
```

C RTP Example: Object Counting

```
1  #include <iostream>
2  #include <cstdlib>
3
4  template <class T>
5  class Counter {
6  public:
7      Counter() {++count_;}
8      Counter(const Counter&) {++count_;}
9      ~Counter() {--count_;}
10     static std::size_t howMany() {return count_;}
11 private:
12     static std::size_t count_;
13 };
14
15 template <class T>
16 std::size_t Counter<T>::count_ = 0;
17
18 // inherit from Counter to count objects
19 class Widget: private Counter<Widget> {
20 public:
21     using Counter<Widget>::howMany;
22     // ...
23 };
24
25 int main() {
26     Widget w1; int c1 = Widget::howMany();
27     Widget w2, w3; int c2 = Widget::howMany();
28     std::cout << c1 << ' ' << c2 << '\n';
29 }
```

Section 2.7.3

Multiple Inheritance and Virtual Inheritance

Multiple Inheritance

- language allows derived class to inherit from more than one base class
- **multiple inheritance (MI)**: deriving from more than one base class
- although multiple inheritance not best solution for most problems, does have some compelling use cases
- one compelling use case is for inheriting interfaces by deriving from abstract base classes with no data members
- when misused, multiple inheritance can lead to very convoluted code
- in multiple inheritance contexts, ambiguities in naming can arise
- for example, if class `Derived` inherits from classes `Base1` and `Base2`, each of which have member called `x`, name `x` can be ambiguous in some contexts
- scope resolution operator can be used to resolve ambiguous names

Ambiguity Resolution Example

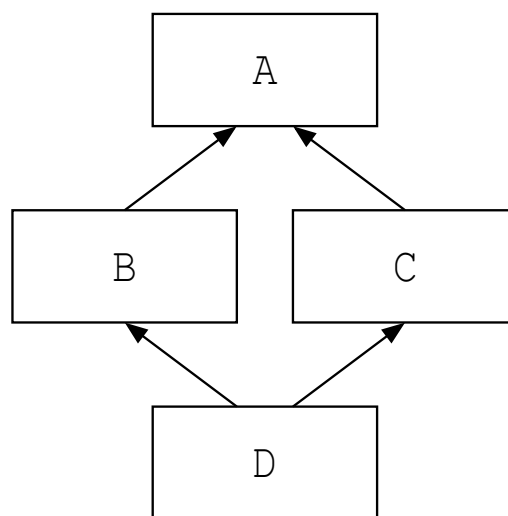
```
class Base1 {  
public:  
    void func ();  
    // ...  
};  
  
class Base2 {  
    void func ();  
    // ...  
};  
  
class Derived : public Base1, public Base2 {  
public:  
    // ...  
};  
  
int main () {  
    Derived d;  
    d.func (); // ERROR: ambiguous function call  
    d.Base1::func (); // OK: invokes Base1::func  
    d.Base2::func (); // OK: invokes Base2::func  
}
```

Multiple Inheritance Example

```
class Input_stream {  
public:  
    virtual ~Input_stream() {}  
    virtual int read_char() = 0;  
    virtual int read(char* buffer, int size) = 0;  
    virtual bool is_input_ready() const = 0;  
    // ... (all pure virtual, no data)  
};  
  
class Output_stream {  
public:  
    virtual ~Output_stream() {}  
    virtual int write_char(char c) = 0;  
    virtual int write(char* buffer, int size) = 0;  
    virtual int flush_output() = 0;  
    // ... (all pure virtual, no data)  
};  
  
class Input_output_stream : public Input_stream,  
    public Output_stream {  
    // ...  
};
```

Dreaded Diamond Inheritance Pattern

- use of multiple inheritance can lead to so called dreaded diamond scenario
- **dreaded diamond** inheritance pattern has following form:



- class D will have *two* subobjects of class A, since class D (indirectly) inherits twice from class A
- situation like one above probably undesirable and often sign of poor design

Dreaded Diamond Example

```
1  class Base {
2  public:
3      // ...
4  protected:
5      int data_;
6  };
7
8  class D1 : public Base { /* ... */ };
9
10 class D2 : public Base { /* ... */ };
11
12 class Join : public D1, public D2 {
13 public:
14     void method() {
15         data_ = 1; // ERROR: ambiguous
16         D1::data_ = 1; // OK: unambiguous
17     }
18 };
19
20 int main() {
21     Join* j = new Join();
22     Base* b;
23     b = j; // ERROR: ambiguous
24     b = static_cast<D1*>(j); // OK: unambiguous
25 }
```


Virtual Inheritance

- when using multiple inheritance, may want to ensure that only one instance of base-class object can appear in derived-class object
- **virtual base class**: base class that is only ever included once in derived class, even if derived from multiple times
- **virtual inheritance**: when derived class inherits from base class that is virtual
- virtual inheritance can be used to avoid situations like dreaded diamond pattern
- order of construction: virtual base classes constructed first in depth-first left-to-right traversal of graph of base classes, where left-to-right refers to order of appearance of base class names

Avoiding Dreaded Diamond With Virtual Inheritance

```
class Base {
public:
    // ...
protected:
    int data_;
};

class D1 : public virtual Base { /* ... */ };

class D2 : public virtual Base { /* ... */ };

class Join : public D1, public D2 {
public:
    void method() {
        data_ = 1; // OK: unambiguous
    }
};

int main() {
    Join* j = new Join();
    Base* b = j; // OK: unambiguous
}
```

Section 2.7.4

References

References I

- 1 N. Meyers. The empty base C++ optimization.
Dr. Dobb's Journal, Aug. 1997.
Available online at <http://www.cantrip.org/emptyopt.html>.
- 2 J. O. Coplien. Curiously recurring template patterns.
C++ Report, pages 24–27, Feb. 1995.
- 3 S. Meyers. Counting objects in C++.
C++ User's Journal, Apr. 1998.
Available online at <http://www.drdobbs.com/cpp/counting-objects-in-c/184403484>.
- 4 A. Nasonov. Better encapsulation for the curiously recurring template pattern.
Overload, 70:11–13, Dec. 2005.

Section 2.8

C++ Standard Library

- C++ standard library provides huge amount of functionality (orders of magnitude more than C standard library)
- uses `std` namespace (to avoid naming conflicts)
- well worth effort to familiarize yourself with all functionality in library in order to avoid writing code unnecessarily

C++ Standard Library (Continued)

- functionality can be grouped into following sublibraries:
 - 1 language support library (e.g., exceptions, memory management)
 - 2 diagnostics library (e.g., assertions, exceptions, error codes)
 - 3 general utilities library (e.g., functors, date/time)
 - 4 strings library (e.g., C++ and C-style strings)
 - 5 localization library (e.g., date/time formatting and parsing, character classification)
 - 6 containers library (e.g., sequence containers and associative containers)
 - 7 iterators library (e.g., stream iterators)
 - 8 algorithms library (e.g., searching, sorting, merging, set operations, heap operations, minimum/maximum)
 - 9 numerics library (e.g., complex numbers, math functions)
 - 10 input/output (I/O) library (e.g., streams)
 - 11 regular expressions library (e.g., regular expression matching)
 - 12 atomic operations library (e.g., atomic types, fences)
 - 13 thread support library (e.g., threads, mutexes, condition variables, futures)

Commonly-Used Header Files

Language-Support Library

Header File	Description
<code>cstdlib</code>	run-time support, similar to <code>stdlib.h</code> from C (e.g., <code>exit</code>)
<code>limits</code>	properties of fundamental types (e.g., <code>numeric_limits</code>)
<code>exception</code>	exception handling support (e.g., <code>set_terminate</code> , <code>current_exception</code>)
<code>initializer_list</code>	<code>initializer_list</code> class template

Diagnostics Library

Header File	Description
<code>cassert</code>	assertions (e.g., <code>assert</code>)
<code>stdexcept</code>	predefined exception types (e.g., <code>invalid_argument</code> , <code>domain_error</code> , <code>out_of_range</code>)

Commonly-Used Header Files (Continued 1)

General-Utilities Library

Header File	Description
<code>utility</code>	basic function and class templates (e.g., <code>swap</code> , <code>move</code> , <code>pair</code>)
<code>memory</code>	memory management (e.g., <code>unique_ptr</code> , <code>shared_ptr</code> , <code>addressof</code>)
<code>functional</code>	functors (e.g., <code>less</code> , <code>greater</code>)
<code>type_traits</code>	type traits (e.g., <code>is_integral</code> , <code>is_reference</code>)
<code>chrono</code>	clocks (e.g., <code>system_clock</code> , <code>steady_clock</code> , <code>high_resolution_clock</code>)

Strings Library

Header File	Description
<code>string</code>	C++ string classes (e.g., <code>string</code>)
<code>cstring</code>	C-style strings, similar to <code>string.h</code> from C (e.g., <code>strlen</code>)
<code>cctype</code>	character classification, similar to <code>ctype.h</code> from C (e.g., <code>isdigit</code> , <code>isalpha</code>)

Commonly-Used Header Files (Continued 2)

Containers, Iterators, and Algorithms Libraries

Header File	Description
<code>array</code>	<code>array</code> class
<code>vector</code>	<code>vector</code> class
<code>deque</code>	<code>deque</code> class
<code>list</code>	<code>list</code> class
<code>set</code>	set classes (i.e., <code>set</code> , <code>multiset</code>)
<code>map</code>	map classes (i.e., <code>map</code> , <code>multimap</code>)
<code>unordered_set</code>	unordered set classes (i.e., <code>unordered_set</code> , <code>unordered_multiset</code>)
<code>unordered_map</code>	unordered map classes (i.e., <code>unordered_map</code> , <code>unordered_multimap</code>)
<code>iterator</code>	iterators (e.g., <code>reverse_iterator</code> , <code>back_inserter</code>)
<code>algorithm</code>	algorithms (e.g., <code>min</code> , <code>max</code> , <code>sort</code>)

Commonly-Used Header Files (Continued 3)

Numerics Library

Header File	Description
<code>cmath</code>	C math library, similar to <code>math.h</code> from C (e.g., <code>M_PI</code> on POSIX-compliant systems, <code>sin</code> , <code>cos</code>)
<code>complex</code>	complex numbers (e.g., <code>complex</code>)
<code>random</code>	random number generation (e.g., <code>uniform_int_distribution</code> , <code>uniform_real_distribution</code> , <code>normal_distribution</code>)

Commonly-Used Header Files (Continued 4)

I/O Library

Header File	Description
<code>iostream</code>	iostream objects (e.g., <code>cin</code> , <code>cout</code> , <code>cerr</code>)
<code>istream</code>	input streams (e.g., <code>istream</code>)
<code>ostream</code>	output streams (e.g., <code>ostream</code>)
<code>fstream</code>	file streams (e.g., <code>fstream</code>)
<code>sstream</code>	string streams (e.g., <code>stringstream</code>)
<code>iomanip</code>	manipulators (e.g., <code>setw</code> , <code>dec</code>)

Regular-Expressions Library

Header File	Description
<code>regex</code>	regular expressions (e.g., <code>basic_regex</code>)

Commonly-Used Header Files (Continued 5)

Atomic-Operations and Thread-Support Libraries

Header File	Description
<code>atomic</code>	atomics (e.g., <code>atomic</code>)
<code>thread</code>	threads (e.g., <code>thread</code>)
<code>mutex</code>	mutexes (e.g., <code>mutex</code> , <code>recursive_mutex</code> , <code>timed_mutex</code>)
<code>condition_variable</code>	condition variables (e.g., <code>condition_variable</code>)
<code>future</code>	futures (e.g., <code>future</code> , <code>shared_future</code> , <code>promise</code>)

Section 2.8.1

Containers, Iterators, and Algorithms

Standard Template Library (STL)

- large part of C++ standard library is collection of class/function templates known as standard template library (STL)
- STL comprised of three basic building blocks:
 - 1 containers
 - 2 iterators
 - 3 algorithms
- containers store elements for processing (e.g., vector)
- iterators allow access to elements for processing (which are often, but not necessarily, in containers)
- algorithms perform actual processing (e.g., search, sort)

Containers

- **container**: class that represents collection/sequence of elements
- usually container classes are template classes
- **sequence container**: collection in which every element has certain position that depends on time and place of insertion
- examples of sequence containers include:
 - `array` (fixed-size array)
 - `vector` (dynamic-size array)
 - `list` (doubly-linked list)
- **associative container**: collection in which position of element in depends on its value or associated key and some predefined sorting/hashing criterion
- examples of associative containers include:
 - `set` (collection of unique keys, sorted by key)
 - `map` (collection of key-value pairs, sorted by key, keys are unique)

Sequence Containers and Container Adapters

Sequence Containers

Name	Description
<code>array</code>	fixed-size array
<code>vector</code>	dynamic-size array
<code>deque</code>	double-ended queue
<code>forward_list</code>	singly-linked list
<code>list</code>	doubly-linked list

Container Adapters

Name	Description
<code>stack</code>	stack
<code>queue</code>	FIFO queue
<code>priority_queue</code>	priority queue

Associative Containers

Ordered Associative Containers

Name	Description
<code>set</code>	collection of unique keys, sorted by key
<code>map</code>	collection of key-value pairs, sorted by key, keys are unique
<code>multiset</code>	collection of keys, sorted by key, duplicate keys allowed
<code>multimap</code>	collection of key-value pairs, sorted by key, duplicate keys allowed

Unordered Associative Containers

Name	Description
<code>unordered_set</code>	collection of unique keys, hashed by key
<code>unordered_map</code>	collection of key-value pairs, hashed by key, keys are unique
<code>unordered_multiset</code>	collection of keys, hashed by key, duplicate keys allowed)
<code>unordered_multimap</code>	collection of key-value pairs, hashed by key, duplicate keys allowed

Typical Container Member Functions

- some member functions typically provided by container classes listed below (where `T` denotes name of container class)

Function	Description
<code>T()</code>	create empty container (default constructor)
<code>T(const T&)</code>	copy container (copy constructor)
<code>T(T&&)</code>	move container (move constructor)
<code>~T</code>	destroy container (including its elements)
<code>empty</code>	test if container empty
<code>size</code>	get number of elements in container
<code>push_back</code>	insert element at end of container
<code>clear</code>	remove all elements from container
<code>operator=</code>	assign all elements of one container to other
<code>operator[]</code>	access element in container

Container Example

- example:

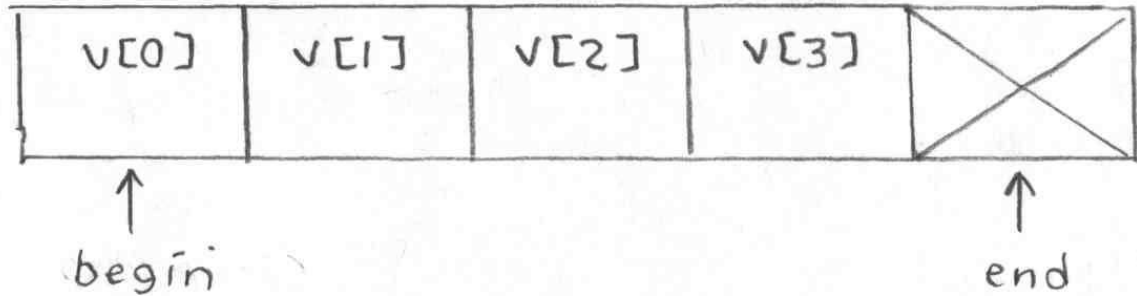
```
1  #include <iostream>
2  #include <vector>
3
4  int main(int argc, char** argv) {
5      std::vector<int> values;
6
7      // append elements with values 0 to 9
8      for (int i = 0; i < 10; ++i)
9          values.push_back(i);
10
11     // print each element followed by space
12     for (int i = 0; i < values.size(); ++i)
13         std::cout << values[i] << ' ';
14     std::cout << '\n';
15
16     return 0;
17 }
```

- program will produce output:

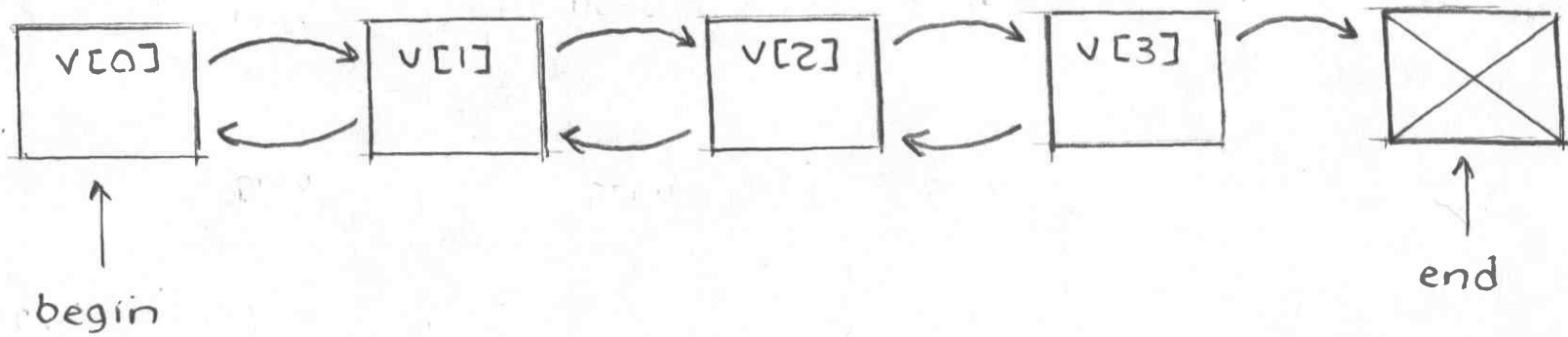
0 1 2 3 4 5 6 7 8 9

Motivation for Iterators

- different containers organize elements (of container) differently in memory
- want uniform manner in which to access elements in any arbitrary container
- organization of elements in array/vector container:

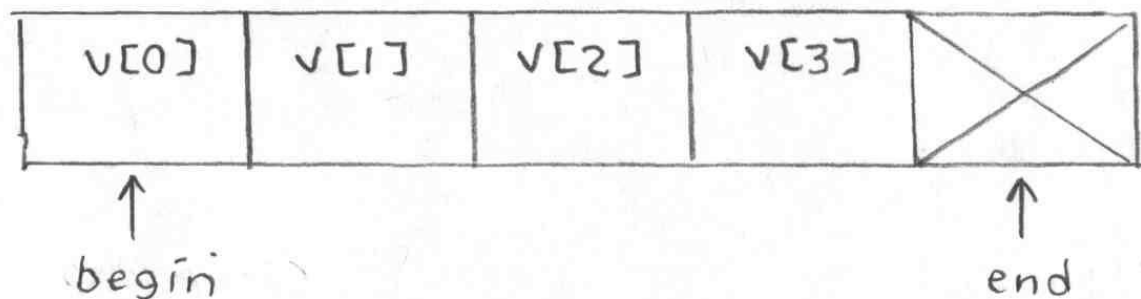


- organization of elements in doubly-linked list container:



Motivation for Iterators (Continued)

- consider array/vector container with `int` elements:



- suppose we want to set all elements in container to zero
- we could use code like:

```
// int* begin; int* end;  
for (int* iter = begin; iter != end; ++iter)  
    *iter = 0;
```

- could we make similar-looking code work for more complicated organization like doubly-linked list?
- yes, create user-defined type that provides all pointer operations used above (e.g., dereference, increment, comparison, assignment)
- this leads to notion of iterator

Iterators

- **iterator**: object that allows iteration over collection of elements, where elements are often (but not necessarily) in container
- iterators support many of same operations as pointers
- in some cases, iterator may actually be pointer; more frequently, iterator is user-defined type
- five different categories of iterators: 1) input, 2) output, 3) forward, 4) bidirectional, and 5) random access
- iterator has particular level of functionality, depending on category
- one of three possibilities of access order:
 - ① forward (i.e., one direction only)
 - ② forward and backward
 - ③ any order (i.e., random access)
- one of three possibilities in terms of read/write access:
 - ① can only read referenced element (once or multiple times)
 - ② can only write referenced element (once or multiple times)
 - ③ can read and write referenced element (once or multiple times)
- const and mutable (i.e., non-const) variants (i.e., read-only or read/write access, respectively)

Abilities of Iterator Categories

Category	Ability	Providers
Input	Reads (once only) forward	<code>istream</code> <code>(istream_iterator)</code>
Output	Writes (once only) forward	<code>ostream</code> <code>(ostream_iterator)</code> , <code>inserter_iterator</code>
Forward	Reads and writes forward	<code>forward_list</code> , <code>unordered_set</code> , <code>unordered_map</code>
Bidirectional	Reads and writes forward and backward	<code>list</code> , <code>set</code> , <code>multiset</code> , <code>map</code> , <code>multimap</code>
Random access	Reads and writes with random access	<code>array</code> , <code>vector</code> , <code>deque</code> , <code>string</code>

Input Iterators

Expression	Effect
<code>T (a)</code>	copies iterator (copy constructor)
<code>*a</code> <code>a->m</code>	dereference as rvalue (i.e., read only); can only be dereferenced once
<code>++a</code>	steps forward (returns new position)
<code>a++</code>	steps forward (returns old position)
<code>a == b</code>	test for equality
<code>a != b</code>	test for inequality

- not assignable (i.e., no assignment operator)

Output Iterators

Expression	Effect
<code>T (a)</code>	copies iterator (copy constructor)
<code>* a</code> <code>a->m</code>	dereference as lvalue (i.e., write only); can only be dereferenced once
<code>++a</code>	steps forward (returns new position)
<code>a++</code>	steps forward (returns old position)

- not assignable (i.e., no assignment operator)
- no comparison operators (i.e., **operator==**, **operator!=**)

Forward Iterators

Expression	Effect
<code>T()</code>	default constructor
<code>T(a)</code>	copy constructor
<code>a = b</code>	assignment
<code>*a</code> <code>a->m</code>	dereference as lvalue (i.e., write only); can only be dereferenced once
<code>++a</code>	steps forward (returns new position)
<code>a++</code>	steps forward (returns old position)
<code>a == b</code>	test for equality
<code>a != b</code>	test for inequality

- must ensure that valid to dereference iterator before doing so

Bidirectional Iterators

- bidirectional iterators are forward iterators that provide additional functionality of being able to iterate backward over elements
- bidirectional iterators have all functionality of forward iterators as well as those listed in table below

Expression	Effect
<code>--a</code>	steps backward (returns new position)
<code>a--</code>	steps backward (returns old position)

Random-Access Iterators

- random access iterators provide all functionality of bidirectional iterators as well as providing random access to elements
- random access iterators provide all functionality of bidirectional iterators as well as those listed in table below

Expression	Effect
$a[n]$	reference element at index n (where n can be negative)
$a += n$	steps n elements forward (where n can be negative)
$a -= n$	steps n elements backward (where n can be negative)
$a + n$	iterator for n th next element
$n + a$	iterator for n th next element
$a - n$	iterator for n th previous element
$a - b$	distance from a to b
$a < b$	test if a before b
$a > b$	test if a after b
$a <= b$	test if a not after b
$a >= b$	test if a not before b

- pointers (built into language) are examples of random-access iterators

Iterator Example

```
1  #include <iostream>
2  #include <vector>
3
4  int main(int argc, char** argv) {
5      std::vector<int> values(10);
6
7      std::cout << "number of elements: " <<
8          (values.end() - values.begin()) << '\n';
9
10     // initialize elements of vector to 0, 1, 2, ...
11     for (std::vector<int>::iterator i = values.begin();
12         i != values.end(); ++i) {
13         *i = i - values.begin();
14     }
15
16     // print elements of vector
17     for (std::vector<int>::const_iterator i =
18         values.begin(); i != values.end(); ++i) {
19         std::cout << *i << ' ';
20     }
21     std::cout << '\n';
22
23     return 0;
24 }
```

Iterator Gotchas

- do not dereference iterator unless it is known to validly reference some object
- some operations on container can *invalidate* some or all iterators referencing elements in container
- critically important to know *which operations invalidate* iterators in order to avoid using iterator that has been invalidated
- incrementing iterator *past end* of container or decrementing iterator *before beginning* of container results in undefined behavior
- input and output iterators can only be dereferenced *once* at each position

- **algorithm**: sequence of computations applied to some generic type
- algorithms use iterators to access elements involved in computation
- often pair of iterators used to specify *range* of elements on which to perform some computation
- what follows only provides brief summary of algorithms
- for more details on algorithms, see:
 - <http://www.cplusplus.com/reference/algorithm>
 - <http://en.cppreference.com/w/cpp/algorithm>

Non-Modifying Sequence Operations

Name	Description
<code>all_of</code>	test if condition true for all elements in range
<code>any_of</code>	test if condition true for any element in range
<code>none_of</code>	test if condition true for no elements in range
<code>for_each</code>	apply function to range
<code>find</code>	find values in range
<code>find_if</code>	find element in range
<code>find_if_not</code>	find element in range (negated)
<code>find_end</code>	find last subsequence in range
<code>find_first_of</code>	find element from set in range
<code>adjacent_find</code>	find equal adjacent elements in range
<code>count</code>	count appearances of value in range
<code>count_if</code>	count number of elements in range satisfying condition
<code>mismatch</code>	get first position where two ranges differ
<code>equal</code>	test whether elements in two ranges differ
<code>search</code>	find subsequence in range
<code>search_n</code>	find succession of equal values in range

Functions (Continued 1)

Modifying Sequence Operations

Name	Description
<code>copy</code>	copy range of elements
<code>copy_if</code>	copy certain elements of range
<code>copy_n</code>	copy n elements
<code>copy_backward</code>	copy range of elements backwards
<code>move</code>	move range of elements
<code>move_backward</code>	move range of elements backwards
<code>swap</code>	exchange values of two objects
<code>swap_ranges</code>	exchange values of two ranges
<code>iter_swap</code>	exchange values of objects referenced by two iterators
<code>transform</code>	apply function to range
<code>replace</code>	replace value in range
<code>replace_if</code>	replace values in range
<code>replace_copy</code>	copy range replacing value
<code>replace_copy_if</code>	copy range replacing value

Functions (Continued 2)

Modifying Sequence Operations (Continued)

Name	Description
<code>fill</code>	fill range with value
<code>fill_n</code>	fill sequence with value
<code>generate</code>	generate values for range with function
<code>generate_n</code>	generate values for sequence with function
<code>remove</code>	remove value from range
<code>remove_if</code>	remove elements from range
<code>remove_copy</code>	copy range removing value
<code>remove_copy_if</code>	copy range removing values
<code>unique</code>	remove consecutive duplicates in range
<code>unique_copy</code>	copy range removing duplicates
<code>reverse</code>	reverse range
<code>reverse_copy</code>	copy range reversed
<code>rotate</code>	rotate elements in range
<code>rotate_copy</code>	copies and rotates elements in range
<code>shuffle</code>	randomly permute elements in range

Functions (Continued 3)

Partition Operations

Name	Description
<code>is_partitioned</code>	test if range is partitioned by predicate
<code>partition</code>	partition range in two
<code>partition_copy</code>	copies range partition in two
<code>stable_partition</code>	partition range in two (stable ordering)
<code>partition_point</code>	get partition point

Sorting

Name	Description
<code>is_sorted</code>	test if range is sorted
<code>is_sorted_until</code>	find first unsorted element in range
<code>sort</code>	sort elements in range
<code>stable_sort</code>	sort elements in range, preserving order of equivalents
<code>partial_sort</code>	partially sort elements in range
<code>partial_sort_copy</code>	copy and partially sort range
<code>nth_element</code>	sort element in range

Functions (Continued 4)

Binary Search (operating on sorted ranges)

Name	Description
<code>lower_bound</code>	get iterator to lower bound
<code>upper_bound</code>	get iterator to upper bound
<code>equal_range</code>	get subrange of equal elements
<code>binary_search</code>	test if value exists in sorted range

Set Operations (on sorted ranges)

Name	Description
<code>merge</code>	merge sorted ranges
<code>inplace_merge</code>	merge consecutive sorted ranges
<code>includes</code>	test whether sorted range includes another sorted range
<code>set_union</code>	union of two sorted ranges
<code>set_intersection</code>	intersection of two sorted ranges
<code>set_difference</code>	difference of two sorted ranges
<code>set_symmetric_difference</code>	symmetric difference of two sorted ranges

Functions (Continued 5)

Heap Operations

Name	Description
<code>is_heap</code>	test if range is heap
<code>is_heap_until</code>	first first element not in heap order
<code>push_heap</code>	push element into heap range
<code>pop_heap</code>	pop element from heap range
<code>make_heap</code>	make heap from range
<code>sort_heap</code>	sort elements of heap

Functions (Continued 6)

Minimum/Maximum

Name	Description
<code>min</code>	get minimum of given values
<code>max</code>	get maximum of given values
<code>minmax</code>	get minimum and maximum of given values
<code>min_element</code>	get smallest element in range
<code>max_element</code>	get largest element in range
<code>minmax_element</code>	get smallest and largest elements in range
<code>lexicographic_compare</code>	lexicographic less-than comparison
<code>is_permutation</code>	test if range permutation of another
<code>next_permutation</code>	transform range to next permutation
<code>prev_permutation</code>	transform range to previous permutation

Functions (Continued 7)

Numeric Operations

Name	Description
<code>iota</code>	fill range with successive values
<code>accumulate</code>	accumulate values in range
<code>adjacent_difference</code>	compute adjacent difference of range
<code>inner_product</code>	compute inner product of range
<code>partial_sum</code>	compute partial sums of range

Algorithms Example

```
1  #include <iostream>
2  #include <vector>
3  #include <algorithm>
4
5  int main(int argc, char** argv) {
6      std::vector<int> values;
7      int x;
8      while (std::cin >> x)
9          values.push_back(x);
10     std::cout << "zero count: " << std::count(
11         values.begin(), values.end(), 0) << '\n';
12     std::random_shuffle(values.begin(), values.end());
13     std::cout << "random order: ";
14     for (std::vector<int>::const_iterator i =
15         values.begin(); i != values.end(); ++i)
16         std::cout << *i << ' ';
17     std::cout << '\n';
18     std::sort(values.begin(), values.end());
19     std::cout << "sorted order: ";
20     for (std::vector<int>::const_iterator i =
21         values.begin(); i != values.end(); ++i)
22         std::cout << *i << ' ';
23     std::cout << '\n';
24     return 0;
25 }
```

Prelude to Functor Example

- consider `std::transform` function template:

```
template <class InputIterator, class OutputIterator,  
         class UnaryOperator>  
OutputIterator transform(InputIterator first,  
                        InputIterator last, OutputIterator result,  
                        UnaryOperator op);
```

- applies `op` to each element in range `[first,last)` and stores each returned value in range beginning at `result`

- `std::transform` might be written as:

```
template <class InputIterator, class OutputIterator,  
         class UnaryOperator>  
OutputIterator transform(InputIterator first,  
                        InputIterator last, OutputIterator result,  
                        UnaryOperator op) {  
    while (first != last) {  
        *result = op(*first);  
        ++first;  
        ++result;  
    }  
    return result;  
}
```

- `op` is object of type that can be used with function call syntax (i.e., function or functor)

Functor Example

```
1  #include <iostream>
2  #include <vector>
3  #include <algorithm>
4
5  struct MultiplyBy { // Functor class
6      MultiplyBy(double factor) : factor_(factor) {}
7      double operator()(double x) const {
8          return factor_ * x;
9      }
10 private:
11     // state information
12     double factor_; // multiplicative factor
13 };
14
15 int main() {
16     MultiplyBy mb(2.0);
17     std::vector<double> v;
18     v.push_back(1);
19     v.push_back(2);
20     v.push_back(3);
21     // v contains 1 2 3
22     std::transform(v.begin(), v.end(), v.begin(), mb);
23     // v contains 2 4 6
24 }
```

Section 2.8.2

The `vector` Class Template

The `vector` Class Template

- one-dimensional array, where type of array elements and storage allocator specified by template parameters

- `vector` declared as:

```
template <class T, class Allocator = allocator<T>>  
class vector;
```

- T: type of elements in vector
- Allocator: type of object used to handle storage allocation (unless custom storage allocator needed, use default `allocator<T>`)
- what follows only intended to provide overview of `vector`
- for additional details on `vector`, see:
 - <http://www.cplusplus.com/reference/stl/vector>
 - <http://en.cppreference.com/w/cpp/container/vector>

Member Types

Member Type	Description
<code>value_type</code>	T (i.e., element type)
<code>allocator_type</code>	Allocator (i.e., allocator)
<code>size_type</code>	type used for measuring size (typically unsigned integral type)
<code>difference_type</code>	type used to measure distance (typically signed integral type)
<code>reference</code>	<code>value_type&</code>
<code>const_reference</code>	const <code>value_type&</code>
<code>pointer</code>	<code>allocator_traits<Allocator>::pointer</code>
<code>const_pointer</code>	<code>allocator_traits<Allocator>::const_pointer</code>
<code>iterator</code>	<i>random-access</i> iterator type
<code>const_iterator</code>	const <i>random-access</i> iterator type
<code>reverse_iterator</code>	reverse iterator type (<code>reverse_iterator<iterator></code>)
<code>const_reverse_iterator</code>	const reverse iterator type (<code>reverse_iterator<const_iterator></code>)

Member Functions

Construction, Destruction, and Assignment

Member Name	Description
constructor	construct vector (overloaded)
destructor	destroy vector
operator=	assign vector

Iterators

Member Name	Description
begin	return iterator to beginning
end	return iterator to end
cbegin	return const iterator to beginning
cend	return const iterator to end
rbegin	return reverse iterator to beginning
rend	return reverse iterator to end
crbegin	return const reverse iterator to beginning
crend	return const reverse iterator to end

Member Functions (Continued 1)

Capacity

Member Name	Description
<code>empty</code>	test if vector is empty
<code>size</code>	return size
<code>max_size</code>	return maximum size
<code>capacity</code>	return allocated storage capacity
<code>reserve</code>	request change in capacity
<code>shrink_to_fit</code>	shrink to fit

Element Access

Member Name	Description
<code>operator []</code>	access element (no bounds checking)
<code>at</code>	access element (with bounds checking)
<code>front</code>	access first element
<code>back</code>	access last element
<code>data</code>	return pointer to start of element data

Member Functions (Continued 2)

Modifiers

Member Name	Description
<code>clear</code>	clear content
<code>assign</code>	assign vector content
<code>insert</code>	insert elements
<code>emplace</code>	insert element, constructing in place
<code>push_back</code>	add element at end
<code>emplace_back</code>	insert element at end, constructing in place
<code>erase</code>	erase elements
<code>pop_back</code>	delete last element
<code>resize</code>	change size
<code>swap</code>	swap content of two vectors

Allocator

Member Name	Description
<code>get_allocator</code>	get allocator used by vector

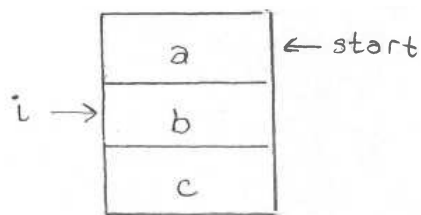
Invalidation of References, Iterators, and Pointers

- **capacity**: total number of elements that vector could hold without requiring reallocation of memory
- any operation that causes reallocation of memory used to hold elements of vector invalidates *all* iterators, references, and pointers referring to elements in vector
- any operation that changes capacity of vector causes reallocation of memory
- any operation that adds or deletes elements can invalidate references, iterators, and pointers
- operations that can potentially invalidate references, iterators, and pointers to elements in vector include:

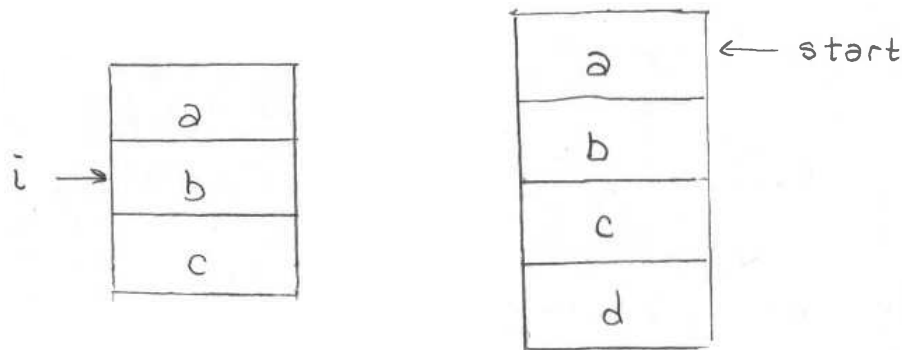
`insert`, `erase`, `push_back`, `pop_back`, `emplace`, `emplace_back`,
`resize`, `reserve`, **`operator=`**, `assign`, `clear`, `shrink_to_fit`, `swap`
(past-the-end iterator only)

Iterator Invalidation Example

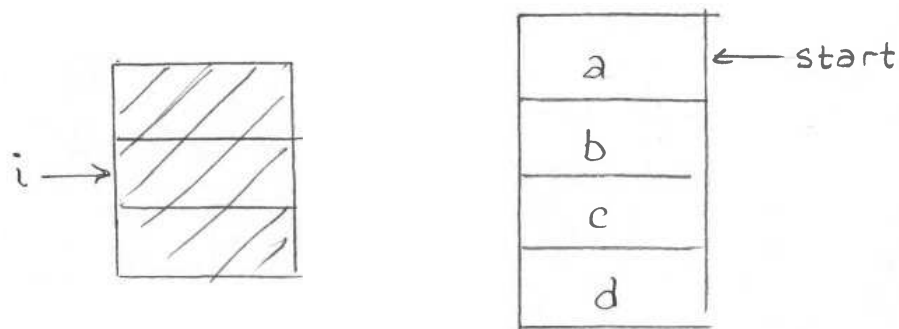
- `start` denotes pointer to first element in array holding vector elements
- `i` is iterator for `vector` (e.g., `vector<T>::const_iterator`, `vector<T>::iterator`)
- initial vector with three elements and capacity of three:



- `push_back(d)` results in new larger array being allocated, contents of old array copied to new one, and then new element added:



- old array is deallocated, iterator `i` is now *invalid*:



vector Example: Constructors

```
std::vector<double> v0;  
    // empty vector  
  
std::vector<double> v1(10);  
    // vector with 10 elements, default constructed  
    // (which for double means uninitialized)  
  
std::vector<double> v2(10, 5.0);  
    // vector with 10 elements, each initialized to 5.0  
  
std::vector<int> v3{1, 2, 3};  
    // vector with 3 elements: 1, 2, 3  
    // std::initializer_list (note brace brackets)
```

vector Example: Iterators

```
1  #include <iostream>
2  #include <vector>
3
4  int main() {
5      std::vector<int> v{0, 1, 2, 3};
6      for (auto& i : v) {++i;}
7      for (auto i : v) {
8          std::cout << ' ' << i;
9      }
10     std::cout << '\n';
11     for (auto i = v.begin(); i != v.end(); ++i) {
12         --(*i);
13     }
14     for (auto i = v.cbegin(); i != v.cend(); ++i) {
15         std::cout << ' ' << *i;
16     }
17     std::cout << '\n';
18     for (auto i = v.crbegin(); i != v.crend(); ++i) {
19         std::cout << ' ' << *i;
20     }
21     std::cout << '\n';
22 }
```

● program output:

```
1 2 3 4
0 1 2 3
3 2 1 0
```

vector Example

```
1  #include <iostream>
2  #include <vector>
3
4  int main() {
5      std::vector<double> values;
6      // ...
7
8      // Erase all elements and then read elements from
9      // standard input.
10     values.clear();
11     double x;
12     while (std::cin >> x) {
13         values.push_back(x);
14     }
15     std::cout << "number of values read: " <<
16         values.size() << '\n';
17
18     // Loop over all elements and print the number of
19     // negative elements found.
20     int count = 0;
21     for (auto i = values.cbegin(); i != values.cend(); ++i) {
22         if (*i < 0.0) {
23             ++count;
24         }
25     }
26     std::cout << "number of negative values: " << count <<
27         '\n';
28 }
```

vector Example: Emplace

```
1  #include <iostream>
2  #include <vector>
3
4  int main() {
5      std::vector<std::vector<int>> v{{1, 2, 3}, {4, 5, 6}};
6      v.emplace_back(10, 0);
7      // The above use of emplace_back is more efficient than:
8      // v.push_back(std::vector<int>(10, 0));
9      for (const auto& i : v) {
10         for (const auto& j : i) {
11             std::cout << ' ' << j;
12         }
13         std::cout << '\n';
14     }
15 }
```

- program output:

```
1 2 3
4 5 6
0 0 0 0 0 0 0 0 0 0
```

Section 2.8.3

The `basic_string` Class Template

The `basic_string` Class Template

- character string type, parameterized on character type, character traits, and storage allocator
- `basic_string` declared as:

```
template <class CharT,  
         class Traits = char_traits<CharT>,  
         class Allocator = allocator<CharT>>  
         class basic_string;
```

- `CharT`: type of characters in string
- `Traits`: class that describes certain properties of `CharT` (normally, use default)
- `Allocator`: type of object used to handle storage allocation (unless custom storage allocator needed, use default)
- `string` is simply abbreviation for `basic_string<char>`
- what follows is only intended to provide overview of `basic_string` template class (and `string` class)
- for more details on `basic_string`, see:
 - http://www.cplusplus.com/reference/string/basic_string
 - http://en.cppreference.com/w/cpp/string/basic_string

Member Types

Member Type	Description
<code>traits_type</code>	Traits (i.e., character traits)
<code>value_type</code>	<code>Traits::char_type</code> (i.e., character type)
<code>allocator_type</code>	Allocator
<code>size_type</code>	<code>allocator_traits<Allocator>::size_type</code>
<code>difference_type</code>	<code>allocator_traits<Allocator>::difference_type</code>
<code>reference</code>	<code>value_type&</code>
<code>const_reference</code>	const <code>value_type&</code>
<code>pointer</code>	<code>allocator_traits<Allocator>::pointer</code>
<code>const_pointer</code>	<code>allocator_traits<Allocator>::const_pointer</code>
<code>iterator</code>	random-access iterator type
<code>const_iterator</code>	const random-access iterator type
<code>reverse_iterator</code>	reverse iterator type (<code>reverse_iterator<iterator></code>)
<code>const_reverse_iterator</code>	const reverse iterator type (<code>reverse_iterator<const_iterator></code>)



Member Functions

Construction, Destruction, and Assignment

Member Name	Description
constructor	construct
destructor	destroy
operator=	assign

Iterators

Member Name	Description
<code>begin</code>	return iterator to beginning
<code>end</code>	return iterator to end
<code>cbegin</code>	return const iterator to beginning
<code>cend</code>	return const iterator to end
<code>rbegin</code>	return reverse iterator to reverse beginning
<code>rend</code>	return reverse iterator to reverse end
<code>crbegin</code>	return const reverse iterator to reverse beginning
<code>crend</code>	return const reverse iterator to reverse end

Member Functions (Continued 1)

Capacity

Member Name	Description
<code>empty</code>	test if string empty
<code>size</code>	get length of string
<code>length</code>	same as <code>size</code>
<code>max_size</code>	get maximum size of string
<code>capacity</code>	get size of allocated storage
<code>reserve</code>	change capacity
<code>shrink_to_fit</code>	shrink to fit

Element Access

Member Name	Description
<code>operator []</code>	access character in string (no bounds checking)
<code>at</code>	access character in string (with bounds checking)
<code>front</code>	access first character in string
<code>back</code>	access last character in string

Member Functions (Continued 2)

Operations

Member Name	Description
<code>clear</code>	clear string
<code>assign</code>	assign content to string
<code>insert</code>	insert into string
<code>push_back</code>	append character to string
<code>operator+=</code>	append to string
<code>append</code>	append to string
<code>erase</code>	erase characters from string
<code>pop_back</code>	delete last character from string
<code>replace</code>	replace part of string
<code>resize</code>	resize string
<code>swap</code>	swap contents with another string

Member Functions (Continued 3)

Operations (Continued)

Member Name	Description
<code>c_str</code>	get nonmodifiable C-string equivalent
<code>data</code>	obtain pointer to first character of string
<code>copy</code>	copy sequence of characters from string
<code>substr</code>	generate substring
<code>compare</code>	compare strings

Search

Member Name	Description
<code>find</code>	find content in string
<code>rfind</code>	find last occurrence of content in string
<code>find_first_of</code>	find character in string
<code>find_first_not_of</code>	find absence of character in string
<code>find_last_of</code>	find character in string from end
<code>find_last_not_of</code>	find absence of character in string from end

Member Functions (Continued 4)

Allocator

Member Name	Description
<code>get_allocator</code>	get allocator

Non-Member Functions

Numeric Conversions

Name	Description
<code>stoi</code>	convert string to int
<code>stol</code>	convert string to long
<code>stoll</code>	convert string to long long
<code>stoul</code>	convert string to unsigned long
<code>stoull</code>	convert string to unsigned long long
<code>stof</code>	convert string to float
<code>stod</code>	convert string to double
<code>stold</code>	convert string to long double
<code>to_string</code>	convert integral or floating-point value to <code>string</code>
<code>to_wstring</code>	convert integral or floating-point value to <code>wstring</code>

string Example

```
1  #include <iostream>
2  #include <string>
3
4  int main() {
5      std::string s;
6      if (!(std::cin >> s)) {
7          s.clear();
8      }
9      std::cout << "string: " << s << '\n';
10     std::cout << "length: " << s.size() << '\n';
11     std::string b;
12     for (auto i = s.cbegin(); i != s.crend(); ++i) {
13         b.push_back(*i);
14     }
15     std::cout << "backwards: " << b << '\n';
16
17     std::string msg = "Hello";
18     msg += ", World!"; // append ", World!"
19     std::cout << msg << '\n';
20
21     const char *cstr = s.c_str();
22     std::cout << "C-style string: " << cstr << '\n';
23 }
```

Numeric/String Conversion Example

```
1  #include <iostream>
2  #include <string>
3
4  int main() {
5      double x = 42.24;
6      // Convert double to string.
7      std::string s = std::to_string(x);
8      std::cout << s << '\n';
9
10     s = "3.14";
11     // Convert string to double.
12     x = std::stod(s);
13     std::cout << x << '\n';
14
15 }
```

Section 2.8.4

Time Measurement

Time Measurement

- time measurement capabilities provided by part of general utilities library (of standard library)
- header file `chrono`
- identifiers in namespace `std::chrono`
- **duration**: time interval
- **time point**: specific point in time (measured relative to epoch)
- **clock**: measures time in terms of time points
- several clocks provided for measuring time
- what follows only intended to provide overview of chrono part of library
- for additional information on chrono part of library, see:
 - <http://www.cplusplus.com/reference/chrono>
 - <http://en.cppreference.com/w/cpp/chrono>

Time Points and Intervals

Name	Description
<code>duration</code>	time interval
<code>time_point</code>	point in time

Clocks

Name	Description
<code>system_clock</code>	system clock (which may be adjusted)
<code>steady_clock</code>	monotonic clock that ticks at constant rate
<code>high_resolution_clock</code>	clock with shortest tick period available

std::chrono Example: Measuring Elapsed Time

```
1  #include <iostream>
2  #include <chrono>
3  #include <cmath>
4
5  double get_result() {
6      double sum = 0.0;
7      for (long i = 0L; i < 10000000L; ++i) {
8          sum += std::sin(i) * std::cos(i);
9      }
10     return sum;
11 }
12
13 int main() {
14     // Get the start time.
15     auto start_time =
16         std::chrono::high_resolution_clock::now();
17     // Do some computation.
18     double result = get_result();
19     // Get the end time.
20     auto end_time = std::chrono::high_resolution_clock::now();
21     // Compute elapsed time in seconds.
22     double elapsed_time = std::chrono::duration<double>(
23         end_time - start_time).count();
24     // Print result and elapsed time.
25     std::cout << "result " << result << '\n';
26     std::cout << "time (in seconds) " << elapsed_time << '\n';
27 }
```

std::chrono Example: Determining Clock Resolution

```
1  #include <iostream>
2  #include <chrono>
3
4  // Get the granularity of a clock in seconds.
5  template <class C>
6  double granularity() {
7      return std::chrono::duration<double>(
8          typename C::duration(1)).count();
9  }
10
11 int main() {
12     std::cout << "system clock:\n" << "period "
13         << granularity<std::chrono::system_clock>() << '\n'
14         << "steady "
15         << std::chrono::system_clock::is_steady << '\n';
16     std::cout << "high resolution clock:\n" << "period "
17         << granularity<std::chrono::high_resolution_clock>()
18         << '\n' << "steady "
19         << std::chrono::high_resolution_clock::is_steady << '\n';
20     std::cout << "steady clock:\n" << "period "
21         << granularity<std::chrono::steady_clock>() << '\n'
22         << "steady "
23         << std::chrono::steady_clock::is_steady << '\n';
24 }
```

Section 2.8.5

Miscellany

std::array Example

```
1  #include <array>
2  #include <iostream>
3  #include <algorithm>
4
5  int main() {
6      // Fixed-size array with 4 elements.
7      std::array<int, 4> a = {{2, 4, 3, 1}};
8
9      // Print elements of array.
10     for (auto i = a.cbegin(); i != a.cend(); ++i) {
11         std::cout << ' ' << *i;
12     }
13     std::cout << '\n';
14
15     // Sort elements of array.
16     std::sort(a.begin(), a.end());
17
18     // Print elements of array.
19     for (auto i = a.cbegin(); i != a.cend(); ++i) {
20         std::cout << ' ' << *i;
21     }
22     std::cout << '\n';
23 }
```

Part 3

More C++

Section 3.1

Exceptions

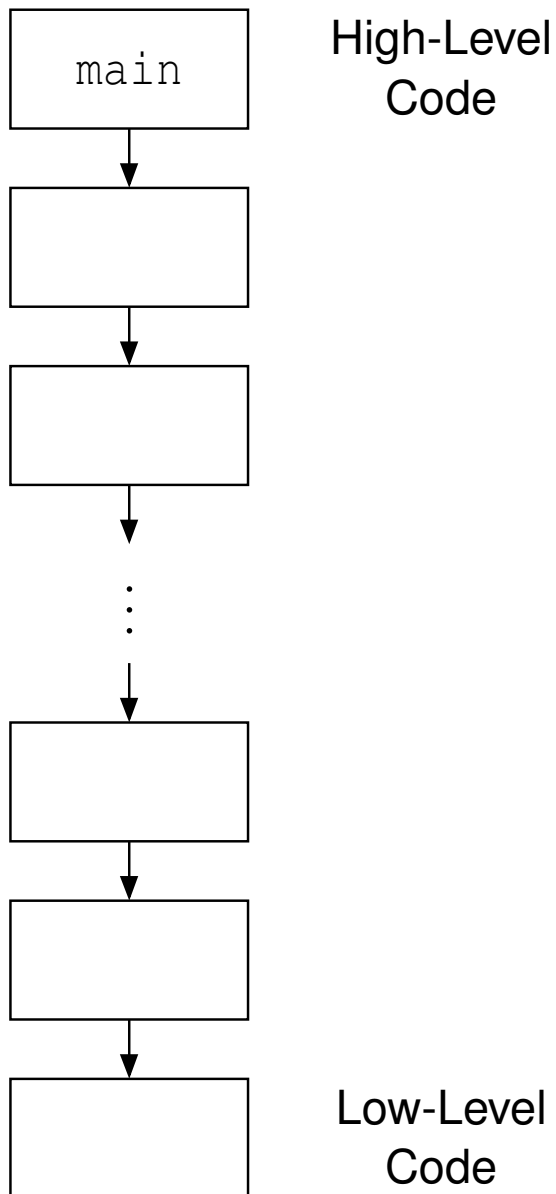
Section 3.1.1

Preliminaries

Exceptions

- exceptions are language mechanism for handling exceptional (i.e., abnormal) situations
- exceptional situation perhaps best thought of as case when code could not do what it was asked to do and usually (but not always) corresponds to error condition
- exceptions often employed for error handling
- exceptions propagate information from point where error *detected* to point where error *handled*
- code that encounters error that it is unable to handle throws exception
- code that wants to handle error catches exception and performs processing necessary to handle error
- exceptions provide convenient way in which to *separate error detection from error handling*

The Problem



- error detected in low-level code
- want to handle error in high-level code
- must propagate error information up call chain

Traditional Error Handling

- if any error occurs, terminate program
 - overly draconian
- pass error code back from function (via return value, reference parameter, or global object) and have caller check error code
 - errors are ignored by default (i.e., explicit action required to check for error condition)
 - caller may forget to check error code allowing error to go undetected
 - code can become cluttered with many checks of error codes, which can adversely affect code readability and maintainability
- call error handler if error detected
 - may not be possible or practical for handler to recover from particular error (e.g., handler may not have access to all information required to recover from error)

Example: Traditional Error Handling

```
1  #include <iostream>
2
3  bool func3() {
4      bool success = false;
5      // ...
6      return success;
7  }
8
9  bool func2() {
10     if (!func3()) {return false;}
11     // ...
12     return true;
13 }
14
15 bool func1() {
16     if (!func2()) {return false;}
17     // ...
18     return true;
19 }
20
21 int main() {
22     if (!func1()) {
23         std::cout << "failed\n";
24         return 1;
25     }
26     // ...
27 }
```


Error Handling With Exceptions

- when error condition detected, signalled by throwing exception (with **throw** statement)
- exception is object that describes error condition
- thrown exception caught by handler (in **catch** clause of **try** statement), which takes appropriate action to handle error condition associated with exception
- handler can be in different function from where exception thrown
- error-free code path tends to be relatively simple, since no need to explicitly check for error conditions
- error condition less likely to go undetected, since uncaught exception terminates program

Example: Exceptions

```
1  #include <iostream>
2  #include <stdexcept>
3
4  void func3() {
5      bool success = false;
6      // ...
7      if (!success) {throw std::runtime_error("Yikes!");}
8  }
9
10 void func2() {
11     func3();
12     // ...
13 }
14
15 void func1() {
16     func2();
17     // ...
18 }
19
20 int main() {
21     try {func1();}
22     catch (...) {
23         std::cout << "failed\n";
24         return 1;
25     }
26     // ...
27 }
```

safe_divide Example: Traditional Error Handling

```
1  #include <iostream>
2  #include <vector>
3  #include <utility>
4
5  std::pair<bool, int> safe_divide(int x, int y) {
6      if (!y) {
7          return std::make_pair(false, 0);
8      }
9      return std::make_pair(true, x / y);
10 }
11
12 int main() {
13     std::vector<std::pair<int, int>> v = {{10, 2}, {10, 0}};
14     for (auto p : v) {
15         auto result = safe_divide(p.first, p.second);
16         if (result.first) {
17             int quotient = result.second;
18             std::cout << quotient << '\n';
19         } else {
20             std::cerr << "division by zero\n";
21         }
22     }
23 }
```

safe_divide Example: Exceptions

```
1  #include <iostream>
2  #include <vector>
3  #include <utility>
4
5  class divide_by_zero {};
6
7  int safe_divide(int x, int y) {
8      if (!y) {
9          throw divide_by_zero();
10     }
11     return x / y;
12 }
13
14 int main() {
15     std::vector<std::pair<int, int>> v = {{10, 2}, {10, 0}};
16     for (auto p : v) {
17         try {
18             std::cout << safe_divide(p.first, p.second) <<
19                 '\n';
20         }
21         catch(const divide_by_zero& e) {
22             std::cerr << "division by zero\n";
23         }
24     }
25 }
```

Exceptions Versus Traditional Error Handling

- advantages of exceptions:
 - exceptions allow for error handling code to be easily separated from code that detects error
 - exceptions can easily pass error information many levels up call chain
 - passing of error information up call chain managed by language (no explicit code required)
- disadvantages of exceptions:
 - writing code that always behaves correctly in presence of exceptions requires great care (as we shall see)
 - although possible to have no execution-time cost when exceptions not thrown, still have memory cost (to store information needed for stack unwinding for case when exception is thrown)

Section 3.1.2

Exceptions

Exceptions

- exceptions are objects
- type of object used to indicate kind of error
- value of object used to provide details about particular occurrence of error
- exception object can have any type (built-in or class type)
- for convenience, standard library provides some basic exception types
- all exception classes in standard library derived (directly or indirectly) from `std::exception` class
- exception object is propagated from one part of code to another by throwing and catching
- exception processing disrupts normal control flow

Standard Exception Classes

Exception Classes Derived from `exception` Class

Type	Description
<code>logic_error</code>	faulty logic in program
<code>runtime_error</code>	error caused by circumstances beyond scope of program
<code>bad_typeid</code>	invalid operand for <code>typeid</code> operator
<code>bad_cast</code>	invalid expression for <code>dynamic_cast</code>
<code>bad_weak_ptr</code>	<code>bad</code> <code>weak_ptr</code> given
<code>bad_function_call</code>	function has no target
<code>bad_alloc</code>	storage allocation failure
<code>bad_exception</code>	use of invalid exception type in certain contexts

Standard Exception Classes (Continued)

Exception Classes Derived from `logic_error` Class

Type	Description
<code>domain_error</code>	domain error (e.g., square root of negative number)
<code>invalid_argument</code>	invalid argument
<code>length_error</code>	length too great (e.g., <code>resize</code> vector beyond <code>max_size</code>)
<code>out_of_range</code>	out of range argument (e.g., subscripting error in <code>vector::at</code>)
<code>future_error</code>	invalid operations on future objects

Exception Classes Derived from `runtime_error` Class

Type	Description
<code>range_error</code>	range error
<code>overflow_error</code>	arithmetic overflow error
<code>underflow_error</code>	arithmetic underflow error
<code>regex_error</code>	error in regular expressions library
<code>system_error</code>	operating-system or other low-level error

Section 3.1.3

Throwing and Catching Exceptions

Throwing Exceptions

- throwing exception accomplished by **throw** statement
- throwing exception transfers control to handler
- object is passed
- type of object determines which handlers can catch it
- handlers specified with **catch** clause of **try** block
- for example

```
throw "OMG!";
```

can be caught by handler of **const char*** type, as in:

```
try {  
    // ...  
}  
catch (const char* p) {  
    // handle character string exceptions here  
}
```

Throwing Exceptions (Continued)

- throw statement initializes temporary object called **exception object**
- type of exception object determined by *static* type of operand of **throw** (so slicing can occur)
- if thrown object is class object, copy/move constructor and destructor must be accessible
- temporary may be moved/copied several times before caught
- advisable for type of exception object to be user defined to reduce likelihood of different parts of code using type in conflicting ways

Catching Exceptions

- exception can be caught by **catch** clause of **try-catch** block
- code that might throw exception placed in **try** block
- code to handle exception placed in **catch** block
- **try-catch** block can have multiple **catch** clauses
- **catch** clauses checked for match in order specified and only first match used
- **catch** (...) can be used to catch any exception
- example:

```
try {  
    // code that might throw exception  
}  
catch (const std::logic_error& e) {  
    // handle logic_error exception  
}  
catch (const std::runtime_error& e) {  
    // handle runtime_error exception  
}  
catch (...) {  
    // handle other exception types  
}
```

Catching Exceptions (Continued)

- catch exceptions by reference in order to:
 - avoid copying, which might throw
 - allow exception object to be modified and then rethrown
 - avoid slicing

Exception During Exception: Catching By Value

```
1  #include <iostream>
2  #include <stdexcept>
3
4  class Error {
5  public:
6      Error(int value) : value_(value) {}
7      Error(Error&& e) : value_(e.value_) {}
8      Error(const Error&) {throw std::runtime_error("copy");}
9      int get() const {return value_;}
10 private:
11     int value_; // error code
12 };
13
14 void func2() {throw Error(42);} // might move
15
16 void func1() {
17     try {func2();}
18     // catch by value (copy throws)
19     catch (Error e) {
20         std::cerr << "yikes\n";
21     }
22 }
23
24 int main() {
25     try {func1();}
26     catch (...) {std::cerr << "exception\n";}
27 }
```

Throwing Polymorphically: Failed Attempt

```
1  #include <iostream>
2
3  class Base {};
4  class Derived : public Base {};
5
6  void func(Base& x) {
7      throw x; // always throws Base
8  }
9
10 int main() {
11     Derived d;
12     try {func(d);}
13     catch (Derived& e) {
14         std::cout << "Derived\n";
15     }
16     catch (...) {
17         std::cout << "not Derived\n";
18     }
19 }
```

- type of exception object determined from *static* type of throw expression

Throwing Polymorphically

```
1  #include <iostream>
2
3  class Base {
4  public:
5      virtual void raise() {throw *this;}
6  };
7  class Derived : public Base {
8  public:
9      virtual void raise() {throw *this;}
10 };
11
12 void func(Base& x) {
13     x.raise();
14 }
15
16 int main() {
17     Derived d;
18     try {func(d);}
19     catch (Derived& e) {
20         std::cout << "Derived\n";
21     }
22     catch (...) {
23         std::cout << "not Derived\n";
24     }
25 }
```

Rethrowing Exceptions

- caught exception can be rethrown by **throw** statement with no operand
- example:

```
try {  
    // code that may throw exception  
}  
catch (...) {  
    throw; // rethrow caught exception  
}
```

Rethrowing Example: Exception Dispatcher Idiom

```
1  void handle_exception() {
2      try {throw;}
3      catch (const exception_1& e) {
4          log_error("exception_1 occurred");
5          // ...
6      }
7      catch (const exception_2& e) {
8          log_error("exception_2 occurred");
9          // ...
10     }
11     // ...
12 }
13
14 void func() {
15     try {operation();}
16     catch (...) {handle_exception();}
17     // ...
18     try {another_operation();}
19     catch (...) {handle_exception();}
20 }
```

- allows reuse of exception handling code

Transfer of Control from Throw Site to Handler

- when exception is thrown, control is transferred to nearest handler (in catch clause) with matching type, where “nearest” means handler for try block most recently entered (by thread) and not yet exited
- if no matching handler found, `std::terminate()` is called
- as control passes from throw expression to handler, destructors are invoked for all automatic objects constructed since try block entered, where automatic objects destroyed in reverse order of construction
- process of calling destructors for automatic objects constructed on path from try block to throw expression called **stack unwinding**
- object not deemed to be constructed if constructor exits due to exception (in which case destructor will not be invoked)
- do not throw exception in destructor since destructors called during exception processing and throwing exception during exception processing will terminate program

Stack Unwinding Example

```
1  void func1() {
2      std::string dave("dave");
3      try {
4          std::string bye("bye");
5          func2();
6      }
7      catch (const std::runtime_error& e) { // Handler
8          std::cerr << e.what() << '\n';
9      }
10 }
11
12 void func2() {
13     std::string world("world");
14     func3(0);
15 }
16
17 void func3(int x) {
18     std::string hello("hello");
19     if (x == 0) {
20         std::string first("first");
21         std::string second("second");
22         throw std::runtime_error("yikes"); // Throw site
23     }
24 }
```

- calling `func1` will result in exception being thrown in `func3`
- during stack unwinding, destructors called in order for `second`, `first`, `hello`, `world`, and `bye` (i.e., reverse order of construction); `dave` unaffected

Function Try Blocks

- function try blocks allow entire function to be wrapped in try block
- function returns when control flow reaches end of catch block (return statement needed for non-void function)

- example:

```
1  #include <iostream>
2  #include <stdexcept>
3
4  int main()
5  try {
6      throw std::runtime_error("yikes");
7  }
8  catch (const std::runtime_error& e) {
9      std::cerr << "runtime error " << e.what() << '\n';
10 }
```

- although function try blocks can be used for any function, most important use cases are for constructors and destructors
- function try block only way to catch exceptions thrown during construction of data members or base objects (which happens before constructor body is entered) or during destruction of data members or base objects (which happens after destructor body exited)

Exceptions and Construction/Destruction

- order of construction:
 - 1 base class objects as listed in type definition left to right
 - 2 data members as listed in type definition top to bottom
 - 3 constructor body
- order of destruction is exact reverse of order of construction, namely:
 - 1 destructor body
 - 2 data members as listed in type definition bottom to top
 - 3 base class objects as listed in type definition right to left
- lifetime of object begins when constructor completes
- constructor might throw in:
 - constructor of base class object
 - constructor of data member
 - constructor body
- need to perform cleanup for constructor body
- will assume destructors do not throw (since very bad idea to throw in destructor)
- any exception caught in function try block of constructor or destructor rethrown implicitly (at end of catch block)

Construction/Destruction Example

```
1  #include <string>
2  #include <iostream>
3
4  struct Base {
5      Base() {}
6      ~Base() {};}
7  };
8
9  class Widget : public Base {
10 public:
11     Widget() {}
12     ~Widget() {}
13     // ...
14 private:
15     std::string s_;
16     std::string t_;
17 };
18
19 int main() {
20     Widget w;
21     // ...
22 }
```


Function Try Block Example

```
1  #include <iostream>
2  #include <stdexcept>
3
4  class Gadget {
5  public:
6      Gadget() {throw std::runtime_error("ctor");}
7      ~Gadget() {}
8  };
9
10 class Widget {
11 public:
12     // constructor uses function try block
13     Widget()
14     try {std::cerr << "ctor body\n";}
15     catch (...) {std::cerr << "exception in ctor\n";}
16     ~Widget() {std::cerr << "dtor body\n";}
17 private:
18     Gadget g_;
19 };
20
21 int main()
22 try {Widget w;}
23 catch (...) {
24     std::cerr << "terminating due to exception\n";
25     return 1;
26 }
```

Section 3.1.4

Exception Specifications

The `noexcept` Specifier

- `noexcept` specifier in function declaration indicates whether or not function can throw exceptions
- `noexcept` specifier with `bool` constant expression argument indicates function does not throw exceptions if expression `true` (otherwise, may throw)
- `noexcept` without argument equivalent to `noexcept (true)`
- except for destructors, not providing `noexcept` specifier equivalent to `noexcept (false)`
- if `noexcept` specifier not provided for destructor, specifier identical to that of implicit declaration (which is, in practice, usually `noexcept`)
- example:

```
void func1(); // may throw anything
void func2() noexcept(false); // may throw anything
void func3() noexcept(true); // does not throw
void func4() noexcept; // does not throw
template <class T>
void func5(T) noexcept(sizeof(T) <= 4);
    // does not throw if sizeof(T) <= 4
```

The `noexcept` Specifier (Continued 1)

- nontrivial `bool` expression for `noexcept` specifier often useful in templates
- example (swap function):

```
1  #include <type_traits>
2  #include <utility>
3
4  // swap two values
5  template <class T>
6  void exchange(T& a, T& b) noexcept (
7      std::is_nothrow_move_constructible<T>::value &&
8      std::is_nothrow_move_assignable<T>::value) {
9      T tmp(std::move(a)); // move construction
10     a = std::move(b);    // move assignment
11     b = std::move(tmp);  // move assignment
12 }
```

The `noexcept` Specifier (Continued 2)

- if function with **`noexcept (true)`** specifier throws exception, `std::terminate` is called immediately

- example:

```
// This function will terminate the program.  
void die_die_die() noexcept {  
    throw 0;  
}
```

- advisable not to use **`noexcept (true)`** specifier unless clear that no reasonable usage of function can throw (in current or *any future* version of code)
- in practice, can often be difficult to guarantee that function will never throw exception (especially when considering *all future* versions of code)

Exceptions and Function Calls

- for some (nonreference) class type `T` and some constant `bool` expression `expr`, consider code such as:

```
T func(T) noexcept (expr);  
T x;  
T y = func(x); // function call
```

- function call can throw exception as result of:
 - ① parameter passing (if pass by value)
 - ② function execution *including return statement*
- in parameter passing, construction and destruction of each parameter happens in context of *calling* function
- consequently, invocation of **noexcept** function can still result in exception being thrown due to parameter passing
- in case of return by value, construction of temporary (if not elided) to hold return value happens in context of *called* function
- if exception due to parameter passing must be avoided, pass by reference or ensure **noexcept** move and/or copy constructor as appropriate
- if exception due to return by value must be avoided, ensure **noexcept** move or copy constructor as appropriate

noexcept Operator

- **noexcept** operator takes expression and returns **bool** indicating if expression can throw exception
- does not actually evaluate expression
- in determining result, only considers **noexcept** specifications for functions involved
- example:

```
1  #include <cstdlib>
2  #include <cassert>
3  #include <utility>
4
5  void increment(int&) noexcept;
6  char* memAlloc(std::size_t);
7
8  // does not throw exception, but not declared noexcept
9  void doesNotThrow() {};
```

```
10
11 int main() {
12     assert(noexcept(1 + 1) == true);
13     assert(noexcept(memAlloc(0)) == false);
14     // Note: does not evaluate expression
15     assert(noexcept(increment(*((int*)0))) == true);
16     assert(noexcept(increment(std::declval<int&>())) ==
17             true);
18     // Note: only uses noexcept specifiers
19     assert(noexcept(doesNotThrow()) == false);
20 }
```

noexcept Operator (Continued)

- **noexcept** operator particularly useful for templates
- example:

```
1  #include <iostream>
2
3  class Int256 { /* ... */ }; // 256-bit integer
4  class BigInt { /* ... */ }; // arbitrary-precision integer
5
6  // function will not throw exception
7  Int256 operator+(const Int256& x, const Int256& y)
8      noexcept;
9
10 // function may throw exception
11 BigInt operator+(const BigInt& x, const BigInt& y);
12
13 // whether function may throw exception depends on T
14 template <class T>
15 T add(const T& x, const T& y) noexcept(noexcept(x + y) &&
16     std::is_nothrow_move_constructible<T>::value)
17 {return x + y;}
18
19 int main() {
20     Int256 i1, i2;
21     BigInt b1, b2;
22     std::cout << "int " << noexcept(add(1, 1)) << '\n'
23         << "Int256 " << noexcept(add(i1, i2)) << '\n'
24         << "BigInt " << noexcept(add(b1, b2)) << '\n';
25 }
```


Dynamic Exception Specifications

- language offers another mechanism for stating exception specifications known as dynamic exception specifications
- dynamic exception specifications are *deprecated* and *should not be used*
- provide exception specification for function using **throw** specifier
- used to specify list of all types of exceptions that can be thrown
- in practice, such a list more of hindrance than help
- if list of all allowable exceptions specified, must check if thrown exception of expected type, which is unnecessary cost
- in terms of compiler optimization, what matters most is whether any exception (regardless of type) can be thrown at all

Section 3.1.5

Storing and Retrieving Exceptions

Storing and Retrieving Exceptions

- might want to store exception and then later retrieve and rethrow it
- exception can be stored using `std::exception_ptr` type
- current exception can be retrieved with `std::current_exception`
- rethrow exception stored in `exception_ptr` object using `std::rethrow_exception`
- provides mechanism for moving exceptions between threads:
 - store exception on one thread
 - then retrieve and rethrow stored exception on another thread
- `std::make_exception_ptr` can be used to make `exception_ptr` object

Example: Storing and Retrieving Exceptions

```
1  #include <exception>
2  #include <stdexcept>
3
4  void yikes () {
5      throw std::runtime_error("Yikes!");
6  }
7
8  std::exception_ptr getException () {
9      try {
10         yikes ();
11     }
12     catch (...) {
13         return std::current_exception ();
14     }
15     return nullptr;
16 }
17
18 int main () {
19     std::exception_ptr e = getException ();
20     std::rethrow_exception (e);
21 }
```

Section 3.1.6

Exception Safety

Resource Management

- **resource**: physical or virtual component of limited availability within computer system
- examples of resources include: memory, files, devices, network connections, processes, threads, and locks
- essential that acquired resource properly released when no longer needed
- when resource not properly released when no longer needed, **resource leak** said to occur
- exceptions have important implications in terms of resource management
- must be careful to avoid resource leaks

Resource Leak Example

```
1 void useBuffer(char* buf) { /* ... */ }
2
3 void doWork() {
4     char* buf = new char[1024];
5     useBuffer(buf);
6     delete[] buf;
7 }
```

- if `useBuffer` throws exception, code that deletes `buf` is never reached

- cleanup operations should always be performed in destructors
- following structure for code is *fundamentally flawed*:

```
void func ()  
{  
    initialize ();  
    do_work ();  
    cleanup ();  
}
```

- code with preceding structure *not exception safe*
- if `do_work` throws, `cleanup` never called and cleanup operation not performed
- in best case, not performing cleanup will probably cause resource leak

Exception Safety and Exception Guarantees

- in order for exception mechanism to be useful, must know what can be assumed about state of program when exception thrown
- operation said to be **exception safe** if it leaves program in valid state when operation is terminated by exception
- several levels of exception safety: basic, strong, nothrow
- **basic guarantee**: all invariants preserved and no resources leaked
- with basic guarantee, partial execution of failed operation may cause side effects
- **strong guarantee**: in addition to basic guarantee, failed operation guaranteed to have no side effects (i.e., commit semantics)
- with strong guarantee, operation can still fail causing exception to be thrown
- **nothrow guarantee**: in addition to basic guarantee, promises not to emit exception (i.e., operation guaranteed to succeed even in presence of exceptional circumstances)

Exception Guarantees

- assume all functions throw if not known otherwise
- code must always provide basic guarantee
- nothrow guarantee should be provided by:
 - destructors
 - move operations (i.e., move constructors and move assignment operators)
 - swap operations
- provide strong guarantee when natural to do so and not more costly than basic guarantee
- examples of strong guarantee:
 - `push_back` for container, subject to certain container-dependent conditions being satisfied (e.g., for `std::vector`, element type has nonthrowing move or is copyable)
 - `insert` on `std::list`
- examples of nothrow guarantee:
 - swap of two containers
 - `pop_back` for container

Resource Acquisition Is Initialization (RAII)

- resource acquisition is initialization (RAII) is programming idiom used to *avoid resource leaks* and *provide exception safety*
- associate resource with owning object (i.e., RAII object)
- period of time over which resource held is tied to lifetime of RAII object
- resource acquired during creation of RAII object
- resource released during destruction of RAII object
- provided RAII object properly destroyed, resource leak cannot occur

Resource Leak Example Revisited

- implementation 1 (not exception safe; has memory leak):

```
1 void useBuffer(char* buf) { /* ... */ }
2
3 void doWork() {
4     char* buf = new char[1024];
5     useBuffer(buf);
6     delete[] buf;
7 }
```

- implementation 2 (exception safe):

```
1 template <class T>
2 class SmartPtr {
3 public:
4     SmartPtr(int size) : ptr_(new T[size]) {}
5     ~SmartPtr() { delete[] ptr_; }
6     operator T* () { return ptr_; }
7     // ...
8 private:
9     T* ptr_;
10 };
11
12 void useBuffer(char* buf) { /* ... */ }
13
14 void doWork() {
15     SmartPtr<char> buf(1024);
16     useBuffer(buf);
17 }
```

Section 3.1.7

Exceptions: Implementation, Cost, and Usage

Implementation of Exception Handling

- standard does not specify how exception handling is to be implemented; only specifies behavior of exception handling
- consider typical implementation here
- potentially significant memory overhead for storing exception object and information required for stack unwinding
- possible to have zero time overhead if no exception thrown
- time overhead significant when exception thrown
- not practical to create exception object on stack, since object frequently needs to be propagated numerous levels up call chain
- exception objects tend to be small
- exception object can be stored in small fixed-size buffer falling back on heap if buffer not big enough

Implementation of Exception Handling (Continued)

- memory required to maintain sufficient information to unwind stack when exception thrown
- two common strategies for maintaining information for stack unwinding: stack-based and table-based strategies
- stack-based strategy:
 - information for stack unwinding is saved on call stack, including list of destructors to execute and exception handlers that might catch exception
 - when exception is thrown, walk stack executing destructors until matching catch found
- table-based strategy:
 - store information to assist in stack unwinding in static tables outside stack
 - call stack used to determine which scopes entered but not exited
 - use look-up operation on static tables to determine where thrown exception will be handled and which destructors to execute
- table-based strategy uses less space on stack but potentially requires considerable storage for tables

Appropriateness of Using Exceptions

- use of exceptions not appropriate in all circumstances
- in practice, exceptions can sometimes (depending on C++ implementation) have prohibitive memory cost for systems with *very limited memory* (e.g., some embedded systems)
- since throwing exception has significant time overhead only use for *infrequently occurring* situations (not common case)
- in code where exceptions can occur, often much more difficult to bound how long code path will take to execute
- since difficult to predict response time of code in presence of exceptions, exceptions often cannot be used in *time critical* component of real-time system (where operation must be guaranteed to complete in specific maximum time)
- considerable amount of code in existence that is *not exception safe*, especially legacy code
- cannot use exceptions in code that is not exception safe

Enforcing Invariants: Exceptions Versus Assertions

- whether invariants should be enforced by exceptions or assertions somewhat controversial
- would recommend only using exceptions for errors from which recovery is likely to be possible
- if error condition detected is indicative of serious programming error, program state may already be sufficiently invalid (e.g., stack trampled, heap corrupted) that exception handling will not work correctly anyhow
- tendency amongst novice programmers is to use exceptions in places where their use is either highly questionable or clearly inappropriate

Section 3.1.8

Smart Pointers and Other RAII Classes

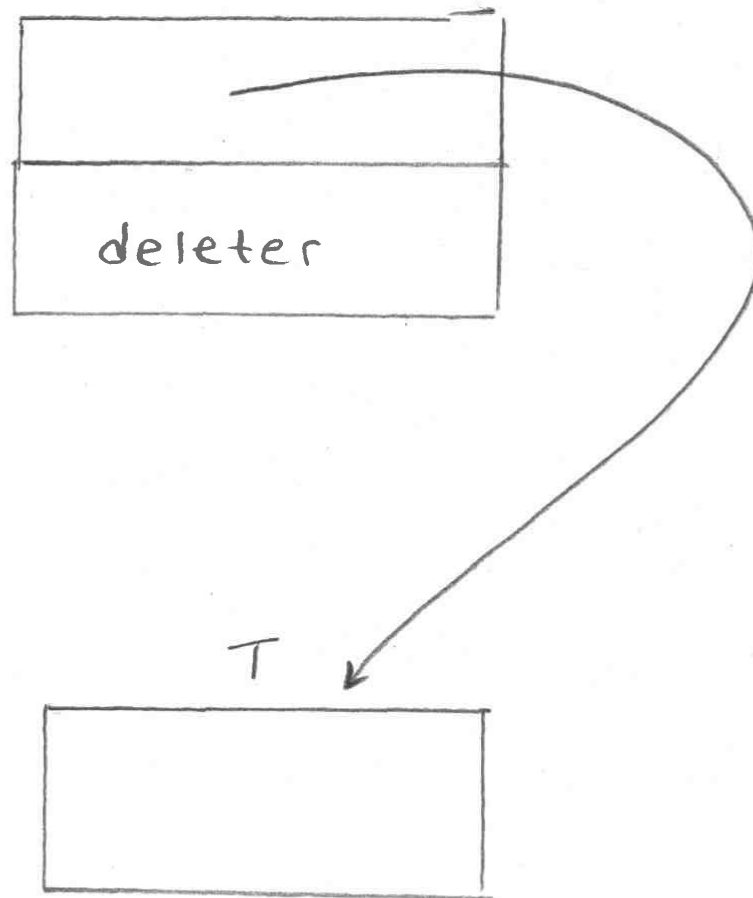
The `std::unique_ptr` Template Class

- `std::unique_ptr` is *smart pointer* that retains *exclusive* ownership of object through pointer
- declaration:

```
template <class T, class Deleter = std::default_delete<T>>  
    class unique_ptr;
```
- T is type of object to be managed (i.e., owned object)
- `Deleter` is callable entity used to delete owned object
- also correctly handles array types via partial specialization (e.g., T could be array of **char**)
- owned object destroyed when `unique_ptr` object goes out of scope
- no two `unique_ptr` objects can own same object
- `unique_ptr` object is movable; move operation transfers ownership
- `unique_ptr` object is not copyable, as copying would create additional owners
- `std::make_unique` template function often used to create `unique_ptr` objects (for exception-safety reasons)

The `std::unique_ptr` Template Class (Continued)

`unique_ptr<T>`



Example: Resource Leak

```
1  #include <cstddef>
2  #include <limits>
3
4  class TwoBufs {
5  public:
6      TwoBufs(std::size_t aSize, std::size_t bSize) :
7          a_(nullptr), b_(nullptr) {
8          a_ = new char[aSize];
9          // If new throws, a_ will be leaked.
10         b_ = new char[bSize];
11     }
12     ~TwoBufs() {
13         delete[] a_;
14         delete[] b_;
15     }
16     // ...
17 private:
18     char* a_;
19     char* b_;
20 };
21
22 void doWork() {
23     // This may leak memory.
24     TwoBufs x(1000000,
25             std::numeric_limits<std::size_t>::max());
26     // ...
27 }
```

Example: `std::unique_ptr`

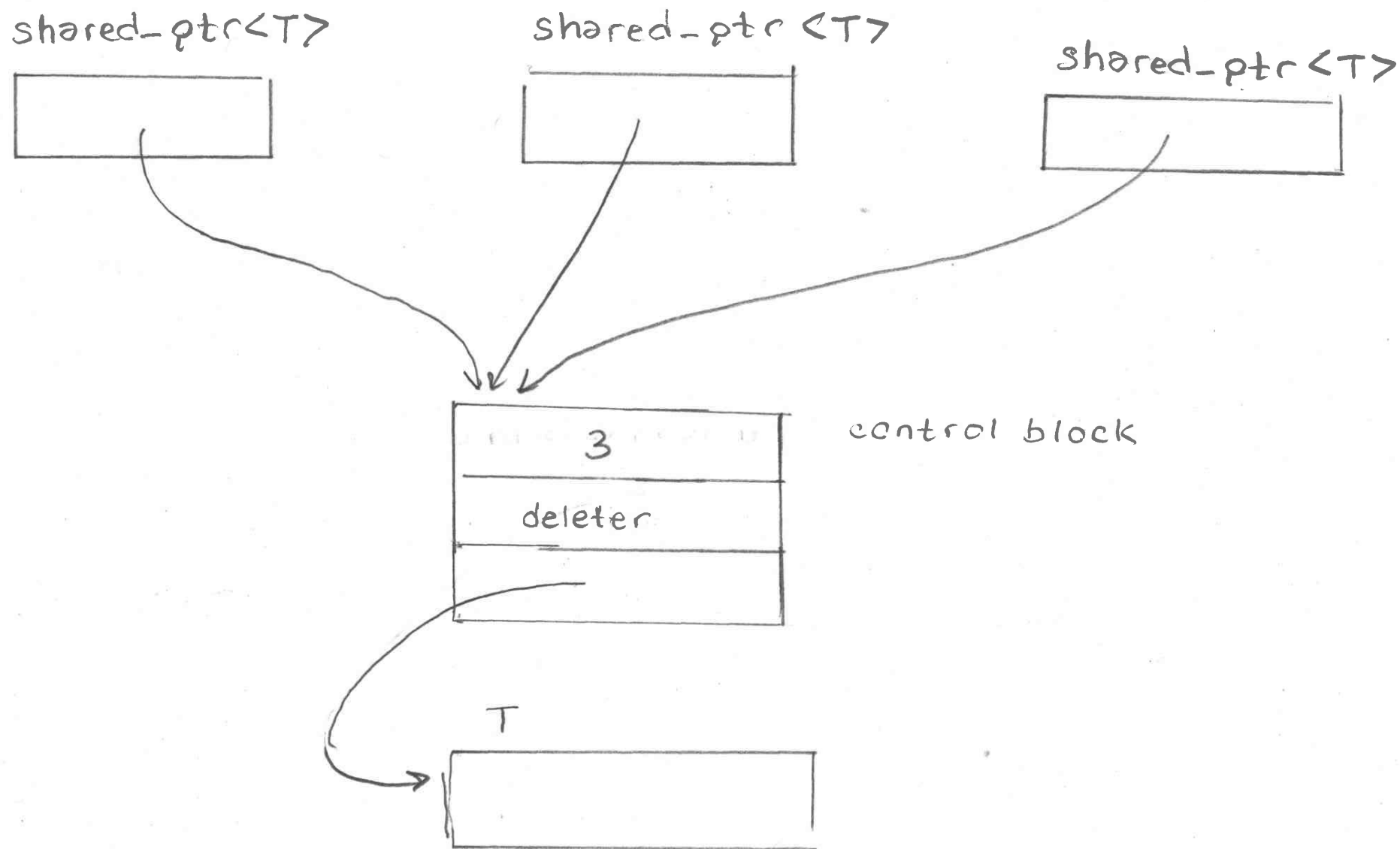
```
1  #include <cstddef>
2  #include <limits>
3  #include <memory>
4
5  class TwoBufs {
6  public:
7      TwoBufs(std::size_t aSize, std::size_t bSize) :
8          a_(std::make_unique<char[]>(aSize)),
9          b_(std::make_unique<char[]>(bSize)) {}
10     ~TwoBufs() {}
11     // ...
12 private:
13     std::unique_ptr<char[]> a_;
14     std::unique_ptr<char[]> b_;
15 };
16
17 void doWork() {
18     // This will not leak memory.
19     TwoBufs x(1000000,
20             std::numeric_limits<std::size_t>::max());
21 }
```

The `std::shared_ptr` Template Class

- `std::shared_ptr` is *smart pointer* that retains *shared* ownership of object through pointer
- declaration:

```
template <class T> class shared_ptr;
```
- T is type of object to be managed (i.e., owned object)
- multiple `shared_ptr` objects may own same object
- owned object is deleted when last remaining owning `shared_ptr` object is destroyed or last remaining owning `shared_ptr` object assigned another pointer via assignment or `reset`
- `shared_ptr` object is movable, where move transfers ownership
- `shared_ptr` object is copyable, where copy creates additional owner
- thread safety guaranteed for `shared_ptr` object itself but not owned object
- `std::make_shared` often used to create `shared_ptr` objects (for both efficiency and exception-safety reasons)
- `shared_ptr` has more overhead than `unique_ptr` so `unique_ptr` should be preferred unless shared ownership required

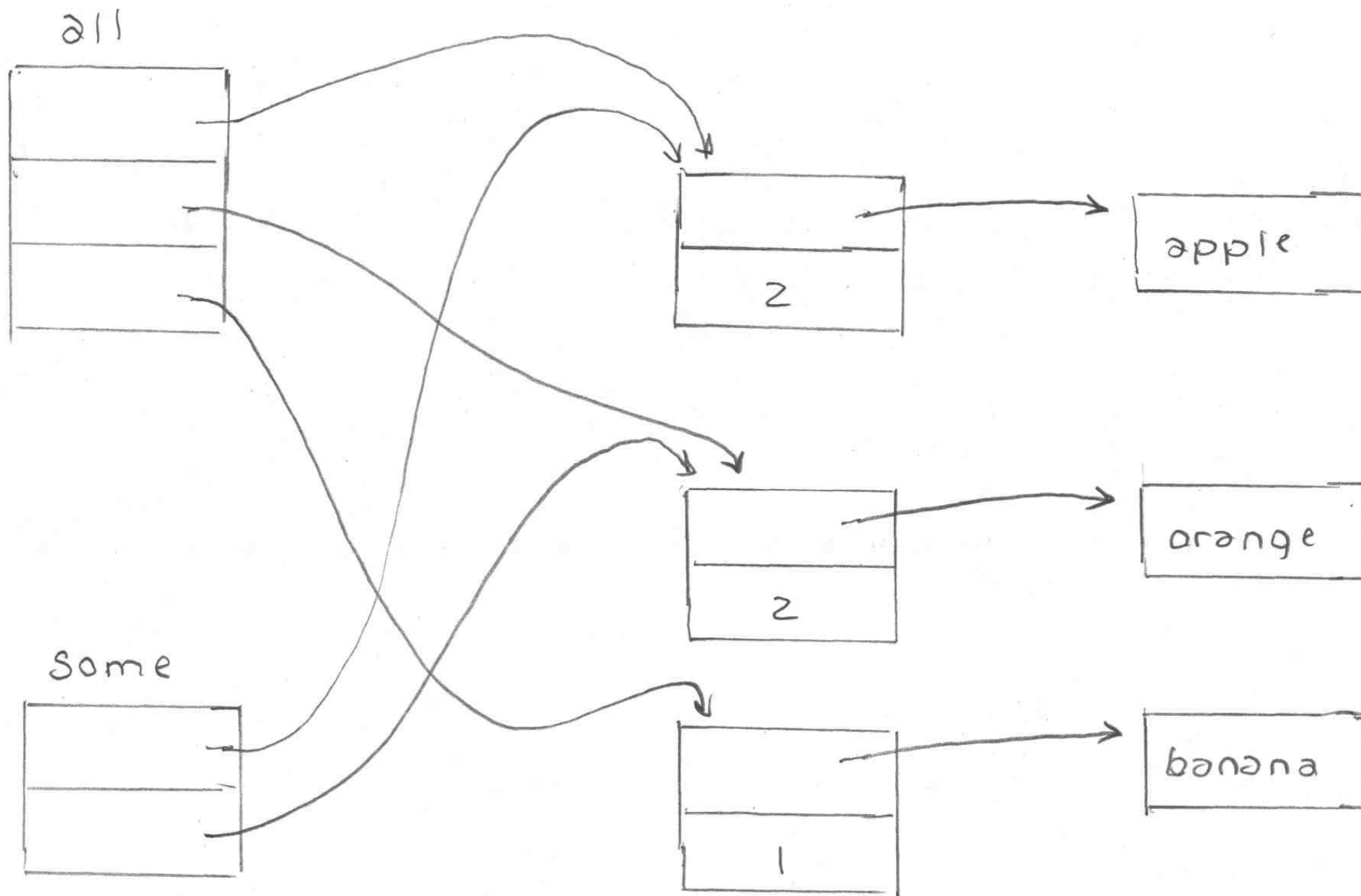
The `std::shared_ptr` Template Class (Continued)



Example: `std::shared_ptr`

```
1  #include <memory>
2  #include <vector>
3  #include <string>
4  #include <iostream>
5
6  using namespace std::literals;
7
8  int main() {
9      std::vector<std::shared_ptr<std::string>> all;
10     all.emplace_back(
11         std::make_shared<std::string>("apple"s));
12     all.emplace_back(
13         std::make_shared<std::string>("orange"s));
14     all.emplace_back(
15         std::make_shared<std::string>("banana"s));
16
17     std::vector<std::shared_ptr<std::string>> some(
18         all.begin(), all.begin() + 2);
19
20     for (auto& x : all) {
21         std::cout << *x << ' ' << x.use_count() << '\n';
22     }
23 }
24
25 /* output:
26 apple 2
27 orange 2
28 banana 1
29 */
```

Example: `std::shared_ptr` (Continued)



RAII Example: Stream Formatting Flags

```
1  #include <iostream>
2  #include <ios>
3  #include <boost/io/ios_state.hpp>
4
5  // not exception safe
6  void unsafeOutput(std::ostream& out, unsigned int x) {
7      auto flags = out.flags();
8      // if exception thrown during output of x, old
9      // formatting flags will not be restored
10     out << std::hex << std::showbase << x << '\n';
11     out.flags(flags);
12 }
13
14 // exception safe
15 void safeOutput(std::ostream& out, unsigned int x) {
16     boost::io::ios_flags_saver ifs(out);
17     out << std::hex << std::showbase << x << '\n';
18 }
```

- RAII objects can be used to save and restore state

Section 3.1.9

Exception Gotchas

shared_ptr Example: Not Exception Safe

```
1  #include <memory>
2
3  class T1 { /* ... */ };
4  class T2 { /* ... */ };
5
6  void func(std::shared_ptr<T1> p, std::shared_ptr<T2> q)
7  { /* ... */ }
8
9  void doWork() {
10     // potential memory leak
11     func(std::shared_ptr<T1>(new T1),
12         std::shared_ptr<T2>(new T2));
13     // ...
14 }
```

- one problematic order:

- 1 allocate memory for T1
- 2 construct T1
- 3 allocate memory for T2
- 4 construct T2
- 5 construct shared_ptr<T1>
- 6 construct shared_ptr<T2>
- 7 call func

- if step 3 or 4 throws, memory leaked

- another problematic order:

- 1 allocate memory for T1
- 2 allocate memory for T2
- 3 construct T1
- 4 construct T2
- 5 construct shared_ptr<T1>
- 6 construct shared_ptr<T2>
- 7 call func

- if step 3 or 4 throws, memory leaked

shared_ptr Example: Exception Safe

```
1  #include <memory>
2
3  class T1 { /* ... */ };
4  class T2 { /* ... */ };
5
6  void func(std::shared_ptr<T1> p, std::shared_ptr<T2> q)
7  { /* ... */ }
8
9  void doWork() {
10     func(std::make_shared<T1>(), std::make_shared<T2>());
11     // ...
12 }
```

- previously problematic line of code now does following:
 - 1 perform following operations in any order:
 - construct `shared_ptr<T1>` via `make_shared<T1>`
 - construct `shared_ptr<T2>` via `make_shared<T2>`
 - 2 call `func`
- each of `T1` and `T2` objects managed by `shared_ptr` at all times so no memory leak possible if exception thrown
- similar issue arises in context of `std::unique_ptr` and can be resolved by using `std::make_unique` in similar way as above

Stack Example

- stack class template parameterized on element type T

```
1  template <class T>
2  class Stack
3  {
4  public:
5      // ...
6      // Pop the top element from the stack.
7      T pop() {
8          // If the stack is empty...
9          if (top_ == start_)
10             throw "stack is empty";
11         // Remove the last element and return it.
12         return *(--top_);
13     }
14 private:
15     T* start_; // start of array of stack elements
16     T* end_;   // one past end of array
17     T* top_;   // one past current top element
18 };
```

- what is potentially problematic about this code with respect to exceptions?

Section 3.1.10

Miscellany

safe_add Example: Traditional Error Handling

```
1  #include <limits>
2  #include <vector>
3  #include <iostream>
4
5  std::pair<bool, int> safe_add(int x, int y) {
6      return ((y > 0 && x > std::numeric_limits<int>::max() - y)
7             || (y < 0 && x < std::numeric_limits<int>::min() - y)) ?
8             std::make_pair(false, 0) : std::make_pair(true, x + y);
9  }
10
11 int main() {
12     constexpr int int_min = std::numeric_limits<int>::min();
13     constexpr int int_max = std::numeric_limits<int>::max();
14     std::vector<std::pair<int, int>> v{
15         {int_max, int_max}, {1, 2}, {int_min, int_min},
16         {int_max, int_min}, {int_min, int_max}
17     };
18     for (auto x : v) {
19         auto result = safe_add(x.first, x.second);
20         if (result.first) {
21             std::cout << result.second << '\n';
22         } else {
23             std::cout << "overflow\n";
24         }
25     }
26 }
```

safe_add Example: Exceptions

```
1  #include <limits>
2  #include <vector>
3  #include <iostream>
4  #include <stdexcept>
5
6  int safe_add(int x, int y) {
7      return ((y > 0 && x > std::numeric_limits<int>::max() - y)
8             || (y < 0 && x < std::numeric_limits<int>::min() - y)) ?
9             throw std::overflow_error("addition") : x + y;
10 }
11
12 int main() {
13     constexpr int int_min = std::numeric_limits<int>::min();
14     constexpr int int_max = std::numeric_limits<int>::max();
15     std::vector<std::pair<int, int>> v{
16         {int_max, int_max}, {1, 2}, {int_min, int_min},
17         {int_max, int_min}, {int_min, int_max}
18     };
19     for (auto x : v) {
20         try {
21             int result = safe_add(x.first, x.second);
22             std::cout << result << '\n';
23         }
24         catch (const std::overflow_error&) {
25             std::cout << "overflow\n";
26         }
27     }
28 }
```

Section 3.1.11

References

References I

- 1 D. Abrahams. [Exception-safety in generic components](#). In *Lecture Notes in Computer Science*, volume 1766, pages 69–79. Springer, 2000.
A good tutorial on exception safety by an expert on the subject.
- 2 T. Cargill. [Exception handling: A false sense of security](#). *C++ Report*, 6(9), Nov. 1994.
[Available online at `http://ptgmedia.pearsoncmg.com/images/020163371x/supplements/Exception_Handling_Article.html`](http://ptgmedia.pearsoncmg.com/images/020163371x/supplements/Exception_Handling_Article.html).
An early paper that first drew attention to some of the difficulties in writing exception-safe code.
- 3 [Exception-Safe Coding in C++](http://exceptionsafecode.com), `http://exceptionsafecode.com`, 2014.
- 4 V. Kochhar, [How a C++ Compiler Implements Exception Handling](http://www.codeproject.com/Articles/2126/How-a-C-compiler-implements-exception-handling), `http://www.codeproject.com/Articles/2126/How-a-C-compiler-implements-exception-handling`, 2002.

- 1 Jon Kalb. Exception-Safe Code, CppCon, Bellevue, WA, USA, Sep 7–12, 2014. (This talk is in three parts.)
- 2 Jon Kalb. Exception-Safe Coding in C++Now, Aspen, CO, USA, May 13–18, 2012. (This talk is in two parts.)

Section 3.2

Rvalue References

Section 3.2.1

Introduction

Motivation Behind Rvalue References

- Rvalue references were added to the language in C++11 in order to provide support for:
 - ① move operations; and
 - ② perfect forwarding.
- A move operation is used to propagate the value from one object to another, much like a copy operation, except that a move operation makes fewer guarantees, allowing for greater efficiency and flexibility in many situations.
- Perfect forwarding relates to being able to pass function arguments from a template function through to another function (called by the template function) while preserving certain properties of those arguments.

Terminology: Named and Cv-Qualified

- A type that includes one or both of the qualifiers **const** and **volatile** is called a **cv-qualified type**.
- A type that is not cv-qualified is called **cv-unqualified**.
- Example:
The types **const int** and **volatile char** are cv-qualified.
The types **int** and **char** are cv-unqualified.
- An object or function that is named by an identifier is said to be **named**.
- An object or function that cannot be referred to by name is said to be **unnamed**.
- Example:

```
std::vector<int> v = {1, 2, 3, 4};  
std::vector<int> w;  
w = v; // w and v are named  
w = std::vector<int>(2, 0);  
    // w is named  
    // std::vector<int>(2, 0) is unnamed
```

Section 3.2.2

Copying and Moving

Propagating Values: Copying and Moving

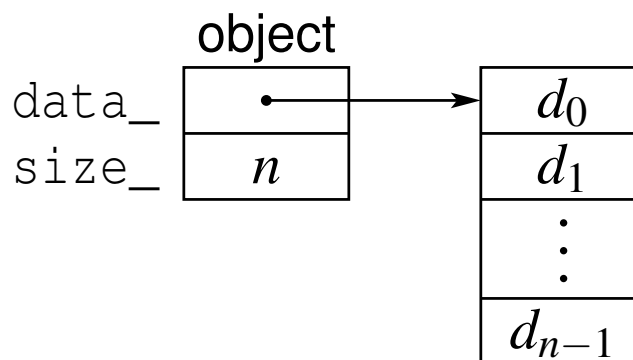
- Suppose that we have two objects of the same type and we want to propagate the value of one object (i.e., the source) to the other object (i.e., the destination).
- This can be accomplished in one of two ways:
 - 1 copying; or
 - 2 moving.
- **Copying** propagates the value of the source object to the destination object *without modifying the source object*.
- **Moving** propagates the value of the source object to the destination object and is *permitted to modify the source object*.
- Moving is always at least as efficient as copying, and for many types, moving is *more efficient* than copying.
- For some types, *copying does not make sense*, while moving does (e.g., `std::ostream`, `std::istream`).

Vector Example: Moving Versus Copying

- Consider a class that represents a one-dimensional array.

```
template <class T>
class Vector {
public:
    // ...
private:
    T* data_; // pointer to element data
              // (allocated with new)
    unsigned int size_; // number of elements
};
```

- Pictorially, the data structure looks like the following:



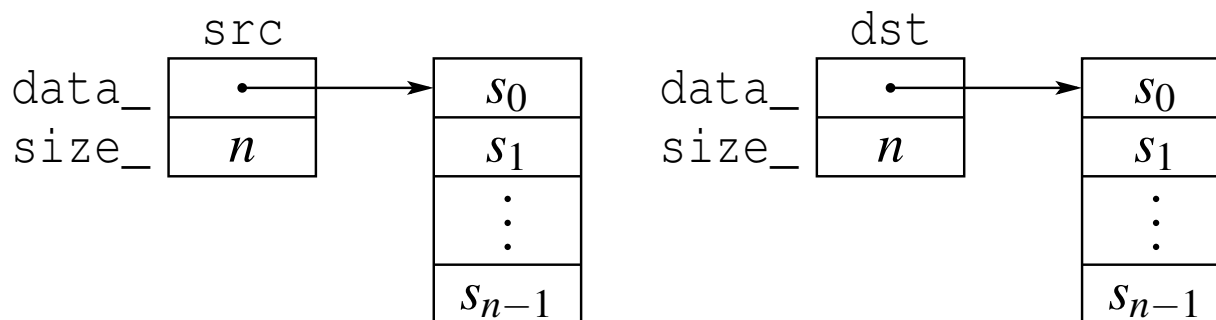
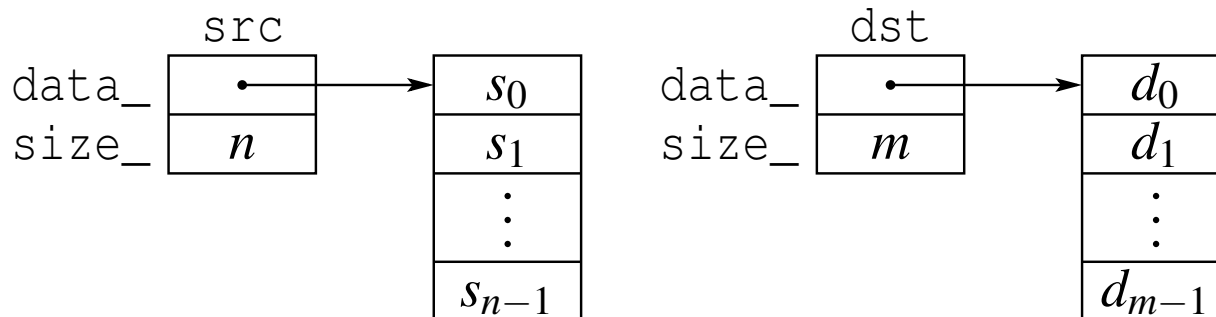
- How would copying be implemented?
- How would moving be implemented?

Vector Example: Copying

- code for copying from source `src` to destination `dst` (not self assignment):

```
delete [] dst.data_;  
dst.data_ = new T[src.size_];  
dst.size_ = src.size_;  
std::copy_n(src.data_, src.size_, dst.data_);
```

- copying requires: one array delete (destruction, memory deallocation), one array new (memory allocation, construction), copying of element data (copy assignment, etc.), and updating `data_` and `size_` data members
- copying proceeds as follows:

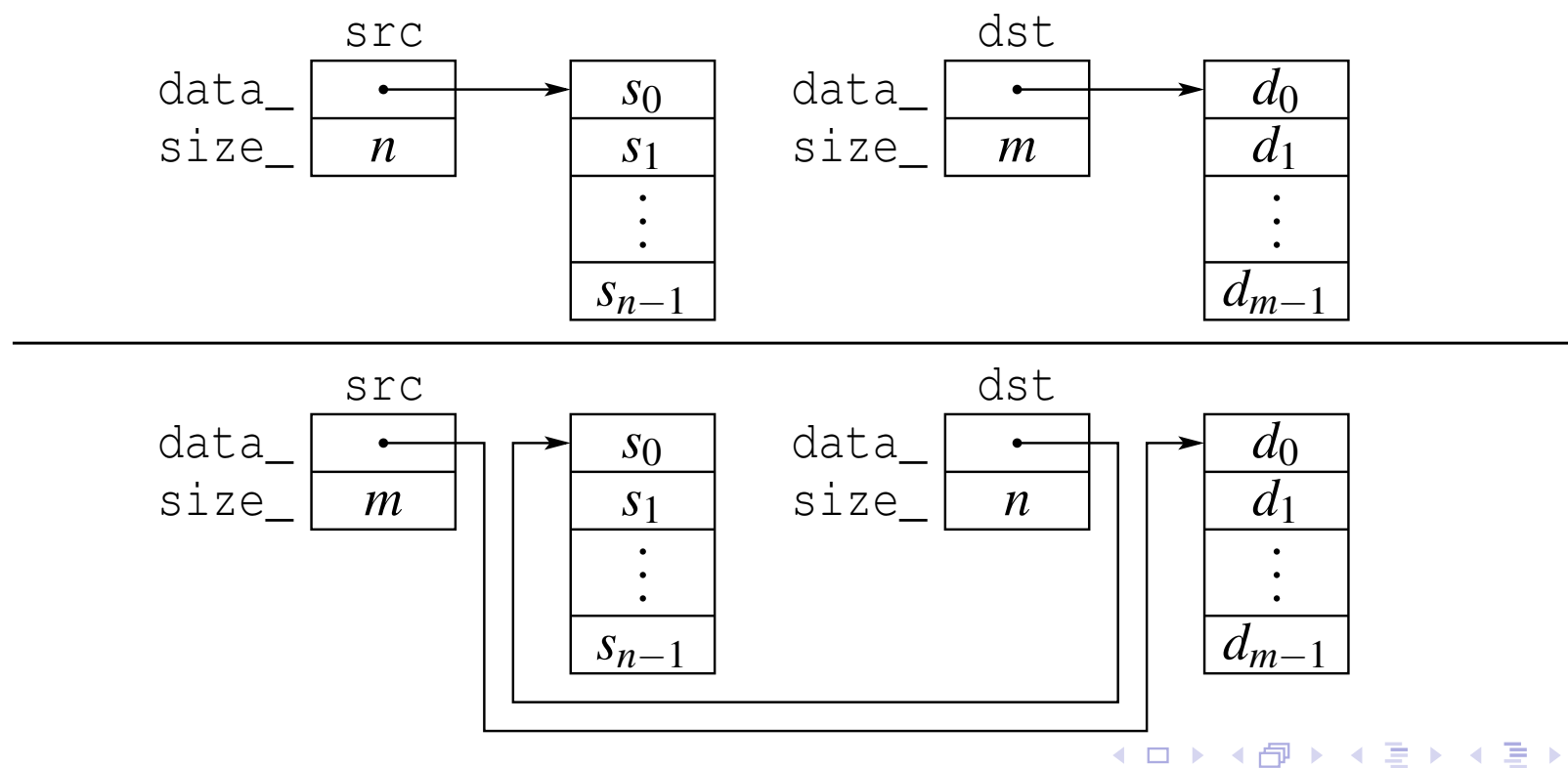


Vector Example: Moving

- code for moving from source `src` to destination `dst`:

```
std::swap(src.data_, dst.data_);  
std::swap(src.size_, dst.size_);
```

- moving only requires updating `data_` and `size_` data members
- although not considered here, could also free data array associated with `src` if desirable to release memory as soon as possible
- moving proceeds as follows:



Moving Versus Copying

- Moving is usually more efficient than copying, often by very large margin.
- So, we should prefer moving to copying.
- We can safely replace a copy by a move when subsequent code does not depend on the value of source object.
- It would be convenient if the language could provide a mechanism for automatically using a move (instead of a copy) in situations where doing so is always guaranteed to be safe.
- For reasons of efficiency, it would also be desirable for the language to provide a mechanism whereby the programmer can override the normal behavior and force a move (instead of a copy) in situations where such a transformation is known to be safe only due to some special additional knowledge about program behavior.
- Rvalue references provide the above mechanisms.

Section 3.2.3

References and Expressions

References

- A **reference** is an alias (i.e., nickname) for an already existing object.
- The language has two kinds of references:
 - 1 lvalue references
 - 2 rvalue references
- An **lvalue reference** is denoted by `&` (often read as “ref”).

```
int i = 5;  
int& j = i; // j is lvalue reference to int  
const int& k = i; k is lvalue reference to const int
```

- An **rvalue reference** is denoted by `&&` (often read as “ref ref”).

```
int&& i = 5; // i is rvalue reference to int  
const int&& j = 17; // j is rvalue reference to const int
```

- The act of initializing a reference is known as **reference binding**.
- Lvalue and rvalue references differ only in their properties relating to:
 - reference binding; and
 - overload resolution.

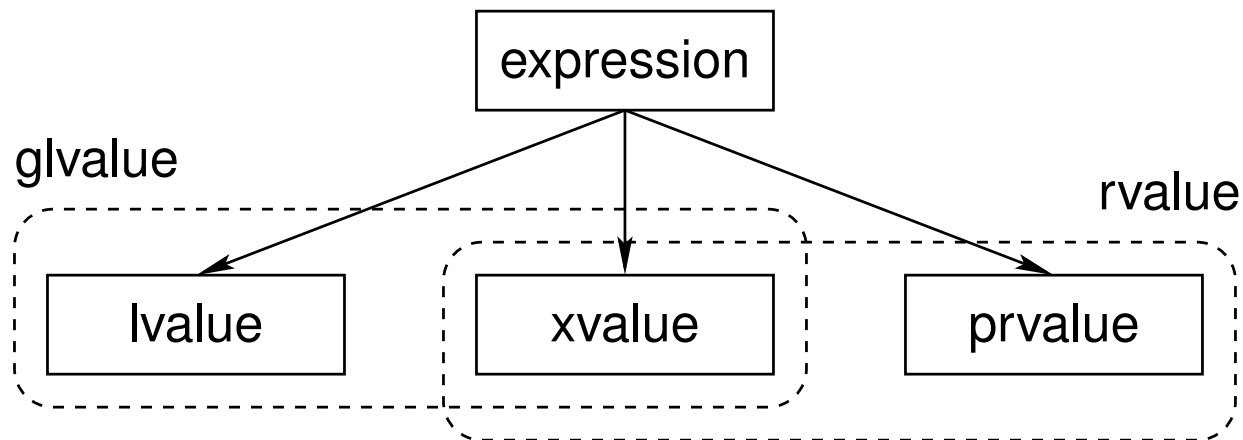
Expressions

- An **expression** is a sequence of operators and operands that specifies a computation.
- An expression has a type and, if the type is not void, a value.
- Example:

```
int x = 0;  
int y = 0;  
int* p = &x;  
double d = 0.0;  
// Evaluate some  
// expressions here.
```

Expression	Type	Value
x	int	0
y = x	int &	reference to y
x + 1	int	1
x * x + 2 * x	int	0
y = x * x	int &	reference to y
x == 42	bool	false
*p	int &	reference to x
p == &x	bool	true
x > 2 * y	bool	false
std::sin(d)	double	0.0

Categories of Expressions



- Every expression can be classified into *exactly one* of the three following categories:
 - 1 lvalue
 - 2 prvalue (pure rvalue)
 - 3 xvalue (expiring value)
- An expression that is an lvalue or xvalue is called a **glvalue** (generalized lvalue).
- An expression that is a prvalue or an xvalue is called an **rvalue**.
- Every expression is either an lvalue or an rvalue (but not both).
- Whether or not it is safe to move (instead of copy) depends on whether an lvalue or rvalue is involved.

- An **lvalue** is an expression that:
 - designates a function or object ; and
 - has an identity (i.e., occupies some *identifiable* location in memory and therefore, in principle, can have its address taken).
- *Named objects* and *named functions* are lvalues. Example:

```
int getValue();  
int i = 0;  
const int j = 1;  
i = j + 1; // i and j are lvalues  
getValue(); // getValue is lvalue [Note: not getValue()]
```

- *Dereferenced pointer*. If *e* is an expression of pointer type, then **e* is an lvalue. Example:

```
char buffer[] = "Hello";  
char* s = buffer;  
*s = 'a'; // *s is lvalue  
*(s + 1) = 'b'; // *(s + 1) is lvalue
```

Lvalues (Continued)

- The result of calling a function whose *return type is an lvalue reference type* is an lvalue. Example:

```
std::vector<int> v = {{1, 2, 3}};  
// int& std::vector<int>::operator[](int);  
int i = v[0]; // v[0] is lvalue
```

- A *string literal* is an lvalue. Example: "Hello World"
- *Named rvalue references* are lvalues. Example:

```
int&& i = 1 + 3;  
int j = i; // i is lvalue
```

- Rvalue references to functions (both named and unnamed) are lvalues.

Moving and Lvalues

- Using a move (instead of a copy) is *not guaranteed to be safe* when the source is an lvalue (since other code can access the associated object by name or through a pointer or reference).
- Example:

```
Vector<int> x;  
Vector<int> y(x);  
    // can we construct by moving (instead of copying)?  
    // source x is lvalue  
    // not safe to move x to y since value of x  
    // might be used below  
y = x;  
    // can we assign by moving (instead of copying)?  
    // source x is lvalue  
    // not safe to move x to y since value of x  
    // might be used below
```

- A **prvalue** (pure rvalue) is an expression that:
 - is a temporary object or subobject thereof, or a value that is not associated with an object; and
 - does not have an identity.
- A prvalue is a kind of rvalue.
- *Temporary objects* are prvalues. Example:

```
std::vector<int> v;  
v = std::vector<int>(10, 2);  
    // std::vector<int>(10, 2) is prvalue  
std::complex<double> u;  
u = std::complex<double>(1, 2);  
    // std::complex<double>(1, 2) is prvalue
```

- A function call whose *return type is not a reference type* is a prvalue.
Example:

```
int func();  
int i = func(); // func() is prvalue
```

Prvalues (Continued)

- All *literals other than string literals* are prvalues. Examples:

```
double pi = 3.1415; // 3.1415 is prvalue
int i = 42; // 42 is prvalue
i = 2 * i + 1; // 2 and 1 are prvalues
char c = 'A'; // 'A' is prvalue
```

- The result yielded by certain built-in operators (e.g., +, -) is a prvalue.

Example:

```
int i, j;
i = 3 + 5; // 3 + 5 is prvalue
j = i * i; // i * i is prvalue
```

- The **this** keyword is a prvalue expression.
- Prvalues need not have any storage associated with them.
- Not requiring prvalue expressions to have storage gives the compiler more freedom in generating code for such expressions.

```
int i = 2;
// 2 is prvalue and need not ever be stored in memory
```


Moving and Prvalues

- Using a move (instead of a copy) is *always safe* when the source is a prvalue (since the prvalue cannot correspond to an object with an identity).
- Example (move from temporary object):

```
Vector<int> getVector();  
Vector<int> x;  
Vector<int> y(getVector());  
    // can we construct by moving (instead of copying)?  
    // source getVector() is prvalue  
    // safe to move since temporary object could not be  
    // used below  
x = getVector();  
    // can we assign by moving (instead of copying)?  
    // source getVector() is prvalue  
    // safe to move since temporary object could not be  
    // used below
```

Xvalues

- An **xvalue** (expiring value) is an expression that:
 - refers to an object (usually near the end of its lifetime);
 - has an identity; and
 - is *deemed to be safe* to use as the source for a move.
- An xvalue is a kind of rvalue.
- An xvalue is the result of certain kinds of expressions involving rvalue references.
- The result of calling a function whose *return type is an rvalue reference type* is an xvalue. Example:

```
std::string s("Hello");  
std::string t = std::move(s); // std::move(s) is xvalue
```

- In the above example, the template function `std::move` converts its argument to an xvalue (since it returns an rvalue reference type).
- *Unnamed rvalue references to objects* are xvalues.

```
std::string s("Hello");  
std::string t;  
t = static_cast<std::string&&>(s);  
// static_cast<std::string&&>(s) is xvalue
```

Moving and Xvalues

- Using a move (instead of a copy) is *deemed to be safe* when the source is an xvalue.
- Example (forced move):

```
Vector<int> v(100, 5);  
Vector<int> u(200, -1);  
for (auto i : v) std::cout << i << '\n';  
for (auto i : u) std::cout << i << '\n';  
v = std::move(u);  
    // std::move(u) is xvalue  
    // safe to force move since later code does  
    // not to use value of u  
for (auto i : v) std::cout << i << '\n';  
    // later code known not to use value of u
```

- The function `std::move` only allows for an object to be treated as if it were safe to use as source of a move, but does not perform a move.

Moving and Lvalues and Rvalues

- With regard to propagating the value from one object to another, we can summarize the preceding results as given below.
- If the source is an *rvalue* (i.e., prvalue or xvalue), using a move instead of a copy is *always safe*.
- If the source is an *lvalue*, using a move instead of a copy is *not guaranteed to be safe*.
- It would be highly desirable if the language would provide a mechanism that would automatically allow a move to be used in the rvalue case and a copy to be employed otherwise.
- In fact, this is exactly what the language does.

More on Lvalues and Rvalues

- Lvalues and rvalues can be either *modifiable or nonmodifiable*.

Example:

```
int i = 0;
const int j = 2;
i = j + 3;
    // i is modifiable lvalue
    // j is nonmodifiable lvalue
    // j + 3 is modifiable rvalue
const std::string getString();
std::string s = getString();
    // getString() is nonmodifiable rvalue
```

- Class rvalues can have cv-qualified types, while non-class rvalues *always have cv-unqualified types*. Example:

```
const int getConstInt(); // const is ignored
const std::string getConstString();
int i = getConstInt();
    // getConstInt() is modifiable rvalue of type int
    // (not const int)
std::string s = getConstString();
    // getConstString() is nonmodifiable rvalue
```

Exercise: Expressions

```
1  #include <iostream>
2  #include <string>
3  #include <utility>
4
5  std::string&& func1(std::string& x) {
6      return std::move(x);
7      // x?  std::move(x)?
8  }
9
10 int main() {
11     const std::string hello("Hello");
12     std::string a;
13     std::string b;
14
15     a = hello + "!";
16     // hello?  hello + "!"?  a = hello + "!"?
17     std::cout << a << '\n';
18     // std::cout?  std::cout << a?
19
20     a = std::string("");
21     // std::string("")?  a = std::string("")?
22     ((a += hello) += "!");
23     // a += hello?
24     b = func1(a);
25     // func1(a)?  b = func1(a)?
26     std::cout << b << '\n';
27 }
```

Built-In Operators, Rvalues, and Lvalues

- Aside from the exceptions noted below, all of the built-in operators *require operands that are rvalues*.
- The operand of each of the following built-in operators must be an lvalue:
 - address of (i.e., unary &),
 - prefix and postfix increment (i.e., ++),
 - prefix and postfix decrement (i.e., --)
- The left operand of the following built-in operators must be an lvalue:
 - assignment (i.e., =)
 - compound assignment (e.g., +=, -=, *=, /=, etc.)
- Aside from the exceptions noted below, all of the built-in operators *yield a result that is an rvalue*.
- The following operators yield a result that is an lvalue:
 - subscript (i.e., [])
 - dereference (i.e., unary *)
 - assignment (i.e., =) and compound assignment (e.g., +=, -=, etc.)
 - prefix increment (i.e., ++) and prefix decrement (i.e., --)
 - function call (i.e., ()) invoking a function that returns a reference type
 - cast to reference type

Operators, Lvalues, and Rvalues

- Whether an operator for a *class type* requires operands that are lvalues or rvalues or yield lvalues or rvalues is determined by the parameter types and return type of the operator function.
- The member selection operator may yield an lvalue or rvalue, depending on the particular manner in which the operator is used. (The behavior is fairly intuitive.)
- The lvalue/rvalue-ness and type of the result produced by the ternary conditional operator depends on the particular manner in which the operator is employed.

Implicit Lvalue-to-Rvalue Conversion

- An implicit conversion from lvalues to rvalues is provided, which can be used in most (but not all) circumstances.
- Example:

```
int i = 1;  
int j = 2;  
int k = i + j;  
    // operands of + must be rvalues  
    // i and j converted to rvalues
```

Section 3.2.4

Reference Binding and Overload Resolution

References: Binding and Overload Resolution

- The kinds of expressions, to which lvalue and rvalue references can *bind*, differ.
- For a nonreference type `T` (such as `int` or `const int`), what kinds of expressions can validly be placed in each of the boxes in the example below?

```
T& r =  ;  
T&& r =  ;
```

- Lvalue and rvalue references also behave differently with respect to *overload resolution*.
- Let `T` be a cv-unqualified nonreference type. Which overloads of `func` will be called in the example below?

```
T operator+(const T&, const T&);  
void func(const T&);  
void func(T&&);  
T x;  
func(x); // calls which version of func?  
func(x + x); // calls which version of func?
```

Reference Binding

- Implicit lvalue-to-rvalue conversion is disabled when binding to references.
- An lvalue reference can bind to an lvalue as long as doing so would not result in the *loss* of any cv qualifiers.

```
const int i = 0;
int& r1 = i; // ERROR: drops const
const int& r2 = i; // OK
const volatile int& r3 = i; // OK
```

- The loss of cv qualifiers must be avoided for *const and volatile correctness*.
- Similarly, an rvalue reference can bind to an rvalue as long as doing so would not result in the *loss* of any cv qualifiers.

```
const std::string getValue();
std::string&& r1 = getValue(); // ERROR: drops const
const std::string&& r2 = getValue(); // OK
```

- Again, the loss of cv qualifiers must be avoided for *const and volatile correctness*.

Reference Binding (Continued)

- An lvalue reference can be bound to an rvalue only if doing so would not result in the *loss* of any cv qualifier and the lvalue reference is *const*.

```
const std::string getConstValue();  
std::string& r1 = getConstValue(); // ERROR: drops const  
const std::string& r2 = getValue(); // OK  
int& ri1 = 42; // ERROR: not const reference  
const int& ri2 = 42; // OK
```

- The requirement that the lvalue reference be const is to prevent temporary objects from being modified in a very uncontrolled manner, which can lead to subtle bugs.
- An rvalue reference can *never* be bound to an lvalue.

```
int i = 0;  
int&& r1 = i; // ERROR: cannot bind to lvalue  
int&& r2 = 42; // OK
```

- Allowing rvalue reference to bind to lvalues would violate the principle of type-safe overloading, which can lead to subtle bugs.

Why Rvalue References Cannot Bind to Lvalues

- In effect, rvalue references were introduced into the language to allow a function to know if one of its reference parameters is bound to an object whose value is safe to change without impacting other code, namely, an rvalue (i.e., a temporary object or xvalue).
- Since an rvalue reference can only bind to an rvalue, any rvalue reference parameter to a function is *guaranteed* to be bound to a temporary object or xvalue.

- Example:

```
class Thing {  
public:  
    // Move constructor  
    // parameter x known to be safe to use as source for move  
    Thing(Thing&& x);  
    // Move assignment operator  
    // parameter x known to be safe to use as source for move  
    Thing& operator=(Thing&& x);  
    // ...  
};  
// parameter x known to be safe to modify  
void func(Thing&& x);
```

- If rvalue references could bind to lvalues, the above guarantee could not be made, as an rvalue reference could then refer to an object whose value cannot be changed safely, namely, an lvalue.

Why Non-Const Lvalue References Cannot Bind to Rvalues

- If non-const lvalue references could bind to rvalues, temporary objects could be modified in many undesirable circumstances.

```
void func(int& x) {  
    // ...  
}  
  
int main() {  
    int i = 1;  
    int j = 2;  
    func(i + j);  
    // ERROR: cannot bind non-const lvalue  
    // reference to rvalue  
    // What would be consequence if allowed?  
}
```

Reference Binding Summary

	Rvalue				Lvalue			
	T	const T	volatile T	const volatile T	T	const T	volatile T	const volatile T
	T&&	✓	C	V	C,V	X	X	X
const T&&	✓	✓	V	V	X	X	X	X
volatile T&&	✓	X	✓	C	X	X	X	X
const volatile T&&	✓	✓	✓	✓	X	X	X	X
T&	X	X	X	X	✓	C	V	C,V
const T&	✓	✓	V	V	✓	✓	V	V
volatile T&	X	X	X	X	✓	C	✓	C
const volatile T&	X	X	X	X	✓	✓	✓	✓

✓: allowed C: strips const V: strips volatile X: other

Reference Binding Example

```
1  #include <string>
2  using std::string;
3
4  string value() {
5      return string("Hello");
6  }
7
8  const string constValue() {
9      return string("World");
10 }
11
12 int main() {
13     string i("mutable");
14     const string j("const");
15
16     string& r01 = i;
17     string& r02 = j; // ERROR: drops const
18     string& r03 = value(); // ERROR: non-const lvalue reference from rvalue
19     string& r04 = constValue(); // ERROR: non-const lvalue reference from rvalue
20
21     const string& r05 = i;
22     const string& r06 = j;
23     const string& r07 = value();
24     const string& r08 = constValue();
25
26     string&& r09 = i; // ERROR: rvalue reference from lvalue
27     string&& r10 = j; // ERROR: rvalue reference from lvalue
28     string&& r11 = value();
29     string&& r12 = constValue(); // ERROR: drops const
30
31     const string&& r13 = i; // ERROR: rvalue reference from lvalue
32     const string&& r14 = j; // ERROR: rvalue reference from lvalue
33     const string&& r15 = value();
34     const string&& r16 = constValue();
35 }
```

Overload Resolution

- Lvalues strongly prefer binding to lvalue references.
- Rvalues strongly prefer binding to rvalue references.
- Modifiable lvalues and rvalues weakly prefer binding to non-const references.

Overload Resolution Summary

	Priority							
	Rvalue				Lvalue			
	T	const T	volatile T	const volatile T	T	const T	volatile T	const volatile T
T&&	1							
const T&&	2	1						
volatile T&&	2		1					
const volatile T&&	3	2	2	1				
T&					1			
const T&	4	3			2	1		
volatile T&					2		1	
const volatile T&					3	2	2	1

Overloading Example 1

```
1  #include <iostream>
2  #include <string>
3
4  void func(std::string& x) {
5      std::cout << "func(std::string&) called\n";
6  }
7
8  void func(const std::string& x) {
9      std::cout << "func(const std::string&) called\n";
10 }
11
12 void func(std::string&& x) {
13     std::cout << "func(std::string&&) called\n";
14 }
15
16 void func(const std::string&& x) {
17     std::cout << "func(const std::string&&) called\n";
18 }
19
20 const std::string&& constValue(const std::string&& x) {
21     return static_cast<const std::string&&>(x);
22 }
23
24 int main() {
25     const std::string cs("hello");
26     std::string s("world");
27     func(s);
28     func(cs);
29     func(cs + s);
30     func(constValue(cs + s));
31 }
32
33 /* Output:
34 func(std::string&) called
35 func(const std::string&) called
36 func(std::string&&) called
37 func(const std::string&&) called
38 */
```

Overloading Example 2

```
1  #include <iostream>
2  #include <string>
3
4  void func(const std::string& x) {
5      std::cout << "func(const std::string&) called\n";
6  }
7
8  void func(std::string&& x) {
9      std::cout << "func(std::string&&) called\n";
10 }
11
12 const std::string&& constValue(const std::string&& x) {
13     return static_cast<const std::string&&>(x);
14 }
15
16 int main() {
17     const std::string cs("hello");
18     std::string s("world");
19     func(s);
20     func(cs);
21     func(cs + s);
22     func(constValue(cs + s));
23 }
24
25 /* Output:
26 func(const std::string&) called
27 func(const std::string&) called
28 func(std::string&&) called
29 func(const std::string&) called
30 */
```

Why Rvalue References Cannot Bind to Lvalues (Revisited)

- If an rvalue reference could bind to an lvalue, this would violate the principle of type-safe overloading.

```
1  #include <iostream>
2  #include <string>
3
4  template <class T>
5  class Container {
6  public:
7      // ...
8      // Forget to provide the following function:
9      // void push_back(const T& value); // Copy semantics
10     void push_back(T&& value); // Move semantics
11 private:
12     // ...
13 };
14
15 int main() {
16     std::string s("Hello");
17     Container<std::string> c;
18     // What would happen here if lvalues
19     // could bind to rvalue references?
20     c.push_back(s);
21     std::cout << s << '\n';
22 }
```

Section 3.2.5

Moving

Move Constructors

- A non-template constructor for class `T` is a **move constructor** if it can be called with one parameter that is of type `T&&`, **const** `T&&`, **volatile** `T&&`, or **const volatile** `T&&`.
- Example (assuming no optimization):

```
struct T {  
    T();  
    T(const T&); // copy constructor  
    T(T&&); // move constructor  
};  
T func(int);  
  
T a(func(1)); // calls T::T(T&&)  
T b = a; // calls T::T(const T&)
```


Move Assignment Operators

- A **move assignment operator** `T::operator=` is a non-static non-template member function of class `T` with exactly one parameter of type `T&&`, `const T&&`, `volatile T&&`, or `const volatile T&&`.
- Example (assuming no optimization):

```
class T {
public:
    T();
    T(const T&); // copy constructor
    T(T&&); // move constructor
    T& operator=(const T&); // copy assignment operator
    T& operator=(T&&); // move assignment operator
    // ...
};

T func(int);

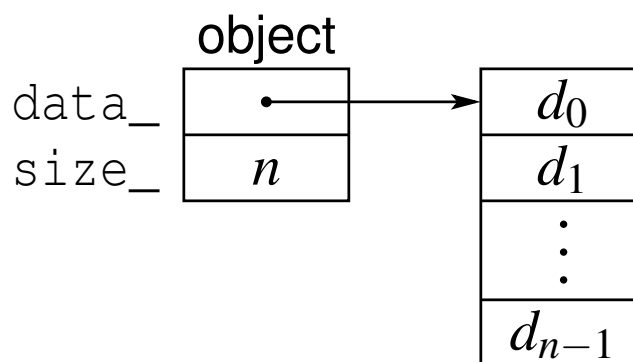
T a;
T b;
a = func(1); // calls T::operator=(T&&)
b = a; // calls T::operator=(const T&)
```

Vector Example Revisited

- Recall the class from earlier that represents a one-dimensional array.

```
template <class T>
class Vector {
public:
    // ...
private:
    T* data_; // pointer to element data
              // (allocated with new)
    unsigned int size_; // number of elements
};
```

- Pictorially, the data structure looks like the following:



Example Without Move Construction/Assignment

```
1  #include <algorithm>
2  #include <complex>
3
4  template <class T>
5  class Vector {
6  public:
7      Vector(unsigned int size, T value = 0) : data_(new T[size]), size_(size)
8          {std::fill_n(data_, size, value);}
9      Vector(const Vector& a) : data_(new T[a.size_]), size_(a.size_)
10         {std::copy_n(a.data_, a.size_, data_);}
11      Vector& operator=(const Vector& a) {
12         if (this != &a) {
13             delete[] data_; size_ = a.size_; data_ = new T[a.size_];
14             std::copy_n(a.data_, a.size_, data_);
15         }
16         return *this;
17     }
18     ~Vector() {delete[] data_;}
19     // ...
20 private:
21     T* data_; // pointer to element data
22     unsigned int size_; // number of elements
23 };
24 typedef Vector<std::complex<double>> Vec;
25 Vec getVector() {return Vec(1000, {0.0, 1.0});}
26
27 int main() {
28     Vec v(0);
29     Vec w = getVector(); // construct from temporary object
30     v = Vec(2000, {1.0, 2.0}); // assign from temporary object
31 }
```

Example With Move Construction/Assignment

```
1  #include <algorithm>
2  #include <complex>
3
4  template <class T>
5  class Vector {
6  public:
7      Vector(unsigned int size, T value = 0) : data_(new T[size]), size_(size)
8          {std::fill_n(data_, size, value);}
9      Vector(const Vector& a) : data_(new T[a.size_]), size_(a.size_)
10         {std::copy_n(a.data_, a.size_, data_);}
11      Vector& operator=(const Vector& a) {
12          if (this != &a) {
13              delete[] data_; size_ = a.size_; data_ = new T[a.size_];
14              std::copy_n(a.data_, a.size_, data_);
15          }
16          return *this;
17      }
18      // Move constructor
19      Vector(Vector&& a) : data_(a.data_), size_(a.size_)
20          {a.size_ = 0; a.data_ = nullptr;}
21      // Move assignment operator
22      Vector& operator=(Vector&& a) {
23          std::swap(size_, a.size_); std::swap(data_, a.data_);
24          return *this;
25      }
26      ~Vector() {delete[] data_;}
27      // ...
28  private:
29      T* data_; // pointer to element data
30      unsigned int size_; // number of elements
31  };
32  typedef Vector<std::complex<double>> Vec;
33  Vec getVector() {return Vec(1000, {0.0, 1.0});}
34
35  int main() {
36      Vec v(0);
37      Vec w = getVector(); // construct from temporary object
38      v = Vec(2000, {1.0, 2.0}); // assign from temporary object
39  }
```

Allowing Move Semantics in Other Contexts

- As we have seen, a reference parameter of a function that is bound to modifiable rvalue can be modified safely (i.e., no observable change in behavior outside of function).
- Sometimes may want to allow a move to be used instead of a copy, when this would not normally be permitted.
- We can allow moves by casting to a non-const rvalue reference.
- This casting can be accomplished by `std::move`, which is declared (in the header file `utility`) as:

```
template <class T>
constexpr typename std::remove_reference<T>::type&&
    move(T&&) noexcept;
```

- For an object `x` of type `T`, `std::move(x)` is similar to `static_cast<T&&>(x)` but saves typing and still works correctly when `T` is a reference type (a technicality yet to be discussed).

Old-Style Swap

- Prior to C++11, a swap function (such as `std::swap`) would typically look like this:

```
1  template <class T>
2  void swap(T& x, T& y) {
3      T tmp(x); // copy x to tmp
4      x = y;    // copy y to a
5      y = tmp;  // copy tmp to y
6  }
```

- In the above code, a swap requires three *copy* operations (namely, one copy constructor call and two copy assignment operator calls).
- For many types `T`, this use of copying is *very inefficient*.
- Furthermore, the above code requires that `T` *must be copyable* (i.e., `T` has a copy constructor and copy assignment operator).
- In C++11, we can write a much better swap function.

Improved Swap

- As of C++11, a swap function would typically look like this:

```
1  template <class T>
2  void swap(T& x, T& y) {
3      T tmp(std::move(x)); // move x to tmp
4      x = std::move(y);   // move y to x
5      y = std::move(tmp); // move tmp to y
6  }
```

- The function `std::move` casts its argument to an rvalue reference.
- Assuming that `T` provides a move constructor and move assignment operator, a swap requires three *move* operations (i.e., one move constructor call and two move assignment operator calls) and *no copying*.
- The use of `std::move` above is essential in order for copying to be avoided.

Why Distinguish Between Lvalues and Rvalues

- By distinguishing between lvalues and rvalues, we can write more efficient code.

- Scenario 1:

```
void doSomething(std::complex<double>& z) {  
    // can the caller detect a change in z?  
}  
  
std::complex<double> z(1.0, 0.0);  
doSomething(z);
```

- Scenario 2:

```
void doSomething(std::complex<double>&& z) {  
    // can the caller detect a change in z?  
}  
  
doSomething(std::complex<double>(1.0, 2.0));
```

- A function parameter that is bound to a modifiable rvalue can be changed without any observable effect outside the function.
- This gives us more freedom in how we deal with the object whose change in value cannot be observed.
- For example, this freedom can be used to replace some copies by moves.

Reference-Qualified Member Functions

- every nonstatic member function has implicit parameter `*this`
- possible to provide reference qualifiers for implicit parameter
- allows overloading member functions on lvalueness/rvalueness of `*this`
- cannot mix reference qualifiers and non-reference qualifiers in single overload set
- provides mechanism for treating lvalue and rvalue cases differently
- useful for facilitating move semantics or preventing operations not appropriate for lvalues or rvalues

Reference-Qualified Member Functions Example

```
1  #include <iostream>
2
3  class Widget {
4  public:
5      void func() const &
6          {std::cout << "const lvalue\n";}
7      void func() &
8          {std::cout << "non-const lvalue\n";}
9      void func() const &&
10         {std::cout << "const rvalue\n";}
11     void func() &&
12         {std::cout << "non-const rvalue\n";}
13 };
14
15 const Widget getConstWidget() {return Widget();}
16
17 int main() {
18     Widget w;
19     const Widget cw;
20     w.func(); // non-const lvalue
21     cw.func(); // const lvalue
22     Widget().func(); // non-const rvalue
23     getConstWidget().func(); // const rvalue
24 }
```

Lvalueness/Rvalueness and the `*this` Parameter

```
1  class Int {
2  public:
3      Int(int x = 0) : value_(x) {}
4      // only allow prefix increment for lvalues
5      Int& operator++() & {++value_; return *this;}
6      // The following allows prefix increment for rvalues:
7      // Int& operator++() {++value_; return *this;}
8      // ...
9  private:
10     int value_;
11 };
12
13 int one() {return 1;}
14
15 int main() {
16     int i = 0;
17     int j = ++i; // OK
18     // int k = ++one(); // ERROR (not lvalue)
19     Int x(0);
20     Int y = ++x; // OK
21     // Int z = ++Int(1); // ERROR (not lvalue)
22 }
```

Move Semantics and the `*this` Parameter

```
1  #include <iostream>
2  #include <vector>
3  #include <utility>
4
5  class Buffer {
6  public:
7      Buffer(char value = 0) : data_(1024, value) {}
8      void data(std::vector<char>& x) const &
9          {x = data_;}
10     void data(std::vector<char>& x) &&
11         {x = std::move(data_);}
12         // ...
13 private:
14     std::vector<char> data_;
15 };
16
17 Buffer getBuffer() {return Buffer(42);}
18
19 int main() {
20     std::vector<char> d;
21     Buffer buffer;
22     buffer.data(d); // copy into d
23     getBuffer().data(d); // move into d
24 }
```

Section 3.2.6

Reference Collapsing and Forwarding References

References to References

- A reference to a reference is not allowed, since such a construct clearly makes no sense.

```
int i = 0;  
int& & j = i; // ILLEGAL: reference to reference
```

- Although one cannot directly create a reference to a reference, a reference to a reference can arise indirectly in several contexts.
- Typedef name:

```
typedef int& RefToInt;  
typedef RefToInt& T; // reference to reference
```

- Template function parameters:

```
template <class T> T func(const T& x) {return x;}  
int x = 1;  
func<int&>(x); // reference to reference
```

- Decltype specifier:

```
int i = 1;  
decltype((i))& j = i; // reference to reference
```

References to References (Continued)

- Auto specifier:

```
int i = 0;  
auto&& j = i; // reference to reference
```

- Class templates:

```
template <class T>  
struct Thing {  
    void func(T&&) {} // reference to reference  
                        // if T is reference type  
};  
Thing<int&> x;
```

- If, during type analysis, a reference to a reference type is obtained, the reference to reference is converted to a simple reference via a process called **reference collapsing**.

Reference Collapsing Rules

- Let TR denote a type that is a reference to type T (where T may be cv qualified).
- The effect of reference collapsing is summarized below. .

Before Collapse	After Collapse
TR&	T&
const TR&	T&
volatile TR&	T&
const volatile TR&	T&
TR&&	TR
const TR&&	TR
volatile TR&&	TR
const volatile TR&&	TR

- In other words:
 - An lvalue reference to any reference yields an lvalue reference.
 - An rvalue reference to an lvalue reference yields an lvalue reference.
 - An rvalue reference to an rvalue reference yields rvalue reference.
 - Any cv qualifiers applied to a reference type are discarded (since cv qualifiers cannot be applied to a reference).

Reference Collapsing Examples

- Due to reference collapsing, T&& syntax may not always be an rvalue reference. Example:

```
typedef int& IntRef;  
int i = 0;  
IntRef&& r = i; // r is int& (i.e., lvalue reference)
```

- Example:

```
typedef int& IntRef;  
typedef int&& IntRefRef;  
typedef const int&& ConstIntRefRef;  
typedef const int& ConstIntRef;  
typedef const IntRef& T1; // T1 is int&  
typedef const IntRefRef& T2; // T2 is int&  
typedef IntRefRef&& T3; // T3 is int&&  
typedef ConstIntRef&& T4; // T4 is const int&  
typedef ConstIntRefRef&& T5; // T5 is const int&&
```

- Example:

```
int i = 0;  
int& j = i;  
auto&& k = j;  
// j cannot be inferred to have type int  
// since rvalue reference cannot be bound to lvalue  
// j inferred to have type int&  
// reference collapsing of int& && yields int&
```

Forwarding References

- A cv-unqualified rvalue reference that appears in a type-deducing context for template parameters is called a **forwarding reference**.
- Type deduction for template parameters of template functions is defined in such a way as to facilitate perfect forwarding.
- Consider the following template-parameter type-deduction scenario:

```
template<class T>
void f(T&& p);

f(expr); // invoke f
```

- Let E denote the type of the expression $expr$. The type T is then deduced as follows:
 - 1 If $expr$ is an **lvalue**, T is deduced as $E\&$, in which case the type of p yielded by reference collapsing is $E\&$.
 - 2 If $expr$ is an **rvalue**, T is deduced as E , in which case p will have the type $E\&\&$.
- Thus, the type $T\&\&$ will be an lvalue reference type if $expr$ is an lvalue, and an rvalue reference type if $expr$ is an rvalue.
- Therefore, the lvalue/rvalue-ness of $expr$ can be determined **inside** f based on whether $T\&\&$ is an lvalue reference type or rvalue reference type.

Forwarding References Example

```
1  #include <utility>
2
3  template <class T> void f(T&& p);
4  int main() {
5      int i = 42;
6      const int ci = i;
7      const int& rci = i;
8      f(i);
9          // i is lvalue with type int
10         // T is int&
11         // p has type int&
12     f(ci);
13         // ci is lvalue with type const int
14         // T is const int&
15         // p has type const int&
16     f(rci);
17         // rci is lvalue with type const int&
18         // T is const int&
19         // p has type const int&
20     f(2);
21         // 2 is rvalue with type int
22         // T is int
23         // p has type int&&
24     f(std::move(i));
25         // std::move(i) is rvalue with type int&&
26         // T is int
27         // p has type int&&
28 }
```

Section 3.2.7

Perfect Forwarding

- **Perfect forwarding** is the act of passing a template function's arguments to another function:
 - without rejecting any arguments that can be passed to that other function
 - without losing any information about the arguments' cv-qualifications or lvalue/rvalue-ness; and
 - without requiring overloading.
- In C++03, for example, the best approximations of perfect forwarding turn all rvalues into lvalues and require at least two (and often more) overloads.

Perfect-Forwarding Example

- Consider a *template* function `wrapper` and another function `func`, each of which takes one argument.
- Suppose that we want to perfectly forward the argument of `wrapper` to `func`.
- The function `wrapper` is to do nothing other than simply call `func`.
- In doing so, `wrapper` must pass its actual argument through to `func`.
- This must be done in such a way that the argument to `wrapper` and argument to `func` have *identical properties* (i.e., match in terms of cv-qualifiers and lvalue/rvalue-ness).
- In other words, the following two function calls must have *identical behavior*, where *expr* denotes an arbitrary expression:

```
wrapper (expr) ;  
func (expr) ;
```

- The solution to a perfect-forwarding problem, such as this one, turns out to be more difficult than it might first seem.

Perfect-Forwarding Example: First Failed Attempt

- For our first attempt, we propose the following code for the (template)

function wrapper:

```
template <class T>
void wrapper(T p) {
    func(p);
}
```

- If `func` takes its parameter by reference, calls to `wrapper` and `func` (with the same argument) can have different behaviors.
- Suppose, for example, that we have the following declarations:

```
void func(int&); // uses pass by reference
int i;
```

- Then, the following two function calls are *not equivalent*:

```
wrapper(i);
// T is deduced as int
// copy of i passed to func
// wrapper cannot change i
```

```
func(i);
// i passed by reference
// func can change i
```

- Problem: The original and forwarded arguments are *distinct objects*.

Perfect-Forwarding Example: Second Failed Attempt

- For our second attempt, we propose the following code for the (template)

function wrapper:

```
template <class T>
void wrapper(T& p) {
    func(p);
}
```

- If, for example, the function argument is an rvalue (such as a non-string literal or temporary object), calls to `wrapper` and `func` (with the same argument) can have different behaviors.
- Suppose, for example, that we have the following declaration:

```
void func(int); // uses pass by value
```

- Then, the following two function calls are *not equivalent*:

```
wrapper(42);
// T is deduced as int
// ERROR: cannot bind rvalue to
// nonconst lvalue reference

func(42);
// OK
```

- Problem: The original and forwarded arguments do not match in terms of *lvalue/rvalue-ness*.

Perfect-Forwarding Example: Third Failed Attempt

- For our third attempt, we propose the following code for the (template) function `wrapper`:

```
template <class T>
void wrapper(const T& p) {
    func(p);
}
```

- If, for example, the function argument is a non-const object, calls to `wrapper` and `func` (with the same argument) will have different behaviors.
- Suppose, for example, that we have the following declaration:

```
void func(int&);
int i;
```

- Then, the following two function calls are *not equivalent*:

```
wrapper(i);
// ERROR: wrapper cannot call func, as this
// would discard const qualifier

func(i);
// OK
```

- Problem: The original and forwarded arguments do not match in terms of *cv-qualifiers*.

Perfect-Forwarding Example: Solution

- Finally, we propose the following code for the (template) function `wrapper`:

```
template <class T>
void wrapper(T&& p) {
    func(static_cast<T&&>(p));
}
```

- Consider now, for example, the following scenario:

```
int i = 42;
const int ci = i;
int& ri = i;
const int& rci = i;
wrapper(expr); // invoke wrapper
```

- The parameter `p` is an alias for the object yielded by the expression `expr`.
- The argument `expr` and argument to `func` match in terms of cv-qualifiers and lvalue/rvalue-ness.

<i>expr</i>	<i>expr</i>		T	argument to <code>func</code>	
	Type	Category		Type (T&&)	Category
<code>i</code>	int	lvalue	int&	int&	lvalue
<code>ci</code>	const int	lvalue	const int&	const int&	lvalue
<code>ri</code>	int&	lvalue	int&	int&	lvalue
<code>rci</code>	const int&	lvalue	const int&	const int&	lvalue
<code>42</code>	int	rvalue	int	int&&	rvalue

Perfect-Forwarding Example: Solution (Continued)

- Although we only considered one specific scenario on the previous slide, the solution works in general.
- That is, the `wrapper` function from the previous slide will perfectly forward its single argument, regardless of what the argument happens to be (or which overload of `func` is involved).
- Thus, we have a general solution to the perfect-forwarding problem in the single-argument case.
- This solution is easily extended to an arbitrary number of arguments.

The `std::forward` Template Function

- To avoid the need for an explicit type-cast operation when forwarding an argument, the standard library provides the `std::forward` function specifically for performing such a type conversion.

- The template function `forward` is defined as:

```
template<class T>
T&& forward(typename std::remove_reference<T>::type& x)
    noexcept {
    return static_cast<T&&>(x);
}
```

- A typical usage of `forward` might look something like:

```
template <class T1, class T2>
void wrapper(T1&& x1, T2&& x2) {
    func(std::forward<T1>(x1), std::forward<T2>(x2));
}
```

- The expression `forward<T>(a)` is an lvalue if `T` is an lvalue reference type and an rvalue otherwise.
- The use of `std::forward` instead of an explicit type cast improves code readability by making the programmer's intent clear.

Perfect-Forwarding Example Revisited

- We now revisit the perfect-forwarding example from earlier.
- In the earlier example, perfect forwarding was performed by the following function:

```
template <class T>
void wrapper(T&& e) {
    func(static_cast<T&&>(e));
}
```

- The above code can be made more readable, however, by rewriting it to make use of `std::forward` as follows:

```
template <class T>
void wrapper(T&& e) {
    func(std::forward<T>(e));
}
```

Forwarding Example

```
1  #include <iostream>
2  #include <string>
3  #include <utility>
4
5  void func(std::string& s) {
6      std::cout << "func(std::string&) called\n";
7  }
8
9  void func(std::string&& s) {
10     std::cout << "func(std::string&&) called\n";
11 }
12
13 template <class T>
14 void wrapper(T&& x) {
15     func(std::forward<T>(x));
16 }
17
18 template <class T>
19 void buggy_wrapper(T x) {func(x);}
20
21 int main() {
22     using namespace std::literals;
23     std::string s("Hi"s);
24     wrapper(s);           // which overload of func called?
25     buggy_wrapper(s);    // which overload of func called?
26     wrapper("Hi"s);     // which overload of func called?
27     buggy_wrapper("Hi"s); // which overload of func called?
28 }
```

Perfect-Forwarding Use Case: Wrapper Functions

- A **wrapper function** is simply a function used to invoke another function, possibly with some additional processing.
- Example:

```
1  #include <iostream>
2  #include <utility>
3  #include <string>
4
5  std::string emphasize(const std::string& s)
6      {return s + "!";}
7
8  std::string emphasize(std::string&& s)
9      {return s + "!!!!";}
10
11 template <class A>
12 auto wrapper(A&& arg) {
13     std::cout << "Calling with argument " << arg << '\n';
14     auto result = emphasize(std::forward<A>(arg));
15     std::cout << "Return value " << result << '\n';
16     return result;
17 }
18
19 int main() {
20     std::string s("Bonjour");
21     wrapper(s);
22     wrapper(std::string("Hello"));
23 }
```

Perfect-Forwarding Use Case: Factory Functions

- A **factory function** is simply a function used to create objects.
- Often, perfect forwarding is used by factory functions in order to pass arguments through to a constructor, which performs the actual object creation.
- Example:

```
1  #include <iostream>
2  #include <string>
3  #include <complex>
4  #include <utility>
5  #include <memory>
6
7  // Make an object of type T.
8  template<typename T, typename Arg>
9  std::shared_ptr<T> factory(Arg&& arg) {
10     return std::shared_ptr<T>(
11         new T(std::forward<Arg>(arg)));
12 }
13
14 int main() {
15     using namespace std::literals;
16     auto s(factory<std::string>("Hello"s));
17     auto z(factory<std::complex<double>>(1.0i));
18     std::cout << *s << ' ' << *z << '\n';
19 }
```


Perfect-Forwarding Use Case: Emplace Operations

- Many container classes provide an operation that creates a new element directly inside the container, often referred to as an **emplace operation**.
- Some or all of the arguments to a member function performing an emplace operation correspond to arguments for a constructor invocation.
- Thus, an emplace operation typically employs perfect forwarding.
- The member function performing the emplace operation forwards some or all of its arguments to the constructor responsible for actually creating the new object.
- Some examples of emplace operations in the standard library include:
 - `std::list` **class**: `emplace`, `emplace_back`, `emplace_front`
 - `std::vector` **class**: `emplace`, `emplace_back`
 - `std::set` **class**: `emplace`, `emplace_hint`
 - `std::forward_list` **class**: `emplace_front`, `emplace_after`

Other Perfect-Forwarding Examples

- `std::thread` constructor uses forwarding to pass through arguments to thread function
- `std::packaged_task` function-call operator uses forwarding to pass through arguments to associated function
- `std::async` uses forwarding to pass through arguments to specified callable entity
- `std::make_unique` forwards arguments to `std::unique_ptr` constructor
- `std::make_shared` forwards arguments to `std::shared_ptr` constructor
- `std::make_pair` forwards arguments to `std::pair` constructor
- `std::make_tuple` forwards arguments to `std::tuple` constructor

Section 3.2.8

References

References I

- 1 S. Meyers, Universal References in C++11, C++ and Beyond, Asheville, NC, USA, Aug. 5–8, 2012.
This talk discusses rvalue/forwarding references.
- 2 S. Meyers. [Universal references in C++11](#). *Overload*, 111:8–12, Oct. 2012.
- 3 S. Meyers, Adventures in Perfect Forwarding, Facebook C++ Conference, Menlo Park, CA, USA, June 2, 2012.
This talk introduces perfect forwarding and discusses matters such as how to specialize forwarding templates and how to address interactions between forwarding and the pimpl idiom.
- 4 T. Becker, C++ Rvalue References Explained, 2013, http://thbecker.net/articles/rvalue_references/section_01.html
- 5 E. Bendersky, Understanding Lvalues and Rvalues in C and C++, 2011, <http://eli.thegreenplace.net/2011/12/15/understanding-lvalues-and-rvalues-in-c-and-c>

- 6 E. Bendersky, Perfect Forwarding and Universal References in C++, 2014, <http://eli.thegreenplace.net/2014/perfect-forwarding-and-universal-references-in-c/>
- 7 M. Kilpelainen. [Lvalues and rvalues](#). *Overload*, 61:12–13, June 2004.
- 8 H. E. Hinnant, Forward, ISO/IEC JTC1/SC22/WG21/N2951, Sept. 27, 2009, <http://www.open-std.org/jtc1/sc22/wg21/docs/papers/2009/n2951.html>
- 9 H. Hinnant and D. Krugler, Proposed Wording for US 90, ISO/IEC JTC1/SC22/WG21/N3143, Oct. 15, 2010, <http://www.open-std.org/jtc1/sc22/wg21/docs/papers/2010/n3143.html>

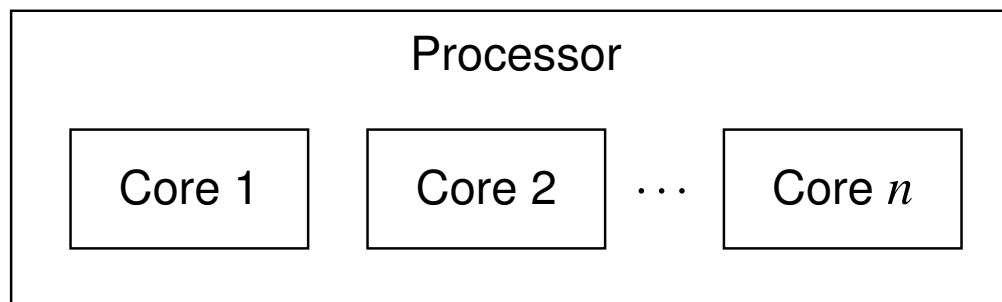
Section 3.3

Concurrency

Section 3.3.1

Preliminaries

Processors

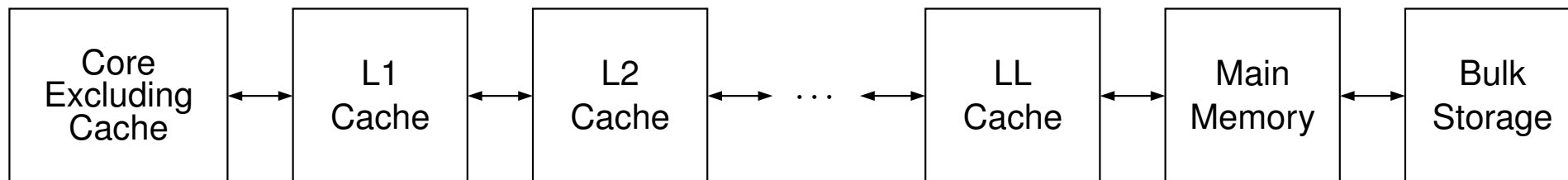


- A **core** is an independent processing unit that reads and executes program instructions, and consists of registers, an arithmetic logic unit (ALU), a control unit, and usually a cache.
- A **processor** is a computing element that consists of one or more cores, an external bus interface, and possibly a shared cache.
- A **thread** is a sequence of instructions (which can be executed by a core).
- At any given time, a core can execute one thread or, if the core supports simultaneous multithreading (such as hyperthreading), multiple threads.
- In the simultaneous multithreading case, the threads share the resources of the core.
- A processor with more than one core is said to be **multicore**.
- Most modern processors are multicore.
- Multicore processors can *simultaneously* execute *multiple* threads.

Processors (Continued)

- A multicore processor said to be **homogeneous** if all of its cores are identical.
- A multicore processor said to be **heterogeneous** if its has more than one type of core.
- Different types of cores might be used in order to:
 - provide different types of functionality (e.g., CPU and GPU)
 - provide different levels of performance (e.g., high-performance CPU and energy-efficient CPU)

Memory Hierarchy



- The component of a system that stores program instructions and data is called **main memory**.
- A **cache** is fast memory used to store copies of instructions and/or data from main memory.
- Main memory is *very slow* compared to the speed of a processor core.
- Due to the latency of main memory, caches are *essential* for good performance.
- Instruction and data caches may be *separate* or *unified* (i.e., combined).
- A cache may be *local* to single core or *shared* between two or more cores.
- The lowest-level (i.e., L1) cache is usually on the core and local to the core.
- The higher-level (i.e., L2, L3, . . . , LL [last level]) caches are usually shared between some or all of the cores.

Examples of Multicore Processors

- Intel Core i7-3820QM Processor (Q2 2012)
 - used in Lenovo W530 notebook
 - 64 bit, 2.7 GHz
 - 128/128 KB L1 cache, 1 MB L2 cache, 8 MB L3 cache
 - *4 cores*
 - *8 threads* (2 threads/core)
- Intel Core i7-5960X Processor Extreme Edition (Q3 2014)
 - targets desktops/notebooks
 - 64 bit, 3 GHz
 - 256/256 KB L1 cache, 2 MB L2 cache, 20 MB L3 cache
 - *8 cores*
 - *16 threads* (2 threads/core)
- Intel Xeon Processor E7-8890 v2 (Q1 2014)
 - targets servers
 - 64 bit, 2.8 GHz
 - 480/480 KB L1 cache, 3.5 MB L2 cache, 37.5 MB L3 cache
 - *15 cores*
 - *30 threads* (2 threads/core)

Examples of Multicore SoCs

- Qualcomm Snapdragon 805 SoC (Q1 2014)
 - used in *Google Nexus 6*
 - 32-bit 2.7 GHz *quad-core* Qualcomm Krait 450 (ARMv7-A)
 - 16/16 KB L1 cache (per core), 2 MB L2 cache (shared)
 - 600 MHz Qualcomm Adreno 420 GPU
- Samsung Exynos 5 Octa 5433 SoC
 - used in *Samsung Galaxy Note 4*
 - high-performance 1.9 GHz *quad-core* ARM Cortex-A57 paired with energy-efficient 1.3 GHz *quad-core* ARM Cortex-A53 (big.LITTLE); both 32-bit (64-bit capable but disabled) (ARMv8-A)
 - Cortex-A57: 48/32 KB L1 cache, 512 KB to 2 MB L2 cache?
 - 700 MHz Mali-T760MP6 GPU
- Apple A8 SoC (2014)
 - used in *Apple iPhone 6, Apple iPhone 6 Plus*
 - 64-bit 1.4 GHz *dual-core* CPU (ARMv8-A)
 - 64/64 KB L1 cache (per core), 1 MB L2 cache (shared), 4 MB L3 cache
 - PowerVR Series 6XT GX6450 (quad-core) GPU

Why Multicore Processors?

- in past, greater processing power obtained through *higher clock rates*
- *clock rates have stopped rising*, topping out at about 5 GHz (little change since about 2005)
- power consumption is linear in clock frequency and quadratic in voltage, but higher frequency typically requires higher voltage; so, considering effect of frequency and voltage together, power consumption grows approximately with *cube* of frequency
- greater power consumption translates into *increased heat production*
- higher clock rates would result in processors *overheating*
- transistor counts *still increasing* (Moore's law: since 1960s, transistor count has doubled approximately every 18 months)
- instead of increasing processing power by raising clock rate of processor core, simply *add more processor cores*
- n cores running at clock rate f use significantly less power and generate less heat than single core at clock rate nf
- going multicore allows for *greater processing power* with *lower power consumption* and *less heat production*

Section 3.3.2

Multithreaded Programming

Concurrency

- A **thread** is a sequence of instructions that can be independently managed by the operating-system scheduler.
- A **process** provides the resources that program needs to execute (e.g., address space, files, and devices) and at least one thread of execution.
- All threads of a process share the *same* address space.
- **Concurrency** is the situation where multiple threads execute over time periods (i.e., from start of execution to end) that *overlap* (but no threads are required to run simultaneously).
- **Parallelism** refers to the situation where multiple threads execute *simultaneously*.
- Concurrency can be achieved with:
 - 1 multiple single-threaded processes; or
 - 2 a single multithreaded process.
- A single multithreaded process is usually preferable, since data can be shared more easily between threads in a single process, due to the threads having a common address space.

Why Multithreading?

- Keep all of the processor cores busy (i.e., *fully utilize* all cores).
 - Most modern systems have multiple processor cores, due to having either multiple processors or a single processor that is multicore.
 - A single thread cannot fully utilize the computational resources available in such systems.
- Keep processes *responsive*.
 - In graphics applications, keep the GUI responsive while the application is performing slow operations such as I/O.
 - In network server applications, keep the server responsive to new connections while handling already established ones.
- *Simplify* the coding of cooperating tasks.
 - Some programs consist of several logically distinct tasks.
 - Instead of having the program manage when the computation associated with different tasks is performed, each task can be placed in a separate thread and the operating system can perform scheduling.
 - For certain types of applications, multithreading can significantly reduce the conceptual complexity of the program.

Section 3.3.3

Multithreaded Programming Models

Memory Model

- A **memory model** (also known as a **memory-consistency model**) is a formal specification of the effect of read and write operations on the memory system, which in effect describes how memory appears to programs.
- A memory model is essential in order for the semantics of a multithreaded program to be well defined.
- The memory model must address issues such as:
 - ordering
 - atomicity
- The memory model affects:
 - programmability (i.e., ease of programming)
 - performance
 - portability

Sequential Consistency (SC)

- The environment in which a multithreaded program is run is said to have **sequential consistency (SC)** if the result of any execution of the program is the same as if the operations of all threads are executed in *some sequential order*, and the operations of each individual thread appear in this sequence in *the order specified by the program*.
- In other words, in a sequentially-consistent execution of a multithreaded program, threads behave as if their operations were simply *interleaved*.
- Consider the multithreaded program (with two threads) shown below, where x , y , a , and b are all integer variables and initially zero.

Thread 1 Code

```
x = 1;  
a = y;
```

Thread 2 Code

```
y = 1;  
b = x;
```

- Some sequentially-consistent executions of this program include:
 - $x = 1; y = 1; b = x; a = y;$
 - $y = 1; x = 1; a = y; b = x;$
 - $x = 1; a = y; y = 1; b = x;$
 - $y = 1; b = x; x = 1; a = y;$

Sequential-Consistency (SC) Memory Model

- Since SC implies that memory must behave in a particular manner, SC implicitly defines a memory model, known as the **SC memory model**.
- In particular, SC implies that each write operation is *atomic* and becomes visible to all threads *simultaneously*.
- Thus, with the SC model, *all* threads see write operations on memory occur *atomically* in the *same* order, leading to all threads having a *consistent view* of memory.
- The SC model precludes (or makes extremely difficult) many hardware optimizations, such as:
 - store buffers
 - caches
 - out-of-order instruction execution
- The SC model also precludes many compiler optimizations, including:
 - reordering of loads and stores
- Although the SC model very is *intuitive*, it comes at a *very high cost* in terms of performance.

Load/Store Reordering Example: Single Thread

- Consider the program with the code below, where x and y are integer variables, all initially zero.

Original Thread 1 Code

```
x = 1;  
y = 1;  
// ...
```

- Suppose that, during optimization, the compiler transforms the preceding code to that shown below, effectively *reordering two stores*.

Optimized Thread 1 Code

```
y = 1;  
x = 1;  
// ...
```

- The execution of the optimized code is *indistinguishable* from a sequentially-consistent execution of the original code.
- The optimized program runs *as if* it were the original program.
- In a *single-threaded* program, loads and stores can be reordered without invalidating the SC model (if data dependencies are correctly considered).

Load/Store Reordering Example: Multiple Threads

- Consider the addition of a second thread to the program to yield the code below.

Original Thread 1 Code

```
x = 1;  
y = 1;  
// ...
```

Thread 2 Code

```
if (y == 1) {  
    assert(x == 1);  
}
```

- Suppose that the compiler makes the same optimization to the code for thread 1 as on the previous slide, yielding the code below.

Optimized Thread 1 Code

```
y = 1;  
x = 1;  
// ...
```

(Unchanged) Thread 2 Code

```
if (y == 1) {  
    assert(x == 1);  
}
```

- Thread 2 can observe x and y being modified in the wrong order (i.e., an order that is inconsistent with SC execution).
- The assertion in thread 2 can never fail in the original program, but can sometimes fail in the optimized program.
- In a *multithreaded* program, the reordering of loads and stores must be avoided *if SC is to be maintained*.

Store-Buffer Example: Without Store Buffer

- Consider the program below, where x , y , a , and b are integer variables, all initially zero.

Thread 1 Code

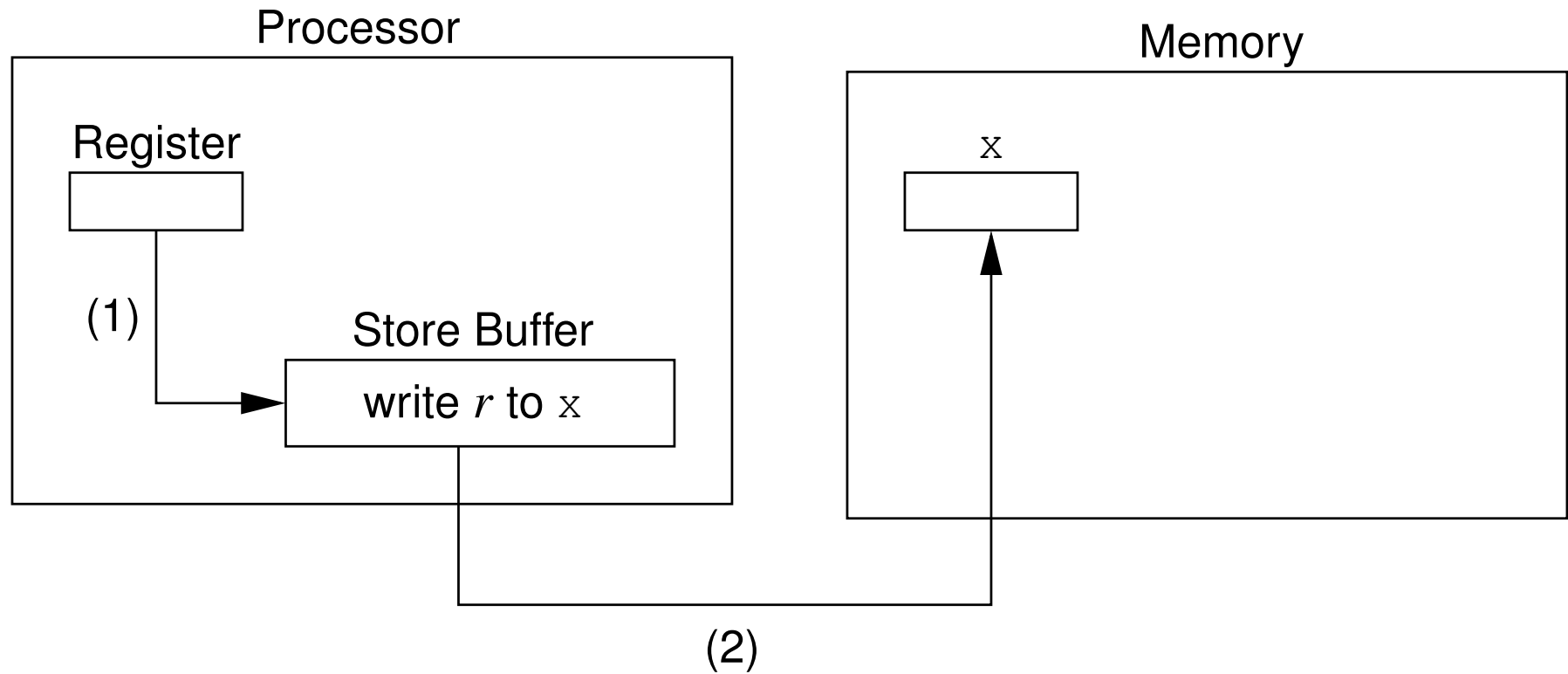
```
x = 1;  
a = y;
```

Thread 2 Code

```
y = 1;  
b = x;
```

- Some possible sequentially-consistent executions of the program include:
 - $x = 1; y = 1; b = x; a = y;$ (a is 1, b is 1)
 - $y = 1; x = 1; a = y; b = x;$ (a is 1, b is 1)
 - $x = 1; a = y; y = 1; b = x;$ (a is 0, b is 1)
 - $y = 1; b = x; x = 1; a = y;$ (a is 1, b is 0)
- In every sequentially-consistent execution of the program, one of “ $x = 1;$ ” or “ $y = 1;$ ” must execute first.
- If “ $x = 1;$ ” executes first, then b cannot be assigned 0.
- If “ $y = 1;$ ” executes first, then a cannot be assigned 0.
- No sequentially-consistent execution can result in a and b *both* being 0.

Store-Buffer Example: Store Buffer



- (1) transfer data from register to store buffer
- (2) flush store buffer to memory

Store-Buffer Example: With Store Buffer (Not SC)

Core 1		Core 2		Memory	
Code	Store Buffer	Code	Store Buffer	x	y
<code>x = 1;</code>	write 1 to x pending			0	0
	no change	<code>y = 1;</code>	write 1 to y pending	0	0
<code>a = y;</code> <code>// a = 0;</code>	no change		no change	0	0
	no change	<code>b = x;</code> <code>// b = 0;</code>	no change	0	0
	write 1 to x completed		no change	1	0
			write 1 to y completed	1	1

- The execution of the program results in `a` and `b` **both** being 0, which **violates SC**.
- The program behaves as if the lines of code in each thread were **reordered** (i.e., reversed), yielding: `a = y;` `b = x;` `x = 1;` `y = 1;`.
- A store buffer (or cache) must be avoided, **if SC is to be maintained**.

Atomicity of Memory Operations

- A fundamental property of SC is that all memory operations are *atomic*.
- Atomic memory operations require *synchronization* between processor cores.
- This synchronization *greatly increases the time required to access memory*, as a result of the time needed by processor cores to communicate and coordinate access to memory.
- Therefore, requiring all memory operations to be atomic is not desirable.
- Allowing non-atomic memory operations, however, would be *inconsistent with a fundamental property of SC*.

Data Races

- If memory operations are *not all atomic*, the possibility exists for something known as a data race.
- Two memory operations are said to **conflict** if they access the *same* memory location and *at least one* of the operations is a write.
- Two conflicting memory operations form a **data race** if they are from different threads and can be executed *at the same time*.
- A program with data races usually has unpredictable behavior (e.g., due to torn reads, torn writes, or worse).
- Example (data race):
 - Consider the multithreaded program listed below, where x , y , and z are (nonatomic) integer variables shared between threads and are initially zero.

```
Thread 1 Code
x = 1;
a = y + z;
```

```
Thread 2 Code
y = 1;
b = x + z;
```

- The program has data races on both x and y .
- Since z is not modified by any thread, z cannot participate in a data race.

Torn Reads

- A **torn read** is a read operation that (due to lack of atomicity) has only partially read its value when another (concurrent) write operation on the same location is performed.
- Consider a two-byte unsigned (big-endian) integer variable x , which is initially 1234 (hexadecimal).
- Suppose that the following (nonatomic) memory operations overlap in time:
 - thread 1 reads x ; and
 - thread 2 writes 5678 (hexadecimal) to x .
- Initially, x is 1234:

Byte 0	Byte 1
12	34
- Thread 1 reads 12 from the first byte of x .
- Thread 2 writes 56 and 78 to the first and second bytes of x , respectively, yielding:

Byte 0	Byte 1
56	78
- Thread 1 reads the second byte of x to obtain the value 78.
- The value read by thread 1 (i.e., 1278) is neither the value of x prior to the write by thread 2 (i.e., 1234) nor the value of x after the write by thread 2 (i.e., 5678).

Torn Writes

- A **torn write** is a write operation that (due to lack of atomicity) has only partially written its value when another (concurrent) read or write operation on the same location is performed.
- Consider a two-byte unsigned (big-endian) integer variable x , which is initially 0.
- Suppose that the following (nonatomic) memory operations overlap in time:
 - thread 1 writes 1234 (hexadecimal) to x ; and
 - thread 2 writes 5678 (hexadecimal) to x .

● Initially, x is 0:

Byte 0	Byte 1
00	00

● Thread 1 writes 12 to the first byte of x , yielding:

Byte 0	Byte 1
12	00

● Thread 2 writes 56 and 78 to the first and second bytes of x , respectively, yielding:

Byte 0	Byte 1
56	78

● Thread 1 writes 34 to the second byte of x , yielding:

Byte 0	Byte 1
56	34

● The resulting value in x (i.e., 5634) is neither the value written by thread 1 (i.e., 1234) nor the value written by thread 2 (i.e., 5678).

SC Data-Race Free (SC-DRF) Memory Model

- From a programmability standpoint, SC is extremely desirable, as it allows one to reason easily about the behavior of a multithreaded program.
- Unfortunately, as we saw earlier, SC precludes almost all useful compiler optimizations and hardware optimizations.
- As it turns out, if we drop the requirement that all memory operations be atomic and then restrict programs to be data-race free, SC can be provided while still allowing most compiler and hardware optimizations.
- This observation is the motivation behind the so called SC-DRF memory model.
- The **sequential-consistency for data-race free programs (SC-DRF) model** provides SC *only for programs that are data-race free*.
- The data-race free constraint is not overly burdensome, since data races will likely result in bugs anyhow.
- Several programming languages have used SC-DRF as the basis for their memory model, including C++, C, and Java.

C++ Memory Model

- The C++ programming language employs, at its default memory model, the *SC-DRF* model.
- Again, with the SC-DRF model, a program behaves as if its execution is sequentially consistent, provided that the program is data-race free.
- Support is also provided for other (more relaxed) memory models.
- For certain memory accesses, it is possible to override the default (i.e., SC-DRF) memory model, if desired.
- The execution of a program that is not data-race free results in *undefined behavior*.

Section 3.3.4

Thread Management

The `std::thread` Class

- `std::thread` class provides means to create new thread of execution, wait for thread to complete, and perform other operations to manage and query state of thread
- `thread` object may or may not be associated with thread (of execution)
- `thread` object that is associated with thread said to be **joinable**
- default constructor creates `thread` object that is **unjoinable**
- can also construct `thread` object by providing callable entity (e.g., function or functor) and arguments (if any), resulting in new thread invoking callable entity
- `thread` function provided with **copies** of arguments so must use reference wrapper class like `std::reference_wrapper` for reference semantics
- `thread` class is movable but **not copyable**
- each `thread` object has ID
- IDs of **joinable** `thread` objects are **unique**
- all **unjoinable** `thread` objects have **same** ID, distinct from ID of every **joinable** `thread` object

The `std::thread` Class (Continued)

- **join operation** waits for `thread` object's thread to complete execution and results in object becoming `unjoinable`
- **detach operation** dissociates thread from `thread` object (allowing thread to continue to execute independently) and results in object becoming `unjoinable`
- using `thread` object as source for move operation results in object becoming `unjoinable`
- if `thread` object joinable when destructor called, exception is thrown
- `hardware_concurrency` member function returns number of hardware threads that can run simultaneously (or zero if not well defined)
- thread creation and join operations establish synchronizes-with relationship (to be discussed later)

Member Types

Member Name	Description
<code>id</code>	thread ID type
<code>native_handle_type</code>	system-dependent handle type for underlying thread entity

Construction, Destruction, and Assignment

Member Name	Description
<code>constructor</code>	construct thread (overloaded)
<code>destructor</code>	destroy thread
<code>operator=</code>	move assign thread

Member Functions

Member Name	Description
<code>joinable</code>	check if thread joinable
<code>get_id</code>	get ID of thread
<code>native_handle</code>	get native handle for thread
<code>hardware_concurrency</code> (static)	get number of concurrent threads supported by hardware
<code>join</code>	wait for thread to finish executing
<code>detach</code>	permit thread to execute independently
<code>swap</code>	swap threads

Example: Hello World With Threads

```
1  #include <iostream>
2  #include <thread>
3
4  void hello ()
5  {
6      std::cout << "Hello World!\n";
7  }
8
9  int main ()
10 {
11     std::thread t(hello);
12     t.join();
13 }
```

```
1  #include <iostream>
2  #include <thread>
3
4  int main ()
5  {
6      std::thread t([]() {
7          std::cout << "Hello World!\n";
8      });
9      t.join();
10 }
```

Example: Thread-Function Argument Passing (Copy Semantics)

```
1  #include <iostream>
2  #include <thread>
3
4  void doWork(int i, int j)
5  {
6      std::cout << i << ' ' << j << '\n';
7  }
8
9  int main()
10 {
11     int i = 42;
12     std::thread t1(doWork, i, 1);
13     t1.join();
14 }
```

Example: Thread-Function Argument Passing (Reference Semantics)

```
1  #include <iostream>
2  #include <vector>
3  #include <functional>
4  #include <thread>
5
6  void doWork(const std::vector<int>& v)
7  {
8      for (auto i : v) {
9          std::cout << i << '\n';
10     }
11 }
12
13 int main()
14 {
15     std::vector<int> v{1, 2, 3, 4};
16
17     // copy semantics
18     std::thread t1(doWork, v);
19     t1.join();
20
21     // reference semantics
22     std::thread t2(doWork, std::ref(v));
23     t2.join();
24 }
```

Example: Thread-Function Argument Passing (Move Semantics)

```
1  #include <iostream>
2  #include <vector>
3  #include <utility>
4  #include <thread>
5
6  void doWork (std::vector<int>&& v)
7  {
8      for (auto i : v) {
9          std::cout << i << '\n';
10     }
11 }
12
13 int main ()
14 {
15     std::vector<int> v{1, 2, 3, 4};
16
17     // move semantics
18     std::thread t1 (doWork, std::move (v));
19     t1.join ();
20 }
```


Example: Moving Threads

```
1  #include <thread>
2  #include <iostream>
3  #include <utility>
4
5  // Return a thread that prints a greeting message.
6  std::thread makeThread() {
7      return std::thread([]() {
8          std::cout << "Hello World!\n";
9      });
10 }
11
12 // Return the same thread that was passed as an argument.
13 std::thread identity(std::thread t) {
14     return t;
15 }
16
17 int main() {
18     std::thread t1(makeThread());
19     std::thread t2(std::move(t1));
20     t1 = std::move(t2);
21     t1 = identity(std::move(t1));
22     t1.join();
23 }
```

Example: Lifetime Bug

```
1  #include <iostream>
2  #include <vector>
3  #include <algorithm>
4  #include <chrono>
5  #include <thread>
6
7  void threadFunc(const std::vector<int>* v) {
8      std::cout << std::accumulate(v->begin(), v->end(), 0)
9          << '\n';
10 }
11
12 void startThread() {
13     std::vector<int> v(1000000, 1);
14     std::thread t(threadFunc, &v);
15     t.detach();
16     // v is destroyed here but detached thread
17     // may still be using v
18 }
19
20 int main() {
21     startThread();
22     // Give the thread started by startThread
23     // sufficient time to complete its work.
24     std::this_thread::sleep_for(std::chrono::seconds(5));
25 }
```

The `std::this_thread` Namespace

Name	Description
<code>get_id</code>	get ID of current thread
<code>yield</code>	suggest rescheduling current thread so as to allow other threads to run
<code>sleep_for</code>	blocks execution of current thread for at least specified duration
<code>sleep_until</code>	blocks execution of current thread until specified time reached

Example: Identifying Threads

```
1  #include <thread>
2  #include <iostream>
3
4  // main thread ID
5  std::thread::id mainThread;
6
7  void func() {
8      if (std::this_thread::get_id() == mainThread) {
9          std::cout << "called by main thread\n";
10     } else {
11         std::cout << "called by secondary thread\n";
12     }
13 }
14
15 int main() {
16     mainThread = std::this_thread::get_id();
17     std::thread t([](){
18         // call func from secondary thread
19         func();
20     });
21     // call func from main thread
22     func();
23     t.join();
24 }
```

Thread Local Storage

- **thread storage duration**: object allocated when thread begins and deallocated when thread ends
- each thread has its own instance of object
- only objects declared **thread_local** have this storage duration
- **thread_local** implies **static** for variable of block scope
- **thread_local** can appear together with **static** or **extern** to adjust linkage
- example:

```
thread_local int counter = 0;  
static thread_local int x = 0;  
thread_local int y;  
  
void func() {  
    thread_local int counter = 0;  
    // equivalent to:  
    // static thread_local int counter = 0;  
}
```

Example: Thread Local Storage

```
1  #include <iostream>
2  #include <vector>
3  #include <thread>
4
5  thread_local int counter = 0;
6
7  void doWork(int id) {
8      static const char letters[] = "abcd";
9      for (int i = 0; i < 10; ++i) {
10         std::cout << letters[id] << counter << '\n';
11         ++counter;
12     }
13 }
14
15 int main() {
16     std::vector<std::thread> workers;
17     for (int i = 1; i <= 3; ++i) {
18         // invoke doWork in new thread
19         workers.emplace_back(doWork, i);
20     }
21     // invoke doWork in main thread
22     doWork(0);
23     for (auto& t : workers) {t.join();}
24 }
```

Section 3.3.5

Sharing Data Between Threads

Shared Data

- In multithreaded programs, it is often necessary to *share resources* between threads.
- Shared resources might include such things as variables, memory, files, devices, and so on.
- The sharing of resources, however, can lead to various problems when multiple threads want access to the same resource simultaneously.
- The most commonly shared resource is *variables*.
- When variables are shared between threads, the possibility exists that one thread may attempt to access a variable while another thread is modifying the same variable.
- Such *conflicting accesses* to variables can lead to data corruption and other problems.
- More generally, when any resource is shared, the potential for problems exists.
- Therefore, mechanisms are needed for ensuring that shared resources can be accessed safely.

Race Conditions

- A **race condition** is a behavior where the outcome depends on the relative ordering of the execution of operations on two or more threads.
- Sometimes, a race condition may be benign (i.e., does not cause any problem).
- Usually, the term “race condition” used to refer to a race condition that is not benign (i.e., breaks invariants or results in undefined behavior).
- A data race is a particularly evil type of race condition.
- A **deadlock** is a situation in which two or more threads are unable to make progress due to being *blocked* waiting for resources held by each other.
- A **livelock** is a situation in which two or more threads are *not blocked* but are unable to make progress due to needing resources held by each other.
- Often, race conditions can lead to deadlocks, livelocks, crashes, and other unpredictable behavior.

Critical Sections

- A **critical section** is a piece of code that accesses a shared resource (e.g., data structure) that must not be simultaneously accessed by more than one thread.
- A synchronization mechanism is needed at the entry to and exit from a critical section.
- The mechanism needs to provide *mutual exclusion* (i.e., prevent critical sections in multiple threads from executing simultaneously).
- Example (FIFO queue):
 - One thread is adding an element to a queue while another thread is removing an element from the same queue.
 - Since both threads modify the queue at the same time, they could corrupt the queue data structure.
 - Synchronization must be employed so that the execution of the parts of the code that add and remove elements are executed in a *mutually exclusive* manner (i.e., cannot run at the same time).

Data-Race Example

Shared (Global) Data

```
double balance = 100.00; // bank account balance
double credit = 50.00; // amount to deposit
double debit = 10.00; // amount to withdraw
```

Thread 1 Code

```
// double tmp = balance;
// tmp = tmp + credit;
// balance = tmp;
balance += credit;
```

Thread 2 Code

```
// double tmp = balance;
// tmp = tmp - debit;
// balance = tmp;
balance -= debit;
```

- above code has data race on `balance` object (i.e., more than one thread may access `balance` at same time with at least one thread writing)

Example: Data Race (Counter)

```
1  #include <iostream>
2  #include <thread>
3
4  unsigned long long counter = 0;
5
6  void func() {
7      for (int i = 0; i < 10000000; ++i) {
8          ++counter;
9      }
10 }
11
12 int main() {
13     std::thread t1(func);
14     std::thread t2(func);
15     t1.join();
16     t2.join();
17     std::cout << counter << '\n';
18 }
```

Example: Data Race and/or Race Condition (IntSet)

```
1  #include <thread>
2  #include <iostream>
3  #include <set>
4
5  class IntSet {
6  public:
7      bool contains(int i) const
8          {return s_.find(i) != s_.end();}
9      void add(int i)
10         {s_.insert(i);}
11 private:
12     std::set<int> s_;
13 };
14
15 IntSet s;
16
17 int main() {
18     std::thread t1([](){
19         for (int i = 0; i < 1000; ++i) s.add(2 * i);
20     });
21     std::thread t2([](){
22         for (int i = 0; i < 1000; ++i) s.add(2 * i + 1);
23     });
24     t1.join(); t2.join();
25     std::cout << s.contains(1000) << '\n';
26 }
```

Section 3.3.6

Mutexes

Mutexes

- A **mutex** is a locking mechanism used to synchronize access to a shared resource by providing *mutual exclusion*.
- A mutex has two basic operations:
 - **acquire**: lock (i.e., hold) the mutex
 - **release**: unlock (i.e., relinquish) the mutex
- A mutex can be held *by only one thread* at any given time.
- If a thread attempts to acquire a mutex that is already held by another thread, the operation will either block until the mutex can be acquired or fail with an error.
- A thread holding a (nonrecursive) mutex *cannot relock* the mutex.
- A thread acquires the mutex before accessing the shared resource and releases the mutex when finished accessing the resource.
- Since only one thread can hold a mutex at any given time and the shared resource is only accessed by the thread holding the mutex, mutually-exclusive access is guaranteed.

The `std::mutex` Class

- `std::mutex` class provides mutex functionality
- *not movable* and *not copyable*
- `lock` member function acquires mutex (blocking as necessary)
- `unlock` member function releases mutex
- thread that owns mutex should not attempt to lock mutex again
- all prior `unlock` operations on given mutex *synchronize with* `lock` operation (on *same* mutex) (synchronizes-with relationship to be discussed later)

Member Types

Name	Description
<code>native_handle_type</code>	system-dependent handle type for underlying mutex entity

Construction, Destruction, and Assignment

Name	Description
constructor	construct mutex
destructor	destroy mutex

Other Member Functions

Name	Description
<code>lock</code>	acquire mutex, blocking if not available
<code>try_lock</code>	try to lock mutex without blocking
<code>unlock</code>	release mutex
<code>native_handle</code>	get handle for underlying mutex entity

Example: Avoiding Data Race Using Mutex (Counter) (mutex)

```
1  #include <iostream>
2  #include <thread>
3  #include <mutex>
4
5  std::mutex m;
6  unsigned long long counter = 0;
7
8  void func() {
9      for (int i = 0; i < 10000000; ++i) {
10         m.lock(); // acquire mutex
11         ++counter;
12         m.unlock(); // release mutex
13     }
14 }
15
16 int main() {
17     std::thread t1(func);
18     std::thread t2(func);
19     t1.join();
20     t2.join();
21     std::cout << counter << '\n';
22 }
```

The `std::lock_guard` Template Class

- `std::lock_guard` is RAII class for mutexes
- declaration:

```
template <class T> class lock_guard;
```
- template parameter `T` specifies type of mutex (e.g., `std::mutex`, `std::recursive_mutex`)
- avoids problem of inadvertently forgetting to release mutex (e.g., due to exception or forgetting `unlock` call)
- constructor takes mutex as argument
- *not movable* and *not copyable*
- acquires mutex in constructor
- releases mutex in destructor
- since language ensures that all objects destroyed at end of lifetime, release of mutex guaranteed (even if some code skipped due to thrown exception)
- advisable to use `lock_guard` instead of calling `lock` and `unlock` explicitly

Member Types

Name	Description
<code>mutex_type</code>	underlying mutex type

Construction, Destruction, and Assignment

Name	Description
constructor	construct lock guard
destructor	destroy lock guard

Example: Avoiding Data Race Using Mutex (Counter) (lock_guard)

```
1  #include <iostream>
2  #include <thread>
3  #include <mutex>
4
5  std::mutex m;
6  unsigned long long counter = 0;
7
8  void func() {
9      for (int i = 0; i < 10000000; ++i) {
10         // lock_guard constructor acquires mutex
11         std::lock_guard<std::mutex> lock(m);
12         ++counter;
13         // lock_guard destructor releases mutex
14     }
15 }
16
17 int main() {
18     std::thread t1(func);
19     std::thread t2(func);
20     t1.join();
21     t2.join();
22     std::cout << counter << '\n';
23 }
```

Example: Avoiding Data Race Using Mutex (IntSet) (lock_guard)

```
1  #include <thread>
2  #include <iostream>
3  #include <set>
4  #include <mutex>
5
6  class IntSet {
7  public:
8      bool contains(int i) const {
9          std::lock_guard<std::mutex> lg(m_);
10         return s_.find(i) != s_.end();
11     }
12     void add(int i) {
13         std::lock_guard<std::mutex> lg(m_);
14         s_.insert(i);
15     }
16 private:
17     std::set<int> s_;
18     mutable std::mutex m_;
19 };
20
21 IntSet s;
22
23 int main() {
24     std::thread t1([](){
25         for (int i = 0; i < 1000; ++i) s.add(2 * i);
26     });
27     std::thread t2([](){
28         for (int i = 0; i < 1000; ++i) s.add(2 * i + 1);
29     });
30     t1.join(); t2.join();
31     std::cout << s.contains(1000) << '\n';
32 }
```

The `std::unique_lock` Template Class

- `std::unique_lock` is another RAII class for mutexes
- declaration:

```
template <class T> class unique_lock;
```
- template parameter `T` specifies type of mutex (e.g., `std::mutex`, `std::recursive_mutex`)
- unlike case of `std::lock_guard`, in case of `unique_lock` do not have to hold mutex over entire lifetime of RAII object
- have choice of whether to acquire mutex upon construction
- also can acquire and release mutex many times throughout lifetime of `unique_lock` object
- upon destruction, if mutex is held, it is released
- since mutex is always guaranteed to be released by destructor, cannot forget to release mutex
- `unique_lock` is used in situations where want to be able to transfer ownership of lock (e.g., return from function) or RAII object needed for mutex but do not want to hold mutex over entire lifetime of RAII object
- *movable* but *not copyable*

std::unique_lock Members

Member Types

Name	Description
<code>mutex_type</code>	underlying mutex type

Construction, Destruction, and Assignment

Name	Description
constructor	construct unique lock
destructor	destroy unique lock
operator=	move assign

Locking Functions

Name	Description
<code>lock</code>	acquire mutex, blocking if not available
<code>try_lock</code>	try to lock mutex without blocking
<code>try_lock_for</code>	try to lock mutex without blocking
<code>try_lock_until</code>	try to lock mutex without blocking
<code>unlock</code>	release mutex

Observer Functions

Name	Description
<code>owns_lock</code>	tests if lock owns associated mutex
<code>operator bool</code>	tests if lock owns associated mutex

Example: Avoiding Data Race Using Mutex (Counter) (unique_lock)

```
1  #include <iostream>
2  #include <thread>
3  #include <mutex>
4
5  std::mutex m;
6  unsigned long long counter = 0;
7
8  void func() {
9      for (int i = 0; i < 1000000; ++i) {
10         // Create a lock object without locking the mutex.
11         std::unique_lock<std::mutex> lock(m, std::defer_lock);
12         // ...
13         // Lock the mutex.
14         lock.lock();
15         ++counter;
16         // The unique_lock destructor releases the mutex.
17     }
18 }
19
20 int main() {
21     std::thread t1(func);
22     std::thread t2(func);
23     t1.join();
24     t2.join();
25     std::cout << counter << '\n';
26 }
```

The `std::lock` Template Function

- `std::lock` variadic template function that can acquire multiple locks simultaneously without risk of deadlock (assuming the only locks involved are ones passed to `lock`)

- declaration:

```
template <class T1, class T2, class... TN>  
void lock (T1&, T2&, TN& ...);
```

- takes as arguments one or more locks to be acquired

Example: Acquiring Two Locks for Swap (Incorrect)

```
1  #include <thread>
2  #include <vector>
3  #include <mutex>
4
5  class BigBuf // A Big Buffer
6  {
7  public:
8      static constexpr int size() {return 16 * 1024 * 1024;}
9      BigBuf() : data_(size()) {}
10     BigBuf& operator=(const BigBuf&) = delete;
11     BigBuf& operator=(BigBuf&&) = delete;
12     void swap(BigBuf& other) {
13         if (this == &other)
14             return;
15         // acquiring the two mutexes in this way can result in deadlock
16         std::lock_guard<std::mutex> lock1(m_);
17         std::lock_guard<std::mutex> lock2(other.m_);
18         std::swap(data_, other.data_);
19     }
20     // ...
21 private:
22     std::vector<char> data_;
23     mutable std::mutex m_;
24 };
25
26 BigBuf a;
27 BigBuf b;
28
29 int main()
30 {
31     std::thread t1([](){
32         for (int i = 0; i < 100000; ++i) a.swap(b);
33     });
34     std::thread t2([](){
35         for (int i = 0; i < 100000; ++i) b.swap(a);
36     });
37     t1.join(); t2.join();
38 }
```

Example: Acquiring Two Locks for Swap

```
1  #include <thread>
2  #include <vector>
3  #include <mutex>
4
5  class BigBuf // A Big Buffer
6  {
7  public:
8      static constexpr int size() {return 16 * 1024 * 1024;}
9      BigBuf() : data_(size()) {}
10     BigBuf& operator=(const BigBuf&) = delete;
11     BigBuf& operator=(BigBuf&&) = delete;
12     void swap(BigBuf& other) {
13         if (this == &other)
14             return;
15         std::unique_lock<std::mutex> lock1(m_, std::defer_lock);
16         std::unique_lock<std::mutex> lock2(other.m_, std::defer_lock);
17         std::lock(lock1, lock2);
18         std::swap(data_, other.data_);
19     }
20     // ...
21 private:
22     std::vector<char> data_;
23     mutable std::mutex m_;
24 };
25
26 BigBuf a;
27 BigBuf b;
28
29 int main()
30 {
31     std::thread t1([](){
32         for (int i = 0; i < 100000; ++i) a.swap(b);
33     });
34     std::thread t2([](){
35         for (int i = 0; i < 100000; ++i) b.swap(a);
36     });
37     t1.join(); t2.join();
38 }
```

The `std::timed_mutex` Class

- `std::timed_mutex` class provides mutex that allows timeout to be specified when acquiring mutex
- if mutex cannot be acquired in time specified, acquire operation fails (i.e., does not lock mutex) and error returned
- adds `try_lock_for` and `try_lock_until` member functions to try to lock mutex with timeout

Example: Acquiring Mutex With Timeout (std::timed_mutex)

```
1  #include <vector>
2  #include <iostream>
3  #include <thread>
4  #include <mutex>
5  #include <chrono>
6
7  std::timed_mutex m;
8
9  void doWork() {
10     for (int i = 0; i < 10000; ++i) {
11         std::unique_lock<std::timed_mutex> lock(m,
12             std::defer_lock);
13         int count = 0;
14         while (!lock.try_lock_for(
15             std::chrono::microseconds(1))) {++count;}
16         std::cout << count << '\n';
17     }
18 }
19
20 int main() {
21     std::vector<std::thread> workers;
22     for (int i = 0; i < 16; ++i) {
23         workers.emplace_back(doWork);
24     }
25     for (auto& t : workers) {t.join();}
26 }
```

Recursive Mutexes

- A **recursive mutex** is a mutex for which a thread may own *multiple* locks *at the same time*.
- After a mutex is first locked by thread *A*, thread *A* can acquire additional locks on the mutex (without releasing the lock already held).
- The mutex is not available to other threads until thread *A* releases all of its locks on the mutex.
- A recursive mutex is typically used when code that locks a mutex must call other code that locks the same mutex (in order to avoid deadlock).
- For example, a function that acquires a mutex and recursively calls itself (resulting in the mutex being relocked) would need to employ a recursive mutex.
- A recursive mutex has *more overhead* than a nonrecursive mutex.
- Code that uses recursive mutexes can often be *more difficult to understand* and therefore *more prone to bugs*.
- Consequently, the use of recursive mutexes should be *avoided if possible*.

Recursive Mutex Classes

- recursive mutexes provided by classes `std::recursive_mutex` and `std::recursive_timed_mutex`
- `recursive_mutex` class similar to `std::mutex` class except allows relocking
- `recursive_timed_mutex` class similar to `std::timed_mutex` class except allows relocking
- implementation-defined limit to number of levels of locking allowed by recursive mutex

Shared Mutexes

- A **shared mutex** (also known as a **multiple-reader/single-writer mutex**) is a mutex that allows both *shared and exclusive* access.
- A shared mutex has *two types of locks*: shared and exclusive.
- **Exclusive lock**:
 - *Only one* thread can hold an *exclusive* lock on a mutex.
 - While a thread holds an exclusive lock on a mutex, no other thread can hold any type of lock on the mutex.
- **Shared lock**:
 - *Any number* of threads (within implementation limits) can take a *shared* lock on a mutex.
 - While any thread holds a shared lock on a mutex, no thread may take an exclusive lock on the mutex.
- A shared mutex would typically be used to protect shared data that is seldom updated but cannot be safely updated if any thread is reading it.
- A thread takes a shared lock for reading, thus allowing *multiple readers*.
- A thread takes an exclusive lock for writing, thus allowing *only one writer with no readers*.
- A shared mutex need not be fair in its granting of locks (e.g., readers could starve writers).

The `std::shared_timed_mutex` Class

- `std::shared_timed_mutex` class provides shared mutex
- `shared_timed_mutex` also allows timeout for acquiring mutex

std::shared_timed_mutex Members

Construction, Destruction, and Assignment

Name	Description
constructor	construct mutex
destructor	destroy mutex
operator= [deleted]	not movable or copyable

Exclusive Locking Functions

Name	Description
lock	acquire exclusive ownership of mutex, blocking if not available
try_lock	try to acquire exclusive ownership of mutex without blocking
try_lock_for	try to acquire exclusive ownership of mutex without blocking
try_lock_until	try to acquire exclusive ownership of mutex without blocking
unlock	release exclusive ownership of mutex

Shared Locking Functions

Name	Description
<code>lock_shared</code>	acquire shared ownership of mutex, blocking if not available
<code>try_lock_shared</code>	try to acquire shared ownership of mutex without blocking
<code>try_lock_shared_for</code>	try to acquire shared ownership of mutex without blocking
<code>try_lock_shared_until</code>	try to acquire shared ownership of mutex without blocking
<code>unlock_shared</code>	release shared ownership of mutex

The `std::shared_lock` Template Class

- `std::shared_lock` is RAII class for shared mutexes
- declaration:

```
template <class T> class shared_lock;
```
- template parameter `T` specifies type of mutex (e.g., `std::shared_timed_mutex`)
- similar interface as `std::unique_lock` but uses shared locking
- constructor may optionally acquire mutex
- may acquire and release mutex many times throughout lifetime of object
- destructor releases mutex if held
- all operations mapped onto shared locking primitives (e.g., `lock` mapped to `lock_shared`, `unlock` mapped to `unlock_shared`)
- for exclusive locking with shared mutexes, `std::unique_lock` and `std::lock_guard` can be used

Example: `std::shared_timed_mutex`

```
1  #include <thread>
2  #include <mutex>
3  #include <iostream>
4  #include <vector>
5  #include <shared_mutex>
6
7  std::mutex coutMutex;
8  int counter = 0;
9  std::shared_timed_mutex counterMutex;
10
11 void writer() {
12     for (int i = 0; i < 10; ++i) {
13         {
14             std::lock_guard<std::shared_timed_mutex> lock(counterMutex);
15             ++counter;
16         }
17         std::this_thread::sleep_for(std::chrono::milliseconds(100));
18     }
19 }
20
21 void reader() {
22     for (int i = 0; i < 100; ++i) {
23         int c;
24         {
25             std::shared_lock<std::shared_timed_mutex> lock(counterMutex);
26             c = counter;
27         }
28         {
29             std::lock_guard<std::mutex> lock(coutMutex);
30             std::cout << std::this_thread::get_id() << ' ' << c << '\n';
31         }
32         std::this_thread::sleep_for(std::chrono::milliseconds(10));
33     }
34 }
35
36 int main() {
37     std::vector<std::thread> threads;
38     threads.emplace_back(writer);
39     for (int i = 0; i < 16; ++i) threads.emplace_back(reader);
40     for (auto& t : threads) t.join();
41 }
```

std::once_flag and std::call_once

- sometimes may want to perform action only once in code executed in multiple threads
- can be achieved through use of `std::once_flag` type in conjunction with `std::call_once` template function
- `std::once_flag` class represents flag used to track if action performed
- declaration of `std::call_once`:

```
template <class Callable, class... Args>  
void call_once(std::once_flag& flag, Callable&& f,  
              Args&&... args);
```

- `std::call_once` invokes `f` only once based on value of `flag` object
- first invocation of `f` is guaranteed to complete before any threads return from `call_once`
- useful for one-time initialization of dynamically generated objects

Example: One-Time Action

```
1  #include <iostream>
2  #include <vector>
3  #include <thread>
4  #include <mutex>
5
6  std::once_flag flag;
7
8  void worker(int id) {
9      std::call_once(flag, [id]() {
10         // This code will be invoked only once.
11         std::cout << "first: " << id << '\n';
12     });
13 }
14
15 int main() {
16     std::vector<std::thread> threads;
17     for (int i = 0; i < 16; ++i) {
18         threads.emplace_back(worker, i);
19     }
20     for (auto& t : threads) {
21         t.join();
22     }
23 }
```

Example: One-Time Initialization

```
1  #include <vector>
2  #include <thread>
3  #include <mutex>
4  #include <cassert>
5  #include <memory>
6
7  std::unique_ptr<int> value;
8  std::once_flag initFlag;
9
10 void initValue() {value = std::make_unique<int>(42);}
11
12 const int& getValue() {
13     std::call_once(initFlag, initValue);
14     return *value.get();
15 }
16
17 void doWork() {
18     const int& v = getValue();
19     assert(v == 42);
20     // ...
21 }
22
23 int main() {
24     std::vector<std::thread> threads;
25     for (int i = 0; i < 4; ++i) {threads.emplace_back(doWork);}
26     for (auto& t : threads) {t.join();}
27 }
```

Static Local Variable Initialization and Thread Safety

- initialization of static local object is thread safe
- object is initialized first time control passes through its declaration
- object deemed initialized upon completion of initialization
- if control enters declaration concurrently while object being initialized, concurrent execution waits for completion of initialization
- code like following is thread safe:

```
const std::string& meaningOfLife() {  
    static const std::string x("42");  
    return x;  
}
```

Section 3.3.7

Condition Variables

Condition Variables

- In concurrent programs, the need often arises for a thread to *wait until a particular event occurs* (e.g., I/O has completed or data is available).
- Having a thread *repeatedly check* for the occurrence of an event can be *inefficient* (i.e., can waste processor resources).
- It is often better to have the thread block and then only resume execution after the event of interest has occurred.
- A **condition variable** is a synchronization primitive that allows threads to *wait (by blocking)* until a particular condition occurs.
- A condition variable corresponds to some event of interest.
- A thread that wants to wait for an event, performs a *wait operation* on the condition variable.
- A thread that wants to notify one or more waiting threads of an event performs a *signal operation* on the condition variable.
- When a signalled thread resumes, however, the signalled condition is not guaranteed to be true (and must be rechecked), since another thread may have caused condition to change.

The `std::condition_variable` Class

- `std::condition_variable` class provides condition variable
- *not movable* and *not copyable*
- `wait`, `wait_for`, and `wait_until` member functions used to wait for condition
- `notify_one` and `notify_all` used to signal waiting thread(s) of condition
- must re-check condition when awaking from wait since:
 - spurious awakenings are permitted
 - between time thread is signalled and time it awakens and locks mutex, another thread could cause condition to change
- concurrent invocation is allowed for `notify_one`, `notify_all`, `wait`, `wait_for`, `wait_until`
- each of `wait`, `wait_for`, and `wait_until` atomically releases mutex and blocks
- `notify_one` and `notify_all` are atomic

Member Types

Name	Description
<code>native_handle_type</code>	system-dependent handle type for underlying condition variable entity

Construction, Destruction, and Assignment

Name	Description
constructor	construct object
destructor	destroy object
operator= [deleted]	not movable or copyable

Notification and Waiting Member Functions

Name	Description
<code>notify_one</code>	notify one waiting thread
<code>notify_all</code>	notify all waiting threads
<code>wait</code>	blocks current thread until notified
<code>wait_for</code>	blocks current thread until notified or specified duration passed
<code>wait_until</code>	blocks current thread until notified or specified time point reached

Native Handle Member Functions

Name	Description
<code>native_handle</code>	get native handle associated with condition variable

Example: Condition Variable (IntStack)

```
1  #include <iostream>
2  #include <vector>
3  #include <thread>
4  #include <mutex>
5  #include <condition_variable>
6
7  class IntStack {
8  public:
9      IntStack() {};
10     IntStack(const IntStack&) = delete;
11     IntStack& operator=(const IntStack&) = delete;
12     int pop() {
13         std::unique_lock<std::mutex> lock(m_);
14         c_.wait(lock, [this]() {return !v_.empty();});
15         int x = v_.back();
16         v_.pop_back();
17         return x;
18     }
19     void push(int x) {
20         std::lock_guard<std::mutex> lock(m_);
21         v_.push_back(x);
22         c_.notify_one();
23     }
24 private:
25     std::vector<int> v_;
26     mutable std::mutex m_;
27     mutable std::condition_variable c_; // not empty
28 };
29
30 constexpr int numIters = 1000;
31 IntStack s;
32
33 int main() {
34     std::thread t1([](){
35         for (int i = 0; i < numIters; ++i) s.push(2 * i + 1);
36     });
37     std::thread t2([](){
38         for (int i = 0; i < numIters; ++i) std::cout << s.pop() << '\n';
39     });
40     t1.join(); t2.join();
41 }
```

The `std::condition_variable_any` Class

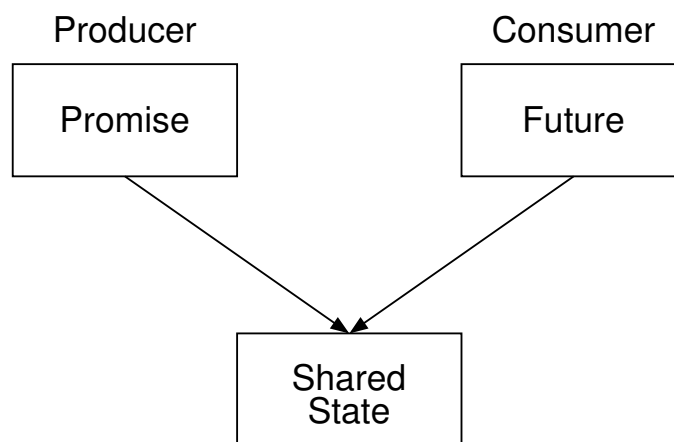
- with `std::condition_variable` class, `std::unique_lock<std::mutex>` class must be used for wait operation
- `std::condition_variable_any` class allows any mutex type (meeting certain basic requirements) to be used
- interface of `std::condition_variable_any` class similar to that of `std::condition_variable` class
- prefer `condition_variable` to `condition_variable_any` since former may be more efficient

Section 3.3.8

Promises and Futures

Promises and Futures

- promise and future together form *one-time* communication channel for passing result (i.e., value or exception) of computation from one thread to same or another thread
- **promise**: object associated with promised result (i.e., value or exception) to be produced
- **future**: object through which promised result later made available
- **shared state**: holds promised result for access through future object (shared by promise object and corresponding future object)
- producer of result uses promise object to store result in shared state
- consumer uses future object (corresponding to promise) to retrieve result from shared state



Promises and Futures (Continued)

- promises and futures useful in both single-threaded and multithreaded programs
- in single-threaded programs, might be used to propagate exception to another part of program
- in multithreaded program, often need arises to do some computation asynchronously and then later get result when ready
- requires synchronization between threads producing and consuming result
- thread consuming result must *wait until result is available*
- must *avoid data races* when accessing result shared between threads
- this type of synchronization can be accomplished via promise and future

The `std::promise` Template Class

- `std::promise` provides access to promise-future shared state for writing result
- declaration:

```
template <class T> class promise;
```
- T is type of result associated with promise (which can be `void`)
- movable but *not copyable*
- `set_value` member function sets result to particular value
- `set_exception` member function sets result to exception
- can set result *only once*
- `get_future` member function retrieves future associated with promise
- `get_future` may be called *only once*
- if `promise` object is destroyed before its associated result is set, `std::future_error` exception will be thrown if attempt made to retrieve result from corresponding `future` object

Construction, Destruction, and Assignment

Name	Description
constructor	construct object
destructor	destroy object
operator=	move assignment

std::promise Members (Continued)

Other Functions

Name	Description
<code>swap</code>	swap two promise objects
<code>get_future</code>	get future associated with promised result
<code>set_value</code>	set result to specified value
<code>set_value_at_thread_exit</code>	set result to specified value while delivering notification only at thread exit
<code>set_exception</code>	set result to specified exception
<code>set_exception_at_thread_exit</code>	set result to specified exception while delivering notification only at thread exit

The `std::future` Template Class

- `std::future` provides access to promise-future shared state for reading result
- declaration:

```
template <class T> class future;
```
- T is type of result associated with future (which can be `void`)
- movable but *not copyable*
- `get` member function retrieves result, blocking if result not yet available
- `get` may be called *only once*
- `wait` member function waits for result to become available without actually retrieving result

Construction, Destruction, and Assignment

Name	Description
constructor	construct object
destructor	destroy object
operator=	move assignment

Other Functions

Name	Description
share	transfer shared state to <code>shared_future</code> object
get	get result
valid	check if <code>future</code> object refers to shared state
wait	wait for result to become available
wait_for	wait for result to become available or time duration to expire
wait_until	wait for result to become available or time point to be reached

Example: Promises and Futures (Without `std::async`)

```
1  #include <future>
2  #include <thread>
3  #include <iostream>
4  #include <utility>
5
6  double computeValue() {
7      return 42.0;
8  }
9
10 void produce(std::promise<double> p) {
11     // write result to promise
12     p.set_value(computeValue());
13 }
14
15 int main() {
16     std::promise<double> p;
17     auto f = p.get_future(); // save future before move
18     std::thread producer(produce, std::move(p));
19     std::cout << f.get() << '\n';
20     producer.join();
21 }
```

The `std::shared_future` Template Class

- `std::shared_future` similar to `future` except object can be copied
- `shared_future` object can be obtained by using `share` member function of `future` class to transfer contents of `future` object into `shared_future` object
- `shared_future` is *copyable* (unlike `future`)
- allows multiple threads to wait for same result (associated with `shared_future` object)
- `get` member can be called multiple times

Example: `std::shared_future`

```
1  #include <iostream>
2  #include <vector>
3  #include <thread>
4  #include <future>
5
6  void consume(std::shared_future<int> f) {
7      std::cout << f.get() << '\n';
8  }
9
10 int main() {
11     std::promise<int> p;
12     std::shared_future<int> f = p.get_future().share();
13     std::vector<std::thread> consumers;
14     for (int i = 0; i < 16; ++i) {
15         consumers.emplace_back(consume, f);
16     }
17     p.set_value(42);
18     for (auto& i : consumers) {
19         i.join();
20     }
21 }
```

The `std::async` Template Function

- `std::async` template function used to launch callable entity (e.g., function or functor) asynchronously
- declaration (uses default launch policy):

```
template <class Func, class... Args>
future<typename result_of<typename decay<Func>::type (
    typename decay<Args>::type...)>::type>
    async(Func&& f, Args&&... args);
```

- declaration (with launch policy parameter):

```
template <class Func, class... Args>
future<typename result_of<typename decay<Func>::type (
    typename decay<Args>::type...)>::type>
    async(launch policy, Func&& f, Args&&... args);
```

- numerous launch policies supported via bitmask `std::launch`
- if `async` bit set, execute on new thread
- if `deferred` bit set, execute on calling thread when result needed
- if multiple bits set, implementation free to choose between them
- in asynchronous execution case, essentially creates promise to hold result and returns associated future; launches thread to execute function/functor and sets promise when function/functor returns

The `std::async` Template Function (Continued)

- `future` (i.e., `future` and `shared_future`) objects created by `async` function have slightly different behavior than `future` objects created in other ways
- in case of `future` object created by `async` function: if `future` object is *last* `future` object referencing its shared state, destructor for `future` object will *block* until result associated with `future` object becomes ready

Example: Promises and Futures (With `std::async`)

```
1  #include <future>
2  #include <iostream>
3
4  double computeValue() {
5      return 42.0;
6  }
7
8  int main() {
9      // invoke computeValue function asynchronously in
10     // separate thread
11     auto f = std::async(std::launch::async, computeValue);
12     std::cout << f.get() << '\n';
13 }
```


Example: Futures and Exceptions

```
1  #include <iostream>
2  #include <vector>
3  #include <cmath>
4  #include <future>
5  #include <stdexcept>
6
7  double squareRoot(double x) {
8      if (x < 0.0) {
9          throw std::domain_error(
10             "square root of negative number");
11     }
12     return std::sqrt(x);
13 }
14
15 int main() {
16     std::vector<double> values{1.0, 2.0, -1.0};
17     std::vector<std::future<double>> results;
18     for (auto x : values) {
19         results.push_back(std::async(squareRoot, x));
20     }
21     for (auto& x : results) {
22         try {
23             std::cout << x.get() << '\n';
24         } catch (const std::domain_error&) {
25             std::cout << "error\n";
26         }
27     }
28 }
```

The `std::packaged_task` Template Class

- `std::packaged_task` template class provides wrapper for callable entity (e.g., function or functor) that makes return value available via future

- declaration:

```
template <class R, class... Args>
    class packaged_task<R(Args...) >;
```

- template parameters `R` and `Args` specify return type and arguments for callable entity
- similar to `std::function` except return value of wrapped function made available via future
- packaged task often used as thread function
- movable but *not copyable*
- `get_future` member retrieves future associated with packaged task
- `get_future` can be called *only once*

std::packaged_task Members

Construction, Destruction, and Assignment

Name	Description
constructor	construct object
destructor	destroy object
operator=	move assignment

Other Functions

Name	Description
valid	check if task object currently associated with shared state
swap	swap two task objects
get_future	get future associated with promised result
operator()	invoke function
make_ready_at_thread_exit	invoke function ensuring result ready only once current thread exits
reset	reset shared state, abandoning any previously stored result

Example: Packaged Task

```
1  #include <iostream>
2  #include <thread>
3  #include <future>
4  #include <utility>
5  #include <chrono>
6
7  int getMeaningOfLife () {
8      // Let the suspense build before providing the answer.
9      std::this_thread::sleep_for (std::chrono::milliseconds (
10         1000));
11     // Return the answer.
12     return 42;
13 }
14
15 int main () {
16     std::packaged_task<int ()> pt (getMeaningOfLife);
17     // Save the future.
18     auto f = pt.get_future ();
19     // Start a thread running the task and detach the thread.
20     std::thread t (std::move (pt));
21     t.detach ();
22     // Get the result via the future.
23     int result = f.get ();
24     std::cout << "The meaning of life is " << result << '\n';
25 }
```

Example: Packaged Task With Arguments

```
1  #include <iostream>
2  #include <cmath>
3  #include <thread>
4  #include <future>
5
6  double power(double x, double y) {
7      return std::pow(x, y);
8  }
9
10 int main() {
11     // invoke task in main thread
12     std::packaged_task<double(double, double)> task(power);
13     task(0.5, 2.0);
14     std::cout << task.get_future().get() << '\n';
15     // reset shared state
16     task.reset();
17     // invoke task in new thread
18     auto f = task.get_future();
19     std::thread t(std::move(task), 2.0, 0.5);
20     t.detach();
21     std::cout << f.get() << '\n';
22 }
```

Section 3.3.9

Atomics

- To avoid data races when sharing data between threads, it is often necessary to employ *synchronization* (e.g., by using mutexes).
- Atomic types are another mechanism for providing synchronized access to data.
- An operation that is indivisible is said to be **atomic** (i.e., no parts of any other operations can interleave with any part of an atomic operation).
- Most processors support atomic memory operations via special machine instructions.
- Atomic memory operations cannot result in torn reads or torn writes.
- The standard library offers the following types in order to provide support for atomic memory operations:
 - `std::atomic_flag`
 - `std::atomic`
- These types provide a uniform interface for accessing the atomic memory operations of the underlying hardware.

Atomics (Continued)

- An atomic type provides guarantees regarding:
 - ① atomicity; and
 - ② ordering.
- An ordering guarantee specifies the manner in which memory operations can become visible to threads.
- Several memory ordering schemes are supported by atomic types.
- The default memory order is sequentially consistent (`std::memory_order_seq_cst`).
- Initially, only this default will be considered.

The `std::atomic_flag` Class

- `std::atomic_flag` provides flag with basic atomic operations
- flag can be in one of two states: set (i.e., true) or clear (i.e., false)
- two operations for flag:
 - **test and set**: set state to true and query previous state
 - **clear**: set state to false
- default constructor initializes flag to *unspecified* state
- *not movable* and *not copyable*
- implementation-defined macro `ATOMIC_FLAG_INIT` can be used to set flag to clear state in (static or automatic) initialization using statement of the form “`std::atomic_flag f = ATOMIC_FLAG_INIT;`”
- guaranteed to be *lock free*
- intended to be used as building block for higher-level synchronization primitives, such as spinlock mutex

Member Functions

Member Name	Description
constructor	constructs object
clear	atomically sets flag to false
test_and_set	atomically sets flag to true and obtains its previous value

Example: Suboptimal Spinlock Mutex

```
1  #include <iostream>
2  #include <thread>
3  #include <atomic>
4  #include <mutex>
5
6  class SpinLockMutex {
7  public:
8      SpinLockMutex() {f_.clear();}
9      void lock() {while (f_.test_and_set()) {}}
10     void unlock() {f_.clear();}
11 private:
12     std::atomic_flag f_; // true if thread holds mutex
13 };
14
15 SpinLockMutex m;
16 unsigned long long counter = 0;
17
18 void doWork() {
19     for (unsigned long long i = 0; i < 1000000ULL; ++i)
20         {std::lock_guard<SpinLockMutex> lock(m); ++counter;}
21 }
22
23 int main() {
24     std::thread t1(doWork), t2(doWork);
25     t1.join(); t2.join();
26     std::cout << counter << '\n';
27 }
```

- default memory order is suboptimal (and will be revisited later)

Example: One-Time Wait

```
1  #include <iostream>
2  #include <atomic>
3  #include <thread>
4  #include <chrono>
5
6  // notReady flag initially not set
7  std::atomic_flag notReady = ATOMIC_FLAG_INIT;
8  int result = 0;
9
10 int main() {
11     notReady.test_and_set(); // indicate result not ready
12     std::thread producer([](){
13         std::this_thread::sleep_for(std::chrono::seconds(1));
14         result = -42;
15         notReady.clear(); // indicate result ready
16     });
17     std::thread consumer([](){
18         // loop until result ready
19         while (notReady.test_and_set()) {}
20         std::cout << result << '\n';
21     });
22     producer.join();
23     consumer.join();
24 }
```

- This is *not* a particularly good use of `atomic_flag`.

The `std::atomic` Template Class

- `std::atomic` class provides types with atomic operations
- declaration:

```
template <class T> struct atomic;
```
- provides object of type `T` with atomic operations
- has partial specializations for integral types and pointer types
- full specializations for all fundamental types
- in order to use class type for `T`, `T` must be trivially copyable and bitwise equality comparable
- not required to be lock free
- on most popular platforms `atomic` is lock free when `T` is built-in type
- *not move constructible* and *not copy constructible*
- assignable but assignment operator returns value not reference
- most operations have memory order argument
- default memory order is SC (`std::memory_order_seq_cst`)

Basic

Member Name	Description
constructor	constructs object
operator=	atomically store value into atomic object
is_lock_free	check if atomic object is lock free
store	atomically replaces value of atomic object with given value
load	atomically reads value of atomic object
operator T	obtain result of <code>load</code>
exchange	atomically replaces value of atomic object with given value and obtain value of previous value
compare_exchange_weak	similar to <code>exchange_strong</code> but may fail spuriously
compare_exchange_strong	atomically compare value of atomic object to given value and perform <code>exchange</code> if equal or <code>load</code> otherwise

Fetch

Member Name	Description
<code>fetch_add</code>	atomically adds given value to value stored in atomic object and obtains value held previously
<code>fetch_sub</code>	atomically subtracts given value from value stored in atomic object and obtains value held previously
<code>fetch_and</code>	atomically replaces value of atomic object with bitwise AND of atomic object's value and given value, and obtains value held previously
<code>fetch_or</code>	atomically replaces value of atomic object with bitwise OR of atomic object's value and given value, and obtains value held previously
<code>fetch_xor</code>	atomically replaces value of atomic object with bitwise XOR of atomic object's value and given value, and obtains value held previously

std::atomic Members (Continued 2)

Increment and Decrement

Member Name	Description
operator++	atomically increment the value of atomic object by one and obtain value after incrementing
operator++ (int)	atomically increment the value of atomic object by one and obtain value before incrementing
operator--	atomically decrement the value of atomic object by one and obtain value after decrementing
operator-- (int)	atomically decrement the value of atomic object by one and obtain value after decrementing

Compound Assignment

Member Name	Description
operator+=	atomically adds given value to value stored in atomic object
operator-=	atomically subtracts given value from value stored in atomic object
operator&=	atomically performs bitwise AND of given value with value stored in atomic object
operator =	atomically performs bitwise OR of given value with value stored in atomic object
operator^=	atomically performs bitwise XOR of given value with value stored in atomic object

Example: Atomic Counter

```
1  #include <iostream>
2  #include <vector>
3  #include <thread>
4  #include <atomic>
5
6  class AtomicCounter {
7  public:
8      AtomicCounter() : c_(0) {}
9      int operator++() {return ++c_;}
10     int get() const {return c_.load();}
11 private:
12     std::atomic<int> c_;
13 };
14
15 AtomicCounter counter;
16
17 void doWork() {
18     for (int i = 0; i < 10000; ++i) {++counter;}
19 }
20
21 int main() {
22     std::vector<std::thread> v;
23     for (int i = 0; i < 10; ++i)
24         {v.emplace_back(doWork);}
25     for (auto& t : v) {t.join();}
26     std::cout << counter.get() << '\n';
27 }
```

Example: Atomic Increment With Compare and Swap

```
1 #include <atomic>
2
3 template <class T>
4 void atomicIncrement (std::atomic<T>& x) {
5     T curValue = x;
6     while (!x.compare_exchange_weak (curValue,
7         curValue + 1)) {}
8 }
```

Example: Counting Contest

```
1  #include <iostream>
2  #include <vector>
3  #include <atomic>
4  #include <thread>
5
6  constexpr int numThreads = 10;
7  std::atomic<bool> ready(false);
8  std::atomic<bool> done(false);
9  std::atomic<int> startCount(0);
10
11 void doCounting(int id) {
12     ++startCount;
13     while (!ready) {}
14     for (volatile int i = 0; i < 20000; i++) {}
15     bool expected = false;
16     if (done.compare_exchange_strong(expected, true))
17         {std::cout << "winner: " << id << '\n';}
18 }
19
20 int main() {
21     std::vector<std::thread> threads;
22     for (int i = 0; i < numThreads; ++i)
23         {threads.emplace_back(doCounting, i);}
24     while (startCount != numThreads) {}
25     ready = true;
26     for (auto& t : threads) {t.join();}
27 }
```

An Obligatory Note on `volatile`

- `volatile` qualifier not useful for multithreaded programming
- `volatile` qualifier makes *no guarantee of atomicity*
- can create object of `volatile`-qualified type whose size is sufficiently large that no current processor can access object atomically
- some platforms may happen to guarantee memory operations on (suitably-aligned) `int` object to be atomic, but in such cases this is normally true *even without `volatile` qualifier*
- `volatile` qualifier *does not adequately address issue of memory consistency*
- `volatile` qualifier does not imply use of memory barriers or other mechanisms needed for memory consistency
- optimizer and hardware might reorder operations (on non-`volatile` objects) across operations on `volatile` objects

Section 3.3.10

Atomics and the Memory Model

Semantics of Multithreaded Programs

- To be able to reason about the behavior of a program, we must know:
 - the *order* in which the operations of the program are performed; and
 - when the effects of each operation become *visible* to other operations in the program, which may be performed in different threads.
- In a single-threaded program, the ordering of operations and when the effects of operations become visible is quite intuitive.
- In a multi-threaded program, this matter becomes *considerably more complicated*.
- In what follows, we examine the above matter more closely (which essentially relates to the memory model).

Happens-Before Relationships

- For two operations A and B performed in the *same or different* threads, A is said to **happen before** B if the effects of A become *visible* to the thread performing B before B is performed.
- The happens-before relationship is *not equivalent* to “happens earlier in time”.
- If operation A happens earlier in time than operation B , this does not imply that the effects of A must be *visible* to the thread performing B before B is performed, due to the effects of caches, store buffers, and so on, which *delay* the visibility of results.
- Happening earlier in time is only a necessary but not sufficient condition for a happens-before relationship to exist.
- Happens-before relationships are *not always transitive*.
- In the absence of something known as a dependency-ordered-before relationship (to be discussed later), which arise relatively less frequently, happens-before relationships are *transitive* (i.e., if A happens before B and B happens before C then A happens before C).

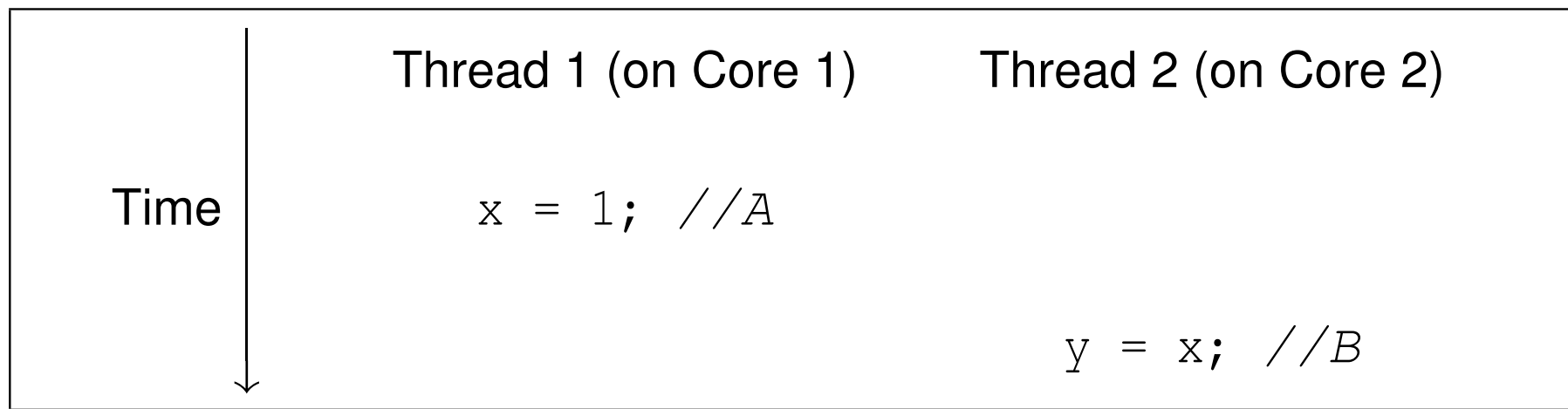
“Earlier In Time” Versus Happens Before

- Consider the multithreaded program (with two threads) shown below, where x and y are integer variables, *initially zero*.

```
Thread 1 Code
x = 1; // A
```

```
Thread 2 Code
y = x; // B
```

- Suppose that the run-time platform is such that memory operations on x are *atomic* so the program is data-race free.
- Consider what happens when the program executes with the particular timing shown below, where *operation A occurs earlier in time than operation B*.



- The value read for x in operation B will not necessarily be 1.

Sequenced-Before Relationships

- Given two operations A and B performed in the *same* thread, the operation A is **sequenced before** B if A precedes B in program order (i.e., source-code order).
- Sequenced-before relationships are *transitive* (i.e., if A is sequenced before B , and B is sequenced before C , then A is sequenced before C).
- Example: In the code below, statement A is sequenced before statement B ; B is sequenced before statement C ; and, by transitivity, A is sequenced before C .

```
x = 1;      // A
y = 2;      // B
z = x + 1;  // C
```

- Example:
 - Consider the line of code below, which performs (in order) the following operations: 1) multiplication, 2) addition, and 3) assignment.

```
y = a * x + b; // (y = ((a * x) + b));
```

- Multiplication is sequenced before addition.
- Addition is sequenced before assignment.
- Thus, by transitivity, multiplication is sequenced before assignment.

Sequenced-Before Relationships (Continued)

- For two operations A and B in the *same* thread, if A is *sequenced before* B then A *happens before* B .
- In other words, program order establishes happens-before relationships for operations *within a single thread*.
- A sequenced-before relationship is essentially an *intra-thread happens-before* relationship. (Note that “intra” means “within”.)
- Example: In the code below, statement A is sequenced before statement B . Therefore, A happens before B . Similarly, B happens before statement C , and A happens before C .

```
x = 1;      // A
y = 2;      // B
z = x + 1;  // C
```

Inter-Thread Happens-Before Relationships

- Establishing whether a happens-before relationship exists between operations in different threads is somewhat more complicated than the same-thread case.
- Inter-thread happens-before relationships establish happens-before relationships for operations in *different* threads.
- For two operations A and B in *different* threads, if A **inter-thread happens before** B then A happens before B .
- Inter-thread happens-before relationships are *transitive* (i.e., if A inter-thread happens before B and B inter-thread happens before C then A inter-thread happens before C).
- Some form of *synchronization* is required to establish an inter-thread happens-before relationship.
- The various forms that this synchronization may take will be introduced on later slides.

Summary of Happens-Before Relationships

- For two operations A and B in either the *same or different* threads, A happens before B if:
 - ① A and B are in the *same* thread and A is sequenced before (i.e., intra-thread happens before) B ; or
 - ② A and B are in *different* threads and A inter-thread happens before B .
- In other words, A happens before B if A either intra-thread happens before or inter-thread happens before B .
- Intra-thread happens-before (i.e., sequenced-before) relationships are *transitive*.
- Inter-thread happens-before relationships are *transitive*.
- Happens-before relationships are *mostly but not always transitive*.
- A happens-before relationship is important because it tells us if the result of one operation *can be seen* by a thread performing another operation.

Synchronizes-With Relationships

- A variety of relationships can imply an inter-thread happens-before relationship, with one being the synchronizes-with relationship.
- For two operations A and B in *different* threads, if A **synchronizes with** B then A *inter-thread happens before* B .
- Example:
 - Consider the two-threaded program shown below, with the shared variable x of type `int`, where x is initially zero.

Thread 1 Code

```
1  x = 1;  
2  // A (call of foo)  
3  foo();
```

Thread 2 Code

```
1  bar();  
2  // B (return from bar)  
3  assert(x == 1);
```

- Suppose that the call of the function `foo` is known to *synchronize with* the return from the function `bar`, which implies that A synchronizes with B .
- Since A synchronizes with B , A must inter-thread happen before B , which implies that *A happens before B* .
- Therefore, the assertion in thread 2 *can never fail*.

Examples of Synchronizes-With Relationships

- **Thread creation.** The completion of the constructor for a `thread` object T synchronizes with the start of the invocation of the thread function for T .
- **Thread join.** The completion of the execution of a thread function for a `thread` object T synchronizes with (the return of) a `join` operation on T .
- **Mutex unlock/lock.** All prior `unlock` operations on a mutex M synchronize with (the return of) a `lock` operation on M .
- **Atomic.** A suitably tagged atomic write operation W on a variable x synchronizes with a suitably tagged atomic read operation on x that reads the value stored by W (where the meaning of “suitably tagged” will be discussed later).

Synchronizes-With Relationship: Thread Create and Join

```
1  #include <thread>
2  #include <cassert>
3
4  int x = 0;
5
6  void doWork() {
7      // A1 (start of thread execution)
8      assert(x == 1); // OK: M1 synchronizes with A1
9      x = 2;
10     // A2 (end of thread execution)
11 }
12
13 int main() {
14     x = 1;
15     std::thread t(doWork); // M1 (completion of constructor)
16     t.join(); // M2 (return from join)
17     assert(x == 2); // OK: A2 synchronizes with M2
18 }
```

- since construction of thread (M1) synchronizes with start of thread function execution (A1), M1 happens before A1 implying that assertion in doWork cannot fail
- since completion of execution of thread function (A2) synchronizes with join operation (M2), A2 happens before M2 implying that assertion in main cannot fail

Synchronizes-With Relationship: Mutex Lock/Unlock

Shared Data

```
std::mutex m;  
int x = 0;  
int y = 0;
```

Thread 1 Code

```
m.lock();  
x = 1;  
m.unlock();
```

Thread 2 Code

```
m.lock();  
y = x;  
m.unlock();
```

Thread 1 Execution

m.lock();

x = 1;

m.unlock();

Ⓐ

synchronizes with

Thread 2 Execution

m.lock();

Ⓑ

y = x;

m.unlock();

Since unlock synchronizes with lock,

Ⓐ happens before Ⓑ

Memory Orders

- Most operations on atomic types allow a memory order to be specified.
- Example:

```
std::atomic<int> x = 0;  
x.store(42, std::memory_order_seq_cst);  
int y = x.load(std::memory_order_seq_cst);
```

- The following memory orders are supported:
 - sequentially consistent (`std::memory_order_seq_cst`)
 - acquire-release (`std::memory_order_acq_rel`)
 - acquire (`std::memory_order_acquire`)
 - release (`std::memory_order_release`)
 - consume (`std::memory_order_consume`)
 - relaxed (`std::memory_order_relaxed`)
- Read operations can use the orders:
 - sequentially consistent, acquire, consume, and relaxed.
- Write operations can use the orders:
 - sequentially consistent, release, and relaxed.
- Read-modify-write operations can use:
 - all of the orders allowed for read and write operations; and
 - acquire-release.

Memory Models

- Although several memory orders can be employed for operations on atomic types, these orders support *four basic models*:
 - 1 sequentially consistent,
 - 2 acquire release,
 - 3 consume release, and
 - 4 relaxed.
- These models differ in the guarantees that they make regarding:
 - whether all writes to all atomic objects become visible to *all* threads *simultaneously* (i.e., total order for all writes to all atomic objects); and
 - whether operations on atomic objects in different threads can establish a *synchronization* relationship (namely, a synchronizes-with or dependency-ordered-before [discussed later] relationship).
- The models listed from strongest (i.e., makes the most guarantees) to weakest (i.e., makes the least guarantees) are:
 - 1 sequentially consistent,
 - 2 acquire release,
 - 3 consume release, and
 - 4 relaxed.

Memory Models (Continued 1)

- These models are *hierarchical* in the sense that each model makes at least all of the same guarantees as its weaker counterparts.
- As we proceed from stronger to weaker models, more guarantees are lost.
- A stronger model may require additional synchronization by hardware, which can *degrade performance*.
- A weaker model *may not provide sufficient guarantees* for the correct functioning of code.
- Using a model that fails to provide sufficient guarantees for correct code behavior will result in *bugs*.
- Also, as the model is weakened, it becomes more difficult to reason about the behavior of code, leading to *incomprehensible code* and an *increased likelihood of (often very subtle) bugs*.

Modification Order

- All writes to a particular atomic object M (over its lifetime) occur in some particular total order, called its **modification order**.
- Each atomic object has its own well-defined modification order.
- For a particular atomic object M , *all* threads in a program are guaranteed to see M change in a manner *consistent with its modification order*.
- Essentially, this guarantee ensures that, once a given thread has seen a particular value of an atomic object, a subsequent read by that thread cannot retrieve an earlier value of the object.
- If such a guarantee were not made, the memory model would be so weak as to be impractical to use.
- Modification order is primarily a *conceptual* tool that is useful for describing memory-model behavior.
- In practice, a thread is unlikely to actually observe every change in the modification order of an object.

Modification Order (Continued)

- For each atomic object M , each thread has its own current position in object's modification order.
- A thread's current position in the modification order of a particular atomic object need not be the same for all threads.
- A read from an atomic object M by a thread T can *optionally* move T 's current position to a later position in the modification order of M and then returns the value at the current position.
- A write to an atomic object M by a thread T appends the value to be written to the modification order of M and updates T 's current position in the modification order of M to correspond to the value written.
- An read-modify-write operation A on an atomic object M reads the *last* value in the modification order of M , modifies the value read appropriately, appends the resulting value to the modification order of M , and updates T 's current position in the modification order of M to correspond to the value written.

Modification Order Example

- Consider an atomic object M with the modification sequence:
 - 0, 1, 2, 3, 4, 5, 6, 7, 8.
- A thread could, for example, legitimately see M undergo any of the following sequences of updates:
 - 0, 4, 8
 - 8
 - 2, 7
 - 0, 1, 2, 5, 7, 8
 - 0, 1, 2, 3, 4, 5, 6, 7, 8
- A thread would, for example, be guaranteed *never* to see M undergo any of the following sequences of updates, as all of these sequences are *inconsistent* with the modification order of M :
 - 1, 0
 - 1, 2, 1
 - 42
 - 0, 1, 2, 3, 4, 5, 6, 7, 6, 8

Relative Ordering of Changes to Different Atomic Objects

- Although each atomic object has its own well-defined modification order, it is not necessarily the case that the modification orders for individual objects can be combined into a single total order over *all* atomic objects.
- Practically speaking, the reason for this is the delay in the visibility of results introduced by store buffers, caches, and so on.
- If a single total order for writes to all atomic objects is not guaranteed, this implies that the relative order of changes to *different* atomic objects need not appear the same to different threads.
- Ensuring the existence of a single total order over all atomic objects would require a significant amount of additional processor synchronization, which can significantly degrade performance.
- Therefore, this guarantee is not required to be made in all cases, the idea being that we only ask for the guarantee when it is needed for correct code behavior.

Modification Order Revisited

- Consider a program with two threads and two shared integer atomic objects x and y , each having the modification order: 0, 1.
- Suppose that no requirement is imposed to guarantee the existence of a single total order on writes to *all* atomic objects.
- Thread 1 could see x and y change in the following manner, consistent with their stated modification order:

Variable	Updates to Value Seen By Thread
x	0 1
y	0 1

- Thread 2 could see x and y change in the following manner, consistent with their stated modification order:

Variable	Updates to Value Seen By Thread
x	0 1
y	0 1

- Observe that thread 1 and thread 2 do not see x and y change in the same order relative to one another (i.e., thread 1 sees x change before y , while thread 2 sees y change before x).

Sequentially-Consistent Model

- The sequentially-consistent model simply corresponds to the default memory model for the language, namely, SC-DRF. (Since data races cannot occur on atomic objects, SC-DRF degenerates into SC for such objects.)
- For the sequentially-consistent model, all memory operations (i.e., read, write, and read-modify-write) must use the sequentially-consistent memory order (`std::memory_order_seq_cst`).
- A *total ordering* is guaranteed on all sequentially-consistent writes to *all* atomic objects.
- All sequentially-consistent writes to atomic objects must become **visible** to all threads *simultaneously*.
- A sequentially-consistent write operation W on an atomic object M (in one thread) *synchronizes with* a sequentially-consistent operation on M (in another thread) that reads the value written by W .
- This model allows for relatively *easy reasoning* about code behavior.

Example: Sequentially-Consistent Model

- shared data:

`x` and `y` are of type `std::atomic<int>` and both are initially zero

- thread 1 code (writes `x`):

```
x.store(1, std::memory_order_seq_cst);
```

- thread 2 code (writes `y`):

```
y.store(1, std::memory_order_seq_cst);
```

- thread 3 code (reads `x` then `y`):

```
int x1 = x.load(std::memory_order_seq_cst);  
int y1 = y.load(std::memory_order_seq_cst);
```

- thread 4 code (reads `y` then `x`):

```
int y2 = y.load(std::memory_order_seq_cst);  
int x2 = x.load(std::memory_order_seq_cst);
```

- memory order guarantees total order for all writes to all atomic objects

- so, thread 3 and thread 4 must agree about order in which `x` and `y` are modified

- not possible to see `x1 == 1` and `y1 == 0` in thread 3 (implying `x` modified before `y`) and `x2 == 0` and `y2 == 1` in thread 4 (implying `y` modified before `x`)

Example: Sequentially-Consistent Model

```
1  #include <atomic>
2  #include <thread>
3  #include <cassert>
4
5  std::atomic<int> x, y, c;
6
7  void w_x() {x.store(1, std::memory_order_seq_cst);}
8
9  void w_y() {y.store(1, std::memory_order_seq_cst);}
10
11 void r_xy() {
12     while (!x.load(std::memory_order_seq_cst)) {}
13     if (y.load(std::memory_order_seq_cst)) {++c;}
14 }
15
16 void r_yx() {
17     while (!y.load(std::memory_order_seq_cst)) {}
18     if (x.load(std::memory_order_seq_cst)) {++c;}
19 }
20
21 int main() {
22     x = 0; y = 0; c = 0;
23     std::thread t1(w_x), t2(w_y), t3(r_xy), t4(r_yx);
24     t1.join(); t2.join(); t3.join(); t4.join();
25     assert(c != 0); // assertion cannot fail
26 }
```

- assertion cannot fail: when **while** loop in `r_xy` terminates, all threads must see `x` as nonzero; when **while** loop in `r_yx` terminates, all threads must see `y` as nonzero; at least one of these must happen before **if** statements in both `r_xy` and `r_yx` executed

Acquire-Release Model

- For the acquire-release model, the memory order is chosen as follows:
 - a read operation uses the acquire order (`std::memory_order_acquire`)
 - a write operation uses the release order (`std::memory_order_release`)
 - a read-modify-write operation uses one of the orders allowed for read and write operations, or the acquire-release order (`std::memory_order_acq_rel`), which results in read acquire and write release.
- *No total ordering* exists on all writes to *all* atomic objects (unlike in the sequentially-consistent model).
- Consequently, threads do not necessarily have to agree on the *relative order* in which different atomics objects are modified.
- A write-release operation W on an atomic object M *synchronizes with* a read-acquire operation on M that reads the value written by W (or a value written by the release sequence headed by W).
- The acquire-release model is useful for situations that involve *pairwise synchronization* of threads, such as with mutexes.
- With the acquire-release model, it is often still possible to reason about code behavior without too much difficulty.

Example: Acquire-Release Model

- shared data:

`x` and `y` are of type `std::atomic<int>` and both are initially zero

- thread 1 code (writes `x`):

```
x.store(1, std::memory_order_release);
```

- thread 2 code (writes `y`):

```
y.store(1, std::memory_order_release);
```

- thread 3 code (reads `x` then `y`):

```
int x1 = x.load(std::memory_order_acquire);  
int y1 = y.load(std::memory_order_acquire);
```

- thread 4 code (reads `y` then `x`):

```
int y2 = y.load(std::memory_order_acquire);  
int x2 = x.load(std::memory_order_acquire);
```

- no ordering relationship between stores to `x` and `y`
- so, thread 3 and thread 4 do not need to agree about order in which `x` and `y` are modified
- possible to see `x1 == 1` and `y1 == 0` in thread 3 (i.e., thread 3 sees `x` change before `y`) and `x2 == 0` and `y2 == 1` in thread 4 (i.e., thread 4 sees `y` change before `x`)

Example: Acquire-Release Model

```
1  #include <atomic>
2  #include <thread>
3  #include <cassert>
4
5  std::atomic<int> x, y, c;
6
7  void w_x() {x.store(1, std::memory_order_release);}
8
9  void w_y() {y.store(1, std::memory_order_release);}
10
11 void r_xy() {
12     while (!x.load(std::memory_order_acquire)) {}
13     if (y.load(std::memory_order_acquire)) {++c;}
14 }
15
16 void r_yx() {
17     while (!y.load(std::memory_order_acquire)) {}
18     if (x.load(std::memory_order_acquire)) {++c;}
19 }
20
21 int main() {
22     x = 0; y = 0; c = 0;
23     std::thread t1(w_x), t2(w_y), t3(r_xy), t4(r_yx);
24     t1.join(); t2.join(); t3.join(); t4.join();
25     assert(c != 0); // assertion can fail
26 }
```

- assertion can fail: one thread seeing x or y being nonzero does not imply other thread sees same

Example: Spinlock Mutex Using `std::atomic_flag`

```
1  #include <iostream>
2  #include <thread>
3  #include <atomic>
4
5  class SpinLockMutex {
6  public:
7      SpinLockMutex() {f_.clear();}
8      void lock() {
9          while (f_.test_and_set(std::memory_order_acquire)) {}
10     }
11     void unlock() {f_.clear(std::memory_order_release);}
12 private:
13     std::atomic_flag f_; // true if thread holds mutex
14 };
15
16 SpinLockMutex m;
17 unsigned long long counter = 0;
18
19 void doWork() {
20     for (unsigned long long i = 0; i < 1000000ULL; ++i)
21         {m.lock(); ++counter; m.unlock();}
22 }
23
24 int main() {
25     std::thread t1(doWork), t2(doWork);
26     t1.join(); t2.join();
27     std::cout << counter << '\n';
28 }
```

● uses acquire-release model



Example: Spinlock Mutex and `std::lock_guard`

```
1  #include <iostream>
2  #include <thread>
3  #include <atomic>
4  #include <mutex>
5
6  class SpinLockMutex {
7  public:
8      SpinLockMutex() {f_.clear();}
9      void lock() {
10         while (f_.test_and_set(std::memory_order_acquire)) {}
11     }
12     void unlock() {f_.clear(std::memory_order_release);}
13 private:
14     std::atomic_flag f_; // true if thread holds mutex
15 };
16
17 SpinLockMutex m;
18 unsigned long long counter = 0;
19
20 void doWork() {
21     for (unsigned long long i = 0; i < 1000000ULL; ++i)
22         {std::lock_guard<SpinLockMutex> lg(m); ++counter;}
23 }
24
25 int main() {
26     std::thread t1(doWork), t2(doWork);
27     t1.join(); t2.join();
28     std::cout << counter << '\n';
29 }
```

Carries-A-Dependency Relationships

- For two operations A and B performed in the *same* thread, A is said to **carry a dependency** to B if the result of A is used as an operand for B (ignoring some special cases).
- Example: In the code below, statement A *carries a dependency* to statement B but not statement C .

```
x = 42;    // A
y = x + 1; // B
z = 0;    // C
```

- Note that “carries a dependency to” is a subset of “is sequenced before” (i.e., the former implies the latter).
- The carries-a-dependency-to relationship is *transitive* (i.e., if A carries a dependency to B and B carries a dependency to C then A carries a dependency to C).
- Example: In the code below, statement A carries a dependency to statement B ; and B carries a dependency to statement C . Therefore, transitively, A carries a dependency to C .

```
x = 42;    // A
y = x + 1; // B
z = 2 * y; // C
```

Dependency-Ordered-Before Relationships

- Another type of synchronization relationship is known as a dependency-ordered-before relationship.
- A write-release operation A is *dependency ordered before* a read-consume operation B if B reads the value written by A (or any side effect in the release sequence headed by A).
- For two operations A and B performed in *different* threads, if A is dependency ordered before B then A *inter-thread happens before* B .
- Thus, dependency-ordered-before relationships can also establish happens-before relationships.

Inter-Thread Happens-Before Relationships Revisited

- The inter-thread happens before relation describes an *arbitrary concatenation* of sequenced-before, synchronizes-with, and dependency-ordered-before relations, *with two exceptions*:
 - ① a concatenation is not permitted to end with dependency ordered before followed by (one or more) sequenced before; and
 - ② a concatenation is not permitted to consist entirely of sequenced-before relations.
- The first restriction is required since a dependency-ordered-before relationship synchronizes *only data dependencies*.
- The second restriction is required since inter-thread happens-before relationship must (by definition) involve operations in *different* threads.

Consume-Release Model

- For the consume-release model, the memory order is chosen as follows:
 - a write operation uses release order (`std::memory_order_release`)
 - a read operation uses the consume order (`std::memory_order_consume`)
- The consume-release model is identical to the acquire-release model with one important difference, namely the type of synchronization relationship established.
- A write-release operation W is *dependency ordered before* a read-consume operation (in a different thread) that reads the value stored by W (or any side effect in the release sequence headed by W).
- In other words, the consume-release model establishes a *dependency-ordered-before* relationship, whereas the acquire-release model establishes a *synchronizes-with* relationship.
- In this sense, the consume-release model is weaker than the acquire-release model (i.e., less data is synchronized).

Example: Consume-Release Model

```
1  #include <thread>
2  #include <atomic>
3  #include <cassert>
4
5  int x = 0;
6  std::atomic<int> y(0);
7
8  void producer() {
9      x = 42;
10     y.store(1, std::memory_order_release);
11 }
12
13 void consumer() {
14     int a;
15     while (!(a = y.load(std::memory_order_consume))) {}
16     assert(x == 42); // data race
17 }
18
19 int main() {
20     std::thread t1(producer);
21     std::thread t2(consumer);
22     t1.join();
23     t2.join();
24 }
```

- program has *data race* on `x`; `a` does not carry dependency to `x` so `x = 42` does not necessarily happen before `x` used in assertion
- if `consume` changed to `acquire`, no data race and assertion cannot fail

Example: Publishing Data Via Pointer

```
1  #include <thread>
2  #include <atomic>
3  #include <cassert>
4  #include <string>
5
6  std::atomic<std::string*> p(nullptr);
7  int x = 0;
8
9  void producer() {
10     std::string* s = new std::string("Hello");
11     x = 42;
12     p.store(s, std::memory_order_release);
13 }
14
15 void consumer() {
16     std::string* s;
17     while (!(s = p.load(std::memory_order_consume))) {}
18     assert(*s == "Hello");
19     // assert(x == 42); would result in data race
20 }
21
22 int main() {
23     std::thread t1(producer), t2(consumer);
24     t1.join(); t2.join();
25 }
```

- assertion cannot fail; store to `p` is dependency ordered before load and load carries dependency to `*s` in assertion

Relaxed Model

- For the relaxed model, all memory operations use the relaxed order (`std::memory_order_relaxed`).
- Like in the acquire-release model, *no total order* exists on updates to *all* atomic objects (collectively).
- Operations on the same variable *within a single thread* satisfy a happens-before relationship (i.e., within a single thread, accesses to a single atomic variable must follow program order).
- Unlike in the acquire-release model, *no inter-thread synchronization* relationship is established.
- No requirement exists on the ordering relative to other threads.
- The relaxed order is sometime suitable for updating counters (e.g., blind event counters).
- Except in very trivial cases, it can be *extremely difficult to reason* about the meaning and/or correctness of code that uses relaxed order.

Behavior of Relaxed Model

- consider atomic memory operations with relaxed order
- for each individual atomic object, all threads have view of updates that is consistent with single modification sequence
- read operation (e.g., `load`):
 - if current position not set, return any element in sequence and set current position to that of returned element
 - otherwise, either leave current position unchanged or move later in sequence and return value at current position
- write operation (e.g., `store`):
 - append value to end of sequence
 - set current position to correspond to appended value
- read-modify-write operation (e.g., `increment`, `decrement`, `exchange`, `compare_exchange`):
 - read last value from sequence
 - modify read value as appropriate to obtain new value
 - append new value to end of sequence
 - set current position to correspond to that of appended value
- considerable flexibility in value returned by read

Example: Relaxed Model

```
1  #include <atomic>
2  #include <thread>
3  #include <cassert>
4
5  std::atomic<int> x, y, c;
6
7  void w_x() {x.store(1, std::memory_order_relaxed);}
8
9  void w_y() {y.store(1, std::memory_order_relaxed);}
10
11 void r_xy() {
12     while (!x.load(std::memory_order_relaxed)) {}
13     if (y.load(std::memory_order_relaxed)) {++c;}
14 }
15
16 void r_yx() {
17     while (!y.load(std::memory_order_relaxed)) {}
18     if (x.load(std::memory_order_relaxed)) {++c;}
19 }
20
21 int main() {
22     x = 0; y = 0; c = 0;
23     std::thread t1(w_x), t2(w_y), t3(r_xy), t4(r_yx);
24     t1.join(); t2.join(); t3.join(); t4.join();
25     assert(c != 0); // assertion can fail
26 }
```

- assertion can fail: one thread seeing x or y being nonzero does not imply other thread sees same

Example: Blind Event Counters

```
1  #include <vector>
2  #include <iostream>
3  #include <thread>
4  #include <atomic>
5
6  std::atomic<unsigned long long> counter(0);
7
8  void doWork() {
9      for (long i = 0; i < 100'000L; ++i) {
10         counter.fetch_add(1, std::memory_order_relaxed);
11     }
12 }
13
14 int main() {
15     std::vector<std::thread> workers;
16     for (int i = 0; i < 10; ++i) {
17         workers.emplace_back(doWork);
18     }
19     for (auto& t : workers) {
20         t.join();
21     }
22     std::cout << "counter " << counter << '\n';
23 }
```

- `fetch_add` can use *relaxed* order, since only incrementing counter *blindly* (i.e., not taking action based on value of counter)
- thread join operations provide synchronization to ensure desired value read for counter when output

Example: Done Flag

```
1  #include <vector>
2  #include <thread>
3  #include <atomic>
4  #include <chrono>
5
6  std::atomic<bool> done;
7
8  void doWork() {
9      while (!done.load(std::memory_order_relaxed)) {
10         // do something here
11     }
12 }
13
14 int main() {
15     std::vector<std::thread> workers;
16     done.store(false, std::memory_order_relaxed); // I hope? ;)
17     for (int i = 0; i < 16; ++i) {
18         workers.emplace_back(doWork);
19     }
20     std::this_thread::sleep_for(std::chrono::seconds(5));
21     done = true; // not relaxed
22     for (auto& t : workers) {
23         t.join();
24     }
25 }
```

- done.store can be relaxed due to synchronization from thread create
- done.load can be relaxed since order not important; different order as if other threads ran at different speeds
- assign to done must be sequentially-consistent to prevent assign from floating past join (due to single-thread optimization)

Example: `std::shared_ptr` Reference Counting

- The copy constructor for `shared_ptr` (which increments a reference count) would look something like:

```
// ...
controlBlockPtr = other->controlBlockPtr;
controlBlockPtr->refCount.fetch_add(1,
    std::memory_order_relaxed);
// ...
```

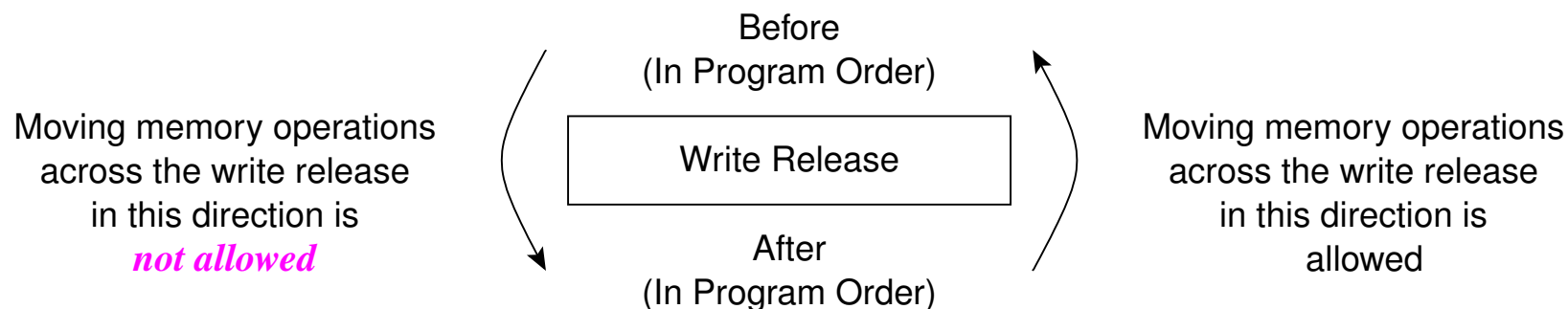
- The destructor for `shared_ptr` (which decrements a reference count) would look something like:

```
// ...
if (!controlBlockPtr->refCount.fetch_sub(1,
    std::memory_order_acq_rel)) {
    delete controlBlockPtr;
}
// ...
```

- The increment operation can use *relaxed* order, since *no action is taken* based on the reference count value.
- The decrement operation needs to use *acquire-release* order so that the decrement cannot float and the correct view of the data is seen by the thread doing the delete (all decrements form a *synchronization chain*).

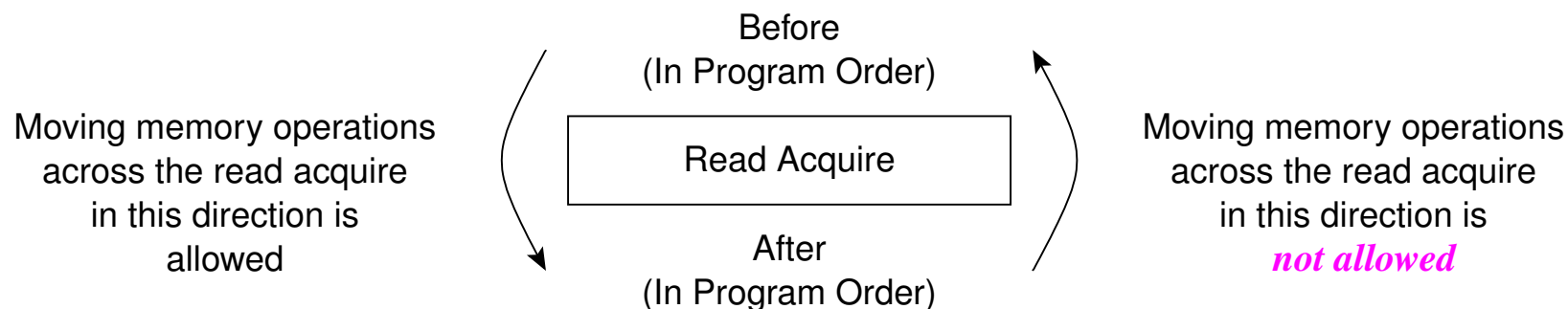
Release Semantics for Memory Operations

- Release semantics is a property that can only apply to operations that *write to memory* (i.e., read-modify-write operations or plain writes).
- A write operation that has release semantics is called a **write release**.
- A write release operation W cannot be reordered with any read or write operation that *precedes* W in program order (i.e., memory operations cannot be moved from before W to after W).
- The term release semantics originates from mutexes.
- In the context of mutexes, the operations prior to a mutex release operation, which correspond to operations in a critical section, must not be moved after the mutex release operation, as operations after the mutex release operation are not protected by the mutex.



Acquire Semantics for Memory Operations

- Acquire semantics is a property that can only apply to operations that *read from memory* (i.e., read-modify-write operations or plain reads).
- A read operation that has acquire semantics is called a **read acquire**.
- A read acquire operation R cannot be reordered with any read or write operation that *follows* R in program order (i.e., memory operations cannot be moved from after R to before R).
- The term acquire semantics originates from mutexes.
- In the context of mutexes, the operations following a mutex acquire operation, which correspond to operations in a critical section, must not be moved before the mutex acquire operation, as operations before the mutex acquire operation are not protected by the mutex.



- A **release sequence** headed by a release operation A on an atomic object M is a maximal contiguous subsequence of side effects in the modification order of M , where the first operation is A , and every subsequent operation
 - is performed by the same thread that performed A , or
 - is an atomic read-modify-write operation.

Release Sequence Example

```
1  #include <thread>
2  #include <atomic>
3  #include <cassert>
4
5  int x = 0;
6  std::atomic<int> y(0);
7
8  int main() {
9      std::thread t1([](){
10         x = 42;
11         y.store(1, std::memory_order_release); // A
12         y.store(2, std::memory_order_relaxed); // B
13     });
14     std::thread t2([](){
15         int r;
16         while ((r = y.load(std::memory_order_acquire)) // C
17             < 2) {}
18         assert(x == 42);
19     });
20     t1.join();
21     t2.join();
22 }
```

- stores to `y` in A and B constitute release sequence headed by store in A
- when while loop terminates, load in C will have read value written by store in B (not store in A)
- A synchronizes with C, since C reads value in release sequence headed by A
- assertion cannot fail, since A happens before C

Fences

- A **memory fence** (also known as a **memory barrier**) is an operation that causes the processor and compiler to enforce an *ordering constraint* on memory operations issued before and after the fence operation.
- Certain types of memory operations before a fence are guaranteed not to be reordered with certain types of memory operations after the fence.
- A fence may also introduce *synchronizes-with* relationships under certain circumstances.
- An **acquire fence** prevents the reordering of any *read or write* following the fence (in program order) with any *read* prior to the fence (in program order). (That is, a memory operation after the fence cannot be moved before any read operation before the fence.)
- A **release fence** prevents the reordering of any *read or write* prior to the fence (in program order) with any *write* following the fence (in program order). (That is, a memory operation before the fence cannot be moved after any write operation after the fence.)
- A fence is *not* a release or acquire operation. as it does not read/write memory.

- memory fences can be inserted via function

`std::atomic_thread_fence`

- declaration:

```
void atomic_thread_fence (std::memory_order order)  
noexcept;
```

- no effect if order is `std::memory_order_relaxed`
- acquire fence if order is `std::memory_order_acquire` or `std::memory_order_consume`
- release fence if order is `std::memory_order_release`
- both acquire and release fence if order is `std::memory_order_acq_rel`
- sequentially consistent acquire and release fence if order is `std::memory_order_seq_cst`

Fences and Synchronizes-With Relationships

- ***Release fence and acquire fence.*** A release fence A synchronizes with an acquire fence B if there exist atomic operations X and Y , both operating on some atomic object M , such that A is sequenced before X , X modifies M , Y is sequenced before B , and Y reads the value written by X or a value written by any side effect in the hypothetical release sequence X would head if it were a release operation.
- ***Release fence and acquire operation.*** A release fence A synchronizes with an atomic operation B that performs an acquire operation on an atomic object M if there exists an atomic operation X such that A is sequenced before X , X modifies M , and B reads the value written by X or a value written by any side effect in the hypothetical release sequence X would head if it were a release operation.
- ***Release operation and acquire fence.*** An atomic operation A that is a release operation on an atomic object M synchronizes with an acquire fence B if there exists some atomic operation X on M such that X is sequenced before B and reads the value written by A or a value written by any side effect in the release sequence headed by A .

Example: Incorrect Code Without Fence

```
1  #include <thread>
2  #include <atomic>
3  #include <iostream>
4
5  std::atomic<bool> ready(false);
6  int data = 0;
7
8  void produce() {
9      data = 42; // write to data can move after store in A
10     // release fence needed here
11     ready.store(true, std::memory_order_relaxed); // A
12 }
13
14 void consume() {
15     while (!ready.load(std::memory_order_relaxed)) {} // B
16     // acquire fence needed here
17     std::cout << data << '\n';
18     // read of data can move before load in B
19 }
20
21 int main() {
22     std::thread t1(produce);
23     std::thread t2(consume);
24     t1.join(); t2.join();
25 }
```

- atomic store (to ready) does not synchronize with atomic load (of ready), due to relaxed order; results in race on data

Example: Correct Code With Fence

```
1  #include <thread>
2  #include <atomic>
3  #include <iostream>
4
5  std::atomic<bool> ready(false);
6  int data = 0;
7
8  void produce() {
9      data = 42;
10     std::atomic_thread_fence(std::memory_order_release);
11     ready.store(true, std::memory_order_relaxed);
12 }
13
14 void consume() {
15     while (!ready.load(std::memory_order_relaxed)) {}
16     std::atomic_thread_fence(std::memory_order_acquire);
17     std::cout << data << '\n';
18 }
19
20 int main() {
21     std::thread t1(produce);
22     std::thread t2(consume);
23     t1.join(); t2.join();
24 }
```

- release fence synchronizes with acquire fence, due to atomic load (of ready) reading from result of atomic store (to ready)

Memory Orders: The Bottom Line

- Use sequentially-consistent order unless there is a compelling case to do otherwise.
- In situations where semantics dictate a clear pairwise synchronization between threads, consider the use of acquire-release order if it can be easily seen to yield correct code.
- Only consider relaxed order in situations where the performance penalty of using a stronger order would be unacceptable.
- *Be very weary of using relaxed order.* Even world experts on the C++ memory model acknowledge that this can be tricky.
- Always have any code using relaxed order thoroughly reviewed by people who are extremely knowledgeable about memory models.

Section 3.3.11

References

References I

- 1 A. Williams. *C++ Concurrency in Action*.
Manning Publications, Shelter Island, NY, USA, 2012.
This is a fairly comprehensive book on concurrency and multithreaded programming in C++. It is arguably the best book available for those who want to learn how to write multithreaded code using C++. Excellent
- 2 M. J. Batty. *The C11 and C++11 Concurrency Model*.
PhD thesis, University of Cambridge, Cambridge, UK, Nov. 2014.
This very well written Ph.D. thesis introduces the C++11/C11 memory model and presents work in mathematically formalizing, refining, and validating this model. Excellent
- 3 M. Batty. *Multicore Programming: C++0x*.
University of Cambridge, Cambridge, UK, Nov. 2010.
This set of slides appear to have been used for part of a course on multicore programming at the University of Cambridge.

- 4 M. Herlihy and N. Shavit. *The Art of Multiprocessor Programming*. Morgan Kaufmann, Burlington, MA, USA, 2008.
A good reference for concurrent programming.

- 1 Herb Sutter. `atomic<>` Weapons: The C++11 Memory Model and Modern Hardware, C++ and Beyond, Asheville, NC, USA, Aug. 5–8, 2012. (This talk is in two parts.)
- 2 Herb Sutter. C++ Concurrency, C++ and Beyond, Asheville, NC, USA, Aug. 5–8, 2012.
- 3 Herb Sutter. Lock-Free Programming (Or, Juggling Razor Blades), CppCon, 2014. (This talk is in two parts.)
- 4 Hans-J. Boehm. Threads and Shared Variables in C++11. Going Native, Redmond, WA, USA, Feb. 2–3, 2012.
- 5 Mike Long. Introducing the C++ Memory Model. Norwegian Developers Conference, Oslo, Norway, Jun. 15–19, 2014.
- 6 Herb Sutter. Machine Architecture and You: Things Your Programming Language Never Told You. Northwest C++ Users' Group, Redmond, WA, USA, Sept. 19, 2007. <http://nwcpp.org/september-2007.html>.

- 7 Pablo Halpern. Overview of Parallel Programming in C++, CppCon, Bellevue, WA, USA, Sept. 8, 2014.
- 8 Valentin Ziegler. C++ Memory Model, Meeting C++, Berlin, Germany, Dec. 6, 2014.

Part 4

Even More C++

Section 4.1

Undefined Behavior and Other Evil Stuff

Undefined, Unspecified, and Implementation-Defined Behavior

- **undefined behavior**: behavior for which standard imposes no requirements (i.e., anything could happen)
- **unspecified behavior**: behavior, for a well-formed program construct and correct data, that depends on the implementation; implementation is not required to document which behavior occurs; range of possible behaviors usually specified in standard
- **implementation-defined behavior**: behavior, for a well-formed program construct and correct data, that depends on the implementation and that each implementation documents (i.e., only know what will happen for a particular implementation)
- *always avoid undefined behavior* and *do not rely on unspecified behavior*; otherwise cannot guarantee correct behavior of program
- *try to avoid relying on implementation-defined behavior*; otherwise cannot guarantee correct behavior of program across all language implementations (i.e., code will not be portable)

Examples of Undefined Behavior

- dereferencing a null pointer; for example:

```
char* p = nullptr;  
char c = *p; // undefined behavior
```

- attempting to modify a string literal or any other const object (excluding mutable data members):

```
const int x = 0;  
const_cast<int&>(x) = 42; // undefined behavior
```

- signed integer overflow
- evaluating an expression that is not mathematically defined; for example:

```
double z = 0.0;  
double x = 1.0 / z; // undefined behavior
```

- not returning a value from a value-returning function (other than `main`)

```
int& increment(int& x) {  
    ++x;  
    // undefined behavior  
}
```

- multiple definitions of the same entity

Examples of Undefined Behavior (Continued)

- performing pointer arithmetic that yields a result before start of or after end (i.e., one past last element) of an array; for example:

```
int v[10];  
int* p = &v[0];  
--p; // undefined behavior
```

- using pointers to objects whose lifetime has ended
- left-shifting values by a negative amount; for example:

```
int i = 1;  
i << (-3); // undefined behavior
```

- shifting values by an amount greater than or equal to the number of bits in the number; for example:

```
int i = 1;  
i << 10000; // undefined behavior
```

- using an automatic variable whose value has not been initialized; for example:

```
void func() {  
    int i; ++i; // undefined behavior  
}
```

Examples of Unspecified Behavior

- order in which arguments to a function are evaluated; for example:

```
1  #include <iostream>
2
3  int count() {
4      static int c = 0;
5      return c++;
6  }
7
8  void func(int x, int y) {
9      std::cout << x << ' ' << y << '\n';
10 }
11
12 int main() {
13     func(count(), count());
14     // what values are passed to func?
15     // 0, 1; or 1, 0?
16 }
```

Examples of Implementation-Defined Behavior

- meaning of **#pragma** directive
- nesting limit for **#include** directives
- search locations for " " and <> headers
- sequence of places searched for header
- signedness of `char`
- `sizeof` built-in types other than **char**, **signed char**, **unsigned char**
- type of `size_t`, `ptrdiff_t`
- parameters to `main` function
- alignment (i.e., restrictions on the addresses at which an object of a particular type can be placed)
- result of right shift of negative value
- precise types used in various parts of C++ standard library (e.g., actual type named by `vector<T>::iterator`)
- meaning of **asm** declaration
- for more examples, see “Index of implementation-defined behavior” section in C++11 standard

Section 4.2

Best Practices, Tips, and Common Pitfalls

Use of `std::istream::eof`

- do not use `std::istream::eof` to determine if earlier input operation has failed, as this will not always work
- `eof` simply returns end-of-file (EOF) flag for stream
- EOF flag for stream can be set during *successful* input operation (when input operation takes places just before end of file)
- when stream extractors (i.e., **operator**>>) used, fields normally delimited by whitespace
- to read all data in whitespace-delimited field, must read *one character beyond* field in order to know that end of field has been reached
- if field followed immediately by EOF without any intervening whitespace characters, reading one character beyond field will cause EOF to be encountered and EOF bit for stream to be set
- in preceding case, however, EOF being set does not mean that input operation failed, only that stream data ended immediately after field that was read

Example: Incorrect Use of eof

- example of *incorrect* use of eof:

```
1  #include <iostream>
2
3  int main() {
4      while (true) {
5          int x;
6          std::cin >> x;
7          // std::cin may not be in a failed state.
8          if (std::cin.eof()) {
9              // Above input operation may have succeeded.
10             std::cout << "EOF encountered\n";
11             break;
12         }
13         std::cout << x << '\n';
14     }
15 }
```

- code incorrectly assumes that eof will only return true if preceding input operation has failed
- last field in stream will be incorrectly ignored if not followed by at least one whitespace character; for example, if input stream consists of three character sequence '1', space, '2', program will output:

```
1
EOF encountered
```

Example: Correct Use of eof

- to determine if input operation failed, simply check if stream in failed state
- if stream *already known to be in failed state* and need to determine specifically if failure due to EOF being encountered, then use eof
- example of correct use of eof:

```
1  #include <iostream>
2
3  int main() {
4      int x;
5      // Loop while std::cin not in a failed state.
6      while (std::cin >> x) {
7          std::cout << x << '\n';
8      }
9      // Now std::cin must be in a failed state.
10     // Use eof to determine the specific reason
11     // for failure.
12     if (std::cin.eof()) {
13         std::cout << "EOF encountered\n";
14     } else {
15         std::cout << "input error (excluding EOF)\n";
16     }
17 }
```

Use of `std::endl`

- `std::endl` is not some kind of string constant
- `std::endl` is stream manipulator and declared as `std::ostream& std::endl(std::ostream&)`
- inserting `endl` to stream always (regardless of operating system) equivalent to outputting single newline character '`\n`' followed by flushing stream
- flushing of stream can incur very substantial overhead; so only flush when strictly necessary

Use of `std::endl` (Continued)

- some operating systems terminate lines with single linefeed character (i.e., `'\n'`), while other operating systems use carriage-return and linefeed pair (i.e., `'\r'` plus `'\n'`)
- existence of `endl` has nothing to do with dealing with handling new lines in operating-system independent manner
- when stream opened in text mode, translation between newline characters and whatever character(s) operating system uses to terminate lines is performed automatically (both for input and output)
- above translation done for all characters input and output and has nothing to do with `endl`

Stream Extraction Failure

- for built-in types, if stream extraction fails, value of target for stream extraction depends on reason for failure
- in following example, what is value of `x` if stream extraction fails:

```
int x;
std::cin >> x;
if (!std::cin) {
    // what is value of x?
}
```

- in above example, `x` may be *uninitialized* upon stream extraction failure
- if failure due to I/O error or EOF, target of extraction is *not modified*
- if failure due to badly formatted data, target of extraction is zero
- if failure due to overflow, target of extraction is closest machine-representable value
- *common error*: incorrectly assume that target of extraction will always be initialized if extraction fails
- for class types, also dangerous to assume target of extraction always written upon failure

The abs Function

- What does the following program output when executed?

```
#include <iostream>
#include <cstdlib>

int main()
{std::cout << abs(-1.5) << '\n';}
```

- In the preceding code, the `abs` function used might be the `abs` function in the C standard library, which is declared as `int abs(int)`. In this case, the program would output a value of 1 (which is probably unexpected).
- What does the following program output when executed?

```
#include <iostream>
#include <cmath>

int main()
{std::cout << std::abs(-1.5) << '\n';}
```

- The `abs` function used in this case is an overload of the function `abs` from the C++ standard library, declared as `double abs(double)`. So, the program outputs a value of 1.5 (as expected).
- Perhaps, the best portable (correct) solution to this problem is to include `cmath` and use `std::abs` instead of `abs`.

Types of Literals

- When specifying a literal, be careful to use a literal of the correct type, as the type can often be quite important.
- For example, what value will be printed by the following code and (more importantly) why:

```
std::vector<double> values;  
values.push_back(0.5);  
values.push_back(0.5);  
// Compute the sum of the elements in the vector values.  
double sum = std::accumulate(values.begin(),  
    values.end(), 0);  
std::cout << sum << '\n';
```

- Hint: The value printed for `sum` is not 1.
- In order to determine what values will be printed, look carefully at the definition of `std::accumulate`.
- Answer: The value printed for `sum` is 0.

Testing Failure State of Streams

- consider `istream` or `ostream` object `s`
- `!s` is equivalent to `s.fail()`
- `bool(s)` is not equivalent to `s.good()`
- `s.good()` is not the same as `!s.fail()`
- do not use `good` as opposite of `fail` since this is wrong

Member Initialization Order

- data members are initialized in order in which declared
- Example:

```
1  #include <cassert>
2
3  class Widget {
4  public:
5      Widget() : y_(42), x_(y_ + 1) {assert(x_ == 43);}
6      int x_;
7      int y_;
8  };
9
10 int main() {
11     Widget w;
12 }
```

- what will above code do when run?
- in constructor, `x_` initialized before `y_`, which results in use of `y_` before its initialization
- strictly speaking, undefined behavior
- in practice, likely `x_` will simply have garbage value when body of constructor executes and assertion will fail

Global Object Initialization Order

- be careful about initialization order of global objects
- Example (program with three source files):

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

```
1 #include <vector>  
2 std::vector<int> v = {1, 2, 3, 4};
```

```
1 #include <vector>  
2 extern std::vector<int> v;  
3 std::vector<int> w = {v[0], v[1]};
```

- no guarantee that v will be constructed before w
- bad things will happen if w is constructed before v
- no guarantee about order of initialization between translation units (i.e., source files [loosely speaking])

Implement Postfix Increment/Decrement via Prefix

- implement postfix increment/decrement in terms of prefix increment/decrement
- ensures that prefix and postfix versions always consistent
- Example:

```
1  class Counter {
2  public:
3      Counter(int count = 0) : count_(count) {}
4      Counter& operator++() {
5          ++count_;
6          return *this;
7      }
8      Counter operator++(int) {
9          Counter old(*this);
10         ++(*this);
11         return old;
12     }
13     // similarly for prefix/postfix decrement
14 private:
15     int count_;
16 };
```

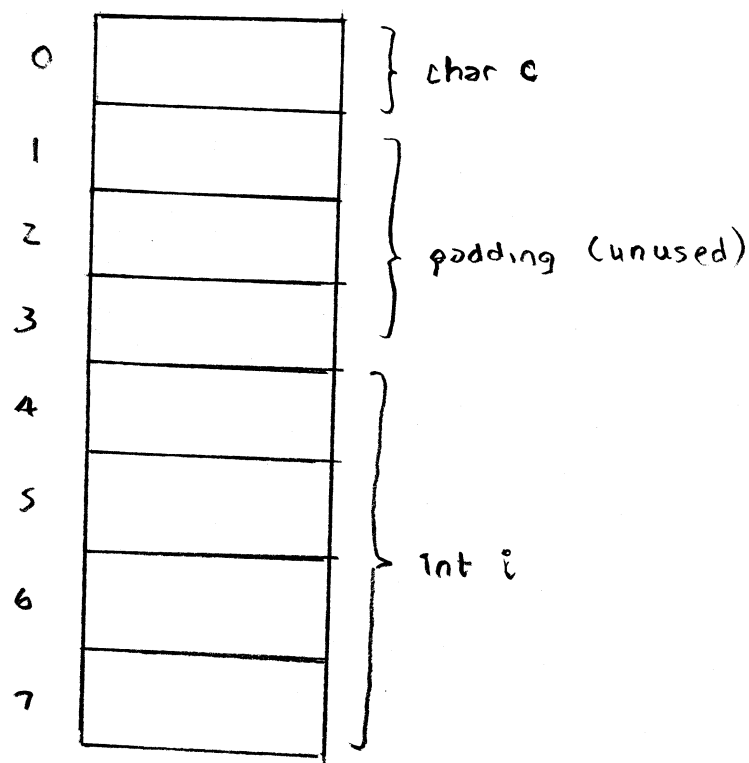
Sizeof Class Versus Sum of Member Sizes

- compilers can (and do) add padding to classes/structs
- Example:

```
1  #include <iostream>
2
3  class Widget {
4  // ...
5  private:
6      char c;
7      int i;
8  };
9
10 int main() {
11     // two numbers printed not necessarily the same
12     std::cout << sizeof(char) + sizeof(int) << ' ' <<
13         sizeof(Widget) << '\n';
14     std::cout << alignof(int) << ' ' <<
15         alignof(Widget) << '\n';
16 }
```

- many processors place alignment restrictions on data (e.g., data type of size n must be aligned to start on address that is multiple of n)
- other factors can also add to size of class/struct (e.g., virtual function table pointer)

Sizeof Class Versus Sum of Member Sizes (Continued)



- **struct** Thing { **char** c; **int** i; };
- suppose **sizeof(int)** is 4 and **alignof(int)** is 4
- implementation adds padding to structure so that **int** data member is suitably aligned (i.e., offset is multiple of 4)

Division/Modulus Operator and Negative Numbers

- for integral operands, division operator yields algebraic quotient with any fractional part discarded (i.e., round towards zero)
- if quotient a / b is representable in type of result, $(a / b) * b + a \% b$ is equal to a
- so, assuming b is not zero and no overflow, $a \% b$ equals $a - (a / b) * b$
- result of modulus operator not necessarily nonnegative
- Example:

```
1  #include <cassert>
2
3  int main() {
4      assert(5 % 3 == 2);
5      assert(5 % (-3) == 2);
6      assert((-5) % 3 == -2);
7      assert((-5) % (-3) == -2);
8  }
```

- What is wrong with the following code?

```
void func(const std::string&);  
std::string s("one");  
const char* p = "two";  
func(std::string(s) + std::string(", ") + std::string(p));  
func(std::string(p) + std::string(", ") + std::string(s));
```

- *Unnecessary temporaries!*

- Fix:

```
func(s + ", " + p);  
func(p + ", " + s);
```

- What is wrong with the following code?

```
std::vector<std::string> v;  
std::string s("one");  
v.push_back(std::string(s));  
v.push_back(std::string(s + ", two"));  
v.push_back(std::string("three"));  
v.push_back(std::string());
```

- Again, *unnecessary temporaries*.

- Fix:

```
v.push_back(s);  
v.push_back(s + ", two");  
v.emplace_back("three");  
v.emplace_back();
```

Classes Holding Multiple Resources

- What is wrong with this code?

```
class TwoResources {
public:
    TwoResources() : x_(nullptr) : y_(nullptr) {
        x_ = new X;
        y_ = new Y;
    }
    ~TwoResources() {
        delete x_;
        delete y_;
    }
private:
    X* x_;
    Y* y_;
};
```

- If an exception is thrown in a constructor, the object being constructed is deemed not to have started its lifetime and no destructor will ever be called for the object.
- So, for example, if `new Y` throws, `x_` will be leaked.
- Fix:

```
class TwoResources {
public:
    TwoResources() : x_(make_unique<X>()),
        y_(make_unique<Y>()) {}
private:
    unique_ptr<X> x_;
    unique_ptr<Y> y_;
};
```

Avoid Returning By Const Value

- What is wrong with the following code?

```
const std::string getMessage() {  
    return "Hello";  
}
```

- The const return value will *interact poorly with move semantics*, as the returned object cannot be used as the source for a move operation (since the source for a move operation must be modifiable).
- Fix:

```
std::string getMessage() {  
    return "Hello";  
}
```


Normally Avoid Using `std::move` When Returning By Value

- What is wrong with the following code?

```
std::vector<int> getVector () {  
    std::vector<int> v;  
    // calculate v  
    return std::move(v);  
}
```

- Due to the use of `std::move`, the type of the expression in the return statement does not match the function return type (i.e., `std::vector<int>` versus `std::vector<int>&&`).
- RVO/NRVO can only be applied if the type of the expression in the return statement matches the function return type.
- So, *RVO/NRVO cannot be applied* in this case.
- If the types *would not have matched* anyways (e.g., a two-element `std::tuple` and a `std::pair`), `std::move` would be reasonable to employ.

- Returning an rvalue reference to an rvalue reference parameter can potentially lead to very subtle bugs.
- Example:

```
std::string&& join(std::string&& s, const char* p) {  
    return std::move(s.append(", ").append(p));  
}  
  
std::string getMessage() {return "Hello";}  
  
void func() {  
    const string& r = join(getMessage(), " World");  
    // lifetime of temporary returned by getMessage  
    // not extended to lifetime of r since not  
    // directly bound to r  
    // r now refers to destroyed temporary  
}
```

- Fix:

```
std::string join(std::string&& s, const char* p) {  
    return std::move(s.append(", ").append(p));  
}
```

- Returning by rvalue reference should probably be avoided, except in very special circumstances (such as `std::forward` and `std::move`).

No Explicit Template Arguments to `std::make_pair`

- Never provide explicit template arguments to `std::make_pair`.
- Let `x` and `y` be objects of type `X` and `Y`, respectively.
- What is wrong with the following code?

```
std::make_pair<X, Y>(x, y)
```

- `make_pair` declared as:

```
template <class T1, class T2>  
    pair<V1, V2> make_pair(T1&& x, T2&& y);
```

where `V1` and `V2` are (except in special case) `std::decay_t<T1>` and `std::decay_t<T2>`, respectively

- If, for example, `X` and `Y` are `int`, then `make_pair` has two rvalue reference parameters which cannot bind to the lvalues `x` and `y`.
- Use `make_pair(x, y)` or sometimes `pair<X, Y>(x, y)`.

Prefer Use of `std::make_shared`

- when creating `std::shared_ptr` objects, prefer to use `std::make_shared` (as opposed to explicit use of **new** with `shared_ptr`)
- more efficient
- control block and owned object can be allocated together
- one memory allocation instead of two; better cache efficiency
- better exception safety (avoid resource leaks)

Section 4.3

Idioms

Proxy Classes

- proxy class provides modified interface to another class

Proxy Class Example

```
1  #include <iostream>
2  #include <utility>
3
4  class BoolVector;
5
6  class Proxy {
7  public:
8      ~Proxy() = default;
9      Proxy& operator=(const Proxy&) = default;
10     operator bool() const;
11     void operator=(bool b);
12 private:
13     friend class BoolVector;
14     Proxy(const Proxy&) = default;
15     Proxy(BoolVector* v, int i) : v_(v), i_(i) {}
16     BoolVector* v_;
17     int i_;
18 };
19
20 class BoolVector {
21 public:
22     BoolVector(int n) : n_(n), d_(new unsigned char[(n + 7) / 8]) {
23         std::fill_n(d_, (n + 7) / 8, 0);
24     }
25     ~BoolVector() {delete [] d_;}
26     int size() const {return n_;}
27     bool operator[](int i) const {return getElem(i);}
28     Proxy operator[](int i) {return Proxy(this, i);}
29 private:
30     friend class Proxy;
31     bool getElem(int i) const {return (d_[i / 8] >> (i % 8)) & 1;}
32     void setElem(int i, bool b) {
33         (d_[i / 8] &= ~(1 << (i % 8))) |= (b << (i % 8));
34     }
35     int n_;
36     unsigned char* d_;
37 };
38
39 inline void Proxy::operator=(bool b) {v_->setElem(i_, b);}
40 inline Proxy::operator bool() const {return v_->getElem(i_);}
```

Proxy Class Example (Continued)

```
1  #include "proxy_class_example_1.hpp"
2
3  int main() {
4      BoolVector v(16);
5      for (int i = 0; i < v.size(); ++i) {
6          v[i] = (i & 1);
7      }
8      for (int i = 0; i < v.size(); ++i) {
9          std::cout << v[i];
10     }
11     std::cout << '\n';
12     const BoolVector& cv = v;
13     for (int i = 0; i < cv.size(); ++i) {
14         std::cout << cv[i];
15     }
16     std::cout << '\n';
17 }
```


Section 4.4

C Compatibility

- Although C++ attempted to maintain compatibility with C where possible, there are numerous incompatibilities between the languages.
- Unfortunately, as C++ and C continue to evolve, the number of incompatibilities between these languages continue to grow.
- In practice, many C programs are valid C++ programs and can therefore be compiled with a C++ compiler.
- Some C programs, however, may require a significant number of changes to be valid C++.
- A few examples of incompatibilities between C++ and C are given in what follows.

Conflicts with New Keywords

```
1  #include <stdio.h>
2  #include <unistd.h>
3
4  /* Delete a file. */
5  int delete(const char* filename) { /* note function name */
6      return unlink(filename);
7  }
8
9  int main(int argc, char** argv) {
10     if (argc >= 2) {
11         if (delete(argv[1])) {
12             printf("cannot delete file\n");
13             return 1;
14         }
15     }
16     return 0;
17 }
```

- C++ introduces many new keywords.
- Some C programs might use some of these keywords as identifiers (e.g., **new**, **delete**).

Function Declarations Without Arguments

```
1  #include <stdio.h>
2
3  int plusOne(); /* no arguments specified */
4
5  int main(int argc, char** argv) {
6      printf("%d\n", plusOne(0));
7      return 0;
8  }
9
10 int plusOne(int i) {
11     return i + 1;
12 }
```

- In C, a function declaration without arguments implies that the arguments are unspecified.
- In C++, a function declaration without arguments implies that the function takes no arguments.

Implicit Return Type

```
1  #include <stdio.h>
2
3  myfunc() { /* implicit return type */
4           return 3;
5  }
6
7  int main(int argc, char **argv) {
8       int i;
9       i = myfunc();
10      printf("%d\n", i);
11      return 0;
12 }
```

- In C, if the return type of a function is not specified, it is treated as **int**.
- In C++, the return type of a function must always be explicitly specified.

More Restrictive Conversions Involving `void*`

```
1  int main(int argc, char** argv) {  
2      int i;  
3      int* ip;  
4      void* vp;  
5      ip = &i;  
6      vp = ip;  
7      ip = vp; /* problematic */  
8      return 0;  
9  }
```

- C provides an implicit conversion from `void*` to any pointer type, while C++ does not.

Scoping Rules for Nested Structs

```
1  struct outer {
2      struct inner {
3          int i;
4      };
5      int j;
6  };
7
8  struct inner a = {1}; /* inner vs. outer::inner */
9
10 int main(int argc, char** argv) {
11     return 0;
12 }
```

- C and C++ both allow nested **struct** types, but the scoping rules differ.

Part 5

Programming

Section 5.1

Good Programming Practices

Formatting, Naming, Documenting

- Be consistent with the *formatting* of the source code (e.g., indentation strategy, tabs versus spaces, spacing, brackets/parentheses).
- Avoid a formatting style that runs against common practices.
- Be consistent in the *naming conventions* used for identifiers (e.g., names of objects, functions, namespaces, types) and files.
- Avoid bizarre naming conventions that run against common practices.
- *Comment* your code. If code is well documented, it should be possible to quickly ascertain what the code is doing without any prior knowledge of the code.
- Use *meaningful names* for identifiers (e.g., names of objects, functions, types, etc.). This improves the readability of code.
- Avoid *magic literal constants*. Define a constant object and give it a meaningful name.

```
const int maxTableSize = 100;  
std::vector<TableEntry> table(maxTableSize);
```

Error Handling

- If a program requires that certain *constraints on user input* be satisfied in order to work correctly, do not assume that these constraints will be satisfied. Instead, always check them.
- Always handle errors *gracefully*.
- Provide *useful* error messages.
- Always *check return codes*. Even if the operation/function theoretically cannot fail (under the assumption of bug-free code), in practice it may fail due to a bug.
- If an operation is performed that can fail, check the *status of the operation* to ensure that it did not fail (even if you think that it should not fail). For example, check for error conditions on streams.
- If a function can fail, always check its *return value*.

- Do not *unnecessarily complicate* code. Use the simplest solution that will meet the needs of the problem at hand.
- Do not impose *bogus limitations*. If a more general case can be handled without complicating the code and this more general case is likely to be helpful to handle, then handle this case.
- Do not *unnecessarily optimize* code. Highly optimized code is often much less readable. Also, highly optimized code is often more difficult to write correctly (i.e., without bugs). Do not write grossly inefficient code that is obviously going to cause performance problems, but do not optimize things beyond avoiding gross inefficiencies that you know will cause performance problems.

Code Duplication

- Avoid *duplication* of code. If similar code is needed in more than one place, put the code in a function. Also, utilize templates to avoid code duplication.
- The avoidance of code duplication has many advantages.
 - ① It simplifies code understanding. (Understand once, instead of n times.)
 - ② It simplifies testing. (Test once, instead of n times.)
 - ③ It simplifies debugging. (Fix bugs in one place, instead of n places.)
 - ④ It simplifies code maintenance. (Change code in one place, instead of n places.)
- Make good use of the available *libraries*. Do not reinvent the wheel. If a library provides code with the needed functionality, use the code in the library.

- Avoid *multiple returns paths* (i.e., multiple points of exit) in functions when they serve to complicate (rather than simplify) code structure.
- Avoid the use of *global objects*. For example, use static data members instead of global objects. In well designed code, global objects are rarely needed.
- Ensure that the code is *const correct*.
- If an object does not need to change, make it const. This improves the readability of code. This also helps to ensure const correctness of code.
- Avoid bringing many unknown identifiers into scope. For example, avoid constructs like:

```
using namespace std;
```

Only bring identifiers into scope if you need them.

- Do not rely on *undefined/unspecified/implementation-defined behavior*. Do not rely on any behavior that is not promised by the language. Do not rely on undocumented features of libraries. That is, do not write code in a way that it may only work on certain computing platforms or when the moon is full.
- Enable *compiler warning messages*. Pay attention to warning messages issued by the compiler.
- Learn how to use a *source-level debugger*. There will be times when you will absolutely need it.
- Be careful to avoid using references, pointers, iterators that do not reference valid data. Always be clear about which operations invalidate references, pointers, and iterators.

Testing: Preconditions and Postconditions

- **precondition**: condition that must be true before function is called
- for example, precondition for function that computes square root of x :
 $x \geq 0$
- **postcondition**: condition that must be true after function is called
- for example, postcondition for function that removes entry from table of size n : new size of table $n - 1$
- whenever feasible, check for violations of preconditions and postconditions for functions
- if precondition or postcondition is violated, terminate program immediately in order to help in localizing bug (e.g., by calling `std::abort` or `std::terminate`)

Testing

- The single most important thing when writing code is that it does the job it was intended to do *correctly*. That is, there should not be any bugs.
- *Test* your code. If you do not spend as much time testing your code as you do writing it, you are likely not doing enough testing.
- Tests should exercise as much of the code as possible (i.e., provide good *code coverage*).
- Design and structure your code so that it is easy to test. In other words, testing should be considered *during design*.
- Your code will have bugs. Design your code so that it will help you to isolate bugs. Use *assertions*. Use *preconditions* and *postconditions*.
- Design your code so that is modular and can be written and tested *in pieces*. The first testing of the software should never be testing the entire software as a whole.
- Often in order to adequately test code, one has to write separate *specialized test code*.

Code Examples

- subscripting operator for 1-D array class:

```
template <class T>
const T& Array_1<T>::operator[(int i)] const {
    // Precondition: index is in allowable range
    assert(i >= 0 && i < data_.size());
    return data_[i];
}
```

- function taking pointer parameter:

```
int stringLength(const char* ptr) {
    // Precondition: pointer is not null
    assert(ptr != 0);
    // Code to compute and return string length.
    // ...
}
```

- function that modifies highly complicated data structure:

```
void modifyDataStructure(Type& dataStructure) {
    // Precondition: data structure is in valid state
    assert(isDataStructureValid(dataStructure));
    // Complicated code to update data structure.
    // ...
    // Postcondition: data structure is in valid state
    assert(isDataStructureValid(dataStructure));
}
```

Section 5.2

Finite-Precision Arithmetic

Code Example

- What do each of the following functions output when executed?

```
void func1() {
    double x = 0.1;
    double y = 0.3;
    double z = 0.4;
    if (x + y == z) {
        std::cout << "true\n";
    } else {
        std::cout << "false\n";
    }
}
```

```
void func2() {
    double x = 1e50;
    double y = -1e50;
    double z = 1.0;
    if (x + y + z == z + y + x) {
        std::cout << "true\n";
    } else {
        std::cout << "false\n";
    }
}
```

```
void func3() {
    for (double x = 0.0; x != 1.0; x += 0.1) {
        std::cout << "hello\n";
    }
}
```

Number Representations Using Different Radixes

- Note: All numbers are base 10, unless explicitly indicated otherwise.
- What is the representation of $\frac{1}{3}$ in base 3?
 $\frac{1}{3} = 0.\overline{3} = 0.1_3$
- What is the representation of $\frac{1}{10}$ in base 2?
 $\frac{1}{10} = 0.1 = 0.0\overline{0011}_2$
- A number may have a representation with a finite number of non-zero digits in one particular number base but not in another.
- Therefore, when a value must be represented with a limited number of significant digits, the number base matters (i.e., affects the approximation error).
- For example, in base 2, $\frac{1}{10}$ cannot be represented exactly using only a finite number of significant digits.
 $0.00011_2 = 0.09375$
 $0.000110011_2 = 0.099609375$
...

Finite-Precision Number Representations

- finite-precision number representation only capable of representing small fixed number of digits
- due to limited number of digits, many values cannot be represented exactly
- in cases that desired value cannot be represented exactly, choose nearest representable value (i.e., round to nearest representable value)
- finite-precision representations can suffer from error due to roundoff, underflow, and overflow
- two general classes of finite-precision representations:
 - 1 fixed-point representations
 - 2 floating-point representations

Fixed-Point Number Representations

- **fixed-point representation**: radix point remains fixed at same position in number
- if radix point fixed to right of least significant digit position, integer format results

Integer Format $a_{n-1} a_{n-2} \cdots a_1 a_0 \blacksquare$

- if radix point fixed to left of most significant digit position, purely fractional format results

Fractional Format $\blacksquare a_{n-1} a_{n-2} \cdots a_1 a_0$

- fixed-point representations *quite limited in range* of values that can be represented
- numbers that vary greatly in magnitude cannot be represented easily using fixed-point representations
- one solution to range problem would be for programmer to maintain scaling factor for each fixed-point number, but this is clumsy and error prone

Floating-Point Number Representations

- **floating-point representation**: radix point is not fixed at particular position within number; instead radix point allowed to move and scaling factor automatically maintained to track position of radix point
- in general, floating-point value represents number x of form

$$x = sr^e,$$

- s is signed integer with fixed number of digits, and called **significand**
- e is signed integer with fixed number of digits, and called **exponent**
- r is integer satisfying $r \geq 2$, and called **radix**
- in practice, r typically 2
- for fixed r , representation of particular x not unique if no constraints placed on s and e (e.g., $5 \cdot 10^0 = 0.5 \cdot 10^1 = 0.05 \cdot 10^2$)

Floating-Point Number Representations (Continued)

- to maximize number of significant digits in significand, s and e usually chosen such that first nonzero digit in significand is to immediate left of radix point (i.e., $1 \leq |s| < r$); number in this form called **normalized**; otherwise called **denormalized**
- other definitions of normalized/denormalized sometimes used but above one consistent with IEEE 754 standard
- Example:

$$\begin{aligned}0.75 &= 0.11_2 = 1.1_2 \cdot 2^{-1} \\1.25 &= 1.01_2 = 1.01_2 \cdot 2^0 \\-0.5 &= -0.1_2 = -1.0_2 \cdot 2^{-1}\end{aligned}$$

IEEE 754 Standard (IEEE Std. 754-1985)

- most widely used standard for (binary) floating-point arithmetic
- specifies four floating-point formats: single, double, single extended, and double extended
- single and double formats called basic formats
- radix 2
- three integer parameters determine values representable in given format:
 - number p of significand bits (i.e., precision)
 - maximum exponent E_{\max}
 - minimum exponent E_{\min}
- parameters for four formats are as follows:

Parameter	Single	Single Extended	Double	Double Extended
p	24	≥ 32	53	≥ 64
E_{\max}	127	> 1023	1023	≥ 16383
E_{\min}	-126	≤ -1022	-1022	≤ -16382
Exponent bias	127	unspecified	1023	unspecified

- with each format, numbers of following form can be represented

$$(-1)^s 2^E (b_0.b_1b_2 \cdots b_{p-1})$$

where $s \in \{0, 1\}$, E is integer satisfying $E_{\min} \leq E \leq E_{\max}$, and $b_i \in \{0, 1\}$

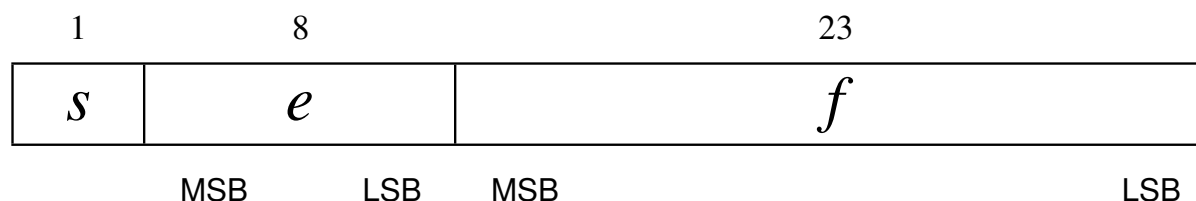
- in addition, can represent four special values: $+\infty$, $-\infty$, signaling NaN, and quiet NaN
- NaNs produced by:
 - operations with at least one NaN operand
 - operations yielding indeterminate forms, such as $0/0$, $(\pm\infty)/(\pm\infty)$, $0 \cdot (\pm\infty)$, $(\pm\infty) \cdot 0$, $(+\infty) + (-\infty)$, and $(-\infty) + (\infty)$
 - real operations that yield complex results, such as square root of negative number, logarithm of negative number, inverse sine/cosine of number that lies outside $[-1, 1]$

IEEE 754 Basic Formats

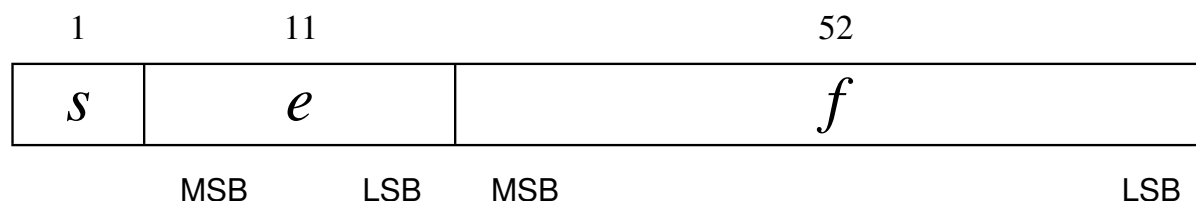
- always represent number in normalized form whenever possible; in such cases, $b_0 = 1$ and b_0 need not be stored explicitly as part of significand
- bit patterns with reserved exponent values (i.e., exponent values that lie outside the range $E_{\min} \leq E \leq E_{\max}$) used to represent ± 0 , $\pm \infty$, denormalized numbers, and NaNs
- each of (basic) formats consist of three fields:
 - a sign bit, s
 - a biased exponent, $e = E + \text{bias}$
 - a fraction, $f = .b_1b_2 \cdots b_{p-1}$
- only difference between formats is size of biased exponent and fraction fields
- value represented by basic format number related to its sign, exponent, and fraction field, but relationship is complicated by the presence of zeros, infinities, and NaNs
- “strange” combination of biased and sign-magnitude formats used to encode floating-point value chosen so that nonnegative floating-point values ordered in same way as integers, allowing integer comparison to compare floating-point numbers

IEEE 754 Basic Formats (Continued)

- single format:



- double format:



- summary of encodings:

Case	Exponent	Fraction	Value
Normal	$E_{\min} \leq E \leq E_{\max}$	—	$(-1)^s 2^E (1 + f)$
Denormal	$E = E_{\min} - 1$	$f \neq 0$	$(-1)^s 2^{E_{\min}} f$
Zero	$E = E_{\min} - 1$	$f = 0$	$(-1)^s 0$
Infinity	$E = E_{\max} + 1$	$f = 0$	$(-1)^s \infty$
NaN	$E = E_{\max} + 1$	$f \neq 0$	NaN

Finite-Precision Arithmetic

- Understand the impact of using finite-precision arithmetic.
- Do not make invalid assumptions about the set of values that can be represented by a particular fixed-point or floating-point type.
- Integer arithmetic can *overflow*. Be careful to avoid overflow.
- Floating-point arithmetic can *overflow and underflow*.
- Perhaps, more importantly, however, floating-point arithmetic has *roundoff error*. If you are not deeply troubled by the presence of roundoff error, you should be as it can cause major problems in many situations.

- 1 D. Goldberg. [What every computer scientist should know about floating-point arithmetic.](#)
ACM Computing Surveys, 23(1):5–48, Mar. 1991
- 2 IEEE Std. 754-1985 — IEEE standard for binary floating-point arithmetic, 1985
- 3 IEEE Std. 754-2008 — IEEE standard for floating-point arithmetic, 2008

- ① John Farrier, Demystifying Floating Point, CppCon, Bellevue, WA, USA, Sept. 24, 2015.

Section 5.3

Documentation for Software Development

Documentation for Software Development

- documentation plays essential role in software development process
- many benefits to formalizing in writing various aspects of software at different points in development process
- consider two types of documents:
 - ① software requirements specification
 - ② software design description
- software requirements specification (SRS): describes what software should do (from external viewpoint)
- software design description (SDD): describes how software works internally

Software Requirements Specification (SRS)

- establishes agreement between consumer and contractors on what software is expected to do as well as what it is not expected to do
- can be thought of as contract between customer and contractor
- functionality: what does software do? (what problem does it solve?)
- external interfaces: how does software interact with external agents, such as humans, hardware, and software (e.g., command-line interface, graphical user interface, application program interface)
- performance: speed, availability, response time, recovery time of various functions
- attributes: considerations regarding reliability, availability, maintainability, portability, security
- design constraints imposed on implementation: implementation language, resource limits, operating environments
- assumptions upon which requirements are based

- distinguish classes of requirements:
 - essential: software will be unacceptable unless requirement met
 - conditional: would enhance software if requirement met, but not unacceptable if requirement not met
 - optional: class of functionality that may or may not be worthwhile
- should not leave details of software requirements to be determined
- only focus on what the software needs to do, not how done (i.e., should not describe any design or implementation details)
- typical use cases
- constraints imposed on software:
 - time constraints
 - memory constraints
- software limitations:
 - restrictions on input data
 - allowable ranges for parameters of methods
 - dependencies on other software (e.g., other programs needed to function)

External Interfaces

- external interfaces: how software interacts with external agents, such as humans, hardware, and software
- command line interface (CLI) (for program)
 - options (e.g., required versus optional, default settings)
 - standard input, output, error
 - exit status
- graphical user interface (GUI) (for program)
 - window layout
 - user interaction (e.g., mouse/keyboard actions)
- application program interface (API) (for library)
 - constants
 - types, classes/methods
 - functions
 - namespaces
- format of all data used by software

Benefits of SRS

- establishes basis for agreement between customer and contractors
- reduces development effort by thoroughly considering all requirements before starting design
- provides basis for estimating costs and schedules
- provides baseline for validation and verification
- facilitates transfer of software product to new users or machines
- serves as basis for enhancement

SRS Example: Sorting Program

- single program that performs sorting
- given records as input, program sorts records and outputs records in sorted order
- record data format (for input and output):
 - records delimited by single newline character
 - each record consists of one or more fields, separated by one or more whitespace characters
- restrictions/constraints:
 - may assume sufficient memory to buffer all records
 - software must work without any modification to source code on any platform with C++ compiler compliant with C++11 standard
- records read from standard input
- sorted records written to standard output
- any error/warning messages written to standard error
- sorts records using n th field in record as key
- can sort in ascending or descending order
- sort key may be numeric or string

SRS Example: Sorting Program (Continued)

- command line interface:

```
sort [-r] [-k $n] [-n]
```

Option	Description
-k \$n	Sort using <i>n</i> th field in record; if not specified, <i>n</i> defaults to 1.
-n	Treat key as real number (instead of string) for sorting purposes; if not specified, key treated as string.
-r	Sort in descending (instead of ascending) order; if not specified, defaults to ascending order.

- give examples illustrating expected use cases

Software Design Description (SDD)

- high-level design: overview of entire system, identifying all its components at some level of abstraction (i.e., overall software architecture)
- detailed design (a.k.a. low-level design): full details of system and its components (e.g., types, functions, APIs, pseudocode, etc.)
- describes high-level and detailed design of software
- some context regarding functionality provided by software
- how design is recursively structured into constituent parts and role of those parts
- types and interfaces (e.g., classes and public members)
- data structures used to represent information to be processed
- internal interfaces (and external interfaces not described in SRS)
- interaction amongst entities
- algorithms

SDD (Continued)

- describe overall structure of software
- carefully consider choice of data structures used to represent information being processed, as choice will almost always have performance implications
- specify any data formats used internally by software
- provide pseudocode for key parts of software
- state any potentially limiting assumptions made

Benefits of SDD

- encourages better planning by forcing design ideas to be more carefully considered and organized
- allows greater scrutiny of design
- captures important design decisions, such as rationale for particular design choices
- allows newcomers to development team to become acquainted with software more easily
- provides point of reference to be used throughout project
- promotes reuse of code (since well documented code more likely to be reused)
- facilitates better software testing (since certain types of testing benefit from understanding of software design)

SDD Example: Sorting Program

- `Key` alias for type that represents sort key (alias for `std::string`)
- `Compare` functor class for comparing `Key` objects
- `Dataset` class represents collection of all records
- specify all class interfaces (i.e., public members)
- `Dataset` class provides:
 - constructor that creates dataset by reading all records from input stream
 - function to output all records in sorted order to output stream
- `Dataset` class to use `std::multimap<Key, std::string, Compare>`
- allows n records to be sorted in $O(n \log n)$ time [n insertions, each requiring $O(\log n)$ time]
- handling n records requires $O(n)$ memory
- only uses C++ standard library

Requirements/Design Document for Degree Project

- document is combination of SRS and SDD with some added information about testing strategies
- briefly introduce problem being addressed by software
- describe each program and library to be developed
- identify parts of any external software (e.g., programs or libraries) that will be used
- describe user interface (e.g., CLI, GUI) for each program
- fully specify all data formats used
- describe overall structure of each program and library
- identify all key data structures and algorithms to be used
- provide pseudocode for key parts of the software
- state any potentially limiting assumptions made by software
- indicate how programs and library code will be tested
- offer any other information that may be helpful (since above list is not exhaustive)
- provide sufficient detail for other people to understand how software is to be structured and how it will be implemented and tested

- 1 IEEE Std. 1016-2009 — IEEE standard for information technology — systems design — software design descriptions, July 2009.
- 2 IEEE Std. 830-1998 — IEEE recommended practice for software requirements specifications, Oct. 1998.

Part 6

Additional Learning Resources

Limits of Knowledge

- Know what you do not know.
- Ask questions when you are uncertain about something and be sure that the person whom you ask is knowledgeable enough to give a *correct* answer.
- Know what information resources can be *trusted*.
- Learn to use reference materials effectively (e.g., documentation on libraries, standards).

- Some good references on various topics related to the C++ programming language, C++ standard library, and other C++ libraries (such as Boost) are listed on the slides that follow.
- Any information on C++ (e.g., books, tutorials, videos, seminars) from the following individuals (who are held in very high regard by the C++ community) is highly recommended:
 - Bjarne Stroustrup (the creator of C++)
 - Scott Meyers
 - Herb Sutter (Convener of ISO C++ standards committee for over 10 years)
 - Andrei Alexandrescu
 - Stephan Lavavej

- 1 ISO/IEC 14882:2011 — information technology — programming languages — C++, Sept. 2011.

This is the definitive specification of the C++ language and standard library. This is an essential reference for any advanced programmer.

- 2 B. Stroustrup. *The C++ Programming Language*. Addison Wesley, 4th edition, 2013.

This is the classic book on the C++ programming language and standard library, written by the creator of the language. This is one of the best references for first learning C++. Excellent

- 3 Standard C++ Foundation web site. <http://www.isocpp.org>, 2014.

This is the web site of a non-profit organization whose purpose is to support the C++ software development community and promote the understanding and use of modern standard C++ on all compilers and platforms. This is an absolutely outstanding source of information on C++. Excellent

- ④ S. Meyers. *Effective Modern C++: 42 Specific Ways to Improve Your Use of C++11 and C++14.*

O'Reilly Media, Cambridge, MA, USA, 2015.

This book covers a list of 42 topics on how to better utilize the C++ language.

Excellent

- ⑤ S. Meyers. *Effective C++: 50 Specific Ways to Improve Your Programs and Designs.*

Addison Wesley, Menlo Park, California, 1992.

This book covers a list of 50 topics on how to better utilize the C++ language.

Excellent

- ⑥ S. Meyers. *More Effective C++: 35 New Ways to Improve Your Programs and Designs.*

Addison Wesley, Menlo Park, California, 1996.

This book covers a list of 35 topics on how to better utilize the C++ language. It builds on Meyers' earlier "Effective C++" book.

Excellent

- 7 S. Meyers. *Effective STL: 50 Specific Ways to Improve Your Use of the Standard Template Library*. Addison Wesley, 2001.

This book covers a list of 50 topics on how to better utilize the Standard Template Library (STL), an essential component of the C++ standard library. Excellent

- 8 N. M. Josuttis. *The C++ Standard Library: A Tutorial and Reference*. Addison Wesley, Upper Saddle River, NJ, USA, 2nd edition, 2012.
This is a very comprehensive book on the C++ standard library. This is arguably the best reference on the standard library (other than the C++ standard). Excellent
- 9 D. Vandevor and N. M. Josuttis. *C++ Templates: The Complete Guide*. Addison Wesley, 2002.
This is a very comprehensive book on template programming in C++. It is arguably one of the best books on templates in C++. Excellent
- 10 A. Williams. *C++ Concurrency in Action*. Manning Publications, Shelter Island, NY, USA, 2012.
This is a fairly comprehensive book on concurrency and multithreaded programming in C++. It is arguably the best book available for those who want to learn how to write multithreaded code using C++. Excellent

- 11 H. Sutter. *Exceptional C++: 47 Engineering Puzzles, Programming Problems, and Solutions.*

Addison Wesley, 1999.

This book covers topics including (but not limited to): proper resource management, exception safety, RAI, and good class design. Excellent

- 12 H. Sutter. *More Exceptional C++: 40 New Engineering Puzzles, Programming Problems, and Solutions.*

Addison Wesley, 2001.

This book covers topics including (but not limited to): exception safety, effective object-oriented programming, and correct use of STL. Excellent

- 13 H. Sutter. *Exceptional C++ Style: 40 New Engineering Puzzles, Programming Problems, and Solutions.*

Addison Wesley, 2004.

This book covers topics including (but not limited to): generic programming, optimization, resource management, and how to write modular code. Excellent

- 14 H. Sutter and A. Alexandrescu. *C++ Coding Standards: 101 Rules, Guidelines, and Best Practices*. Addison Wesley, 2004.

This book presents 101 best practices, idioms, and common pitfalls in C++ in order to allow the reader to become a more effective C++ programmer. Excellent

- 15 A. Langer and K. Kreft. *Standard C++ IOStreams and Locales*. Addison Wesley, 2000.

This book provides a very detailed look at C++ I/O streams and locales.

Said-To-Be Excellent

- 16 V. A. Punathambekar. *How to interpret complex C/C++ declarations*. <http://www.codeproject.com/Articles/7042/How-to-interpret-complex-C-C-declarations>, 2004.

This is a detailed tutorial on how to interpret complex C/C++ type declarations.

This tutorial explains how type declarations are parsed in the language, which is essential for all programmers to understand clearly. Excellent

Other C++ References I

- 1 S. B. Lippman, J. Lajoie, and B. E. Moo. *C++ Primer*. Addison Wesley, Upper Saddle River, NJ, USA, 4th edition, 2005.
- 2 A. Koenig and B. E. Moo. *Accelerated C++: Practical Programming by Example*. Addison Wesley, Upper Saddle River, NJ, USA, 2000.
- 3 B. Eckel. *Thinking in C++—Volume 1: Introduction to Standard C++*. Prentice Hall, 2nd edition, 2000.
- 4 B. Eckel and C. Allison. *Thinking in C++—Volume 2: Practical Programming*. Prentice Hall, 1st edition, 2003.
- 5 B. Stroustrup. *Programming: Principles and Practice Using C++*. Addison Wesley, Upper Saddle River, NJ, USA, 2009.
An introduction to programming using C++ by the creator of the language.

Other C++ References II

- 6 A. Alexandrescu. *Modern C++ Design*.
Addison Wesley, Upper Saddle River, NJ, USA, 2001.
- 7 D. Abrahams and A. Gurtovoy. *C++ Template Metaprogramming: Concepts, Tools, and Techniques from Boost and Beyond*.
Addison Wesley, Boston, MA, USA, 2004.
- 8 D. D. Gennaro. *Advanced C++ Metaprogramming*.
CreateSpace Independent Publishing Platform, 2011.
- 9 Boost web site. <http://www.boost.org>, 2014.
The web site for the Boost C++ libraries.
- 10 B. Karlsson. *Beyond the C++ Standard Library: An Introduction to Boost*.
Addison Wesley, Upper Saddle River, NJ, USA, 2005.
An introduction to (some parts of) the Boost library.

- 11 B. Schaling. *The Boost C++ Libraries*. XML Press, 2nd edition, 2014.
An introduction to the Boost library. Online version at <http://theboostcpplibraries.com>.
- 12 M. Kilpelainen. *Overload resolution — selecting the function*. *Overload*, 66:22–25, Apr. 2005.
Available online at <http://accu.org/index.php/journals/268>.

Yet More C++ References I

- 1 Herb Sutter's Web Site: <http://herbsutter.com>
- 2 Herb Sutter's Guru of the Week: <http://www.gotw.ca/gotw/>
- 3 Bjarne Stroustrup's Web Site: <http://www.stroustrup.com>
- 4 ISO C++ Working Group web site: <http://www.open-std.org/jtc1/sc22/wg21/>
- 5 C++ FAQ: <http://www.parashift.com/c++-faq/>
- 6 Newsgroup comp.lang.c++.moderated: <https://groups.google.com/forum/#!forum/comp.lang.c++.moderated>
- 7 <http://en.cppreference.com>
- 8 <http://www.cplusplus.com>
- 9 Stackoverflow: <http://stackoverflow.com>
- 10 Cpp Reddit (C++ discussions, articles, and news): <https://www.reddit.com/r/cpp>

Yet More C++ References II

- ⑪ **Cplusplus Reddit (C++ questions, answers, and discussion):** <https://www.reddit.com/r/cplusplus>
- ⑫ **ACCU Overload Journal:** <http://accu.org/index.php/journals/c78/>
- ⑬ **The C++ Source:** <http://www.artima.com/cppsource>

- 1 Scott Schurr. `constexpr`: Introduction, CppCon, Bellevue, WA, USA, Sept 19–25, 2015.
- 2 Scott Schurr. `constexpr`: Applications, CppCon, Bellevue, WA, USA, Sept 19–25, 2015.

C++ Programming Competitions

1 Google Code Jam

<https://code.google.com/codejam/>

2 Topcoder

<https://www.topcoder.com/>

3 IEEEExtreme 24-Hour Programming Competition

<http://www.ieee.org/xtreme>

4 ACM International Collegiate Programming Contest (ICPC)

<http://icpcnews.com/>

5 CodeChef

<https://www.codechef.com/>

- Use as many information resources as you can to learn as much as you can about C++.
- Read books, articles, and other documents.
- Watch videos.
- Attend lectures and seminars.
- Participate in programming competitions.
- But most importantly:

Write code!

Write lots and lots and lots of code!

- The only way to truly learn a programming language well is to use it heavily (i.e., write lots of code using the language).