

Electronic Edition

This file is part of the electronic edition of *The Unicode Standard, Version 5.0*, provided for online access, content searching, and accessibility. It may not be printed. Bookmarks linking to specific chapters or sections of the whole Unicode Standard are available at

<http://www.unicode.org/versions/Unicode5.0.0/bookmarks.html>

Purchasing the Book

For convenient access to the full text of the standard as a useful reference book, we recommend purchasing the printed version. The book is available from the Unicode Consortium, the publisher, and booksellers. Purchase of the standard in book format contributes to the ongoing work of the Unicode Consortium. Details about the book publication and ordering information may be found at

<http://www.unicode.org/book/aboutbook.html>

Joining Unicode

You or your organization may benefit by joining the Unicode Consortium: for more information, see [Joining the Unicode Consortium](http://www.unicode.org/consortium/join.html) at

<http://www.unicode.org/consortium/join.html>

This PDF file is an excerpt from *The Unicode Standard, Version 5.0*, issued by the Unicode Consortium and published by Addison-Wesley. The material has been modified slightly for this electronic edition, however, the PDF files have not been modified to reflect the corrections found on the Updates and Errata page (<http://www.unicode.org/errata/>). For information on more recent versions of the standard, see <http://www.unicode.org/versions/enumeratedversions.html>.

Many of the designations used by manufacturers and sellers to distinguish their products are claimed as trademarks. Where those designations appear in this book, and the publisher was aware of a trademark claim, the designations have been printed with initial capital letters or in all capitals.

The Unicode® Consortium is a registered trademark, and Unicode™ is a trademark of Unicode, Inc. The Unicode logo is a trademark of Unicode, Inc., and may be registered in some jurisdictions.

The authors and publisher have taken care in the preparation of this book, but make no expressed or implied warranty of any kind and assume 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.

The *Unicode Character Database* and other files are provided as-is by Unicode®, Inc. No claims are made as to fitness for any particular purpose. No warranties of any kind are expressed or implied. The recipient agrees to determine applicability of information provided. *Dai Kan-Wa Jiten*, used as the source of reference Kanji codes, was written by Tetsuji Morohashi and published by Taishukan Shoten.

Cover and CD-ROM label design: Steve Mehallo, www.mehallo.com

The publisher offers excellent discounts on this book when ordered in quantity for bulk purchases or special sales, which may include electronic versions and/or custom covers and content particular to your business, training goals, marketing focus, and branding interests. For more information, please contact U.S. Corporate and Government Sales, (800) 382-3419, corpsales@pearsoned.com. For sales outside the United States please contact International Sales, international@pearsoned.com

Visit us on the Web: www.awprofessional.com

Library of Congress Cataloging-in-Publication Data

The Unicode Standard / the Unicode Consortium ; edited by Julie D. Allen ... [et al.]. — Version 5.0.
p. cm.

Includes bibliographical references and index.

ISBN 0-321-48091-0 (hardcover : alk. paper)

1. Unicode (Computer character set) I. Allen, Julie D.

II. Unicode Consortium.

QA268.U545 2007

005.7'22—dc22

2006023526

Copyright © 1991–2007 Unicode, Inc.

All rights reserved. Printed in the United States of America. This publication is protected by copyright, and permission must be obtained from the publisher prior to any prohibited reproduction, storage in a retrieval system, or transmission in any form or by any means, electronic, mechanical, photocopying, recording, or likewise. For information regarding permissions, write to Pearson Education, Inc., Rights and Contracts Department, 75 Arlington Street, Suite 300, Boston, MA 02116. Fax: (617) 848-7047

ISBN 0-321-48091-0

Text printed in the United States on recycled paper at Courier in Westford, Massachusetts.

First printing, October 2006

Appendix B

Unicode Publications and Resources

This appendix provides information about the Unicode Consortium and its activities, particularly regarding publications other than the Unicode Standard. The Unicode Consortium publishes a number of technical standards and technical reports, and the current list of those, with abstracts of their content, is included here for convenient reference.

The Unicode Web site also has many useful online resources. *Section B.6, Other Unicode Online Resources*, provides a guide to the kinds of information available online.

B.1 The Unicode Consortium

The Unicode Consortium was incorporated in January 1991, under the name Unicode, Inc., to promote the Unicode Standard as an international encoding system for information interchange, to aid in its implementation, and to maintain quality control over future revisions.

To further these goals, the Unicode Consortium cooperates with the Joint Technical Committee 1 of the International Organization for Standardization and the International Electrotechnical Commission (ISO/IEC JTC1). It holds a Class C liaison membership with ISO/IEC JTC1/SC2; it participates in the work of both JTC1/SC2/WG2 (the technical working group for the subcommittee within JTC1 responsible for character set encoding) and the Ideographic Rapporteur Group (IRG) of WG2. The Consortium is a member company of the InterNational Committee for Information Technology Standards, Technical Committee L2 (INCITS/L2), an accredited U.S. standards organization. Many members of the Unicode Consortium have representatives in many countries who also work with other national standards bodies. In addition, a number of organizations are Liaison Members of the Consortium. For a list, see “Unicode Consortium Liaison Members” on page xlix.

Membership in the Unicode Consortium is open to organizations and individuals anywhere in the world who support the Unicode Standard and who would like to assist in its extension and widespread implementation. Full, Institutional, Supporting, and Associate Members represent a broad spectrum of corporations and organizations in the computer and information processing industry. For a list, see “Unicode Consortium Members” on page xlvi. The Consortium is supported financially solely through membership dues.

The Unicode Technical Committee

The Unicode Technical Committee (UTC) is the working group within the Consortium responsible for the creation, maintenance, and quality of the Unicode Standard. The UTC follows an open process in developing the Unicode Standard and its other technical publications. It coordinates and reviews all technical input to these documents and decides their contents. For more information on the UTC and the process by which the Unicode Standard and the other technical publications are developed, see:

<http://www.unicode.org/consortium/utc.html>

Other Activities

Going beyond developing technical standards, the Unicode Consortium acts as registration authority for the registration of script identifiers under ISO 15924, and it has a technical committee dedicated to the maintenance of the Common Locale Data Repository (CLDR). The repository contains a large and rapidly growing body of data used in the locale definition for software internationalization. For further information about these and other activities of the Unicode Consortium, visit:

<http://www.unicode.org>

B.2 Unicode Publications

In addition to the Unicode Standard, the Unicode Consortium publishes Unicode Technical Standards and Unicode Technical Reports. These materials are published as electronic documents only and, unlike Unicode Standard Annexes, do not form part of the Unicode Standard.

A *Unicode Standard Annex* (UAX) forms an integral part of the Unicode Standard. The Unicode Standard may require conformance to normative content in a Unicode Standard Annex, if so specified in the Conformance chapter of that version of the Unicode Standard. The version number of a UAX document corresponds to the version number of the Unicode Standard at the last point that the UAX document was updated. The Unicode Standard Annexes are printed in the back of this book, following the indices.

A *Unicode Technical Standard* (UTS) is a separate specification with its own conformance requirements. A UTS may include a requirement for an implementation of the UTS to also conform to a specific, base level of the Unicode Standard, but conformance to the Unicode Standard as such does not require conformance to any UTS.

A *Unicode Technical Report* (UTR) contains informative material. Unicode Technical Reports do not contain conformance requirements of their own, nor does conformance to the Unicode Standard require conformance to any specifications contained in any of the UTRs. Other specifications, however, are free to cite the material in UTRs and to make any level of conformance requirements within their own context.

In the past, some normative material was published as Unicode Technical Reports. Currently, however, such material is published either as a Unicode Technical Standard or a Unicode Standard Annex.

The Unicode Web site is the source for the most current version of all three categories of technical reports:

<http://www.unicode.org/reports/>

The following sections provide lists of abstracts for current Unicode Technical Standards and Unicode Technical Reports. They are listed numerically within each category. There are gaps in the numerical sequence because some of the reports have been superseded or have been incorporated into the text of the standard.

B.3 Unicode Technical Standards

UTS #6: A Standard Compression Scheme for Unicode

This report presents the specifications of a compression scheme for Unicode and sample implementation.

UTS #10: Unicode Collation Algorithm

This report provides the specification of the Unicode Collation Algorithm, which provides a specification for how to compare two Unicode strings while remaining conformant to the requirements of the Unicode Standard.

UTS #18: Unicode Regular Expression Guidelines

This document describes guidelines for how to adapt regular expression engines for use with the Unicode Standard.

UTS #22: Character Mapping Markup Language (CharMapML)

This document specifies an XML format for the interchange of mapping data for character encodings. It provides a complete description for such mappings in terms of a defined mapping to and from Unicode code points, and a description of alias tables for the interchange of mapping table names.

UTS #35: Locale Data Markup Language (LDML)

This document describes an XML format (*vocabulary*) for the exchange of structured locale data.

UTS #37: Ideographic Variation Database

This document describes the organization of the Ideographic Variation Database and the procedure to add sequences to that database.

UTS #39: Unicode Security Mechanisms

Because Unicode contains such a large number of characters and incorporates the varied writing systems of the world, incorrect usage can expose programs or systems to possible security attacks. This report specifies mechanisms that can be used in detecting possible security problems.

B.4 Unicode Technical Reports

UTR #16: UTF-EBCDIC

This document presents the specifications of UTF-EBCDIC: EBCDIC Friendly Unicode (or UCS) Transformation Format.

UTR #17: Character Encoding Model

This document clarifies a number of the terms used to describe character encodings and indicates where the different encoding forms of the Unicode Standard fit in. It elaborates the Internet Architecture Board's (IAB) three-layer "text stream" definitions into a five-layer structure.

UTR #20: Unicode in XML and Other Markup Languages

This document contains guidelines on the use of the Unicode Standard in conjunction with markup languages such as XML.

UTR #23: The Unicode Character Property Model

This document presents a conceptual model of character properties defined in the Unicode Standard.

UTR #25: Unicode Support for Mathematics

The Unicode Standard includes virtually all of the standard characters used in mathematics. This set supports a variety of math applications on computers, including document presentation languages like $\text{T}_\text{E}\text{X}$, math markup languages like MathML and OpenMath, internal representations of mathematics in systems like Mathematica, Maple, and MathCAD, computer programs, and plain text. This document describes the Unicode mathematics character groups and gives some of their imputed default math properties.

UTR #26: Compatibility Encoding Scheme for UTF-16: 8-Bit (CESU-8)

This document specifies an 8-bit Compatibility Encoding Scheme for UTF-16 (CESU) that is intended for internal use within systems processing Unicode to provide an ASCII-compatible 8-bit encoding that is similar to UTF-8 but preserves UTF-16 binary collation. *It is not intended or recommended as an encoding used for open information exchange.* The Unicode Consortium does not encourage the use of CESU-8, but does recognize the existence of data in this encoding and supplies this technical report to clearly define the format and to distinguish it from UTF-8. This encoding does not replace or amend the definition of UTF-8.

UTR #30: Character Foldings

This report identifies a set of operations that map similar characters to a common target. Such operations, called character foldings, are used to ignore certain distinctions between similar characters. The report also provides an algorithm for applying these operations to searching, plus additional guidelines.

UTR #36: Unicode Security Considerations

Because Unicode contains such a large number of characters and incorporates the varied writing systems of the world, incorrect usage can expose programs or systems to possible security attacks. This document describes some of the security considerations that programmers, system analysts, standards developers, and users should take into account, and it provides specific recommendations to reduce the risk of problems.

B.5 Unicode Technical Notes

Unicode Technical Notes provide information on a variety of topics related to Unicode and internationalization technologies.

These technical notes are independent publications, not approved by any of the Unicode Technical Committees, nor are they part of the Unicode Standard or any other Unicode specification. Publication does not imply endorsement by the Unicode Consortium in any way. These documents are not subject to the Unicode Patent Policy. Unicode Technical Notes can be found on the Unicode Web site at:

<http://www.unicode.org/notes/>

The technical notes cover the following topics (among others):

- Algorithms
- Collation
- Compression and code set conversions
- Language identification

- Migration of software
- Modern and historical scripts
- Text layout and rendering
- Tutorials
- Social and cultural issues

B.6 Other Unicode Online Resources

The Unicode Consortium provides a number of online resources for obtaining information and data about the Unicode Standard as well as updates and corrigenda.

Unicode Online Resources

Unicode Web Site

<http://www.unicode.org>

Unicode Anonymous FTP Site

<ftp://ftp.unicode.org>

Charts. The charts section of the Web site provides online charts for all of the Unicode characters, plus specialized charts for normalization, collation, case mapping, script names, and Unified CJK Ideographs.

<http://www.unicode.org/charts/>

Common Locale Data Registry (CLDR). Machine-readable repository, in XML format, of locale information for use in application and system development.

<http://www.unicode.org/cldr/>

Conferences. The Internationalization and Unicode Conferences are of particular value to anyone implementing the Unicode Standard or working on internationalization. A variety of tutorials and conference sessions cover current topics related to the Unicode Standard, the World Wide Web, software, internationalization, and localization.

<http://www.unicode.org/conference/>

E-mail Discussion List. Subscription instructions for the public e-mail discussion list are posted on the Unicode Web site.

FAQ (Frequently Asked Questions). The FAQ pages provide an invaluable resource for understanding the Unicode Standard and its implications for users and implementers.

<http://www.unicode.org/faq/>

Online Unicode Character Database. This page supplies information about the online Unicode Character Database (UCD), including links to documentation files and the most up-to-date version of the data files, as well as instructions on how to access any particular version of the UCD.

<http://www.unicode.org/ucd/>

Online Unihan Database. The online Unihan Database provides interactive access to all of the property information associated with CJK ideographs in the Unicode Standard.

<http://www.unicode.org/chart/unihan.html>

Policies. These pages describe Unicode Consortium policies on stability, patents, and Unicode Web site privacy. The stability policies are particularly important for implementers, documenting invariants for the Unicode Standard that allow implementations to be compatible with future and past versions. Accordingly, the stability policies are also reprinted in this book in *Appendix F, Unicode Encoding Stability Policies*, for easy reference.

<http://www.unicode.org/policies/>

Updates and Errata. This page lists periodic updates with corrections of typographic errors and new clarifications of the text.

<http://www.unicode.org/errata/>

Versions. This page describes the version numbering used in the Unicode Standard, the nature of the Unicode character repertoire, and ways to cite and reference the Unicode Standard, the Unicode Character Database, and Unicode Technical Reports. It also specifies the exact contents of each and every version of the Unicode Standard, back to Unicode 1.0.0.

<http://www.unicode.org/versions/>

Where Is My Character? This page provides basic guidance to finding Unicode characters, especially those whose glyphs do not appear in the charts, or that are represented by sequences of Unicode characters.

<http://www.unicode.org/standard/where/>

How to Contact the Unicode Consortium

The best way to contact the Unicode Consortium to obtain membership information or order additional copies of this book is via the Web site:

<http://www.unicode.org/contacts.html>

The Web site also lists the current telephone, fax, and courier delivery address. The Consortium's postal address is:

P.O. Box 391476
Mountain View, CA 94039-1476
USA