

This PDF file is an excerpt from *The Unicode Standard, Version 4.0*, issued by the Unicode Consortium and published by Addison-Wesley. The material has been modified slightly for this online 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/standard/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 Addison-Wesley was aware of a trademark claim, the designations have been printed in initial capital letters. However, not all words in initial capital letters are trademark designations.

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 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, <http://www.mehallo.com>

The publisher offers discounts on this book when ordered in quantity for bulk purchases and special sales. For more information, customers in the U.S. please contact U.S. Corporate and Government Sales, (800) 382-3419, corpsales@pearsontechgroup.com. For sales outside of the U.S., please contact International Sales, +1 317 581 3793, international@pearsontechgroup.com

Visit Addison-Wesley on the Web: <http://www.awprofessional.com>

Library of Congress Cataloging-in-Publication Data

The Unicode Standard, Version 4.0 : the Unicode Consortium /Joan Aliprand... [et al.].

p. cm.

Includes bibliographical references and index.

ISBN 0-321-18578-1 (alk. paper)

1. Unicode (Computer character set). I. Aliprand, Joan.

QA268.U545 2004

005.7'2—dc21

2003052158

Copyright © 1991–2003 by Unicode, Inc.

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, recording or otherwise, without the prior written permission of the publisher or Unicode, Inc. Printed in the United States of America. Published simultaneously in Canada.

For information on obtaining permission for use of material from this work, please submit a written request to the Unicode Consortium, Post Office Box 39146, Mountain View, CA 94039-1476, USA, Fax +1 650 693 3010 or to Pearson Education, Inc., Rights and Contracts Department, 75 Arlington Street, Suite 300 Boston, MA 02116, USA, Fax: +1 617 848 7047.

ISBN 0-321-18578-1

Text printed on recycled paper

1 2 3 4 5 6 7 8 9 10—CRW—0706050403

First printing, August 2003

Chapter 16

Code Charts

Disclaimer

Character images shown in the code charts are not prescriptive. In actual fonts, considerable variations are to be expected.

The code charts that follow present the characters of the Unicode Standard. Characters are organized into related groups called *blocks*. In the Unicode Standard, character blocks generally contain characters from a single script. In many cases, a script is fully represented in its character block. There are, however, important exceptions, most notably in the area of punctuation characters.

A character names list follows each character chart, except for CJK ideographs and Hangul syllables, as discussed in *Section 16.2, CJK Unified Ideographs*, and *Section 16.3, Hangul Syllables*, respectively. The character names list itemizes every character in the block and provides supplementary information in many cases.

An index to distinctive character names is found at the back of this book; a full set of character names (including earlier Version 1.0 names) is in the Unicode Character Database.

16.1 Character Names List

The following illustration identifies the components of typical entries in the character names list.

| <i>code</i> | <i>image</i> | <i>entry</i> | |
|-------------|--------------|--|---|
| 00AE | ® | REGISTERED SIGN = REGISTERED TRADE MARK SIGN | (Version 1.0 name) |
| 00AF | - | MACRON = overline, APL overbar • this is a spacing character → 02C9 ¯ modifier letter macron → 0304 ◌ combining macron → 0305 ◌ combining overline ≈ 0020 ☐ 0304 ◌ | (Unicode name) (alternative names) (informative note) (cross reference) (compatibility decomposition) |
| 00E5 | å | LATIN SMALL LETTER A WITH RING ABOVE • Danish, Norwegian, Swedish, Walloon ≡ 0061 a 030A ◌ | (canonical decomposition) |

Images in the Code Charts and Character Lists

Each character in these code charts is shown with a representative glyph. A representative glyph is not a prescriptive form of the character, but one that enables recognition of the intended character to a knowledgeable user and facilitates lookup of the character in the code charts. In many cases, there are more or less well-established alternative glyphic representations for the same character.

Designers of high-quality fonts will do their own research into the preferred glyphic appearance of Unicode characters. In addition, many scripts require context-dependent glyph shaping, glyph positioning, or ligatures, none of which is shown in the code charts.

The representative glyphs in the code charts are based on a serified, Times-like font. For example, even the ASCII character U+0061 LATIN SMALL LETTER A has two common alternative forms, the “a” used in Times and the “a” that occurs in many other font styles. In a Times-like font, the character U+03A5 GREEK CAPITAL LETTER UPSILON looks like “Υ”; the form Υ is common in other font styles.

A different case is U+010F LATIN SMALL LETTER D WITH CARON, which is commonly typeset as d' instead of \acute{d} . In such cases, the code charts show the more common variant in preference to a more didactic archetypical shape.

Many characters have been unified and have different appearances in different language contexts. The shape shown for U+2116 № NUMERO SIGN is a fullwidth shape as it would be used in East Asian fonts. In Cyrillic usage, № is the universally recognized glyph.

In certain cases, characters need to be represented by more or less condensed, shifted, or distorted glyphs to make them fit the format of the code charts. For example, U+0D10 ഐ MALAYALAM LETTER AI is shown in a reduced size to fit the character cell.

Sometimes characters need to be given artificial shapes to make them recognizable in the code charts. Examples are U+00AD ⎯ SOFT HYPHEN and U+2011 ⎯ NON-BREAKING HYPHEN, where the special behavior of the hyphen is indicated by the dashed box and the letters.

When characters are used in context, the surrounding text gives important clues as to identity, size, and positioning. In the code charts, these clues are absent. For example, U+2075 ⁵ SUPERSCRIPT FIVE is shown much smaller than it would be in a Times-like text font.

Combining characters are shown with a dotted circle—for example, U+0940 ॠ DEVANAGARI VOWEL SIGN II. The relative position of the dotted circle gives an approximate indication of the location of the base character in relation to the combining mark. During rendering, additional adjustments are necessary. Accents such as U+0302 COMBINING CIRCUMFLEX ACCENT are adjusted vertically and horizontally based on the height and width of the base character, as in “ \hat{i} ” versus “ \hat{W} ”.

For non-European scripts, typical typefaces were selected that allow as much distinction as possible among the different characters.

The Unicode Standard contains many characters that are used in writing minority languages or that are historical characters, often used primarily in manuscripts or inscriptions. Where there is no strong tradition of printed materials, the typography of a character may not be settled.

Character Names

The character names in the code charts precisely match the normative character names in the Unicode Character Database. Character names are unique and stable. By convention

they are in uppercase. Because character names are stable, mistaken names will not be revised, but may be annotated. For example:

2118 ⅈ SCRIPT CAPITAL P
 = Weierstrass elliptic function
 • actually this has the form of a lowercase calligraphic p, despite its name

Aliases

An alias (preceded by =) is an alternate name for a character. Characters may have several aliases, and aliases for different characters are not guaranteed to be unique. Aliases are informative and may be updated. By convention, aliases are in lowercase, except where they contain proper names. Where an alias matches the name of a character in *The Unicode Standard, Version 1.0*, it is listed first and given in all caps. Because the formal character names may differ in unexpected ways from commonly used names (for example, PILCROW SIGN = paragraph sign), some aliases may be useful alternate choices for indicating characters in user interfaces. In the Hangul Jamo block, U+1100..U+11FF, the normative short jamo names are given as aliases in uppercase.

Cross References

Cross references (preceded by →) are used to indicate a related character of interest, but without indicating the nature of the relation. Possibilities are a different character of similar appearance or name, the other member of a case pair, or some other linguistic relationship.

Explicit Inequality. The two characters are not identical, although the glyphs that depict them are identical or very close.

003A : COLON
 → 0589 : armenian full stop
 → 2236 : ratio

Other Linguistic Relationships. These relationships include transliterations (such as between Serbian and Croatian), typographically unrelated characters used to represent the same sound, and so on.

01C9 lj LATIN SMALL LETTER LJ
 → 0459 љ cyrillic small letter lje
 ≈ 006C l 006A j

Cross references are neither exhaustive nor symmetric. Typically a general character would have cross references to more specialized character, but not the other way around.

Information About Languages

An informative note may include a list of the language(s) using that character where this information is considered useful. For case pairs, the annotation is given only for the lowercase form, to avoid needless repetition. An ellipsis “...” indicates that the listed languages cited are merely the principal ones among many.

Case Mappings

When a case mapping corresponds *solely* to a difference based on SMALL versus CAPITAL in the names of the characters, the case mapping is not given in the names list but only in the Unicode Character Database.

0041 A LATIN CAPITAL LETTER A

01F2 Dz LATIN CAPITAL LETTER D WITH SMALL LETTER Z
 ≈ 0044 D 007A z

When the case mapping cannot be predicted from the name, the information is given in a note.

00DF ß LATIN SMALL LETTER SHARP S
 = Eszett
 • German
 • uppercase is “SS”
 • in origin a ligature of 017F f and 0073 s
 → 03B2 β greek small letter beta

Decompositions

The decomposition sequence (one or more letters) given for a character is either its canonical mapping or its compatibility mapping. The canonical mapping is marked with an *identical to* symbol ≡.

00E5 â LATIN SMALL LETTER A WITH RING ABOVE
 • Danish, Norwegian, Swedish, Walloon
 ≡ 0061 a 030A ̊

212B Å ANGSTROM SIGN
 ≡ 00C5 Ä latin capital letter a with ring above

Compatibility mappings are marked with an *almost equal to* symbol ≈. Formatting information may be indicated inside angle brackets.

01F2 Dz LATIN CAPITAL LETTER D WITH SMALL LETTER Z
 ≈ 0044 D 007A z

FF21 A FULLWIDTH LATIN CAPITAL LETTER A
 ≈ <wide> 0041 A

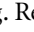
The following compatibility formatting tags are used in the Unicode Character Database:

| | |
|------------|---|
| | A font variant (for example, a blackletter form) |
| <noBreak> | A no-break version of a space, hyphen, or other punctuation |
| <initial> | An initial presentation form (Arabic) |
| <medial> | A medial presentation form (Arabic) |
| <final> | A final presentation form (Arabic) |
| <isolated> | An isolated presentation form (Arabic) |
| <circle> | An encircled form |
| <super> | A superscript form |
| <sub> | A subscript form |
| <vertical> | A vertical layout presentation form |
| <wide> | A wide (or zenkaku) compatibility character |
| <narrow> | A narrow (or hankaku) compatibility character |
| <small> | A small variant form (CNS compatibility) |
| <square> | A CJK squared font variant |
| <fraction> | A vulgar fraction form |
| <compat> | Otherwise unspecified compatibility character |

In the character names list accompanying the code charts, the “<compat>” label is suppressed, but all other compatibility formatting tags are explicitly listed in the compatibility mapping.


Decompositions are not necessarily full decompositions. For example, the decomposition for U+212B Å ANGSTRÖM SIGN can be further decomposed using the canonical mapping for U+00C5 Å LATIN SMALL LETTER A WITH RING ABOVE. (For more information on decomposition, see *Section 3.7, Decomposition*.)

Reserved Characters


Character codes that are marked “<reserved>” are unassigned and reserved for future encoding. Reserved codes are indicated by a  glyph. To ensure readability, many instances of reserved characters have been suppressed from the names list.



060D  <reserved>

Reserved codes may also have cross references to assigned characters located elsewhere.

2073  <reserved>
→ 00B3³ superscript three

Noncharacters

Character codes that are marked “<not a character>” refer to noncharacters. They are designated code points that will never be assigned to a character. These codes are indicated by a  glyph. Noncharacters are only shown in the code charts where they occur together with other characters in the same block. For a complete list of noncharacters, see *Section 15.8, Noncharacters*.

FFFF  <not a character>
• the value FFFF  is guaranteed not to be a Unicode character at all

Subheads

The character names list contains a number of informative subheads that help divide up the list into smaller sublists of similar characters. For example, in the Miscellaneous Symbols block, U+2600..U+26FF, there are subheads for “Astrological symbols,” “Chess symbols,” and so on. Such subheads are editorial and informative, and should not be taken as providing any definitive, normative status information about characters in the sublists they mark or about any constraints on what characters could be encoded in the future at reserved code points within their ranges. The subheads are subject to change.

16.2 CJK Unified Ideographs

A character names list is not provided for any of the CJK Unified Ideograph character blocks because the name of a unified ideograph simply consists of its Unicode code point preceded by CJK UNIFIED IDEOGRAPH-.

As with other character charts, each Unicode character in these blocks is shown with its Unicode code point and a single representative glyph. Note that varying typographic practices throughout East Asia may require glyphs other than the representative one to be used so that the display is correct for a particular country or language.

A table providing mappings between the CJK ideographs included in the Unicode Standard and those in other character set standards is included in Unihan.txt in the Unicode Character Database.

A radical-stroke index to CJK ideographs is in *Chapter 17, Han Radical-Stroke Index*.

16.3 Hangul Syllables

A character names list is not provided for characters in the Hangul Syllables block, U+AC00..U+D7AF, because the name of a Hangul syllable can be determined by algorithm as described in *Section 3.12, Conjoining Jamo Behavior*. The short names used in that algorithm are listed in the code charts as aliases in the Hangul Jamo block, U+1100..U+11FF.

The code charts, pages 420-1186 in the book, are omitted here. Please see the online charts at:

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

Note: The online code charts are continuously updated and may contain characters added after the publication of *The Unicode Standard, Version 4.0*. To find out whether a particular character was part of the Unicode Standard, Version 4.0, please consult either the printed edition of the standard (ISBN 0-321-18578-1) or version 4.0.0 of the Unicode Character Database. Normative references to the Unicode Standard, Version 4.0, should use the printed edition.