

## Chapter 10

# *East Asian Scripts*

All of the writing systems of East Asia have as their root the ideographic script developed in China in the second millennium BCE. Originally intended to write the various Chinese dialects, this script was exported throughout China's sphere of influence and was used for centuries to write even non-Chinese languages such as Japanese, Korean, and Vietnamese. Speakers of other languages, such as Yi, created their own ideographs in imitation of Chinese.

The Chinese language is largely monosyllabic and noninflecting, and an ideographic writing system suits it well. Ideographs are less well suited for unrelated languages. In Japan, this problem was solved by creating two syllabic scripts, *hiragana* and *katakana*. In Korea, an alphabetic system was invented, with letters grouped into ideograph-like syllabic blocks called *hangul*.

Ideographs continued to be the exclusive way of writing Vietnamese until the twentieth century, when they were replaced with an alphabet based on the Latin script. They are still extensively used in Japanese, where they are called *kanji*, and somewhat more rarely in Korean, where they are called *hanja*. In mainland China, the government has promoted the use of more modern, simplified forms of the ideographs over the older, more traditional forms used in Taiwan and overseas Chinese communities.

*Appendix A, Han Unification History*, describes how the diverse typographic traditions of mainland China, Taiwan, Japan, Korea, and Vietnam have been reconciled to provide a common set of ideographs in the Unicode Standard for all of these languages and regions.

The Unicode Standard includes a complete set of Korean *hangul* syllables, as well as the individual letters (*jamos*), which can be also be used to write Korean. *Section 3.11, Conjoining Jamo Behavior*, describes how to use the conjoining *jamos* and how to convert between two methods for representing Korean.

In all East Asian scripts, individual characters are written within uniformly sized squares. Diacritical marks are rarely used, although phonetic annotations are not uncommon. Traditionally, East Asian scripts were written from the top of the page downward, with columns running right to left across the page. Under the influence of Western practices, a horizontal left-to-right writing sequence is now common.

Many older character sets include characters intended to simplify the implementation of East Asian scripts, such as variant punctuation forms for text written vertically, halfwidth forms (which occupy only half a square), and fullwidth forms (which allow Latin letters to occupy a full square). These characters are included in the Unicode Standard for compatibility with these older standards.

---

## 10.1 Han

### CJK Unified Ideographs

These blocks contain a set of Unified Han ideographic characters used in the written Chinese, Japanese, and Korean languages.<sup>1</sup> The term *Han*, derived from the Chinese Han Dynasty, refers generally to Chinese traditional culture. The Han ideographic characters make up a coherent script, which was traditionally written vertically, with the vertical lines ordered from right to left. In modern usage, especially in technical works and in computer-rendered text, the Han script is written horizontally from left to right and is freely mixed with Latin or other scripts. When used in writing Japanese or Korean, the Han characters are interspersed with other scripts unique to those languages (Hiragana and Katakana for Japanese; Hangul syllables for Korean).

The Han ideographic characters constitute a very large set, numbering in the tens of thousands. They have a long history of use in East Asia. Enormous compendia of Han ideographic characters exist because of a continuous, millennia-long scholarly tradition of collecting all Han character citations, including variant, mistaken, and nonce forms, into annotated character dictionaries.

Because of the large size of the Han ideographic character repertoire, and because of the particular problems the characters pose for standardizing their encoding, this character block description is more extended than that for other scripts and is divided into subsections. The first three subsections (CJK Standards, Blocks, and Mapping to Standards) describe the character set standards used as sources, the way in which the Unicode Standard divides ideographs into blocks, and some of the issues involved in mapping these characters to other character sets. These subsections are followed by an extended discussion of the characteristics of Han characters, with particular attention being paid to the problem of unification of encoding for characters used for different languages. Then there is a formal statement of the principles behind the Unified Han character encoding adopted in the Unicode Standard and order of their arrangement. For a detailed account of the background and history of development of the Unified Han character encoding, see also *Appendix A, Han Unification History*.

#### **CJK Standards**

The Unicode Standard draws its Han character repertoire of 27,484 characters from a number of character set standards. These standards are grouped into six sources, as indicated in *Table 10-1*. The primary work of unifying and ordering the characters from these sources was done by the Ideographic Rapporteur Group (IRG), a subgroup of ISO/IEC JTC1/SC2/WG2.

The G, T, J, K, and V sources represent the characters submitted to the IRG by its member bodies. The G source consists of submissions from mainland China, the Hong Kong SAR, and Singapore. The other four are the submissions from Taiwan, Japan, Korea, and Vietnam, respectively. The U source represents character set standards that were not submitted to the IRG by any member body but which were used by the Unicode Consortium.

---

1. Although the term “CJK”—Chinese, Japanese, and Korean—is used throughout this text to describe the languages that currently use Han ideographic characters, it should be noted that earlier Vietnamese writing systems were based on Han ideographs. Consequently, the term “CJKV” would be more accurate in a historical sense. Han ideographs are still used for historical, religious, and pedagogical purposes in Vietnam.

**Table 10-1. Sources for Unified Han**

|           |           |  |                          |
|-----------|-----------|--|--------------------------|
| G source: | G0        | GB2312-80  |                          |
|           | G1        | GB12345-90 with 58 Hong Kong and 92 Korean “Idu” characters  |                          |
|           | G3        | GB7589-87 unsimplified forms   |                          |
|           | G5        | GB7590-87 unsimplified forms   |                          |
|           | G7        | General Purpose Hanzi List for Modern Chinese Language, and General List of Simplified Hanzi   |                          |
|           | GS        | Singapore Characters   |                          |
|           | G8        | GB8565-88  |                          |
|           | GE        | GB16500-95   |                          |
|           | T source: | T1   | CNS 11643-1992 1st plane |
|           |           | T2   | CNS 11643-1992 2nd plane |
| T3        |           | CNS 11643-1992 3rd plane with some additional characters   |                          |
| T4        |           | CNS 11643-1992 4th plane   |                          |
| T5        |           | CNS 11643-1992 5th plane   |                          |
| T6        |           | CNS 11643-1992 6th plane   |                          |
| T7        |           | CNS 11643-1992 7th plane   |                          |
| TF        |           | CNS 11643-1992 15th plane  |                          |
| J source: | J0        | JIS X 0208-1990  |                          |
|           | J1        | JIS X 0212-1990  |                          |
|           | JA        | Unified Japanese IT Vendors Contemporary Ideographs, 1993  |                          |
| K source: | K0        | KS C 5601-1987 (unique ideographs)   |                          |
|           | K1        | KS C 5657-1991   |                          |
|           | K2        | PKS C 5700-1 1994  |                          |
|           | K3        | PKS C 5700-2 1994  |                          |
| V source: | V0        | TCVN 5773:1993   |                          |
|           | V1        | TCVN 6056:1995   |                          |
| U source: |           | KS C 5601-1987 (duplicate ideographs)<br>ANSI Z39.64-1989 (EACC)<br>Big-5 (Taiwan)<br>CCCII, level 1<br>GB 12052-89 (Korean)<br>JEF (Fujitsu)<br>PRC Telegraph Code<br>Taiwan Telegraph Code (CCDC)<br>Xerox Chinese<br>Han Character Shapes Permitted for Personal Names (Japan)<br>IBM Selected Japanese and Korean Ideographs |                          |

In some cases, the entire ideographic repertoire of the original character set standards was *not* included in the corresponding source. Three reasons explain this decision:

1. Where the repertoires of two of the character set standards within a single source have considerable overlap, the characters in the overlap might be included only once in the source. This approach is used, for example, with GB 2312-80 and GB 12345-90, which have many ideographs in common. Characters in GB 12345-90 that are duplicates of characters in GB 2312-80 are not included in the G source.
2. Where a character set standard is based on unification rules that differ substantially from those used by the IRG, many variant characters found in the character set standard will not be included in the source. This situation is the case with CNS 11643-1992, EACC, and CCCII. It is the only case where full round-trip compatibility with the Han ideograph repertoire of the relevant character set standards is not guaranteed.
3. KS C 5601-1987 contains numerous duplicate ideographs included because they have multiple pronunciations in Korean. These multiply-encoded ideographs are not included in the K source but are included in the U source to provide full round-trip compatibility with KS C 5601-1987.

## Blocks

Ideographs are found in three blocks of the Unicode Standard: the CJK Unified Ideographs block (U+4E00–U+9FFF, common ideographs), the CJK Unified Ideographs Extension A block (U+3400–U+4DFF, rare ideographs), and the CJK Compatibility Ideographs block (U+F900–U+FAFF, duplicates, unifiable variants, and corporate characters).

Characters in the CJK Unified Ideographs and CJK Unified Ideographs Extension A blocks are defined by the IRG and are derived entirely from the G, T, J, K, and V sources.

The CJK Unified Ideographs block represents characters submitted to the IRG prior to 1992 and consists of commonly used characters. Characters in the CJK Unified Ideographs Extension A block are rarer, were submitted to the IRG between 1992 and 1998, and are not unifiable with characters in the CJK Unified Ideographs block.

The only difference in the unification work done by the IRG on these two blocks is that the source separation rule was applied to the CJK Unified Ideographs block and not to the CJK Unified Ideographs Extension A block. This rule states that ideographs that are distinct in a source must not be unified. (For further discussion, see the subsection “Principles” later in this section.)

Characters unique to the U source are found in the CJK Compatibility Ideographs block. There are 12 of these characters: U+FA0E, U+FA0F, U+FA11, U+FA13, U+FA14, U+FA1F, U+FA21, U+FA23, U+FA24, U+FA27, U+FA28, and U+FA29. The remaining characters in the CJK Compatibility Ideographs block are either duplicates or unifiable variants of characters in the one of the other blocks and are included in the Unicode Standard for reasons of round-trip compatibility.

## General Characteristics of Han Ideographs

The authoritative Japanese dictionary *Kouzien* defines Han characters to be

characters that originated among the Chinese to write the Chinese language. They are now used in China, Japan, and Korea. They are logographic (each character represents a word, not just a sound) characters that developed from pictographic and ideographic principles. They are also used phonetically. In Japan they are generally called *kanzi* (Han, that is, Chinese, characters) including the “national characters” (*kokuzi*) such as *touge* (mountain pass), which have been created using the same principles. They are also called *mana* (true names, as opposed to *kana*, false or borrowed names).<sup>1</sup>

For many centuries, written Chinese was the accepted written standard throughout East Asia. The influence of the Chinese language and its written form on the modern East Asian languages is similar to the influence of Latin on the vocabulary and written forms of languages in the West. This influence is immediately visible in the mixture of Han characters and native phonetic scripts (*kana* in Japan, *hangul* in Korea) as now used in the orthographies of Japan and Korea (see *Table 10-2*).

The evolution of character shapes and semantic drift over the centuries has resulted in changes to the original forms and meanings. For example, the Chinese character 湯 *tang* (Japanese *tou* or *yu*, Korean *thang*), which originally meant “hot water,” has come to mean “soup” in Chinese. “Hot water” remains the primary meaning in Japanese and Korean, whereas “soup” appears in more recent borrowings from Chinese, such as “soup noodles”

---

1. Lee Collins' translation from the Japanese, *Kouzien*, Izuru, Shinmura, ed. (Tokyo: Iwanami Syoten, 1983).

**Table 10-2. Common Han Characters**

| <i>Han Character</i> | <i>Chinese<sup>a</sup></i> | <i>Japanese</i> | <i>Korean</i> | <i>English Translation</i> |
|----------------------|----------------------------|-----------------|---------------|----------------------------|
| 天                    | tian <sup>1</sup>          | ten, ame        | chen          | heaven, sky                |
| 地                    | di <sup>4</sup>            | ti, tuti        | ci            | earth, ground              |
| 人                    | ren <sup>2</sup>           | zin, hito       | in            | man, person                |
| 山                    | shan <sup>1</sup>          | san, yama       | san           | mountain                   |
| 水                    | shui <sup>3</sup>          | sui, mizu       | swu           | water                      |
| 上                    | shang <sup>4</sup>         | zyou, ue        | sang          | above                      |
| 下                    | xia <sup>4</sup>           | ka, sita        | ha            | below                      |

a. The superscripted numbers in this table represent Chinese (Mandarin) tone marks.

(Japanese *tanmen*; Korean *thangmyen*). Still, the identical appearance and similarities in meaning are dramatic and more than justify the concept of a unified Han script that transcends language.

The “nationality” of the Han characters became an issue only when each country began to create coded character sets (for example, China’s GB 2312-80, Japan’s JIS X 0208-1978, and Korea’s KS C 5601-87) based on purely local needs. This problem appears to have arisen more from the priority placed on local requirements and lack of coordination with other countries, rather than out of conscious design. Nevertheless, the identity of the Han characters is fundamentally independent of language, as shown by dictionary definitions, vocabulary lists, and encoding standards.

**Terminology.** Several standard romanizations of the term used to refer to East Asian ideographic characters are commonly used. They include *hanzi* (Chinese), *kanzi* (Japanese), *kanji* (colloquial Japanese), *hanja* (Korean), and Chữ hán (Vietnamese). The standard English translations for these terms are interchangeable: Han character, Han ideographic character, East Asian ideographic character, or CJK ideographic character. For the purpose of clarity, the Unicode Standard uses some subset of the English terms when referring to these characters. The term *Kanzi* is used in reference to a specific Japanese government publication. The unrelated term *KangXi* (which is a Chinese reign name, rather than another romanization of “Han character”) is used only when referring to the dictionary on which the Unified Repertoire and Ordering, Version 2.0, was based.

**Distinguishing Han Character Usage Between Languages.** There is some concern that unifying the Han characters may lead to confusion because they are sometimes used differently by the various East Asian languages. Computationally, Han character unification presents no more difficulty than employing a single Latin character set that is used to write languages as different as English and French. Programmers do not expect the characters “c”, “h”, “a”, and “t” alone to tell us whether *chat* is a French word for cat or an English word meaning “informal talk.” Likewise, we depend on context to identify the American hood (of a car) with the British bonnet. Few computer users are confused by the fact that ASCII can also be used to represent such words as the Welsh word *ynghyd*, which are strange looking to English eyes. Although it would be convenient to identify words by language for programs such as spell-checkers, it is neither practical nor productive to encode a separate Latin character set for every language that uses it.

Similarly, the Han characters are often combined to “spell” words whose meaning may not be evident from the constituent characters. For example, the two characters “to cut” and “hand” mean “postal stamp” in Japanese, but the compound may appear to be nonsense to a speaker of Chinese or Korean (see *Figure 10-1*).

### Figure 10-1. Han Spelling

|        |   |      |   |                       |
|--------|---|------|---|-----------------------|
| 切      | + | 手    | = | 1. Japanese “stamp”   |
| to cut |   | hand |   | 2. Chinese “cut hand” |

Even within one language, a computer requires context to distinguish the meanings of words represented by coded characters. The word *chuugoku* in Japanese, for example, may refer to China or to a district in central west Honshuu (see *Figure 10-2*).

### Figure 10-2. Context for Characters

|        |   |         |   |                                 |
|--------|---|---------|---|---------------------------------|
| 中      | + | 国       | = | 1. China                        |
| middle |   | country |   | 2. Chuugoku district of Honshuu |

Coding these two characters as four so as to capture this distinction would probably cause more confusion and still not provide a general solution. The Unicode Standard leaves the issues of language tagging and word recognition up to a higher level of software and does not attempt to encode the language of the Han characters.

**Sorting Han Ideographs.** The Unicode Standard does not define a method by which ideographic characters are sorted; the requirements for sorting differ by locale and application. Possible collating sequences include phonetic, radical-stroke (*KangXi*, *Xinhua Zidian*, and so on), four-corner, and total stroke count. Raw character codes alone are seldom sufficient to achieve a usable ordering in any of these schemes; ancillary data are usually required. (See *Table 10-5*.)

**Character Glyphs.** In form, Han characters are monospaced. Every character takes the same vertical and horizontal space, regardless of how simple or complex its particular form is. This practice follows from the long history of printing and typographical practice in China, which traditionally placed each character in a square cell. When written vertically, there are also a number of named cursive styles for Han characters, but the cursive forms of the characters tend to be quite idiosyncratic and are not implemented in general-purpose Han character fonts for computers.

There may be a wide variation in the glyphs used in different countries and for different applications. The most commonly used typefaces in one country may not be used in others.

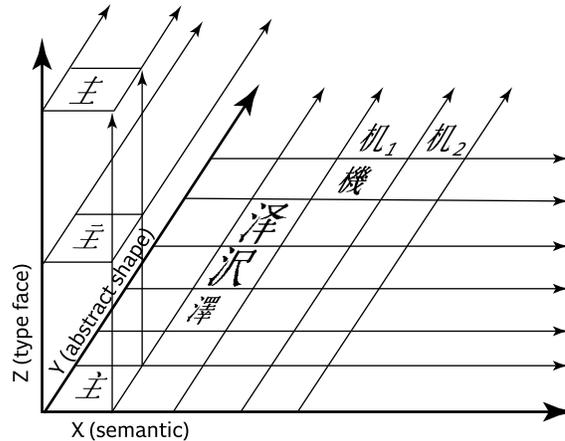
The types of glyphs used to depict characters in the Han ideographic repertoire of the Unicode Standard have been constrained by available fonts. Users are advised to consult authoritative sources for the appropriate glyphs for individual markets and applications. It is assumed that most Unicode implementations will provide users with the ability to select the font (or mixture of fonts) that is most appropriate for a given locale.

## Principles

**Three-Dimensional Conceptual Model.** To develop the explicit rules for unification, a conceptual framework was developed to model the nature of Han ideographic characters. This model expresses written elements in terms of three primary attributes: semantic

(meaning, function), abstract shape (general form), and actual shape (instantiated, type-face form). These attributes are graphically represented in three dimensions according to the *X*, *Y*, and *Z* axes (see *Figure 10-3*).

**Figure 10-3. Three-Dimensional Conceptual Model**



The semantic attribute (represented along the *X* axis) distinguishes characters by meaning and usage. Distinctions are made between entirely unrelated characters such as 澤 (marsh) and 機 (machine) as well as extensions or borrowings beyond the original semantic cluster such as 机<sub>1</sub> (a phonetic borrowing used as a simplified form of 機) and 机<sub>2</sub> (table, the original meaning).

The abstract shape attribute (the *Y* axis) distinguishes the variant forms of a single character with a single semantic attribute (that is, a character with a single position on the *X* axis).

The actual shape (typeface) attribute (the *Z* axis) is for differences of type design (the actual shape used in imaging) of each variant form.

Only characters that have the same abstract shape (that is, occupy a single point on the *X* and *Y* axes) are potential candidates for unification. *Z* axis typeface and stylistic differences are generally ignored.

**Unification Rules.** The following rules were applied during the process of merging Han characters from the different source character sets:

**R1 Source Separation Rule.** *If two ideographs are distinct in a primary source standard, then they are not unified.*

- This rule is sometimes called the *round-trip rule* because its goal is to facilitate a round-trip conversion of character data between an IRG source standard and the Unicode Standard without loss of information.
- This rule was applied only for the work on the original Unified Repertoire and Ordering (URO). The IRG dropped this rule in 1992 and will not use it in future work.

For example, the ideographs from the URO shown in *Figure 10-4* would normally be subject to unification by rule R3; however, their unification is prevented because they are distinct in the primary source standard J<sub>0</sub> (JIS X 0208-1990).

**R2 Noncognate Rule.** *In general, if two ideographs are unrelated in historical derivation (noncognate characters), then they are not unified.*

### Figure 10-4. Preserving Variants

劍 劍 劔 劔 劔 劔

“sword”

For example, the following ideographs (in *Figure 10-5*), although visually quite similar, are nevertheless not unified because they are historically unrelated and have distinct meanings.

### Figure 10-5. Not Cognates, Not Unified

土 ≠ 士  
earth warrior, scholar

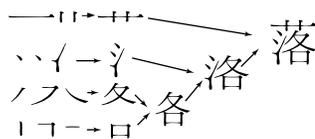
**R3** *By means of a two-level classification (described next), the abstract shape of each ideograph is determined. Any two ideographs that possess the same abstract shape are then unified provided that their unification is not disallowed by either the source separation rule or the noncognate rule.*

**Two-Level Classification.** Using the three-dimensional model, characters are analyzed in a two-level classification. The two-level classification distinguishes characters by abstract shape (Y axis) and actual shape of a particular typeface (Z axis). Variant forms are identified based on the difference of abstract shapes.

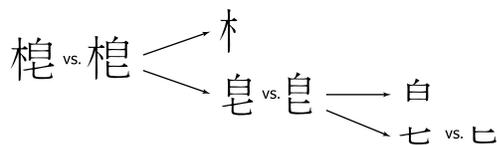
To determine differences in abstract shape and actual shape, the structure and features of each component of an ideograph are analyzed as follows.

**Ideograph Component Structure.** The component structure of each ideograph is examined. A component is a geometrical combination of primitive elements. Various ideographs can be configured with these components used in conjunction with other components. Some components can be combined to make a component more complicated in its structure. Therefore, an ideograph can be defined as a component tree with the entire ideograph as the root node and with the bottom nodes consisting of primitive elements (see *Figure 10-6* and *Figure 10-7*).

### Figure 10-6. Component Structure



### Figure 10-7. The Most Superior Node of a Component



**Ideograph Features.** The following features of each ideograph to be compared are examined:

- Number of components
- Relative position of components in each complete ideograph
- Structure of a corresponding component
- Treatment in a source character set
- Radical contained in a component

**Uniqueness.** If one or more of these features are different between the ideographs compared, the ideographs are considered to have different abstract shapes and therefore are considered unique characters and are not unified.

**Unification.** If all of these features are identical between the ideographs, the ideographs are considered to have the same abstract shape and are therefore unified.

The examples in *Table 10-3* represent some typical differences in abstract character shape. The ideographs are therefore *not* unified.

**Table 10-3. Ideographs Not Unified**

| Characters | Reason   |
|------------|--|
| 崖 ≠ 厓      | Different Number of Components   |
| 峰 ≠ 峯      | Same Number of Components Placed in Different Relative Position                                      |
| 擴 ≠ 擴      | Same Number and Same Relative Position of Components, Corresponding Components Structure Differently |
| 區 ≠ 區      | Characters Treated Differently in a Source Character Set   |
| 祕 ≠ 秘      | Characters with Different Radical in a Component   |
| 爲 ≠ 為      | Same Abstract Shape, Difference in Actual Shape  |

Differences in actual shape of ideographs that *have* been unified are illustrated in *Table 10-4*.

**Table 10-4. Ideographs Unified**

| Characters | Reason  |
|------------|---|
| 周 ≈ 周      | Different Writing Sequence                                |
| 雪 ≈ 雪      | Differences in Overshoot at the Stroke Termination        |
| 酉 ≈ 酉      | Differences in Contact of Strokes                         |
| 鉅 ≈ 鉅      | Differences in Protrusion at the Folded Corner of Strokes |
| 堙 ≈ 堙      | Differences in Bent Strokes                               |
| 朱 ≈ 朱      | Differences in Stroke Termination                         |
| 父 ≈ 父      | Differences in Accent at the Stroke Initiation            |
| 八 ≈ 八      | Difference in Rooftop Modification                        |
| 說 ≈ 說      | Difference in Rotated Strokes/Dots <sup>a</sup>           |

- a. These ideographs (having the same abstract shape) would have been unified except for the source separation rule.

**Han Ideograph Arrangement.** The arrangement of the Unicode Han characters is based on the position of characters as they are listed in four major dictionaries. The *KangXi Zidian* was chosen as primary because it contains most of the source characters and because the dictionary itself and the principles of character ordering it employs are commonly used throughout East Asia.

The Han ideograph arrangement follows the index (page and position) of the dictionaries listed in *Table 10-5* with their priorities.

**Table 10-5. Han Ideograph Arrangement**

| Priority | Dictionary              | City    | Publisher                         | Version         |
|----------|-------------------------|---------|-----------------------------------|-----------------|
| 1        | <i>KangXi Zidian</i>    | Beijing | Zhonghua Bookstore, 1989          | 7th edition     |
| 2        | <i>Dai Kan-Wa Jiten</i> | Tokyo   | Taishuukan Shoten, 1986           | Revised edition |
| 3        | <i>Hanyu Da Zidian</i>  | Chengdu | Sichuan Cishu Publishing, 1986    | 1st edition     |
| 4        | <i>Dae Jaweon</i>       | Seoul   | Samseong Publishing Co. Ltd, 1988 | 1st edition     |

When a character is found in the *KangXi Zidian*, it follows the *KangXi Zidian* order. When it is not found in the *KangXi Zidian* and it is found in *Dai Kan-Wa Jiten*, it is given a position extrapolated from the *KangXi* position of the preceding character in *Dai Kan-Wa Jiten*. When it is not found in either *KangXi* or *Dai Kan-Wa*, then the *Hanyu Da Zidian* and *Dae Jaweon* dictionaries are consulted in a similar manner.

Ideographs with simplified *KangXi* radicals are placed in a group following the traditional *KangXi* radical from which the simplified radical is derived. For example, characters with the simplified radical 讠 corresponding to *KangXi* radical 言 follow the last nonsimplified character having 言 as a radical. The arrangement for these simplified characters is that of the *Hanyu Da Zidian*.

The few characters that are not found in any of the four dictionaries are placed following characters with the same *KangXi* radical and stroke count.

The radical-stroke order that results is a culturally neutral order. It does not exactly match the order found in common dictionaries. Information for sorting all CJK ideographs by the radical-stroke method is found on the CD-ROM. It should be used if characters from the CJK Unified Ideographs and CJK Unified Ideographs Extension A blocks and compatibility ideographs are to be properly interleaved.

The form of the charts for the CJK Unified Ideographs block is described in the introduction to *Chapter 14, Code Charts*. A full radical-stroke index is also provided in *Section 15.1, Han Radical-Stroke Index*, to help users locate characters in the main charts.

### Mapping to Standards

The mappings defined by the IRG between the ideographs in the Unicode Standard and the IRG sources are included on the CD-ROM. These mappings are considered to be normative parts of ISO/IEC 10646-1; that is, the characters are *defined* to be the targets for conversion of these characters in these character set standards.

These mappings have as their source editions of the relevant national standards that were provided directly to the IRG by its member bodies. In some cases, these editions differ from the published editions generally available.

The mappings defined by the IRG are considered normative for ISO/IEC 10646-1, and must be considered normative for the Unicode Standard as well. Because they may not match the mappings derived from published editions of the standards, developers of individual mapping tables, which are intended to handle real-life intercharacter set mappings,

may choose to use alternative mappings more directly correlated with published character set editions. The Unicode Consortium provides tables indicating where the variant mappings may be desirable.

Specialized conversion systems may also choose more sophisticated mapping mechanisms—for example, semantic conversion, variant normalization, or conversion between simplified and traditional Chinese.

The Unicode Consortium also provides mapping information that extends beyond the normative mappings defined by the IRG. These additional mappings include mappings to character set standards included in the U source, including duplicate characters from KS C 5601-1987, mappings to portions of character set standards omitted from IRG sources, references to standard dictionaries, and suggested character/stroke counts.

## CJK Compatibility Ideographs: U+F900–U+FAFF

The Korean national standard KS C 5601-1987, which served as one of the primary source sets for the Unified CJK Ideograph Repertoire and Ordering, Version 2.0, contains 268 duplicate encodings of identical ideograph forms to denote alternative pronunciations. That is, in certain cases, the standard encoded a single character multiple times to denote different linguistic uses. This approach is like encoding the letter “a” five times to denote the different pronunciations it has in the words *hat*, *able*, *art*, *father*, and *adrift*. They are in all ways identical in shape to their nominal counterparts, and so were excluded by the IRG from its sources. For round-trip conversion with KS C 5601-1987, they are encoded separately from the primary CJK Unified Ideographs block.

In addition, another 34 ideographs from various regional and industry standards were encoded in this block, primarily to achieve round-trip conversion compatibility. Twelve of these 34 ideographs (U+FA0E, U+FA0F, U+FA11, U+FA13, U+FA14, U+FA1F, U+FA21, U+FA23, U+FA24, U+FA27, U+FA28, and U+FA29) are not encoded in the CJK Unified Ideographs Areas. These 12 characters are not duplicates and should be treated as a small extension to the set of unified ideographs.

## Kanbun: U+3190–U+319F

This block contains a set of Kanbun marks used in Japanese texts to indicate the Japanese reading order of classical Chinese texts. They are not encoded in any current character encoding standards but are widely used in literature. They are typically written in an annotation style to the left of each line of vertically rendered Chinese text.

See also enclosed CJK letters and months (U+3200..U+32FF) and CJK compatibility (U+3300..U+33FF).

## CJK and KangXi Radicals: U+2E80–U+2FD5

East Asian ideographic *radicals* are ideographs or fragments of ideographs used to index dictionaries and word lists, and as the basis for creating new ideographs. The term *radical* comes from the Latin *radix*, meaning *root*, and refers to the part of the character under which the character is classified in dictionaries. *Section 15.1, Han Radical-Stroke Index*, provides information on how to use radical-stroke lookup to locate ideographs encoded in the Unicode Standard.

There is no single radical set in general use throughout East Asia; however, the set of 214 radicals used in the eighteenth-century *KangXi* dictionary is universally recognized.

The visual appearance of radicals is often very different when they are used as radicals from what it is when they are stand-alone ideographs. Indeed, many radicals have multiple graphic forms when used as parts of characters. A standard example is the water radical, which is written 水 when an ideograph and generally 氵 when part of an ideograph.

The Unicode Standard includes two blocks of encoded radicals: The KangXi Radicals block (U+2F00 through U+2FD5), which contains the base forms for the 214 radicals, and the CJK Radicals Extension block (U+2E80 through U+2EF3), which contains a set of variant shapes taken by the radicals either when they occur as parts of characters or are used for simplified Chinese. These variant shapes are commonly found as independent and distinct characters in dictionary indices—such as for the radical-stroke charts in the Unicode Standard. As such, they have not been subject to the usual unification rules used for other characters in the standard.

Most of the radicals in the CJK and KangXi Radicals blocks are also part of the CJK Unified Ideographs block of the Unicode Standard. Radicals that have one graphic form as an ideograph and another as part of an ideograph are generally encoded in both forms in the CJK Unified Ideographs block (such as U+6C34 and U+6C35 for the water radical).

**Standards.** CNS 11643-1992 includes a block of radicals separate from its ideograph block. This block includes 212 of the 214 KangXi radicals. These characters are included in the KangXi Radicals block.

Those radicals that are ideographs in their own right have a definite meaning and are usually referred to by that meaning. Accordingly, most of the characters in the KangXi Radicals block have been assigned names reflecting their meaning. The other radicals have been given names based on their shape.

**Semantics.** Characters in the CJK and KangXi Radicals blocks should never be used as ideographs. They have different properties and meaning. U+2F00 KANGXI RADICAL ONE IS not equivalent to U+4E00 CJK UNIFIED IDEOGRAPH 4E00. The former is to be treated as a symbol, the latter as a word or part of a word.

It is necessary to make a semantic distinction between a character used as an ideograph and the same character used as a radical. To emphasize this difference, radicals may be given a distinct font style from their ideographic counterparts.

## Ideographic Description: U+2FF0–U+2FFB

Although the Unicode Standard includes over 27,000 ideographs, many thousands of extremely rare ideographs were nevertheless left unencoded. As an example, nearly half of the characters in the *KangXi* dictionary are unencoded. Research into cataloging additional ideographs for encoding continues, but it is anticipated that at no point will the entire set of potential, encodable ideographs be completely exhausted. In particular, ideographs continue to be coined and such new coinages will invariably be unencoded.

The 12 characters in the Ideographic Description block provide a mechanism for the standard interchange of text that must reference unencoded ideographs. Unencoded ideographs can be described using these characters and encoded ideographs; the reader can then create a mental picture of the ideographs from the description.

This process is different from a formal *encoding* of an ideograph. There is no canonical description of unencoded ideographs; there is no semantic assigned to described ideographs; there is no equivalence defined for described ideographs. Conceptually, ideograph

descriptions are more akin to the English phrase, “an ‘e’ with an acute accent on it,” than to the character sequence “U+006E U+0301.”

In particular, support for the characters in the Ideographic Description block does *not* require the rendering engine to recreate the graphic appearance of the described character.

Note also that many of the ideographs that users might represent using the Ideographic Description characters will be formally encoded in future versions of the Unicode Standard.

The Ideographic Description algorithm depends on the fact that virtually all CJK ideographs can be broken down into smaller pieces which are themselves ideographs. The broad coverage of the ideographs already encoded in the Unicode Standard implies that the vast majority of unencoded ideographs can be represented using the Ideographic Description characters.

***Ideographic Description Sequences.*** Ideographic Description Sequences are defined by the following grammar. See *Section 0.2, Notational Conventions*, for the notational conventions used here.

*IDS ::= UnifiedIdeograph | Radical | BinaryDescriptionOperator IDS IDS  
| TertiaryDescriptionOperator IDS IDS IDS*

*BinaryDescriptionOperator ::= U+2FF0 | U+2FF1 | U+2FF4 | U+2FF5 | U+2FF6 | U+2FF7  
| U+2FF8 | U+2FF9 | U+2FFA | U+2FFB*

*TertiaryDescriptionOperator ::= U+2FF2 | U+2FF3*

*Radical ::= U+2E80 | U+2E81 | ... | U+2EF2 | U+2EF3 | U+2F00 | U+2F01 | ... | U+2FD4  
| U+2FD5*

*UnifiedIdeograph ::= U+3400 | U+3401 | ... | U+4DB4 | U+4DB5 | U+4E00 | U+4E01 | ...  
| U+9FA4 | U+9FA5 | U+FA0E | U+FA0F | U+FA11 | U+FA13 | U+FA14  
| U+FA1F | U+FA21 | U+FA23 | U+FA24 | U+FA27 | U+FA28 | U+FA29*

In addition to the above grammar, Ideographic Description Sequences have two additional length constraints:

- No sequence can be longer than 16 Unicode scalar values in length.
- No sequence can contain more than six *UnifiedIdeographs* in a row without an intervening Ideographic Description character.

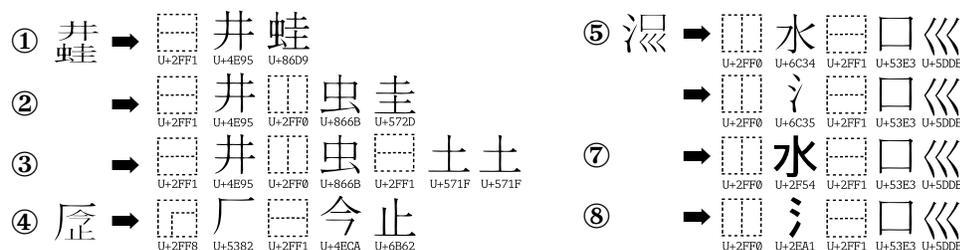
A sequence of characters that includes Ideographic Description characters but does not conform to the above grammar and length constraints is not an Ideographic Description Sequence.

The operators indicate the relative graphic positions of the operands running from left to right and from top to bottom.

Note that non-unique compatibility ideographs (U+F900 through U+FA2D, but not U+FA0E, U+FA0F, U+FA11, U+FA13, U+FA14, U+FA1F, U+FA21, U+FA23, U+FA24, U+FA27, U+FA28, or U+FA29) are not counted as unified ideographs for the purposes of this grammar, although they do have the ideographic property (see *Section 4.10, Letters and Other Useful Properties*). Non-unique compatibility ideographs are excluded from Ideographic Description Sequences to incrementally reduce the ambiguity of such sequences.

*Figure 10-8* illustrates the use of this grammar to provide descriptions of unencoded ideographs.

A user wishing to represent an unencoded ideograph will need to analyze its structure to determine how to describe it using an Ideographic Description Sequence. As a rule, it is best to use the natural radical-phonetic division for an ideograph if it has one and to use as

**Figure 10-8. Using the Ideographic Description Characters**

short a description sequence as possible, but there is no requirement that these rules be followed. Beyond that, the shortest possible Ideographic Description Sequence is preferred.

The length constraints allow random access into a string of ideographs to have well-defined limits. Only a small number of characters need to be scanned backward to determine whether those characters are part of an Ideographic Description Sequence.

The fact that Ideographic Description Sequences can contain other Ideographic Description Sequences means that implementations may need to be aware of the *recursion depth* of a sequence and its *back-scan length*. The recursion depth of an Ideographic Description Sequence is the maximum number of pending operations encountered in the process of parsing an Ideographic Description Sequence. In *Figure 10-8*, the maximum recursion depth is shown in the third example, where three operations are still pending at the end of the Ideographic Description Sequence.

The back-scan length is the maximum number of ideographs unbroken by Ideographic Description characters in the sequence. None of the examples in *Figure 10-8* has more than two ideographs in a row; for all, the back-scan depth is two.

The Unicode Standard places no formal limits on the recursion depth of Ideographic Description Sequences. It does, however, limit the back-scan depth for valid Ideographic Description Sequences to be six or less.

**Equivalence.** Many unencoded ideographs can be described in more than one way using this algorithm, either because the pieces of a description can themselves be broken down further (examples one through three in *Figure 10-8*), or because duplications appear within the Unicode Standard (examples five and six in *Figure 10-8*).

The Unicode Standard does not define equivalence for two Ideographic Description Sequences that are not identical. *Figure 10-8* contains numerous examples illustrating how different Ideographic Description Sequences might be used to describe the same ideograph.

In particular, Ideographic Description Sequences are not to be used to provide alternative graphic representations of encoded ideographs. Searching, collation, and other content-based text operations would then fail.

**Interaction with the Ideographic Variation Mark.** As with ideographs proper, the Ideographic Variation Mark (U+303E) may be placed before an Ideographic Description Sequence to indicate that the description is only an approximation of the original ideograph desired. A sequence of characters that includes an Ideographic Variation Mark is not an Ideographic Description Sequence.

**Rendering.** Ideographic Description characters are visible characters. They are not to be treated as control characters. The sequence U+2FF1 U+4E95 U+86D9 must have a distinct appearance from U+4E95 U+86D9.

An implementation may render a valid Ideographic Description Sequence either by rendering the individual characters separately, or by parsing the Ideographic Description Sequence and drawing the ideograph so described. In the latter case, the Ideographic Description Sequence should be treated as a ligature of the individual characters for purposes of hit testing, cursor movement, and other user interface operations. (See *Section 5.12, Editing and Selection.*)

**Character Boundaries.** Ideographic Description characters are not combining characters, and there is no requirement that they affect character or word boundaries. Thus U+2FF1 U+4E95 U+86D9 may be treated as a sequence of three characters or even three words.

Implementations of the Unicode Standard may choose to parse Ideographic Description Sequences when calculating word and character boundaries, but such a decision will make the algorithms involved significantly more complicated and slower.

**Standards.** The Ideographic Description characters are found in GBK—an extension to GB 2312-80 that adds all Unicode ideographs not already in GB 2312-80. GBK is defined as a normative annex of GB 13000.1-93.

---

## 10.2 Hiragana

### Hiragana: U+3040–U+309F

Hiragana is the cursive syllabary used to write Japanese words phonetically and to write sentence particles and inflectional endings. It is also commonly used to indicate the pronunciation of Japanese words. Hiragana syllables are phonetically equivalent to corresponding Katakana syllables.

**Standards.** The Hiragana block is based on the JIS X 0208-1990 standard, extended by the nonstandard syllable U+3094 VU, which is included in some Japanese corporate standards.

**Combining Marks.** Hiragana and the related script Katakana use U+3099 COMBINING KATAKANA-HIRAGANA VOICED SOUND MARK and U+309A COMBINING KATAKANA-HIRAGANA SEMI-VOICED SOUND MARK to generate voiced and semi-voiced syllables from the base syllables, respectively. All common precomposed combinations of base syllable forms using these marks are already encoded as characters, and use of these precomposed forms is the predominant JIS usage. These combining marks must follow the base character to which they apply. As most implementations and JIS standard treat these marks as spacing characters, the Unicode Standard also contains two corresponding noncombining (spacing) marks at U+309B and U+309C.

**Iteration Marks.** The two characters U+309D HIRAGANA ITERATION MARK and U+309E HIRAGANA VOICED ITERATION MARK are punctuation-like characters that denote the iteration (repetition) of a previous syllable according to whether the repeated syllable has an unvoiced or voiced consonant, respectively.

---

## 10.3 Katakana

### Katakana: U+30A0–U+30FF

Katakana is the noncursive syllabary used to write non-Japanese (usually Western) words phonetically in Japanese. It is also used to write Japanese words with visual emphasis. Katakana syllables are phonetically equivalent to corresponding Hiragana syllables. Katakana contains two characters, U+30F5 KATAKANA LETTER SMALL KA and U+30F6 KATAKANA LETTER SMALL KE, that have no direct correspondent in Hiragana; they are used in special Japanese spelling conventions (for example, the spelling of placenames that include archaic Japanese connective particles).

**Standards.** The Katakana block is based on the JIS X 0208-1990 standard.

**Punctuation-like Characters.** U+30FB KATAKANA MIDDLE DOT is used to separate words when writing non-Japanese phrases. U+30FC KATAKANA-HIRAGANA PROLONGED SOUND MARK is used predominantly with Katakana and occasionally with Hiragana to denote a lengthened vowel of the previously written syllable. The two iteration marks, U+30FD KATAKANA ITERATION MARK and U+30FE KATAKANA VOICED ITERATION MARK, serve the same function in Katakana writing that the two Hiragana iteration marks serve in Hiragana writing.

### Halfwidth and Fullwidth Forms: U+FF00–U+FFEF

In the context of East Asian coding systems, a double-byte character set (DBCS) such as JIS X 0208-1990 or KS C 5601-1987 is generally used together with a single-byte character set (SBCS), such as ASCII or a variant of ASCII. Text that is encoded with both a DBCS and SBCS is typically displayed such that the glyphs representing DBCS characters occupy two display cells where a display cell is defined in terms of the glyphs used to display the SBCS (ASCII) characters. In these systems, the two-display-cell width is known as the *fullwidth* or *zenkaku* form, and the one-display-cell width is known as the *halfwidth* or *hankaku* form.

Because of this mixture of display widths, certain characters often appear twice, once in fullwidth form in the DBCS repertoire and once in halfwidth form in the SBCS repertoire. To achieve round-trip conversion compatibility with such mixed encoding systems, it is necessary to encode both fullwidth and halfwidth forms of certain characters. This block consists of the additional forms needed to support conversion for existing texts that employ both forms.

In the context of conversion to and from such mixed width encodings, all characters in the General Scripts Area should be construed as halfwidth (*hankaku*) characters if they have a fullwidth equivalent elsewhere in the standard or if they do not occur in the mixed width encoding; otherwise, they should be construed as fullwidth (*zenkaku*). Specifically, most characters in the CJK Phonetics and Symbols Area and the Unified CJK Ideograph Area, along with the characters in the CJK Compatibility Ideographs, CJK Compatibility Forms, and Small Form Variants blocks, should be construed as fullwidth (*zenkaku*) characters. For a complete description of the East Asian Width property, see Unicode Technical Report #11, “East Asian Width,” on the CD-ROM or the up-to-date version on the Unicode Web site.

The characters in this block consist of fullwidth forms of the ASCII block (except SPACE), certain characters of the Latin-1 Supplement, and some currency symbols. In addition, this

block contains halfwidth forms of the Katakana and Hangul Compatibility Jamo characters. Finally, a number of characters from the Symbols Area are replicated here (U+FFE8..U+FFEE) with explicit halfwidth semantics.

As with other compatibility characters, the preferred Unicode encoding is to use the nominal counterparts of these characters and use rich text font or style bindings to select the appropriate glyph size and width.

**Unifications.** The fullwidth form of U+0020 SPACE is unified with U+3000 IDEOGRAPHIC SPACE.

## 10.4 Hangul

### Hangul Jamo: U+1100–U+11FF

Korean Hangul may be considered to be a syllabic script. As opposed to many other syllabic scripts, the syllables are formed from a set of alphabetic components in a regular fashion. These alphabetic components are called *jamo*.

The Unicode Standard contains both the complete set of precomposed modern Hangul syllable blocks and the set of conjoining Hangul jamo in this block. This set of conjoining Hangul jamo can be used to encode all modern and ancient syllable blocks. For a description of conjoining jamo behavior and precomposed Hangul Syllables, see *Section 3.11, Conjoining Jamo Behavior*, and the Hangul Syllables character block description (U+AC00..U+D7A3).

The Hangul jamo are divided into three classes: *choseong* (leading consonants, or syllable-initial characters), *jungseong* (vowels, or syllable-peak characters), and *jongseong* (trailing consonants, or syllable-final characters). In the following discussion, these classes are abbreviated as *L* (leading consonant), *V* (vowel), and *T* (trailing consonant).

For use in composition, two invisible filler characters act as placeholders for choseong or jungseong: U+115F HANGUL CHOSEONG FILLER and U+1160 HANGUL JUNGSEONG FILLER.

**Collation.** The unit of collation in Korean text is normally the Hangul syllable block. Because of the arrangement of the conjoining jamo, their sequences may be collated with a binary comparison. For example, in comparing (a) *LVTLV* with (b) *LVLV*, the first syllable block (*LVT*) should be compared with the second (*LV*). Supposing the first two characters are identical—because all trailing consonants have binary values greater than all leading consonants—the *T* would compare as greater than the second *L* in (b). This result produces the correct ordering between the strings. The positions of the fillers in the code charts were also chosen with this condition in mind.

- As with any coded characters, collation cannot depend simply on a binary comparison. Odd sequences such as superfluous fillers will produce an incorrect sort, as will cases where a non-jamo character follows a sequence (such as comparing *LVT* with *LVx*, where *x* is a Unicode character above U+11FF, such as U+3000 IDEOGRAPHIC SPACE).

If mixtures of precomposed syllable blocks and jamo are collated, the easiest approach is to decompose the precomposed syllable blocks into conjoining jamo before comparing.

### Hangul Compatibility Jamo: U+3130–U+318F

This block consists of spacing, nonconjoining Hangul consonant and vowel (jamo) elements. These characters are provided solely for compatibility with the KS C 5601 standard. Unlike the characters found in the Hangul Jamo block (U+1100..U+11FF), the jamo characters in this block have no conjoining semantics.

The characters of this block are considered to be fullwidth forms in contrast with the halfwidth Hangul Compatibility Jamo found at U+FFA0..U+FFDF.

**Standards.** The Unicode Standard follows KS C 5601 for Hangul Jamo elements.

## Hangul Syllables: U+AC00–U+D7A3

The Hangul script used in the Korean writing system consists of individual consonant and vowel letters (jamo) that are visually combined into square display cells to form entire syllable blocks. Hangul syllables may be encoded directly as precomposed combinations of individual jamo or as decomposed sequences of conjoining jamo. The latter encoding is supported by the Hangul Jamo block (U+1100..U+11FF). The syllabic encoding method is described here.

Modern Hangul syllable blocks can be expressed with either two or three jamo, either in the form *consonant + vowel* or in the form *consonant + vowel + consonant*. There are 19 possible leading (initial) consonants (choseong), 21 vowels (jungeong), and 27 trailing (final) consonants (jongseong). Thus there are 399 possible two-jamo syllable blocks and 10,773 possible three-jamo syllable blocks, for a total of 11,172 modern Hangul syllable blocks. This collection of 11,172 modern Hangul syllables encoded in this block is known as the *Johab* set.

**Standards.** The Hangul syllables are taken from KS C 5601-1992, representing the full Johab set. This group represents a superset of the Hangul syllables encoded in earlier versions of Korean standards (KS C 5601-1987, KS C 5657-1991).

**Equivalence.** Each of the Hangul syllables encoded in this block may be encoded by an equivalent sequence of conjoining jamo; however, the converse is not true because thousands of archaic Hangul syllables may be encoded only as a sequence of conjoining jamo. Implementations that use a conjoining jamo encoding are able to represent these archaic Hangul syllables.

**Hangul Syllable Composition.** The Hangul syllables can be derived from conjoining jamo by a regular process of composition. The algorithm that maps a sequence of conjoining jamo to the encoding point for a Hangul syllable in the Johab set is detailed in *Section 3.11, Conjoining Jamo Behavior*.

**Hangul Syllable Decomposition.** Conversely, any Hangul syllable from the Johab set can be decomposed into a sequence of conjoining jamo characters. The algorithm that details the formula for decomposition is provided in *Section 3.11, Conjoining Jamo Behavior*.

**Hangul Syllable Name.** The character names for Hangul syllables are derived algorithmically from the decomposition. (For full details, see *Section 3.11, Conjoining Jamo Behavior*.)

**Hangul Syllable Representative Glyph.** The representative glyph for a Hangul syllable can be formed from its decomposition based on the categorization of vowels shown in *Table 10-6*.

**Table 10-6. Line-Based Placement of Jungseong**

| Vertical |     | Horizontal |    | Horizontal and Vertical |     |
|----------|-----|------------|----|-------------------------|-----|
| 1161     | A   | 1169       | O  | 116A                    | WA  |
| 1162     | AE  | 116D       | YO | 116B                    | WAE |
| 1163     | YA  | 116E       | U  | 116C                    | OE  |
| 1164     | YAE | 1172       | YU | 116F                    | WEO |
| 1165     | EO  | 1173       | EU | 1170                    | WE  |
| 1166     | E   |            |    | 1171                    | WI  |
| 1167     | YEO |            |    | 1174                    | YI  |
| 1168     | YE  |            |    |                         |     |
| 1175     | I   |            |    |                         |     |

If the vowel of the syllable is based on a vertical line, place the leading consonant to its left. If the vowel is based on a horizontal line, place the preceding consonant above it. If the vowel is based on a combination of vertical and horizontal lines, place the preceding consonant above the horizontal line and to the left of the vertical line. In either case, place a following consonant, if any, below the middle of the resulting group.

In any particular font, the exact placement, shape, and size of the components will vary according to the shapes of the other characters and the overall design of the font.

See also enclosed CJK letters and months (U+3200..U+32FF), CJK compatibility (U+3300..U+33FF), and halfwidth and fullwidth forms (U+FF00..U+FFEF).

## 10.5 Bopomofo

### Bopomofo: U+3100–U+312F

*Bopomofo* constitute a set of characters used to annotate or teach the phonetics of Chinese, primarily the standard Mandarin language. The characters are used in dictionaries and teaching materials, but not in the actual writing of Chinese text. The formal Chinese names for this alphabet are *Zhuyin-Zimu* (“phonetic alphabet”) and *Zhuyin-Fuhao* (“phonetic symbols”), but the informal term “Bopomofo” (analogous to “ABCs”) provides a more serviceable English name and is also used in China. The Bopomofo were developed as part of a populist literacy campaign following the 1911 revolution; thus they are acceptable to all branches of modern Chinese culture, although in the People’s Republic of China their function has been largely taken over by the Pinyin romanization system.

**Standards.** The standard Mandarin set of Bopomofo is included in the People’s Republic of China standard GB 2312 and in the Republic of China (Taiwan) standard CNS 11643.

**Mandarin Tone Marks.** Small modifier letters used to indicate the five Mandarin tones are part of the Bopomofo system. In the Unicode Standard they have been unified into the Modifier Letter range, as shown in *Table 10-7*.

**Table 10-7. Mandarin Tone Marks**

|             |        |                              |
|-------------|--------|------------------------------|
| first tone  | U+02C9 | MODIFIER LETTER MACRON       |
| second tone | U+02CA | MODIFIER LETTER ACUTE ACCENT |
| third tone  | U+02C7 | CARON                        |
| fourth tone | U+02CB | MODIFIER LETTER GRAVE ACCENT |
| light tone  | U+02D9 | DOT ABOVE                    |

**Standard Mandarin Bopomofo.** The order of the Mandarin Bopomofo letters U+3105..U+3129 is standard worldwide. The code offset of the first letter U+3105 BOPOMOFO LETTER B from a multiple of 16 is included to match the offset in the ISO-registered standard GB 2312. The character U+3127 BOPOMOFO LETTER I is usually written as a vertical stroke when Bopomofo text is set vertically. In the Unicode Standard, this representation is considered to be a rendering variation; the variant is not assigned a separate character code.

**Extended Bopomofo.** To represent the sounds of Chinese dialects other than Mandarin, the basic Bopomofo set U+3105..U+3129 has been augmented by additional phonetic characters. These extensions are much less broadly recognized than the basic Mandarin set. The three extended Bopomofo characters U+312A..U+312C are cited in some standard reference works, such as the encyclopedia *Xin Ci Hai*. Another set of 24 extended Bopomofo, encoded at U+31A0..U+31B7, was designed in 1948 to cover additional sounds of the Minnan and Hakka dialects. The extensions are used together with the main set of Bopomofo characters to provide a complete phonetic orthography for those dialects. There are no standard Bopomofo letters for the phonetics of Cantonese or several other Southern Chinese dialects.

The small characters encoded at U+31B4..U+31B7 represent syllable-final consonants not present in standard Mandarin or Mandarin dialects. They have the same shapes as Bopomofo “p”, “t”, “k”, and “h”, respectively, but are rendered smaller than the initial consonants; they are also generally shown close to the syllable medial vowel character. These final letters are encoded separately so that Minnan and Hakka dialects can be represented unambiguously in plain text without having to resort to subscripting or other fancy text mechanisms to represent the final consonants.

**Extended Bopomofo Tone Marks.** In addition to the Mandarin tone marks enumerated in Table 10-7, additional tone marks appropriate for use with the extended Bopomofo transcriptions of Minnan and Hakka can be found in the Modifier Letter range, as shown in Table 10-8. The “departing tone” refers to the *qusheng* in traditional Chinese tonal analysis, with the *yin* variant historically derived from voiceless initials and the *yang* variant from voiced initials. Southern Chinese dialects in general maintain more tonal distinctions than Mandarin.

**Table 10-8. Minnan and Hakka Tone Marks**

|                     |        |                          |
|---------------------|--------|--------------------------|
| yin departing tone  | U+02EA | YIN DEPARTING TONE MARK  |
| yang departing tone | U+02EB | YANG DEPARTING TONE MARK |

**Rendering of Bopomofo.** Bopomofo is rendered left to right in horizontal text, but also commonly appears in vertical text. It may be used by itself in either orientation, but most commonly appears in interlinear annotation of Chinese (Han character) text. It is not uncommon for children’s books to be completely annotated with Bopomofo pronunciations for every character. This interlinear annotation is structurally quite similar to the system of Japanese *ruby* annotation, but it has additional complications that result from the explicit usage of tone marks with the Bopomofo letters.

In horizontal interlineation, the Bopomofo is generally placed above the corresponding Han character(s); tone marks, if present, appear at the end of each syllabic group of Bopomofo letters. In vertical interlineation, the Bopomofo is generally placed on the right side of the corresponding Han character(s); tone marks, if present, appear in a separate interlinear row to the right side of the vowel letter. When using extended Bopomofo for Minnan and Hakka, the tone marks may also be mixed with Latin digits 0–9 to express changes in actual tonetic values resulting from juxtaposition of basic tones.

---

## 10.6 Yi

### Yi: U+A000–U+A4CF

The Yi syllabary is used to write the Yi language, a member of the Sino-Tibetan language family. The script is also known as Cuan or Wei.

The Yi, also known as Lolo and Nuo-su, are one of the largest non-Han minorities in the People's Republic of China (PRC). Most live in southwestern China, but others live in Myanmar, Laos, and Vietnam. Yi is one of the official languages of PRC.

The earliest surviving samples of classical Yi, an ideographic script, date from about 500 years ago. Unlike other Sinoform scripts, the ideographs themselves appear not to be derived from Han ideographs. There are some 8,000 to 10,000 characters in the classical Yi script, although the exact ideographs used varied from region to region.

To improve literacy in Yi, the Yi syllabary was introduced in the 1970s. This syllabary is encoded in the Unicode Standard; the classical ideographic Yi script is not encoded at this time.

Each Yi syllable consists of a consonantal initial, a final, and a tone. The core Yi syllabary consists of 819 signs for syllables with the first three tones (high, low, and middle low), plus a mark added to the form for the middle low tone to indicate a fourth tone (middle high).

**Standards.** In 1991, a national standard for Yi was adopted by the PRC as GB 13134-91. This encoding includes all 1,165 Yi syllables and is the basis for the encoding used by the Unicode Standard, which includes all 1,165 Yi syllables.

**Naming Conventions and Order.** The Yi syllables are named on the basis of their romanized sound values. The tone is indicated by appending a letter to the romanization: “t” for the high tone, “p” for the low tone, “x” for the middle high tone, and no letter for the middle low tone.

**Rendering.** Yi follows the writing rules for Han ideographs. Characters are generally written left to right or occasionally top to bottom. There is no typographic interaction between individual characters of the Yi script.

**Yi Radicals.** To facilitate the lookup of Yi characters in dictionaries, a set of radicals has been invented. The Yi repertoire is divided into several subsets, each of which shares a common stroke (radical). The name used for the radical is that of the corresponding Yi character closest to it in shape, with a “b” added as a suffix.

This PDF file is an excerpt from *The Unicode Standard, Version 3.0*, 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 (see <http://www.unicode.org/unicode/uni2errata/UnicodeErrata.html>). More recent versions of the Unicode standard exist (see <http://www.unicode.org/unicode/standard/versions/>).

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 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.

ISBN 0-201-61633-5

Copyright © 1991-2000 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.

This book is set in Minion, designed by Rob Slimbach at Adobe Systems, Inc. It was typeset using FrameMaker 5.5 running under Windows NT. ASMUS, Inc. created custom software for chart layout. The Han radical-stroke index was typeset by Apple Computer, Inc. The following companies and organizations supplied fonts:

Apple Computer, Inc.  
Atelier Fluxus Virus  
Beijing Zhong Yi (Zheng Code) Electronics Company  
DecoType, Inc.  
IBM Corporation  
Monotype Typography, Inc.  
Microsoft Corporation  
Peking University Founder Group Corporation  
Production First Software

Additional fonts were supplied by individuals as listed in the *Acknowledgments*.

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.

All other company and product names are trademarks or registered trademarks of the company or manufacturer, respectively.

The publisher offers discounts on this book when ordered in quantity for special sales. For more information please contact:

Corporate, Government, and Special Sales  
Addison Wesley Longman, Inc.  
One Jacob Way  
Reading, Massachusetts 01867

Visit A-W on the Web: <http://www.awl.com/cseng/>

First printing, January 2000.