

Universal Multiple-Octet Coded Character Set
International Organization for Standardization
Organisation Internationale de Normalisation
Международная организация по стандартизации

Doc Type: Ideographic Rapporteur Group Document
Title: Request to clarify some IDS issues
Source: Eiso Chan (陈永聪, Culture and Art Publishing House)
Status: Individual Contribution to IRG #53, Shenzhen, PRC
Action: For consideration by IRG
Date: 2019-10-24

In the IRG WS works and the Jianzi encoding works, there are some unclear issues related to IDS, so I request IRG to clarify them.

1. Scope of Character Description Component (CDC)

In IRG PnP, Glossary, it shows “It includes all coded CJK unified ideographs, Kangxi Radicals, CJK Radical Supplements, and coded CJK Compatibility ideographs.” Coded CJK unified ideographs mean all encoded characters in URO and all CJK extensions, Kangxi Radicals mean all encoded characters in *Kangxi Radicals* block (U+2F00 to U+2FDF), CJK Radical Supplements mean all encoded characters in *CJK Radical Supplement* block (U+2E80 to U+2EFF) and coded CJK Compatibility ideographs mean all encoded characters in *CJK Compatibility Ideographs* block and *CJK Compatibility Ideographs Supplement* block.

When we check the script property values in UCD for the characters mentioned above as below, we will find that they are related to Han.

Characters	Script property values
CJK Unified Ideographs	Han
Kangxi Radicals	Han
CJK Radical Supplements	Han
CJK Compatibility Ideographs	Han

However, there are also other Han characters in CJK Symbols and Punctuation block, they are the Ideographic Iteration Marks, Ideographic Number and the Suzhou Numerals shown as below.

UCS	Character	Character Name
U+3005	々	IDEOGRAPHIC ITERATION MARK
U+3007	〇	IDEOGRAPHIC NUMBER ZERO
U+3021	丨	HANGZHOU NUMERAL ONE
U+3022	𠄎	HANGZHOU NUMERAL TWO
U+3023	𠄎	HANGZHOU NUMERAL THREE

- <Binary_Symbol> ::= ☐ | ☑ | ☒ | ☓ | ☔ | ☕ | ☖ | ☗ | ☘ | ☙
- <Ternary_Symbol> ::= ☚ | ☛

If IRG has confirmed the CJK strokes, fullwidth question mark, the CJK supplementary components or the Suzhou numerals are not allowed for IDS in the IRG encoding works, it's better to clarify in IRG PnP, which will be helpful for the submitters and reviewers.

The following IDS syntax is cited from the latest version of Unicode Core Specification.

Ideographic Description Sequences. Ideographic Description Sequences are defined by the following grammar. The list of characters associated with the Ideographic and Radical properties can be found in the Unicode Character Database. In particular, the Ideographic property is intended to apply to other siniform ideographic systems, in addition to CJK ideographs. [...]

IDS := *Ideographic* | *Radical* | *CJK_Stroke* | *Private Use* | U+FF1F

| *IDS_BinaryOperator* *IDS* *IDS*

| *IDS_TernaryOperator* *IDS* *IDS* *IDS*

CJK_Stroke := U+31C0 | U+31C1 | ... | U+31E3

IDS_BinaryOperator := U+2FF0 | U+2FF1 | U+2FF4 | U+2FF5 | U+2FF6 | U+2FF7 |

U+2FF8 | U+2FF9 | U+2FFA | U+2FFB

IDS_TernaryOperator := U+2FF2 | U+2FF3

In the Unicode IDS syntax, *Ideographic* means all the characters which their scripts are belong to *Han*.

2. Encoding Model for Jianzi

I, the source showed Culture and Art Publishing House, submitted WG2 N5041 to request to encode Jianzi in Unicode and ISO/IEC 10646 in future, and now this work is analyzing and ongoing. Prof. Tang, the source showed China National Database (aka Chinese Characters Repertoire), pointed out there are something not better in the encoding model based on musical notation meanings. I read Prof. Tang's document and discussed with the Guqin experts and designers in my team, we thank the comment, but we still don't think the encoding model based on IDS provided by Prof. Tang is well.

2.1. Unstable Sequences

After WG2 #68, Mr. Chen Zhuang, Prof. Tang Yingmin and I met at PKU to discuss the Jianzi encoding model. I pointed out why the IDS method is not better.

In Unicode core specification, it shows the important description as below. (p. 725 in Unicode, 12.0.0)

The Unicode Standard does not define equivalence for two Ideographic Description Sequences that are not identical. *Figure 18-9* contains numerous examples illustrating how different Ideographic Description Sequences might be used to describe the same ideograph. In particular, Ideographic Description Sequences should not be used to provide alternative graphic representations of encoded ideographs in data interchange. Searching, collation, and other content-based text operations would then fail.

In IRG PnP, Annex B, it shows both of 𠃉 傾 and 𠃉 化頁 should be acceptable for the IDS for 傾 (U+50BE) and so on.

For the big-sized characters in the Jianzi system, there are also the similar situations. For example, the IDS for 𦉳 (散勾三弦) could be used 𦉳 𦉳 𦉳 𦉳 𦉳 or 𦉳 𦉳 𦉳 𦉳 𦉳 as their IDS. The first one is easy for the end user if he or she doesn't know Guqin, but more Guqin players will

approximation of the ideograph desired. The IVI is not considered a part of the Ideographic Description Sequence and does not invalidate the sequence.

I request for consideration to add this paragraph in IRG PnP, Glossary as well.

4. Additional Issue

Yifan pointed out the so-called “para-ideographs” in L2/19-346 to the Script Ad-hoc Group and UTC, but there is no comment in the Script Ad-hoc Group meeting in Sept. 27th, 2019, see L2/19-343.

We reviewed this document which raises a number of points about the Gongche characters and other “para-ideographic” characters. Because the Script Ad Hoc reviews proposals for non-CJK characters, we recommend the document be submitted to the UTC and to IRG for discussion.

As the original submitter of the Gongche proposal, I oppose to move the seven Gongche characters out of the end of CJK Ext. B to destabilize Unicode, Version 13.0, but it’s necessary for all the IRG expert to consider how to encode other so-called “para-ideographs” in Yifan’s document.

I understand the “para-ideographs” means the hybrid characters and the characters used in the ideograph running text but the stroke shapes are not like the common ones. In the list provided by Yifan, the first character in Table 2 could be unified with U+211A0 (𠄠), which has been added to BabelStoneHan via `ss01` and `cv01`; the first character in Table 5 could be unified with U+2CF36 (𠄠) because of cognate; the second character in Table 5 is the cursive form of U+2E574 (𠄠).

Andrew West and I also collected more hybrid character like Yifan said. Andrew added them to his famous font BabelStoneHan in PUA.

PUA	Char.	Pseudo-IDS	Note
U+F3E2	𠄠	𠄠タカ	Kanji-Katakana hybrid = 鷹 taka in the Tokyo place name 墨砧 ~ = ボク きぬた たか (Boku kinuta taka).
U+F8C0	𠄠	𠄠木 K	Sawndip-Latin hybrid = ge (?) "pine"
U+F8C1	𠄠	𠄠纒 K	Sawndip-Latin hybrid = gej "to untie" (= 擲)
U+F8C2	𠄠	𠄠老 K	Sawndip-Latin hybrid = geq "old" (= 齧)
U+F8C3	𠄠	𠄠疒 A	Sawndip-Latin hybrid = ae "cough" (= 痲)
U+F8C4	𠄠	𠄠丫 E	Sawndip-Latin hybrid = ngez (?) "branch"
U+F8C5	𠄠	𠄠身 N	Sawndip-Latin hybrid = enj "to stick out one's chest or stomach" (= 𠄠益先 GZ-1482101 in Ext. G)
U+F8C6	𠄠	𠄠艹 M	Sawndip-Latin hybrid = em "kind of grass" [芭芒] (= 菩)
U+F8C7	𠄠	𠄠口 ㇇	Hanja-Hangul hybrid = 圖 do "map"

PUA	Char.	Pseudo-IDS	Note
U+F3E2	鷹	𠩺广𠩺タカ	Kanji-Katakana hybrid = 鷹 taka in the Tokyo place name 墨帖～ = ボク きぬた たか (Boku kinuta taka).
U+F8C8	腐	𠩺广ふ	Kanji-Hiragana hybrid = 腐 in the word 豆腐 (とふ tofu)
U+F8C9	𠩺	𠩺广 K	Kanji-Latin hybrid = 慶 in the name of Keiō University (慶應大学)
U+F8CA	𠩺	𠩺广 O	Kanji-Latin hybrid = 應/応 in the name of Keiō University (慶應大学)
U+F8CB	𠩺	𠩺广𠩺KO	Kanji-Latin hybrid = 慶應/慶応 in the name of Keiō University (慶應大学)

In the current CJK encoding method, the hybrid characters mentioned above can't be encoded in CJKUI, but it's not better to encode in another new block. It's hard to use IDS to describe based on the current IDS syntax.

I request the IRG experts to consider how to handle these hybrid characters.

5. Acknowledgement

Thanks for the feedback comments from Andrew West (魏安), Lee Collins (康立論), Dr. Ken Lunde (小林劍), William Nelson, John Knightley (李忠仕), Kushim Jiang (姜兆勤) and Henry Chan (陳輝恒).

The new definition of IVI is written by Ken Whistler.

(End of Document)