ISO/IEC JTC1/SC2/WG2/IRG

Ideographic Rapporteur Group

(IRG)

Source/Contribution Identifier :    Hong Kong Special Administrative Region, China

Meeting :    IRG Meeting #49 in San Jose, US

Title :    Feedback to IRGN2234 on IRG Standing Document Summary Version 2

Status :    Member feedback for discussion

The HKSAR has reviewed IRGN2234 on IRG Standing Document Summary Version 2.  Please refer to the parts highlighted in green in the marked-up draft attached to this document for details of proposed refinements.

END OF DOCUMENT

Universal Multiple-Octet Coded Character Set
UCS

ISO/IEC JTC 1/SC 2/WG 2/IRGN2234 Draft
Date: 2017-06-26

| | |
|---|---|
| Source: | IRG |
| Meeting: | After IRG Meeting #48, Seoul, Korea |
| Title: | IRG Standing Document Summary Version 2. |
| References | IRGN1648 |
| Status: | Member's submission |
| Actions required: | Feedback Requested |
| Distribution: | IRG |
| Medium: | Electronic |
| Pages: | 6 |

## Introduction:

IRG Working Document Series (IWDS) is a set of IRG maintained documents to keep up-to-date examples of CJK unification related example cases to supplement the published Annex S of ISO/IEC 10646 for IRG unification work. IRG also decided in its meeting #48 that a list of submitters' Printed Character Normalization Guidelines are to be kept in IWDS for keeping track of the transformation rules used for handwritten character to printed character conversions.

The maintenance of the IRG Working Document Series should comply with the operational procedures established in Annex E of the *IRG Principles and Procedures*.

The Standing Document Series consists of the following documents:

Series 1: Summary of unification rules and sample examples.
Series 2: List of UCV (Unifiable Component Variations) (UCV) of Ideographs.
Series 3: List of NUC (Non-Unifiable Components (NUC) of Ideographs and Overly-Unified Ideographs)
Series 4: List of Possibly Mis-Unified Ideographs (MUI).
Series 5: List of documents, each is used to describinge onea submitter's

normalization guidelines (SNG) for conversion of one particular script to the printed form of ideographs(SNG).

**File format of IWDS:**

Each of the IWDS file is named as IWDS_*SSS_II* where *SSS* is the name of the specific series (may be more than 3 letters for Series 5) and *II* refers to the IRG meeting where the list is confirmed. Thus, the Standing dDocument sSeries has 5 threads of documents. The first 4 threads are specified as follows:

IWDS_SUM: Summary document (Series 1) … This document.
IWDS_UCV: List of UCV list (Series 2)
IWDS_NUC: List of Non-Uunifiable Components (Series 3)
IWDS_MUI: Possibly Mis-Uunified Ideographs (Series 4). _

In Series 5, the guidelines are submitter dependent as well as script dependent. Thus, the file names need to include a submitter idID followed by a script idID. For example, the ROK's current guidelines is are for conversion of cursive Kai to printed Kai style. So the file name should be: IWDS_SNG_ROK_Kai_48.

**Detailed Specification of the Standing Document Series**
This section explains the nature of each series as well as the format and the information contained in each series.

1. Summary of IRG Standing Document Series (SUM)

The summary document (this document) is a definitive document giving detailed specification for each of the data files including specification on the nature of the data, the data format, and the examples.

2. List of Unifiable Component Variations (UCV)

The UCV list provides the a list of component variations to bewhich are unifiable, as observed from existing UCS multi-column charts, or proposed and agreed among IRG members to be unifiable.

If two ideographs differ only in terms of the components in the UCV list, but satisfy the requirement for dis-unification according to dis-unification rules, these ideographs may be encoded differently. However, these cases are exceptional and should be exhaustively listed in this document under the related components to avoid confusion for consideration of other characters.

Unification is meant to be at the component level only. In other words, if the

components themselves are also ideographs proper, this list does not imply that the corresponding ideographs proper are unifiable.

The following is the format for each entry in the file:

a. No.:  The serial number of this entry for reference that is unique throughout the standardization works.

b. Criteria: List of actual glyphs.

c. References: Excerpt from existing document (e.g. JIS X 0213 and HYDZD).

d. Exceptions: The exhaustive list of dis-unification examples.

e. Compatible/Duplicate/Examples: The example list of unified ideographs and compatibility ideographs, and notes if necessary.

The following is an example of a typical entry in the UCV list.



3.  List of Non-Unifiable Components (NUC) of Ideographs(NUC)

The NUC list provides the a list of component variations which are not to be unified. This list should be kept as minimal as possible. Components that are not obviously unifiable will not be listed here. That is, it should only list those that are close in glyph shapes and can be confusing cognitively. In other words, this list should only contain the components which are (possibly inappropriately)

unified by precedence during the IRG working process, ~~or~~and components that are stated to be unifiable by some local national standards~~,~~ but not in the UCS.

Furthermore, this list should not contain components which are ~~either~~ (1) KangXi radicals (such as 工 vs. 土 ) or (2) simplified vs. traditional components with no precedence of unification (such as 門 vs 门).

The following is the format for each entry in the file:

a. Components: List of non-unifiable glyphs that is unique throughout the standardization work~~s~~.

b. Analysis: Reasons for dis-unification and each reason will be listed separately. Typical ~~ones are~~reasons include: the ideographs in question are already separated ~~ones and ideographs which are~~or they are encoded by one~~-~~ side only.

c. ~~Ex~~Samples: List of exhaustive possibly over~~ly~~-unified ideographs (if exists).

The following is an example of a typical entry in the ~~U~~NUC list.~~:~~

| Components | Analysis | Samples |
|---|---|---|
| 麻 厤 | Unifications | |
| | Separations | 麻厤, 厤歷, 厤曆, 歷曆, 瀝瀝, 曆厤, 厤歷 |
| 麻-only | none | |
| 麻-only | 厤灰歷厤厤歷歷灰攦歷厤厤展厤厤歷厤厤鷹魔鷹鷹 |
| 閒 間 | Unifications | 㗵 053/244 35F4 5-3856 5-2454 4-492D 4-4113 3-226F 3-0279 0-3338 0-1924 2-4-34 2442 a4c2 f1c0 35F4 㗵 30 (口) 12 4277 ガン人名 マ 㗵口(マグチ・姓) \| 㗵 㗵 \| 㗵 |
| | Separations | 僩僴, 嫺嫻, 憪憪, 憪憪, 捫捫, 椚椚, 洞澗, 燗燗, 癇痫, 瞯瞷, 磵磵, 簡簡, 網網, 萠莔, 裥裥, 覸覸, 調調, 鐗鐗, 開開, 閒閒, 鷴鷴, 嵧嵧, 徣徣, 琱琱, 酐酐, 鎇鎇, 鮰鮰, 鷗鷗 |
| 間-only | 㗵塥㗵琱椚㗵㗵琱瞷 |
| 閒-only | 胢殑獝鮰轀琱 |

4. List of Possibly Mis-~~Uu~~nified Ideographs (MUI).

The MUI list provides the possibly mis-unified ideographs as pairs of CJK compatibility ideograph~~s~~s and their corresponding CJK unified ideographs, which have different ~~semantics~~ meanings and pronunciations with the supplied related reference information in a single document (possibly a dictionary).

It is possible that the coded CJK compatibility ideographs listed in this document will be proposed as new CJK unified ideographs. However, extreme care must be taken to ~~assure the~~ensure compatibility with existing standards in

accordance with Annex I of WG2's *Principles and Procedures*.

The following is the format for each entry in the file:

a. U-code: The UCS code point.
b. Characters: List of possibly non-unifiable ideographs.
c. References: Excerpts of their usage from a single document source.

The following is an example of a typical entry in the MUI list.



## 5. Set of Printed Character Normalization Guidelines

As the normalization guidelines are submitter/culture/language dependent as well as script dependent, each document should provide an overview of the scope of the guidelines, the major references and authoritative document sources from where the guidelines are derived. A set of rules can be described using text descriptions with ample examples for people in other language/culture environments to follow and to help with the review and acceptance of evidences. It can also serve as a possible unification/dis-unification guide for other submitters. Content of the normalization table should include, but not limited to, the following data for each entry:

a. Serial numbers: A numbering system used internally for indexing and searching.
b. Variant glyph(s): Actual glyph shapes of the components in the source script.
c. Normalized glyph: The corresponding normalized glyph shape of the component.
d. Evidences: Examples of actual character glyphs with reference to their sources.
e. Comments: Any remark that may be helpful to IRG review.

The following is an example of some typical entries in a normalization table (from ROK: IRGN2154V1.1).

| SN (Serial Number) | Variant component shape | Normalized component shape | Evidences, examples, source | comments |
|---|---|---|---|---|
| 1−1 | 步(少) | 步(少) | 頻−頻[A50;14d;5(甲)] 賓−賓[A272;27b;8(丙)] | ※ refer to UCV_N009. |
| 1−2 | 少 | 少 | 砂−砂[A272;27c;4(丙)] 陟−陝[A104;175d;6(丙)] | |
| 1−3 | 少 | 少 | 少−少[A50;30b;7(甲)] | |
| 1−4 | 少 | 少 | 少−少[A40;445b;6(甲)] | |
| 1−5 | 少 | 少 | 少−少[B061;322d;8(乙)] | |
| 1−6 | 少 | 少 | 少−少[B094;516b;1(丙)] | |
| 1−7 | 少 | 少 | 雋−雋[B081;501a;4] | |
| 1−8 | 示(小) | 尐(少) | 歳−歳[B069;097c;1(乙)] | ※ refer to UCV_N015. |

(End of document)