## ISO/IEC JTC1/SC2/WG2/IRG N1757

## Preliminary Proposal for an Ideographic Variation Database Registration

John H. Jenkins 井作恆 Richard Cook 曲理查 Ken Lunde 小林劍 Unicode Technical Committee 6 April 2011

#### Summary

This documented is being presented to IRG members for information only. No action is required.

The Unicode Technical Committee maintains a database of Han ideographs which have been brought to its attention as potential candidates for encoding. This database is found in Unicode Technical Report #45, U-source Ideographs (UTR 45, found at <a href="http://www.unicode.org/reports/tr45/">http://www.unicode.org/reports/tr45/</a>). All U-source indices refer to this document, but a fuller description of the sources for these ideographs is found in Unicode Standard Annex #38, Unicode Han Database (UAX 38, found at <a href="http://www.unicode.org/reports/tr38/">http://www.unicode.org/reports/tr45/</a>).

The UTC is currently investigating the possibility of representing a number of the unencoded ideographs in UTR 45 via registered ideographic variation sequences (IVSs) as described in Unicode Technical Standard #37, Unicode Ideographic Variation Database (UTS 37, found at <a href="http://www.unicode.org/reports/tr37/">http://www.unicode.org/reports/tr37/</a>). The actual database of registered variation sequences is referred to as the Ideographic Variation Database (IVD).

This document summarizes the current state of that investigation.

#### Rationale

The forms included in this proposal are all completely synonymous with the base characters in question. Usage of one form rather than another reflects locale-specific preference, and any differences in glyph shape are otherwise irrelevant to UCS encoding.

The proposed forms are all regular modern simplifications of traditional Chinese characters, produced by applying well-known simplification rules. The cases presented here involve only one-to-one mappings, and conversion from one form to the other is completely reversible.

There is currently an expectation that registered variation sequences cover only forms which are unifiable under Annex S rules (see UTS 37 §2). Registering the current collection would involve a broadening of current guidelines to explicitly allow unification of so-called "y-variants" under certain conditions.

(We note that Annex S is itself an informative part of ISO/IEC 10646 and not normative. The recommendation to use Annex S as a guideline for determining whether or not two ideographs can be unified is also informative within UTS 37.)

Annex S describes variation of ideographs using a three axis model: the x-axis relates to differences in meaning, the y-axis relates to relatively large differences in shape, and the z-axis relates to relatively minor differences in shape. The process of determining whether two synonymous ideographs are y- or z-variants of one another is described in detail in Annex S.

Characters with very similar shapes but very different usage (differing in meaning, along the x-axis) are not unifiable. For example, the characters  $\pm$  (U+672A) and  $\pm$  (U+672B) may be written with very similar forms, but have different usage. They are divergent of n the x-axis (clearly separated in dictionaries), and so are not unifiable. Even though they may easily be confused with one another, that alone has not prevented them from being separately encoded.

Ideographs with only very minor shape differences (such as  $\tilde{\pi}$  and  $\tilde{\pi}$ –U+82B1 as drawn with two different "kai" typefaces) are z-variants as defined by Annex S and may be unified.

Characters with larger shape differences but identical meanings are termed "y-variants." An example would be  $\blacksquare$  (U+8C9D) and  $\square$  (U+8D1D), which are a traditional Chinese/ simplified Chinese pair. The general approach in developing the UCS, applied to every script other than Han, would be to unify these two forms. Annex S currently includes no guidelines for how it could be done now for CJKV ideographs. Adding such guidelines to a future revision of Annex S seems desirable.

As a rule, CJKV y-variants have hitherto been separately encoded, largely for historical reasons. This unnecessarily complicates the encoding, overburdens the IRG and computing processes, and has no practical benefit for anyone. Indeed, the encoding of duplicate characters adds significant cost to implementers and frustrates end-users.

#### Example 1

To illustrate, let's start by looking at the screen shot below.

#### <u>藍若非為科幻小說努力-明報網上書店</u>Q

猶幸明窗出版社繼續支持科幻小說,除了膾灸人口的「衛斯理系列」,更有香港科幻學會主席李 逆熵所寫的《泰拉文明消失之謎》,可洛的「女媧之門系列」,而藍若非的「藍... books.mingpao.com/cfm/books.cfm?Path=editor\_149.htm - 頁庫存檔

#### 香港版:倪匡科幻小说系列原振侠传奇《灵椅》(重200克)-图书价格:15.00 ... ♀ - [轉為繁體網頁] 2010年12月28日 ... 网上书店[明灯楼]在线销售图书香港版:倪匡科幻小说系列原振侠传奇《灵 椅》(重200克), 图书价格:15.00;图书品相:8.5成品相;孔夫子旧书网汇集全国 ... book.kongfz.com/10468/102804575/ - 中華人民共和國 - 頁庫存檔

This screen shot is from a Google search for the phrase "香港科幻小說書店" ("Hong Kong science fiction bookstores"). Two characters in the search phrase (說 and 書) have standard simplifications (说 and 书), and Google has accommodated this, with a search result containing the traditional forms on top and a search result containing the simplified forms on the bottom.

In order to do this, Google needs to have a table of character equivalents. Having to double-check equivalents via a table has a performance impact. It also means that the table of variants has to be continually updated whenever a new CJK Unified Ideographs Extension is encoded.

Moreover, since there is no authoritative source for variant data, different vendors will have different tables and different degrees of support for this feature. This detracts from the end-user experience.

If y-variants are represented via IVSs, however, the variation relationship is automatically encoded in the text itself. If you want one y-variant to match any of the others, you simply ignore the variation selector. If you want one y-variant *not* to match any of the others, you take the variation selector into account. This doesn't even require that you explicitly add support for any particular variation sequence; it's simply a matter filtering out variation selectors via a range check or not filtering them out.

#### Example 2

Searching is not the only process to benefit. Even though fonts do need to be updated whether you use separate encoding or a registered variation sequence, display of text also works better with the use of variation sequences.

For example, consider the line 訏謨定命 from the Chinese classic, the *Shi Jing*.<sup>1</sup> The first two characters both have potential simplifications: 訏 to 玗, and 謨 to 谟, but only the latter is currently encoded. The simplified Chinese form of the line is therefore

<sup>&</sup>lt;sup>1</sup> From Ode 256 (蕩之什,抑); it's translated as "He who takes counsel widely, is final in his commands" by Waley.

usually given as 訏谟定命. The form 玗 is attested, however, and a simplified Chinese book may want to write it 玗谟定命.

If 诗 is separately encoded and a particular user hasn't updated their fonts, what they will see is something like 図谟定命. If 诗 is represented using a variation sequence and the user hasn't updated their fonts, they will see 訏谟定命, which is distinctly better.

(And a text-to-speech engine would know to pronounce 计谟定命 as *heoi1 mou4 ding6 ming6* without needing to be specifically updated as well.)

#### Summary

In general, representation of y-variants with IVSs provides for a better user experience and a more robust and flexible representation of text. It goes without saying that it also reduces the IRG's workload since fewer characters need be submitted for consideration. It would also reduce the workload of specific IRG members—most notably China—because the amount of effort they expend producing and tracking their encoding proposals is reduced.

The characters in this proposal are not common, and use of one y-variant or the other is not mandatory in current practice. Whereas 書店 would be considered "wrong" in a simplified Chinese text, 訏谟定命 is perfectly acceptable, even though 迂谟定命 would be preferred if it's available. There is therefore no significant difficulty in using variation sequences to represent these forms—but there is a distinct benefit.

#### Sources

Detailed descriptions of the sources for these forms are found in UAX 38 and UTR 45. These sources include:

- *ABC Chinese-English Comprehensive Dictionary*. John DeFrancis, ed. University of Hawai'i Press, 2003. ISBN: 0-8248-2766-X.
- •《漢語大字典》,湖北辭書出版社, Wuhan, 1988. ISBN: 7-5403-0030-2/H.16. [<http://www.unicode.org/reports/tr38/#kHanYu>.]
- 《现代汉语词典》 [Xiàndài Hànyǔ Cídiǎn 'Modern Chinese Dictionary']. 中国社会科学院语言研究所词典编辑室编 [Chinese Academy of Social Sciences, Linguisitics Research Institute, Dictionary Editorial Office, eds.]. 北京: 商务印书馆, 2007 [第 5 版; 2007 年 11 月北京第 377 次印刷. ISBN: 7-100-04385-9/H.1100.]
- 《现代汉语词典》 [Xiàndài Hànyǔ Cídiǎn 'Modern Chinese Dictionary']. 中国社会科学 院语言研究所词典编辑室编 [Chinese Academy of Social Sciences, Linguisitics Research Institute, Dictionary Editorial Office, eds.]. 北京: 商务印书馆, 1983 [1978 年

12 月第 1 版; 1983 年 1 月第 2 版; 1984 年 1 月北京第 49 次印刷印张 54; 统一书号: 17017.91]. [<http://www.unicode.org/reports/tr38/#kXHC1983>.]

- 《形音義規範字典》 Xíng-yīn-yì Guīfàn Zìdiǎn. 主编:李行健. 臺北市:五南圖書出版股份 有限公司;語文出版社 (《現代漢語規範字典》), 2003. ISBN: 957-11-3317-5.
  [Taiwanese stroke order; also contains PRC simplified characters, but without stroke diagrams; cp. XHG.]
- 《现代汉语规范字典》 Xiàndài Hànyǔ Guīfàn Zìdiǎn. 主编:李行健. 北京: 語文出版社 (《現代漢語規範字典》), 1998. ISBN: 7-80126-346-4/H.76. [1998 年 1 月; earlier PRC version of ROC 《形音義規範字典》.]
- 文林 Wénlín Software for Learning Chinese, Version 4.0.1. Wenlin Institute: Eureka, California. [includes electronic ABC English-Chinese/Chinese-English Dictionary, DeFrancis et al.]

#### References

IRG N1468 Recommendation For IRG To Use IVD Collections http://appsrv.cse.cuhk.edu.hk/~irg/irg/irg30/IRGN1468IVS\_Recommendation.pdf

IRG N1646 Principles and Procedures <a href="http://appsrv.cse.cuhk.edu.hk/~irg/irg/irg34/IRGN1646Confirmed.doc">http://appsrv.cse.cuhk.edu.hk/~irg/irg/irg34/IRGN1646Confirmed.doc</a>

IRG N1765 Current Status of IVS Support in OSes & Application

UAX #38: Unicode Han Database (Unihan) http://www.unicode.org/reports/tr38/

UTR #45: U-Source Ideographs http://www.unicode.org/reports/tr45/

## Appendix A Text Data to Add to IVD Database Files

We append below a summary of the data to be added to the two text files included in the IVD.

Addition to IVD\_Collections.txt

UTC; UTC-[0-9]{5}; <u>http://www.unicode.org/reports/tr45/</u>

Additions to IVD\_Sequences.txt

5D19	E0101;	UTC;	UTC-00668
7A68	E0100;	UTC;	UTC-00669
7D41	E0101;	UTC;	UTC-00029
7D9D	E0101;	UTC;	UTC-00914
8A0F	E0100;	UTC;	UTC-00071
8B30	E0101;	UTC;	UTC-00030
8B46	E0101;	UTC;	UTC-00675
8B54	E0101;	UTC;	UTC-00676
8F36	E0101;	UTC;	UTC-00024
91B2	E0101;	UTC;	UTC-00038
9265	E0101;	UTC;	UTC-00052
96A4	E0101;	UTC;	UTC-00674
982B	E0101;	UTC;	UTC-00677
992C	E0101;	UTC;	UTC-00678
99BC	E0101;	UTC;	UTC-00842
9A23	E0101;	UTC;	UTC-00679
9D4F	E0100;	UTC;	UTC-00117
9DB1	E0101;	UTC;	UTC-00680
9DC3	E0101;	UTC;	UTC-00061
9DC7	E0101;	UTC;	UTC-00068
9F6E	E0101;	UTC;	UTC-00013

## Appendix B Mappings and Glyphs

We append below a summary of the proposed IVSs, showing the glyphs for the composed form as shown in UTR 45 and, in parentheses, the glyph for the base form.

<u+5d19< th=""><th>U+E0101&gt;</th><th>==</th><th>UTC-00668</th><th>芲</th><th>(崙)</th></u+5d19<>	U+E0101>	==	UTC-00668	芲	(崙)
<u+7a68< th=""><th>U+E0100&gt;</th><th>==</th><th>UTC-00669</th><th>穮</th><th>(穨)</th></u+7a68<>	U+E0100>	==	UTC-00669	穮	(穨)
<u+7d41< th=""><th>U+E0101&gt;</th><th>==</th><th>UTC-00029</th><th>纯</th><th>(絁)</th></u+7d41<>	U+E0101>	==	UTC-00029	纯	(絁)
<u+7d9d< th=""><th>U+E0101&gt;</th><th>==</th><th>UTC-00914</th><th>綝</th><th>(綝)</th></u+7d9d<>	U+E0101>	==	UTC-00914	綝	(綝)
<u+8a0f< th=""><th>U+E0100&gt;</th><th>==</th><th>UTC-00071</th><th>讶</th><th>(訏)</th></u+8a0f<>	U+E0100>	==	UTC-00071	讶	(訏)
<u+8b30< th=""><th>U+E0101&gt;</th><th>==</th><th>UTC-00030</th><th>涟</th><th>(謰)</th></u+8b30<>	U+E0101>	==	UTC-00030	涟	(謰)
<u+8b46< th=""><th>U+E0101&gt;</th><th>==</th><th>UTC-00675</th><th>语</th><th>(譆)</th></u+8b46<>	U+E0101>	==	UTC-00675	语	(譆)
<u+8b54< th=""><th>U+E0101&gt;</th><th>==</th><th>UTC-00676</th><th>谍</th><th>(譔)</th></u+8b54<>	U+E0101>	==	UTC-00676	谍	(譔)
<u+8f36< th=""><th>U+E0101&gt;</th><th>==</th><th>UTC-00024</th><th>緧</th><th>(輶)</th></u+8f36<>	U+E0101>	==	UTC-00024	緧	(輶)
<u+91b2< th=""><th>U+E0101&gt;</th><th>==</th><th>UTC-00038</th><th>畞</th><th>(醲)</th></u+91b2<>	U+E0101>	==	UTC-00038	畞	(醲)
<u+9265< th=""><th>U+E0101&gt;</th><th>==</th><th>UTC-00052</th><th>铽</th><th>(鉥)</th></u+9265<>	U+E0101>	==	UTC-00052	铽	(鉥)
<u+96a4< th=""><th>U+E0101&gt;</th><th>==</th><th>UTC-00674</th><th>隤</th><th>(隤)</th></u+96a4<>	U+E0101>	==	UTC-00674	隤	(隤)
<u+982b< th=""><th>U+E0101&gt;</th><th>==</th><th>UTC-00677</th><th>粄</th><th>(頫)</th></u+982b<>	U+E0101>	==	UTC-00677	粄	(頫)
<u+992c< th=""><th>U+E0101&gt;</th><th>==</th><th>UTC-00678</th><th>翖</th><th>(餬)</th></u+992c<>	U+E0101>	==	UTC-00678	翖	(餬)
<u+99bc< th=""><th>U+E0101&gt;</th><th>==</th><th>UTC-00842</th><th>玟</th><th>(馼)</th></u+99bc<>	U+E0101>	==	UTC-00842	玟	(馼)
<u+9a23< th=""><th>U+E0101&gt;</th><th>==</th><th>UTC-00679</th><th>騣</th><th>(騣)</th></u+9a23<>	U+E0101>	==	UTC-00679	騣	(騣)
<u+9d4f< th=""><th>U+E0100&gt;</th><th>==</th><th>UTC-00117</th><th>鹤</th><th>(鵏)</th></u+9d4f<>	U+E0100>	==	UTC-00117	鹤	(鵏)
<u+9db1< th=""><th>U+E0101&gt;</th><th>==</th><th>UTC-00680</th><th>骞</th><th>(鶱)</th></u+9db1<>	U+E0101>	==	UTC-00680	骞	(鶱)
<u+9dc3< th=""><th>U+E0101&gt;</th><th>==</th><th>UTC-00061</th><th>斟</th><th>(鷃)</th></u+9dc3<>	U+E0101>	==	UTC-00061	斟	(鷃)
<u+9dc7< th=""><th>U+E0101&gt;</th><th>==</th><th>UTC-00068</th><th>高及</th><th>(鷇)</th></u+9dc7<>	U+E0101>	==	UTC-00068	高及	(鷇)
<u+9f6e< th=""><th>U+E0101&gt;</th><th>==</th><th>UTC-00013</th><th>止<u>大</u> 囚可</th><th>( 齮 )</th></u+9f6e<>	U+E0101>	==	UTC-00013	止 <u>大</u> 囚可	( 齮 )

### Korea JTC1/SC2, Committee on Character Codes

KIM, Kyongsok, Chairperson of Korea JTC1/SC2 email: gimgsO AT gmail DOT com, phone: +82-51-510-2292 address: Division of Computer Science and Engineering, Pusan National University 2 Busandaehagro 63 Beon Gil, Geumjeonggu, BUSAN 609-735, Rep. of KOREA

Author: KIM, Kyongsok Date: 2011.06.01. Status: NB position Subject: R.O.Korea's comments RE: IRG N1757 (UTC Preliminary Proposal for an IVD Registration) Relevant documents: IRG N1757,

1. A related IRG resolution (M36.2):

# **Resolution IRG M36.2: UTC Preliminary Proposal for an IVD Registration** (IRGN1757)

Unanimous

The IRG has reviewed IRGN1757 from the UTC and is concerned that an IVD registration from the UTC might be treated as a de facto encoding. The IRG requests that the UTC not proceed with the registration until this issue is addressed.

The IRG encourages its members to do further review of IRGN1757 and provide feedback to the UTC at any time.

#### 2. R.O. Korea's feedback:

 $R.\,0.\,Korea$  discussed this issue within domestic committee and concluded as follows:

- In CJKU, y-variants such as simplified Chinese characters were encoded separately from the traditional ones (base characters).

- If y-variants should be registered with IVD, some y-variants (simplified) will be using the same code position as traditional ones while other y-variants (simplified) will be using code position different from traditional ones.

- R.O. Korea is concerned that this will cause confusion to the users.

- Furthermore, in the future, there will be confusion as to whether y-variants should be given a separate code position or registered with IVD.

- As a conclusion, R.O.Korea suggests that y-variants be given a separate code position instead of registering y-variants with IVD.

\* \* \*