IRG N2176 (v2 modified 2016/10/20)

Title: New UCV Suggestions

Author: Henry Chan

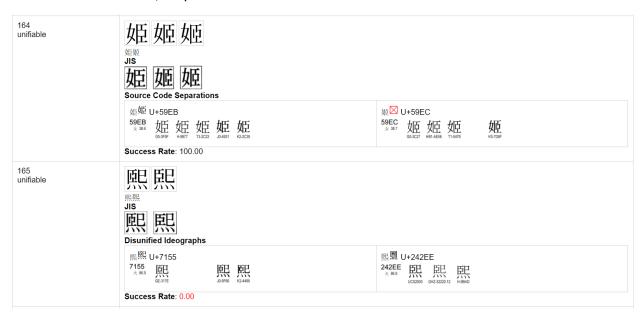Type: Individual Contribution to IRG #47

## Issue 1A

### 臣 vs 臣

According the current NUCV, this pair cannot be unified:

In reality, only U+268F1 is cognate with U+8CFE. (U+24C7B, U+2981D, U+29821, U+2981E and U+29822 cannot be verified). The other characters are non-unifiable.

Also, following this rule, there are also rules 164 and 165 which indicate when 臣 or 臣 are in ［□女 ⬛］ or ［曰□⬛巳灬］, they are unifiable:

| 164 unifiable | 姬 姬 姬<br>姬姬<br>**JIS**<br>姬 姬 姬<br>**Source Code Separations** | | |
|---|---|---|---|
| | 姬姬 U+59EB<br>59EB 姬 姬 姬 姬 姬<br>女 38.6   G5-3F5F   H-9B77   T3-2C23   J0-4931   K2-2C35 | | 姬⊠ U+59EC<br>59EC 姬 姬 姬    姬<br>女 38.7   G0-3C27   HB1-AE56   T1-5478    K0-7D6F |
| | **Success Rate**: 100.00 | | |
| 165 unifiable | 熙 熙<br>熙熙<br>**JIS**<br>熙 熙<br>**Disunified Ideographs** | | |
| | 熙熙 U+7155<br>7155 熙    熙 熙<br>火 86.9   GE-317E    J0-5F66   K2-4466 | | 熙⊠ U+242EE<br>242EE 熙 熙 熙<br>火 86.9   UCS2003   GHZ-32220.12   H-9BAD |
| | **Success Rate**: 0.00 | | |

臣 and 臣 are completely different in pronunciation, so if two characters containing these component, it should be easy to identify if they are non-cognate. In handwriting, the two components are somewhat frequently messed up. Such a rule may unnecessarily hinder the unification of "improper" variants.
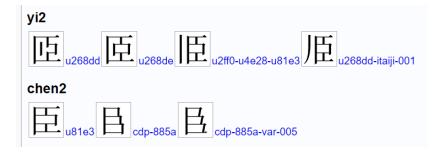
This is similar to these other UCV rules in nature:

| 126 · unifiable | 243 · unifiable |
|---|---|
| 本 本 | 免 兔 |

Their use in characters is often as a phonetic component, but it is easy to distinguish.

It is suggested that this rule be removed from the NUCV.

**Issue 1B**

Per handwriting conventions, it is suggested that these variants be (independently) added to the UCV:
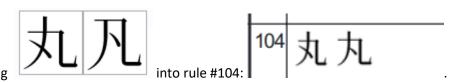


Existing disunified characters are as follows:





The Kangxi entry for U+205f1 is:

 which indicates they are cognate, and could have been unified, when complemented by rule

364: 

This rule substitutes UCV rule 163, 164, 165, and 167, and complements rule 364.

**Issue 2**

Suggest merging  into rule #104:  .

**Rationale**

U+2007D is the ancient form of U+4E38.  According to Kangxi dictionary, Kangxi thinks that 丸 is a 俗體 (corrupted form) of 凡:



The origin of the ancient shape is from seal script, where the "人" does not protrude out of the container:



『說文解字』

【卷九】【丸部】
圜，傾側而轉者。从反仄。凡丸之屬皆从丸。胡官切

『說文解字注』

(丸) 圜也。也字各本無。今依韵會補。以疊韵爲訓也。今丸藥其一崗也。商頌。松栢丸丸。傳曰。丸丸、易直也。按謂其滑易而調直也。丸義之引伸也。大雅。松柏斯兌。傳亦云。兌、易直也。兌與丸、古蓋音同而義同矣。傾側而轉者。从反仄。圜則不能平立。故从反仄以象之。仄而反復、是爲丸也。胡官切。十四部。凡丸之屬皆从丸。

The shape of U+2007D 凡 is easily confused with U+51E1 凡. As such, as mentioned in IRGN2174 Part 1, U+2F8FA has been mistakenly unified to U+6C4E 汎 instead of U+6C4D 汍.

The small top protrusion is significant in distinguishing the etymology of the character in the Kangxi Dictionary.  In modern life,  the shape of U+2007D 凡 is rare. The shape of U+4E38 丸 is often used, as a single character or as a character component.  It is recommended that U+2007D and U+4E38 be unifiable components to reduce any confusion.

Existing Disunified Examples:

#2 丸 ancient

丸 u4e38　凡 u2007d
汃 u6c4d　汎 u2f8fa
骪 u9aab　骫 u9aaa

## 4E38
丶 3.2

丸 G0-4D68　丸 HB1-A459　丸 T1-443A　丸 J0-345D　丸 K0-7C2F

## 2007D
丶 3.2

凡 UCS2003　凡 GKX-0080.14　凡 T6-2132

## 9AAA
骨 188.3

骫 GE-4643　骫 T4-4221　骫 J1-6956　骫 K1-6966

## 9AAB
骨 188.3

骪 G3-796D　骪 HB2-E0E9　骪 T2-486F　骪 K2-6F54

The glyphs for U+9AAA from G-source, T-source and K-source may considered for changing to better reflect the etymology.

**Issue 3**

Suggest to add this pair:

永 u6c38-ue0100    永 u6c38-ue0101    汞 zihai-101007

Characters with the 永 phonetic are frequently written as the shape "zihai-101007" in Kangxi dictionary.

Existing Unification example:

9721
雨  173.10

霡  霡  霡  霡

G3-744B    H-A066    T4-5E38    J1-6727

Existing Disunification (mis-disunification) example from Extension-B:

錄 (u92a2)

錄 (u92a2) (=jmj-026837) (=u92a2

錄 (u92a2-ue0101) (=jmj-026838

錄 (koseki-458370) (=juki-92a2)

錄 (hkcs_m92a2)

錄 (u2896d)
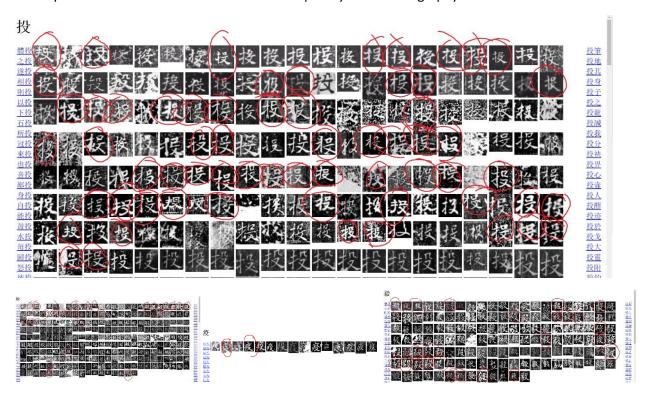[関連字（その他）]

錄 (u2896d) (=juki-bdde)

錄 (hkcs_m2896d)

Per the discussion in IRG#47, this difference is usually normalized away by ROK. We should make them unifiable as well, as not all regions can carry out normalization to their glyphs.

**Issue 4**



#4 殳 handwriting deviation

殳 u6bb3　殳 u2ff1-u20003-u4e42　旻 u2ff1-u53e3-u53c8

These pairs are handwritten deviation forms frequently seen in calligraphy:



Allowing the disunification of could lead to disastrous results.

However, it may be too late to do anything with the existing characters in Extension F. The disunification of the JMJ characters in Extension F should not be regarded as a "general rule".

**Existing Disunified cases:**

U+2A832 旻: variant of U+20B1B 殳. 殳 is unifiable with 旻.

U+23CDA 浸: variant of U+6C92 没. 殳 component is unifiable with 殳.

U+22919 懿: variant of U+2289E 懿.

U+253CF 殿: variant of U+6BB9 殿

U+22936 懿: Variant of U+6128 愨.

**Issue 5**

Propose adding the following unification:

亞 u4e9e    亜

The shape difference is minor and could be handled via IVS. The two shapes a purely calligraphic stroke ordering differences.

Current disunified characters:

亞 u4e9e   亜 u4e9c

啞 u555e   唖 u5516

惡 u60e1   悪 u60aa

婭 u5a6d   娅 u2bc2b

椏 u690f   桠 u2c0b2

噁 u5641   噁 u2bad6

A unification example:

2827C
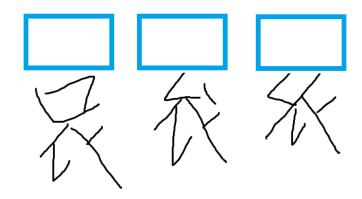身 158.11    軀惡    軀惡    軀惡    2
UCS2003    TF-6127    H-8BAE

Previous discussion about disunifying this pair remarked that they should not be disunified, but this could not be regarded as a general unification example.

The proposal is to regard make it the preferred unification, given the comments by Whistler to reduce "unnecessary variants".

**Issue 6**

Add new entry into UCV:



Disunified examples (that are cognate):

| | |
|---|---|
| 袁 u8881 | 袠 u212ae |
| 遠 u9060 | 遠 u903a |
| 猿 u733f | 猿 u2b7a4 |
| 袞 u889e | 袞 u2c844 |
| 罷 u7758 | 罷 u262b7 |

This proposal supersedes #293:

**Issue 7**

Add this rule:

羊 u7f8a  芈 u2634b  芉 zihai-021410  羊 u7f8a-03-var-002  芈 u2634b-03

These shapes arise from the different Kaishu-fication of the seal shape:

『說文解字』

【卷四】【羊部】羊
祥也。从丫，象頭角足尾

『說文解字注』

羊

（羊）祥也。疊韵。考工記注
曰。牛羊之字。以形舉也。許多引

(Screenshot from zdic.net)

Existing Disunified (mis-disunified) cognates:

羊 u7f8a  芈 u2634b
攘 u3a3e  㨇 u22d47
嬙 u21842  嬚 u21818
繕 u7e55  繕 u26187
撍 u22d48  撍 u22d2f
羞 u7f9e  羞 u2636e  羞 u2635f
薈 u22431  薈 u2242f
蕎 u26c36  蕎 u26c86

**Issue 8**

We should expand existing UCV rule #393 to cover more scenarios.

UCV Rule #393:

**393 · unifiable**

走走

Rule 393 should be expanded to cover the following scenarios:

足 u8db3  昆 u20bc1  走 u8d70  赱 u8d71  縦 u7e26  縱 u260b5
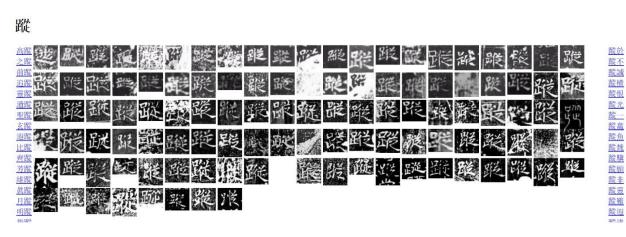
Proposed modified UCV Rule:

#393:

**Issue 9**

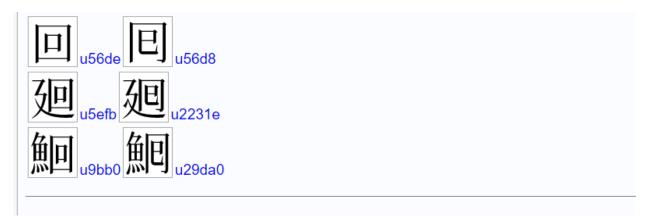Consider adding rule to UCV for these variants:



These variants are cognates.

This variation is quite common.  In fact, it is often combined with modified #393 UCV:



However, the shape difference might be "too significant" to unify?

**Issue 10**

Consider adding rule to UCV for these variants:

回 u56de 㔾 u56d8
廻 u5efb 㘞 u2231e
鮰 u9bb0 𩶠 u29da0

**Issue 11A**

Consider adding rule to UCV for these variants:

#11 丑刃 handwriting deviation

丑 u4e11 刄 u4e12
杻 u677b 枫 u2b788

The Zhuang characters submission by PRC has also normalized these shapes:

**Issue 11B**

2CED8 　�088
一 1.7

JMJ-057477　　　　　is a variant of 芻.

In fact, all of these are cognates according to Hanyu Dazidian and Zhonghua Zihai:

嚃☒嚃 (U+35D9)　　嚃☒嚃 (U+20D3E)
傷☒傷 (U+3473)　　僲☒僲 (U+202A3)
搊☒搊 (U+640A)　　揕☒揕 (U+22BA2)
榖☒榖 (U+3BB2)　　橲☒橲 (U+23516)
燭☒燭 (U+717C)　　煊☒煊 (U+2429B)
犓☒犓 (U+7293)　　牷☒牷 (U+24684)　　犓☒犓 (U+246C3)
皺☒皺 (U+76BA)　　皷☒皷 (U+24FFF)
穊☒穊 (U+4173)　　稑☒稑 (U+257F7)
篘☒篘 (U+7BD8)　　篁☒篁 (U+25BA2)
糫☒糫 (U+25EE4)　　糧☒糧 (U+25EA3)
膓☒膓 (U+4422)　　朣☒朣 (U+2673B)
艣☒艣 (U+447C)　　艤☒艤 (U+26A59)
蒭☒蒭 (U+84AD)　　薑☒薑 (U+26C6A)
謅☒謅 (U+8B05)　　譆☒譆 (U+27A56)
蹰☒蹰 (U+280D8)　　躋☒躋 (U+28093)
鄒☒鄒 (U+9112)　　郪☒郪 (U+28708)
鞒☒鞒 (U+29304)　　韜☒韜 (U+292D3)
騶☒騶 (U+9A36)　　驔☒驔 (U+298FF)
鸀☒鸀 (U+29C48)　　囍☒囍 (U+29C42)
齺☒齺 (U+9F7A)　　齟☒齟 (U+2A634)

However, the difference between 芻 and JMJ-057477 is too big in my opinion.  Consider add to NUCV.

**Issue 12**

Consider adding a UCV rule for these variants:



#12 □ 丙丙 handwriting variant

| | |
|---|---|
| 麗 u9e97 | 麗 u2a2d8 |
| 廲 u5ef2 | 麗 u22312 |
| 酈 u9148 | 酈 u287eb |
| 鸝 u9e1d | 鸝 u2a239 |
| 藶 u457b | 藶 u27173 |
| 覿 u4695 | 覿 u278ae |
| 艫 u27956 | 艫 u27957 |

They are cognate.

**Issue 13:**

Consider adding UCV for these variants:

| | | | |
|---|---|---|---|
| 犮 u72ae | 叐 u53d0 | | |
| 拔 u62d4 | 抜 u39de | 拔 u629c | 㨌 u22b0c |
| 跋 u8dcb | 跜 u47e6 | | |
| 妭 u59ad | 媛 u2a972 | | |
| 髮 u9aee | 髪 u9aea | 髳 u28c73 | |
| 𩊊 u2928a | 鞍 u292a4 | | |
| 柭 u67ed | 枛 u2342a | | |
| 𤤒 u24912 | 玅 u24923 | | |
| 紱 u7d31 | 綏 u25fc8 | | |
| 䮂 u4b82 | 駁 u2989a | | |

These examples are consistent with ROK's normalization.

ROK normalization also includes this pair.   However, such normalization generally requires semantic decomposition to determine the correct glyph shape.  Consider adding to UCV with a note.
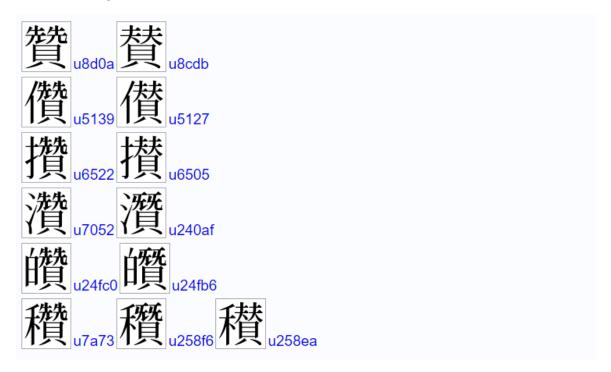
**Issue 14:**

Consider adding UCV for these variants:

#14 朁 handwriting variant

| | | |
|---|---|---|
| 朁 u6701 | 㬱 u3b31 | |
| 僭 u50ed | 僣 u50e3 | |
| 劖 u3506 | 劆 u207c6 | |
| 嘈 u5646 | 嘈 u2110b | 嘈 u20fb1 |
| 嶜 u5d9c | 嶜 u21f2b | 巆 u21fa1 |
| 憯 u61af | 憯 u39a7 | 憯 u2285a |
| 槮 u6a6c | 槮 u236bd | |
| 潛 u6f5b | 潜 u6ff3 | 潜 u6f5c |
| 熷 u71b8 | 熷 u243d6 | |
| 譖 u8b56 | 譖 u27b82 | 譛 u8b5b |
| 瞶 u406e | 瞶 u25333 | 瞨 u252cb |
| 羳 u437c | 羳 u263cb | |
| 膪 u4436 | 膪 u233af | 膪 u2681b |
| 篸 u7c2a | 篸 u2c582 | 簪 u7c2e |
| 鐕 u9415 | 鐕 u28be9 | 錯 u941f |
| 蕳 u26ef3 | 蕈 u26ed8 | |
| 譖 u8b56 | 譖 u8b5b | |
| 蹧 u28154 | 蹧 u28155 | |
| 鄪 u48df | 鄪 u48e0 | |
| 雦 u2903f | 雦 u29040 | |
| 霅 u29168 | 霅 u2916e | |
| 顅 u4aec | 顅 u29557 | |

The more etymologically correct glyph is on the left.

**Issue 15:**

Consider adding UCV rule for these characters:

贊 u8d0a　贊 u8cdb

儹 u5139　儧 u5127

攒 u6522　攢 u6505

瓉 u7052　瓚 u240af

曊 u24fc0　曊 u24fb6

穳 u7a73　穳 u258f6　穳 u258ea

The more etymologically correct form is on the left.


**Issue 16**

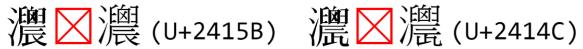Consider adding UCV rule for these characters:

#16 蠻 handwriting variant

鬵 u9b35　鬵 u29c55　鬵 u29c3f

儹 u203b1　儹 u20437

蠽 u27329　蠽 u2746f

蠿 u2757d　蠿 u2759f

The more etymologically correct form is on the left.

**Issue 17:**

Consider adding these glyphs to a UCV rule:

#17 辰/畏 ancient

辰 u8fb0　辰 u28443　瓰 u20a37　辰 u28444

畏 u754f　𤱖 u24c56

灐⊠灐 (U+2415B)　灐⊠灐 (U+2414C)

These glyphs are derived from different Song/Ming typeface interpretations of the similar in Shuowenjiezi.


**Issue 18:**

Consider adding these glyphs to a UCV rule:
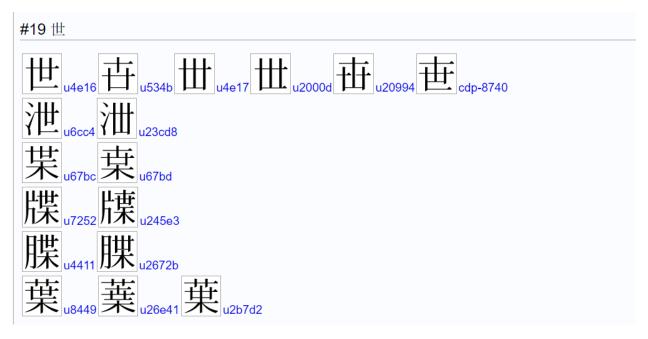
#18 奔奔

奔 u5954　奔 u22343

踌 u2807c　踌 u280e6

A "running person" is represented in oracle script as a person throwing his arms up and down with his head slanted.  It is modernized as 大 (big) (two arms downwards) or 夭 (devil) (head slanted and two arms downwards), which are both characters of unrelated origin.  The presence of an additional slant is not contributive to the meaning of the character.  (Similar to issue 13)

**Issue 19:**

Consider adding UCV to cover the following cognates:

#19 世

世 u4e16 　 卋 u534b 　 世 u4e17 　 丗 u2000d 　 丗 u20994 　 丗 cdp-8740

泄 u6cc4 　 泄 u23cd8

枼 u67bc 　 某 u67bd

牒 u7252 　 牒 u245e3

𤐑 u4411 　 脒 u2672b

葉 u8449 　 葉 u26e41 　 葉 u2b7d2

The included variations are cognate.

**Issue 20:**

Consider noting these two cases as mis-disunification of "missing of minor dot".

#20 Taboo Ommitted Dots:

瓜 u74dc 　 瓜 u244f0

玄 u7384 　 玄 u248e5

**Issue 21:**

Consider covering the following cognates:

#21 丩 卭 cdp-88b5

丩 u4e29 卩 u200cf
杊 u673b 枛 u233b9

**Issue 22A:**

272 · unifiable

皋 臯 臯 臯 臬

臯

Existing UCV:

Consider covering these cognates:

#22 臯臯臯臯

臯 u768b 臯 u7690 臯 u81ef 臬 u2690e

噪 u55e5 嘩 u5651 嘷 u5637

暤 u66a4 暉 u66ad 暤 u66cd 暤 u2328a

槹 u69d4 槹 u69f9 槹 u6a70 槹 u23636

滜 u6edc 潿 u23f4e

獋 u7346 獋 u7354 獋 u734b 獋 u2c342

皞 u769e 皞 u76a1 皞 u76a5

翱 u7ff1 翱 u7ff6 翱 u7ffa

韓 u97df 韓 u293e4

鶴 u9dce 鶴 u9df1

嶂 u37f8 嶂 u21ec7 嶂 u21f17

瘇 u3fc1 瘒 u24e83

蓇 u450c 蓇 u26e86

覾 u468c 覾 u2788c

�campered u4702 �campered u27b01

鎯 u4730 鎯 u27bcc

顆 u4ae7 顆 u2954d

嫂 u217f7 嫂 u21816

稾 u25845 稾 u25890

腺 u267b5 腺 u267de

轖 u28380 轖 u283a6

The shape difference between 衤 and 丰 are too significant. However, the top part difference of 白 and 自 should be unifiable, when the bottom part is consistent.

**Issue 22B:**

Considering unifying 白 and 自 in this case as well:

| 2EBB0 | 鼻 209.0 | 臭 JMJ-059275 |
| 2EBB1 | 鼻 209.0 | 臰 JMJ-059277 |
| 2EBB2 | 鼻 209.0 | 鼻 JMJ-059273 |
| 2EBB3 | 鼻 209.0 | 鼻 JMJ-059272 |
| 2EBB4 | 鼻 209.0 | 鼻 JMJ-059276 |
| 2EBB5 | 鼻 209.0 | 鼻 JMJ-059274 |

**Issue 23:**

#withdrawn

**Issue 24:**

These glyphs are the same character. Consider marking them as mis-disunifications due to rule #328 and #329:

**328** · unifiable
丞 从 瓜 豖

**329** · unifiable
眾 眔 眾 衆 象

#24 乑

巫 u2b851　乑 u4e51

臮 u81ee　臮 u2690c

濕 u23f84　濕 u23f4d

眾 u773e　衆 u8846　眾 u2c454

溼 u6f40　溔 u6f68　澩 u6f48

遝 u28599　遝 u285bc

霂 u29149　霂 u29167

鱳 u29ef4　鰥 u29ec4

#25 莁糸

**Issue 25:**

Consider adding UCV rule for these variants:





The glyph design of u26f09 should be  in the middle but the PRC normalization conventions convert it into grass (  )

**Issue 26:**

Consider new UCV rule for these cognates:

| | |
|---|---|
| 榮 u7162 | 熒 u712d |
| 檺 u23727 | 檓 u6a69 |
| 償 u203bd | 僗 u348c |
| 蕶 u26f53 | 蕶 u26eba |

**Issue 27 & 28:**

Ditto as above.

| | | |
|---|---|---|
| 宜 u5b9c | 宂 u5b90 | 冝 u519d |
| 疊 u758a | 疊 u7589 | 疊 u24d01 |
| 曡 u66e1 | 曡 u3b2a | |
| 氊 u3cb2 | 氊 u6c0e | |
| 攝 u3a79 | 攝 u3a78 | |

| | |
|---|---|
| 示 u793a | 禾 u25605 |
| 社 u793e | 社 u21279 |
| 禰 u256dc | 禰 u271d8 |
| 祈 u7948 | 斫 u23098 |

**Issue 29:** Ditto.

#39 辻 u758c 辵 u758c-itaiji-001

婕 u5a55 媫 u5aab
健 u5022 偯 u507c
捷 u6377 揵 u3a17
葏 u84f5 蒇 u26ef4
蠌 u27425 蠌 u2747b
嗹 u5551 嚏 u20e1d
楗 u234c9 楗 u23579
睫 u776b 瞂 u25224
菨 u8410 菨 u26d49
諜 u8ab1 諜 u27a8c

**Issue 30:**

287 · unifiable
卑 畀

Update #287 to cover all following characters and mark as mis-disunification:

#30 卑 u5351 畀 u24c1e

卑 u5351 畀 u24c1e
淠 u6e12 淠 u23d2a
碑 u7891 碑 u254d3
稗 u7a17 稗 u257d1
粺 u7cba 粺 u25e9b
脾 u813e 脾 u26709
睥 u41d1 睥 u25a8e
痺 u75fa 痺 u24dd2
髀 u9ac0 髀 u29a59
俾 u4ffe 俾 u20237

**Issue 31 (added per discussion in IRG#47 1st day)**

Add to UCV:

詹 詹 詹

This variation is very frequently seen:

詹

子詹
鄭詹
少詹
叔詹
夫詹
士詹
倫詹
使詹
陪詹
道詹
速詹
逮詹
與詹
約詹



擔

荷擔
武擔
塞擔
後擔
錫擔
負擔



膽

肝膽
是膽
之膽
曰膽
張膽
賭膽
魂膽



As discussed, the allowed disunification could be catastrophic given the high number of characters containing the 詹 phonetic.

No disunification example exist.

**Issue 32 (added per discussion of IRG #47 1ˢᵗ day)**

Add new NUCV:

襃 裵

As discussed, shape difference is significant: number of components may be counted differently.

Disunification in URO:

U+61F7 懷 vs U+61D0 懐

U+58DE 壞 vs U+58CA 壊

**Issue 33 (added per discussion on IRG #47 1ˢᵗ day)**

Modify #268 and #269:



To



To cover a unification case discussed by IRG regarding U+24261

**Discussion:**

Traditionally, in certain East Asian regions, a heavy emphasis is placed on "correct" or "proper" form in official contexts. Variants frequently seen in calligraphy is because calligraphy is "an art". Handwritten characters are often "normalized" to a certain glyph shape before they are added into national encoded character sets.

Each region may have different preferred forms. The ISO/IEC10646 is a standard that encodes via a character basis, not a glyph basis. Therefore, similar forms are usually unified to the same codepoint. This unification across all regions (GHTJKV) is most significant in URO. The current Annex S and hence UCV rules are based on these inter-region unifications precedent.

However, due to legal reasons, some regions may need to assign a new code-point in their national character set for every glyph variant that occurs, in handwriting or print, no matter how small the variance. This is usually particular to family register computer systems.

Unlike the normal character sets submitted to IRG, these family register character sets may contain many variants of the same character. Now, the IRG has decided that IVS is the better solution. This "unification" (intra-region unification) by encoding variants via IVS is also currently called "unification", but it is different in nature to the unification (inter-region unification) carried out previously. This new type of unification is in nature a kind of "normalization" because the variant shapes are discarded from UCS. The only difference is where it is carried out: at the level of IRG (e.g. TCA's postponed unifications in IRG #47), or before the submission to IRG (e.g. ROK's normalization in IRG #47).

In the intra-region unification, often semantic analysis is required, and the definition of "minor difference" is rather arbitrary. More often, it depends on the evidence supplied from the submitter to determine the correct base character it should unify to. In inter-region unification, if the glyph shape difference is rather large, it has been simply disunified (e.g. such as the separate coding of 壞 and 壊 in URO).

The current UCV is not constructed to handle these two situations well. Also, the current UCV is very long and quite hard to use. Many characters that are mis-disunifications are currently categorized as "Disunified", because analyzing them is hard.

## Proposal

I propose that we split the UCV into two main subgroups: Assimilation and Variance.

**Assimilation** should concern where two similar shape but phonetically distinct characters are often mixed up when used inside a character.

Examples include:



Whether two glyphs containing these components can be unified (same character) or cannot be unified (non-cognate, different character) heavily depends on semantic analysis. A lot of attention should be placed – these are cases where previous unifications have gone (horribly) wrong.

My additional proposal is, IRG should specify a virtual "proper" shape for the separate etymologies, for determining a SC and FS for multi-source and single-source characters. The "proper" shape need not follow Kangxi closely, but should use forms traditionally used to denote different etymology. An example is, in Taiwan and Hong Kong,  are distinguished components which reflect their etymology, which is generally not distinguished in Kangxi. Having a distinguished virtual shape will more intuitively identify any non-cognates.

**Variance** would concern variances in glyph that do not incur phonetic differences for characters than contain them. Unifications in this category should be nearly 100%. Complex components which swap out part of their component with another component with a completely different semantic meaning, but have no added effect on glyphs containing this complex component also belong to this category.

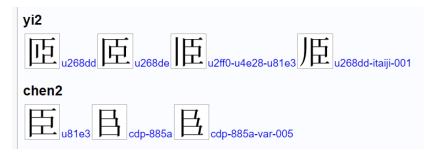Variances that do not concern any etymological differences:

Variant components with components which differ in etymology, but have no added effect:



A single "proper" form should be specified for SC and FS counting across all regions.

For this proposal, some rules would need to be split into two or more rules.

Assimilation rules may also include variance rules. (Refer to Issue 1B:)



For the top four glyphs, they are a "variance set 1". For the bottom three glyphs, they are a "variance set 2". The compound of these two variance sets is another "assimilation" rule.

**Discussion Item 2:**

To aid in duplicate removal and identification of similar shape non-cognates, a virtual normalized glyph can be generated for each glyph submitted according to the proper forms specified in the Assimilation rules and the proper form for the Variation rules. By specifying a normalized virtual glyph, SC can also be easily calculated easily.

The SC does not have to follow anyone's convention closely; it only needs to be consistent. In my opinion, IRGN954AR is a preliminary specification of "proper" forms because the first stroke and stroke count is generally from the glyph of the head character.

The generated virtual glyphs, if necessary, can be included directly in the ISO/IEC10646 standard as well. This is similar to UCS2003 glyphs in Extension B, which effectively act as the glyph that other characters are unified to. The UCS2003 glyphs in the past suffer from the problem that the glyph normalization was not consistent. With "proper" forms specified, we can also avoid this problem.