

Title:	HKSCS and GB 18030 PUA characters, background document
Source:	UTC/US
Authors:	Michel Suignard, Eric Muller, John Jenkins
Action:	For consideration by UTC and IRG

Summary

This document describes characters still encoded in the Private Use Area of ISO/IEC 10646/Unicode as commonly found in the mapping information for Chinese coded characters such as HKSCS and GB-18030. It describes a new encoding proposal to eliminate these Private Use Area allocations, so that the PUA can really be used for its true purpose. Doing so would tremendously improve interoperability between the East Asian market platforms because support for Government related encoded repertoire would not interfere with local comprehensive usage of the PUA area.

Hong Kong Supplementary Character Set (HKSCS)

According to <http://www.info.gov.hk/digital21/eng/hkscs/download/big5-iso.txt> there are a large number of HKSCS-2001 characters still encoded in the Private Use Area (PUA). A large majority of these characters looks like CJK Basic stroke that could be used to describe the appearance of CJK characters. Although there are already collections of various CJK fragments (such as CJK Radicals Supplement, Kangxi Radical) and methods to describe their arrangement using the Ideographic Description Characters, these 'stroke' elements stand on their own merit as an interesting mechanism to describe CJK characters and corresponding glyphs.

Most of these characters have been proposed for encoding on the CJK Extension C. However that extension is not yet mature, but at the same time removing characters from the PUA is urgent. A better solution seems to create a new block containing these Basic CJK Unified Stroke characters and to fill the current CJK Unified Ideographs blocks with the other characters that cannot be considered as BASIC stroke elements but should just be encoded as regular CJK Unified Ideographs.

This covers the large majority of the HKSCS PUA characters. The remaining cases are as following:

- 𠄎 (HKSCS code 8862), 𠄎 (HKSCS code 8864), 𠄎 (HKSCS code 88A3) and 𠄎 (HKSCS code 88A5) can already be encoded as sequence of base characters followed by combining characters. Therefore they don't need a new encoding. Furthermore, should they be encoded, they could never be used in any process involving character normalization. For example they would never be usable as part of an International Domain Name (IDN). Instead, the four composite sequences identified by the UCS Sequence identifiers should be formally identified in ISO 10646 as such.
- 𠄎 (HKSCS code 88A9: ground symbol) and 𠄎 (HKSCS code 88AA: fuse symbol) should be encoded as addition to the Miscellaneous Technical block.
- 𠄎 (HKSCS code C87E) looks very similar to the CJK Unified Ideograph U+4491 艹 which is a simplified view of the grass radical encoded itself as U+2F8B 艹. There are many grass radicals already encoded in the standard, either as part of the CJK Unified set or in other parts:

- 1) U+2F8B 艹 Kangxi Radical Grass
- 2) U+2EBE 艹 CJK Radical Grass One
- 3) U+2EBF 艹 CJK Radical Grass Two

- 4) U+2EC0 ^{††} CJK Radical Grass Three
- 5) U+8278 艸 CJK Unified Ideograph
- 6) U+8279 ^{††} CJK Unified Ideograph
- 7) U+4491 [≠] CJK Unified Ideograph
- 8) U+FA5D ^{††} CJK Compatibility Ideograph
- 9) U+FA5E ^{††} CJK Compatibility Ideograph

Fonts used for HKSCS such as “MingLiu_HKSCS” do not seem to contain a character for U+4491. Therefore it seems appropriate to encode that private use character as U+4491, unless the intended usage is a CJK radical which would have different character properties. However, encoding a fourth variant of the CJK Radical Grass seems unwise unless strong evidence is provided.

The following table shows the following information:

- HKSCS Code
- UCS PUA Code
- Representative glyphs as in MingLiu_HKSCS
- Proposed Code (xx represents the block position that would be used for the proposed Unified CJK Basic Stroke, all proposed code values are new characters except for U+4491 and composite sequences)
- Comment as appropriate

(The shaded area in the glyph cells shows the character bounding box, underlined proposed code shows addition)

HKSCS Code	PUA Code	Glyph	Proposed Code	Comment
8840	F303	↗	<u>31C0</u>	Proposed for CJK Unified Ext C-0001
8841	F304)	<u>31C1</u>	Proposed for CJK Unified Ext C-0518
8842	F305	↘	<u>31C2</u>	Proposed for CJK Unified Ext C-0519
8843	F306	↖	<u>9FB2</u>	Proposed for CJK Unified Ext C-0520, not encoded as a stroke character because it is a variant of HKSCS 8843
8844	F307	└	<u>31C3</u>	
8846	F309	┐	<u>31C4</u>	Proposed for CJK Unified Ext C-0400
8849	F30C	┘	<u>31C5</u>	HKSCS C87A already maps to similarly shaped U+200CC ㄣ
884A	F30D	フ	<u>31C6</u>	Proposed for CJK Unified Ext C-0402
884D	F310	┘	<u>31C7</u>	Proposed for CJK Unified Ext C-0399
884F	F312	ㄣ	<u>31C8</u>	Proposed for CJK Unified Ext C-0403
8850	F313	┐	<u>31C9</u>	Proposed for CJK Unified Ext C-0404
8851	F314	フ	<u>31CA</u>	Proposed for CJK Unified Ext C-0405
8852	F315	ㄣ	<u>31CB</u>	Proposed for CJK Unified Ext C-0406
8854	F317	┘	<u>31CC</u>	Proposed for CJK Unified Ext C-0407
8855	F318	ㄣ	<u>31CD</u>	Proposed for CJK Unified Ext C-0408
8862	F325	Ⓔ	<00CA,0304>	Already encoded as composite sequence

8864	F327	Ě	<00CA,030C>	Already encoded as composite sequence
88A3	F344	ē	<00EA,0304>	Already encoded as composite sequence
88A5	F346	ě	<00EA,030C>	Already encoded as composite sequence
88A9	F34A	𠄎	<u>23DA</u>	Miscellaneous Technical Block
88AA	F34B	𠄏	<u>23DB</u>	Miscellaneous Technical Block
8C43	F57A	熿	<u>9FA6</u>	Proposed for CJK Unified Ext C-13160
8C6D	F5A4	曙	<u>9FA7</u>	Proposed for CJK Unified Ext C-10543
8C74	F5AB	匡	<u>9FA8</u>	Proposed for CJK Unified Ext C-3169
8CB7	F5CC	蕲	<u>9FA9</u>	Proposed for CJK Unified Ext C-18342
8CB9	F5CE	𠄑	<u>9FAA</u>	Proposed for CJK Unified Ext C-19996
8CBB	F5D0	幹	<u>9FAB</u>	Proposed for CJK Unified Ext C-20303
8CC0	F5D5	鋼	<u>9FAC</u>	Proposed for CJK Unified Ext C-21171
8CD7	F5EC	騏	<u>9FAD</u>	Proposed for CJK Unified Ext C-22927
8CD8	F5ED	騏	<u>9FAE</u>	Proposed for CJK Unified Ext C-22952
8CDA	F5EF	鉞	<u>9FAF</u>	Proposed for CJK Unified Ext C-21016
C879	F7E5	ㄣ	<u>31CE</u>	Proposed for CJK Unified Ext C-0211
C87E	F7EA	𠄒	4491	Already encoded as CJK Unified Ideograph
C8A1	F7EB	𠄓	<u>9FB0</u>	Proposed for CJK Unified Ext C-3309
C8A3	F7ED	𠄔	<u>9FB1</u>	Proposed for CJK Unified Ext C-5373

Hong Kong HKSCS extension

In addition to the current version of the HKSCS, the Hong Kong Government has a process of character ‘approval’ by which new characters get added to a list for future inclusion in HKSCS. If the characters do not exist yet in ISO/IEC 10646, they are supposed to be submitted to the appropriate channel, for more details see <http://www.info.gov.hk/digital21/eng/hkscs/applcn.html>. Unfortunately in the meantime, such characters are allocated to the PUA. The current list at <http://www.info.gov.hk/digital21/eng/hkscs/download/newchar.pdf> contains three such characters.

Of those three characters, one seems to be already encoded in CJK Unified Extension B block as U+23551 𠄑.

The following table shows the following information:

- HKSCS Code
- UCS PUA Code
- Representative glyphs as in the HK web site.
- Proposed Code
- Comment as appropriate

(Underlined proposed code shows addition)

HKSCS Code	PUA Code	Glyph	Proposed Code	Comment
8CEB	F600	榊	23551	Serial Number 16, already encoded as U+23551 榊
8CED	F602	鐘	<u>9FB3</u>	Serial Number 17
8D48	F61C	壘	<u>9FB4</u>	Serial Number 44

It should be noted that the character coded in 23551 榊 is slightly different. It could be argued that they should be dis-unified.

GB-18030 Coded Character Set

GB-18030 also contains some Private Use Area characters and many of these characters can be categorized along the same criteria as the HKSCS set: CJK Basic Stroke, CJK Unified Ideographs, Latin letter and various symbols. Symbols and Latin letter can be categorized as follows:

- Vertical Variants (GB Codes A6D9-A6DF, A6EC-A6ED, A6F3). Should they be encoded, they should become part of a supplement block to the CJK Compatibility Block U+FE30-4F which is fully populated. The new block could be called ‘Vertical Forms’ and use the range U+FE10-U+FE1F.
- 𠄎 (GB Code A8BC) is a Latin character already encoded in U+1E3F.

The following table shows the following information:

- GB-18030 Code
- UCS PUA Code
- Representative glyphs as in SimSun-18030
- Proposed Code (xx represents the block position that would be used for the proposed Unified CJK Basic Stroke, yyy represents the block position that would be used for a possible CJK Compatibility Forms Supplement)
- Comment as appropriate

(The shaded area in the glyph cells shows the character bounding box, underlined proposed code shows addition)

GB18030 Code	PUA Code	Glyph	Proposed Code	Comment
A6D9	E78D	’	<u>FE10</u>	
A6DA	E78E	◦	<u>FE12</u>	Vertical variant of U+3002 ◦
A6DB	E78F	、	<u>FE11</u>	Vertical variant of U+3001 、
A6DC	E790	∴	<u>FE13</u>	
A6DD	E791	∵	<u>FE14</u>	
A6DE	E792	!	<u>FE15</u>	
A6DF	E793	?	<u>FE16</u>	
A6EC	E794	𠄎	<u>FE17</u>	Vertical variant of U+3016 𠄎

A6ED	E795	𠄎	<u>FE18</u>	Vertical variant of U+3017 𠄎
A6F3	E796	⋮	<u>FE19</u>	Vertical ellipsis, but can't use U+22EE as it does not decompose in <vertical>U+2026
A8BC	E7C7	𠄏	1E3F	
FE51	E816	𠄐	20087	
FE52	E817	𠄑	20089	
FE53	E818	𠄒	200CC	
FE59	E81E	𠄓	<u>9FB5</u>	
FE61	E826	𠄔	<u>9FB6</u>	
FE66	E82B	𠄕	<u>9FB7</u>	
FE67	E82C	𠄖	<u>9FB8</u>	Proposed for CJK Unified Ext C-0537
FE6C	E831	𠄗	215D7	
FE6D	E832	𠄘	<u>9FB9</u>	
FE76	E83B	𠄙	2298F	
FE7E	E843	𠄚	<u>9FBA</u>	Similar to U+20509 尖 but as high half component
FE90	E854	𠄛	<u>9FBB</u>	Similar to U+2099D 卓 but as left half component
FE91	E855	𠄜	241FE	
FEA0	E864	𠄝	<u>9FBC</u>	Similar to U+470C 𠄝 but as high half component

Conclusion

Addressing the PUA issues as suggested for these two important coded character set (HKSCS and GB-18030) would go a long way in facilitating data interchange across East Asia. It would also be extremely beneficial to bring the requirements for any new characters similar to these to the attention of the Unicode Technical Committee and ISO SC2/WG2 before taking the decision to encode them in the Private Use Area.