

Source:	Yi Bai(白易), CheonHyeong Sim(沈天珩)
Title:	Proposal to consider adding CodeCharts support for kIRG_KPSource representative glyphs in Unicode
Date:	2022-10-16
Action:	To be considered by UTC

This document proposes the addition of kIRG_KPSource representative glyphs by a font from DPRK source to the CodeCharts, which is missing in current Unicode version, and will complete representative glyphs from all sources.

Background

Unlike other sources, very limited DPRK-source ideographs are included in current CodeCharts (8 in the Extension C block, 106 in the CJK Compatibility Ideographs block, 50 in the CJK Compatibility Ideographs Supplement block and 1 in the Extension H block). The difficulty comes from that DPRK has been inactive in Unicode community for more nearly two decades. Recently, we acquire a font from an Android application, namely Okpyon, from DPRK source, which includes all glyphs that has kIRG_KPSource.

After careful examination of the font, we found that the font includes exact the same number of ideographs as in the kIRG_KPSource, and the glyph shape in the font is identical with the ones in CodeCharts, i.e. it is the same font included in the current CodeCharts for DPRK's glyphs. Thus we propose the font to be included in future Unicode CodeCharts as representative glyphs of kIRG_KPSource.

Proposal

We propose including the font in the code chart, which will be an improvement to the code charts as it will complete the task of displaying representative glyphs from all sources.

The font is already used for CJK Ext-C, CJK Compatibility Ideographs, CJK Compatibility Ideographs Supplement, and CJK Ext-H block. Thus we propose the font to be included in other blocks, namely CJK Unified Ideographs (URO), CJK Ext-A and CJK Ext-B in future Unicode versions.

	KP0	KP1	<i>KPU</i>	CodeCharts Coverage (as of Unicode15.0)
URO	4652	10359		No
ExtA	1	3188		No
ExtB		5766		No
ExtC		8		Yes
ExtH		1		Yes
Comp.		106	<i>1</i>	Yes
Comp. Supp.		49	<i>1</i>	Yes
<i>Unencoded</i>		83		/

Reconstruction of font's mapping

As the order and blocks in the font is different from normal Unicode orders, the glyphs in the font are not a simple mapping to Unicode codepoints. To give a better understand of the font, we will introduce the reconstruction of font's mapping to Unicode here briefly.

1. Layout of Hanja–part in the font

The Hanja–part in the font can be divided into three parts:



The first part is CJK Radical Supplement, there is 1 Hanja as shown in the picture above; (Figure 1)

The second part is CJK Unified Ideographs and CJK Unified Ideographs Extension A, there are 18815 Hanjas, which has a full coverage of all the KP1–characters in BMP except for U+431B(鰲) and all the KP0–characters according to Unihan. There are 617 Hanjas(included in the 18815 Hanjas) without a KP–source, and it seems to be out of scope of our reconstruction;

The third part is U+AXXX and Private Use Area. There are 6014 Hanjas(range: U+A000–U+ABFF and U+E000–U+EB7D), including U+431B which is actually mapped to U+E211 in the font, the characters from CJK–ExtB, CJK–ExtC, CJK–Comp, CJK–Comp Supplement, and the characters unincluded in Unicode. These characters are all from KP1. And there are also some characters(not necessarily Hanja) in PUA after U+EB7E, but these are also useless for our reconstruction. The sketch maps of the Hanjas from the second and the third part will be shown below. (Figures 2–4)

\$33AD	\$33AE	\$33AF	\$33B0	\$33B1	\$33B2	\$33B3	\$33B4	\$33B5	\$33B6	\$33B7	\$33B8	\$33B9	\$33BA	\$33BB	\$33BC	\$33BD	\$33BE	\$33BF	\$33C0	\$33C1	\$33C4	\$33C7	\$33C8
rad	rad%	rad%	ps	ns	μs	ms	pV	nV	μV	mV	kV	MV	pW	nW	μW	mW	kW	MW	kΩ	MΩ	cc	Co.	dB
\$33CA	\$33CB	\$33D2	\$33D6	\$33DD	\$33DE	\$33DF	\$340C	\$341C	\$3425	\$342D	\$3431	\$3433	\$343A	\$343C	\$3441	\$3444	\$3445	\$344A	\$344B	\$344C	\$3457	\$3459	
ha	HP	log	mol	Wb	Ym	A%	gal	苞	孰	鬯	𠙴	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇
\$345B	\$345C	\$345D	\$345E	\$345F	\$3463	\$3465	\$3466	\$3473	\$3474	\$3476	\$3477	\$3479	\$3478	\$347E	\$347F	\$3480	\$3481	\$3484	\$3485	\$348A	\$348D	\$3491	\$3492
𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇
\$3493	\$3494	\$3495	\$3499	\$349A	\$349B	\$349C	\$349D	\$349F	\$34A0	\$34A4	\$34A6	\$34A9	\$34B3	\$34B9	\$34BF	\$34D1	\$34D6	\$34DA	\$34E0	\$34E8	\$34F6	\$34F8	
𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇

(Figure 2)

\$9F7E	\$9F8D	\$9F8E	\$9F8F	\$9F90	\$9F91	\$9F92	\$9F93	\$9F94	\$9F95	\$9F96	\$9F98	\$9F9C	\$9F9D	\$9FA0	\$9FA1	\$9FA2	\$9FA3	\$9FA4	\$9FA5	\$A000	\$A001	\$A002	\$A003
麌	龍	龐	龔	龐	龔	龔	龔	龔	龔	龔	龔	龜	龜	龜	龜	龜	龜	龜	龜	丘	丈	甫	並
\$A004	\$A005	\$A006	\$A007	\$A008	\$A009	\$A00A	\$A00B	\$A00C	\$A00D	\$A00E	\$A00F	\$A010	\$A011	\$A012	\$A013	\$A014	\$A015	\$A016	\$A017	\$A018	\$A019	\$A01A	\$A01B
盟	壺	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇
\$A01C	\$A01D	\$A01E	\$A01F	\$A020	\$A021	\$A022	\$A023	\$A024	\$A025	\$A026	\$A027	\$A028	\$A029	\$A02A	\$A02B	\$A02C	\$A02D	\$A02E	\$A02F	\$A030	\$A031	\$A032	\$A033
𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇
\$A034	\$A035	\$A036	\$A037	\$A038	\$A039	\$A03A	\$A03B	\$A03C	\$A03D	\$A03E	\$A03F	\$A040	\$A041	\$A042	\$A043	\$A044	\$A045	\$A046	\$A047	\$A048	\$A049	\$A04A	\$A04B
𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇

(Figure 3)

\$F2CE	\$F2CF	\$F2D0	\$F2D1	\$F2D2	\$F2D3	\$F2D4	\$F2D5	\$F2D6	\$F2D7	\$F2D8	\$F2D9	\$F2DA	\$F2DB	\$F2DC	\$F2DD	\$F2DE	\$F2DF	\$F2E0	\$F2E1	\$F2E2	\$F2E3	\$F2E4	\$F2E5
𦨑	𦨒	𦨓	𦨔	𦨕	𦨖	𦨗	𦨘	𦨙	𦨚	𦨛	𦨜	𦨝	𦨞	𦨟	𦨟	𦨟	𦨟	𦨟	𦨟	𩶻	𩶻	𩶻	𩶻
\$F2E6	\$F2E7	\$F2E8	\$F2E9	\$F2EA	\$F2EB	\$F2EC	\$F2ED	\$F2EE	\$F2EF	\$F2FO	\$F2F1	\$F2F2	\$F2F3	\$F2F4	\$F2F5	\$F2F6	\$F2F7	\$F2F8	\$F2F9	\$F2FA	\$F2FB	\$F2FC	\$F2FD
𩶻	鷓	鶴	𩶻	𩶻	𩶻	𩶻	𩶻	𩶻	𩶻	𩶻	𩶻	𩶻	𩶻	𩶻	𩶻	𩶻	𩶻	𩶻	𩶻	𩶻	𩶻	𩶻	𩶻
\$F2FE	\$F2FF	\$F300	\$F301	\$F302	\$F303	\$F304	\$F305	\$F306	\$F307	\$F308	\$F309	\$F30A	\$F30B	\$F30C	\$F30D	\$F30E	\$F30F	\$F310	\$F311	\$F312	\$F313	\$F314	\$F315
𠂇	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	𠂇	𠂇	𠂇	𠂇
\$F316	\$F317	\$F318	\$F319	\$F31A	\$F31B	\$F31C	\$F31D	\$F31E	\$F31F	\$F320	\$F321	\$F322	\$F323	\$F324	\$F325	\$F326	\$F327	\$F328	\$F329	\$F32A	\$F32B	\$F32C	\$F32D
𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇
\$F32E	\$F32F	\$F330	\$F331	\$F332	\$F333	\$F334	\$F335	\$F336	\$F337	\$F338	\$F339	\$F33A	\$F33B	\$F33C	\$F33D	\$F33E	\$F33F	\$F340	\$F341	\$F342	\$F343	\$F344	\$F345
𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇
\$F346	\$F347	\$F348	\$F349	\$F34A	\$F34B	\$F34C	\$F34D	\$F34E	\$F34F	\$F350	\$F351	\$F352	\$F353	\$F354	\$F355	\$F356	\$F357	\$F358	\$F359	\$F35A	\$F35B	\$F35C	\$F35D
𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇	𠂇

(Figure 4)

As can be seen from the foregoing, there are $1 + (18815 - 617) + 6014 = 24213$ Hanjas;

In KPS10721:2000, the first codepoint is 0x3400, and the last one is 0x9294. There are $0x9294 - 0x3400 + 1 = 24213$ Hanjas.

According to Unihan, there are 4736 "holes"(codepoints are not continuous) in KP1. Meanwhile, there are 4653 Hanjas in KP0. If these Hanjas are filled into KP1, there will be $4736 - 4653 = 83$ holes left;

For the Hanjas mapped to U+AXXX and PUA, 106 of them are encoded in CJK-Comp, 1 of them is encoded in CJK-ExtA, 5766 of them are encoded in CJK-ExtB, 8 of them are encoded in CJK-ExtC, 49 of them are encoded in CJK-Comp Supplement,

1 of them is encoded in CJK-ExtH. So there should be 6014-106-1-5766-8-49-1=83 holes left;

There are 94 characters in IRGN897, and 13 of them are later encoded to CJK-ExtC or horizontally extended to CJK-ExtA or CJK-ExtB. And there are also 2 characters withdrawn from the initial proposal of DPRK-Comp Characters(KP1-441D and KP1-510B, see WG2 N2573 for reference). IRGN897 came after that so these 2 characters was certainly not included in. So there should be 94-13+2=83 holes.

All the numbers listed above can match up, which shows the credibility of the font.

2. Procedure of the reconstruction

After knowing the credibility of the font, and we are able to conjecture that all the KP0-Hanjas are included in KP1 by the number of the holes, we can start doing the reconstruction.

First, we can easily discover that the Hanjas in KP1 are sorted by the radical and the residual strokes. When both the radical and the RS are same, the Hanjas are sorted by the Korean pronunciation(DPRK-Order). For example: (Radical Grass+6, 0x6D82-0x6DC5)

0x6D8X	萸간	茳강	荅격	苗곡	茆공	芨교	芨교	芨구	芨규	茆기	芨길	茆과	芨괄	茆꽝
0x6D9X	芨뇨													
0x6DAX	苧자	苧전	苧족	苧종	苧증	苧지	苧천	苧천	苧천	苧천	苧종	苧치	苧치	苧채
0x6DBX	芨표	芨환	芨환	芨환	芨환									
0x6DCX	茆율	茆이	茆이	茆이	茆이									

The Hanjas with Radical Grass and RS 6 in KP0 are as follows: (Hanjas have already been arranged in order by the pronunciation in KP0)

荼다 茫망 苓명 茁복 茈수 茁순 茈자 草초 莉形 莉行 莉회 荒황 茁여 茁용 莉이 莉인 茁임

They can be filled into the holes just right.

0x6D8X	萸간	茳강	荅격	苗곡	茆공	芨교	芨교	芨구	芨규	茆기	芨길	茆과	芨괄	茆꽝
0x6D9X	芨뇨													
0x6DAX	苧자	苧전	苧족	苧종	苧증	苧지	苧천	苧천	苧천	苧천	苧종	苧치	苧치	苧채
0x6DBX	芨표	芨환	芨환	芨환	芨환									
0x6DCX	茆율	茆이	茆이	茆이	茆이									

Then we can do the filling for all the KP0-Hanjas.

Several known issues

1. About U+249D6(環)

There is no glyph for this Hanja but U+746F(璐) appears twice. It seems to be a mapping issue which led to a glyph mistake. See below:

0x594X	瑣감 瑣거 瑣개 瑣노 瑣단 瑣돌 瑣대 瑣모
0x595X	瑂미 瑔민 瑔변 瑔서 瑔선 瑶성 瑟슬 瑶전 瑶체 瑶창 瑶춘 瑶하 瑶호 瑶훈 瑶해
0x596X	瑂환 瑶황 瑶야 瑶연 瑶영 瑶우 瑶유 瑶위 瑶원 瑶각 瑶괴 瑶당 瑶도
0x597X	瑤류 瑶류 瑶률 瑶마 瑶방 瑶쇄 瑶쇄 瑶장 瑶조 瑶진 瑶진 瑶차 瑶침 瑶퇴 瑶영
0x598X	瑤오 瑶온 瑶요 瑶용

According to Unihan, KP1–594E should be U+249D6(瑣,환), but according to its pronunciation, it should be mapped to KP1–5961; According to Unihan, KP1–596F should be U+249E8(瑂), but there is not a same glyph in the font, we only have two "瑂" and a "瑣". While the glyph "瑣" refers to KP0–D7D7, we can conjecture that the glyph "瑂" is mapped to KP1–596F and U+249E8 for the cognition, but according to the RS, "瑂" seems to be more appropriate to be mapped to KP1–594E while "瑣" to KP1–596F.

2. Mapping issues

We discovered some mapping issues when doing the reconstruction. Some of them are completely wrong, and the rest of them are cognate to the current mappings but there are probably some better mappings. And we will divide these Hanjas into two parts to illustrate.

2.1. Wrong mappings

KP1–50FB(浑) is now mapped to U+23CD9(浑), but they are completely different.

0x50FX 汎술 汎색 汎생 汎자 汎沮 汎저 汎전 汎정 注주 汎지 浑집 沸제 沸찰 泉천 沾침
From the pronunciation, we could know that the glyph is not wrong but the mapping is wrong, because U+23CD9(浑) is read as 소. So it should be fixed to U+23CC0;

KP1–5B5D(瘡) is now mapped to U+24D6A(瘡), but they are completely different.

0x5B5X 瘡瘍 瘡疫 瘡우 瘡가 瘡감 瘡거 瘡고 瘡구 瘡근 瘡과 瘡닐 瘡단 瘡동 瘡凶 痘령 瘡법
From the pronunciation, we could know that the glyph is not wrong but the mapping is wrong, because U+24D6A(瘡) is a character used for a person's name in Taiwan or a Chữ Nôm in Việt Nam, its pseudo-Korean-pronunciation should be 소 or 조. So it should be fixed to U+3F94;

KP1–7EF4(酸) is now mapped to U+9166(酸), but they are completely different. Although the pronunciation for "酸" and "酸" are both 빨, firstly, from the two Hanjas illustrated before, it is more likely to be a mapping issue but not a glyph issue; secondly, "酸" is the simplified form of "醕" in the Mainland of China, it is unlikely to be used in DPRK, while "酸" is a traditional character without being simplified anywhere. So it should be fixed to U+48EE.

2.2. Suboptimal mappings

Some Hanjas are encoded separately with the different actual shape in Unicode, but they are cognate. For example, KP1–4A00(矇) is now mapped to U+66DA(矇), while its glyph is more close to U+232E1. But we cannot say it is wrong due to the cognition. They should probably be unified under the current unification rules. For these characters, we maintain a list of them in other sources^[1].

(Remark: The table lists all the mapping issues regardless of whether it is a wrong mapping or a suboptimal mapping.

For KP1–5450 and KP1–8346, they have a similar case as "鄜" mentioned before, whose current mapping is not the best but the best mapping has already been occupied by another Hanja.)

(End of document)

Ref.

- [1] CheonHyeong Sim. KP1–reconstitution.pdf[DB/OL]. <http://cheonhyeong.com/PDF/KP1–reconstitution.pdf>, 2022–06–19: PP.74–75.