

**** Formulas ****

The compositional character formulas represent the core of the Jizura project. I became aware of the 漢字データベース project sometime around 2007 or 2008; at the time, being able to download thousands upon thousands of formulas under the GPL license was what really jumpstarted my project.

The first years I spent mainly with sorting out the data, trying to make sense of it, correcting it where I found it to be faulty or problematic, and writing first attempts at software that would help me to handle the data, compile collections of derivative data, make it searchable and so on.

u-cjk/4e00 一 ●
 u-cjk/4e01 丁 ●一
 u-cjk/4e02 𠂇 ●一
 u-cjk/4e03 七 ①一
 u-cjk/4e04 上 ①一
 u-cjk/4e05 十 ①一
 u-cjk/4e06 ア ①一
 u-cjk/4e07 万 ①一
 u-cjk/4e08 丈 ①一
 u-cjk/4e09 三 ①一
 u-cjk/4e0a 上 ①一
 u-cjk/4e0b 下 ①一
 u-cjk/4e0c 𠂇 ①一
 u-cjk/4e0d 不 ①一
 u-cjk/4e0e 与 ①一
 u-cjk/4e0e 与 ①一
 u-cjk/4e0f 𠂇 ①一
 u-cjk/4e10 𠂇 ①一
 u-cjk/4e11 丑 ①一
 u-cjk/4e12 刃 ①一
 u-cjk/4e13 𠂇 ①一

[...]

u-cjk/4f9c 俯 ①一
 u-cjk/4f9d 依 ①一
 u-cjk/4f9e 伽 ①一
 u-cjk/4f9f 俯 ①一
 u-cjk/4fa0 俠 ①一
 u-cjk/4fa1 伽 ①一
 u-cjk/4fa2 伽 ①一
 u-cjk/4fa3 伽 ①一
 u-cjk/4fa4 伽 ①一
 u-cjk/4fa5 伽 ①一
 u-cjk/4fa6 伽 ①一
 u-cjk/4fa7 伽 ①一
 u-cjk/4fa8 伽 ①一
 u-cjk/4fa9 伽 ①一
 u-cjk/4faa 伽 ①一
 u-cjk/4fab 伽 ①一
 u-cjk/4fac 伽 ①一
 u-cjk/4fad 伽 ①一
 u-cjk/4fae 伽 ①一
 u-cjk/4faf 伽 ①一
 u-cjk/4fb0 伽 ①一
 u-cjk/4fb1 伽 ①一
 u-cjk/4fb2 伽 ①一

[...]

u-cjk/513d 儼 ①一
 u-cjk/513e 儼 ①一
 u-cjk/513f 儼 ①一
 u-cjk/5140 兀 ①一
 u-cjk/5141 允 ①一
 u-cjk/5142 允 ①一
 u-cjk/5143 元 ①一
 u-cjk/5144 兄 ①一

[...]

u-cjk/53ad 厭 ①一
 u-cjk/53ae 斯 ①一

u-cjk/53af 厭 ①一
 u-cjk/53af 厭 ①一
 u-cjk/53b0 廠 ①一
 u-cjk/53b1 廠 ①一
 u-cjk/53b2 廠 ①一
 u-cjk/53b3 廠 ①一
 u-cjk/53b4 廠 ①一
 u-cjk/53b5 廠 ①一
 u-cjk/53b6 廠 ①一
 u-cjk/53b7 廠 ①一

[...]

u-cjk-xb/20194 貌 ①一
 u-cjk-xb/20195 貌 ①一
 u-cjk-xb/20196 貌 ①一
 u-cjk-xb/20197 貌 ①一
 u-cjk-xb/20198 貌 ①一
 u-cjk-xb/20199 貌 ①一
 u-cjk-xb/2019a 貌 ①一
 u-cjk-xb/2019b 貌 ①一
 u-cjk-xb/2019c 貌 ①一
 u-cjk-xb/2019d 貌 ①一
 u-cjk-xb/2019e 貌 ①一
 u-cjk-xb/2019f 貌 ①一
 u-cjk-xb/201a0 貌 ①一
 u-cjk-xb/201a1 貌 ①一
 u-cjk-xb/201a2 貌 ①一
 u-cjk-xb/201a3 貌 ①一
 u-cjk-xb/201a4 貌 ①一
 u-cjk-xb/201a5 貌 ①一
 u-cjk-xb/201a6 貌 ①一
 u-cjk-xb/201a7 貌 ①一
 u-cjk-xb/201a8 貌 ①一
 u-cjk-xb/201a9 貌 ①一

[...]

u-cjk-xd/2b81b 齧 ①一
 u-cjk-xd/2b81c 齧 ①一
 u-cjk-xd/2b81d 齧 ①一

u-cjk-sym/3005 々 ①一
 u-cjk-sym/3006 々 ①一
 u-cjk-sym/3007 々 ①一
 u-cjk-sym/3021 一 ①一
 u-cjk-sym/3022 一 ①一
 u-cjk-sym/3023 一 ①一
 u-cjk-sym/3024 一 ①一
 u-cjk-sym/3025 一 ①一
 u-cjk-sym/3026 一 ①一
 u-cjk-sym/3027 一 ①一
 u-cjk-sym/3027 一 ①一
 u-cjk-sym/3028 一 ①一
 u-cjk-sym/3028 一 ①一
 u-cjk-sym/3029 一 ①一
 u-cjk-sym/3038 一 ①一
 u-cjk-sym/3039 一 ①一
 u-cjk-sym/303a 一 ①一
 u-cjk-sym/303b 一 ①一
 u-cjk-sym/303d 一 ①一

u-cjk-rad1/2f00 一 ①一
 u-cjk-rad1/2f01 一 ①一
 u-cjk-rad1/2f02 一 ①一

[...]

jzr/e100 申 ①一
 jzr/e100 申 ①一
 jzr/e100 申 ①一
 jzr/e101 鬼 ①一
 jzr/e102 艮 ①一

[...]

hzk1/b3d0 &hzk1#xb3d0; ①一
 hzk1/d044 &hzk1#xd044; ①一
 hzk1/d143 &hzk1#xd143; ①一
 hzk1/d2ea &hzk1#xd2ea; ①一

u-cjk/4e00	一	<1>
jzr/e167	/	<1>
u-cjk/4e8c	二	<11>
u-cjk-xb/2011f	二	<11>
u-cjk-xb/20120	二	<11>
u-cjk/4e09	三	<111>
u-cjk/4e96	三	<1111>
jzr/e15a	手	<11112>
u-cjk/8a00	言	<1111251>
u-cjk-xb/229d4	截	<1111342511534>
u-cjk/5f0e	式	<111154>
u-cjk/5f10	式	<111154>
u-cjk/4e30	圭	<1112>

札子五筆法	手	<11112>
jzr/e13f	𠂇	<11112>
jzr/e14d	𠂇	<11121112>
jzr/e225	𠂇	<111211121234>
u-cjk-xb/234f5	𠂇	<111211121251431>
u-cjk-xb/27bee	𠂇	<11121112125351>
u-cjk-xb/277f7	𠂇	<111211121234>
u-cjk-xa/3ece	𠂇	<1112111213434112431>
u-cjk-xb/28af9	𠂇	<111211121343412112431>
u-cjk-xb/24acd	𠂇	<111211121343445412121>
u-cjk-xb/24ad2	𠂇	<111211121344345412121>
u-cjk-xb/24ad2	𠂇	<1112111213445434>
u-cjk-xb/21689	𠂇	

1	→	1			→	1
1	→	1			→	1
1	→	1			→	1
1	→	11			→	11
1	→	11			→	11
1	→	11			→	11
1	→	11			→	111
1	→	11			→	1111
1	→	11			→	1111111111111111
1	→	11			→	11112
1	→	11			→	11112225211
1	→	11			→	1112
1	→	11			→	1112
1	→	11			→	1112
1	→	11			→	1112

jzr-e24d 三
jzr-e15a 丰
jzr-e2b3 萬
jzr-e14d 丰
jzr-e13f 丰
丰x2h jzr-e225 丰

$$[\dots]$$

1	—	12	jzr-e21a	↓	12
1	—	12	jzr-e21a	↓	125111152
1	—	12	jzr-e1a6	↓	12
1	—	12	jzr-e1a6	↓	1215

jzr-e21a 𠄎
jzr-e1a6 𠄎

1 — 12	丁 12
1 — 12	丁 12
1 — 12	丁 121
1 — 12	丁 121
# 1 — 12	丁 1212
1 — 12	丁 12121
1 — 12	丁 12121
1 — 12	丁 12122111

	丁		
	丁v1		丁
	工		
jzr-e2d6	工x4q		玨
	平		
	正		
jzr-e21c	正xlr		五
	章		

$$[\dots]$$

2	25	□	25
2	25	□	25
2	25	□	25
2	25	□	251
2	25	□	251
2	25	□	251
2	25	□	251
2	25	□	251
2	25	□	251
2	25	□	251
2	25	□	251
2	25	□	251
2	25	□	251
2	25	□	251
2	25	□	251
2	25	□	251
2	25	□	2511
2	25	□	2511

	v1	jzr-e2e3	
	x2c		
jzr-e352			
jzr-e365			
	x2c		
	x2h		
	x2v		
	x3h		
	x3p		
	x4q		
jzr-e2a9			
jzr-e2aa			
jzr-e2dd			

萬字一覽置換索引

5(34)		53	jzr-e1cb	↘	531		女x2v	
5(34)		53	jzr-e1cb	↘	531		女x3p	姦姦
5(34)		55		↘	55			
5(34)		55		↘	555			
5(34)		55		↘	555			
5(34)		55		↘	5551		x6q	𠂇𠂇
5(37))	5)	5		
5(47))	5)	5		
5(47)		\	5		\	5		
5(47)		\	53	jzr-e27c	↘	53	jzr-e27c	
5(47)		\	53	jzr-e27c	↘	534	jzr-e27a	
5(47)		\	5		\	5		
5(8)		=	5		=	5		
# ideographic space:								
5(8)		5			5			
5(8)		▣	5		▣	5		
5(8)		↗	5		↗	5		
5(8)		ε	5		ε	5		
5(8)		~	5		~	5		
5(8)		ø	5		ø	5		
5(8)		\$	5		\$	5		
5(8)		○	5		○	5		

(over 4,000 lines)

**** Figural Themes ****

[illegible][illegible]

** Shape Identity Mappings (SIMs) **

(22,000 lines)

To keep u-cjk-xb/2011e 二 ‘archaic form of 下 上’ as distinct from u-cjk/4e8c 二 ‘two’, which has the exact same shape, may be justified by the 《康熙字典》, which emulates the 《說文解字》 in this respect; however, doing so would be incompatible with our ‘shapes-only’ policy. I do not know what motivated the IRG to encode u-cjk-xb/201a2 人 (‘人-like, somewhat flattened shape that appears as top component’) as distinct from u-cjk/4eba 人 ‘man; human’; anyway, u-cjk-xb/201a2 人 does not seem to appear in 《康熙字典》 and is likewise incompatible with the ‘shapes-only’ policy.

u-cjk-xb/2011f 二 ‘archaic form of 下’ is in fact (if minimally) distinguishable from both u-cjk-xb/20120 二 ‘archaic form of 二 “two” and u-cjk/4e8c 二 ‘two’, so we keep these codepoints separate—except when writing out character formulas, where we always use u-cjk/4e8c 二 in an attempt to prevent becoming overly specific.

u-cjk/4eba 人 u-cjk-rad1/2f08 人
u-cjk/4eba 人 u-cjk-xb/201a2 人
u-cjk/4e8c 二 u-cjk-rad1/2f06 二
u-cjk/4e8c 二 u-cjk-xb/2011e 二
u-cjk/4e8c 二 u-cjk-xb/2011f 二!components
u-cjk/4e8c 二 u-cjk-xb/20120 二!components

The amount of duplication in Unicode is certainly surprising to the layman; some of this is due to Unicode’s ‘roundtrip encoding compatibility with legacy encodings’, some is deliberate, and some is due to oversight.

As far as CJK Ideographs are concerned, the two biggest single sources for duplicate codepoints are Compatibility Codepoints and the (I believe misguided) attempt to treat well-known ‘radicals’ (字典部首) and some characters explicitly labeled as ‘symbols’ apart from the ‘ideographs proper’. To do so, the Unicode Blocks u-cjk-rad1/2foo ... u-cjk-rad1/2fdf *Kangxi Radicals*, u-cjk-rad2/2e80 ... u-cjk-rad2/2eff *CJK Radicals Supplement* and u-cjk-sym/3000 ... u-cjk-sym/303f *CJK Symbols and Punctuation* have been set up; the latter contains, inter alia, u-cjk-sym/3007 〇 ‘zero’ (to me undoubtedly a ‘character’, just as 一, 二, 三, 萬 are characters, not symbols), duplicated as u-cjk-strk/31e3 〇 ‘(circular stroke, used for some ideographs created in Korea)’.

u-cjk-xb/2967f 食 u-cjk-rad2/2ede 食
u-cjk/98df 食 u-cjk-rad1/2fb7 食
u-cjk/98df 食 u-cjk-rad2/2edd 食
u-cjk/98e0 食 u-cjk-rad2/2edf 食
u-cjk/9963 𠂇 u-cjk-rad2/2ee0 𠂇
u-cjk/98e0 食 u-cjk-xb/2967f 食!components
u-cjk/91d1 金 u-cjk-cmpil/f90a 金
u-cjk/91d1 金 u-cjk-rad1/2fa6 金
u-cjk/91d1 金 u-cjk/91d2 金!components
u-cjk/8fb6 辶 u-cjk-cmpil/fa66 辶
u-cjk/8fb6 辶 u-cjk-rad2/2ecc 辶
u-cjk/8fb6 辶 u-cjk-rad2/2ecd 辶!components
u-cjk/8fb6 辶 u-cjk-rad2/2ece 辶!components
u-cjk/5ef4 辶 u-cjk-rad1/2f35 辶
u-cjk-sym/3023 𠂇 jzr-fig/e177 𠂇
u-cjk/4e28 | u-cjk-sym/3021 |

u-cjk/4e28 | u-cjk-rad1/2f01 |
u-cjk/4e28 | u-cjk-strk/31d1 |
u-cjk/4e85 | u-cjk-rad1/2f05 |
u-cjk/4e85 | u-cjk-strk/31da |
u-cjk/4e28 | u-cjk/4e85 |!components/search
u-cjk/4ea0 彳 u-cjk-rad1/2f07 彳
u-cjk/4ea0 彳 u-cjk-sym/3026 彳
u-cjk/5341 十 u-cjk-rad1/2f17 十
u-cjk/5341 十 u-cjk-sym/3038 十
u-cjk/535d 卩 u-cjk-cmpil/fa5d 卩
u-cjk/535d 卩 u-cjk-rad2/2ec0 卩
u-cjk-rad2/2ebf 卩 u-cjk-cmpil/fa5e 卩
u-cjk/5344 卩 u-cjk/8279 卩
u-cjk/5344 卩 u-cjk-rad2/2ebe 卩
u-cjk/5344 卩 u-cjk-sym/3039 卩
u-cjk-xb/25ad7 𠂇 u-cjk-rad2/2eae 𠂇
u-cjk-xb/2099d 卓 u-cjk/9fba 卓
u-cjk-xd/2b740 𠂇 jzr/e14f 𠂇
u-cjk-xc/2b1e6 𠂇 jzr/e182 𠂇
u-cjk/6075 患 u-cjk-cmpil/fa6b 患
u-cjk/8218 館 u-cjk-cmpil/fa6d 館
u-cjk/6ed1 滑 u-cjk-cmpil/f90d 滑
u-cjk/5951 契 u-cjk-cmpil/f909 契
u-cjk/61f6 懶 u-cjk-cmpil/f90d 懶
u-cjk/7669 癩 u-cjk-cmpil/f90e 癩
u-cjk/908f 邏 u-cjk-cmpil/f913 邏
u-cjk/4e82 亂 u-cjk-cmpil/f91b 亂
u-cjk/6feb 濫 u-cjk-cmpil/f922 濫
u-cjk/85cd 藍 u-cjk-cmpil/f923 藍
u-cjk/8964 檻 u-cjk-cmpil/f924 檻
u-cjk/5eca 廊 u-cjk-cmpil/f928 廊
u-cjk-xb/20628 𠂇 u-cjk-rad2/2e87 𠂇
u-cjk-xb/206a4 々 u-cjk-sym/3005 々
u-cjk-xb/208de 𠂇 u-cjk-cmpil/2f9dd 𠂇
u-cjk-xb/21018 𠂇 u-cjk-xb/2103c 𠂇
u-cjk-xb/21d0b 𠂇 u-cjk-cmpil/2f8f8 𠂇
u-cjk-xb/21f37 𠂇 u-cjk-xb/2439a 𠂇
u-cjk-xb/22331 𠂇 u-cjk-cmpil/2f891 𠂇
u-cjk-xb/232ad 𠂇 u-cjk-xb/232ab 𠂇
u-cjk-xb/23d1e 𠂇 u-cjk-cmpil/2f906 𠂇
u-cjk-xb/23f41 𠂇 u-cjk-xb/23f9e 𠂇
u-cjk-xb/23f5e 𠂇 u-cjk-cmpil/2f910 𠂇
u-cjk-xb/243ab 𠂇 u-cjk-cmpil/2f91f 𠂇
u-cjk-xb/24425 𠂇 u-cjk-xb/2444b 𠂇
u-cjk-xb/249bc 𠂇 u-cjk-xb/249e9 𠂇
u-cjk-xb/24c36 𠂇 u-cjk-cmpil/2f935 𠂇
u-cjk-xb/24d14 𠂇 u-cjk-rad2/2eaa 𠂇
u-cjk-xb/24fb8 𠂇 u-cjk-cmpil/2f93c 𠂇
u-cjk-xb/25249 𠂇 u-cjk-cmpil/fad5 𠂇
u-cjk/9091 邑 u-cjk-rad1/2fa2 邑
u-cjk/9149 酉 u-cjk-rad1/2fa3 酉

漢字データベース

jzr/e1ef 馬 cdp/896a &cdp#x896a;
jzr/e1f0 甘 cdp/88c8 &cdp#x88c8;
jzr/e1f1 𠂇 cdp/85fd &cdp#x85fd;
jzr/e1f8 𠂇 cdp/8569 &cdp#x8569;
u-cjk/5e80 𠂇 c2/215a &c2#x215a;
u-cjk/68b5 梵 c1/5b31 &c1#x5b31;
u-cjk/904b 運 c1/672a &c1#x672a;
u-cjk/9db3 𠂇 c2/6d34 &c2#x6d34;
u-cjk-xb/26ddf 𠂇 cb/15d5 &cb#x15d5;
u-cjk/5e97 店 c1/4d33 &c1#x4d33;
u-cjk-xa/387b 𠂇 c3/2866 &c3#x2866;
u-cjk/7385 𠂇 c2/2c28 &c2#x2c28;
u-cjk/9d19 𠂇 c2/5c39 &c2#x5c39;
u-cjk/6b63 正 c1/465f &c1#x465f;
u-cjk/9117 𠂇 c1/673a &c1#x673a;
u-cjk/7380 𠂇 c1/7b66 &c1#x7b66;

yadda yadda yadda yadda yadda yadda yadda yadda
yadda yadda yadda yadda yadda yadda yadda yadda
yadda yadda yadda yadda yadda yadda yadda yadda

yadda yadda yadda yadda yadda yadda yadda yadda yadda
yadda yadda yadda yadda yadda yadda yadda yadda yadda
yadda yadda yadda yadda yadda yadda yadda

**** Variants and Usage ****

This file contains around 8,400 records that connect ‘variant characters’ and ‘regions of usage’.

The regions as used here are C for the PRC, J for Japan, T for Taiwan, H for Hong Kong, M for Macau, and K for North and South Korea; this reflects both how these countries and territories are represented in the Unicode Consortium, and, in particular, in the data presented as a result of the Unicode IICore (International Ideographs Core, 國際表意文字核心, 東アジアの諸国で一般に使用される漢字集合)¹ effort.

The term ‘variant character’ is understood as an umbrella term for what is variously called 俗字, 古字, 略字, 異體字, 簡/繁體字, 新字体 and so on in the traditional literature; no effort has been made to differentiate beyond making simple statements like ‘glyph A [which is used in regions X...] is a variant of glyph B [which is used in regions Y...]’. In the software, variant relationships are modeled as both symmetric and transitive²; further, there is no temporal aspect whatsoever encoded. All of these assumptions make the display and the handling of the data simpler, but they also oversimplify somewhat.

Here are some samples:

01 也CJKTTHM 𠂇 亡 芒

02 鷗JKTHM 鷗C 鷗J
 03 飢JKTHM 饑KTHM 飢C
 04 個JKTHM 箇JKTh 个CJ 个
 05 團JKTHM 糰T 团C 团J 糰
 06 紂 紂
 07 龜JKTHM 龟C 龟J 𪚩 𪚩 𪚩 龜 龜 龜 龜
 龜 龜 龜 龜 龜 龜
 08 亞JKTHM 亚C 亜J
 09 畝JKTHM 亩C 畝 畝 畝 畝 畝
 10 儼cTHM 限cJt 混
 11 台CJKTTHM 𪚩JKTHM 臺KThM 檯TM 枱th 儼 龔
 𠂇 允 壺 壺 壺 壺 壺 壺 壺 壺
 12 元CJKTTHM 圓JKTHM 円JK 圓C 圓

Line 01 tells us that 𠂇, 亡 and 芒 are variants of 也. Of these, only 也 has a ‘usagecode’—CJKTTHM in this case—attached to it from which we can immediately see that the IICore team did not include any of 𠂇, 亡, 芒 in their listing of important characters; in other words, these glyphs are presumably of minor importance for everyday communication. Conversely, CJK-THM implies that the glyph 也 is used in all six regions under consideration.

In line 02, 鷗 is marked as used in the PRC and 鷗 as used in Japan; the other regions use 鷗, as it is marked JKTHM. Both 鷗 and 鷗 show up with J, the implication being that both are current in Japan, possibly for different usages (I guess 鷗 is used in names).

¹ <https://zh.wikipedia.org/wiki/國際表意文字核心>

² that is, when ‘A is a variant of B’ holds, then ‘B is a variant of A’ also holds, and when additionally ‘B is a variant of C’ is true, then ‘A is a variant of C’ is also true.