

- ★ a list of compositional formulas of glyphs;
- ★ a catalogue of ‘factors’, i.e. basic building blocks;

- ★ a way of ordering factors (using strokeorders);
- ★ an algorithm to derive XXXXXXXXXXXX

**** KWIC ****

version 1

占		占	<
占戈		占战	<
占义		占卤	<
占※戍口		占鹹	<
占※食		占鯨	<
占灬		占点	<
貝 占		貼	<
巾 占		帖	<
贝 占		贴	<
黑 占		點	<
钅 占		钻	<
亻 占		佔	<
魚 占		鮎	<
广 占		店	<
立 占		站	<
米 占		粘	<
氵 占		沾	<
氵 占※		滷	<
臣亼 占※皿		鹽	<
禾人水 占		黏	<

version 2

貝	占	貼	>
巾	占	帖	>
貝	占	貼	>
黑	占	點	>
钅	占	钻	>
亻	占	佔	>
魚	占	鮎	>
广	占	店	>
立	占	站	>
米	占	粘	>
氵	占	沾	>
禾人水	占	黏	>
	占戈	战	
	占乂	卤	
氵	占𣎵	灑	
	占𣎵戌口	鹹	
臣亠	占𣎵皿	鹽	
	占𣎵食	齏	
	占𣎵	點	

version 3

		占	占
禾人水		占	占
	貝	占	占
	巾	占	占
	貝	占	占
	黑	占	占
	钅	占	占
	亻	占	占
	魚	占	占
	广	占	占
	立	占	占
	米	占	占
	彳	占	占
		戈	占
		又	占
	彳	占	𠂔
		𠂔戍口	占
臣一		𠂔皿	占
		𠂔食	占
		𠂔	占

version 1

水	
𠂔 水	水
𠂔 水	泰
𠂔 水リ	录
𠂔 水リ	剥
尸 水 ^四 勺虫	剝
尸 水牛	屬
尸 水牛 ^レ	犀
石 水	遲
日共 水	碌
金 水	暴
金 ^レ 水	録
衤 水	錄
衤 ^レ 水	禄
纟 水	祿
纟 ^レ 水	绿
糸 ^レ 水	緑
禾人 水占	綠
廿月 ^共 水	黏
日日共 水	藤
禾勿人 水	曝
月木人 水	黎
火日共 水	膝
彡木人 水	爆
	漆

version 2

		水	
𠂔		水	
𠂔		水	
石	𠂔		水
日	共		水
金	𠂔		水
金	亼		水
衤	𠂔		水
衤	亼		水
彡	𠂔		水
糸	𠂔		水
糸	亼		水
廿	月	𠂔	
日	日	共	
禾	勹		人
月	木		人
火	日	共	
彳	木		人
彳	日	共	
禾	人		占
	𠂔		水
	亼		水
尸	水	勹	虫
尸	水		牛

version 3

		水		水
𠂔	𠂔	水		泰
日共	日共	水		暴
日日共	日日共	水		曝
火日共	火日共	水		爆
彡日共	彡日共	水		瀑
月木人	月木人	水		膝
彡木人	彡木人	水		漆
禾勾人	禾勾人	水		黎
廿月爻	廿月爻	水		藤
𠂔	𠂔	水		录
石𠂔	石𠂔	水		碌
金𠂔	金𠂔	水		録
衤𠂔	衤𠂔	水		禄
纟𠂔	纟𠂔	水		绿
糸𠂔	糸𠂔	水		緑
金亼	金亼	水		錄
衤亼	衤亼	水		祿
纟亼	纟亼	水		綠
禾人	禾人	占		黏
𠂔	𠂔	水リ		剥
亼	亼	水リ		剝
尸	尸	水四虫		屬
尸	尸	水牛		犀

彳|日共|水

瀑 <||

尸|水牛讠

遲 ||

尸|水牛讠

遲

version 1

方	方
方 ^人 其	方旗 <
方 ^人 矢	方族 <
方 ^人 氏	方旅 <
方 ^人 也	方施 <
方 ^人 疋	方旋 <
方 ^人 子讠	方遊 <
方令	方於 <
方攴	方放 <
土 方	方坊 <
卅 方	方芳 <
亻 方	方仿 <
彳 方	方仿 <
月 方	方肪 <
言 方	方訪 <
厶 方	方旁 <
讠 方	方訪 <
尸 方	方房 <
卩 方	方防 <
纟 方	方紡 <
女 方	方妨 <
钅 方 ^人 矢	方鏃 <
王 方 ^人 疋	方璇 <
彳 方 ^人 子	方游 <
甫 方攴	方敷 <
亻 方攴	方傲 <
白 方攴讠	方邀 <
木 ^四 方	方楞 <
木 ^厶 方	方榜 <
石 ^厶 方	方磅 <
虫 ^厶 方	方螃 <
钅 ^厶 方	方傍 <
亻 ^厶 方	方傍 <
彳 ^厶 方	方傍 <
月 ^厶 方	方膀 <
言 ^厶 方	方謗 <
亻 ^四 方	方悞 <
讠 ^厶 方	方謗 <
彳 方攴	方激 <
穴 方攴	方竅 <
纟 方攴	方繳 <
糸 方攴	方繳 <
自 穴 方讠	方邊 <

version 2

方	方
土 方	方坊
卅 方	方芳
亻 方	方仿 >
彳 方	方仿 >
月 方	方肪 >
言 方	方訪 >
厶 方	方旁 >
讠 方	方訪 >
尸 方	方房 >
卩 方	方防 >
纟 方	方紡 >
女 方	方妨 >
木 ^四 方	方楞 >
木 ^厶 方	方榜 >
石 ^厶 方	方磅 >
虫 ^厶 方	方螃 >
钅 ^厶 方	方傍 >
亻 ^厶 方	方傍 >
彳 ^厶 方	方傍 >
月 ^厶 方	方膀 >
言 ^厶 方	方謗 >
亻 ^四 方	方悞 >
讠 ^厶 方	方謗 >
方 ^人 其	方旗
方 ^人 矢	方族
钅 方 ^人 矢	方鏃
方 ^人 氏	方旅
方 ^人 也	方施
方 ^人 疋	方旋
王 方 ^人 疋	方璇
彳 方 ^人 子	方游
方 ^人 子讠	方遊
方令	方於
方攴	方放
甫 方攴	方敷
亻 方攴	方傲 >
彳 方攴	方傲 >
白 方攴	方邀 >
穴 方攴	方邀 >
纟 方攴	方邀 >
糸 方攴	方邀 >
白 方攴讠	方邀 >
自 穴 方讠	方邊

version 3

方	方
土 方	方坊
卅 方	方芳
木 ^四 方	方楞
亻 ^四 方	方悞
亻 方	方仿
彳 方	方仿
月 方	方肪
言 方	方訪
厶 方	方旁
木 ^厶 方	方榜
石 ^厶 方	方磅
虫 ^厶 方	方螃
钅 ^厶 方	方傍
亻 ^厶 方	方傍
彳 ^厶 方	方傍
月 ^厶 方	方膀
言 ^厶 方	方謗
讠 ^厶 方	方謗
尸 方	方房
卩 方	方防
纟 方	方紡
女 方	方妨
方 ^人 其	方旗
方 ^人 矢	方族
钅 方 ^人 矢	方鏃
方 ^人 氏	方旅
方 ^人 也	方施
方 ^人 疋	方旋
王 方 ^人 疋	方璇
彳 方 ^人 子	方游
方 ^人 子讠	方遊
方令	方於
方攴	方放
甫 方攴	方敷
彳 方攴	方傲
彳 方攴	方傲
白 方攴	方邀
穴 方攴	方邀
纟 方攴	方邀
糸 方攴	方邀
亻 方攴	方邀
白 方攴讠	方邀
自 穴 方讠	方邊

True to the original ideas of the KWIC principle, each glyph appears as many times in the index as it has factors. Within the 15,000 or so most common characters, the maximum number of factors per glyph is 6, while the average is a little under 3.

|方^人子讠
方|^人子讠
方^人|子讠
方^人子|讠

遊
遊
遊
遊

|尸米四勺虫
 尸|米四勺虫
 尸米|四勺虫
 尸米四|勺虫
 尸米四勺|虫

屬屬屬屬屬

		木	缶	木	一	𣪠	𣪠	鬱
	木		缶	木	一	𣪠	𣪠	鬱
	木	缶		木	一	𣪠	𣪠	鬱
	木	缶	木		一	𣪠	𣪠	鬱
	木	缶	木	一		𣪠	𣪠	鬱
	木	缶	木	一	𣪠		𣪠	鬱

檉木
 檉木
 檉木
 檉木
 檉木
 檉木

**** 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 this.

The formula language used here is based on the Ideographic Description Language (IDL) as proposed by the Unicode Consortium,¹ with a few extensions; compared to other approaches,² IDL has the advantage of being at the right level of abstraction for our purposes and being both human-readable and syntactically straightforward.

u-cjk/4e00 一 ●
 u-cjk/4e01 丁 ●一丿
 u-cjk/4e02 乚 ●一乚
 u-cjk/4e03 乚 ●一乚
 u-cjk/4e04 乚 ●一乚
 u-cjk/4e05 乚 ●一乚
 u-cjk/4e06 乚 ●一乚
 u-cjk/4e07 万 ●一乚
 u-cjk/4e08 丈 ●ナ、
 u-cjk/4e09 三 ●一二
 u-cjk/4e0a 上 ●ト一
 u-cjk/4e0b 下 ●ト一
 u-cjk/4e0c 丌 ●一乚
 u-cjk/4e0d 丌 ●アト
 u-cjk/4e0e 与 ●与一
 u-cjk/4e0e 与 ●●与一一
 u-cjk/4e0f 𠂇 ●●丁丿
 u-cjk/4e10 𠂇 ●下与
 u-cjk/4e11 丑 ●丿土
 u-cjk/4e12 𠂇 ●刃一
 u-cjk/4e13 𠂇 ●二●乚、

[...]

u-cjk/4f9c 俯 ●一舟
 u-cjk/4f9d 依 ●一衣
 u-cjk/4f9e 伽 ●一如
 u-cjk/4f9f 俯 ●一存
 u-cjk/4fa0 侠 ●一夹
 u-cjk/4fa1 伽 ●一面
 u-cjk/4fa2 伽 ●一再
 u-cjk/4fa3 伽 ●一吕
 u-cjk/4fa4 伽 ●一考
 u-cjk/4fa5 伽 ●一尧
 u-cjk/4fa6 伽 ●一贞
 u-cjk/4fa7 伽 ●一则
 u-cjk/4fa8 伽 ●一乔
 u-cjk/4fa9 伽 ●一会
 u-cjk/4faa 伽 ●一齐
 u-cjk/4fab 伽 ●一妄
 u-cjk/4fac 伽 ●一农
 u-cjk/4fad 伽 ●一尽
 u-cjk/4fae 伽 ●一每
 u-cjk/4faf 侯 ●一●コ矢
 u-cjk/4fb0 伽 ●一君
 u-cjk/4fb1 伽 ●一呈
 u-cjk/4fb2 伽 ●一辰

[...]

u-cjk/513d 儼 ●一彙
 u-cjk/513e 儼 ●一囊
 u-cjk/513f 儿 ●ノ乚
 u-cjk/5140 兀 ●一儿
 u-cjk/5141 允 ●厶儿
 u-cjk/5142 允 ●
 u-cjk/5143 元 ●一兀
 u-cjk/5144 兄 ●口儿

[...]

u-cjk/53ad 厭 ●尸獸
 u-cjk/53ae 廝 ●尸斯
 u-cjk/53af 厖 ●麻心
 u-cjk/53af 厖 ●麻&cdp#x8962;
 u-cjk/53b0 廠 ●尸敞
 u-cjk/53b1 廋 ●尸僉
 u-cjk/53b2 厲 ●尸萬
 u-cjk/53b3 廠 ●尸敢
 u-cjk/53b4 厖 ●厭甲
 u-cjk/53b5 廋 ●原原
 u-cjk/53b6 厶 ●厶、
 u-cjk/53b7 厶 ●ナム

[...]

u-cjk-xb/20194 覯 ●享兒
 u-cjk-xb/20195 𠂇 ●并夜
 u-cjk-xb/20196 𠂇 (●一●方氏合)
 u-cjk-xb/20197 𠂇 ●旅合
 u-cjk-xb/20198 𠂇 (●一八口一衣)
 u-cjk-xb/20199 𠂇 ●享夜
 u-cjk-xb/2019a 𠂇 (●一𠂇死)
 u-cjk-xb/2019b 𠂇 ●一𠂇
 u-cjk-xb/2019c 𠂇 ●旅妻
 u-cjk-xb/2019d 𠂇 (●一●日日八人戊)
 u-cjk-xb/2019e 𠂇 ●享宜
 u-cjk-xb/2019f 𠂇 (●一𠂇𠂇𠂇)
 u-cjk-xb/201a0 𠂇 (●一●𠂇同百𠂇)
 u-cjk-xb/201a1 𠂇 (●一●𠂇同然)
 u-cjk-xb/201a2 𠂇 ▽
 u-cjk-xb/201a3 𠂇 ●一人
 u-cjk-xb/201a4 𠂇 ●人丁
 u-cjk-xb/201a5 𠂇 ●人ノ
 u-cjk-xb/201a6 𠂇 ●一●人一
 u-cjk-xb/201a7 𠂇 ●一人
 u-cjk-xb/201a8 𠂇 ●一了
 u-cjk-xb/201a9 𠂇 ●一凡

[...]

¹ Unicode Specification V6.0, Ch. 12 <http://www.unicode.org/versions/Unicode6.0.0/ch12.pdf>
² https://en.wikipedia.org/wiki/Chinese_character_description_languages

u-cjk-xd/2b81b 𪔐 𪔐齒星
 u-cjk-xd/2b81c 𪔑 𪔑齒兒
 u-cjk-xd/2b81d 𪔒 𪔒敵龜

u-cjk-sym/3005 𠂇 ▽
 u-cjk-sym/3006 𠂈 ●
 u-cjk-sym/3007 〇 ●
 u-cjk-sym/3021 | ▽
 u-cjk-sym/3022 𠂉 〇 | |
 u-cjk-sym/3023 𠂊 〇 | |
 u-cjk-sym/3024 𠂋 ▽
 u-cjk-sym/3025 𠂌 ●
 u-cjk-sym/3026 𠂍 ▽
 u-cjk-sym/3027 𠂎 〇 一
 u-cjk-sym/3027 𠂎 〇 二
 u-cjk-sym/3028 𠂏 〇 二
 u-cjk-sym/3028 𠂏 〇 三
 u-cjk-sym/3029 𠂐 〇 又
 u-cjk-sym/3038 十 ▽
 u-cjk-sym/3039 𠂑 ▽
 u-cjk-sym/303a 𠂒 ▽
 u-cjk-sym/303b 𠂓 ●
 u-cjk-sym/303d 𠂔 ●

u-cjk-rad1/2f00 一 ▽
 u-cjk-rad1/2f01 | ▽
 u-cjk-rad1/2f02 丶 ▽

[...]

jzr/e100 𠂕 〇 甲
 jzr/e100 𠂕 〇 白 𠂕
 jzr/e100 𠂕 〇 由
 jzr/e101 𠂖 〇 𠂖
 jzr/e102 𠂗 〇 𠂗

[...]

hzk1/b3d0 &hzk1#xb3d0; 𠂙 𠂙 𠂙 土 方 文
 hzk1/d044 &hzk1#xd044; 𠂚 𠂚 𠂚 佳 乃
 hzk1/d143 &hzk1#xd143; 𠂛 𠂛 𠂛 𠂛 𠂛 用
 hzk1/d2ea &hzk1#xd2ea; 𠂜 𠂜 黃

[...]

** Strokeorders (札字五筆法) **		
u-cjk/4e00	一	<1>
jzr/e167	丿	<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	𠂇	<11112>
u-cjk/8a00	言	<1111251>
u-cjk-xb/229d4	𦇧	<1111342511534>
u-cjk/5f0e	𠂇	<111154>
u-cjk/5f10	𠂇	<111154>
u-cjk/4e30	丰	<1112>
jzr/e13f	𠂇	<1112>
jzr/e14d	𠂇	<1112>
jzr/e225	𠂇	<11121112>
u-cjk-xb/234f5	𠂇	<111211121234>
u-cjk-xb/27bee	𠂇	<111211121251431>
u-cjk-xb/277f7	𠂇	<11121112125351>
u-cjk-xa/3ece	𠂇	<11121112134>
u-cjk-xb/28af9	𠂇	<1112111213434112431>
u-cjk-xb/24acd	𠂇	<111211121343412112431>
u-cjk-xb/24ad2	𠂇	<111211121343445412121>
u-cjk-xb/24ad2	𠂇	<111211121344345412121>
u-cjk-xb/21689	𠂇	<1112111213445434>

萬字一覽置換索引

2 25	冂 2511	日x4q	𠂔
2 25	冂 2511	日	

[...]

3 / 31	厶 31	厶	
3 / 31	厶 31115	𠂔	
3 / 31	厶 3112	午	
3 / 31	厶 3112	牛	
3 / 31	厶 3112	牛x2v	𠂔
3 / 31	厶 3112	牛x4q	𠂔
3 / 31	厶 31121	生	
3 / 31	厶 31121	生x2h	𠂔
3 / 31	厶 311212	年	
3 / 31	厶 3112121	𠂔	
3 / 31	厶 31122221	jzr-e18d 無	
3 / 31	厶 311234	朱	
3 / 31	厶 31125	jzr-e353 用	
3 / 31	厶 311252	jzr-e175 𠂔	
3 / 31	厶 311252	𠂔	
3 / 31	厶 31134	失	
3 / 31	厶 31134	矢	
3 / 31	厶 3115	气	
3 / 31	厶 312	jzr-e13d 𠂔	
3 / 31	厶 312	jzr-e13e 𠂔	
3 / 31	厶 312	jzr-e13d 𠂔x2h	竹
3 / 31	厶 312	jzr-e13e 𠂔x3p	𠂔
3 / 31	厶 312	jzr-e13e 𠂔x6p	𠂔
3 / 31	厶 312	jzr-e13e 𠂔x8q	𠂔
3 / 31	厶 3121	生	
3 / 31	厶 3121	𠂔	
3 / 31	厶 31211	𠂔	

[...]

4 \ 41	ㄣ 41	ㄣ	
4 \ 41	ㄣ 4111251	言	
4 \ 41	ㄣ 4111251	言x2h	𠂔
4 \ 41	ㄣ 4111251	言x2v	𠂔
4 \ 41	ㄣ 4111251	言x4q	𠂔
4 \ 41	ㄣ 41121	主	
4 \ 41	ㄣ 412234	jzr-e11a 亦	
4 \ 41	ㄣ 4124	𠂔	
4 \ 41	ㄣ 4125125251	高	
4 \ 41	ㄣ 413	广	

[...]

5(12)	冂 5	冂 5	冂
5(12)	冂 51 jzr-e346	冂 51 jzr-e346	冂
5(12)	冂 51 jzr-e346	冂 5121	𠂔
5(12)	冂 51 jzr-e350	冂 51	jzr-e350 𠂔
5(12)	冂 51 jzr-e350	冂 5121121211	jzr-e31e 𠂔
5(12)	冂 51 jzr-e1c2	冂 51	jzr-e1c2 𠂔
5(12)	冂 51 jzr-e1c2	冂 5122111	jzr-e21b 𠂔
5(12)	冂 51 jzr-e1c2	冂 5134	𠂔
5(12)	冂 51 jzr-e341	𠂔 51 jzr-e341	𠂔
5(12)	冂 51 jzr-e341	𠂔 511	𠂔
5(12)	冂 51 jzr-e341	𠂔 511	𠂔x2v jzr-e26d 𠂔

5(12)	冂	51	jzr-e341	𠂇	511112	聿
5(12)	冂	51	jzr-e341	𠂇	51112	聿
5(12)	冂	51	jzr-e341	𠂇	511121	聿
5(12)	冂	51	jzr-e341	𠂇	5112	聿
5(12)	冂	51	jzr-e341	𠂇	5112	聿
5(12)	冂	51	jzr-e341	𠂇	51121	聿
5(12)	冂	51	jzr-e341	𠂇	5112234	聿
5(12)	冂	51	jzr-e341	𠂇	51123134	聿

[...]

5(34)	𠂇	5	𠂇	5	𠂇	
5(34)	𠂇	53	jzr-e1cb	𠂇	53	jzr-e1cb
5(34)	𠂇	53	jzr-e1cb	𠂇	531	女
5(34)	𠂇	53	jzr-e1cb	𠂇	531	女x2h
5(34)	𠂇	53	jzr-e1cb	𠂇	531	女x2v
5(34)	𠂇	53	jzr-e1cb	𠂇	531	女x3p
5(34)	𠂇	55		𠂇	55	𠂇
5(34)	𠂇	55		𠂇	555	𠂇
5(34)	𠂇	55		𠂇	555	𠂇x6q
5(34)	𠂇	55		𠂇	5551	𠂇
5(37)	𠂇	5		𠂇	5	𠂇
5(47)	𠂇	5		𠂇	5	𠂇
5(47)	𠂇	5		𠂇	5	𠂇
5(47)	𠂇	53	jzr-e27c	𠂇	53	jzr-e27c
5(47)	𠂇	53	jzr-e27c	𠂇	534	jzr-e27a
5(47)	𠂇	5		𠂇	5	𠂇
5(8)	𠂇	5		𠂇	5	𠂇
# ideographic space:						
5(8)	𠂇	5		𠂇	5	𠂇
5(8)	𠂇	5		𠂇	5	𠂇
5(8)	𠂇	5		𠂇	5	𠂇
5(8)	𠂇	5		𠂇	5	𠂇
5(8)	𠂇	5		𠂇	5	𠂇
5(8)	𠂇	5		𠂇	5	𠂇
5(8)	𠂇	5		𠂇	5	𠂇
5(8)	𠂇	5		𠂇	5	𠂇

(over 4,000 lines)

**** Figural Themes ****

hemes **

(22,000 lines)

**** Shape Identity Mappings (SIMs) ****

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/2012a 人 ‘(人-like, somewhat flattened shape that appears as top component)’ as distinct from u-cjk/4eba 人 ‘man; human’; anyway, u-cjk-xb/2012a 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/2012a 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/2012a 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-cmp11/f90a 金
u-cjk/91d1 金 u-cjk-rad1/2fa6 金

u-cjk/91d1 金 u-cjk/91d2 金!components
u-cjk/8fb6 𠂇 u-cjk-cmp11/fa66 𠂇
u-cjk/8fb6 𠂇 u-cjk-rad2/2ecc 𠂇
u-cjk/8fb6 𠂇 u-cjk-rad2/2ecd 𠂇!components
u-cjk/8fb6 𠂇 u-cjk-rad2/2ece 𠂇!components
u-cjk/5ef4 𠂇 u-cjk-rad1/2f35 𠂇
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 J u-cjk-rad1/2f05 J
u-cjk/4e85 J u-cjk-strk/31da J
u-cjk/4e28 | u-cjk/4e85 J!components/search
u-cjk/4ea0 𠂇 u-cjk-rad1/2f07 𠂇
u-cjk/4ea0 𠂇 u-cjk-sym/3026 𠂇
u-cjk/5341 + u-cjk-rad1/2f17 +
u-cjk/5341 + u-cjk-sym/3038 +
u-cjk/535d 𠂇 u-cjk-cmp11/fa5d 𠂇
u-cjk/535d 𠂇 u-cjk-rad2/2ec0 𠂇
u-cjk-rad2/2ebf 𠂇 u-cjk-cmp11/fa5e 𠂇
u-cjk/5344 𠂇 u-cjk/8279 𠂇
u-cjk/5344 𠂇 u-cjk-rad2/2ebe 𠂇
u-cjk/5344 𠂇 u-cjk-sym/3039 𠂇
u-cjk-xb/25ad7 𠂇 u-cjk-rad2/2eae 𠂇
u-cjk-xb/2099d 𠂇 u-cjk/9fba 𠂇
u-cjk-xd/2b740 𠂇 jzr/e14f 𠂇
u-cjk-xc/2b1e6 𠂇 jzr/e182 𠂇
u-cjk/6075 惠 u-cjk-cmp11/fa6b 惠
u-cjk/8218 𠂇 u-cjk-cmp11/fa6d 𠂇
u-cjk/6ed1 滑 u-cjk-cmp11/f904 滑
u-cjk/5951 契 u-cjk-cmp11/f909 契
u-cjk/61f6 懶 u-cjk-cmp11/f90d 懶
u-cjk/7669 癩 u-cjk-cmp11/f90e 癩
u-cjk/908f 邏 u-cjk-cmp11/f913 邏
u-cjk/4e82 亂 u-cjk-cmp11/f91b 亂
u-cjk/6feb 濫 u-cjk-cmp11/f922 濫
u-cjk/85cd 藍 u-cjk-cmp11/f923 藍
u-cjk/8964 檻 u-cjk-cmp11/f924 檻
u-cjk/5eca 廊 u-cjk-cmp11/f928 廊
u-cjk-xb/20628 𠂇 u-cjk-rad2/2e87 𠂇
u-cjk-xb/206a4 𠂇 u-cjk-sym/3005 𠂇
u-cjk-xb/208de 𠂇 u-cjk-cmp12/2f9dd 𠂇
u-cjk-xb/21018 𠂇 u-cjk-xb/2103c 𠂇
u-cjk-xb/21d0b 𠂇 u-cjk-cmp12/2f8f8 𠂇
u-cjk-xb/21f37 𠂇 u-cjk-xb/2439a 𠂇
u-cjk-xb/22331 𠂇 u-cjk-cmp12/2f891 𠂇
u-cjk-xb/232ad 𠂇 u-cjk-xb/232ab 𠂇
u-cjk-xb/23d1e 𠂇 u-cjk-cmp12/2f906 𠂇
u-cjk-xb/23f41 𠂇 u-cjk-xb/23f9e 𠂇
u-cjk-xb/23f5e 𠂇 u-cjk-cmp12/2f910 𠂇
u-cjk-xb/243ab 𠂇 u-cjk-cmp12/2f91f 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-cmp12/2f935 𠂇
u-cjk-xb/24d14 𠂇 u-cjk-rad2/2eaa 𠂇
u-cjk-xb/24fb8 𠂇 u-cjk-cmp12/2f93c 𠂇
u-cjk-xb/25249 𠂇 u-cjk-cmp11/fad5 𠂇
u-cjk/9091 邑 u-cjk-rad1/2fa2 邑
u-cjk/9149 酉 u-cjk-rad1/2fa3 酉

萬字一覽置換索引

漢字データベース

jzr/e1ef 馬 cdp/896a &cdp#x896a;
jzr/e1f0 甘 cdp/88c8 &cdp#x88c8;
jzr/e1f1 亞 cdp/85fd &cdp#x85fd;
jzr/e1f8 兩 cdp/8569 &cdp#x8569;
u-cjk/5e80 庀 c2/215a &c2#x215a;
u-cjk/68b5 梵 c1/5b31 &c1#x5b31;
u-cjk/904b 運 c1/672a &c1#x672a;
u-cjk/9db3 𨔵 c2/6d34 &c2#x6d34;
u-cjk-xb/26ddf 蒹 cb/15d5 &cb#x15d5;
u-cjk/5e97 店 c1/4d33 &c1#x4d33;
u-cjk-xa/387b 廊 c3/2866 &c3#x2866;

u-cjk/7385 紗 c2/2c28 &c2#x2c28;
u-cjk/9d19 鴿 c2/5c39 &c2#x5c39;
u-cjk/6b63 正 c1/465f &c1#x465f;
u-cjk/9117 鄙 c1/673a &c1#x673a;
u-cjk/7380 羅 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 yadda yadda yadda yadda yadda

**** Dictionary Data ****

The relatively little dictionary data I have comes largely from Unicode's UniHan.txt data and Jim Breene's EDICT project³. Overall, it is of very mixed quality and would be in need of a thorough overhaul.

@glyphs 冰冰

u-cjk/51b0 冰 py:bīng

u-cjk/51b0 冰 ka:ヒョウ

u-cjk/51b0 冰 hi:こおり, ひ, こおる

u-cjk/51b0 冰 hg:빙

u-cjk/51b0 冰 gloss:ice; ice-cold

@glyphs 氫氢

u-cjk/6c2b 氢 py:qīng

u-cjk/6c2b 氢 gloss:amonia; hydrogen nitride

@glyphs 況况

u-cjk/6cc1 況 py:kuàng

u-cjk/6cc1 況 ka:キョウ

u-cjk/6cc1 況 hi:まし・て, いわ・んや, おもむき

u-cjk/6cc1 況 hg:황

u-cjk/6cc1 況 gloss:condition, situation;

furthermore

@glyphs 渾渾

u-cjk/6d51 浑 py:hún

u-cjk/6e3e 渾 ka:コン

u-cjk/6e3e 渾 hi:すべ・て, にご・る

u-cjk/6e3e 渾 hg:혼

u-cjk/6e3e 渾 gloss:muddy, turbid; blend, merge, mix

@glyphs 游遊

u-cjk/6e38 游 py:yóu

u-cjk/6e38 游 ka:ユウ, リュウ

u-cjk/904a 遊 ka:ユウ, ユ

u-cjk/6e38 游 hi:あそ・(び|ぶ), およ・ぐ

u-cjk/904a 遊 hi:あそ・(ぶ|ばす)

u-cjk/6e38 游 hg:유

u-cjk/6e38 游 gloss:to swim; float, drift; wander, roam

³ http://www.edrdg.org/wwwjdic/wwwjdicinf.html#dicfil_tag

**** 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)⁴ effort.

The term ‘variant character’ is understood as an umbrella term for what is variously called 俗字, 古字, 本字, 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 transitive⁵; 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CJKTm 𠂇 𠂈 𠂉
02 鷗JKThm 鷗C 鷗J

03 飢JKThm 饑KThm 飢C
04 個JKThm 箇JKThm 个CJ 𠂇 𠂈 𠂉
05 團JKThm 糰T 团C 团J 糰
06 紂 紂
07 龜JKThm 龟C 龜J 𪚩 𪚪 𪚫 龜 龜 龜 龜
龜 龜 龜 龜 龜 龜
08 亞JKThm 亞C 亜J
09 畝JKThm 亩C 畝 畝 畝 畝 畝
10 假cThm 限cJt 混
11 台CJKTm 廳JKThm 臺KThm 檯TM 枱th 儼 龔
𡗗 允 壺 壺 壺 壺 壺 壺 壺 壺 壺
12 元CJKTm 圓JKThm 円JK 圓C 圓

Line 01 tells us that 𠂇, 𠂈 and 𠂉 are variants of 也. Of these, only 也 has a ‘usagecode’—CJKTm in this case—attached to it from which we can immediately see that the IICore team did not include any of 𠂇, 𠂈, 𠂉 in their listing of important characters; in other words, these glyphs are presumably of minor importance for everyday communication. Conversely, CJKThm 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).

⁴ <https://zh.wikipedia.org/wiki/國際表意文字核心>

⁵ 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.

**** Frequency and Rankings ****

Using data from the Leeds Corpus⁶ for Chinese (PRC) and Japanese character usages, a frequency list compiled by Chih-Hao Tsai⁷ and a NodeJS module⁸ that promises to perform a reasonable ranking using Bayesian statistics, a little over 15,600 characters mentioned in at least one of the above sources were given a ranking index, starting from #00,001 for the most frequent character, 的, to #15,677 for the least common one, 气.

It must be said that the procedure as outlined above is probably not very scientific; it does not take the IICore selection of characters into account and neglects potentially useful material such as counts derived from Wikipedia articles, the lists published in connection with language proficiency tests such as the JLPT⁹ or the HSK¹⁰. On the other hand, there is no way to make any meaningful usage statistic without handling huge corpuses and making decisions about which texts to include and exclude from the counts.

As it stands, our listing does seem to make some intuitive sense, and has been a valuable tool to facilitate the sifting out of tens of thousands of Unicode CJK code points describing characters which virtually never appear in modern written material.

00,001 的
00,002 人
00,003 一
00,004 中
00,005 上
00,006 要
00,007 大
00,008 在
00,009 出
00,010 以
00,011 自
00,012 他
00,013 年
00,014 可
00,015 多
00,016 家
00,017 能
00,018 生
00,019 好
00,020 本
00,021 得
00,022 日
00,023 前
00,024 子
00,025 用

00,026 方
[...]
01,008 漂
01,009 爸
01,010 实
01,011 陳
01,012 电
01,013 融
01,014 飲
01,015 架
01,016 歲
01,017 預
01,018 籍
01,019 陣
01,020 关
01,021 憲
01,022 県
01,023 貼
01,024 狗
01,025 統
01,026 妙
01,027 漢
01,028 減
01,029 查
01,030 藏
01,031 搞
01,032 变
01,033 见
01,034 緊
01,035 棒
01,036 癸
01,037 樂
[...]
15,664 泽
15,665 樞
15,666 掘
15,667 熬
15,668 扎
15,669 艱
15,670 秒
15,671 炸
15,672 齣
15,673 糈
15,674 稈
15,675 塌
15,676 蝶
15,677 气

⁶ <http://corpus.leeds.ac.uk/list.html> ⁷ <http://technology.chtsai.org/charfreq/> ⁸ <https://github.com/mumme/smart-ranking> ⁹ <http://www.jlpt.jp/e/> ¹⁰ <http://www.chinesetest.cn/index.do>