# Annex S
## (informative)
## Procedure for the unification and arrangement of CJK ideographs

The graphic character collections of CJK Unified ideographs in this document are specified in Clause 34. They are derived from many more ideographs which are found in various different national and regional standards for coded character sets (the "sources").

Annex S describes how the ideographs in this standard are derived from the sources by applying a set of unification procedures. It also describes how the ideographs in this standard are arranged in the sequence of consecutive code points to which they are assigned.

The source references for CJK Unified ideographs are specified in Clause 24.

Within the context of this document a unification process is applied to the ideographic characters taken from the codes in the source groups. In this process, single ideographs from two or more of the source groups are associated together, and a single code point is assigned to them in this standard. The associations are made according to a set of procedures that are described below. Ideographs that are thus associated are described here as "unified".

> NOTE – The unification process does not apply to the following collections of ideographic characters:
>
> CJK RADICALS SUPPLEMENT (2E80 - 2EFF)
>
> KANGXI RADICALS (2F00 - 2FDF)
>
> CJK COMPATIBILITY IDEOGRAPHS (F900 - FAFF with the exception of FA0E, FA0F, FA11, FA13, FA14, FA1F, FA21, FA23, FA24, FA27, FA28 and FA29)
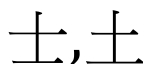>
> CJK COMPATIBILITY IDEOGRAPHS SUPPLEMENT (2F800-2FA1F).

## S.1    Unification procedure

### S.1.1    Scope of unification

Ideographs that are unrelated in historical derivation (non-cognate characters) have not been unified.

> EXAMPLE

土,土

> NOTE – The difference of shape between the two ideographs in the above example is in the length of the lower horizontal line. This is considered an actual difference of shape. Furthermore, these ideographs have different meanings. The meaning of the first is "Soldier" and of the second is "Soil or Earth".

An association between ideographs from different sources is made here if their shapes are sufficiently similar, according to the following system of classification.

### S.1.2    Two level classification

A two-level system of classification is used to differentiate (a) between abstract shapes and (b) between actual shapes determined by particular typefaces. Variant forms of an ideograph, which can not be unified, are identified based on the difference between their abstract shapes.

### S.1.3    Procedure

A unification procedure is used to determine whether two ideographs have the same abstract shape or different ones. The unification procedure has two stages, applied in the following order:

a)   Analysis of component structure;

b)   Analysis of component features;

### S.1.3.1 Analysis of component structure

In the first stage of the procedure the component structure of each ideograph is examined. A component of an ideograph is a geometrical combination of primitive elements. Alternative ideographs can be configured from the same set of components. Components can be combined to create a new component with a more complicated structure. An ideograph, therefore, can be defined as a component tree, where the top node is the ideograph itself, and the bottom nodes are the primitive elements. This is shown in Figure S.1.
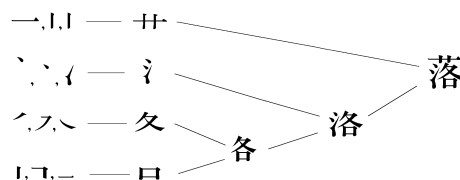


**Figure S.1 - Component structure**

### S.1.3.2 Analysis of component features

In the second stage of the procedure, the components located at corresponding nodes of two ideographs are compared, starting from the top level, as shown in Figure S.2.



**Figure S.2 - The most superior node of a component**

The following features of each ideograph to be compared are examined:

a)  the number of components,

b)  the relative position of the components in each complete ideograph,

c)  the structure of corresponding components.

If one or more of the features a) to c) above are different between the ideographs in the comparison, the ideographs are considered to have different abstract shapes and are therefore not unified.

If all of the features a) to c) above are the same between the ideographs, the ideographs are considered to have the same abstract shape and are therefore unified.

### S.1.4 Examples of differences of abstract shapes

To illustrate rules derived from a) to c) in S.1.3.2, some typical examples of ideographs that are not unified, owing to differences of abstract shapes, are shown below.

### S.1.4.1 Different number of components

The examples below illustrate rule a) since the two ideographs in each pair have different numbers of components.

崖•厓, 肱•厷, 降•夅

### S.1.4.2 Different relative positions of components

The examples below illustrate rule b). Although the two ideographs in each pair have the same number of components, the relative positions of the components are different.

峰•峯, 荊•荆

When the ideographs consists of two horizontally aligned components, a difference of the last stroke of the left-hand component going beneath the right-hand component should not warrant separate encoding, as in the case of the source glyphs for U+34F3:

㓳•㓳

### S.1.4.3 Different structure of a corresponding component

The examples below illustrate rule c). The structure of one (or more) corresponding components within the two ideographs in each pair is different.

拡•擴, 策•筞, 𭥆•燮, 圣•巠, 佥•僉, 区•區, 夹•夾,
单•單, 雈•萑, 戋•戔, 賛•贊, 襄•襄, 聿•聿, 間•閒,
朵•朵, 雋•隽, 恒•恆, 奂•奐, 从•从, 枭•枭, 夏•夏

### S.1.5 Differences of actual shapes

To illustrate the classification described in S.1.2, some typical examples of ideographs that are unified are shown below. The two or more ideographs in each group below have different actual shapes, but they are considered to have the same abstract shape and are therefore unified on their own or as component in larger ideographs. The differences are classified according to the following examples.

a) Differences in rotated strokes/dots

半•半, 勺•勺, 羽•羽羽, 酋•酋, 兼•兼, 益•益, 每•每

b) Differences in overshoot at the stroke initiation and/or termination

身•身, 雪•雪, 拐•拐, 不•不, 非•非, 周•周, 告•告

c) Differences in contact of strokes

奥•奧, 酉•酉, 児•児, 查•查, 奔•奔

d) Differences in protrusion at the folded corner of strokes

巨•巨, 成•成

e) Differences in bent strokes

西•西

f) Differences in folding back at the stroke termination

朱•朱

g)   Differences in accent at the stroke initiation

父•父, 丈•丈, 乏•乏

h)   Differences in "rooftop" modification

八•八, 穴•穴

i)   Addition or omission of a minor stroke

步•步, 者•者, 臭•臭, 呂•吕, 単•単

j)   Combinations of the above differences

刃•刃•刃, 直•直, 県•県

k)   Miscellaneous

辶•辶•辶, 示•示•礻, 且•旦•皀, 食•食•飠, 黄•黄, 盈•昷, 曷•曷, 包•包,
青•青, 册•冊, 爭•争, 备•孟•孟, 彔•录, 幵•并, 骨•骨, 吴•吳•呉,
眞•眞•真, 爲•為, 曾•曾•曽, 專•専, 内•内, 晉•晋, 龜•龜, 艹•艹

> NOTE – Some of the group items are unified when used as components in more complex ideographs, but are not unified themselves for other reasons, such as the source separation rule.

These differences in actual shapes of a unified ideograph are presented in the corresponding source columns for each code point entry in the code charts in Clause 34 of this document.

## S.1.6    Source separation rule

To preserve data integrity through multiple stages of code conversion (commonly known as "round-trip integrity"), any ideographs that are separately encoded in any one of the source standards listed below have not been unified.

| | |
|---|---|
| G-source: | GB2312-80, GB12345-90, GB7589-87*, GB7590-87*, GB8565-88*, General Purpose Hanzi List for Modern Chinese Language* |
| T-source: | TCA-CNS 11643-1986/1st plane, TCA-CNS 11643-1986/2nd plane, TCA-CNS 11643-1986/14th plane* |
| J-source: | JIS X 0208-1990, JIS X 0212-1990 |
| K-source: | KS X 1001:2004 (previously KS C 5601-1989), KS X 1002:2001 (previously KS C 5657-1991) |

> NOTE 1 – The characters from the J source: JIS X 0212-1990 encoded in this document are listed in the collection 372 JAPANESE IDEOGRAPHICS SUPPLEMENT.

> NOTE 2 – A " * " after the reference number of a standard indicates that some of the ideographs included in that standard are not introduced into the unified collection.

However, some ideographs encoded in two standards belonging to the same source group (e.g. GB2312-80 and GB12345-90) have been unified during the process of collecting ideographs from the source group.

The source separation rule described in S.1.6  only applies to the CJK UNIFIED IDEOGRAPHS block specified in the Basic Multilingual Plane.

> NOTE 3 – CJK Compatibility ideographs are created following a rule very similar to the source separation rule. However, the end result is the combination of a single CJK Unified ideograph and one or several CJK Compatibility ideographs. When the source separation rule is applied, all 'similar' source CJK ideographs result in separate CJK Unified ideographs.

## S.2    Arrangement procedure

### S.2.1    Scope of arrangement

The arrangement of the CJK Unified ideographs in the code charts of Clause 34 of this document is based on the filing order of ideographs in the following dictionaries.

| Priority | Dictionary | Edition |
|---|---|---|
| 1 | Kangxi Dictionary 康熙字典 | Beijing 7<sup>th</sup> edition |
| 2 | Daikanwa Jiten 大漢和辞典 | 9<sup>th</sup> edition |
| 3 | Hanyu Dazidian 漢語大字典 | 1<sup>st</sup> edition |
| 4 | Daejaweon 大字源 | 1<sup>st</sup> edition |

The dictionaries are used according to the priority order given in the table above. Priority 1 is highest. If an ideograph is found in one dictionary, the dictionaries of lower priority are not examined.

### S.2.2    Procedure

#### S.2.2.1    Ideographs found in the dictionaries

a)   If an ideograph is found in the Kangxi Dictionary, it is positioned in the code chart in accordance with the Kangxi Dictionary order.

b)   If an ideograph is not found in the Kangxi Dictionary but is found in the Daikanwa Jiten, it is given a position at the end of the radical-stroke group under which is indexed the nearest preceding Daikanwa Jiten character that also appears in the Kangxi dictionary.

c)   If an ideograph is found in neither the Kangxi nor the Daikanwa, the Hanyu Dazidian and the Daejaweon dictionaries are referred to with a similar procedure.

#### S.2.2.2    Ideographs not found in the dictionaries

If an ideograph is not found in any of the four dictionaries, it is given a position at the end of the radical-stroke group (after the characters that are present in the dictionaries) and it is indexed under the same radical-stroke count.

## S.3    Source separation examples

The pairs (or triplets) of ideographs shown below are exceptions to the unification rules described in S.1 . They are not unified because of the source separation rule described in S.1.6.

> NOTE 1 – The particular source group (or groups) that causes the source separation rule to apply is indicated by the letter (G, J, K, or T) that appears to the right of each pair (or triplet) of ideographs. The source groups that correspond to these letters are identified in S.1.6.

> NOTE 2 – Seven pairs described below have J sources which were originally part of the JIS X 212-1990 level-3, therefore covered by the Source separation rule, which are also part of the JIS X 213:2004 with source identified as J13. The pairs (J13 character code points emphasized) are: **U+5861**-U+586B, U+6483-**U+64CA**, U+75E9-**U+7626**, **U+83D1**-U+8458, U+848B-**U+8523**, U+91A4-**U+91AC**, and **U+985A**-U+985B.

| 丟 丢 | T | 刏 刐 | J | 俁 俣 | TJK | 値 值 | T |
|---|---|---|---|---|---|---|---|
| 4E1F 4E22 | | 4EDE 4EED | | 4FC1 4FE3 | | 5024 503C | |

| 么 幺 | GT | 併 併 | T | 俞 兪 | T | 偷 偸 | T |
|---|---|---|---|---|---|---|---|
| 4E48 5E7A | | 4F75 5002 | | 4FDE 516A | | 5077 5078 | |

| 争 爭 | GTJ | 侣 侶 | T | 俱 俱 | T | 偽 僞 | TJ |
|---|---|---|---|---|---|---|---|
| 4E89 722D | | 4FA3 4FB6 | | 4FF1 5036 | | 507D 50DE | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 兑 兑 | T | 剝 剥 | T | 唧 唧 | T | 壯 壯 | GTJ |
| 514C 5151 | | 525D 5265 | | 5527 559E | | 58EE 58EF | |
| 兎 兔 | TJ | 劒 劔 | J | 喩 喻 | T | 壽 壽 | T |
| 514E 5154 | | 5292 5294 | | 55A9 55BB | | 58FD 5900 | |
| 兖 兗 | T | 匀 匀 | T | 噓 噓 | T | 夐 夐 | T |
| 5156 5157 | | 52FB 5300 | | 5618 5653 | | 5910 657B | |
| 冊 册 | TJ | 单 単 | T | 噏 噏 | GTJ | 夲 本 | GTJ |
| 518A 518C | | 5355 5358 | | 568F 5694 | | 5932 672C | |
| 净 淨 | G | 即 卽 | TK | 囯 国 | T | 奧 奧 | J |
| 51C0 51C8 | | 5373 537D | | 56EF 56FD | | 5965 5967 | |
| 几 几 | T | 卷 巻 | TJ | 圈 圏 | TJ | 奨 奬 奨 | TJ |
| 51E2 51E3 | | 5377 5DFB | | 5708 570F | | 5968 596C 734E | |
| 刃 刄 | TJ | 叁 参 | GT | 圎 圓 | T | 妆 妝 | GT |
| 5203 5204 | | 53C1 53C2 | | 570E 5713 | | 5986 599D | |
| 刊 刋 | TJ | 參 叄 | T | 圖 圗 | T | 妍 姸 | T |
| 520A 520B | | 53C3 53C4 | | 5716 5717 | | 598D 59F8 | |
| 删 删 | T | 吕 呂 | T | 坙 坙 | T | 姍 姗 | T |
| 5220 522A | | 5415 5442 | | 5759 5DE0 | | 59CD 59D7 | |
| 別 别 | T | 吞 呑 | T | 埒 埓 | J | 姬 姬 | GT |
| 5225 522B | | 541E 5451 | | 57D2 57D3 | | 59EB 59EC | |
| 券 劵 | TJ | 吳 吴 吴 | TJ | 墈 墍 | T | 娛 娯 娱 | T |
| 5238 52B5 | | 5433 5434 5449 | | 5848 588D | | 5A1B 5A2F 5A31 | |
| 刹 剎 | T | 呐 吶 | T | 填 填 | TJ | 婕 婫 | T |
| 5239 524E | | 5436 5450 | | 5861 586B | | 5A55 5AAB | |
| 剏 剙 | T | 告 吿 | T | 增 增 | T | 媾 媮 | T |
| 524F 5259 | | 543F 544A | | 5897 589E | | 5A7E 5AAE | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 媪 媼 | TK | 尪 尫 | T | 彔 录 | T | 戩 戬 | GT |
| 5AAA 5ABC | | 5C2A 5C2B | | 5F54 5F55 | | 6229 622C | |
| 媯 嬀 | T | 屶 屷 | T | 彙 彚 | T | 戲 戱 | T |
| 5AAF 5B00 | | 5C36 5C37 | | 5F59 5F5A | | 622F 6231 | |
| 媎 嬈 | T | 屏 屛 | T | 彛 彜 | J | 戶 户 戸 | T |
| 5B0E 5B14 | | 5C4F 5C5B | | 5F5B 5F5C | | 6236 6237 6238 | |
| 孄 孷 | GT | 峥 崢 | GT | 彝 彝 | T | 戻 戾 | T |
| 5B24 5B37 | | 5CE5 5D22 | | 5F5D 5F5E | | 623B 623E | |
| 孳 孶 | T | 巓 巔 | T | 彦 彥 | T | 抛 拋 | T |
| 5B73 5B76 | | 5DD3 5DD4 | | 5F65 5F66 | | 629B 62CB | |
| 宫 宮 | T | 岼 幷 | T | 德 徳 | T | 拔 拔 | TJ |
| 5BAB 5BAE | | 5E21 5E32 | | 5FB3 5FB7 | | 629C 62D4 | |
| 寛 寬 | T | 带 帶 | TJ | 徴 徵 | T | 挩 挽 | T |
| 5BDB 5BEC | | 5E2F 5E36 | | 5FB4 5FB5 | | 6329 635D | |
| 寜 寧 | T | 并 幵 | T | 惠 惠 | TJ | 揷 插 插 | TJ |
| 5BDC 5BE7 | | 5E76 5E77 | | 6075 60E0 | | 633F 63D2 63F7 | |
| 寝 寢 | GTJ | 廄 廏 | T | 悅 悦 | T | 捏 捏 | TJ |
| 5BDD 5BE2 | | 5EC4 5ECF | | 6085 60A6 | | 634F 63D1 | |
| 専 專 | J | 弑 弒 | T | 惧 悞 | T | 搜 搜 | TJ |
| 5C02 5C08 | | 5F11 5F12 | | 609E 60AE | | 635C 641C | |
| 将 將 | GTJ | 强 強 | T | 悳 惪 | T | 揭 楬 | T |
| 5C06 5C07 | | 5F37 5F3A | | 60B3 60EA | | 63B2 63ED | |
| 尔 尒 | T | 弹 彈 | T | 慍 愠 | T | 摇 搖 摇 | TJ |
| 5C13 5C14 | | 5F39 5F3E | | 6120 614D | | 63FA 6416 6447 | |
| 尙 尚 | T | 彐 彑 | TJ | 慎 愼 | TJ | 搵 搵 | T |
| 5C19 5C1A | | 5F50 5F51 | | 613C 614E | | 63FE 6435 | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 擊擊 | TJ | 概概 | T | 污污 | T | 潛潛 | GTJK |
| 6483 64CA | | 6982 69EA | | 6C5A 6C61 | | 6F5B 6FF3 | |
| 敎教 | T | 楹楹 | T | 沒没 | TJ | 瀨瀨 | T |
| 654E 6559 | | 6985 69B2 | | 6C92 6CA1 | | 7028 702C | |
| 敓敓 | T | 橵橵 | T | 淨淨 | TJ | 為爲 | GTJ |
| 6553 655A | | 699D 6A27 | | 6D44 6DE8 | | 70BA 7232 | |
| 既旣 | T | 槇槙 | J | 涉涉 | T | 熒熒 | GTJK |
| 65E2 65E3 | | 69C7 69D9 | | 6D89 6E09 | | 712D 7162 | |
| 昂昻 | T | 樣樣 | TJ | 涗涚 | T | 熙熙 | J |
| 6602 663B | | 69D8 6A23 | | 6D97 6D9A | | 7155 7199 | |
| 晚晚 | T | 橫橫 | T | 淚淚 | T | 熅熅 | T |
| 665A 6669 | | 6A2A 6A6B | | 6D99 6DDA | | 7174 7185 | |
| 曁曁 | T | 步步 | T | 淥淥 | T | 狀狀 | GT |
| 66A8 66C1 | | 6B65 6B69 | | 6DE5 6E0C | | 72B6 72C0 | |
| 曾曾 | J | 歲歲 | T | 清清 | T | 瑤瑶 | TJ |
| 66FD 66FE | | 6B72 6B73 | | 6DF8 6E05 | | 7464 7476 | |
| 枴枴 | T | 歿歿 | T | 渴渴 | T | 瓶瓶 | T |
| 67B4 67FA | | 6B7F 6B81 | | 6E07 6E34 | | 74F6 7501 | |
| 查查 | T | 殼殼 | GTJ | 溫溫 | T | 產産 | T |
| 67E5 67FB | | 6BBB 6BBC | | 6E29 6EAB | | 7522 7523 | |
| 柵栅 | T | 毀毀 | T | 潙潙 | T | 瘦瘦 | J |
| 67F5 6805 | | 6BC0 6BC1 | | 6E88 6F59 | | 75E9 7626 | |
| 梲梲 | T | 每每 | T | 溉溉 | T | 皋皞 | T |
| 68B2 68C1 | | 6BCE 6BCF | | 6E89 6F11 | | 76A1 76A5 | |
| 楡榆 | T | 氳氳 | T | 滾滾 | T | 真真 | TJ |
| 6961 6986 | | 6C32 6C33 | | 6EDA 6EFE | | 771E 771F | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 眾 衆 | TJK | 緣 縁 | T | 葍 蔅 | T | 諫 諫 | TJ |
| 773E 8846 | | 7DE3 7E01 | | 8480 8495 | | 8ACC 8AEB | |
| 研 硏 | T | 縕 縕 | T | 蒋 蔣 | GJ | 謠 謡 | J |
| 7814 784F | | 7DFC 7E15 | | 848B 8523 | | 8B20 8B21 | |
| 祿 禄 | TJ | 繈 繦 | T | 蒍 蔿 | T | 豜 豣 | T |
| 797F 7984 | | 7E48 7E66 | | 848D 853F | | 8C5C 8C63 | |
| 禿 禿 | T | 羮 羹 | TJ | 薀 薀 | T | 走 走 | TJ |
| 79BF 79C3 | | 7FAE 7FB9 | | 8570 8580 | | 8D70 8D71 | |
| 稅 税 | T | 翶 翺 | T | 薫 薰 | T | 輤 軧 | T |
| 7A05 7A0E | | 7FF6 7FFA | | 85AB 85B0 | | 8EFF 8F27 | |
| 穗 穂 | TJ | 胼 胼 | T | 蘊 蘊 | T | 輜 輺 | J |
| 7A42 7A57 | | 80FC 8141 | | 85F4 860A | | 8F1C 8F3A | |
| 筝 箏 | GJ | 脫 脱 | T | 虛 虚 | T | 輼 輀 | T |
| 7B5D 7B8F | | 812B 8131 | | 865A 865B | | 8F3C 8F40 | |
| 簳 簳 | T | 腽 膃 | T | 蛻 蜕 | T | 达 迖 | T |
| 7BB3 7C08 | | 817D 8183 | | 86FB 8715 | | 8FBE 8FD6 | |
| 簒 篡 | T | 舃 舄 | GT | 衛 衞 | TJK | 迸 迸 | TJ |
| 7BE1 7C12 | | 8203 8204 | | 885B 885E | | 8FF8 902C | |
| 粤 粵 | T | 舍 舎 | TJ | 袞 裦 | TK | 遙 遥 | J |
| 7CA4 7CB5 | | 820D 820E | | 886E 889E | | 9059 9065 | |
| 絕 絶 | T | 舖 舗 | J | 裝 装 | GJK | 邢 郉 | T |
| 7D55 7D76 | | 8216 8217 | | 88C5 88DD | | 90A2 90C9 | |
| 綠 緑 | T | 荘 莊 | TJ | 訮 訮 | T | 郎 郞 | T |
| 7DA0 7DD1 | | 8358 838A | | 8A2E 8A7D | | 90CE 90DE | |
| 緖 緒 | T | 蓄 蓄 | TJ | 說 説 | T | 鄉 鄕 鄊 | T |
| 7DD2 7DD6 | | 83D1 8458 | | 8AAA 8AAC | | 90F7 9109 9115 | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 醖 醞 | T | 陻 陻 | G | 餅 餅 | TJ | 鳳 鳳 | T |
| 9196 919E | | 9667 9689 | | 9905 9920 | | 9CEF 9CF3 | |
| 醤 醬 | J | 靑 青 | T | 馱 馱 | TJK | 鶇 鶫 | J |
| 91A4 91AC | | 9751 9752 | | 99B1 99C4 | | 9D87 9DAB | |
| 鈃 鈒 | T | 静 靜 | GTJ | 駢 騈 | TK | 鷆 鷏 | J |
| 9203 9292 | | 9759 975C | | 99E2 9A08 | | 9DC6 9DCF | |
| 鋭 鋭 | T | 靭 靱 | J | 骩 骫 | T | 麪 麫 | T |
| 92B3 92ED | | 976D 9771 | | 9AA9 9AAB | | 9EAA 9EAB | |
| 錄 録 | T | 頹 頽 | T | 高 髙 | T | 麼 麽 | T |
| 9304 9332 | | 9839 983D | | 9AD8 9AD9 | | 9EBC 9EBD | |
| 鍊 鍊 | TK | 顏 顔 | TJ | 髪 髮 | TJ | 黃 黄 | T |
| 932C 934A | | 984F 9854 | | 9AEA 9AEE | | 9EC3 9EC4 | |
| 鎭 鎮 | TJ | 顚 顛 | J | 鬪 鬭 | T | 黑 黒 | T |
| 93AD 93AE | | 985A 985B | | 9B2C 9B2D | | 9ED1 9ED2 | |
| 閱 閲 | T | 飮 飲 | J | 鰛 鰮 | TJ | | |
| 95B1 95B2 | | 98EE 98F2 | | 9C1B 9C2E | | | |

## S.4 Non-unification examples

In accordance with the unification procedures described in S.1, the pairs (or triplets) of ideographs shown below are not unified. The reason for non-unification is indicated by the reference which appears to the right of each pair (or triplet). For "non-cognate" see S.1.1.

NOTE – The reason for non-unification in these examples is different from the source separation rule described in S.1.6.

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 冑 冑 | non cognate | 况 況 | S.1.4.3 | 寶 寶 | S.1.4.3 | 叕 叕 | S.1.4.3 |
| 5191 80C4 | | 51B5 6CC1 | | 5BF3 5BF6 | | 6560 656A | |
| 冲 沖 | S.1.4.3 | 垛 垜 | S.1.4.3 | 廳 廰 | S.1.4.1 | 朌 肦 | non cognate |
| 51B2 6C96 | | 579B 579C | | 5EF0 5EF3 | | 670C 80A6 | |
| 决 決 | S.1.4.3 | 孼 孽 | S.1.4.2 | 懐 懷 | S.1.4.1 | 朏 胐 | non cognate |
| 51B3 6C7A | | 5B7C 5B7D | | 61D0 61F7 | | 670F 80D0 | |

| | | | |
|---|---|---|---|
| 朐朐 non cognate<br>6710 80CA | 朣朣 non cognate<br>6723 81A7 | 稻稻 S.1.4.3<br>7A32 7A3B | 聴聽聽 S.1.4.1<br>8074 807C 807D |
| 朓朓 non cognate<br>6713 8101 | 朵朶 S.1.4.3<br>6735 6736 | 翱翱 S.1.4.3<br>7FF1 7FF6 | 荆荊 S.1.4.2<br>8346 834A |
| 朘朘 non cognate<br>6718 8127 | 灔灧 S.1.4.3<br>7054 7067 | 耇耈耉 S.1.4.3<br>8007 8008 8009 | 躱躲 S.1.4.3<br>8EB1 8EB2 |