# Annex S
## (informative)

# Procedure for the unification and arrangement of CJK Ideographs

The graphic character collections of CJK unified ideographs in ISO/IEC 10646-1 are specified in clause 27. They contain almost 27,500 ideographs, and are derived from over 66,000 ideographs which are found in various different national and regional standards for coded character sets (the "source codes").

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

The source code standards are shown in clause 27 in five groups according to their origins. The groups are identified as the G-, T-, J-, K- and V-sources.

For the purposes of ISO/IEC 10646-1 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 position 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 in the Basic multilingual Plane:
> - CJK RADICALS SUPPLEMENT (2E80 - 2EFF)
> - KANGXI RADICALS (2F00 - 2FDF)
> - CJK COMPATIBILITY IDEOGRAPHS (F900 - FAFF with the exception of FA1F and FA23).

## 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:

# 士, 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.
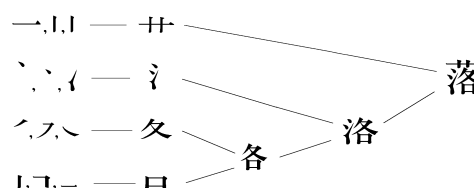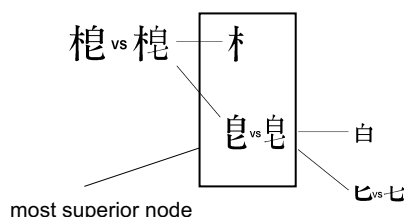


**Figure S.1  - Component structure**

## 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 most superior node, 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•夆

### 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, 荆•荊

### 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, 𭕄•燚, 圣•巠,
僉•僉, 区•區, 夹•夾, 単•單,
雚•雚, 戋•戔, 贊•贊, 襄•襄,
韭•韮, 間•閒, 朵•朵, 隽•隽,
恒•恆, 奐•奐, 𠆢人•𠆢入, 柰•柰,
叞•叞

### 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 three ideographs in each group below have different actual shapes, but they are considered to have the same abstract shape, and are therefore unified.

辶•辶•辶, 示•示•礻, 且•𠃓•皀, 食•食•𩙿,
黄•黃, 盘•盁, 曷•曷, 包•包,
靑•青, 每•每, 册•冊, 爭•争,
畬•畬•畬, 彔•录, 步•步, 者•者,
臭•臭, 幵•并, 骨•骨, 呂•吕,
直•直, 県•県, 吴吴吳, 眞•真•真,
爲•為, 単•単, 會•曽•曽, 成•成,
專•専, 內•内, 晉•晋, 龜•龜,
卅•卄

The differences are further classified according to the following examples.

a) Differences in rotated strokes/dots

半•半, 勾•勺, 羽•羽•羽, 酓•酉,
兼•兼, 益•益

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・奥，酉・酉，児・児，査・査，
奔・奔

d) Differences in protrusion at the folded corner of strokes

巨・巨

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，爻・爻

h) Differences in "rooftop" modification

八・八，穴・穴

j) Combinations of the above differences

刃・刃・刃

These differences in actual shapes of a unified ideograph are presented in the corresponding source columns for each code position entry in the code table in clause 27 of this International Standard.

### 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 C 5601-1989, KS C 5657-1991

(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.

## S.2. Arrangement procedure

### S.2.1 Scope of arrangement

The arrangement of the CJK UNIFIED IDEOGRAPHS in the code table of clause 27 of this International Standard is based on the filing order of ideographs in the following dictionaries.

| Priority | Dictionary | | Edition | |
|---|---|---|---|---|
| 1 | Kangxi Dictionary | 康熙字典 | Beijing 7th | edition |
| 2 | Daikanwa Jiten | 大漢和辞典 | 9th | edition |
| 3 | Hanyu Dazidian | 汉语大字典 | 1st | edition |
| 4 | Daejaweon | 大字源 | 1st | 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 table 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 code separation examples

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

NOTE - The particular source code group (or groups) that causes the source code 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 code groups that correspond to these letters are identified at the beginning of this Annex.

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 丟丟 | T | 兎兔 | TJ | 劒劔 | J | 喩喻 | T |
| 4E1F 4E22 | | 514E 5154 | | 5292 5294 | | 55A9 55BB | |
| 么幺 | GT | 兖充 | T | 勻匀 | T | 噓嘘 | T |
| 4E48 5E7A | | 5156 5157 | | 52FB 5300 | | 5618 5653 | |
| 争爭 | GTJ | 冊册 | TJ | 単单 | T | 嚏嚔 | GTJ |
| 4E89 722D | | 518A 518C | | 5355 5358 | | 568F 5694 | |
| 仍仭 | J | 净淨 | G | 即卽 | TK | 国囯 | T |
| 4EDE 4EED | | 51C0 51C8 | | 5373 537D | | 56EF 56FD | |
| 併併 | T | 兂兂 | T | 卷巻 | TJ | 圈圏 | TJ |
| 4F75 5002 | | 51E2 51E3 | | 5377 5DFB | | 5708 570F | |
| 侣侶 | T | 刃双 | TJ | 叁参 | GT | 圎圓 | T |
| 4FA3 4FB6 | | 5203 5204 | | 53C1 53C2 | | 570E 5713 | |
| 俁俣 | TJK | 刊刋 | TJ | 參叄 | T | 圖圗 | T |
| 4FC1 4FE3 | | 520A 520B | | 53C3 53C4 | | 5716 5717 | |
| 俞兪 | T | 删刪 | T | 吕呂 | T | 坙坙 | T |
| 4FDE 516A | | 5220 522A | | 5415 5442 | | 5759 5DE0 | |
| 俱俱 | T | 别別 | T | 吞呑 | T | 垒垳 | J |
| 4FF1 5036 | | 5225 522B | | 541E 5451 | | 57D2 57D3 | |
| 值値 | T | 券劵 | TJ | 吴吳吴 | TJ | 墅堅 | T |
| 5024 503C | | 5238 52B5 | | 5433 5434 5449 | | 5848 588D | |
| 偸偷 | T | 刹剎 | T | 呐呐 | T | 填填 | TJ |
| 5077 5078 | | 5239 524E | | 5436 5450 | | 5861 586B | |
| 偽僞 | TJ | 刱剏 | T | 告吿 | T | 增增 | T |
| 507D 50DE | | 524F 5259 | | 543F 544A | | 5897 589E | |
| 兊兑 | T | 剝剥 | T | 喞唧 | T | 壮壯 | GTJ |
| 514C 5151 | | 525D 5265 | | 5527 559E | | 58EE 58EF | |

| | | | |
|---|---|---|---|
| 壽 壽　T<br>58FD 5900 | 孳 孳　T<br>5B73 5B76 | 并 幷　T<br>5E76 5E77 | 憊 憊　T<br>60B3 60EA |
| 复 复　T<br>5910 657B | 宮 宮　T<br>5BAB 5BAE | 廄 廄　T<br>5EC4 5ECF | 慍 慍　T<br>6120 614D |
| 本 本　GTJ<br>5932 672C | 寬 寬　T<br>5BDB 5BEC | 弑 弑　T<br>5F11 5F12 | 慎 慎　TJ<br>613C 614E |
| 奧 奧　J<br>5965 5967 | 寧 寧　T<br>5BDC 5BE7 | 強 強　T<br>5F37 5F3A | 戩 戩　GT<br>6229 622C |
| 奬 奬 奬　TJ<br>5968 596C 734E | 寢 寢　GTJ<br>5BDD 5BE2 | 弹 弹　T<br>5F39 5F3E | 戲 戲　T<br>622F 6231 |
| 妆 妝　GT<br>5986 599D | 專 專　J<br>5C02 5C08 | 彐 彑　TJ<br>5F50 5F51 | 戶 户 戸　T<br>6236 6237 6238 |
| 妍 妍　T<br>598D 59F8 | 将 將　GTJ<br>5C06 5C07 | 彔 录　T<br>5F54 5F55 | 戾 戾　T<br>623B 623E |
| 姍 姍　T<br>59CD 59D7 | 尒 尔　T<br>5C13 5C14 | 彙 彙　T<br>5F59 5F5A | 抛 抛　T<br>629B 62CB |
| 姬 姬　GT<br>59EB 59EC | 尚 尚　T<br>5C19 5C1A | 彛 彜　J<br>5F5B 5F5C | 拔 拔　TJ<br>629C 62D4 |
| 娛 娛 娛　T<br>5A1B 5A2F 5A31 | 尫 尪　T<br>5C2A 5C2B | 彝 彝　T<br>5F5D 5F5E | 挩 挩　T<br>6329 635D |
| 婕 婙　T<br>5A55 5AAB | 尶 尷　T<br>5C36 5C37 | 彥 彦　T<br>5F65 5F66 | 插 插 插　TJ<br>633F 63D2 63F7 |
| 婾 婾　T<br>5A7E 5AAE | 屏 屏　T<br>5C4F 5C5B | 德 德　T<br>5FB3 5FB7 | 捏 捏　TJ<br>634F 63D1 |
| 媼 媼　TK<br>5AAA 5ABC | 峥 峥　GT<br>5CE5 5D22 | 徵 徵　T<br>5FB4 5FB5 | 搜 搜　TJ<br>635C 641C |
| 媯 嫣　T<br>5AAF 5B00 | 巓 巓　T<br>5DD3 5DD4 | 惠 惠　TJ<br>6075 60E0 | 揭 揭　T<br>63B2 63ED |
| 孍 孍　T<br>5B0E 5B14 | 帋 帋　T<br>5E21 5E32 | 悅 悅　T<br>6085 60A6 | 搖 搖 搖　TJ<br>63FA 6416 6447 |
| 孏 孏　GT<br>5B24 5B37 | 帶 帶　TJ<br>5E2F 5E36 | 惝 惝　T<br>609E 60AE | 搵 搵　T<br>63FE 6435 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 擊擊 | TJ | 槇槙 | J | 淥淥 | T | 產産 | T |
| 6483 64CA | | 69C7 69D9 | | 6DE5 6E0C | | 7522 7523 | |
| 敎教 | T | 樣樣 | TJ | 淸清 | T | 瘦瘦 | J |
| 654E 6559 | | 69D8 6A23 | | 6DF8 6E05 | | 75E9 762 | |
| 敓敓 | T | 橫橫 | T | 渴渴 | T | 皡皥 | T |
| 6553 655A | | 6A2A 6A6B | | 6E07 6E34 | | 76A1 76A5 | |
| 既旣 | T | 步歩 | T | 溫温 | T | 眞真 | TJ |
| 65E2 65E3 | | 6B65 6B69 | | 6E29 6EAB | | 771E 771F | |
| 昂昻 | T | 歲歳 | T | 溈潙 | T | 眾衆 | TJK |
| 6602 663B | | 6B72 6B73 | | 6E88 6F59 | | 773E 8846 | |
| 晚晚 | T | 歿殁 | T | 溉漑 | T | 硏研 | T |
| 665A 6669 | | 6B7F 6B81 | | 6E89 6F11 | | 7814 784F | |
| 暨曁 | T | 殼殻 | GTJ | 滾滾 | T | 祿禄 | TJ |
| 66A8 66C1 | | 6BBB 6BBC | | 6EDA 6EFE | | 797F 7984 | |
| 曾曽 | J | 毀毁 | T | 潛潜 | GTJK | 禿禿 | T |
| 66FD 66FE | | 6BC0 6BC1 | | 6F5B 6FF3 | | 79BF 79C3 | |
| 枴枴 | T | 每每 | T | 瀨瀬 | T | 稅税 | T |
| 67B4 67FA | | 6BCE 6BCF | | 7028 702C | | 7A05 7A0E | |
| 查査 | T | 氲氳 | T | 為爲 | GTJ | 穗穗 | TJ |
| 67E5 67FB | | 6C32 6C33 | | 70BA 7232 | | 7A42 7A57 | |
| 柵柵 | T | 污汚 | T | 熒熒 | GTJK | 箏箏 | GJ |
| 67F5 6805 | | 6C5A 6C61 | | 712D 7162 | | 7B5D 7B8F | |
| 梲梲 | T | 沒没 | TJ | 熙熙 | J | 箅箅 | T |
| 68B2 68C1 | | 6C92 6CA1 | | 7155 7199 | | 7BB3 7C08 | |
| 楡榆 | T | 淨淨 | TJ | 煜熅 | T | 簒篡 | T |
| 6961 6986 | | 6D44 6DE8 | | 7174 7185 | | 7BE1 7C12 | |
| 槪概 | T | 涉涉 | T | 狀狀 | GT | 粵粤 | T |
| 6982 69EA | | 6D89 6E09 | | 72B6 72C0 | | 7CA4 7CB5 | |
| 榅榲 | T | 涗涚 | T | 瑤瑶 | TJ | 絕絶 | T |
| 6985 69B2 | | 6D97 6D9A | | 7464 7476 | | 7D55 7D76 | |
| 橻橻 | T | 淚淚 | T | 瓶瓶 | T | 綠綠 | T |
| 699D 6A27 | | 6D99 6DDA | | 74F6 7501 | | 7DA0 7DD1 | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 緒緒 | T | 蔿蔿 | T | 輜輜 | J | 陧陧 | G |
| 7DD2 7DD6 | | 848D 853F | | 8F1C 8F3A | | 9667 9689 | |
| 緣緣 | T | 薀薀 | T | 輼輼 | T | 靑青 | T |
| 7DE3 7E01 | | 8570 8580 | | 8F3C 8F40 | | 9751 9752 | |
| 縕縕 | T | 薫薫 | T | 达达 | T | 静靜 | GTJ |
| 7DFC 7E15 | | 85AB 85B0 | | 8FBE 8FD6 | | 9759 975C | |
| 繈繈 | T | 蘊蘊 | T | 迸迸 | TJ | 靭靱 | J |
| 7E48 7E66 | | 85F4 860A | | 8FF8 902C | | 976D 9771 | |
| 羮羹 | TJ | 虚虛 | T | 遙遥 | J | 頹頹 | T |
| 7FAE 7FB9 | | 865A 865B | | 9059 9065 | | 9839 983D | |
| 翶翺 | T | 蛻蛻 | T | 邢邢 | T | 顔顔 | TJ |
| 7FF6 7FFA | | 86FB 8715 | | 90A2 90C9 | | 984F 9854 | |
| 胼胼 | T | 衛衞 | TJK | 郎郎 | T | 顚顛 | J |
| 80FC 8141 | | 885B 885E | | 90CE 90DE | | 985A 985B | |
| 脫脱 | T | 袠袞 | TK | 鄉鄕鄊 | T | 飮飲 | J |
| 812B 8131 | | 886E 889E | | 90F7 9109 9115 | | 98EE 98F2 | |
| 膃膃 | T | 裝裝 | GJK | 醖醞 | T | 餅餠 | TJ |
| 817D 8183 | | 88C5 88DD | | 9196 919E | | 9905 9920 | |
| 鳥鳥 | GT | 訮訮 | T | 醬醬 | J | 馱馱 | TJK |
| 8203 8204 | | 8A2E 8A7D | | 91A4 91AC | | 99B1 99C4 | |
| 舍舎 | TJ | 說説 | T | 鈃鈃 | T | 駢駢 | TK |
| 820D 820E | | 8AAA 8AAC | | 9203 9292 | | 99E2 9A08 | |
| 舖舗 | J | 諫諫 | TJ | 銳鋭 | T | 觖觖 | T |
| 8216 8217 | | 8ACC 8AEB | | 92B3 92ED | | 9AA9 9AAB | |
| 莊荘 | TJ | 謠謡 | J | 錄録 | T | 高髙 | T |
| 8358 838A | | 8B20 8B21 | | 9304 9332 | | 9AD8 9AD9 | |
| 蓿蓿 | TJ | 豼豼 | T | 鍊錬 | TK | 髮髮 | TJ |
| 83D1 8458 | | 8C5C 8C63 | | 932C 934A | | 9AEA 9AEE | |
| 萓薝 | T | 走赱 | TJ | 鎭鎮 | TJ | 鬬鬭 | T |
| 8480 8495 | | 8D70 8D71 | | 93AD 93AE | | 9B2C 9B2D | |
| 蔣蔣 | GJ | 軒軒 | T | 閲閲 | T | 鰛鰮 | TJ |
| 848B 8523 | | 8EFF 8F27 | | 95B1 95B2 | | 9C1B 9C2E | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 鳳鳳 | T | 鷶鶊 | J | 麼麼 | T | 黑黒 | T |
| 9CEF 9CF3 | | 9DC6 9DCF | | 9EBC 9EBD | | 9ED1 9ED2 | |
| 鶇鶇 | J | 麨麪 | T | 黃黄 | T | | |
| 9D87 9DAB | | 9EAA 9EAB | | 9EC3 9EC4 | | | |

In accordance with the unification procedures described in S.1 of this Annex 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 code separation rule described in clause S.1.6.

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 冑冑 | non cognate | 寶寶 | S.1.4.3 | 胸胷 | non cognate | 稻稻 | S.1.4.3 |
| 5191 80C4 | | 5BF3 5BF6 | | 6710 80CA | | 7A32 7A3B | |
| 冲沖 | S.1.4.3 | 廳廰 | S.1.4.1 | 朓朓 | non cognate | 翱翶 | S.1.4.3 |
| 51B2 6C96 | | 5EF0 5EF3 | | 6713 8101 | | 7FF1 7FF6 | |
| 决決 | S.1.4.3 | 懷懐 | S.1.4.1 | 腠腠 | non cognate | 耇耈耉 | S.1.4.3 |
| 51B3 6C7A | | 61D0 61F7 | | 6718 8127 | | 8007 8008 8009 | |
| 況况 | S.1.4.3 | 敠敪 | S.1.4.3 | 朣瞳 | non cognate | 聴聼聽 | S.1.4.1 |
| 51B5 6CC1 | | 6560 656A | | 6723 81A7 | | 8074 807C 807D | |
| 垛垜 | S.1.4.3 | 朌肦 | non cognate | 朵朶 | S.1.4.3 | 荊荊 | S.1.4.2 |
| 579B 579C | | 670C 80A6 | | 6735 6736 | | 8346 834A | |
| 孼孽 | S.1.4.2 | 朏胐 | non cognate | 灔灧 | S.1.4.3 | 躱躲 | S.1.4.3 |
| 5B7C 5B7D | | 670F 80D0 | | 7054 7067 | | 8EB1 8EB2 | |