

1. Introduction

This document is a standing document of ISO/IEC JTC 1/SC 2/WG 2/IRG ([Ideographic Rapporteur Group](#)). It consists of a set of principles and procedures on a number of items relevant to the preparation, submission and development of a repertoire of [Chinese-Japanese-Korean \(CJK\) Unified Ideographs](#) extensions for additions to the standard ([ISO/IEC 10646](#)). The document also contains procedures and guidelines [governing the work of the IRG](#). Submitters should check the standard documents (including all the amendments and corrigenda) before preparing new submissions. Submitters are encouraged to visit the "[Unihan Database](#)" page on the Unicode web site for more information on checking whether a CJK [ideograph](#) is already encoded in the standard [or not](#).

For anything not explicitly [covered](#) in this document, [the IRG](#) will follow the Principles and Procedures of WG2 and [other](#) higher level directives.

2. IRG's Scope of Work

The IRG [deals with](#) CJK ideograph-related tasks under the supervision of WG2 (SC2 Resolution M13-05). The following [is a](#) list of IRG projects:

- CJK Unified Ideograph [Repertoire](#) and its extensions
- Kangxi Radicals and CJK Radical Supplements
- Ideographic Description Characters
- IICORE ([International Ideographs Core](#))
- CJK Strokes
- Old Hanzi

The following sections are dedicated [to](#) CJK Unified Ideographs and [a](#) set of principles and procedures to be applied in the development of a new repertoire of CJK Unified Ideographs. Standardizing CJK Compatibility Characters maintained in UCS for the purpose of round-trip integrity with other standards is out of IRG's [work](#) scope.

3. Development of CJK Unified Ideographs

(TBD) When and under which conditions will a new extension of CJK Unified Ideographs be developed as an IRG project?

3.1 General principles - [characters not glyphs](#)

A. UCS encodes abstract characters. A member of CJK Unified Ideographs is such an abstract character that should be determined by its own abstract shape:

A CJK [ideographic](#) character can be written in many actual forms depending on [the](#) writing style [adopted](#). [Common writing styles include](#) Song or Ming style (typical print form), Kai style (hand written form) [and](#) Cao style (cursive form). Stylistically different forms of the same character can be represented by different number of different type of strokes and/or components, which could affect identification of the same abstract shape. In order to [reach](#) a common ground to identify those abstract shapes to be encoded as distinct CJK Unified Ideographs, [the IRG only](#) accepts submissions [using](#) print form of glyphs (usually Song or Ming style).

B. Unification procedures of CJK ideographs:

Standard print forms of CJK ideographs are constructed with a combination of known components and/or stroke types. Most of them are determined by two components - a radical chosen to classify the character in [dictionaries](#) and possibly reflect the meaning of the character and a phonetic component which represents the [pronunciation](#) of the character [\[to be revisited\]](#). Basically, two submitted print forms of glyphs with the same phonetic component are distinct characters if they have different radicals. For non trivial cases, further shape analysis will be conducted. Two similar glyphs shall be decomposed into radicals, components and/or stroke types and evaluated by following the unification procedures described in Annex S of ISO/IEC 10646.

C. Non-cognate rule:

Deleted: for ###

Deleted: Ideograph

Deleted: written

Deleted: works for

Deleted: are the

Deleted: .

Deleted: s

Deleted: .

Deleted: for

Deleted: the

Deleted: Characters

Deleted: Not

Deleted: G

Deleted: Ideographic

Deleted: a

Deleted: such as

Deleted: ,

Deleted: , etc., and those s

Deleted: s

Deleted: facilitate

Deleted: consisting only from

Deleted: dictionary

Deleted: reading

Deleted: 36

Deleted: 6

Deleted: 9

Deleted: 9

Deleted: 9

No matter how similar two ideographs is in actual shape, non-cognate or semantically different glyphs shall be considered to have different abstract shapes.

'戍'(U+620C) and '戌'(U+620D) differ only in rotated strokes/dots (S.1.5 a).

'日'(U+66F0) and '日'(U+65E5) differ only in contact of strokes (S.1.5 c). [TCA to provide a relevant example for this case]

'于'(U+4E8E) and '干'(U+5E72) differ only in folding back at the stroke termination (S.1.5 f).

Because a shape analysis might not tell non-cognateness or semantic differences, it is the submitter's responsibility to provide supporting evidence, and information in order to invoke the non-cognate rule.

D. Enhancement of Annex S with new submission:

Examples in Annex S shall be continuously updated. In reviewing character submissions, the IRG shall consider whether or not a new submission is worth including in Annex S as a new example for unification or disunification.

3.2 Preparation for submission to the IRG

A. Required data to be submitted:

● A glyph image with a specified dimension and filename in bitmap format (128 x 128 image) for each proposed ideograph in Song or Ming style.

The following data for each proposed ideograph must be submitted in the specified text format (usually in UTF-8) together with its glyph image.

- KXI (Kangxi Index with a flag to indicate real or virtual).
- KX Radical Code (Kangxi Radical with a flag to indicate simplified or traditional).
- Stroke Count (SC) of the Non-radical Component.
- First Stroke (FS) Code of the Non-radical Component.
- Ideographic Description Sequence.
- Unique ID to indicate source and the name of the glyph image for track-keeping.
- Evidence, and supporting information to support the proposed glyph shape and the usage and context with pronunciations, meanings, etc. of the proposed ideograph to convince the IRG that it is actually being used and/or non-cognate with other similar ideographs.

B. Optional data:

- For questionable characters especially those with possible unification questions, member bodies are encouraged to supply more detailed evidence of actual usage from authoritative sources and additional information on other related characters, variants and characters similar in shape or meaning already encoded in UCS for review.
- TrueType font for the glyph of the proposed ideographs (as specified under point 5 of A.1 – Submitter's Responsibilities in Annex A, WG2N3452).

C. 5 % rule:

For any character encoding standards, a common general principle is to encode the same character once and only once. It is the submitter's responsibility to filter out already encoded characters from its proposal. In assessing the suitability of a proposed ideograph for encoding, the IRG shall evaluate the credibility and quality of the submitter's proposal. If the IRG should find more than 5 % of duplicated characters in the latest UCS from the submitter's source set during the review process, the whole submission will be removed from the subsequent IRG working drafts for that particular IRG project.

3.3 Production and review process of IRG working drafts

A. Production of IRG working drafts:

After the IRG has accepted all submissions, the IRG Technical Editor will produce a set of IRG working drafts.

(TBD) The following key points should be noted:

- All working drafts should be registered with an IRG document number.

Deleted: are

Deleted: .

Deleted: the

Deleted: cognacy

Deleted: s

Deleted: to

Deleted: S

Deleted: The

Deleted: examples

Deleted: s

Deleted: the

Deleted: t

Deleted: with

Deleted: KangXi

Deleted: KangXi

Deleted: keep

Deleted: s

Deleted: readings

Deleted: for

Deleted: candidates

Deleted: use

Deleted: to

Deleted: her

Deleted: IRG

Deleted: s

Deleted: of

Deleted: t

Deleted: e

Deleted: Describe the following:

Deleted: 36

Deleted: 6

Deleted: 9

Deleted: 9

Deleted: 9

- All editors should request an IRG document number from the Rapporteur and comments should be submitted with the IRG document number assigned.
- Consolidated comments should be prepared with an IRG document number.
- Unique Character ID: once given, must not change across all versions of the same project.
- M set (i.e. Main set containing stable candidate characters without any queries raised by member bodies during the review process), D set (i.e. Discussion set containing questionable candidate characters with queries raised by member bodies during the review process) and other sets (e.g. G set (Glyph set) containing candidate characters with glyph problems) should be generated for the purpose of discussion. [(Note: add explanation about D Set, M set and other sets for reference.) Criteria for putting characters into these sets should also be stated.]
- Machine generated duplicate lists according to IDS data.
- The file name should follow the format of "IRGNnnnnXXXX" where "n" is assigned document number and "X" are alphabets for easy identification. No spaces are allowed but use of underscore "_" for separation is allowed. Use short form "Vn", e.g. V3 for version 3. Use shorter forms s as far as possible for convenience's sake.

B. Review process of IRG working drafts:

(TBD) Describe the following:

- how to split review work.
 - The Project Editor can split and assign the review work to member bodies depending on the size of the submission to be reviewed.
- what to look for.
 - duplicate characters
 - attribute errors (glyph shape/quality, KXI, Rad, SC, FS)
- how to review
 - use listed examples in IRGN954R to check attribute
 - use known patterns in Annex S
 - use the updated list of characters of unification examples in IRG standing documents.
- how to give feedback.
 - Each review cycle has its schedule. Members missing the review deadline will not have their comments considered.

3.4 Preparation for discussions at IRG meetings

A. Unification issues:

After filtering out obvious cases from machine generated duplication reports, submitters must prepare arguments with further evidence and information (such as dictionaries, publications and legal documents), supporting the use of those proposed ideographs which have been questioned as being possibly unifiable with existing UCS or other proposed ideographs in the same working drafts. The questioned ideographs with no counter arguments shall be automatically marked as unified and the IRG will move on.

For questionable characters, member bodies must supply more detailed evidence of use from authoritative sources s and additional information on other related characters, variants and characters similar in shape or meaning already encoded in UCS for review.

Further examples on the relationship with the other characters that are possibly unifiable can speed up the review and enhance quality of work.

B. Data issues:

(TBD) Describe the following: [pending Anan San to clarify the purpose of this section]

- Different choice of Rad, SC, FS etc may or may not affect KXI. When making a different choice of the radical, other attributes will be affected and should be changed accordingly.

3.5 Recording of unification arguments and decisions

The IRG should maintain all records s of unification arguments and decisions and publish them on its website. A search engine will be provided to facilitate the searching of such information for reference.

Deleted: id
Deleted: do

Deleted:

Deleted: use

Deleted: e

Formatted: Bullets and Numbering

Deleted: amount

Deleted: data

Formatted: Bullets and Numbering

Deleted: Use

Deleted: of

Deleted: Use

Deleted: of

Deleted: the

Deleted: return

Formatted: Bullets and Numbering

Deleted: s

Deleted: , e.g. dictionaries, legal documents, publications, etc. for all ...

Deleted: to be

Deleted: to

Deleted: to

Deleted: the

Deleted: , which

Deleted: In case of

Deleted: R

Deleted: may

Deleted: it at the IRG

Deleted: S

Deleted: adopted

Deleted: these

Deleted: 36

Deleted: 6

Deleted: 9

Deleted: 9

Deleted: 9

(TBD) Recording format and useful indices for easy search.

3.6 Preparation for submission to WG2

(TBD) Describe the following:

- Preparation of TrueType fonts (fonts have to be available in accordance with the requirements stated in point 5 of A.1 – Submitter’s Responsibilities in Annex A, WG2N3452).
- Source references.
- Packed multi-column format.
- The IRG should at least conduct one round of review of the table generated with TrueType font before submission.
- Member bodies are encouraged to review and comment on IRG submissions to WG2. The Rapporteur will forward member bodies’ comments to WG2.

Deleted: Multi

Deleted: IRG

Deleted: d

Deleted: r

4. Handling Defect Reports

- The IRG will follow WG2 procedures on reporting of defects according to Annex I and J of WG2 P&P document (attached to this document).

5. IRG Website

The IRG maintains its own website at <http://www.cse.cuhk.edu.hk/~irg/>, hosted by the Department of Computer Science and Engineering of the Chinese University of Hong Kong. IRG meeting notices, minutes, resolutions, document register, documents and standing documents are made available at this site. Hyperlinks to WG2 websites will be provided for member bodies’ easy access.

Deleted: w

Deleted:

Deleted:

Deleted: in

Deleted: 36

Deleted: 6

Deleted: 9

Deleted: 9

Deleted: 9

Annex A: Information accompanying submissions

[cf Sections 3.2A Required data to be submitted and 3.2B Optional data]

Annex B: IDS matching [Japan's Mr Taichi Kawabata is invited to contribute to this part]

(TBD)

B.1 Simple IDS

(TBD)

B.2 Handling of complex or incomplete IDS

(TBD)

Annex C: Work flow and stages of progression [The IRG Chief Editor and Technical Editor are invited to contribute to this part](TBD)

C.1 The IRG working drafts

(TBD)

C.2 Stages of progression

(TBD)

C.3 Dealing with urgent requests

(TBD)

- For submissions with the status of "National" or "Regional" standards, the IRG will consider giving priority in processing them with regard to the workload to be incurred.

C.3 Dealing with individual submissions to WG2

(TBD)

Guidelines to deal with individual submissions to WG2:

- small enough set.
- Urgent.
- the proposal is sound and stable after exercising due diligence.
- Member bodies' submissions to WG2 for encoding characters in the compatibility zone have to go through the same unification review of CJK ideographs by the IRG.
- the same proposal should be submitted to the IRG, with additional information if it might introduce any potential conflicts with IRG work projects.

Deleted: ¶

Deleted: to

Deleted: e

Deleted: for

Deleted: consideration of

Deleted:

Deleted: require

Deleted: ing

Deleted: 36

Deleted: 6

Deleted: 9

Deleted: 9

Deleted: 9

WG2 PnP Annex I: Guideline for handling of CJK ideograph unification and/or disunification error

(Source: [ISO/IEC JTC 1/SC 2/WG 2 N2576R](#) – 2003-10-21)

There are two kinds of errors that may be encountered related to coded CJK unified ideographs.

Case 1: *to be unified* error - Ideographs that should have been unified are assigned separate code points.

Case 2: *to be disunified* error - Ideographs that should not have been unified are unified and assigned a single code point. An example of this is the request from TCA in document [N2271](#).

When such errors are found, the following guidelines will be used by WG 2 to deal with them.

1.1 Guideline for “to be unified” errors

- A. The “to be unified” pair will be left disunified. Once a character is assigned a code position in the standard, it will not be removed from the standard.
- B. If necessary, an additional note may be added to an appropriate section in the standard.

1.2 Guideline for “to be disunified” errors

- A. The ideographs to be disunified should be disunified and should be given separate code positions as soon as possible (disunification in some sense, and character name change in some sense also). These ideographs will have two separate glyphs and two separate code positions. One of these ideographs will stay at its current encoded position. The other one will have a new glyph and a new code position.
- B. For the ideographs that are encoded in the BMP, the code charts in ISO/IEC 10646 are presented in multiple columns, with possibly differing glyph shapes in each column. The question of which glyph shall be used for the currently encoded ideograph will be resolved as follows. In the interest of synchronization between ISO/IEC 10646 and the Unicode standard, the ideograph with the glyph shape that is similar to the glyph that is published in the “[Unicode Charts](#)” will continue to be associated with its current code position. For the ideographs outside the BMP, the glyph shape in ISO/IEC 10646 and the Unicode Charts are identical and will be used with its current code position.
- C. The disunified ideograph will have a glyph that is different from the one that retains the current code position.
- D. The net result will be an addition of new ideograph character and a correction and an additional entry to the source reference table.

1.3 Discouragement of new disunification request

There is a possibility of “pure true disunification” request. This is almost like the new source code separation request. This kind of request shall not be accepted disregarding the reasoning behind. Key difference between “TO BE DISUNIFIED” and “SHALL NOT BE DISUNIFIED” is as follows.

- a. If character pair is non-cognate (meanings are different), that pair of characters is TO BE DISUNIFIED.
- b. If a character pair is cognate (means the same but different shape), that pair of characters SHALL NOT BE DISUNIFIED.

Disunification request with reason of mis-application (over-application usually) of unification rule should NOT be accepted due to the principle in resolution [M41.11](#).

Deleted: 36

Deleted: 6

Deleted: 9

Deleted: 9

Deleted: 9

WG2 PnP Annex J: Guideline for correction of CJK ideograph mapping table errors

(Source: [ISO/IEC JTC 1/SC 2/WG 2 N2577](#) – 2003-09-02)

In principle, mapping table or reference to code point of existing national/regional standard (in the source reference tables) must not be changed. But once a fatal error is found it should be corrected as early as possible, under following guidelines:

J.1 Priority of error correction procedure

- A. Consider adding new code position and source-reference mapping for the character in question rather than changing the mapping table.
- B. If change of mapping table is unavoidable, correction should be done as soon as possible.

J.2 Announcement of addition or correction of mapping table

Once any addition or correction of mapping table is made, an announcement of the change should be made immediately. Usually this will be in the form of a resolution of a WG 2 meeting, followed by subsequent process resulting in an appropriate amendment to the standard.

J.3 Collection and maintenance of mapping tables that are not owned by WG 2

There are many mapping tables, which are included in national/regional standards or developed by third parties. These are out of WG 2's scope. Any organization (such as Unicode Consortium) that collects mapping information, maintains it consistently and makes this information widely available is invited and encouraged to do so.

Deleted: 36

Deleted: 6

Deleted: 9

Deleted: 9

Deleted: 9

References

Document numbers in the first column in the following table refer to IRG working documents (ISO/IEC JTC 1/SC 2/WG 2/IRG Nxxxx), except where noted otherwise. For those documents for which a link is not given, you may try <http://www.cse.cuhk.edu.hk/~irg/>; some of the older documents are available only in paper form (contact the IRG Rapporteur of JTC1/SC 2/WG 2/IRG – Prof. Lu Qin).

| Doc. No. | Title | Source | Date |
|---------------------------|--|-----------------------------------|------------|
| WG2 N3201 | Principles and Procedures for Allocation of New Characters and Scripts and handling of Defect Reports on Character Names Annex S | WG2 | 2007-03-14 |
| N681 | | Bruce Peterson and IRG Rapporteur | 1999-11-18 |
| N881 | CJK Extension C Submission Format | IRG | 2001-12-04 |
| N953 | Minutes of the Adhoc meeting on submitted documents: N941, N942, N944, N945, N948, N949 | CJK ad hoc group | 2002-11-22 |
| N954 | Report on first stroke/stroke count by ad hoc group | CJK ad hoc group | 2002-11-22 |
| N954AR | N954 Appendix: First Stroke / Stroke Count Chart | CJK ad hoc group | 2002-11-21 |
| N955 | IRG Radical Classification | Ideograph Radical Ad Hoc | 2002-11-21 |
| N956 | Ideograph Unification | Ideograph Radical Ad Hoc | 2002-11-21 |
| N1105 | Amendments to IRG N954AR | Macao | 2005-01-03 |
| N1183 | IDS decomposition principles(Revised by IRG) | KAWABATA, Taichi | 2005-12-28 |
| N1197 | Sample evidences for CJK C1 candidates | Japan | 2006-05-22 |

Deleted: 36

Deleted: 6

Deleted: 9

Deleted: 9

Deleted: 9