

Universal Multiple-Octet Coded Character Set
International Organization for Standardization
Organisation Internationale de Normalisation
Международная организация по стандартизации

Doc Type: Working Group Document
Title: Review of Revised Proposal to Encode Jurchen (N3788)
Source: Andrew West
Status: Individual Contribution
Action: For consideration by JTC1/SC2/WG2 and UTC
Date: 2010-04-15

1. Character Names

A. Errors in the Source Dictionary

The following corrections of errors in the phonetic transcriptions in the source dictionary should be reflected in the character names:

- 012-02 X-0051 𐰇 (湯古) *tanŋgu* should be *taŋgu*
- 021-05 X-0108 𐰇 (湯古) *tanŋgu* should be *taŋgu*
- 023-01 X-0114 𐰇 (希兒) *çir* should be *çir*
- 041-05 X-0194 𐰇 (只里) *džiri* should be *džiri*
- 045-03 X-0216 𐰇 (哈称) *xatfi* should be *xatfi*
- 084-07 X-0405 𐰇 (間) *giyen* should be *gijen*
- 118-01 X-0557 𐰇 (先) *fiɛn* should be *çien*
- 138-02 X-0639 𐰇 (謙) *fien* should be *teien*
- 168-02 X-0766 𐰇 (申) *fin* should be *fin*
- 169-03 X-0771 𐰇 *futici?* should be *fufiçi*
- 178-01 X-0809 𐰇 *futici?* should be *fufiçi*
- 178-02a X-0810a 𐰇 *futici?* should be *fufiçi*
- 207-01 X-0952 𐰇 (只速) *džisu* should be *džisu*
- 223-01 X-1040 𐰇 (厥) *kiye* should be *kije*

B. Simplification of the Phonetic Transcription

The phonetic reconstructions used in the source dictionary are strictly phonetic, and so use different phonetic symbols to represent the same phoneme in different contexts.

The following phonetic symbols only occur before back and mid vowels:

- *f* before *a, o, u, ə, ɨ*
- *ʃ* before *a, o, u, ə, ɨ*
- *dʒ* before *a, o, u, ə, ɨ*
- *ʒ* before *u*

Whereas the following corresponding phonetic symbols only occur before front vowels:

- *ç* before *i, y*
- *tɕ* before *i*
- *ʤ* before *i, y*

The vowels *e* and *ə* are also in complementary distribution:

- *e* only occurs as *ei, ie* and *je*
- *ə* never occurs before/after *i* or after *j*

We believe that there is no need to preserve such non-phonemic distinctions in the character names, and suggest representing pairs of phonetic symbols that represent the same phoneme and that are in complementary distribution (i.e. *f/ç, ʃ/tɕ, dʒ/ʤ* and *e/ə*) using the same ASCII letter or letters, for example using 'E' for both *e* and *ə*.

C. Use of Accented Letters in Character Names

The character names proposed in N3788 include extended characters (e.g. *Š, Ṧ, Ž, Ž̇, È*), which are disallowed according to the character naming rules. We propose the following scheme for translating the phonetic transcription in the source dictionary to legal character names (all unlisted characters are unchanged):

- *e/ə* = E
- *ɛ* = AE
- *ɨ* = NG
- *ʤ/dʒ* = J
- *tɕ/ʃ* = C
- *ç/f* = SH
- *ʒ* = ZH
- *j* = Y (in complementary distribution with 'Y' representing the vowel [y])

We prefer "J" and "C" to "DZH" and "TSH" as Manchu romanization uses "j" and "c", so the resultant Jurchen names appear much closer to their Manchu cognates.

2. Wrongly Ordered Characters

The following characters appear to be misordered:

- X-0663b (Radical 7 / 4 strokes) should be moved to after X-0218
- X-0894a (Radical 32 / 6 strokes) should be moved to after X-0931d
- X-0454 (Radical 22 / 8 strokes) and X-0455 (Radical 21 / 8 strokes) should be swapped
- X-0993 (214-08) and X-0994 (214-09) should be swapped

3. Duplicate Characters

N3788 proposes encoding the character in entry 048-03 (X-0231a and X-0231b) twice as it is a variant form of two different characters (048-02 and 051-01). We believe that this character should not be encoded twice.

4. Radical 21

The glyph form of Radical 21 卅 does not exactly match the shape of the radical in the characters under radical 21. The radical glyph has crossing horizontal strokes, whereas the characters under this radical do not:

- 092-05 X-0452 𠂇
- 092-06 X-0453 𠂈
- 093-02 X-0455 𠂉

We believe that although the difference is minor, it would be best to change the glyph for Radical 21 to match its actual appearance so that it is more distinct from Radical 20:

Radical 20	Radical 21
卅	卅

5. Glyph Variants

The revised proposal (N3788) proposes 1,430 characters for encoding, of which:

- 523 characters have no variant forms (including nearly a hundred characters with unknown readings that may in fact be variants of other characters proposed for encoding)
- 368 characters are the primary form of a character with multiple glyph variants
- 539 characters are glyph variants

Thus, almost 40% of the proposed Jurchen characters are glyph variants, which is a considerably higher proportion than for any other script encoded or proposed for encoding in the UCS. Moreover, very many of the glyph variants proposed for encoding show very insignificant differences, often just reflecting a slightly different way of writing the same stroke in different manuscript sources, as for example:

AI	𠂇 𠂇 𠂇
AMBA	𠂈 𠂈
AN	𠂉 𠂉
AXU	𠂊 𠂊

BIRA	侑侑
BUXA	汰汰
CII	戛戛
CUEN	仔仔
DAI	米米
ESE	𠂇𠂇
FUN	𠂇𠂇
GE	秀秀
FI	宋宋
FO	玫玫
GE	屈屈
GU	𠂇𠂇
I	𠂇𠂇
INDA	𠂇𠂇
IU	𠂇𠂇
JAL	𠂇𠂇
JO	𠂇𠂇
JU	𠂇𠂇
KI	𠂇𠂇
MA	𠂇𠂇
MI	𠂇𠂇
MINGGAN	五五
MO	𠂇𠂇
MUA	呆呆

NADAN	𐰇𐰇
NAN	𐰇𐰇
O	𐰇𐰇
SHI	𐰇𐰇
SHIA	𐰇𐰇𐰇
SHII	𐰇𐰇𐰇𐰇
SHIIR	𐰇𐰇
TAIYI	𐰇𐰇
TU	𐰇𐰇
U	𐰇𐰇
UJE	𐰇𐰇
XA	𐰇𐰇
YA	𐰇𐰇

It may be appropriate to separately encode some significant character variants, but it may be more appropriate to represent insignificant glyph variants by means of variation sequences.

In several respects, the nature of the Jurchen character variants is significantly different to that of Tangut character variants, which means that although variation sequences are not appropriate in the case of Tangut, they may be the most appropriate solution for Jurchen:

Tangut	Jurchen
Relatively few variant characters (about 120 out of the 6,054 characters proposed for the main Tangut block)	Relatively many variant characters (539 out of 1,430 proposed characters)
Most variants are attested in more than one modern Tangut dictionary	Most variants are only used in Jin Qicong's dictionary
Most Tangut variants are distinguished by a different structural composition (one or more different component elements, or the same component elements arranged differently)	Most Jurchen variants are distinguished only by a different way of writing the same stroke (e.g. the angle or length of a stroke) or a different placement of a stroke or a difference in stroke count