

Date: 2007-04-22

A response to N3240 (dated 2007-04-11), the National Proposal (China) entitled "Proposal on Adding 3 Tibetan Characters and a symbol for ISO/IEC 10646 in BMP"

Submitted by: Robert R. Chilton (individual expert contributor)
Technical Director, Asian Classics Input Project
Email: acip@well.com

The three new characters proposed in N3240 are:

1. "(TIBETAN ABOVE TRANSFORMED LETTER RAGU) to 0F6B"
2. "(TIBETAN ABOVE TRANSFORMED LETTER LA) to 0F6C"
3. "(TIBETAN ABOVE LETTER SA) to 0F6D"

These correspond to the three "head" letters/glyphs, of RA, LA, and SA, defined by Tibetan orthography and grammar. I will refer to these three as RA-head, LA-head, and SA-head.

The proposal N3240 asserts three reasons for why these three glyph-characters should be added to the BMP. In what follows I examine these three reasons and explain why the stated reasons are **not** sufficient to justify adding the proposed new characters.

It should be noted that international experts on Tibetan-language computing, including members of the Chinese national delegation, already considered whether or not to encode these three "head" letter glyphs as separate characters when the Tibetan script was added to the ISO/IEC 10646 BMP. The consensus of these experts was to adopt an encoding model for Tibetan script that represents the elements of Tibetan script according to their relative spatial position rather than their lexical meaning.¹ It was therefore decided at that time (1995) that no specifically defined lexical "head" characters are needed.

The above-referenced proposal, N3240, calls into question the earlier decision reached by international experts regarding the fact that there is no need for encoding specific Tibetan characters representing RA-head, LA-head and SA-head.

Here I will explain each of the reasons given in the N3240 proposal as to why these three characters are needed, together with analysis showing why these reasons are neither wholly valid nor sufficient to justify adding these proposed Tibetan characters.

¹ In forming syllable-like units, Tibetan letters can be combined in both left-to-right sequence and top-to-bottom sequence. For this reason, the BMP for Tibetan includes two sets of the consonant characters: in nominal position the letter follows to the right of the preceding character(s); in subjoined position the letter combines with and appears below the preceding character(s).

N3240 Proposal Reason #1. *It is suggested that the proposed new characters are needed for demonstrating (i.e., in Tibetan-language textbooks) the orthographic structure of the script.²* Specifically, in teaching students how to write the script, there is a need to present, in isolation, the abbreviated "head" forms of the letters RA, LA and SA. Typically (i.e., in the "u-chen" script), the RA-head has a different shape from a nominal letter RA whereas the LA-head and SA-head are of the same shape as the corresponding nominal letters LA and SA but smaller in size.

While it might be tempting to define separate characters for this purpose, this is not strictly necessary. Similar needs for showing contextual presentation forms of various characters are encountered in other combining scripts such as Arabic. In such cases, the characters Zero-Width Joiner (U+200D) and Zero-Width Non-Joiner (U+200C) are used in combination with the existing characters of the script to produce contextual variant forms as needed.

Therefore, what is needed here are not three new Tibetan characters but rather a consensus on how to implement display of these needed "head" letter glyphs in their isolated presentation forms. This would be a simple matter of defining how Tibetan fonts should display the isolated presentation forms via a sequence of the desired letter, whether TIBETAN LETTER RA [U+0F62], TIBETAN LETTER LA [U+0F63], or TIBETAN LETTER SA [U+0F66], in combination with ZERO-WIDTH JOINER [U+200D].

Alternatively, these three "head" letter glyphs could be newly encoded explicitly and solely for such isolated presentation purposes. In this case, they should be clearly defined as standalone symbols and not as alphabetic letters. That is, they should be described as individual, spacing, non-combining symbols. This would prevent any misunderstanding or misuse that would wrongly treat these presentation-form glyphs (symbols) as alphabetic (letter) characters. Likewise, the names should clearly show that these characters are intended as standalone symbolic glyphs, i.e., as:

TIBETAN SYMBOL RA MGO
TIBETAN SYMBOL LA MGO
TIBETAN SYMBOL SA MGO

N3240 Proposal Reason #2. *It is suggested that because the three "head" letters occur with high frequency in Tibetan, separate head-letter characters are desirable to support data entry such as keyboarding.³* The apparent rationale is that Tibetan keyboarding would be improved if there are separate dedicated keys for these head-letter characters.

² This is the presumed intent of passages in the N3240 proposal that mention "description of glyph" and "describe the characteristics of Modern Tibetan character".

³ This is the presumed intent of passages in the N3240 proposal that mention "application software development" and "these 3 characters are high frequency character in the Modern Tibetan and usually used as a combining element to form Tibetan composite character".

It is easy to show that even if such a keyboard behavior is desired, there is no need to encode specific head-letter characters in the BMP in order to achieve it. Various systems for keyboard entry of Tibetan are already in use and deemed acceptable; and virtually any desired keyboard layout can be implemented without any need for defining new Tibetan characters.

Moreover, there are at least two strong arguments **against** this particular rationale:

Firstly, if the proposed new characters are encoded as combining characters, they would need to be typed *after* the "root" letter that they surmount or modify. This goes against the traditional manner in which words are spelled verbally and also against the order in which the letters are written out by hand. Traditionally the "head" letters are spelled and written *before* the "root" letters that they surmount or modify.

Secondly, from the perspective of data processing and interchange in digital environments, adding these proposed characters as Tibetan letters will create many unnecessary problems and complications. Current fonts and rendering systems for Tibetan follow the specification for the Tibetan script that has been in place since 1996. These fonts and rendering systems are now in widespread use in the Microsoft Windows, Macintosh OS X, and Linux operating environments. They operate on the assumption that vertical stacks of Tibetan letters will be encoded as a sequence of a single nominal letter followed by one or more combining subjoined letters. Adding the three proposed characters will force significant changes, and add increased complexity, for all Tibetan fonts and rendering systems since they would have to additionally recognize alternate (and lexically equivalent) sequences for Tibetan letter-stacks which could be encoded as a nominal letter, followed by zero or one "head" combining letter, followed by one or more combining subjoined letters.

Even more disconcerting, adding these proposed new characters will require additional steps of normalization between Tibetan data encoded using the current methodology (character sequences ordered according to relative spatial position) and Tibetan data encoded using the proposed new methodology (character sequences ordered according to both lexical and spatial parameters). Such normalization will require the re-ordering of characters in the data sequence. The current methodology encodes the lexical "head" letter in sequence prior to the "root" letter that it modifies; the "root" letter in turn is encoded as a combining character in subjoined position. In contrast, the proposed new methodology would encode the "root" letter first—in nominal position—followed by the lexical "head" letter which attaches as a combining character.

N3240 Proposal Reason #3. *It is suggested that use of these proposed new characters will help in sorting (collation) of Tibetan data.⁴* Sorting Tibetan-script data in culturally expected order is not simple. Fortunately, the problem of machine collation of Tibetan data has already

⁴ This is the clear intent of passages in the N3240 proposal mentioning "solve the problems with Tibetan sorting based on traditional Tibetan order" and "sort order of Tibetan character base on the root letter, so the rightly find out root letter of every Tibetan glyph, then Tibetan character has sorted is possible based on BMP".

been solved (more than once).⁵ LDML collation data for sorting Tibetan is attached as Appendix A.

This supposed rationale has little if any merit. Beyond the fact that the sorting "problem" for Tibetan-character data has already been solved, adding the proposed new characters would do little to simplify the "problem". The general difficulty in sorting Tibetan is due to the fact that the "root" letter, which is given primary weight in determining the relative sort order of Tibetan syllables, can be preceded by one or two pre-radical letters.

Listed below are the 11 combinations of letter(s) that can occur ahead of the "root" letter:

༄	<TIBETAN LETTER GA [U+0F42]>
༅	<TIBETAN LETTER DA [U+0F51]>
༅	<TIBETAN LETTER BA [U+0F56]>
༅	<TIBETAN LETTER MA [U+0F58]>
༅	<TIBETAN LETTER -A [U+0F60]>
༅	<TIBETAN LETTER RA [U+0F62]>
༅	<TIBETAN LETTER LA [U+0F63]>
༅	<TIBETAN LETTER SA [U+0F66]>
༅	<TIBETAN LETTER BA [U+0F56], TIBETAN LETTER RA [U+0F62]>
༅	<TIBETAN LETTER BA [U+0F56], TIBETAN LETTER LA [U+0F63]>
༅	<TIBETAN LETTER BA [U+0F56], TIBETAN LETTER SA [U+0F66]>

Although one might argue that adding the three proposed new characters would simplify sorting for three of the cases listed above, it would do nothing to address the other eight cases. Further, some of the most difficult syllables to parse with regard to determining the "root" letter do not involve any "head" letters. Here are two examples:

ବ୍ୟାକ <BA[U+0F56], GA[U+0F42], SA U+0F66>
ମୁଦ୍ରଣ <MA[U+0F58], NGA[U+0F44], SA[U+0F66]>

The ambiguity lies in the fact that, by the rules of Tibetan grammar, either the first or the second letter could reasonably serve as the "root" letter in these syllables.

⁵ For example, culturally expected sort order for Tibetan-character data has been implemented within Microsoft Windows Vista, Linux, and Mimer SQL. (LDML data for Tibetan collation is attached as Appendix A.) The schema described for sorting Tibetan-character data is universal and general in the sense that it works equally well whether the Tibetan-script data contains modern Tibetan, classical Tibetan, Dzongkha (Bhutanese), Ladakhi, Balti (Yige), foreign words transliterated into Tibetan script, or any combination of these. For an overview, see my presentation to the Tibetan Information Technology Panel at the 2003 meeting of the *International Association for Tibetan Studies* (IATS) entitled "Sorting Unicode Tibetan using a Multi-Weight Collation Algorithm" <http://www.columbia.edu/~ph2046/iats/it/Chilton_slides.pdf>.

Any solution to the "problem" of sorting Tibetan that proposes the encoding of new characters as a means of explicitly identifying the "root" letter would need to define at least five additional new characters, for the five prefix letters, in addition to the three head letters proposed in N3240. One can imagine other proposals, such as defining separate lexical characters for various of the letters in the Tibetan alphabet—depending on whether they represent a "prefix", "head", "root", "subjoined" or "suffix" letter—that might be set forth. However, such "solutions" are neither necessary nor desirable.

Finally, adding the three new characters proposed in N3240—or any other characters similarly defined according to lexical category—would actually hinder sorting to the extent that existing solutions for sorting Tibetan would need to be modified in order to account for the new and alternative way(s) to represent Tibetan syllables that would be introduced with such additions to the character repertoire.

Summary. The reasons given in N3240 do not support the contention that the current character repertoire for Tibetan is lacking and that three new "head" characters are required. However, the N3240 proposal reminds us that there is occasional need for displaying contextual presentation forms, in isolation, of the Tibetan letters RA, LA and SA. It remains to be decided whether this requirement can best be accomplished through (1) combining, in sequence, the nominal character for the letter together with Zero-Width Joiner, or else by way of (2) encoding three new **symbol** characters which would serve solely and explicitly for the purpose of displaying these contextual glyph-forms.

Appendix A

LDML collation for sorting Tibetan-character data

```
<collation>
<rules>
<reset>&#xF40;</reset>
    <s>&#xF88;&#xF90;</s>
    <p>&#xF51;&#xF40;</p>
    <p>&#xF56;&#xF40;</p>
    <p>&#xF62;&#xF90;</p>
    <p>&#xF63;&#xF90;</p>
    <p>&#xF66;&#xF90;</p>
    <p>&#xF56;&#xF62;&#xF90;</p>
    <p>&#xF56;&#xF66;&#xF90;</p>
<reset>&#xF41;</reset>
    <s>&#xF88;&#xF91;</s>
    <p>&#xF58;&#xF41;</p>
    <p>&#xF60;&#xF41;</p>
<reset>&#xF42;</reset>
    <p>&#xF51;&#xF42;&#xF42;</p>
    <p>&#xF51;&#xF42;&#xF44;</p>
    <p>&#xF51;&#xF42;&#xF51;</p>
    <p>&#xF51;&#xF42;&#xF56;</p>
    <p>&#xF51;&#xF42;&#xF5D;</p>
    <p>&#xF51;&#xF42;&#xF60;</p>
    <p>&#xF51;&#xF42;&#xF62;</p>
    <p>&#xF51;&#xF42;&#xF63;</p>
    <p>&#xF51;&#xF42;&#xF66;</p>
    <p>&#xF51;&#xF42;&#xF74;</p>
    <p>&#xF51;&#xF42;&#xF7A;</p>
    <p>&#xF51;&#xF42;&#xF7C;</p>
    <p>&#xF51;&#xF42;&#xFB1;</p>
    <p>&#xF51;&#xF42;&#xFB2;</p>
    <p>&#xF56;&#xF42;&#xF42;</p>
    <p>&#xF56;&#xF42;&#xF51;</p>
    <p>&#xF56;&#xF42;&#xF58;</p>
    <t>&#xF56;&#xF42;&#xF7E;</t>
    <p>&#xF56;&#xF42;&#xF5D;</p>
    <p>&#xF56;&#xF42;&#xF60;</p>
    <p>&#xF56;&#xF42;&#xF62;</p>
    <p>&#xF56;&#xF42;&#xF7A;</p>
    <p>&#xF56;&#xF42;&#xF7C;</p>
    <p>&#xF56;&#xF42;&#xFB1;</p>
    <p>&#xF56;&#xF42;&#xFB2;</p>
    <p>&#xF56;&#xF42;&#xFB3;</p>
```

<p>མགར</p>

<p>མགལ</p>

<p>མགུ</p>

<p>མགེ</p>

<p>མགོ</p>

<p>མགྱ</p>

<p>མགྲ</p>

<p>འགག</p>

<p>འགང</p>

<p>འགད</p>

<p>འགན</p>

<p>འགབ</p>

<p>འགམ</p>

<t>འགཾ</t>

<p>འགའ</p>

<p>འགར</p>

<p>འགལ</p>

<p>འགས</p>

<p>འགི</p>

<p>འགུ</p>

<p>འགེ</p>

<p>འགོ</p>

<p>འགྱ</p>

<p>འགྲ</p>

<p>རྒ</p>

<p>ལྒ</p>

<p>སྒ</p>

<p>བརྒ</p>

<p>བསྒ</p>

<reset>ང</reset>

<p>དངག</p>

<p>དངང</p>

<p>དངན</p>

<p>དངར</p>

<p>དངུ</p>

<p>དངོ</p>

<p>མངག</p>

<p>མངན</p>

<p>མངའ</p>

<p>མངར</p>

<p>མངལ</p>

<p>མངོ</p>

<p>རྔ</p>

<p>ལྔ</p>

<p>སྔ</p>

<p>བརྔ</p>

<p>བསྔ</p>

<reset>ཅ</reset>

<p>གཅ</p>

<p>བཅ</p>

<p>ལྕ</p>

<p>བལྕ</p>

<reset>ཆ</reset>

<p>མཆ</p>

<p>འཆ</p>

<reset>ཇ</reset>

<p>མཇ</p>

<p>འཇ</p>

<p>རྗ</p>

<p>ལྗ</p>

<p>བརྗ</p>

<reset>ཉ</reset>

<s>ྋྙ</s>

<p>གཉ</p>

<p>མཉ</p>

<p>རྙ</p>

<p>སྙ</p>

<p>བརྙ</p>

<p>བསྙ</p>

<reset>ཏ</reset>

<t>ཊ</t>

<p>གཏ</p>

<p>བཏ</p>

<p>རྟ</p>

<p>ལྟ</p>

<p>སྟ</p>

<p>བརྟ</p>

<p>བལྟ</p>

<p>བསྟ</p>

<reset>ཐ</reset>

<t>ཋ</t>

<p>མཐ</p>

<p>འཐ</p>

<reset>ད</reset>

<t>ཌ</t>

<p>གདག</p>

<p>གདང</p>

<p>གདན</p>

<p>གདབ</p>

<p>གདམ</p>

<t>གདཾ</t>

<p>གདའ</p>

<p>གདར</p>

<p>གདལ</p>

<p>གདས</p>

<p>གདི</p>

<p>གདུ</p>

<p>གདེ</p>

<p>གདོ</p>

<p>བདག</p>

<p>བདམ</p>

<t>བདཾ</t>

<p>བདའ</p>

<p>བདར</p>

<p>བདལ</p>

<p>བདས</p>

<p>བདུ</p>

<p>བདེ</p>

<p>བདོ</p>

<p>མདག</p>

<p>མདང</p>

<p>མདན</p>

<p>མདའ</p>

<p>མདར</p>

<p>མདུ</p>

<p>མདེ</p>

<p>མདོ</p>

<p>འདག</p>

<p>འདང</p>

<p>འདད</p>

<p>འདན</p>

<p>འདབ</p>

<p>འདམ</p>

<t>འདཾ</t>

<p>འདཝ</p>

<p>འདའ</p>

<p>འདར</p>

<p>འདལ</p>

<p>འདས</p>

<p>འདི</p>

<p>འདུ</p>

<p>འདེ</p>

<p>འདོ</p>

<p>འདྲ</p>

<p>རྡ</p>

<p>ལྡ</p>

<p>སྡ</p>

<p>བརྡ</p>

<p>བལྡ</p>

<p>བསྡ</p>

<reset>ན</reset>

<t>ཎ</t>

<p>གནག</p>

<p>གནང</p>

<p>གནད</p>

<p>གནན</p>

<p>གནམ</p>

<t>གནཾ</t>

<p>གནཝ</p>

<p>གནའ</p>

<p>གནས</p>

<p>གནུ</p>

<p>གནོ</p>

<p>མནག</p>

<p>མནང</p>

<p>མནན</p>

<p>མནབ</p>

<p>མནམ</p>

<t>མནཾ</t>

<p>མནའ</p>

<p>མནར</p>

<p>མནལ</p>

<p>མནུ</p>

<p>མནེ</p>

<p>མནོ</p>

<p>རྣ</p>

<p>སྣ</p>

<p>བརྣ</p>

<p>བསྣ</p>

<reset>པ</reset>

<s>ྉྤ</s>

<p>དཔག</p>

<p>དཔང</p>

<p>དཔད</p>

<p>དཔའ</p>

<p>དཔར</p>

<p>དཔལ</p>

<p>དཔས</p>

<p>དཔུ</p>

<p>དཔེ</p>

<p>དཔོག</p>

<p>དཔོང</p>

<p>དཔོད</p>

<p>དཔོན</p>

<p>དཔོར</p>

<p>དཔྱ</p>

<p>དཔྲ</p>

<p>ལྤ</p>

<p>སྤ</p>

<reset>ཕ</reset>

<s>ྉྥ</s>

<p>འཕ</p>

<reset>བ</reset>

<p>དབག</p>

<p>དབང</p>

<p>དབད</p>

<p>དབན</p>

<p>དབབ</p>

<p>དབའ</p>

<p>དབར</p>

<p>དབལ</p>

<p>དབས</p>

<p>དབུ</p>

<p>དབེ</p>

<p>དབོ</p>

<p>དབྱ</p>

<p>དབྲ</p>

<p>འབག</p>

<p>འབང</p>

<p>འབད</p>

<p>འབན</p>

<p>འབབ</p>

<p>འབམ</p>

<t>འབཾ</t>

<p>འབའ</p>

<p>འབར</p>

<p>འབལ</p>

<p>འབི</p>

<p>འབུ</p>

<p>འབེ</p>

<p>འབོ</p>

<p>འབྱ</p>

<p>འབྲ</p>

<p>རྦ</p>

<p>ལྦ</p>

<p>སྦ</p>

<reset>མ</reset>

<t>ཾ</t>

<t>ྂ</t>

<t>ྃ</t>

<p>དམག</p>
<p>དམང</p>
<p>དམན</p>
<p>དམཝ</p>
<p>དམའ</p>
<p>དམར</p>
<p>དམས</p>
<p>དམི</p>
<p>དམུ</p>
<p>དམེ</p>
<p>དམོད</p>
<p>དམྱ</p>
<p>རྨ</p>
<p>སྨ</p>

<reset>ཙ</reset>
 <p>གཙ</p>
 <p>བཙ</p>
 <p>རྩ</p>
 <p>སྩ</p>
 <p>བརྩ</p>
 <p>བསྩ</p>

<reset>ཚ</reset>
 <p>མཚ</p>
 <p>འཚ</p>

<reset>ཛ</reset>
 <p>མཛ</p>
 <p>འཛ</p>
 <p>རྫ</p>
 <p>བརྫ</p>

<reset>ཞ</reset>
 <p>གཞ</p>
 <p>བཞ</p>

<reset>ཟ</reset>
 <p>གཟ</p>
 <p>བཟ</p>

<reset>ཡ</reset>
 <p>གཡ</p>

<reset>ར</reset>

<t>ཪ</t>

<p>བརླ</p>

<reset>ཤ</reset>

<t>ཥ</t>

<p>གཤ</p>

<p>བཤ</p>

<reset>ས</reset>

<p>གསག</p>

<p>གསང</p>

<p>གསད</p>

<p>གསན</p>

<p>གསབ</p>

<p>གསའ</p>

<p>གསར</p>

<p>གསལ</p>

<p>གསས</p>

<p>གསི</p>

<p>གསུ</p>

<p>གསེ</p>

<p>གསོ</p>

<p>བསག</p>

<p>བསང</p>

<p>བསད</p>

<p>བསབ</p>

<p>བསམ</p>

<t>བསཾ</t>

<p>བསར</p>

<p>བསལ</p>

<p>བསི</p>

<p>བསུ</p>

<p>བསེ</p>

<p>བསོ</p>

<p>བསྭ</p>

<p>བསྲ</p>

<p>བསླ</p>

<reset>ཧ</reset>

<p>ལྷ</p>

<reset>ི</reset>

<t>ྀ</t>

<reset>ེ</reset>

<t>ཻ</t>

<reset>ོ</reset>

<t>ཽ</t>

<p>ྐ</p>

<p>ྑ</p>

<p>ྒ</p>

<p>ྔ</p>

<p>ྕ</p>

<p>ྖ</p>

<p>ྗ</p>

<p>ྙ</p>

<p>ྟ</p>

<t>ྚ</t>

<p>ྠ</p>

```
<t>&#x0F9B;</t>
<p>&#x0FA1;</p>
<t>&#x0F9C;</t>
<p>&#x0FA3;</p>
<t>&#x0F9E;</t>
<p>&#x0FA4;</p>
<p>&#x0FA5;</p>
<p>&#x0FA6;</p>
<p>&#x0FA8;</p>
<p>&#x0FA9;</p>
<p>&#x0FAA;</p>
<p>&#x0FAB;</p>
<p>&#x0FAD;</p>
<t>&#x0FBA;</t>
<p>&#x0FAE;</p>
<p>&#x0FAF;</p>
<p>&#x0FB0;</p>
<p>&#x0FB1;</p>
<t>&#x0FBB;</t>
<p>&#x0FB2;</p>
<t>&#x0FBC;</t>
<p>&#x0FB3;</p>
<p>&#x0FB4;</p>
<t>&#x0FB5;</t>
<p>&#x0FB6;</p>
<p>&#x0FB7;</p>
<p>&#x0FB8;</p>
<reset>&#x0F39;</reset>
<s>&#x0F84;</s>
<s>&#x0F71;</s>
<s>&#x0F39;</s>
<s>&#x0F7F;</s>
<s>&#x0F85;</s>
<s>&#x0F88;</s>
<s>&#x0F89;</s>
<s>&#x0F8A;</s>
<s>&#x0F8B;</s>
<reset>&#x0FB2;&#x0F71;&#x0F80;</reset>
    <i>&#x0F77;</i>
<reset>&#x0FB3;&#x0F71;&#x0F80;</reset>
    <i>&#x0F79;</i>
<reset>&#x0F51;&#x0F42;&#x0F42;&#x0F66;</reset>
    <t>&#x0F51;&#x0F42;&#x0F4A;</t>
    <t>&#x0F51;&#x0F42;&#x0F4C;</t>
<reset>&#x0F56;&#x0F42;&#x0F42;&#x0F66;</reset>
    <t>&#x0F56;&#x0F42;&#x0F4A;</t>
```



```
<reset>&#x0FBF;</reset>
<p>&#x0F00;</p>
<reset>&#x0EC6;</reset>
<p>&#x0F0B;</p>
<t>&#x0F0C;</t>
<s>&#x0F0D;</s>
<s>&#x0F0E;</s>
<s>&#x0F0F;</s>
<s>&#x0F10;</s>
<s>&#x0F11;</s>
<s>&#x0F14;</s>
</rules>
</collation>
```
