

Collation of Transliterating Tibetan Characters

Heming Huang^{1,2,*} and Feipeng Da¹

¹ School of Automation, Southeast University, Nanjing, Jiangsu 210096, China

² School of computer Science, Qinghai Normal University, Xining, Qinghai 810008, China
 huang-heming@sohu.com, daftp@seu.edu.cn

Abstract. The transliterating Tibetan characters used specially to transliterate foreign scripts have two collations: collation with the rules of native Tibetan dictionary and with that of transliterating Tibetan dictionary. This paper proposes two general structures for transliterating characters. Based on these general structures, a collation scheme is developed so that all transliterating characters can be collated correctly and effectively.

Keywords: Tibetan, character, collation, structure.

1 Introduction

The Tibetan script is an alphasyllabary, a segmental writing system in which consonant-vowel sequences are written as a unit. Tibetan has two alphabets: the native Tibetan alphabet used in daily life of Tibetan people and the transliterating Tibetan alphabet used specially to transliterate foreigner scripts especially the Sanskrit.

[illegible]

The transliterating Tibetan is different from the native Tibetan in many ways. One difference is that the transliterating Tibetan has two kinds of collation. The first kind is that all the characters need to be collated are just the transliterating characters and

* Corresponding author.

native Tibetan syllables under this circumstance. It should be judged with the native Tibetan orthography. Generally, a pre-composed character is a transliterating character if it meets one of the following conditions.

- 1) A pre-composed character has the transliterating vowel འ, སྟེ, རྩེ, ལྷེ, ཤྭེ, བྱེ, མྱེ, ཐྱེ, ཏྱེ, ཅྱེ, or ཇྱེ.
- 2) A pre-composed character has the diacritic sigh ི, ུ, ྲྀ, ོ, ེ, ཻ, ཽ, or ཿ.
- 3) A pre-composed character has the transliterating consonants རྒྱ, ལྗ, ལྡ, ལླ, ལྰ, ལྱ, ལྵ, ལྶ, ལྷ, ལྸ, ལྐྵ, ལྺ, ལྻ, ལྼ, or ལ྽.
- 4) A pre-composed character has two consonants, but the first consonant is none of འ, རྩེ, and སྟེ while the second consonant is none of འ, རྩེ, ལྷེ, and ཤྭེ. Examples of such characters are རྒྱེ, ལྗེ, ལྡེ, and ལླེ.
- 5) A pre-composed character has three consonants, but the first one is none of འ, རྩེ, and སྟེ while the third one is none of འ, རྩེ, ལྷེ, and ཤྭེ. Examples of such characters are རྒྱལྡ, ལྗལྡ, and ལྡལྡ.
- 6) A pre-composed character has more than three consonants. Examples of such characters are རྒྱལྡལྡ, ལྗལྡལྡ, and ལྡལྡལྡ.
- 7) A horizontal combination of several consonants, but there is no prefix consonant or suffix consonant according to the restriction rules of native Tibetan Standard orthography to these positions. Examples of such combinations are གཁལལ, གཁལ, གལལ, and གལལལ.
- 8) A horizontal combination of a consonant and a pre-composed character, but the consonant is neither the prefix consonant nor the suffix consonant. Examples of such combinations are གལྡ, གལྡེ, གལྡེེ, གལྡེེེ, and ངལྡ.
- 9) A horizontal combination of several pre-composed characters, but the last one is none of རྩེ, རྩྥ, and རྩྭ. Examples of such combinations are རྩྭེེ, རྩྭེེེ, and རྩྭེེེེ.

3 The General Structure of Transliterating Characters

The collation of a transliterating character is not decided by its component letters directly. A transliterating character may be decomposed into several syllables firstly and then its collation is decided by those syllable series. Therefore, it is necessary to describe the syllable of transliterating characters.

3.1 The Collation Rules of the Transliterating Tibetan Dictionary

A transliterating character may be the vertical composition of basic consonant, foot consonant, and vowel and there are no concepts of prefix consonant, suffix consonant, and superscript consonant. Therefore, the phrases གཤི, མཐུར, and བཅོན belong to the chapters ག, མ, and བ respectively; and the phrases རྟམ, ལྟམ, and སྟམ belong to the chapters ར, ལ, and ས respectively. Furthermore, a transliterating syllable may have two foot consonants and two vowels. For example, the syllable ལྟམ has two vowels. The first one is ལ.

4.1 Collated with the Rules of the Transliterating Character Dictionary

When two transliterating characters are collated with the rules of the transliterating character dictionary, the scheme of the transliterating character collation consists of the following five steps as shown in Fig. 4.

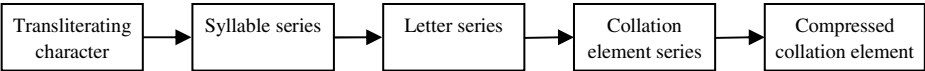


Fig. 4. The scheme of the transliterating character collation

- Step 1: Decompose each transliterating character into syllable series first.
- Step 2: Expand each syllable further into the letter series according to the sort order shown in Fig. 3. If there is no letter in a certain position, a space ‘□’ is used instead.
- Step 3: Replace each letter in the letter series with the corresponding collation element.
- Step 4: Compress the collation element series.
- Step 5: Compare the two compressed collation element series and we have got the collation result of two transliterating characters.

However, this paper just focuses on the first three steps.

To collate the two characters རྩུ་ and རྩུ་, for example, they should firstly be expanded into syllable series ‘རྩུ་’ and ‘རྩུ་’ respectively; and then each syllable is further expanded into letter series, thus we have got ‘རྩུ་ཨྱུ་ཨྱུ་ཨྱུ་ཨྱུ་’ and ‘རྩུ་ཨྱུ་ཨྱུ་ཨྱུ་ཨྱུ་’; Finally, compare the two letter series as we compare two English strings and we have got the collation result of two characters རྩུ་ and རྩུ་. Table 1 gives some more examples of this kind collation.

Table 1. The collation of the transliterating characters with the rules of the transliterating character dictionary

Characters	Syllable series	Letter series
རྩུ་	རྩུ་	རྩུ་ཨྱུ་ཨྱུ་ཨྱུ་ཨྱུ་
རྩུ་	རྩུ་	རྩུ་ཨྱུ་ཨྱུ་ཨྱུ་ཨྱུ་
རྩུ་	རྩུ་	རྩུ་ཨྱུ་ཨྱུ་ཨྱུ་ཨྱུ་
རྩུ་	རྩུ་	རྩུ་ཨྱུ་ཨྱུ་ཨྱུ་ཨྱུ་
རྩུ་	རྩུ་	རྩུ་ཨྱུ་ཨྱུ་ཨྱུ་ཨྱུ་
རྩུ་	རྩུ་	རྩུ་ཨྱུ་ཨྱུ་ཨྱུ་ཨྱུ་
རྩུ་	རྩུ་	རྩུ་ཨྱུ་ཨྱུ་ཨྱུ་ཨྱུ་

4.2 Collated with the Rules of Native Tibetan Syllable Dictionary

A typical Tibetan syllable, such as རྩུ་, is a two-dimensional combination of its letters. To syllable རྩུ་, the letter རྩུ་ at the center is the base consonant, the letter རྩུ་ above the base consonant is the head consonant, the letter རྩུ་ in the prefix position is the

5 Conclusion

Compared with the native Tibetan characters, the transliterating characters are used not so popularly; however, there are more than six thousands of them. Therefore, it is necessary to study the collation of these transliterating characters. The paper proposes two structures that can deal with the two kinds of collation of transliterating characters: collated with rules of native Tibetan dictionaries and with the rules of transliterating dictionaries. Based on the proposed structures, all transliterating characters can be collated successfully and effectively with the rules of two different dictionaries.

Acknowledgment. This work is partially supported by NSFC under Grant No.60963016 and Key laboratory of Tibetan Information Processing, Ministry of Education of the People's Republic of China. The authors also thank the anonymous reviewers for their invaluable comments and suggestions.

References

1. An, S.: Sanskrit-Tibetan-Chinese dictionary. Nationalities Publishing House, Beijing (1991)
2. Zhang, Y.: Tibetan-Chinese Dictionary. Nationalities Publishing House, Beijing (1985)
3. Jiang, D., Kang, C.: The sorting mathematical model and algorithm of written Tibetan language. *Chinese Journal of Computers* 27(4), 524–529 (2004)
4. Huang, H., Da, F.: General Structure Based Collation of Tibetan Syllables. *J. Inf. Comput.* 6(5), 1693–1703 (2010)
5. Huang, H., Da, F.: Discussion on Collation of Tibetan Syllables. In: *IALP 2010*, pp. 35–38 (December 2010)
6. National Standard of PRC, Information Technology-Tibetan Coded Character Sets for Information Interchange-Extension A (GB/T 20542-2006). Standards Press of China, Beijing (May 2007)
7. National Standard of PRC, Information Technology-Tibetan Coded Character Sets for Information Interchange-Extension B (GB/T 22238-2008). Standards Press of China, Beijing (January 2009)