

Calculation analysis on Consonant and character for Corpus study of Gesar epic "HorLing"

DuoLa¹,TashiGyal^{2*}

1.Northwest University for Nationalities,Lanzhou,Gansu,China,730030

2.Tibet University,Lhasa,Tibet,China,850000

DuoLa67@126.com zzzx.77@163.com

ABSTRACT. We made an econometric analysis on consonants and characters after the establishment of Gesar epic classic version HorLing corpus. Firstly, we set up a 2-million-consonants corpus for further verification and comparison of the character frequency of HorLing. Secondly, we established the theory of Tibetan consonants combination rules and Tibetan theory consonants and wish the Tibetan theory consonants coverage on HorLing. By the analysis, we not only understood the status of the Gesar epic consonants and characters, but also clarified the application of Tibetan consonants and characters in the real life.

Keywords: Gesar; Tibetan corpus; consonant and characters

1 Introduction

It is well known that HorLing is the classic version of Gesar epic and the body of the epic. Gesar is the world's longest epic and "living epic". It is still widely talked and sung so far. In spoken and written Tibetan language and dialect differences, it's between spoken and written language. It broke through each dialect barriers and became a universal language. It is important to study the language, which have a restricted view for a nation's way of thinking and attitude of observing.

This paper makes the following contributions: 1) make a qualitative calculation research on vowel sounds; 2) calculate Tibetan theory characters by rule description and software generation; 3) compare with the consonant and characters of Gesar epic, which helps master the language system of the epic; 4) provide a reference data for Tibetan language research and Tibetan information processing. Therefore, it is a very meaningful experiment.

2 Vowels and frequency in HorLing

In ISO/IEC10646 "Tibetan coding Basic Collection" (hereinafter referred to as the "Basic Collection"), the Tibetan is coded according to the dynamic superposition

* Corresponding author

roots in the scope of letters. They represent not only latter roots, but appear as roots and act as latter root. The other five could also act as previous root. "རྩལ་ལྷན་" could be used as ko can, but they cannot used as ko can roots. The diversity of the role makes them stand out in more than 1 million consonants and become a top ten high frequency consonants. More impressive is that the cumulative frequency of these 10 consonants and 4 vowels achieved 76.0236%. That is to say, 1014827 consonants in the full text of Horling except syllables and punctuation, 76% of the 14 consonants were repeated and composed in the various roles.

The highest level of the art is not to make the simple things complicated, but to describe the complicated things by the most simple and fluent language, as is not only the art and science. If this view is correct, we could say that the epic of "Gesar" is a perfect combination of art and science. We could also make this question, of course, the characteristics of the phonetic writing is to realize complex expression by a few letter combinations. There is no doubt that it should be acknowledged. What make us surprising is that 30 letters of alphabet writing should only use 10 of them and reached 76% of expression. This work is a miracle.

In 168 coded consonants version of "Basic collection", there has 77 consonant characters [3], including Tibetan single consonant, 49 dok can root character, 28 Sanskrit consonant character and dok can characters. In the epic Gesar, statistics of consonant and components of "HorLing". There are 49 pureTibetan characters except seven vowels and six Sanskrit characters. It shows that in the level of consonant and component, "Gesar" epic "HorLing" covers 100% of the modern Tibetan consonant and component, effectively prove "HorLing" has important value as a template for Tibetan language research.

4 About the character

Tibetan has 206 basic character in traditional grammar (ལྷན་ལྷན་ལྷན་)[4] if "body" understood as character, the statistics might not consider other factors. But the basic character is put forward, which gives us a consideration and evolves "character". For the attention of the "character" and put forward again of this concept due to the rise of Tibetan information technology, the problem of character were raised at the beginning of encoding for Tibetan.

In terms of Chinese characters, a word can be classified as a character, such as 6763 Chinese characters were collected in GB 2312, including 3755 primary/first-level Chinese characters, 3008 secondary characters, which basically meet the requirements of regular users of computer processing of Chinese characters, and its coverage to mainland China 99.75% of Chinese characters operating frequency. From which we can see that every word is a syllable unit, at the same time it is a character.

acter, and appear some Sanskrit categories, but the frequency of the Sanskrit is very low, just appear while reciting om mani padme hung and ali's four great rivers, and several animal and plant names such as "water lotus", "lion"etc. 34 Sanskrit and frequency table (table 1) as follows:

Table 1. Frequency table of the Sanskrit and frequency

	charac-	time	frequen-	No	charac-	time	frequen-	Amount
1	ॐ	142	0.0205	18	ॐ	2	0.0003	69335
2	ॐ	30	0.0043	19	ॐ	1	0.0001	69335
3	ॐ	22	0.0032	20	ॐ	1	0.0001	69335
4	ॐ	17	0.0025	21	ॐ	1	0.0001	69335
5	ॐ	13	0.0019	22	ॐ	1	0.0001	69335
6	ॐ	12	0.0017	23	ॐ	1	0.0001	69335
7	ॐ	7	0.0010	24	ॐ	1	0.0001	69335
8	ॐ	6	0.0009	25	ॐ	1	0.0001	69335
9	ॐ	5	0.0007	26	ॐ	1	0.0001	69335
10	ॐ	5	0.0007	27	ॐ	1	0.0001	69335
11	ॐ	5	0.0007	28	ॐ	1	0.0001	69335
12	ॐ	4	0.0006	29	ॐ	1	0.0001	69335
13	ॐ	3	0.0004	30	ॐ	1	0.0001	69335
14	ॐ	2	0.0003	31	ॐ	1	0.0001	69335
15	ॐ	2	0.0003	32	ॐ	1	0.0001	69335
16	ॐ	2	0.0003	33	ॐ	1	0.0001	69335
17	ॐ	2	0.0003	34	ॐ	1	0.0001	69335

(Note: not including in the total number of syllable point)

As you can see, in 34 Sanskrit character, a total of 17 character appear twice or more, the rest is only appear one time. The highest frequency in Sanskrit is "ॐ", a total of 142 times, in fact the character and its combination are both "ॐ", it is a character with different form of "ॐ" (hero)".

We already know that the number of modern Tibetan character is 456, which appear 411 times in "HorLing", there are still 45 Tibetan character does not appear, and sill 40 characters did not appear in the epic except 5 new character "ॐ". 411 Horling character accounted for 90.1316% in modern Tibetan characters, there are nearly 10% of the them has not been covered. If some high-frequency character did not appear, then typicality of the HorLing language will be questioned. Therefore, we need to know which character does not appear, and the function of those does not appear. Through comparison and selection, we found that the 40 characters are: "ॐ"

2	ག	48251	6.9591	101920	14.6995	693356
3	ད	46848	6.7567	148768	21.4562	693356
4	ཅ	44149	6.3674	192917	27.8237	693356
5	ཆ	44031	6.3504	236948	34.1741	693356
6	ཇ	34726	5.0084	271674	39.1825	693356
7	ཉ	34056	4.9118	305730	44.0942	693356
8	ཏ	31894	4.5999	337624	48.6942	693356
9	ཐ	27888	4.0222	365512	52.7164	693356
10	འ	25557	3.6860	391069	56.4023	693356

In characters, it is still a"འ" in the first place, the second and third places are "ག" and "ད" same as characters statistics layout, and then change on the ordinal, but on the whole, still the 10characters.

In the top 100characters, mainly is function word, such as five the auxiliary case are in the top 50. They are, "འ" ranks 12th, appear 10408 times, frequency is 1.5011%; "ལ" in the 18th, 5054 times, the frequency is 0.7289%; "ཞ" in the 33th, 3171 times, frequency is 0.4573%; "ཀ" in the 44th, 2568 times, frequency is 0.3704%; "ཡ" is 46th, 2417 times, frequency is 0.3486%. The frequency of ordering shows that it is obvious for the adhesion phenomenon(function word cohere with national word) of the case.

The cumulative frequency of the first 100 characters is up to 625379 words, accounted for 89.94% of total corpus, the lowest frequency is 945 times. The 100 characters appear 1033 times except the end of the last three characters. So they can be called thousand-times character.

From the number of 101to 200, character root has taken main status. The highest frequency of this phase is 935 times, the lowest is 265 times of the 200th character. Which is relatively balanced, unlike the frequency of the first 100 characters, the first appeared 53669 times, and 100th appear 945 times, they are differ more than 52000 times.

The cumulative frequency of the first 200 characters is up to 675826 times, accounted for 97.47% of the total number. From the number of 101 to 200, mostly is superposition character, followed by consonant and vowel, illustrate the major component of Tibetan alphabet combination in this frequency. The character without superposition and vowel is only one "ེ", its number is 107, still in the first place, appear 872 times, frequency is 0.1258%.

The frequency and distribution of 447 characters shows that the reflects the usage and distribution of Tibetan characters in epic, on the other hand it also reflects a general situation of the entire Tibetan character. All 447 characters including Sanskrit, there are 29 of which frequency is only once, 10 of which frequency is 2, seven of which frequency is 3, five of which frequency is 4, three of which frequency is 6, only one of which frequency is 7, two of which frequency is 8, there of which frequency is 9, two of which frequency is 10, there of which frequency is 11, there of which frequency is 12, there of which frequency is 13.As you can see from these frequency, 2-

3characters distributed for each level in low-frequency ones equably. Throughout the 447 characters, there are 115 characters which frequency are below 30.The rest 33 Sanskrit except "ཨ" are mixed In this frequency channel. It can be seen that the character of this frequency channel is not "uncommon character".

7 Conclusion

"HorLing", commonly known as "HorLing wars "or" the war of HorLing", mainly describes the war events between tribe of Hore and Ling, shows the changes of the Tibetan tribal society. Mr Qian mintz wrote in his article: "HorLing wars is a historical picture scroll of Tibetan hero's grand, majestic and grand, it is permeated with enthusiasm upward, positive enterprising and fight passion. This is very similar to the Iliad"[9] He also asserts that "Tibetan HorLing wars is the Chinese nation's "Iliad""[10].

It is earlier to study of A Dream of Red Mansions by corpus method, which inspired us for the first time to study Gesar epic by establishing Tibetan tagging corpus. Though the study we not only get the number of Tibetan theory character, but also verified the frequency of character in the practical application, and classified the characters. Moreover, we verified the 40 low-frequency characters in a larger range (2 million characters) the verification of the corpus to further confirm the reliability of the frequency of characters in "HorLing", it makes sense.

ACKNOWLEDGEMENTS

This research project was financially supported by the National Natural Science Foundation of China (NSFC)(No. 61262053) and The ministry of education philosophy and social sciences research project for major project(No. 13JZD028), the National Natural Science Foundation of China (NSFC)(No. 61163043).

REFERENCES

1. Di Jiang,Cong Jun Long.:Study on Tibetan character, Social sciences academic press, Aug, pp.88-89(2010)
2. Qian mintz, HorLing wars, the Iliad, Foreign literature, pp.125-251 (1986)
3. Di Jiang,Cong Jun Long.: Study on Tibetan character, Social sciences academic press, Aug, pp.77-78(2010)
4. Caidanxiarong.:The Tibetan language grammar, Qinghai minorities press, pp.25-26(1998)
5. Edited by Luo Bingfen, Tubo Medical Literature Essential, Ethnic Publishing House, pp. 127-300(2002)
6. Edited by Chenjian, Wangyao.:Duanguang Manuscript Tibetan Literature, ethnic publishing house , pp.12-321(1983)
7. Compiled by Tsering Tso etc.: Atomic Physics (Tibetan-Chinese Languages), Gansu Ethnic Publishing House, pp.7-29(2005)

8. Compiled by Zenlha Ngawang Tsutrin, Ancient Tibetan Dictionary, Ethnic Publishing House, pp.15-16(1997)
9. Qian mintz, HorLing wars, the Iliad, Foreign literature, pp.103-105 (1986)
10. Qian mintz, HorLing wars, the Iliad, Foreign literature, pp.114-120(1986)