

# CTTC: A Collection of Tibetan Text Corpora

Huidan Liu<sup>a</sup>, Congjun Long<sup>b</sup>, Longlong Ma<sup>a</sup>, Jian Wu<sup>a</sup>, Le Sun<sup>a</sup>

<sup>a</sup>.Institute of Software, Chinese Academy of Sciences, Beijing, China, 100190

<sup>b</sup>.Institute of Ethnology and Anthropology, Chinese Academy of Social Sciences, Beijing, China, 100081

huian@iscas.ac.cn, longcj@cass.org.cn, wujian@iscas.ac.cn, sunle@iscas.ac.cn

## Abstract

The Chinese Academy of Sciences launched the Multi-Layer MultiLingual Resource Database (MLLRD) project which aims to collect language resources for natural language processing tasks for low resource languages used in China, such as Mongolian, Tibetan, Uyghur and so on. Tibetan text corpus building is one of the sub projects, in which we have built a Collection of Tibetan Text Corpora(CTTC), including: (1) Tibetan web article corpus which has 440,900 documents. (2)Tibetan text classification corpus. (3) Chinese-Tibetan parallel text corpus which has 773,068 sentence pairs. (4) Part-Of-Speech tagged corpus which has 52,041 sentences. (5) Tibetan tree bank which has 6,040 trees. The paper reports the methods to build these corpora, the contents and scales of each corpus, and applications of them.

**Keywords:** Tibetan, Corpus, Machine Translation, Tree Bank

## 1. Introduction

Corpora are the basic and necessary materials for natural language processing. The Chinese Academy of Sciences launched the Multi-Layer MultiLingual Resource Database (MLLRD) project which aims to collecting language resources for natural language processing tasks for low resource languages used in China. The project collects resources generally for three tasks, namely machine translation, speech recognition, hand written character recognition. For machine translation task, bilingual sentence level aligned parallel text are collected for Chinese-Tibetan, Chinese-Uyghur and Chinese-Mongolian. There are more than 300 thousand sentence pairs for any of the three language pairs. For speech recognition task, text-speech aligned sentences are collected for Tibetan, Uyghur and Mongolian. There is 360GB speech data in total. For hand written character recognition, hand written characters from 300 writers are collected for each of the three languages. Tibetan text corpus building is one of the sub projects. In the sub project we have built a Collection Tibetan Text Corpora(CTTC), including: (1) Tibetan web article corpus. (2)Tibetan text classification corpus. (3) Chinese-Tibetan parallel text corpus. (4) Part-Of-Speech tagged corpus. (5) Tibetan tree bank. We introduce these corpora in the following sections.

## 2. Tibetan Web Article Corpus

### 2.1. Sources of the Corpus

Previously Liu et al. (2012b) proposed an approach to build a large scale text corpus for Tibetan natural language processing. We adopt the method to build our corpus. We use a web crawler initialized with a seed URL list, which includes some well-known Tibetan websites. Then we check the crawled web page whether it contains Tibetan text with a Tibetan examiner, and if a page has Tibetan text in it, all URLs which it links to are appended to the fetching list of the crawler. The procedure continues until no new Tibetan web pages are found. After that we know where to get Tibetan text. For Tibetan web article corpus, we crawled arti-

cles from 19 Tibetan websites which mainly focus on news and broadcastings(Table 1).

1	<a href="http://blog.amdotibet.cn">http://blog.amdotibet.cn</a>
2	<a href="http://epaper.chinatibetnews.com">http://epaper.chinatibetnews.com</a>
3	<a href="http://tb.chinatibetnews.com">http://tb.chinatibetnews.com</a>
4	<a href="http://tb.tibet.cn">http://tb.tibet.cn</a>
5	<a href="http://tb.xzxw.com">http://tb.xzxw.com</a>
6	<a href="http://ti.gzznews.com">http://ti.gzznews.com</a>
7	<a href="http://ti.kbcmw.com">http://ti.kbcmw.com</a>
8	<a href="http://ti.tibet3.com">http://ti.tibet3.com</a>
9	<a href="http://tibet.cpc.people.com.cn">http://tibet.cpc.people.com.cn</a>
10	<a href="http://tibet.people.com.cn">http://tibet.people.com.cn</a>
11	<a href="http://tibetan.qh.gov.cn">http://tibetan.qh.gov.cn</a>
12	<a href="http://www.amdotibet.cn">http://www.amdotibet.cn</a>
13	<a href="http://www.qhtb.cn">http://www.qhtb.cn</a>
14	<a href="http://www.tbmgar.com">http://www.tbmgar.com</a>
15	<a href="http://www.tibet3.com">http://www.tibet3.com</a>
16	<a href="http://www.tibetnr.com">http://www.tibetnr.com</a>
17	<a href="http://www.tibetology.ac.cn">http://www.tibetology.ac.cn</a>
18	<a href="http://www.vtibet.com">http://www.vtibet.com</a>
19	<a href="http://xizang.news.cn">http://xizang.news.cn</a>

Table 1: Sources of Tibetan web article corpus.

### 2.2. URL Filtering

Web pages can be classified into two kinds, namely “topic” and “hub”. A topic page contains long text in it while a hub page contains many links to the topic pages. As our target is to extract Tibetan web articles from the web pages, We only care about the topic pages rather than the hub pages. Topic pages rather than hub pages are selected with a rule based method by checking the url.

Table 2 and Table 3 show some URLs of topic pages and hub pages of the three Tibetan web sites respectively. Comparing tens of thousands of URLs of the three web sites, we find the following rules:

- The topic URLs of “Tibetan Web of China” have the pattern of “{host}/{column}/{year}-

Site	Example URLs
China Tibet Online	http://tibet.people.com.cn/141101/15137028.html http://tibet.people.com.cn/141101/15199715.html http://tibet.people.com.cn/15143391.html
Tibetan's Web of China	http://ti.tibet3.com/folkways/2008-12/10/content_370366.htm http://ti.tibet3.com/medicine/2009-10/27/content_99171.htm

Table 2: Example URLs of topic pages.

Site	Example URLs
China Tibet Online	http://tibet.people.com.cn/140827/141059/index3.html http://tibet.people.com.cn/96372/125163/index.html http://tibet.people.com.cn/141101/index11.html
Tibetan's Web of China	http://ti.tibet3.com/culture/index.htm http://ti.tibet3.com/tour/node_701.htm http://ti.tibet3.com/economy/index.htm

Table 3: Example URLs of hub pages.

{month}/{date}/content\_{articleid}.htm". Everyone of them contains the string "content\_".

- The hub URLs of "Tibetan Web of China" contain the string "index" or "node".
- The topic URLs of "China Tibet Online" have the pattern of "{host}/{columnid}/{articleid}.html". Characters between the host URL "{host}" and the file suffix name "html" are numbers or slash.
- The hub URLs of "China Tibet Online" contain the string "index".

With these rules, we make text extraction only on the topic pages.

### 2.3. Text extraction

We analysed the layout structure of the web pages from each web site and get clues to build templates to extract topic title, publishing date, author, topic content and some other topic related informations. Figure 1 shows the structure of a web page<sup>1</sup>. From the figure, we see that there are some HTML tags giving the boundaries of different text blocks, and we can find the corresponding HTML tags of the article title, publishing date, author, article content and so on.

### 2.4. Content of the Corpus

At present, the corpus has about 440.9 thousands documents, 9.50 million sentences, 228 million syllables in total. Each Article is saved as an XML file. Figure 2 shows an article from the corpus.

### 2.5. Quality of the Corpus

Some predefined rules are used to check whether there are spelling errors in a syllable in a previous of the corpus. The statistical data show that there are 9700 misspelt ones out of the 20743 Tibetan syllables occurred in the corpus, which shares 46.7628%. But their occurrence is only 27,427 in the 93 million syllable in total, sharing only 0.0308% (Liu et al., 2017), which shows that the corpus has a very high quality.

<sup>1</sup>http://tibet.people.com.cn/15260188.html

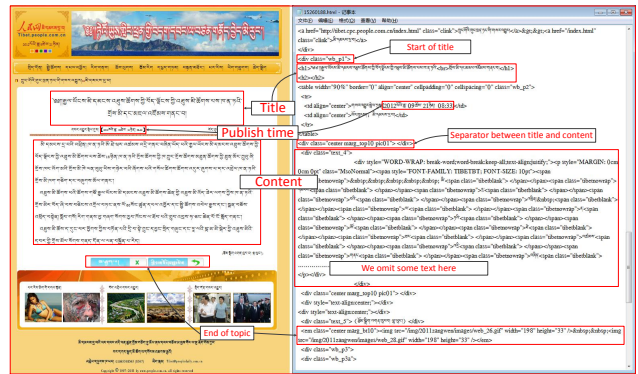


Figure 1: The structure of a web page from "China Tibet Online".

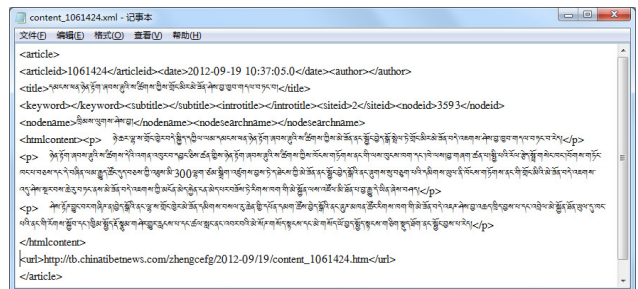


Figure 2: The text extracted from a page from "China Tibet news", in XML format.

## 3. Tibetan Text Classification Corpus

The web article corpus is further processed to build a text classification corpus. It's a heavy task to manually classify those document into domains. However, we can get the domain information for a certain subsets of the web article corpus. For some web sites listed above, we can get the domain information from the URL of each web page. For instance, the URL "http://tb.chinatibetnews.com/xzmeishi/2011-12/05/content\_831210.htm" shows it belongs to a column called "xzmeishi". so it must be a page about Tibetan foods, because "xz" is the abbreviated form of Chinese word "xizang" (西藏), which means the Tibetan Autonomous Region, while "meishi" means "delicious food". So we can classify the documents in the corpus into domains. Table 4 and 5 list the domains of two subsets of the articles from two web sites named "China Tibet News" and "Tibetan's web of China" respectively. Obviously, a large part of the documents in the corpus are news as expected, because the two web sites are both hold by news agencies.

## 4. Chinese-Tibetan Parallel Corpus

### 4.1. Sources of the Corpus

We get documents for the Chinese-Tibetan parallel corpus from two types of sources. The first source of the corpus is translation agencies. A large part of documents in our corpus are collected from several translating agencies. As most of them are translated from Chinese to Tibetan, we

	Domain	#doc	(%)	#sent	(%)
1	Art	3,240	4.76	112,642	8.71
2	Economy	712	1.05	12,477	0.96
3	History	2,897	4.25	19,627	1.52
4	News	25,247	37.08	576,842	44.59
5	Picture	12,732	18.70	51,088	3.95
6	Politics	3,230	4.74	63,437	4.90
7	Rural Life	2,402	3.53	35,535	2.75
8	Social Life	1,153	1.69	9,881	0.76
9	Specials	9,986	14.67	268,003	20.72
10	Technology	1,988	2.92	38,321	2.96
11	Buddhism	1,983	2.91	48,832	3.77
12	Food	215	0.32	2,963	0.23
13	Medicine	720	1.06	36,676	2.84
14	Tour	1,588	2.33	17,296	1.34
Total		68,093	100.00	1,293,620	100.00

Table 4: Domains of a subset of the documents from “China Tibet News”.

Order	Domain	#doc	(%)	#sent	(%)
1	Art	92	0.35	3,021	0.45
2	Culture	885	3.40	109,749	16.18
3	Economy	78	0.30	7,749	1.14
4	Education	15	0.06	695	0.10
5	Music	323	1.24	3,169	0.47
6	News	24,055	92.45	519,576	76.61
7	Photo	80	0.31	2,548	0.38
8	Policy	116	0.45	7,062	1.04
9	Politics	124	0.48	7,668	1.13
10	Medicine	107	0.41	11,417	1.68
11	Tour	145	0.56	5,563	0.82
Total		26,020	100.00	678,217	100.00

Table 5: Domains of a subset of the documents from “Tibetan’s web of China”.

know the correspondence between the Chinese part and the Tibetan part when we got them. We have nearly 600 thousand sentence pairs from the first source.

The second source of the corpus is the web. We collected articles from two web sites as listed in Table 6 which publish articles in both Chinese and Tibetan. They mainly focus on news and broadcastings. We have about 202 thousand sentence pairs from the second source.

	Host	Language
1	http://tb.xzxw.com	Tibetan
2	http://www.xzxw.com	Chinese
3	http://ti.tibet3.com	Tibetan
4	http://www.tibet3.com	Chinese

Table 6: Sources of the bilingual corpus.

## 4.2. Document Alignment

We use a feature based method to find the Chinese correspondence for each Tibetan article. Three kinds of features are used: numbers, common punctuations and geographic names in the context of each document. Numbers

and some punctuations have same presentation in Chinese and Tibetan while geographic names are translated fixedly. Thus we regard them as good clues to make the document alignment.

### 4.2.1. Number Extraction

Table 7 shows three ways to present Tibetan numbers. In our method, we extract first two forms of numbers in Tibetan documents, and transfer the numbers presented as Tibetan symbol digits to Arabic numbers.

Form	Description	Example
Arabic numbers	consist of Arabic number (0 to 9)	“2012”
Tibetan digital numbers	alike Arabic numbers, consist of Tibetan digital character	“ཨ་ཨ་ཨ་” (2010)
Tibetan syllable numbers	consist of one or several Tibetan syllables	“ཨ་ཨ་ཨ་ཨ་” (15)

Table 7: Three ways to present Tibetan numbers

Table 8 shows two ways to present Chinese numbers. In our method, we extract first form of numbers in Chinese documents.

Form	Description	Example
Arabic numbers	consist of Arabic digit (0 to 9)	‘2012’
Chinese digital numbers	consist of one or several Chinese syllables	“十五” (15)

Table 8: Two ways to present Chinese numbers

### 4.2.2. Punctuation Extraction

Chinese and Tibetan have their own punctuation marking system respectively. However, Tibetan borrows some Chinese punctuations, such as parentheses “()”, book title mark “” and double quotation mark. In our method, we extract these three punctuation marks as features for they will be preserved in the same form when an article is translated from Chinese to Tibetan or from Tibetan to Chinese.

### 4.2.3. Geographic Names Extraction

We use a bilingual dictionary of Chinese and Tibetan, which consist of most of place of interest and administrative division in Tibet to extract geographic names in articles with maximum matching method. Tibetan geographic names are translated to Chinese which is taken as the features to make document alignment.

### 4.2.4. Candidate Document Pair Generation

In the Internet, there are millions of Chinese and Tibetan documents, so it’s necessary to filter document pairs that are impossible to be parallel. As the number of extracted Chinese articles are much larger than that of Tibetan ones,









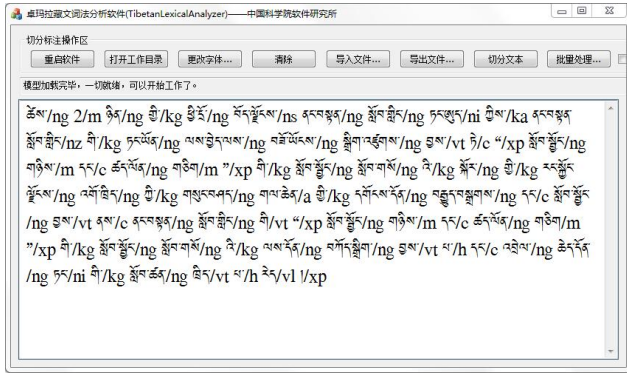


Figure 8: Tibetan lexical analyser.

### 7.7. Tibetan Parsing

We use the Berkeley parser to train a Tibetan parser with a former version of the Tibetan tree bank. The training set and the test set include 3,746 and 354 trees respectively. Experiments show that if the POS tags are provided, the parser achieves a better performance. The precision, recall and F1 are 0.9251, 0.9273 and 0.9262 respectively.

### 8. Conclusion

As a low-resource language, Tibetan language processing is facing a big challenge. In this paper, we introduced our work on building a collection of Tibetan text corpora(CTTC), including: (1) Tibetan web article corpus. (2)Tibetan text classification corpus. (3) Chinese-Tibetan parallel text corpus. (4) Part-Of-Speech tagged corpus. (5) Tibetan tree bank. The corpora are applied in many research tasks such language modelling, machine translation, lexical analysis, text classification, parsing and so on.

In the future, we will collect more web text to increase the scales of these corpora, especially for Tibetan tree bank as its scale is still small. We will also make more annotations based on the existing corpora to build corpora for other Tibetan NLP tasks. CTTC is available for academic researches by contacting the authors.

### 9. Acknowledgements

We thank the reviewers for their critical and constructive comments and suggestions that helped us improve the quality of the paper. The research is partially supported by Informationization Project of the Chinese Academy of Sciences (No.XXH12504-1-10) and Research Project of the National Language Committee(ZDI135-17).

### 10. Bibliographical References

Ai, J., Yu, H., and Li, Y. (2009). Statistical analysis on tibetan shaped structure. *Journal of Computer Applications*, 29(7):2029–2031.

Gao, D. and Gong, Y. (2005). A statistically study on the qualities of all modern tibetan character set. *Journal of Chinese Information Processing*, 19(1):71–75.

Jiang, D. and Dong, Y. (1994). Statistical analysis on linear processing of tibetan clustered structures. *Chinese Information Processing*, (4):44–46.

Jiang, D. and Dong, Y. (1995). Research on property of tibetan characters as information processing. *Journal of Chinese Information Processing*, 9(2):37–44.

Jiang, D. and Kong, J. (2006). *Advances on the Minority Language Processing of China*. Social Sciences Academic Press, Beijing, China.

Jiang, D. and Long, C. (2010). *On Characters of Tibetan Writing System: Alpbabetic Characters, Pronunciations, ISO Codes, Frequencies, Sorting Orders, Picture Symbols and Transliterations*. Social Sciences Academic Press, Beijing, China.

Jiang, D. (2006). History and development of tibetan text information processing. In *Frontiers of Chinese Information Processing - Proceedings of the 25th Anniversary Conference of Chinese Information Processing Society of China*, pages 83–97. Tsinghua University Press, Beijing, China.

Liu, H., Nuo, M., Ma, L., and et al. (2011). Tibetan Word Segmentation as Syllable Tagging Using Conditional Random Fields. In *Proceedings of the 25th Pacific Asia Conference on Language, Information and Computation (PACLIC 2011)*, pages 168–177.

Liu, H., Nuo, M., Ma, L., and et al. (2012a). Segt: A pragmatic tibetan word segmentation system. *Journal of Chinese Information Processing*, 26(1):97–103.

Liu, H., Nuo, M., Wu, J., and He, Y. (2012b). Building large scale text corpus for tibetan by extracting text from web pages. In *Proceedings of the 10th asian language resources at COLING 2012*, pages 8–17.

Liu, H., Long, C., Nuo, M., and Wu, J. (2015). Tibetan word segmentation as sub-syllable tagging with syllables part-of-speech property. In Maosong Sun, et al., editors, *Chinese computational linguistics and natural language processing based on naturally annotated big data (LNAI 9427)*, pages 189–201.

Liu, H., Hong, J., Nuo, M., and Wu, J. (2017). Statistics and analysis on spell errors of tibetan syllables based on a large scale web text corpus. *Journal of Chinese Information Processing*, 31(2):61–70.

Lu, Y., Ma, S., Zhang, M., and Luo, G. (2003). Researches of calculations of tibetan characters, pieces, syllables, vocabulary and universal frequency and its applications. *Journal of Northwest Minorities University(Natural Science)*, 24(48):32–42.

Yu, X., Wu, J., and Hong, J. (2011). Research and realization of dictionary-based chinese-tibetan sentence alignment. *Journal of Chinese Information Processing*, 25(4):57–62.