

# himalayan linguistics

A free refereed web journal and archive devoted to the study of the  
languages of the Himalayas

## Himalayan Linguistics

---

*A Chinese to Tibetan machine translation system with multiple translating strategies*

**Huidan Liu**

Chinese Academy of Sciences

**Weina Zhao**

Chinese Academy of Sciences; Qinghai Normal University

**Minghua Nuo**

Chinese Academy of Sciences; Inner Mongolia University

**Jinling Hong**

Chinese Academy of Sciences

**Xin Yu**

Chinese Academy of Sciences

**Jian Wu**

Chinese Academy of Sciences

### ABSTRACT

This paper proposes a Chinese to Tibetan machine translation system with multiple translating strategies. The key corpora and technologies are explained in detail. Experiments show the sub-systems output the translation of each phrase in the same order as they are in the Chinese sentence rather than in a Tibetan sentence, which leads to unsatisfactory translation quality. Consequently, an order adjusting model is essential to Chinese to Tibetan translation system. The phrase recall of the SMT (MaxProb) sub-system makes an improvement of 9.71% over the popular off-the-shelf language neutral statistical machine translation programme Moses. Our translation system achieves a speed of about 0.175s per sentence on a PowerEdge R710 server with 2 4-core Intel E5620 (2.40 GHz) CPUs and 16 GB memories, which meets the requirement of a computer aided translation system.

### KEYWORDS

Tibetan, NLP, machine translation

This is a contribution from *Himalayan Linguistics*, Vol. 15(1): 149–166.

ISSN 1544-7502

© 2016. All rights reserved.

This Portable Document Format (PDF) file may not be altered in any way.

Tables of contents, abstracts, and submission guidelines are available at  
[escholarship.org/uc/himalayanlinguistics](http://escholarship.org/uc/himalayanlinguistics)

# *A Chinese to Tibetan machine translation system with multiple translating strategies*

**Huidan Liu**

Chinese Academy of Sciences

**Weina Zhao**

Chinese Academy of Sciences; Qinghai Normal University

**Minghua Nuo**

Chinese Academy of Sciences; Inner Mongolia University

**Jinling Hong**

Chinese Academy of Sciences

**Xin Yu**

Chinese Academy of Sciences

**Jian Wu**

Chinese Academy of Sciences

## **1 Introduction**

The Tibetan language is the main carrier of Tibetan culture; it plays a very important role in Tibetan people's daily communication (Chen 2003). According to a previous investigation, there are about 100 organizations which have Chinese-Tibetan translation business in Tibet and about 1,000 translators who make translation their profession. But the annual amount of Chinese text to be translated into Tibetan is more than 50 million Chinese characters every year (Luo et al. 2010). Thus, there is a big gap between the translating ability and the market requirement. Machine translation technology is an urgent requirement, which contributes to the improvement of translation speed as well as its quality.

In the paper, we propose a Chinese to Tibetan machine translation system with multiple translation strategies. The paper is organized as follows: In Section 2 we recall related work on machine translation in general and Tibetan related machine translation in particular. In Section 3, we introduce the architecture of our machine translation system. The core corpora and technologies are explained in Section 4. Then, in Section 5, we experimentally compare the performances of the system with some other methods. Section 6 offers some conclusions.

## 2 Related Work

There are mainly three types of machine translation methods, namely rule-based machine translation (RBMT), example-based machine translation (EBMT), and statistical machine translation (SMT). Research of Chinese-Tibetan machine translation focused on rule-based methods before 2010 due to a lack of parallel corpora and other basic Tibetan language resources (Degai 2001; Cai 2005; Zhaluo 2005; Kanzhuo et al. 2006). In recent years, statistical Chinese-Tibetan machine translation methods are playing a more and more important role.

Generally speaking, a large rule set is needed in a RBMT system. However, it takes huge human resources and requires a high level of language knowledge. An EBMT system takes less human resources and requires a lower level of language knowledge. But if we can't find a very similar instance for the input text, the translation can't be satisfactory. A parallel corpus is essential to a SMT system. As long phrase takes more memory and will confront the problem of data sparseness, the well-known Moses statistical translation programme sets the limit of phrase length to 7 words or less (Koehn 2007).

Hou (2007) presents an example based Chinese-Mongolian machine translation method. The method consists of several parts, including example searching, segment splitting, matching and recombining. The method is based on word alignment. It uses word alignment information for segment matching, and computing the similarities by the number of matching words and length, and selects the best example. Kang et al. (2007) proposed a hybrid method which combines a statistical method with linguistic rules to extract Chinese multi-word chunks for translation purposes.

In Chinese-Tibetan machine translation, there are several RBMT systems reported. Cai (2005) built a Chinese-Tibetan machine translation system for government documents, which is based on dictionary and grammar rule templates. Degai et al. (2001) also built a machine translation system based on knowledge rules. Zhaluo (2006) proposed a method to build translation rules for complex sentences. Kanzhuo (2006) discussed the classification of verbs, different forms of verbs and the changeable regular pattern of verbs when they are in different tenses, and suggested several methods to improve the quality of translation.

There are two crucial problems to solve in Chinese-Tibetan machine translation. First, at the sentence level, all Chinese sentences have SVO structure, but Tibetan has SOV structure. So, long distant order adjusting is essential in Chinese-Tibetan machine translation. Second, verbs in Chinese do not inflect for tense, but in Tibetan a verb may have many variants. For instance, the forms of the verb "complete" are listed in Table 1.

	present tense	future tense	past tense	imperative mood
inactive voice	ལྷོབ་	བལྷོབ་	བལྷོབས་	ལྷོབས་
active voice	འལྷོབ་	འབལྷོབ་	ལྷོབ་	

Table 1. Forms of the Tibetan verb 'complete'

In this paper, we focus on building a Chinese to Tibetan machine translation system with multiple translating strategies to improve the translation quality subject to a poor language resource. It is expected that the system can find the translation if a sentence itself or a very similar sentence is in the parallel corpus, while it uses phrase-based machine translation decoding to generate a

translation in the many cases in which no similar sentences can be found, and a dictionary based decoding makes the greatest effort to assure that at least one translation is found for every word in the sentence.

### 3 System Structure

The proposed system is a combination of three different types of machine translation models, namely EBMT model, SMT model and RBMT model, as shown in Figure 1.

In the training phase of the system, Chinese articles and their translations are collected from several government translation organizations. These articles are processed into bilingual sentence pairs to form the bilingual sentence level parallel corpus. The sentence pairs are segmented into words by ICTCLAS (Zhang 2003) using Chinese and Tibetan word segmentation tools (Liu 2011; Liu 2012), and indexed to form the EBMT model. The segmented sentence pairs are word aligned and used to train the SMT model. Bilingual phrase pairs are extracted from the segmented sentence pairs. They form the RBMT model with word pairs from bilingual dictionaries (Zhang 1993; The ethnic publishing house 2002).

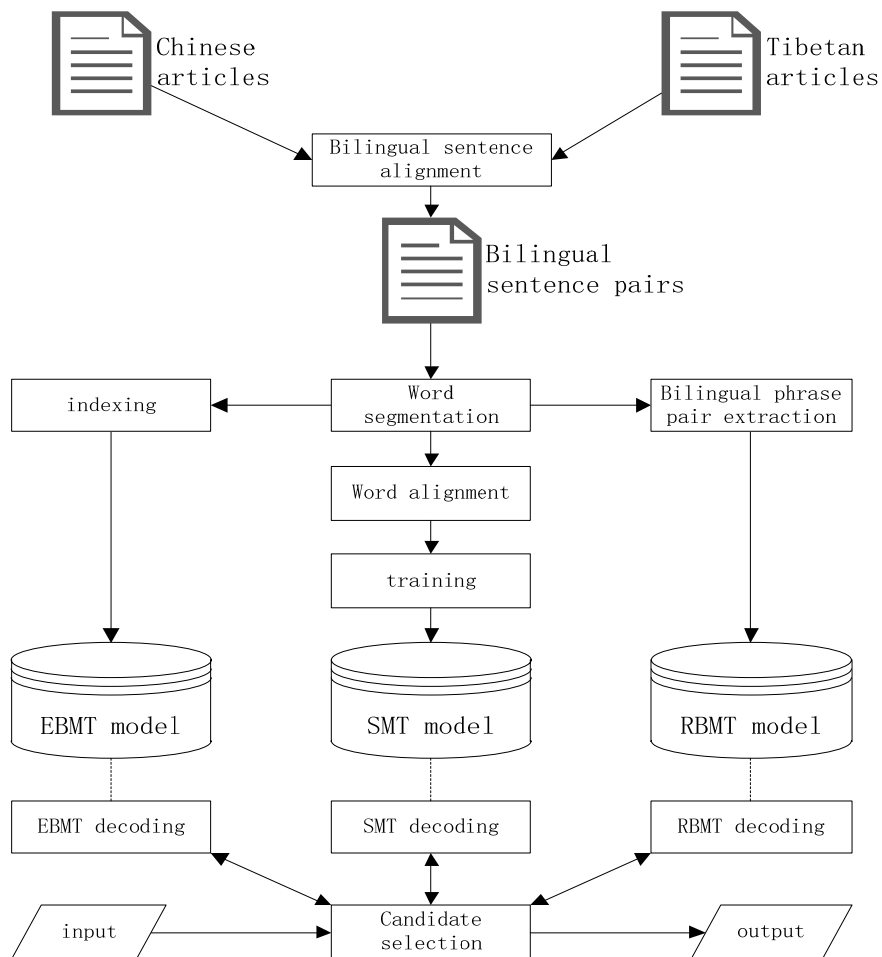


Figure 1. Data flow of the proposed system.

In the application phase of the system, as in the input of the system, a Chinese sentence is sent to the three sub-systems to be translated. The EBMT decoder segments the sentence into words, and tries to find similar sentences in the indexed bilingual sentence pairs by Levenshtein distance (Hirschberg 1975). If the EBMT decoder finds a Chinese sentence similar to the input one, and the similarity is larger than a predefined threshold, its Tibetan counterpart is extracted and taken as the output of the EBMT sub-system. Otherwise, the input sentence will be sent to the SMT decoder. The SMT decoder splits the input sentence into words and phrases, and computes the probability for every potential candidate. The probability determines how good each candidate is. The best candidate is selected as the output of the SMT sub-system. If the former two decoders cannot offer a good enough translation, the Chinese sentence is sent to the RBMT decoder, word to word translation is performed with a maximum matching algorithm with a trie structure (Fredkin 1960).

In the following section, we discuss the key corpora and technologies in both of the training phase and the application phase in detail.

## **4 Key corpora and technologies**

### **4.1 Tibetan word segmentation**

#### **4.1.1 Rule based method**

We design and implement a Tibetan word segmentation system named "SegT". It identifies critical words with a fast algorithm while segmenting each Tibetan sentence to chunks with case-auxiliary words, such as ཞེ, ཞུ, ཞུ, ཡི, ཞེས, ཞུས, ཞུས, and ཡིས. Each chunk is segmented into words by both forward maximum matching and backward maximum matching with a Trie tree structure<sup>1</sup>. It detects ambiguities by bidirectional segmentation, and disambiguates making use of pre-determined word frequencies, such that the segmentation yielding the more frequent words is selected.

In the procedure, the structure of each syllable is analyzed to identify abbreviated syllables while segmenting each block into words. In Tibetan text, some syllables, including འེ, ས, ར, འང, འམ, འོ (We call them abbreviation marker (AM) in this paper), can adhere to the previous word without a syllable delimiter "tsheg". The combination of these abbreviation markers with their preceding syllables produce 'abbreviated syllables'. For example, when the genitive case word འེ follows the word རྒྱལ་པོ་ (king), no "tsheg" appears between them and they are fused to form རྒྱལ་པོ་འེ (king[+genitive], king's), in which འེ is an abbreviated syllable. When the ergative case word ས follows the word འང (we), it forms འངས (we[+ergative]), in which མས is an abbreviated syllable. Table 2 shows more examples.

---

<sup>1</sup> <http://www.chasen.org/~taku/software/darts/>

word	AM	result	explanation
ང་	མ་	ངམ་	Tsheg is omitted.
གལ་ཆེ་	འི་	གལ་ཆེའི་	Tsheg is omitted.
གོ་	འང་	གོའང་	Tsheg is omitted.
རྣ་བ་	འམ་	རྣ་བའམ་	Tsheg is omitted.
ལྷ་	འམ་	ལྷའམ་	Tsheg is omitted.
དབང་པོ་ལྷ་	འོ་	དབང་པོ་ལྷའོ་	Tsheg is omitted.
སྤྲུལ་མཐའ་	འི་	སྤྲུལ་མཐའི་	འ(/a/) and tsheg are omitted.
ནམ་མཐའ་	འི་	ནམ་མཐའི་	འ(/a/) and tsheg are omitted.
བཤད་པ་	ར་	བཤད་པར་	Tsheg is omitted.
རྒྱལ་པོ་	འི་	རྒྱལ་པོའི་	Tsheg is omitted.

Table 2. Examples of Tibetan abbreviated syllables.

In the rule based word segmentation system, Tibetan words are collected from several dictionaries, namely "Tibetan Chinese General Dictionary" (Zhang 1993), "The Antitheses Chinese-Tibetan Dictionary" (The ethnic publishing house 2002), and "The Antitheses Chinese-Tibetan Oral Dictionary" (Yu 1983), as well as some digital dictionaries. About 220 thousand words are included in the segmenting dictionary. Word frequencies are collected by segmenting a Tibetan text corpus with 230 thousand sentences and 2 million words in total, which includes newspaper articles, law text, political papers and books. In a previous study, experiments show that the precision of the system reaches 96.98% (Liu 2012).

#### 4.1.2 Statistical based method

As Statistical based methods detect out of vocabulary items much more effectively than rule based methods, we also developed a Tibetan word segmentation tool base on the Conditional Random Fields (CRFs) machine learning model. We reformulate Tibetan word segmentation as a syllable tagging problem, and propose an approach using conditional random fields (CRFs) for Tibetan word segmentation.

We convert the segmented words in the corpus into a tagged sequence of Tibetan syllables (or sub-syllables). We tag each syllable with one of the four tags, B (Begin), M (Middle), E (End) and S (Single) depending on its position within a word. Two additional tags, namely ES (End and Single) and SS (Single and Single) are used when we take the syllable rather than the sub-syllable as the tagging unit (Liu 2011; Liu 2015). For each syllable:

- (1) It is tagged B if it is the left boundary of a word.
- (2) It is tagged M if it is at middle of a word.
- (3) It is tagged E if it is the right boundary of a word.
- (4) It is tagged S if it is a word by itself.
- (5) It is tagged ES if it comes from a multiple-syllable word and an AM.
- (6) It is tagged SS if it comes from a single-syllable word and an AM.

Then, for the Tibetan sentence in (a), which means (b), it's segmented into (c) manually. Consequently, it's converted into (d) or (e) by applying the aforementioned tags to form a word

segmentation corpus to be used as the training set for the CRFs. Earlier research shows that this method achieves precisions higher than 94.43% (Liu 2015).

- (a) ང་ཚོས་སྤྱི་ཚོགས་རིང་ལུགས་ཀྱི་སློབ་དབང་བའི་ལམ་ལུགས་དང་ཚོལ་བསྐྱོན་ཐོབ་སྤོང་གི་ཚད་ནུབ་མཐའ་འཁྲུངས་བས་ཡིད།
- (b) We have always followed the principles of socialist public ownership and distribution according to work.
- (c) ང་ཚོ་ས་/སྤྱི་ཚོགས་རིང་ལུགས་/གི་/སློབ་དབང་བའི་ལམ་ལུགས་/དང་/ཚོལ་/བསྐྱོན་/ཐོབ་/སྤོང་/གི་/ཚད་ནུབ་/མཐའ་འཁྲུངས་/བས་/ཡིད་/།
- (d) ང་/B ཚོས་/ES སྤྱི་/B ཚོགས་/M རིང་/M ལུགས་/E ཀྱི་/S སློབ་/B ལ་/M དབང་/M བའི་/M ལམ་/M ལུགས་/E དང་/S ཚོས་/S བསྐྱོན་/S ཐོབ་/S སྤོང་/S ཀྱི་/S ར་/B རྩོམ་/E མཐའ་/B འཁྲུངས་/E བས་/S ཡིད་/S །/S
- (e) ང་/B ཚོ་/E ས་/S སྤྱི་/B ཚོགས་/M རིང་/M ལུགས་/E ཀྱི་/S སློབ་/B ལ་/M དབང་/M བ་/M འི་/M ལམ་/M ལུགས་/E དང་/S ཚོས་/S བསྐྱོན་/S ཐོབ་/S སྤོང་/S ཀྱི་/S ར་/B རྩོམ་/E མཐའ་/B འཁྲུངས་/E བ་/B ས་/E ཡིད་/S །/S

### 4.2 Bilingual sentence level parallel corpus building

Constructing corpora is a basic necessity for Natural Language Processing. For Chinese-Tibetan machine translation, large scale bilingual parallel corpora are still basic resources which are urgently needed. We collected Tibetan text and Chinese text from several translating organizations. A bilingual sentence level parallel corpus was built as part of our project, which includes 571 thousand bilingual sentence pairs. As shown in Table 3 and Table 4, two versions are included. Set A is a long sentence version and Set B is a short sentence version. Each sentence pair in Set A has a complete Chinese sentence which ends with a period, while each sentence pair in Set B has a shorter Chinese sentence which may end with comma. Both of the sets are used in the machine translation system.

Long version	Chinese	西藏自治区面积 122 万平方公里, 平均海拔在 4000 米以上, 有着独特的自然生态和地理环境。
	Tibetan	བོད་རང་སྐྱོང་ལྗོངས་ཀྱི་ས་ཁྲུན་ལ་སྤྱི་ལེ་གྲ་བཞི་མ་ཁྲི་122ཡིད་པ་དང་། རྒྱ་མཚོའི་ངོ་ས་ལས་མཐོ་ཚད་ཆ་སྟོ་མས་སུ་སྤྱི་4000ཡན་བེན་པ་དང་། རང་བྱུང་གི་སྤྱི་ཚོགས་དང་ས་གཤིས་ཀྱི་ཡུལ་ཁམས་དམིགས་བསལ་ཡིན།
Short version	Chinese	西藏自治区面积 122 万平方公里,
	Tibetan	བོད་རང་སྐྱོང་ལྗོངས་ཀྱི་ས་ཁྲུན་ལ་སྤྱི་ལེ་གྲ་བཞི་མ་ཁྲི་122ཡིད་པ་དང་།
	Chinese	平均海拔在 4000 米以上,
	Tibetan	རྒྱ་མཚོའི་ངོ་ས་ལས་མཐོ་ཚད་ཆ་སྟོ་མས་སུ་སྤྱི་4000ཡན་བེན་པ་དང་།
	Chinese	有着独特的自然生态和地理环境。
Tibetan	རང་བྱུང་གི་སྤྱི་ཚོགས་དང་ས་གཤིས་ཀྱི་ཡུལ་ཁམས་དམིགས་བསལ་ཡིན།	

Table 3. Long sentence pair and the corresponding short sentence pairs.

Domain	#Sentence (Set A)	#Sentence (Set B)
Law text	115,299	68,535
Leader's book	53,292	96,181
News	26,613	4,270
Government reports	72,849	102,795
Dictionary	31,234	
Total	299,287	271,781

Table 4. Bilingual sentence level parallel corpus.

The Chinese part of the corpus is processed by some rules and open source tools. Some rules are used to segment Chinese text into sentences. A Chinese word segmentation tool named “ICTCLAS”<sup>2</sup> is used to segment Chinese sentence into words. Tibetan word segmentation methods described in the previous subsection are used to segment Tibetan sentence into words. There are also two problems to solve. The first is how to find the boundary of a Tibetan sentence. The second is how to align Tibetan sentences with Chinese sentences.

#### *4.2.1 Tibetan sentence boundary detection*

Tibetan sentence boundary detection is a nontrivial problem because punctuation marks in Tibetan are not used exclusively to mark sentence breaks, in other words the existence of a punctuation marker does not necessarily suggest the boundary of a sentence. In particular, “། (SHAD)” is one of the most common and significant punctuation marks in Tibetan, which functions like a period, a caesura sign, and a comma.

Tibetan is an ergative language whose structure is SOV (subjective, objective, and verb), so a predicate almost always found at the end of a sentence. Therefore, we choose the predicate to help us disambiguate the end-of-sentence patterns. The diverse constituents of predicates can be divided into four types:

- (1) Verb. For example: ང་དགོ་ཚན་ཡིན།
- (2) Verb with Auxiliary. For example: ངས་རང་བོད་ཡིག་འབྲི་ཤིང་།
- (3) Verb with Aspect-Evidentially Mark. For example: ངས་ལྟོད་རང་གི་ལྷག་ལས་ལ་གཞོན་པ་མ་ཡོང་།
- (4) Verb with Modal Particle. For example: ལྟོད་རང་ལྷ་ས་བ་ཡིན་ནམ།

We observed from the corpus that in a majority of cases the last word of a Tibetan sentence is an auxiliary. Specifically, 78.79% of sentences end with auxiliaries. Considering the importance of auxiliaries, we built an Auxiliary List. A Tibetan verb lexicon is also used in our algorithm. They are used to find the sentence boundaries in a Tibetan article. A previous paper shows that the method is robust and efficient, and its accuracy reaches 99.26% (Zhao 2010).

---

<sup>2</sup> <http://ictclas.nlpir.org/newsdownloads?DocId=389>

#	positive	negative	#	positive	negative
1	ལྷན	མི་ལྷན	16	མོད	མི་མོད
2	མི	མི་མི	17	འདོད	མི་འདོད
3	ཤེས	མི་ཤེས	18	པར་འདོད	པར་མི་འདོད
4	མཉམས	མི་མཉམས	19	བར་འདོད	བར་མི་འདོད
5	འོས	མི་འོས	20	དང	མི་དང
6	བར་འོས	མི་བར་འོས	21	པར་དང	པར་མི་དང
7	བར་འོས	མི་བར་འོས	22	བར་དང	བར་མི་དང
8	བྱང	མི་བྱང	23	བཏུབ	མི་བཏུབ
9	བྱང་དང་མ་བྱང		24	ལྷ་བཏུབ	ལྷ་མི་བཏུབ
10	ཡུང་བྱང		25	ཏུ་བཏུབ	ཏུ་མི་བཏུབ
11	ཚོག	མི་ཚོག	26	ད་བཏུབ	ད་མི་བཏུབ
12	ཉན	མི་ཉན	27	བྱ་བཏུབ	བྱ་མི་བཏུབ
13	ལྷོ	མི་ལྷོ	28	ར་བཏུབ	ར་མི་བཏུབ
14	དགོས	མི་དགོས	29	ལ་བཏུབ	ལ་མི་བཏུབ
15	ལུས	མི་ལུས	30	དཀའ	དཀའ་མིན

Table 5. A part of the Auxiliary List

#### 4.2.2 Bilingual sentence alignment

A Chinese-Tibetan dictionary with 137,873 items was collected by combining several published dictionaries (Liu 2011; Liu 2012; Liu 2015). Bilingual articles are respectively segmented into monolingual sentences. They are further segmented into words. As the correspondences of some words in Tibetan sentence to their Chinese translations in Chinese sentence exist, a dynamic programming algorithm is applied to find the correspondence of the sentences in each pair of Bilingual articles. A previous study shows that the aligning precision of this approach is 84.8% (Yu 2010). We implemented a tool for further proofreading to correct alignment errors (Yu 2012).

#### 4.3 Bilingual word alignment

As Bilingual word pairs are core resources in statistical machine translation, each bilingual sentence pair is segmented into word sequence pairs and aligned. Word alignment finds their correspondences in the target language for the words in the source language. Figure 2 shows a word alignment example. In the system, the alignment is denoted by  $\alpha = \{1 \rightarrow 14; 2 \rightarrow 1; 3 \rightarrow 2; 4 \rightarrow 3; 5 \rightarrow 4; 6 \rightarrow 6, 7; 7 \rightarrow 10; 8 \rightarrow 11; 9 \rightarrow 16\}$

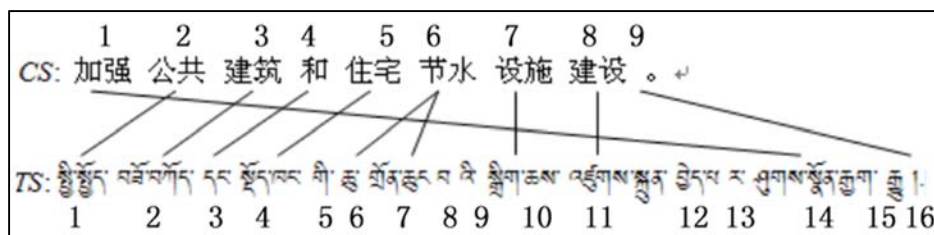


Figure 2. Word alignment result of an example sentence pair.

In the training procedure of statistical machine translation, the word pairs are extracted and the translation probabilities from each word in the source language to each of its correspondences in the target language are fitted, which results in the translation model that will be used in the translation decoding procedure.

We also use the Giza++ toolkit<sup>3</sup> (Och 2003), which is widely used in many machine translation systems, to enact the word alignment.

#### 4.4 Bilingual phrase pair extraction

As mentioned above, we have a Chinese-Tibetan dictionary with 137,873 items. This scale is too small for many NLP purposes. In particular, we need many more word or phrase pairs to build the dictionary based machine translation sub system. As we have a bilingual sentence level parallel corpus, we can extract phrase pairs from it. Such a procedure has two stages.

##### 4.4.1 Stage 1: extracting Chinese phrases

A phrase with more than one word is also called Multi Word Expression (MWE). A collocation measure is used to find the left and right boundaries of a Chinese phrase. For an adjacent pair of words ( $w_1, w_2$ ), the collocation is defined by the following formula (Nuo 2011):

$$Collocation(w_1, w_2) = \frac{VMI(w_1, w_2)}{H(w_1) + H(w_2)}$$

where  $w_1$  and  $w_2$  represent the occurrence of two words,  $H(w_1)$  is the entropy of the word  $w_1$ , and  $VMI(w_1, w_2)$  is the average mutual information of the two words (Nuo 2011), defined as follows:

$$VMI(w_1, w_2) = P(w_1, w_2) \log \frac{P(w_1, w_2)}{P(w_1) * P(w_2)} + P(\overline{w_1}, \overline{w_2}) \log \frac{P(\overline{w_1}, \overline{w_2})}{P(\overline{w_1}) * P(\overline{w_2})} \\ + P(\overline{w_1}, w_2) \log \frac{P(\overline{w_1}, w_2)}{P(\overline{w_1}) * P(w_2)} + P(w_1, \overline{w_2}) \log \frac{P(w_1, \overline{w_2})}{P(w_1) * P(\overline{w_2})}$$

where  $P(w_1)$  is the occurrence probability of the word  $w_1$ ,  $P(w_1, w_2)$  is the occurrence probability of the adjacent pair  $(w_1, w_2)$ , and  $P(w_1, \overline{w_2})$  is the occurrence probability of adjacent pairs starting with  $w_1$  but followed by any word other than  $w_2$ .

For the word sequence of  $w_1 w_2 w_3$ , if we denote  $x = Collocation(w_1, w_2)$  and  $y = Collocation(w_2, w_3)$ , then, the BindingDegree( $x, y$ ) is defined as follow (Nuo 2011):

$$BindingDegree(x, y) = \begin{cases} x/y, & \text{if } y \geq x \\ y/x, & \text{if } y < x \end{cases}$$

<sup>3</sup> <http://www.statmt.org/moses/giza/GIZA++.html>

So the  $\text{Collocation}(w_1, w_2)$  tells whether  $w_1 w_2$  forms a bi-word phrase, while the  $\text{BindingDegree}(x, y)$  tells whether  $w_3$  can be appended to the bi-word phrase  $w_1 w_2$  to form a tri-word phrase  $w_1 w_2 w_3$ .

#### *4.4.2 Stage 2: extracting the Tibetan correspondences for Chinese phrase*

Generally, if a phrase occurs in the source language in every pair in a certain set of bilingual sentence pairs, its translation into the target language occurs also in the same set. When a Chinese phrase is extracted, we extract all the sentence pairs which it occurs in to form a candidate sentence pair set A. Then, words occurring in every Tibetan sentences in set A are extracted, which forms the Tibetan translation of the Chinese phrase.

### **4.5 Example based machine translation**

When a Chinese sentence is to be translated, it is segmented into word sequences. If there is a similar sentence (or an identical sentence) in the bilingual sentence level parallel corpus, the Tibetan counterpart of the sentence pair can be taken as the translation of the Chinese sentence. This is the advantage of example based machine translation.

In the proposed system, for a Chinese sentence, all sentence pairs which have any word in common with in the Chinese sentence to be translated are extracted as the candidate set B. Then the Levenshtein Distances (Hirschberg 1975) (LD, word as the unit) between the Chinese sentence and all those in set B are computed. The similarity of two Chinese sentences CS\_A and CS\_B as word sequences is defined as follows:

$$\text{Sim}(\text{CS}_A, \text{CS}_B) = 1 - \frac{\text{LD}(\text{CS}_A, \text{CS}_B)}{\max(\text{len}(\text{CS}_A), \text{len}(\text{CS}_B))}$$

The sentence with the maximum similarity to the input sentence is selected. If the similarity is larger than the predefined threshold (0.7 in our system, as determined by a test on some typical sentences), it is selected and the Tibetan counterpart in the sentence pair is taken as the translation of the input sentence.

#### 4.6 Statistical machine translation

phrase	Candidates	Phrase	candidates
我们	ང་ཚོ་ལ་	民族	མི་རིགས་
我们	ང་ཚོ་	民族	མི་རིགས་ཁག་
我们要	ང་ཚོ་ལ་	民族地区	མི་རིགས་ས་ཁུལ་གྱི་
我们要	ང་ཚོ་ལ་ བྱ་དགོས་	民族地区	མི་རིགས་ས་ཁུལ་ དྲ་
我们要不断	ང་ཚོ་ལ་ རྒྱུ་ཚད་མེད་པ་ ར་	民族地区的	མི་རིགས་ས་ཁུལ་གྱི་
我们要不断	ང་ཚོ་ལ་ རྒྱུ་མི་འཚད་པ་ ར་	民族地区的	མི་རིགས་ས་ཁག་གི་
要	བྱ་དགོས་	地区	ས་ཁུལ་
要	དགོས་	地区	ས་ཁུལ་ཁག་ དང་
不断加强	ཁྱེད་རྒྱ་ ར་ རྒྱལ་སྤོན་ ཟམ་ མི་ ཚད་པ་ བརྒྱབ་ བ་	干部	ལས་བྱེད་པ་
不断加强	ར་ རྒྱལ་སྤོན་ ཟམ་ མི་ ཚད་པ་ བརྒྱབ་ བ་	干部	འགོ་ཁྲིད་ལས་བྱེད་པ་
加强	རྒྱལ་སྤོན་	干部队伍	ལས་བྱེད་པ་ འི་ དབྱེད་ཁག་
加强	ཁྱེད་རྒྱལ་ས་	干部队伍	ལས་བྱེད་པ་ འི་ དབྱེད་ཇེ་
加强民族	མི་རིགས་མཐུན་སྦྱིལ་ ལ་ རྒྱལ་སྤོན་	干部队伍建设	ལས་བྱེད་པ་ འི་ དབྱེད་ཇེ་ འཛུགས་སྐྱོང་ ཁྱེད་པ་
加强民族	མི་རིགས་མཐུན་སྦྱིལ་ ལ་ རྒྱལ་སྤོན་རྒྱལ་ བ་	干部队伍建设	ལས་བྱེད་པ་ འི་ དབྱེད་ཁག་ ལེགས་སྐྱོང་

Table 6. Each phrase in the sentence has multiple translation candidates.

After word alignment, we can train a SMT model. The probability of any translation candidate for a phrase in the source language is stored in the model. The SMT decoder will make the decision to select the overall best candidate for each phrase in a sentence to be translated. We take an example to explain the procedure.

Sentence to be translated : 我们要不断加强民族地区的干部队伍建设。

English translation: We have to keep strengthening the construction of cadres in ethnic areas.

The decoder will find multiple candidates for nearly every phrase in the sentence. Table 6 shows some candidates for some phrases in the sentence to be translated.

Let's denote the Chinese sentence by  $f$ , and the Tibetan sentence by  $e$ , and denote the  $i$ th phrase in  $f$  by  $f_i$ , and the  $i$ th phrase in  $e$  by  $e_i$ . The probability of the Tibetan sentence is the translation of the Chinese sentence can be measured by the following formula:

$$P(e|f) = \exp \left[ \lambda_\phi \sum_{i=1}^I \log \phi(f_i|e_i) + \lambda_{LM} \sum_{i=1}^{|e|} \log p_{LM}(e_i | e_1 \cdots e_{i-1}) \right]$$

Note the  $\phi(f_i|e_i)$  is the translation probability to translate a Tibetan phrase to a Chinese phrase, and  $P_{LM}(e_i | e_1 \cdots e_{i-1})$  is the probability of how likely the Tibetan phrases forms a Tibetan sentence, which is computed by the language modeling technology (Ponte 1998).  $\lambda_\phi$  and  $\lambda_{LM}$  are the weights of the translation model and the language model respectively in the computation, and they are determined by training on the bilingual parallel corpus.



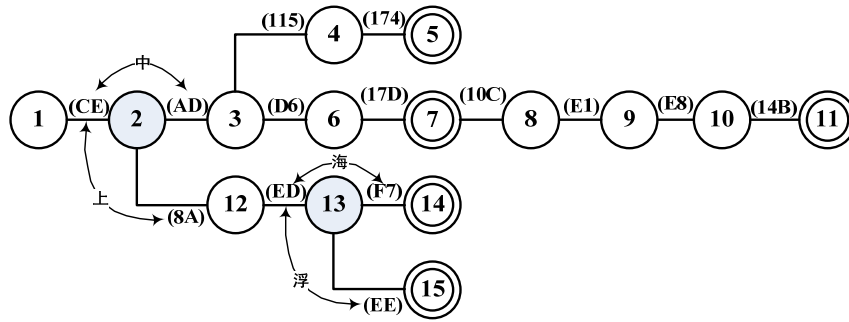


Figure 5. The new trie of key set  $K$

In each terminal node, a pointer to the Tibetan translations of the Chinese word is indexed, which will be used to extract the corresponding translation. With this structure, the forward maximum matching method can be performed to translate a Chinese sentence to Tibetan. Backward maximum matching method can be performed similarly (Aoe 1989).

## 5 System Evaluation

In this section the performances of the different sub-systems are compared, and the well-known Moses<sup>4</sup> statistical machine translation system is also included in the comparison. 69,756 sentence pairs are included in the experiments, in which 429 pairs are randomly selected as the test set. The other pairs are included in the training set. The domain distributions of the training set and test set are listed in Table 7.

In the proposed system, two SMT sub systems are implemented, the first is maximum probability decoding model (MaxProb). The well-known Moses system is also taken as the second SMT sub system. Two RBMT sub systems are implemented too, which use Forward Maximum Matching (FMM) decoding and Backward Maximum Matching (BMM) decoding respectively.

The performances of different MT sub-systems are listed in Table 8. The performance data are collected on a PowerEdge R710 server with 2 4-core Intel E5620 (2.40 GHz) CPUs and 16 GB of memory. All sub systems (including Moses) are running on the same server. The two evaluation metrics BLEU4 and NIST are widely used in machine translation research (Papineni 2002; George 2002). BLEU4 simply calculates the geometric average of n-gram precisions for n=1 to 4 adding equal weight to each one, while NIST also calculates how informative a particular n-gram is. Obviously, Moses gives the best translations comparing with all sub-systems of the proposed system. However, it takes nearly 10 times more time. The EBMT sub-system gives the fewest correct translations. The other three sub systems have similar performance in both translation quality and time cost. However, 25,512 phrases are extracted by the SMT(MaxProb) sub system, in which 2258 long phrases are not extracted by Moses with default configuration. So at the phrase level, the recall of the SMT(MaxProb) sub system improves the results by 9.71% over Moses.

<sup>4</sup> <http://www.statmt.org/moses/>

	Domain	# of sentence pair
Training set	Law text	8,595
	Leader's book	25,112
	Government reports	34,610
Total of training set.		67,327
Test set	Law text	63
	Leader's book	151
	Government reports	215
Total of test set.		429
Total		69,756

Table 7. Domain distributions of the training set and test.

Table 9 shows the output of different sub-systems for the sentence mentioned in the former section. It shows that the EBMT sub system successfully find the sentence itself in the RBMT model because it's included in the training set. The other three sub-systems find translation candidates for every phrase in the sentence, but they are arranged in the same order as they are in the Chinese sentence rather than in a Tibetan sentence. That is why these sub-systems have a worse translation quality than Moses. So, an order adjusting model is essential to a Chinese to Tibetan translation system.

MT (sub) system	BLEU4	NIST	Time(s)
EBMT	0.0410	3.1472	41
SMT(MaxProb)	0.2168	5.6882	75
SMT(Moses)	0.2771	6.2338	633
RBMT(FMM)	0.2040	5.6235	68
RBMT(BMM)	0.1986	5.5944	67

Table 8. Performances of different MT sub systems.

Source	我们(we) 要(have to) 不断(keep) 加强(strengthening) 民族(ethnic) 地区(are) 的('s) 干部(cadres) 队伍(group) 建设(construction) 。
English	We have to keep strengthening the construction of cadres group in ethnic areas.
Reference	ང་ཚོ་ལ་ མི་རིགས་ ས་ཁུལ་ གྱི་ ལས་བྱེད་པ་ འི་ དཔུང་ལག་ འཇུགས་སྐྱོང་བྱེད་ཤགས་ ཆེ་ཏུ་ ཟམ་མི་ཚད་པ་ གཏོང་དགོས།
EBMT	ང་ཚོ་ལ་ མི་རིགས་ ས་ཁུལ་ གྱི་ ལས་བྱེད་པ་ འི་ དཔུང་ལག་ འཇུགས་སྐྱོང་བྱེད་ཤགས་ ཆེ་ཏུ་ ཟམ་མི་ཚད་པ་ གཏོང་དགོས།
SMT	ང་ཚོ་ལ་ བྱེད་རྒྱ་ར་ ཤགས་སྐྱོན་ ཟམ་མི་ཚད་པ་ བརྒྱབ་པ་ མི་རིགས་ ས་ཁུལ་ གྱི་ ལས་བྱེད་པ་ འི་ དཔུང་ལྗེ་ འཇུགས་སྐྱོང་ བྱེད་པ།
RBMT(FMM)	ང་ཚོ་ལ་ རྒྱན་ཚད་མེད་པ་ ར་ མི་རིགས་ ས་ཁུལ་ གྱི་ བྱེད་ཤགས་ ཆེ་ཏུ་གཏོང་ དགོས་ ལས་བྱེད་པ་ འི་ དཔུང་ལྗེ་ ལེགས་ སྐྱོང་ བྱ་རྒྱུ་ འོ།
RBMT(BMM)	ང་ཚོ་ལ་ ལུ་མཐུན་ ཤགས་སྐྱོན་རྒྱག་ དགོས་ མི་རིགས་ གངས་ལུང་ ས་ཁུལ་ གྱི་ ལས་བྱེད་པ་ འི་ དཔུང་ལག་ ལེགས་ སྐྱོང་ འོ།

Table 9. A comparison of sub systems' outputs.

## 6 Conclusion

There is a big gap between the capacity for Chinese to Tibetan professional human translation and the market requirement. Machine translation technology is an urgent requirement, which contributes to the improvement of translation speed as well as the quality. We made great effort to build a Chinese to Tibetan machine translation system and discuss the key corpora and technologies in the paper. Experiments show the sub systems output the translation of each phrase in the same order as they are in the Chinese sentence rather than in a Tibetan sentence, which leads to a low quality translation. So an order adjusting model is essential to a Chinese to Tibetan translation system. As the translation quality is not good enough, we will make efforts toward order adjusting in the future. A computer aided tool is another requirement to make full use of the output of the system, and generates a better translation by interacting with the translator.

## 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 National Science Foundation (No.61202219, No.61202220, and No.61303165) and Informationization Project of the Chinese Academy of Sciences (No. XXH12504-1-10).

## REFERENCES

- Aoe, Junichi. 1989. "An efficient digital search algorithm by using a double-array structure". *IEEE Transactions on Software Engineering* 15.9, 1066-1077.
- Cai, Rangjia. 2011. "Tibetan corpus processing method". *Computer Engineering and Application* 47.6: 138-139, 146. (才让加.藏语语料库加工方法研究[J].计算机工程与应用,2011,47(6):138-139,146.)
- Cai, Zangtai; and Hua, Guanjia. 2005. "Research of Banzhida Chinese-Tibetan document translation system based on the dichotomy of syntax analysis". *Journal of Chinese Information Processing* 19.6: 9-14. (才藏太,华关加.班智达汉藏公文翻译系统中基于二分法的句法分析方法研究 [J].中文信息学报,2005, 19 (6): 9-14.)

- Chang, Pi-Chuan; Galley, Michel; and Manning, Chris. 2008. "Optimizing Chinese word segmentation for machine translation performance". Paper presented at ACL 3th Workshop on Statistical Machine Translation.
- Chen, Yuzhong; and Yu, Shiwen. 2003. "The current status and future of Tibetan information processing technologies". *China Tibetology* 4: 97-107. (陈玉忠,俞士汶. 藏文信息处理技术的研究现状与展望[J]. 中国藏学, 2003, (4): 97-107.)
- Degai, Cailang; Li, Yanfu; and Hangqing Chaojia, et al. 2001. "The design and implementation of a practical Chinese-Tibetan machine translation system". In: *Proceedings of the 863 high-tech project "Intelligent Computer Workshop"*, 405-411. (德盖才郎,李延福,项清朝加等. 实用化汉藏机器翻译系统的设计与实现[C].//863 计划智能计算机主题学术会议论文集.2001:405-411.)
- Doddington, George. 2002. "Automatic evaluation of machine translation quality using n-gram Co-occurrence Statistics". In: *Proceedings of the Second International Conference on Human Language Technology Research*, 138-145. San Francisco: Morgan Kaufmann Publishers Inc.
- Fredkin, Edward. 1960. "Trie memory". *Communications of the ACM*, 3.9: 490-500.
- Hirschberg, D. S. 1975. "A linear space algorithm for computing maximal common subsequences". *Communications of the ACM*. 18.6: 341-343.
- Hou, Hongxu; Liu, Qun; and Nasun Urt. 2007. "Example based Chinese-Mongolian machine translation". *Journal of Chinese Information Processing* 21.4: 65-72. (侯宏旭,刘群,那顺乌日图等. 基于实例的汉蒙机器翻译[J]. 中文信息学报, 2007, 21(4): 65-72.)
- Kang, Byeong-kwu; Zhang, Qinlong; Chen, Yirong; and Chang, Bao-bao. 2007. "Chinese multi-word chunks extraction for computer aided translation". *Journal of Chinese Information Processing* 21.1: 9-16. (姜柄圭,张秦龙,谌贻荣等. 面向机器辅助翻译的汉语语块自动抽取研究[J]. 中文信息学报, 2007, 21 (1): 9-16.)
- Kanzhuo Caidan; Jin, Weixun; and Luo, Zhihua et al. 2006. "To search and solve problems about verbs in Chinese-Tibetan translation system". *Terminology Standardization and Information Technology* 3. (看卓才旦,金为勋,李延福等. 汉藏翻译系统中的动词处理研究 [J]. 术语标准化与信息技术, 2006, (3).)
- Koehn, Philip; Hoang, Hieu; Birch, Alexandra; Callison-Burch, Chris; Federico, Marcello; Bertoldi, Nicola; Cowan, Brooke; Shen, Wade; Moran, Christine; Zens, Richard; Dyer, Chris; Bojar, Ondrej; Costantin, Alexandra; and Herbst, Evan. 2007. "Moses: Open source toolkit for statistical machine translation". In: *Proceedings of the ACL 2007 (demo and poster sessions)* Prague, June 2007, 177-180.
- Liu, Huidan; Long, Congjun; Nuo, Minghua; and Wu, Jian. 2015. "Tibetan word segmentation as sub-syllable tagging with syllable's Part-Of-Speech property". In: Maosong Sun; Zhiyuan Liu; Min Zhang; and Yang Liu (eds.), *Chinese computational linguistics and natural language processing based on naturally annotated big data* (LNAI 9427), 189-201.
- Liu, Huidan; Nuo, Minghua; Ma, Longlong; Wu, Jian; and He, Yeping. 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)*, 168-177.
- Liu, Huidan, Minghua Nuo, Weina Zhao, Longlong Ma. 2012. "SegT: A pragmatic Tibetan word segmentation system". *Journal of Chinese Information Processing* 26.1: 97-103. (刘汇丹, 诺明花, 赵维纳, 吴健, 贺也平. SegT: 一个实用的藏文分词系统[J]. 中文信息学报, 2012, 26(1): 97-103.)

- Liu, Huidan; Zhao, Weina; Nuo, Minghua; Jiang, Li; Wu, Jian; and He, Yeping. 2010. "Tibetan number identification based on classification of number components in Tibetan word segmentation". Paper presented at International Conference on Computational Linguistics 2010, pp. 719-724.
- Luo, Aijun; Gelang; Wujin Jiacan; and Ling, Yongxiao. 2010. "Investigation and analysis on the Chinese to Tibetan translator group in Tibet". *Tibet Technology* 5: 21-13. (罗爱军, 格朗, 伍金加参等. 西藏汉藏翻译队伍状况调查与分析 [J]. 西藏科技, 2010.)
- Nuo, Minghua; Liu, Huidan; and Ma, Longlong et al. 2011. "Automatic acquisition of Chinese-Tibetan multi-word equivalent pair from bilingual corpora". In: *Proceedings of the International Conference on Asian Language Processing (IALP2011)*, 177-180.
- Nuo, Minghua; Wu, Jian; Liu, Huidan; and Ding, Zhiming. 2011. "Research on phrase translation extraction methods in the process of Chinese-Tibetan phrase pair extraction". *Journal of Chinese Information Processing* 25.3: 112-117. (诺明花, 吴健, 刘汇丹等. 汉藏短语对抽取中短语译文获取方法研究[J]. 中文信息学报, 2011, 25(3): 112-117.)
- Och, Franz Josef; and Ney, Hermann. 2003. "A systematic comparison of various statistical alignment models". *Computational Linguistics* 29.1: 19-51
- Ponte, Jay M.; and Croft, W. Bruce. 1998. "A Language modeling approach to information retrieval". *SIGIR '98 Proceedings of the 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, 275-281. New York: ACM.
- Papineni, Kishore; Roukos, Salim; Ward, Todd; and Zhu, Weijing. 2002. "BLEU: A method for automatic evaluation of machine translation". In: *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, 311-318.
- The Ethnic Publishing House. 2002. *The antitheses Chinese-Tibetan dictionary*. Beijing: The ethnic publishing house. 2002.7. (民族出版社编, 汉藏对照词典. 民族出版社. 2002.7)
- Yu, Daoquan. 1983. *The antitheses Chinese-Tibetan oral dictionary*. Beijing: The Ethnic Publishing House. (于道泉, 汉藏对照词典. 民族出版社. 1983.10)
- Yu, Xin; Zhang, Liqiang; and Liu, Huidan. 2012. "Mechanism for Large-scale Chinese-Tibetan Bilingual Corpus Construction". In: *Proceedings of the 4th China National Youth Conference on Minority Language Information Processing*, 37-42. (于新, 张立强, 刘汇丹等. 大规模汉藏双语语料库构建机制[C]. //第四届全国少数民族青年自然语言信息处理论文集. 2012: 37-42.)
- Yu, Xin; Zhao, Weina; and Wu, Jian. 2010. "Dictionary-based Chinese-Tibetan sentence alignment". In: *Proceedings of the 2010 IEEE International Conference on Intelligent Computing and Integrated Systems*.
- Zhao, Weina; Liu, Huidan; Yu, Xin; Wu, Jian; and Zhang, Pu. 2010. "The construction of the Chinese-Tibetan bilingual parallel corpora for the Chinese-Tibetan bilingual computer-aided translation systems". In: *Proceedings of the 3rd China National Youth Conference on Minority Language Information Processing*, 43-46. (赵维纳, 刘汇丹等. 面向汉藏辅助翻译系统的平行语料库建设. 第三届全国少数民族青年自然语言信息处理学术研讨会, 2010.)
- Zhao, Weina; Yu, Xin; and Liu, Huidan. 2010. "Sentence boundary detection based on auxiliary verbs in Modern Tibetan". Paper presented at the Conference on Language Investigation and Information Processing (LIIP2010).

Zhang, Huaping; Liu, Qun; Cheng, Xueqi; Zhang, Hao; Yu, Hongkui. 2003. "Chinese lexical analysis using Hierarchical Hidden Markov Model". Paper presented at the Second SIGHAN Workshop Affiliated with 41st ACL; Sapporo Japan, July, pp. 63-70.

Zhaluo; Suonan Renqian. 2006. "Research on the translation rules for complex sentences in Chinese-Tibetan machine translation". In: *Chinese information processing Frontiers - Proceedings of the 25th Anniversary Conference of the Chinese Information Processing Society*. (扎洛,索南仁欠. 汉藏机器翻译中复句的翻译规则研究 [A]. 中文信息处理前沿进展——中国中文信息学会二十五周年学术会议[C]. 2006.)

Zhang, Yisun. 1993. *Tibetan Chinese general dictionary*. Beijing: The Ethnic Publishing House. 1993.12. (张怡荪编, 藏汉大辞典. 民族出版社.1993.12)

Huidan Liu  
huidan@iscas.ac.cn

Weina Zhao  
zhaoweina1999@yahoo.com.cn

Minghua Nuo  
csnmh@imu.edu.cn

Jinling Hong  
jinling@iscas.ac.cn

Xin Yu  
yuxin08@iscas.ac.cn

Jian Wu  
wujian@iscas.ac.cn