

Free Linguistic and Speech Resources for Tibetan

Guanyu Li^{*}, Hongzhi Yu^{*}, Thomas Fang Zheng[†], Jinghao Yan^{*}, Shipeng Xu^{*}

^{*}Key Laboratory of National language Intelligent Processing Gansu Province, Northwest Minzu University, Lanzhou, China
Email: guanyu-li@163.com, yuhongzhi@hotmail.com, yjh527a@163.com, xspxkrs@aliyun.com Tel: +86-13809316272

[†]Center for Speech and Language Technologies, Tsinghua University, Beijing, China
E-mail: fzheng@tsinghua.edu.cn

Abstract—Tibetan is an important low-resource language in China. A key factor that hinders the speech and language research for Tibetan is the lack of resources, particularly free ones. This paper describes our recent progression on Tibetan resource construction supported by the NSFC M2ASR project, including the phone set, lexicon, as well as the transcription of a large scale speech corpus. Following the M2ASR free data program, all the resources are publicly available and free for researchers. We also release a small Tibetan speech database that can be used to build a proto type Tibetan speech recognition system.

I. INTRODUCTION

Tibetan language is a key member in the family of minor languages in China. It belongs to the Sino-Tibetan language family, the Tibeto-Burman subgroup. The speakers are about 6 million people, mainly distributed in China (Tibet, Qinghai, Gansu, Sichuan and Yunnan provinces), India, Bhutan, and Nepal. Compared to the major languages such as English and Mandarin, the research for Tibetan is far from extensive, on both linguistics and speech processing. A key factor that hinders the research is that the resources are very limited and far from being standard. For example, the lexicon is still in a small scale, and large-scale speech databases are very rare. Most seriously, most of the resources are held by individual institutes, with very limited sharing and openness.

The Multilingual Minorlingual Automatic Speech Recognition (M2ASR) project aims to change the situation. An ambition of this project is to construct a full set of language and speech resources for 5 minor languages (Tibetan, Mongolia, Uyghur, Kazak and Kirgiz), and make the resources open and free for research purposes. In this paper, we report our progress on Tibetan resource construction, including the phone set, the lexicon, transcriptions and speech databases. All the resources are available on the project webpage (<http://m2asr.cslt.org>), and can be obtained by either free download or delivery on request.

Note that there are 3 Tibetan dialect areas in China: U-Tsang, Amdo, Kham. People in the three areas use the same written form, but pronounce very differently. In U-Tsang, the most popular dialect is the Lhasa Tibetan, and in Amdo, the Xiahe Tibetan (or Labrang Tibetan) is the mostly

influential. The M2ASR project focuses on the two dialects as they are spoken by most of Tibetan people. In the following sections, we will first briefly summarize the written and pronunciation system of Tibetan, and then propose our work on resource construction for the two dialects respectively.

II. CHARACTERS OF TIBETAN

Tibetan scripts are written in alphabets. From view of written form, there are 30 consonant letters and 4 vowel signs in Tibetan (note all dialects are the same in writing). Each syllable is a combination of several consonant letters and a vowel sign. Words are comprised of one or several syllables. In the Tibetan script, syllables and words are written from left to right, and are separated by the same delimiter “.” (called ཚེག་ (/tsheg/) in Tibetan).

Each syllable involves a radical consonant letter, and other consonant letters could be appended to the radical consonant as superscript, subscript, prescript, postscript and post-postscript to form a syllable (Fig 1). A syllable must contain a vowel sign, but a vowel sign corresponds to a sound /a/ can be omitted. In general, the vowel signs ཨ, ཨ, ཨ, ཨ sound /i/, /u/, /e/, /o/ respectively, but exceptions also exist, as their pronunciations can be changed following some regular rules. Note that in all the dialects of Tibetan, two syllables may be pronounced the same but each syllable has only a single pronunciation. In other words, there are many homophones but no polyphones in Tibetan.

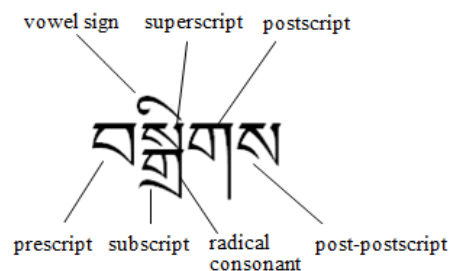


Fig. 1 Constitution of syllable.

The radical consonant, the prescript and superscript consonants together form the initial part of a syllable, and

the vowel sign, the postscript and post-postscript consonants altogether form the final part. In ancient Tibetan, there are many consonant compositions. These consonant compositions are largely preserved in the modern Xiahe dialect, however in the Lhasa dialect, most consonant compositions sound just like a single constant. Another distinction between the Lhasa dialect and the Xiahe dialect is that the former is tonal (four tones in total: 43, 44, 12 and 113[1]) while the latter is toneless. For these reasons, the two dialects sound very different and should be treated different in resource construction.

III. RESOURCES FOR LHASA TIBETAN

In this section, we describe our work on resources construction for the Lhasa dialect. The resources include the phones set, the syllable lexicon, the word lexicon, text database and speech database.

A. Phone set

The phone set involves small and distinct pronunciation units. These units are related to the consonants and vowels in the written form, and more reflect the true pronunciation. We follow the seminal work by Ge Sang Ju Mian [1] and define 29 consonants and 8 cardinal vowels in Lhasa Tibetan. The phone clusters of the consonants are presented in Table I, where the consonant /f/ only appears in foreign words. The vowels in Lhasa Tibetan are listed in table II [1][2][3]. There is a long vowel form for each of the 8 cardinal vowels; 5 vowels (/ɛ/, /e/, /ø/, /i/ and /y/) have a glottalized form and 3 vowels (/e/, /i/ and /y/) have a nasalized form. Additionally, there are 4 consonants (/k/, /m/ and /p/) that can be augmented to the end of a vowel to form a coda, and there are two compound vowels: /au/ and /iu/.

TABLE I
CLUSTERS OF LHASA TIBETAN CONSONANTS

		bilabial		labiodental		Apical alveolar		retroflex		Palatal		Velar		Labial-velar		glottal	
Plosive	voiceless	unaspirated	p		t		c	K									ʔ
		aspirated	p ^h		t ^h		c ^h	k ^h									
Affricate	voiceless	unaspirated		f	ts	tʂ	tc										
		aspirated			tʂ ^h	tʂ ^h	tc ^h										
Fricative		unaspirated			s	ʂ	ɕ										h
nasal	voiced	unaspirated	m		n		ɲ	ŋ									
approximant	voiced	unaspirated						j					w				
Lateral fricative						l											
Lateral approximant	voiced	unaspirated				l											
trill	voiced	unaspirated				r											

TABLE II
VOWELS OF LHASA TIBETAN

characters			vowels				Post consonant					
tongue position	rounded	long	glottalized	nasalized			k	m	p	u	ŋ	
												front
front	low											

front	medium low		ɛ	ɛ:	ɛʔ												
front	medium high		e	e:	eʔ	ẽ			em	ep							ej
front	medium high	rounded	ø	ø:	øʔ												
front	high		i	i:	iʔ	ĩ		ik	im	ip							ij
front	high	rounded	y	y:	yʔ	ỹ											
back	medium	rounded	o	o:					ok	om	op						oj
back	high	rounded	u	u:					uk	um	up						uj

For speech recognition, we made a slight modification to construct the ASR phone set. The first modification is that we merge a cardinal vowel with its corresponding long vowel form. The second modification is to treat /au/ and /iu/ as two single phones rather than splitting them into ingredient phones. To ease the text processing with computers, all these phones (IPAs) are transformed into Latin letter, as shown in Table III.

TABLE III
PHONES AND LATIN TRANSFORMATION OF LHASA TIBETAN

IPA	Latin	IPA	Latin	IPA	Latin	IPA	Latin	IPA	Latin	IPA	Latin	IPA	Latin	IPA	Latin
c	c	l	l	s	s	ŋ	ng	te	tx	ɛ	Ec	øʔ	edb	ỹ	yy
c ^h	ch	m	m	t	t	ɲ	nn	te ^h	tx ^h	ɛʔ	Ecb	i	i	o	o
h	h	n	n	t ^h	th	ɕ	x	ʔ	ab	e	E	iʔ	ib	u	u
j	j	p	p	tʂ	ts	ʂ	ss	f	f	eʔ	Eb	ĩ	ii	au	au
k	k	p ^h	ph	tʂ ^h	tsh	tʂ	ds	l	lh	ẽ	Ee	y	y	iu	iu
k ^h	kh	r	r	w	w	tʂ ^h	dsh	a	a	ø	Ed	yʔ	yb		

B. Syllable lexicon

A syllable lexicon involves a set of syllables whose pronunciations are defined. There are more than 8000 possible syllables in Tibetan, including syllables for foreign words. We construct the syllable lexicon by constructing a text corpus involving 420,000 sentences (including both written and spoken), and then selected the most frequent syllables from this corpus. After removing some syllables that are for transliterating Sanskrit words only, we obtained a syllable lexicon consisting of 6013 syllables. By applying the pronunciation rules, these syllables were segmented into initials and finals, and the initials and finals were further split into phones [4]. All these syllables and their phone sequence forms were manually checked to ensure the quality.

C. Word lexicon

The word lexicon translates words into syllable sequences. To construct the lexicon, the same text corpus used in the syllable lexicon construction is used to form the word list. This is performed by a word segmentation followed by a frequency-based filtering. By this approach, we obtained 27000 frequent words. This primary set was further extended by adding two extra sets: a set of 10000 nouns and

Download English Version:

<https://daneshyari.com/en/article/4997052>

Download Persian Version:

<https://daneshyari.com/article/4997052>

[Daneshyari.com](https://daneshyari.com)