

Rule-Based Pāli Romanization System for Myanmar Language

Lei Lei Win

* Faculty of Computer Science
** University of Computer Studies, Meiktila

DOI: 10.29322/IJSRP.8.9.2018.p8158
<http://dx.doi.org/10.29322/IJSRP.8.9.2018.p8158>

Abstract- Typically, Myanmar is the most religious Buddhist country with regard to the percentage of the population living as monks and the amount of money spent on religion. Pāli is a language that has been widely used in the Buddhist scriptures. Generally, the Pāli words are expressed with stacked consonant so that there may be some difficulties to pronounce the Myanmar Pāli word. Therefore, this paper presents the text to speech system for Myanmar Pāli word by using Romanization rules. Firstly, the input words or sentences in Pāli are accepted. Then, these words are checked as it is Pāli words or not by using rules for Pāli. After that, these Pāli words are converted into their corresponding roman symbols by using Romanization rules for Pāli. Finally, the system generates the speech of Pāli words. The aim of this paper is to help the Buddhists and the new Buddhist students who are unfamiliar with some of the Pāli words offer used in the study of Buddhism. According the experimental result, the system achieved the acceptable level of accuracy.

Index Terms- Pāli language, Text-to-speech, Romanization rules

I. INTRODUCTION

Text to speech system is the conversion of the input text into their corresponding output speech. Myanmar language is the official language and it is tonal, pitch-register and syllable time language. Moreover, the nature of Myanmar language is monosyllabic and analytic language. The sentence order is “subject-object-verb”. It also known as Burmese and it is a member of the Lolo-Burmese grouping of the Sino-Tibetan language family. The language uses a Brahmic script called the Burmese script. The Burmese alphabet consists of 33 letters and 12 vowels, and is written from left to right. It requires no spaces between words, although modern writing usually contains spaces after each clause to enhance readability. Besides the Burmese language, the Burmese alphabet is also used for the liturgical languages of Pāli and Sanskrit. Pāli is the language used to preserve the Buddhist canon of the Theravada (ထေရဝါဒ) Buddhist tradition (is a branch of Buddhism that uses the teaching of the Pāli Canon, a collection of the oldest recorded Buddhist texts, as its doctrinal core), which is regarded as the oldest complete collection of Buddhist texts surviving in an Indian language. Pāli is closely related to Sanskrit, but its grammar and structure are simpler. Traditional Theravadins regard Pāli as the language spoken by the Buddha himself, but in

the opinion of leading linguistic scholars, Pāli was probably a synthetic language created from several vernaculars to make the Buddhist texts comprehensible to Buddhist monks living in different parts of northern India. As Theravada Buddhism spread to other parts of southern Asia, the use of Pāli as the language of the texts spread along with it, and thus Pāli became a sacred language in Sri Lanka, Myanmar, Thailand, Laos, Cambodia and Vietnam. Pāli has been used almost exclusively for Buddhist teachings, although many religious and literary works related to Buddhism were written in Pāli at a time. So Pāli is a spoken language, written in the script of the land where it is used: for example, in Myanmar, it is written in Myanmar script. Strictly speaking Myanmar and Indic scripts are abugidas (alpha-syllabaries) and not alphabets [1].

Myanmar language has been greatly influenced by the Pāli language due to the widespread practice of Buddhism and the study of Buddhist scriptures in Myanmar. [2] expressed that the essentially Indian genius, the psychological subtleties, and high thoughts of Buddhism have forced the Burmese language to grow, deepen and expand continually. As a consequence of Pāli influence on Myanmar language, usages of Pāli and Pāli derived words are wide and frequent in Myanmar text. Some Pāli words were directly incorporated into Myanmar language. Later era, the re-researchers focus on speech processing. Therefore, this paper presents text to speech system for Myanmar Pāli words by using Romanization rules to help the people who are interested in Pāli that is writing in roman script. It can also help the students and teachers of oriental studies.

[3] proposed the system to use for geographical names in the democratic people's republic of Korea. Its contains the Romanization of Hangeul. The way to write Thai language using roman alphabets are proposed by [4]. It could be performed on the basis of orthographic form (transliteration) or pronunciation (transcription) or both. Romanization system for Japanese Kana is proposed by [5]. It has been in use by the U.S. Board on Geographic names and the U.K. [6] is the system of the Romanization of Shan. It has been developed for use in Romanization names written in the Shan script.

The rest of paper is organized as follow: in Section II, Pali language is explained. The nature of Pali word in Myanmar language is presented in Section III. The architecture of the proposed system is explained in Section IV with the step by step explanation. The paper is concluded in Section V.

II. PĀLI LANGUAGE

Language is the speech, spoken by the people for communication, composed of letters (akkharā) or alphabet. Pāli is the language in which is composed the Tipiṭaka. The word Pāli is used in the sense of "Text", sacred Text and the same thing for the etymology of Pāli is the Holy Text, the Scriptures or the canon. Pāli language is a branch of Indo-European fami-y and a sister language of Sanskrit. Pāli was first committed to writing in Sṛlāṅkā in the 1st century AD for the Buddhist Canon. It is the spoken language. It has no own script but only sound. So, Pāli is transliterated into various local scripts. Pāli is an inflectional language (declension, conjugation, assimilation). Pāli had contributed mostly to the growth of Myanmar as a national language. Brahmanism, Hinduism, Buddhism, Sanskrit and Pāli have been well introduced to Myanmar from the beginning of our history of 4th century AD. The Myanmar invented Myanmar Script using Pallava Script. Moreover, phonetics and ideas were taken from Pāli language and literature. As Theravada Buddhism flourished in Myanmar, Pāli became the medium of writing [7].

III. PĀLI WORD NATURE IN MYANMAR LANGUAGE

Every language has its own alphabet which contains letters of that language. They are called akkharā, lipi, script or writing. The Akkahara (alphabet) means one which is eternam or imperishable, however, they are pronounced or used. Therefore, the original meaning of alphabet is "sound". As they describe the quality of sound, they are also called syllable (vanna). Ka, kha, ga, gha, na, etc. are characters of alpha-bet, these characters are called script (lipi). Generally, Pāli is the name of a language and it has no own script. When its alphabets are written, various scripts are used: in Myanmar the Myanmar, in India Devanāgarī, in Srilāṅkā the Sinha-lese, in Thailand the Siamese, etc. The speech spoken in a language can be written with symbol that are so called alphabet or script. Languages such as English, German used Roman script in writing. Because of Roman origin, it is called Roman script. Romanization is the representation of a language written in a non-Roman script using the Roman alphabet. In the Myanmar writing system, as example, ကံ၊ခါ၊ဂါ၊ဃါ၊ဇါ are Myanmar script and the corresponding Ro-man scrips are k, kh, g, gh, ṅ. The Pāli alphabet contains 8 vowels, 33 consonants, and 1 nasal sound. All vowels can produce their sounds by themselves. Consonants cannot produce their sounds by themselves. So they are called mutes which produce their respective sounds only in combination with the vowels. Pāli language contains 41 letters. The Pāli words are pronounced with defined Ṭhān [8]. Ṭhān means the organ of articulation. There are six organ of articulation in Pāli. They are –

- Kaṅṭha- Gutturals throat
- Tālu- Palatals, hard plate
- Muddha- Cerebrals, soft plate
- Danta- Dentals, teeth
- Oṭṭha- Labials, lips
- Nāsā- Nasals, nose

The 33 consonant can be classified into six group according to the above six Ṭhān. Therefore, the Pāli alphabet for both Myanmar script and Roman script classified with Ṭhān are

described in Table I and II. Moreover, it can also be classified with Vagga and Avagga sound as shown in Table III and IV.

TABLE I. Vowel in Myanmar script

Vowels							
အ	အာ	ဣ	ဤ	ဥ	ဦ	ဧ	ဩ

TABLE II. Consonant in Myanmar script

Consonants				
က	ခ	ဂ	ဃ	င
စ	ဆ	ဇ	ဈ	ည
ဏ	တ	ထ	ဒ	ဏ
ပ	ဖ	ဗ	ဗ	မ
ယ	ရ	လ	ဝ	သ
	ဟ	ဌ	အံ (၀)	

TABLE III. Vowel in Roamn script

Ṭhān	Short (Rassa Sara)	Long (Dhīgha Sara)
Gutturals	a	ā
Palatals	i	ī
Labials	u	ū
Gutturals+ Palatals	e	
Gutturals+ Labials		o

TABLE IV. Vowel in Roman script

Ṭhān	Vagga					Avagga (Unclassified)
	(Classified)					
	Unaspi-rate	Aspi-rate	Unaspi-rate	Aspi-rate	Nasals	h
First	Second	Third	Fourth	Fifth		
Kaṅṭha	k	kh	g	gh	ṅ	y
Tālu	c	ch	j	jh	ñ	r, l
Muddha	ṭ	ṭh	ḍ	ḍh	ṇ	l, s
Danta	t	th	d	dh	n	v
Oṭṭha	p	ph	b	bh	m	y
Nāsā					ṃ	

IV. PROPOSED SYSTEM ARCHITECTURE

There are four main modules in the proposed Pāli Romanization system: (i) text normalization, (ii) check the input text is Pāli word or not, (iii) Romanized the Pāli word and finally (iv) generate the speech of the Romanized Pāli words. Firstly, the input Myanmar Pāli words are preprocessed as syllable segmentation for the next processing. Then, the normalized words are checked which are Pāli words or not. If it is Pāli words, these are transformed into roman script by using proposed Romanization rules. Finally, the system generated the speech output of the Myanmar Pāli words. The system architecture of the proposed system is shown in Fig 1.


```

22. Else If ("o" is over (ဝဲထဲမှထဲဝဲ)) then
23.     ISPALI = true;
24. Else If (current syllable does not contain (excludeWord)) then
25.     ISPALI=true;
26. Else
27.     ISPALI=false;
28. End if
29. End procedure
End
    
```

Fig.3. The algorithm for checking Myanmar Pāli words

C. Romanized Pāli Words

In linguistics, Romanizaion is the conversion of writing from a different writing system to the Roman (Latin) script, or a system for doing so. Methods of Romanization include transliteration, for representing written text, and transcription, for representing the spoken word, and combinations of both. Transcription methods can be subdivided into phonemic transcription, which records the phonemes or units of semantic meaning in speech, and stricter phonetic transcription, which records speech sounds with precision [11].

Myanmar language is one of few alphabets capable of transcribing Pāli text with 100% orthographic fidelity. However, because Pāli is no longer a spoken, but a written language, the standard pronunciation of Pāli text occurs in agreement with the phonetic values and inherent rules of the corresponding alphabets used. Accordingly, like other language such as Thai, Sinhala, Lao and Khmer, speakers, the Myanmar have a very distinct accent when using Pāli words.

After checking the input word is a Pāli or non-Pāli word, if the input word is Pāli, the system converted into their corresponding Romanized symbol. Therefore, we created the Romanized table as shown in Table V. Then, the input Pāli words are transformed into Roman script by using the following 10 Romanization rules with example words.

TABLE V. Romanization table

Consonant	Unicode	Roman_Symbol
က	U+1000	k
ခ	U+1001	kh
ဂ	U+1002	g
...		
ဥ	U+102f	u
ဦ	U+1030	ū
ဧ	U+1031	e

Rule 1. Combination of Consonant with Vowels

k+ a= ka (က), k+ i= ki (ကီ), k+ u= ku (ကု), k+ e= ke (ကေ)
k+ ā= kā (ကာ), k+ ī= kī (ကီ), k+ ū= kū (ကူ), k+ o= ko (ကေ)

Rule 2. Combination of the vowels with the niggahita

a+ ṁ= aṁ (အံ), i+ ṁ= iṁ (အိ), u+ ṁ= uṁ (အူ)
ā+ ṁ= āṁ (အံ), ī+ ṁ= īṁ (အိ), ū+ ṁ= ūṁ (အူ)

Rule 3. Combination of the consonants, vowels with the niggahita

k+ a+ ṁ= kaṁ (ကံ), k+ i+ ṁ= kiṁ (ကိ), k+ u+ ṁ= kuṁ (ကူ)
k+ ā+ ṁ= kaṁ (ကံ), k+ ī+ ṁ= kiṁ (ကိ), k+ ū+ ṁ= kuṁ (ကူ)

Rule 4. First Group alphabet+ First Group alphabet

k+ k= kk sakka (သကုက)
c+ c= cc sacca (သဗ္ဗစ)
t+ t= tt vatta (ဝဋ္ဌ)
t+ t= tt satta (သတုတ)
p+ p= pp sappa (သပုပ)

Rule 5. First Group alphabet+ Second Group alphabet

k+ kh= kkh yakkha (ယကုခ)
c+ ch= cch accha (အစုဆ)
t+ th= tth sattha (သတုထ)
p+ ph= pph puppha (ပုပုဖ)

Rule 6. Third group alphabet + Third group alphabet

g+ g= gg agga (အဂ္ဂ)
j+ j= jj ajja (အဇ္ဈ)
d+ d= dd bhadda (ဘဒ္ဒ)
b+ b= bb sabba (သဗ္ဗ)

Rule 7. Third group alphabet+ Fourth group alphabet

g+ gh= ggh byaggha (ဗျဂ္ဂ)
j+ jh= jjh majjha (မဇ္ဈ)
d+ dh= ddh buddha (ဗုဒ္ဓ)
b+ bh= bbh labbhati (လဗ္ဗဘတိ)

8. Fifth group alphabet+ Consonant of same Thān

ñ+ g= ṅg maṅgala (မင်္ဂလ)
ñ+ c= ñc pañca (ပဉ္စ)
ṇ+ ḍ= ṇḍ kaṇḍa (ကဏ္ဍ)
n+ t= nt ananta (အနန္တ)
m+ bh= mbh sambhūla (သမုဘူလ)

9. Combination of Avagga letters

y+ y= yy seyya (သေယျ)

y+ v= yv	yvāham (ယွာဟံ)
y+ h= yh	paggayha (ပဂ္ဂယှ)
v+ h= vh	avhayati (အဝါ ယတိ)
l+ l= ll	vallari (ဝလ္လရိ)
l+ y= ly	kalyā (ကလျှ)
s+ y= sy	nisya (နိဿှ)
s+ v= sv	svāham (သွာဟံ)

10. Combination of the Vagga with Avagga

k+ y= ky	sakyamuni (သကျမုနိ)
k+ r= kr	cakra (စကရ)
d+ v= dv	dvāra (ဒွါရ)
n+ h= nh	nhāna (နာနာ)

The example Romanized Pāli sentence is shown in Table by using the mentioned Romanization rules.

TABLE VI. Sample Romanization result

Input Pāli Sentence	နမော တဿ ဘဂဝတော အရဟတော သမုမာ သမုဗုဒ္ဓဓဿ
Romanized Sentecne	Namo tassa bhagavato arahato sammā ssaṃbuddhassa

D. Generated Pāli Speech

In the speech generation step, the converted roman scripts are transformed again into speech output. MaryTTS speech engine is used in this potion. MaryTTS is a multilingual Text-to-Speech Synthesis platform written in Java and it is an open-source platform. It was originally developed as a collaborative project of DFKI’s Language Technology Lab and the Institute of Phonetics at Saarland University. Now, Multimodal Speech Processing Group in the Cluster of Excellence MMCI and DFK maintain the MaryTTS. As of version 5.2, MaryTTS supports German, British and American English, French, Italian, Luxembourgish, Russian, Swedish, Telugu, and Turkish; more languages are in preparation.

MaryTTS comes with toolkits for quickly adding support for new languages and for building unit selection and HMM-based synthesis voices [12].

V. EXPERIMENTAL RESULT

Generally, the performance of Myanmar Pāli word Romanization can be calculated in different ways. In this system, the goodness of transformation is measured by four types of outcomes: (1) Correct Transformed (CT): A Pāli word was converted correctly and is detected to be correct; (2) Correct Rejection (CR): A Pāli word was transformed incorrectly and is

detected to be incorrect; (3) False Transformed (FA): A Pāli word was converted incorrectly and is detected to be correct; (4) False Rejection (FR): A Pāli word was transformed correctly and is detected to be incorrect. For this four measure, 500 Pāli sentences are tested. The experimental results are shown in Fig 4. Typically, the performance of an error detection algorithm can be calculated in different ways. One way is to measure the scoring accuracy (SA), which is calculated by formula shown below:

$$SA = ((CT+CR) / (CT + CR + FT + FR)) * 100;$$

The ratio of CTs and CRs can be calculated by the classification algorithm: precision, recall, and F-measure metrics. These measurements are as follows:

$$\text{Precision of CT} = (CT / (CT + FA)) * 100;$$

$$\text{Precision of CR} = (CR / (CR + FR)) * 100;$$

$$\text{Recall of CT} = (CT / (CT + FR)) * 100;$$

$$\text{Recall of CR} = (CR / (CR + FA)) * 100;$$

$$F\text{-measure} = \frac{2*(Precision*Recall)}{(Precision+Recall)}$$

VI. CONCLUSION

This paper presented the rule based Myanmar Pāli word Romanization system. Therefore, the ten Romanization rules are discussed. Before the Pāli words are Romanized, the input words are checked that they are Pāli word or not. Consequently, the Myanmar Pāli word checking algorithm is presented. This paper is mainly focus on Myanmar Pāli word checking and Romanized these Pāli words so that the MaryTTS engine is used for speech generation. According to the experimental result for Romanization, the system achieved the overall accuracy is 89.6. For some words, such as “အောင်မင်္ဂလာ”. In this word, although “မင်္ဂလာ” is Pāli word, the syllalbe “အောင်” is not Pāli word. In this case, the system may wrongly check as the Pāli words. In this case, the accuracy may be decreased. Nowadays, in Myanmar, the researchers focus on the speech processing in Myanmar natural language processing. In the future, the high quality speech output will be generated by using other speech synthesis methods like concatenation speech synthesis approach by recording own voice for Pāli words.

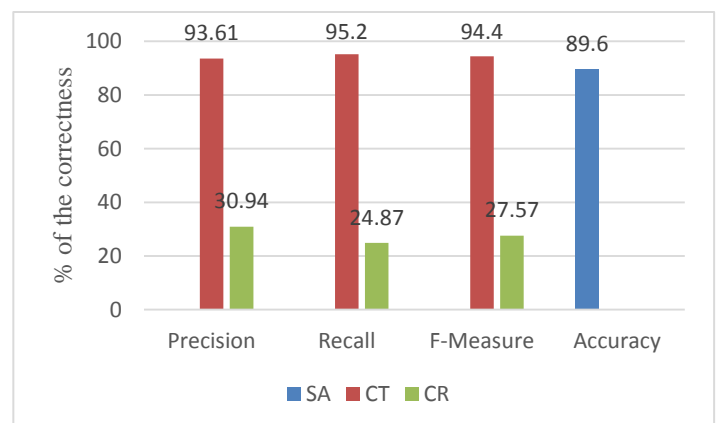


Fig.4. The experimental result of the Myanmar Pāli word
Romanization

REFERENCES

- [1] C. Duroiselle, "Practical Grammar of Pali Language", Third Edition 1997
- [2] M. H. Bode, The Pali Literature of Burma, The Royal Asiatic Society of Great Britain and Ireland, first published 1909, reprinted 1966
- [3] Ryang, Sonia, ed. Koreans in Japan: Critical voices from the margin. Routledge, 2013.
- [4] Aroonmanakun, Wirote. "A unified model of Thai romanization and word segmentation." In Proceedings of The 18th Pacific Asia Conference on Language, Information and Computation, pp. 205-214. 2004.
- [5] Kudo, Yoko. "Modified Hepburn Romanization System in Japanese Language Cataloging: Where to Look, What to Follow." Cataloging & Classification Quarterly 49, no. 2 (2011): 97-120.
- [6] "Romanization System for Shan", BGN/PCGN 2011 System.
- [7] "Fundamental of the Pāli Language", University of Mandalay, Department of Oriental Studies.
- [8] C. Duroiselle, A Practical Pali Grammar of the Pali Language, Third Edition, 1997, originally printed at the British Burma Press, 1921
- [9] C.S. Hlaing and A. Thida. "Phoneme based Myanmar text to speech system." International Journal of Advanced Computer Research 8, no. 34 (2018): 47-58.
- [10] Maung, Zin Maung, and Yoshiki Mikami. "A rule-based syllable segmentation of Myanmar text." In Proceedings of the IJCNLP-08 Workshop on NLP for Less Privileged Languages. 2008
- [11] <https://en.wikipedia.org/wiki/Romanization>
- [12] <http://mary.dfki.de/index.html>