

# Complexity of Letter to Sound Conversion (LTS) in Khmer Language: under the context of Khmer Text-to-Speech (TTS)

(Annanda th. Rath, Long S. Meng, Heng. Samedi, Long Nipaul, Sok K. heng)<sup>1</sup>

<sup>1</sup> NLP lab, Department of Computer and Communication Engineering, Institute of Technology of Cambodia, Cambodia, PAN10 and IDRC Canada

rannanda@itc.edu.kh, seangmeng@itc.edu.kh, sheng@itc.edu.kh, nipaul@itc.edu.kh, sok.kimheng@itc.edu.kh

## Abstract

Text to sound conversion is an important step in the process of creating the Text-to-Speech application for any languages. The complexity is generally based on the way the script is written. Some languages such as Khmer, Thai or Lao, one word may be pronounced to more than one ways depending on the syllable boundary in the word. This leads to the great difficulty in conversing correctly the text. The correctness of conversing text-to-sound in Khmer language is hard to archive 100%; however, it is not impossible.

This paper presents the steps in Khmer LTS, the difficulties and our research result.

**Index Terms:** LTS, TTS

## 1. Introduction

Speech synthesis systems use two basic approaches to determine the pronunciation of a word based on its spelling, a process which is often called text-to-phoneme, text-to-sound or grapheme-to-phoneme conversion. The simplest approach to text-to-phoneme conversion is the dictionary-based approach, where a large dictionary, also called lexicon, containing all the words of a language and their correct pronunciations is stored by the program. Determining the correct pronunciation of each word is a matter of looking up each word in the dictionary and replacing the spelling with the pronunciation specified in the dictionary. The other approach is rule-based, in which pronunciation rules are applied to words to determine their pronunciations based on their spellings.

Each approach has advantages and drawbacks. The dictionary-based approach is quick and accurate, but completely fails if it is given a word which is not in dictionary. As dictionary size grows, so too does the memory space requirements of the synthesis system, another drawback is that we have to update our dictionary whenever the new word comes. On the other hand, the rule-based approach works on any input, but the complexity of the rules grows substantially as the system takes into account irregular spellings or pronunciations. As a result, nearly all speech synthesis systems use a combination of these approaches.

Some languages, like Spanish, have a very regular writing system, and the prediction of the pronunciation of words based on their spellings is quite successful. Speech synthesis systems for such languages often use the rule-based method extensively, resorting to dictionaries only for those few words, like foreign names and borrowings, whose pronunciations are not obvious from their spellings. On the other hand, speech synthesis systems for languages like English, which have extremely irregular spelling systems, are more likely to rely on

dictionaries and to use rule-based methods only for unusual words, or words that aren't in their dictionaries. [1]

Like English, our future Khmer Text To Speech System will rely on lexicon and use rule-based methods only for unusual words, or words that are not present in the lexicon.

## 2. Background of Khmer Language

The Cambodian scripts (called Khmer letters) are derived from various forms of the ancient Brahmi script of South India, and the ancient Khmer script itself. The Cambodian script has symbols for thirty-three consonants, twenty-four dependent vowels, twelve independent vowels, and several diacritic symbols. Most consonants have reduced or modified forms, called sub-consonants or subscript, when they occur as the second member of a consonant cluster. Vowels may be written before, after, over, or under a consonant symbol. [2]

The consonants are divided into 2 groups or series:  $\alpha$ -series and  $\beta$ -series (see fig. 1 for the group of consonants along with their pronunciation). Most of consonants (in  $\alpha$ -series resp.  $\beta$ -series) have their counterparts (in  $\beta$ -series resp.  $\alpha$ -series). For consonants in one group that don't have their counterpart in another group, diacritic mark  $\ddot{\circ}$  (MUUSIKATOAN) and  $\tilde{\circ}$

(TRIISAP) are used to produce their counterpart.  $\ddot{\circ}$  changes

the consonant to  $\alpha$ -series while  $\tilde{\circ}$  changes the consonant to  $\beta$ -series. Among the 24 vowels there is 1 abstract (inherent) vowel which is embedded in a consonant. The pronunciation of a vowel, including the inherent vowel, is determined by the series of the initial consonant or consonant cluster that it accompanies. Khmer syllable in written form can contain 1 or 2 vowels. In the later case, last vowel is a vowel that has embedded consonant (see fig. 2 for the list of vowels and sequences of vowels that come together along with their pronunciation for each group). Independent vowels incorporate with both an initial consonant and a vowel (see fig. 3 for list of independent vowels along with their pronunciations). Khmer words are predominantly of one or two syllables. Words with three or more syllables exist, particularly those pertaining to science, the arts, and religion. These words are loanwords, usually derived from Pali, Sanskrit, or more recently, French. [3] Among the loanwords, about 90% are from Pali and Sanskrit. [4] Here, we classify words in Khmer language into 3 groups: original Khmer words (group 1), words borrowed from Pali and Sanskrit (group 2) and words borrowed from other languages like English, French and Chinese (group 3). The majority of words in Khmer language are in the first 2 groups. The last group contains only a small number of words.

Consonants and subscripts				
Letter				Sound
α-series		ɔ-series		
Cons.	Sub.	Cons.	Sub.	
ក	ក	គ	គ	k
ខ	ខ	ឃ	ឃ	k <sup>h</sup>
ង	ង	ង	ង	ŋ
ច	ច	ជ	ជ	c
ឆ	ឆ	ឈ	ឈ	c <sup>h</sup>
ញ	ញ	ញ	ញ	ɲ
ដ	ដ	ឌ	ឌ	d
ត, ថ	ត, ថ	ណ, ណ	ណ, ណ	t <sup>h</sup>
ណ	ណ	ន	ន	n
តិ	តិ	ទ	ទ	t
ប	ប	បី	បី	b
ផ	ផ	ភ	ភ	p <sup>h</sup>
ប៉	ប៉	ព	ព	p
ម	ម	ម	ម	m
យ	យ	យ	យ	j
រ	រ	រ	រ	r
ឡ	ឡ	ល	ល	l
វ	វ	វ	វ	w
ស	ស	ស	ស	s
ហ	ហ	ហ	ហ	h
អ	អ	អ	អ	ʔ

Figure 1. Text to sound mapping for consonants and subscripts

Vowels		
Letter	Sound	
	α-series	ɔ-series
Inherent vowel	ɑ	ɔ
ា	a:	i:ə
ិ	e	i
េ	ej	i:
ុ	ə	ɨ
ូ	œ:	ɨ:
ុ	o	u
ួ	ɔ:ɔ	u:
ុ	u:ə	u:ə
ើ	a:ə	ə:
ើ	ɨ:ə	ɨ:ə
ើ	i:ɜ	i:ɜ
ើ	e:	ɛ:
ើ	a:ɛ	ɛ:
ើ	aj	ej
ើ	a:o	o:
ើ	aw	əw
ើ	om	um
ើ	am	um
ើ	am	oam
ើ	ah	ɛah
ើ	eh	ih
ើ	əh	ɨh
ើ	oh	uh
ើ	eh	ih
ើ	ɔh	uəh

Figure 2. Text to sound mapping for vowels and sequence of vowels

Independent vowels	
Letter	Sound
ា	ʔa:
េ	ʔe
េ	ʔej
ុ	ʔu
ុ	ʔu:
ុ	ʔo:w
ុ	ʔɨ
ុ	ʔɨ:
ុ	li
ុ	li:
ុ	ʔa:ɛ
ុ	ʔaj
ុ	ʔa:ɔ
ុ	ʔa:ɔ
ុ	ʔaw

Figure 3. Text to sound mapping for independent vowels

### 3. Pronunciation of Khmer Word

As we can see in fig. 4, the third consonant រី can be mapped to /kɑ/ or /k/ depending on whether it is in the onset or the coda of a syllable. Note that the consonant រ is deleted if it is in the coda. So without knowing syllable boundary, pronouncing Khmer word is relatively difficult if not impossible. But once we have identified syllable boundary, the pronunciation is quite easy. So our approach to find pronunciation is in three steps. The first step is to do preprocessing to find syllable boundary and to replace diacritic

marks with appropriate text. The second step is to find the pronunciation of each syllable. And the last step is to apply sound change rules. Note that for our approach, we don't need syllabification at phoneme level.

We will firstly present the rule to pronounce each syllable and after that we will show how the preprocessing is done.

ក	.	ក	ៃ	.	ក	.	ក	រ	ក	.	ក	ៃ	ក	.	ក	រ
ka:	.	k	aɛ	.	ka:	.	ka:		ka:	.	k	aɛ	k	.	ka:	

Figure 4. Pronunciation of the word កៃកករី

### 3.1 Rule and algorithms for syllable pronunciation

In written form, a syllable in Khmer word can be one of this: C, CC(្ក), CS, CSC(្ក), I, IC, CV, CVC(្ក), CSV, CSVC(្ក), CSSV where C is a consonant, S is a subscript, I is an independent vowel, V is a vowel, ្ក (called Bantoc) is diacritic mark to make vowel a short vowel and parenthesis means optional. **Remark:** C, S, I and V are all in script form not in phoneme form.

#### Algorithm:

For each character in the syllable:

If the character is a consonant or subscript

Take corresponding sound in the text to sound mapping for the consonant and subscript

If it is not followed by a vowel or a subscript

Add inherent vowel (see below)

If the character is a vowel (including inherent vowel which may have been added in previous step)

Take the corresponding sound in text to sound mapping for the vowel where the group is the group of the preceding consonant or consonant cluster. To find group the consonant is in, just look in text to sound mapping table for the consonant. For consonant cluster see section 4 below.

If the character is an independent vowel

Take the corresponding sound in text to sound mapping table for independent vowel.

If there is ្ក at the end then make the vowel a short vowel

### 3.2 Find the group of the consonant cluster

Consonant cluster in Khmer word is one of the forms CS or CSS. For the second form CSS there is only 2 possibilities:

[ក្ករ] /str/ and [ក្ករ] /skt/ which are in α-group.

For CS form, the algorithm to determine the group is as follow:

If the consonant and subscript are in the same group

The group of the cluster is the same as the group of consonant or subscript

Else If the subscript is not in [ក្ក ក្ក ក្ក ក្ក ក្ក ក្ក]

The group of the cluster is the same as the group of the subscript with two exceptions: ក្ក្ក (solidly packed) and ក្ក្ក (to tame)

Else the group of the cluster is the same as the group of the consonant

### 3.3 Preprocessing

For words with more than one syllable, the group (α or ɔ) of the consonant or consonant cluster of the second or higher



2. Insert syllable boundary after vowels ័្រ ័្រ and diacritic marks

័

3. Insert syllable boundary before and after the form  $C\overset{\circ}{\circ}$  and remove diacritic mark ័

4. for vowel ័

If it is preceded by vowel ័ and followed by consonant ឃ

Insert syllable boundary after ឃ

Else Insert syllable boundary after ័ and change it to ័

5. Insert syllable boundary after diacritic mark ័: and replace it by

័

6. Insert syllable boundary before independent vowels

7. For the form CV, insert syllable boundary before C

8. Insert syllable boundary before the form  $CC\overset{\circ}{\circ}$

9. If the form  $C_1C_2\overset{\circ}{\circ}$  is at the end of the word, insert syllable boundary before it and change  $C_2\overset{\circ}{\circ}$  to ័

10. If the form  $C_1C_2\overset{\circ}{\circ}$  is not at the end of the word, insert syllable boundary before it and change it to  $C_1\overset{\circ}{\circ}.C_2$

11. If the form  $C_1\overset{\circ}{\circ}C_2$  is at the end of the word, insert syllable boundary before it and change it to ័.  $C_1C_2\overset{\circ}{\circ}$  ('.' denotes syllable boundary)

12. If the form  $VC\overset{\circ}{\circ}$  is at the end of the word, change it to V ័

13. If the form  $\overset{\circ}{\circ}CS$  appears at the end of the word, remove S

14. If the form  $\overset{\circ}{\circ}CV$  appears at the end of the word, remove V

15. If the form V ័ appears at the end of the word, remove vowel ័

16. For consonant cluster CS which is not at the beginning of the word and no syllable boundary has been added before C

If CS is possible consonant cluster (see fig. 6 for the list of possible consonant clusters)

Insert syllable boundary before C

Else Insert syllable boundary after C

If C is preceded by a consonant

Insert ័ before C Replace subscript S with the corresponding consonant

17. For the form  $\overset{\circ}{\circ}C$ , insert syllable boundary after C

18. If the form  $C\overset{\circ}{\circ}$  appears at the end of the word, insert syllable boundary before C and change ័ to ័

19. If the form  $C\overset{\circ}{\circ}$  appears at the end of the word, insert syllable boundary after C

20. If the form  $C\overset{\circ}{\circ}$  appears at the end of the word, insert syllable boundary before C and change ័ to ័

21. If the form  $C\overset{\circ}{\circ}$  appears at the end of the word, insert syllable boundary before C and change ័ to ័

22. If the form  $\overset{\circ}{\circ}$  appears at the end of the word, change it to ័

23. If the form  $\overset{\circ}{\circ}CS$  appears at the end of the word, remove S

24. For the form  $\overset{\circ}{\circ}S$ , insert syllable boundary after ័ and insert ័ before S

25. For the form  $\overset{\circ}{\circ}C$ , insert syllable boundary after C

26. If the form  $C_1C_2$  is not at the end of the word and it is not followed by diacritic mark ័, insert syllable boundary after  $C_1$

27. If the form VC is not at the end of the word, insert syllable boundary after V

## 4. Evaluation and discussion

We implemented this algorithm using Python. To test the accuracy, we have built a lexicon containing 9445 words. The accuracy is about 72%. Most of the incorrect results are due to the group change we mentioned at the beginning of section 2. There are also some incorrect results due to incorrect Pali / Sanskrit word detection. Some of the errors are Pali / Sanskrit words. We haven't tested the accuracy of our Pali / Sanskrit word detection yet because we don't have information on word group in the lexicon we have built for the test.

The result obtained after the evaluation is acceptable, but not well. Further works have to be done to improve our Text to Sound Conversion algorithms. In our method, the improvement can be done by adding the module to predict group change of consonant or consonant cluster of the second or higher syllable in the word with more than one syllable. The statistical approach as described in [5] may be also be used.

## 5. Complexities of LTS in Khmer language

The difficulties that we found are like others languages such as Thai or Lao, where one word may be pronounced to more than one ways depending on the syllable boundary. This leads to the ambiguity of the meaning while conversing to the sound. For example, in Khmer language, the word ័័័័័័ (to mumble) can have 2 different pronunciations depending on syllable boundary. first one is /kɑ:kɑ:kɑ:kɑ:/ corresponding to ័.័.័.័ and another is /kɑ:kɑ:kɑ:/ corresponding to ័.័័.័.័. Some methods such content analysis technique is needed to improve the performance of the TTS application creating. In Khmer language, there are many most frequently used words that fall into the category mentioned above. As the result of our experimentation, we found that most of the error came from the mis-identifying of the meaning of the words as those words have the ambiguity form.

## 6. Conclusions

Our method to do text to sound conversion is by doing some preprocessing at script level to determine syllable boundary and replace diacritic marks with appropriate text. After the preprocessing, we convert each syllable to sound using some simples' rules. The last step is the application of sound change rules. The output of our text to sound conversion also includes syllable boundary. To measure the accuracy of our method, we implemented our algorithm using Python and built a lexicon containing about 10000 words. The test yields about 72% accuracy. This result is the first outcome of our research on Khmer language. The further research needs to be done in order to improve the accuracy of the algorithms that we have developed. This work will be carried on in next phrase.

## References

- [1] [http://en.wikipedia.org/wiki/Speech\\_synthesis](http://en.wikipedia.org/wiki/Speech_synthesis), 09/10/2008
- [3] [http://en.wikipedia.org/wiki/Khmer\\_language](http://en.wikipedia.org/wiki/Khmer_language), 10/10/2008
- [2][http://www.seasite.niu.edu/khmer/writingsystem/writing\\_introduction/intro\\_set.htm](http://www.seasite.niu.edu/khmer/writingsystem/writing_introduction/intro_set.htm), 10/10/2008
- [4] Khmer grammar
- [5] Issues in Building General Letter to Sound Rules, Alan W Black – Kevin Lenzo – Vincent Pagel