

# Thai to Khmer Rule-Based Machine Translation using Reordering Word to Phrase

Sukchatri Prasomsuk<sup>1</sup> and Sukchatri Prasomsuk<sup>2</sup>

<sup>1</sup> University of Phayao

*Received: 8 December 2016 Accepted: 5 January 2017 Published: 15 January 2017*

---

## Abstract

In this paper, an effective machine translation system from Thai to Khmer language on a website is proposed. To create a web application for a high performance Thai- Khmer machine translation (ThKh-MT), the principles and methods of translation involve with lexical base. Word reordering is applied by considering the previous word, the next word and subject-verb agreement. The word adjustment is also required to attain acceptable outputs. Additional steps related to structure patterns are added in a combination with the classical methods to deal with translation issues. PHP is implemented to build the application with MySQL as a tool to create lexical databases. For testing, 5,100 phrases and sentences are selected to evaluate the system. The result shows 89.25 percent of accuracy and 0.84 for F-Measure which infers to a higher efficiency than that of Google and other systems.

---

*Index terms*— thai khmer translation, machine translation (MT), rule based, pattern-based.

## 1 Introduction

Association of Southeast Asian Nations (ASEAN) consists of ten countries with various cultures and languages. Thailand and Cambodia are included in ASEAN, and the eastern border of Thailand is adjacent to Cambodia. Therefore, efficient communication is significant for international relations between these two countries. Cambodian natives have Khmer as a national language while formal language in Thailand is Thai. The linguistic differences of Thai and Khmer in both writing and speaking contribute to a translation barrier. For instance, since Thai language has been adapted partly from Pali, Sanskrit and Old Khmer, Thai vocabulary is relatively diverse. Thai language also contains complex orthography and relational markers. Furthermore, standard written Thai is complicated due to various combinations of syllabic alphabets, which consists of 44 basic consonants, 21 vowel symbols and 4 tone diacritics, applied under the rule that all diacritics appear in front of, above or below the consonants. Furthermore, Thai syntax has a noun classifier system as well as conforms to a basic sentence structure called subject-verb-object (SVO) Khmer contains 33 consonants, 23 dependent vowels and 15 independent vowels; however, no tone is presented. Due to the linguistic differences, current Thai-Khmer translation systems have scarcely achieved complete and accurate outputs. Moreover, the existent systems have rarely been created and developed. There is also a shortage of intellectuals who are competent in both languages and able to convey knowledge for creating a system of translation. As a result, the improvement of the Thai-Khmer translation system has been disrupted. Document translation between Thai and Khmer which requires high accuracy has consequently encountered difficulties. To solve the issues, machine translation (MT) from Thai to Khmer language requires development.

The proposed system in this paper implements translation techniques including rule-based algorithm with verification of sentence patterns to improve translation quality. The overview operation of the translation system is to input a Thai language text in a web application and then convert it into a desired output in Khmer. A lexical analyzer is first applied in the process to divide Thai sentences or phrases into individual syllabic words so that the separated words are analyzed and processed in the following steps resulting in Khmer sentences.

## 2 II. Related and Previous Works

There have been many attempts to research on machine translation between Thai and other languages. English-Thai machine translation was developed in 1998 with regard to the sentence-based technique which combines the rule-based and the example-based method to establish a system for English to Thai sentence translation [1]. However, the research result of performance evaluation and comparison was not indicated. In 2012, a technique called generalized patterns is presented to improve machine translation from Japanese to Thai language [2]. The method was compared to the others implemented in Google and Bing translators by executing 3,107 Japanese sentences in testing. F-Measure score was applied to assess performance of the translator.

Machine translation between Khmer and other language has also been researched. One of the studies selected Moses DoMY CE, which is statistical machine Abstract-In this paper, an effective machine translation system from Thai to Khmer language on a website is proposed. To create a web application for a high performance Thai-Khmer machine translation (ThKh-MT), the principles and methods of translation involve with lexical base. Word reordering is applied by considering the previous word, the next word and subject-verb agreement. The word adjustment is also required to attain acceptable outputs. Additional steps related to structure patterns are added in a combination with the classical methods to deal with translation issues. PHP is implemented to build the application with MySQL as a tool to create lexical databases. For testing, 5,100 phrases and sentences are selected to evaluate the system. The result shows 89.25 percent of accuracy and 0.84 for F-Measure which infers to a higher efficiency than that of Google and other systems.

with a horizontal and vertical writing direction from left to right and from top to bottom, respectively. Similarly, translation (SMT), as a tool to create an online system for English -Khmer translation based on Python, XML and HTML language in 2013 [3]. There is also research in 2014 on developing a French-Khmer dictionary called 'MotàMot' [4]. In 2015, an automatic machine translation was created to provide translation between Khmer and other 20 languages by using three statistical methods: the phrase-based approach, the hierarchical phrasebased approach and the operation sequence model (OSM) as well as selecting BLEU and RIBES to evaluate translation quality [5].

There is, furthermore, research specifically on Thai-Khmer machine translation. For example, Thai -Khmer machine translation on a website has been developed based on Java (JSP) and SQL (Appserv) with 4,000 words from a Thai -Khmer dictionary as a database [6]. In testing, 212 sample sentences have been processed, and the result has shown 72.16% of accuracy which is higher than that of Google translator. In 2014, the rule-based machine translation (RBMT) combined with statistical methods was recognized to be widely applied in automated translation [7]. The technique has shown the potential to improve translation between Thai and Khmer. Even though such classical technique is applied, the research has rarely a result with high performance.

## 3 III.

### 4 Background of thai to Khmer Translation

Sentences in Thai and Khmer language are similarly formed; on the other hand, ordering and semantic structure are different. With regard to the existent methods, the newly presented one for the proposed system is expected to balance between advantages and disadvantages of the classical techniques and be straightforward for implementation. In this paper, a process to translate Thai to Khmer language is composed of six main steps including 1) Input process: reading Thai text into the system from a website screen, 2) Word segmentation: applying Lex To and the longest matching approach to divide Thai sentences into words, 3) Word search: retrieving data from the database of Thai-Khmer dictionary to find a matched-meaning word in Khmer for each Thai word, 4) Boundary check: considering a boundary of each Thai word such as conjunction, verb, adjective and surrounding nouns to inspect parts of speech, 5) Pattern verification: examining Thai sentence patterns by using the rule-based algorithm, and 6) Khmer word rearrangement: reordering Khmer words in phrases or sentences.

To build a Thai-Khmer dictionary for testing in this paper, approximately 37,052 Thai words from the Royal Institute Dictionary (RID, 1999) are translated according to the existent Thai-Khmer dictionary [8]. In the process of examining word boundaries, patterns and various conditions of grammar rules are taken into account to solve translation mistakes.

The classification of machine translation architecture which is regularly implemented on clientweb server for online translation is the direct model shown in Figure 1. The direct machine translation architecture transforms a source language sentence (Thai) into a target language sentence (Khmer). Besides, the proposed system applies the indirect architecture which is demonstrated as a diagram in Figure ???. A sample screen of the program is also provided Figures 3. IV.

## 5 Web Page

### 6 Methodology and Proposed Algorithm of Reordering

In general, Thai and Khmer sentences are sorted verbatim. Regarding to the verbatim characteristics of these two languages, the classical algorithm of word reordering could appear to be a proper tool to cope with phrase and sentence arrangement. On the other hand, the reordering method is unable to suit all cases of input phrases



161 (k?un-maa-t?ii ?nii-tham ?may)Khmer: ?? ????? ?á???????? ? ? ? (neak-mÉ??"É?"k-tii nih-haet Ê?"v?y) ??  
 162 ????? ?á???????? ? ? ? ? (neak-mÉ??"É?"k-tii nih-tv?? Ê?"v?y) B. Sample 2: [ThPattern 1 ] [KhPattern 12 ]  
 163 Thwd{1} + Thwd{2} + Thwd{3} + ? + Thwd{ x n } Khwd{1}+Khwd{2}+Khwd{4}+Khwd{3}+Khwd{5} +  
 164 ? + Khwd{ x n }

165 If the word "?????" (moo?) follows a number (of time indications), the Khmer word "???? ?á?" "á?"?" ? ?"  
 166 (pram bu?n-mao?) is replaced by drow eht "?? ? ????? ?á?" "á?"?" or swap the position with that of the word "??  
 167 ? ?-????á?" "á?"?". In this case, the example sentences are provided below. Reordering words and translating  
 168 are in the final step to diminish the translation issue. After the pattern each word which is then rearranged to  
 169 be in a proper position. As a result, a Khmer sentence is attained as the output.

170 V.

## 11 Performance Evaluation

171 The proposed system is assessed for translation performance from Thai to Khmer by sentences from various  
 172 sample documents as the input. In the testing process, the total phrases and sentences The algorithm is applied  
 173 for a Thai phrase consisting of a word (Thwd{2}) ?? ?.

174 Thai: "?? ?????à,?"? ???? ??????? ? ? ?"  
 175 (iik -sÉ??"É?"?d?an-c?a?-jaÉ?"-pay-kam.p?uu?c?aa) (Next two month I will go to Cambodia.) Khmer:

176 ?????? ????? ? ? á?"? ?????? ? ? ? (ti?t-pii-k?ae-k?om-n??-t?v-kampuÊ?"cie) ?? ????????? ? ?  
 177 á?"? ?????? ? ? ? (pii-k?ae-ti?t-k?om-n??-t?v-kampuÊ?"cie) mapping is completed, around 37,000 words from  
 178 Thai -Khmer dictionary database are retrieved to match Table I: Sample of Phrase/Sentences for testing Three

179 translation systems including Google translator [11], Chhun's translation system and the proposed system in  
 180 this paper are assessed through translating the sample phrases and sentences. The translated outputs of each  
 181 translation system are categorized into three groups consisting of accuracy (correct), acceptance (acceptable) and  
 182 mistake (wrong). According to 5,100 Thai sentences selected for testing, the proposed system is able to translate  
 183 4,083 words correctly (80.06%), reach the acceptable level of translation for 469 sentences (9.189%) and produce  
 184 errors only in 548 sentences (10.75%). The total translation accuracy of the proposed system becomes 89.25 %  
 185 which is a sum of its accuracy and acceptance value. On the other hand, Chhun's translation contributes to 3,590  
 186 correct sentences (70.38%) which is less than those of the proposed system, 658 acceptable sentences (12.9%) and  
 187 857 mistakes (16.81%). Google translation also achieve less accuracy compared to the proposed system: 1,067  
 188 correct sentences (20.9%). whereas it acquires 798 acceptable sentences (15.64%) and 3,230 mistakes (63.34%),  
 189 respectively, higher than those of the proposed one.

190 Moreover, performances of all systems are compared with regard to system precision, recall and efficiency by  
 191 implementing F-measure as shown in Table ?? The result in Tables 2 reveals that the proposed system attains  
 192 the highest score in all evaluations: the precision is 0.89, the recall is 0.80 and the efficiency (F-Measure) is 0.84.

## 12 VI.

## 13 Conclusion

195 The methodology in this paper is presented for creating Thai to Khmer machine translation system by using  
 196 syntactic and semantic analysis to transform and structure patterns as well as implementing the rulebased  
 197 translation. The presented processes can also simplify compound sentences into simple ones based on predefined  
 198 sentence structures. The previous word, the next word and the subject-verb agreement are also considered. In  
 199 addition, switching with more suitable words, reordering words and adjusting output sentences are also performed  
 200 with regard to Thai and Khmer grammar. As a result, the proposed system is apparently able to improve the  
 201 quality of source texts and translated outputs as well as assist Thai-Khmer language learners. Nevertheless, a  
 202 larger amount of sample sentences in the corpus than that which is currently applied in the proposed system is  
 203 necessary to achieve higher performance in Thai-Khmer translation. Furthermore, the larger dictionary database  
 204 as well as the higher diversity of sample sources would be added to the process. Other methods or tools would  
 205 also be considered to develop Thai-Khmer translation in future research.

206 VII. <sup>1 2</sup>

<sup>1</sup>© 20 7 Global Journa ls Inc. (US)

<sup>2</sup>© 20 7 Global Journa ls Inc. (US) 1



This car is red.

Thai sentence  
S + V + O  
?? ? + ?? ? + ?? ?? c?a? + kin+  
k?aâw

Year 2017

31

)

( H

Khmer sentence

S + V + O

= ?? ? ? + ?? ? +?? k?om + ?am +  
bay

Figure 3:

Example:

Thai sentence S + Khmer sentence S + (V) + O

(V) + O

?????? + ?? ??? ? + ?? ?à\_?"? = ???á???"? + ?á???"? + ??? ???? rot ?yon + k?an-ní+ si? ?d??? r

Example:

[ThPattern  
x ]

Year Thai sentence S + Khmer sentence S + V + O + unit(s)

2017 V + O + unit(s)

32 + ?? + ??? + 5 + ??? = +?á???"+?? ? ?+??? +???á???" + mii + k?a ?y + haâ + fÉ?"É?"? + mien

Volume

Note: ? =

XVII

incorrect

Issue

ordering

III

words, ? =

Ver-

Correct A.

sion

Sample 1:

I

[ThPattern

1 ]

[KhPattern

11 ]

Thwd{1} +

Thwd{2} +

Thwd{3} +

? + Thwd{

x n }

)

( H

Khwd{1}+Khwd{2}+(non,

Khwd{3-1}

or Khwd{3-

2})

Global [ThPattern 1 ]

[KhPattern 11 ] : : : [KhPattern 1n ] [KhPattern 21 ] [KhPattern 22 ] [KhPattern

Journal of  
Computer  
Science  
and  
Technology

Translation Methods	Precision	Recall	F-Measure
Google	0.57	0.21	0.31
Chhun	0.84	0.70	0.76
Proposed System	0.89	0.80	0.84

Figure 5:

208 .1 Acknowledgment

209 This research was funded by School of Information and Communication Technology, University of Phayao. We  
210 would like to show our appreciation to Cambodian students who have assisted the research by correcting Khmer  
211 words, phrases and sentences. We are also immensely grateful to other researchers for their supports in this  
212 research project.

213 .2 Global Journals Inc. (US) Guidelines Handbook 2017

214 www.GlobalJournals.org

215 [Vanavong et al. ()] , Y Vanavong , R Saravut , Thai-Khmer Dictionary . 2012. Cambodia: Nokothom Book  
216 Shop. p. .

217 [Thu et al. (2015)] ‘A Large-scale Study of Statistical Machine Translation Methods for Khmer Language’. Y K  
218 Thu , V Chea , A Finch , M Utiyama , E Sumita . *The 29th Pacific Asia Conference on Language, Information  
219 and Computation*, (Shanghai, China) October 30 -November 1, 2015. p. .

220 [Available ()] <http://translate.google.com/> Available, 2016.

221 [ S ()] ‘Color in Khmer : Perception and Grammatical Construction’. S . *Thesis, Silpakorn Univ* 2005. p. .

222 [Loy] *Fundamental Khmer 2*, P S Loy . p. . Ramkhamhaeng University

223 [Google web site, Google Translator] *Google web site, Google Translator*,

224 [Jabin et al. (2013)] ‘How to Translate from English to Khmer using Moses’. S Jabin , S Samak , K Sokphyrum  
225 . *International Journal of Engineering Inventions* Sep.2013. p. .

226 [Sae-Tang and Prayote (2012)] ‘Japanese-Thai Machine Translation with Generalized Patterns’. P Sae-Tang , A  
227 Prayote . *Computing and Convergence Technology (ICCCT), 7th International Conference*, (Seoul, Korea)  
228 Dec. 2012. IEEE. p. .

229 [Mangeot] ‘MotàMot project: conversion of a French-Khmer published Dictionary for building a multilingual  
230 lexical system’. M Mangeot . <http://arxiv.org/abs/1405.5674> *Languag-es Resources and Evaluation  
231 Conference*, p. .

232 [Malézieux et al. ()] ‘Rule-Based machine translation (RBMT) with statistical knowledge’. G Malézieux , A Bosc  
233 , V Berment . *Proceedings of the 5th Workshop on South and Southeast Asian NLP, 25th International  
234 Conference on Computational Linguistics*, (the 5th Workshop on South and Southeast Asian NLP, 25th  
235 International Conference on Computational LinguisticsDublin, Ireland) August 23-29, 2014. p. .

236 [Chancharoen et al. (1998)] ‘Sentence-based machine translation for English-Thai’. K Chancharoen , N Tannin ,  
237 B Sirinaovakul . *The 1998 IEEE Asia-Pacific Conference on Circuits and Systems, Chiangmai Plaza Hotel*,  
238 (Chiangmai, Thailand) Nov.1998. p. .

239 [Table II: F-Measure Results for Thai into Khmer] *Table II: F-Measure Results for Thai into Khmer*,

240 [Chhun ()] ‘Thai-Khmer simple sentence translation on Web’. C Chhun . *the 1st RMUTL Chiangrai National  
241 Conference, (RCCON)*, (Thailand) 2015.