

Introduction &  
Background

Objective

Methodology

Result

Discussion

Conclusion

Thanks



**UTM**  
UNIVERSITI TEKNOLOGI MALAYSIA

**Study of Machine Translation  
to Build Pattern-based Model  
For Vietnamese – Cham**



Mohamad Budi

Ph.D Student. Mohamad Budi  
Prof.Dr. Mohamad Bin Bilal Ali  
Dr. Noor Dayana Abd Halim

# 1. Introduction & Background

- There are various efforts have been made in developing Machine translation (MT) systems and many approaches on MT research: Pattern-based, transfer-based, interlingua-based, and etc.
- Pattern-based MT is a very traditional MT method that uses translation patterns and translation word (phrase). In this method the translation get high-quality translation results if the input sentence matches the translation pattern and this translation pattern is correct.
- Preliminary studies, we found that this method suitable for Vietnam Cham language, because Cham grammar structure is unclear, hence, we need to create the bilingual corpus, bilingual dictionary and translation sample to support for this method.

Next

## 2. Background

- In the pioneering work of Nagao, machine translation (MT) is the translation of text by a computer with no human involvement. (Nagao, 1984; Sato & Nagao, 1990).
- Today with various efforts are involved in developing MT systems for practical use. A few different types of main techniques for machine version; Statistics Based Machine Translation (SBMT), Rules Based Machine Translation (RBMT), and Hybrid Systems which combine RBMT and SBMT are being used. (Mamta, 2015).
- Pattern-based MT is a traditional method that was proposed in the 1960s (H.Maruyama, 1993).



### 3. Objective

The objectives of the research to be investigated are as follows:

- To build machine translation model for Vietnamese - Cham language.
- To build data source for bilingual corpus, bilingual dictionary, translation sample to support for this method.
- To develop application program for proposed model.



## 4.1 Vietnamese and Cham Simple Sentence Structure

- Many linguists classify that Cham language is a member of the Malayo-Polynesian and Cham language is polysyllabic word forms.
- Vietnamese language is a part of Austroasiatic language. Used Latin alphabet with additional diacritics for tones and certain letters. Vietnamese language is Isolating language with single syllabic.

# Vietnamese and Cham Simple Sentence Structure

The similarities both Vietnamese and Cham language were relatively orderly. Grammatical structure: Subject (S), Verb (V) and Object (O). It can be divided into 6 difference patterns such as SVO, SOV, VSO, VOS, OVS and OSV. Almost all the language sentences are SVO or SOV (Russel, 1986). However, the structure of Vietnamese and Cham language is similar: SOV.

No.	Bilingual Sentences			
(1)	Tom	nhớ	Jakei	Vietnamese
	Tom	miss	Jakei	
	Ớ	Ớ	Ớ	Cham
	S	V	O	
(2)	Jakei	nhớ	Tom	Vietnamese
	Jakei	miss	Tom	
	Ớ	Ớ	Ớ	Cham
	S	V	O	
(3)	Nhớ	Jakei	Tom	Vietnamese
	Miss	Jakei	Tom	
	Ớ	Ớ	Ớ	Cham
	V	S	O	

Figure 1. Simple Structure of Vietnamese and Cham Languages

## 4.2 Vietnamese Cham Pattern-Based Machine Translation<sup>7</sup>

The model of pattern-based MT system as proposed has three main resources: Bilingual corpus, bilingual Vietnamese-Cham dictionary, and translation sample.

Step 1. Prepare Bilingual corpus, Bilingual dictionary, and Translation sample.

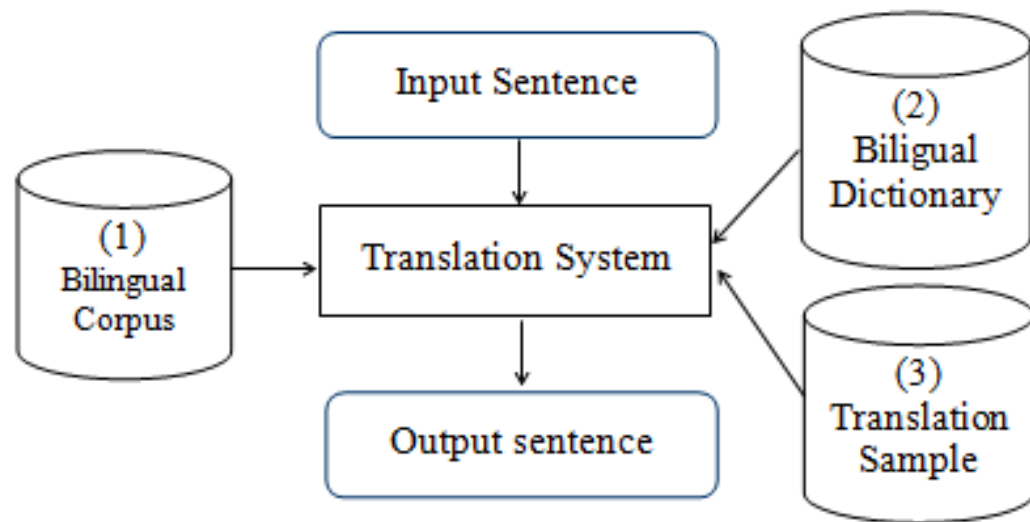
Step 2. Input Vietnamese sentence.

Step 3. Search for a Vietnamese pattern that matches the input of step 2.

Step 4. Output Cham Pattern corresponding to the Cham pattern made in Step 3.

Step 5. Generate a Cham sentence by using Vietnamese-Cham in Bilingual corpus, Bilingual dictionary and Cham pattern in Translation sample.

Next



**Figure 2:** Pattern-based Vietnamese Cham Machine Translation Model

## 4.2.1 Bilingual Corpus (cont.)

Bilingual corpus is a pairs of Vietnamese-Cham sentences.

Each pair of sentences are separated by a blank line.

Bilingual corpus are simple and standard sentences.

Corpus does not contain the text in confused sentences. (Dinh Dien, 2006)

No.	Bilingual Sentence
1.	Anh ở đâu ? អ្នកនៅ ឲ្យក្រី ក្រី ណាទី ?
	← blank line
2.	Tôi ở Sài Gòn . ខ្ញុំនៅ ឲ្យក្រី ក្រី នៅសៃកុង .

Where are you?

Abang daok pak halei?

**Figure 3.** Vietnamese- Cham Sample Sentence Construction

Next

## 4.2.2 Translation Sample Construction

+ Case sample sentence with a variable

No.	Sentence one variable	
1.	Anh học chữ <b>X</b> khi nào ? អ្នក ហ៊ុន អក្សរ <b>Y</b> នៅពេល ?	When did you learn <b>X</b> language? Abang bac akhar <b>Y</b> habier?
	“Anh học chữ <b>Cham</b> khi nào?” អ្នក ហ៊ុន អក្សរ <b>ចំ</b> នៅពេល?	(when did you learn <b>Cham</b> language?).

Translate sentence: “Anh học chữ **Cham** khi nào?” then the word "**Cham**" is replaced by "**X**" and its meaning “ចំ” is also replaced by "**Y**" and the translation result is:

អ្នក ហ៊ុន អក្សរ **ចំ** នៅពេល?

Next

## 4.2.2 Translation Sample Construction

+ Case sample sentence with a variable

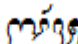
No.	Sentence one variable	
1.	Anh học chữ <b>X</b> khi nào ? អ្នក ហ៊ុន អក្សរ <b>Y</b> ពេលណា ?	When did you learn <b>X</b> language? Abang bac akhar <b>Y</b> habier?
2.	X nay anh ở đâu? Y នៅ ឯណា បច្ចុប្បន្ន ពេល ឥឡូវ	


**Figure 4.** Sample Sentence With a Variable

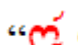
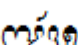
Sentence (1), translate sentence: “Anh học chữ **Cham** khi nào?” then the word “**Cham**” is replaced by “**X**” and its meaning “អក្សរ” is also replaced by “**Y**” and the translation result is: អ្នក ហ៊ុន អក្សរ **អ** ពេលណា? (when did you learn **Cham** language?).





Next

## Case sample sentence with many variable

Assuming exist sentence “X like Y”. In Vietnamese “**X1** thích **X2**”. and the corresponding translation “**Y1**  **Y2**”

To translate the sentence “**I** like **Cat**”, in Vietnamese “**Tôi** thích **mèo**” then “**Tôi**” is replaced by “**X1**” and its meaning “” was replaced by “**Y1**”.

Step 1: “  **Y2**”

Similarly, “**mèo**” was replaced by “**X2**” and its meaning “” was replaced by “**Y2**”. Meaning of complete sentences was “  ” (I like cat).

(V) Tôi muốn X1 một X2 (I want X1 a X2)  
 (C) ကို နှစ်ခု Y1 နှင့် Y2  
 (V1) Tôi muốn mua một X (I want buy a X)  
 (C1) ကို နှစ်ခု ဖြစ် နှင့် Y  
 (V2) Tôi muốn X một con mèo (I want X a Cat)  
 (C2) ကို နှစ်ခု Y နှင့် (ဖြစ်) မိလ္လာ

**Figure 5. Sample Sentence With Many Variable**

## 4.3 String Edit Distance

In order to extract the string, we use method of string edit distance (Computing Levenshtein distance) to measure distance between strings.

For turning a word (a) = “SPAKE” into (b)= “PARK”

		P	A	R	K
	c <sub>00</sub>	c <sub>01</sub>	c <sub>02</sub>	c <sub>03</sub>	c <sub>04</sub>
S	c <sub>10</sub>	c <sub>11</sub>	c <sub>12</sub>	c <sub>13</sub>	c <sub>14</sub>
P	c <sub>20</sub>	subst c <sub>21</sub>	delete c <sub>22</sub>	c <sub>23</sub>	c <sub>24</sub>
A	c <sub>30</sub>	insert c <sub>31</sub>	???		
K					
E					

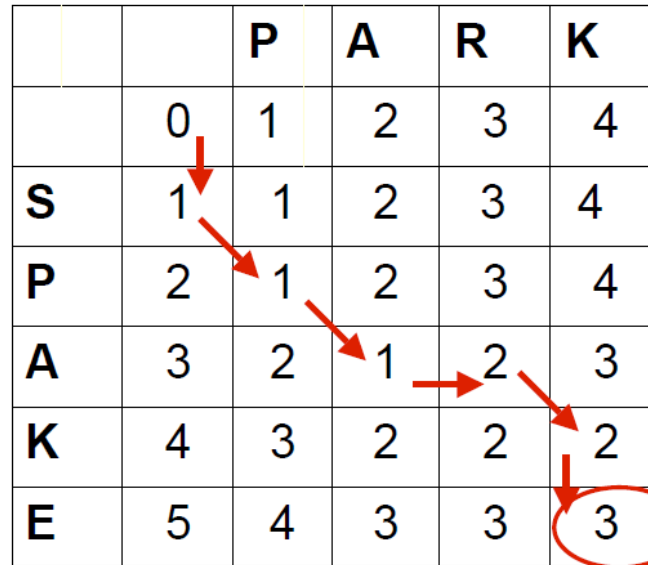
First delete “S”, second insert “R” and third delete “E” in sentence (a).

$$D_{ij} = \begin{cases} D(i,j) = \text{score of best alignment from } s1..s_i \text{ to } t1..t_j \\ \min \begin{cases} D(i-1, j-1) & , \text{ if } s_i=t_j & // \text{ Copy} \\ D(i-1, j-1) + 1, & \text{ if } s_i \neq t_j & // \text{ Substitute} \\ D(i-1, j) + 1 & // \text{ Insert} \\ D(i, j-1) + 1 & // \text{ Delete} \end{cases} \end{cases}$$

Next

## 4.3 String Edit Distance (cont.)

		P	A	R	K
	0	1	2	3	4
S	1	1	2	3	4
P	2	1	2	3	4
A	3	2	1	2	3
K	4	3	2	2	2
E	5	4	3	3	3



**Figure 7.** The Final Score of Both Strings Edit Distance

## 5. Result

The size of the bilingual dictionary is approximately 4,500 sentences, bilingual corpus is 2500 pairs of Vietnam Cham sentence, 1950 pairs of translation sample, 324 function words, 768 vocabulary, 57 prepositions, and all are stored in text file.

**Table 1.** Size of Data Warehouse in Testing

No.	Data	Size
1.	Bilingual dictionary	4,500
2.	Bilingual corpus	2500
3.	Pattern mapping	1950
4.	Function words	325
5.	Vocabulary	768
6.	Preposition	57



## 5. Result (cont.)

**Table 2.** Sample Sentence of Multiple Meaning

Sentence	Meaning
1. Vietnamese Cham	Tôi thích Mèo 𑜉𑜂𑜫 𑜉𑜂𑜫 𑜇𑜨𑜂𑜫 <i>I like Cat</i> (correct)
2. Vietnamese Cham	Tôi thích Mèo và Chó 𑜉𑜂𑜫 𑜉𑜂𑜫 𑜇𑜨𑜂𑜫 𑜉𑜂𑜫 𑜇𑜨𑜂𑜫 <i>I like Cat and Dog</i> (correct)
3. Vietnamese Cham	Cô ấy đẹp 𑜉𑜂𑜫 𑜉𑜂𑜫 𑜉𑜂𑜫 / 𑜉𑜂𑜫 𑜉𑜂𑜫 𑜉𑜂𑜫 <i>She is beautiful / She is nice female</i>
4. Vietnamese Cham	Nhà này đắt 𑜉𑜂𑜫 𑜉𑜂𑜫 𑜉𑜂𑜫 (correct) <i>This house is expensive</i> 𑜉𑜂𑜫 𑜉𑜂𑜫 𑜉𑜂𑜫 (incorrect) <i>This house is high</i>

## 6. Discussion

Based on Table 2 the results showed that, sentence (1) the sample sentence one variable “Tôi thích **X**” or “**X** thích Mèo”, or can be two variables as “**X** thích **Y**”. This sentence translation is correct because the input sentence “**Tôi** thích **Mèo**” exists in bilingual corpus or these words exist in bilingual dictionary.

Similar in sentence (2), this sentence translation is correct because it exists in bilingual corpus or these words exist in bilingual dictionary, and it can be generated into three variables: “**X** thích **Y** và **Z**”, or two variables “Tôi thích **X** và **Y**” and these samples can be converted into sample with one variable as “**X** thích Mèo và Chó”, “Tôi thích **X** và Chó”, or “Tôi thích Mèo và **X**”.



## 6. Discussion (cont.)

In our experiments, the characteristic of this translation is high quality translation results if the input sentence matches the translation pattern and this translation pattern is correct.

However, translation patterns and translation word dictionaries are usually made manually. Therefore, there are many costs in making a pattern-based machine translation system.

Anyhow, this machine translation for Vietnamese Cham has been important implication for studying, teaching, translating as well as preservation of Cham script and Cham language.



## 7. CONCLUSIONS

We proposed a technique pattern-based model to translate Vietnamese into Cham language. Initially we tested program based on a bilingual corpus, bilingual dictionaries and translation sample with limited resources, through observations the results achieved were relatively satisfying. Therefore, to develop a machine translation application for Vietnamese-Cham is necessary.

For future work, in order to make the application well translated, we are interested in building large enough resources, good quality for bilingual corpus, bilingual dictionaries and translation sample. Besides, we propose a new model for a MT system that will combine rule-based and example-based approach and will be applied to the Vietnamese - Cham translation.

Thanks for listening

Terima kasih

ကျေးဇူးတင်ပါတယ်

Next



Prof. Dr. Mohamad Bin Bilal Ali  
Dr. Noor Dayana Abd Halim  
PhD. Student. Mohamad Budi

University Technology Malaysia (UTM)- Malaysia

Back