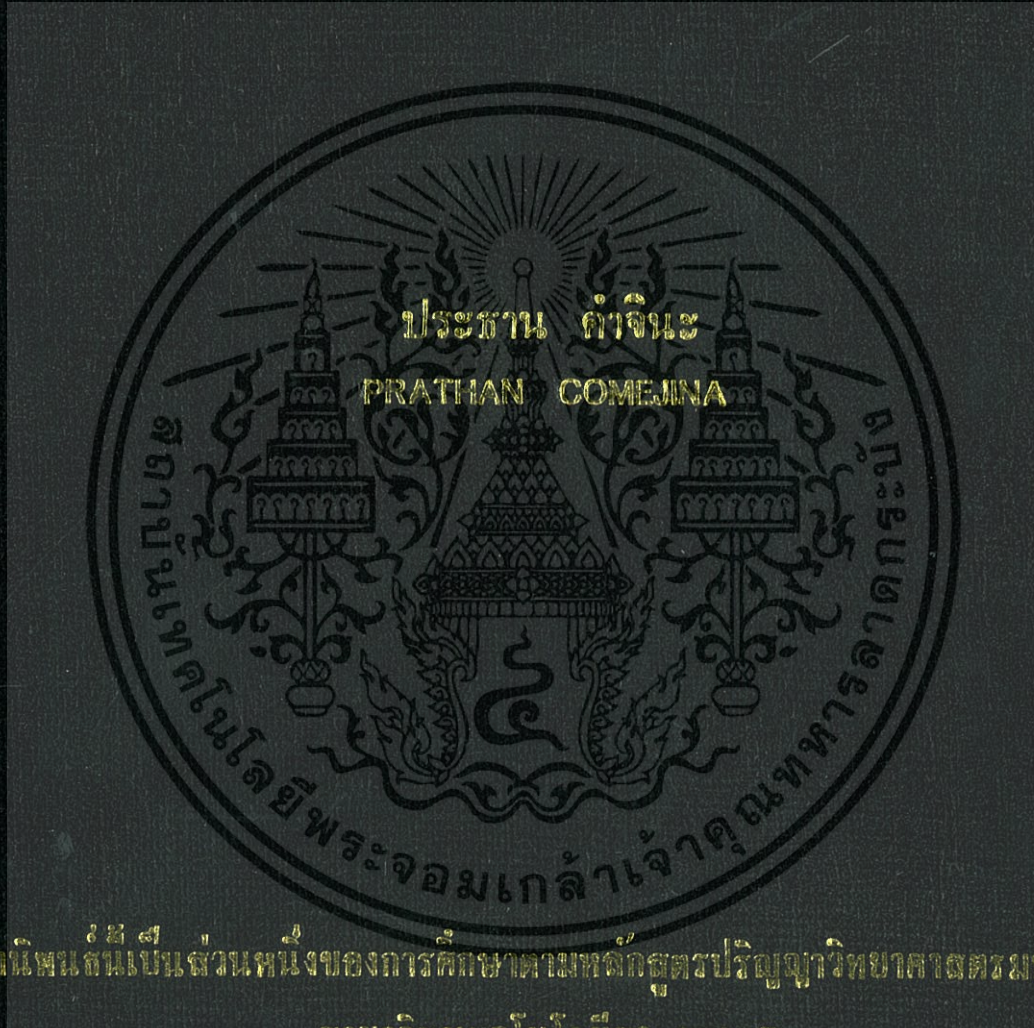


การแปลภาษาไทยเป็นภาษาไทยล้านนา

THAI TO THAI-LANNA MACHINE TRANSLATION



วิทยานิพนธ์นี้เป็นส่วนหนึ่งของการศึกษาค้นคว้าตามหลักสูตรปริญญาวิทยาศาสตรมหาบัณฑิต

สาขาวิชาเทคโนโลยีสารสนเทศ

คณะเทคโนโลยีสารสนเทศ

สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

พ.ศ. 2556

KMITL-2013-IT-M-001-002

ห้องสมุดคณะเทคโนโลยีสารสนเทศ พระจอมเกล้าลาดกระบัง

การแปลภาษาไทยเป็นภาษาไทยล้านนา

THAI TO THAI-LANNA MACHINE TRANSLATION



ประธาน คำจินะ

PRATHAN COMEJINA

๒๐๑๗ ๒๐๑๗ ๒๐๑๗

๗

เลขหมู่.....
เลขทะเบียน..... 7163
วัน,เดือน,ปี..... 17 ต.ค. 2556

b. ๗๒๕๕๙๐๕๕
i.

วิทยานิพนธ์นี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญาวิทยาศาสตรมหาบัณฑิต

สาขาวิชาเทคโนโลยีสารสนเทศ

คณะเทคโนโลยีสารสนเทศ

สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

พ.ศ. 2556

KMITL-2013-IT-M-001-002

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

THAI TO THAI-LANNA MACHINE TRANSLATION



**A THESIS SUBMITTED IN PATIAL FULFILLMENT
OF THE REQUIREMENT FOR THE DEGREE OF
MASTER OF SCIENCE IN INFORMATION TECHNOLOGY
FACULTY OF INFORMATION TECHNOLOGY
KING MONGKUT'S INSTITUTE OF TECHNOLOGY LADKRABANG**

2013

KMITL-2013-IT-M-001-002

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



COPYRIGHT 2013


FACULTY OF INFORMATION TECHNOLOGY

KING MONGKUT'S INSTITUTE OF TECHNOLOGY LADKRABANG

เอกสารนี้เป็นเอกสารลิขสิทธิ์สงวนไว้สำหรับใช้เพื่อการศึกษาเท่านั้น เมื่อเผยแพร่เห็นไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

คณะเทคโนโลยีสารสนเทศ
สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง
ใบรับรองวิทยานิพนธ์

หัวข้อวิทยานิพนธ์ การแปลภาษาไทยเป็นภาษาไทยล้านนา
Thai to Thai-lanna Machine Translation
นักศึกษา นายประธาน คำจិនะ
รหัสประจำตัว 51066409
ปริญญา วิทยาศาสตรมหาบัณฑิต
สาขาวิชา เทคโนโลยีสารสนเทศ
อาจารย์ที่ปรึกษาวิทยานิพนธ์ รองศาสตราจารย์ ดร.พรฤดี เนติโสภาคกุล

| คณะกรรมการสอบวิทยานิพนธ์ | ลายมือชื่อ |
|--|---|
| รองศาสตราจารย์ ดร.วราภรณ์ กริสรุเดช |  |
| รองศาสตราจารย์ ดร.นุชรี เปรมชัยสวัสดิ์ | |
| รองศาสตราจารย์ ดร.อาริต ธรรมโน | |
| รองศาสตราจารย์ ดร.พรฤดี เนติโสภาคกุล | |
| ผู้ช่วยศาสตราจารย์ ดร.กันต์พงษ์ วรรัตน์ปัญญา | |

สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

KING MONGKUT'S INSTITUTE OF TECHNOLOGY LADKRABANG

วัน/เดือน/ปี ที่สอบ วันอังคารที่ 9 เมษายน 2556 เวลา 09.30 น.

สถานที่สอบ ณ ห้อง 335 (ชั้น 3) คณะเทคโนโลยีสารสนเทศ

คณะเทคโนโลยีสารสนเทศรับรองแล้ว



(รองศาสตราจารย์ ดร.กันทรบูรณ์ สัตตวิริยวงศ์)

คณบดีคณะเทคโนโลยีสารสนเทศ

วันที่...14...เดือน...พฤษภาคม...พ.ศ. 2556

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

| | |
|------------------|--------------------------------|
| ชื่อวิทยานิพนธ์ | การแปลภาษาไทยเป็นภาษาไทยล้านนา |
| นักศึกษา | นายประธาน คำจันะ |
| รหัสนักศึกษา | 51066409 |
| ปริญญา | วิทยาศาสตรมหาบัณฑิต |
| สาขาวิชา | เทคโนโลยีสารสนเทศ |
| แขนงวิชา | วิทยาการสารสนเทศ |
| พ.ศ. | 2556 |
| อาจารย์ที่ปรึกษา | รศ.ดร.พรฤดี เนติโสภาคกุล |

บทคัดย่อ

วิทยานิพนธ์ฉบับนี้นำเสนอกระบวนการแปลภาษาไทยเป็นภาษาไทยล้านนาในระดับประโยค มีกระบวนการประมวลผลอยู่ 3 ขั้นตอนหลักได้แก่ ขั้นตอนก่อนการประมวลผล ขั้นตอนการแปล ประโยคไทยล้านนาและขั้นตอนการสร้างประโยคไทยล้านนา ในขั้นตอนก่อนการประมวลผลจะ ประกอบด้วยขั้นตอนย่อย 2 ขั้นตอนคือ การตัดคำในงานวิจัยนี้ใช้เครื่องมือ KU-CUT ของ มหาวิทยาลัยเกษตรศาสตร์และขั้นตอนวิเคราะห์โครงสร้างไวยากรณ์ไม่พึงบริบทประยุกต์เทคนิค ทางสถิติ(PCFG) ขั้นตอนการแปลประโยคไทยล้านนาในขั้นตอนนี้จะเป็นการนำคำจากขั้นตอน ก่อนการประมวลผลไปค้นหาในพจนานุกรมหากไม่พบคำแปลจะทำการแปลงรูปคำไทยล้านนา ตามกฎการเขียนคำไทยล้านนา ในขั้นตอนสุดท้ายการสร้างประโยคไทยล้านนา เป็นการปรับ โครงสร้างประโยคไทยให้เป็นโครงสร้างประโยคไทยล้านนา โดยผ่านไวยากรณ์ปริวรรตเพิ่มพูน และลดความกำกวมของโครงสร้างด้วยวิเทอร์บีอัลกอริทึม เมื่อผ่านกระบวนการทั้งหมดแล้วจะได้ ประโยคไทยล้านนาที่เรียงลำดับตามความน่าจะเป็นของการเกิดประโยค สิ่งที่ผู้วิจัยได้จัดสร้าง ได้ได้แก่ (1) พจนานุกรมอิเล็กทรอนิกส์ไทย-ไทยล้านนา (2) เครื่องจักรวิเคราะห์ไวยากรณ์ไม่พึง บริบท บริบทประยุกต์เทคนิคทางสถิติ (3) กฎการแปลงรูปภาษาไทยเป็นภาษาไทยล้านนา (4) กฎ ไวยากรณ์ปริวรรตเพิ่มพูนเพื่อจัดโครงสร้างประโยคไทยล้านนา สุดท้ายเรียงลำดับคำแปลที่น่าจะ เป็นตามลำดับเมื่อคำนวณจากวิเทอร์บีอัลกอริทึม การทดสอบประสิทธิภาพในการแปลโดยใช้นิ ทิตานพื้นบ้านและประโยคทดสอบจำนวน 206 ประโยค นำผลลัพธ์มาเปรียบเทียบกับประโยคไทย ล้านนาจากผู้เชี่ยวชาญภาษาไทยล้านนา ปรากฏว่ามีประโยคที่แปลถูกต้อง 162 ประโยคจาก 206 ประโยคคิดเป็นอัตราส่วนร้อยละ 78.64

| | |
|-----------------------|--|
| Thesis | Thai to Thai-lanna Machine Translation |
| Student | Mr.Prathan Comejina |
| Student ID. | 51066409 |
| Degree | Master of Science |
| Program | Information Technology |
| Major | Information Science |
| Year | 2013 |
| Thesis Advisor | Assoc.Prof. Dr. Pornrudee Netisopakul |

ABSTRACT

This thesis proposes to translate Thai language to Thai-lanna language in a sentence level. There are three main natural language processes: preprocessing, mapping Thai to Thai-lanna, and Thai-lanna syntactic analysis. Preprocessing steps consist of two sub-steps, Thai word segmentation using a KUCUT tool and Thai sentence analysis using Probabilistic Context Free Grammar (PCFG). Mapping steps translate Thai words to Thai-lanna words using Thai-Thai-lanna dictionary. If a word is not found in the dictionary, we use Thai-lanna written rules to directly translate the Thai word. The final process constructs a Thai-lanna sentence using Thai-Thai lanna transformation rule. If ambiguity occurs – there can be more than one sentence structure- we sort the sentences based on their probabilities given by a viterbi algorithm. Our contributions are construction of (1) Thai-Thai lanna electronic dictionary in a trie structure, (2) a Thai PCFG grammar for parsing Thai sentences, (3) rules for Thai-Thai lanna word translation, (4) transformation rules for Thai-Thai lanna sentence use during Thai-lanna sentence generation, and finally, implementation of viterbi algorithm to sort a set of possible Thai-lanna sentences. We use our system to translate three folk stories and a set of general sentences, totally 206 sentences. The first sentence in the sorted translated list is judged by a Thai lanna expert to be translated correctly 162 sentences out of 206 sentences, or about 78.64%.

กิตติกรรมประกาศ

วิทยานิพนธ์ฉบับนี้สำเร็จลุล่วงไปได้ด้วยดีด้วยความกรุณาจาก รศ.ดร.พรฤดี เนติโสภากุล อาจารย์ที่ปรึกษาวิทยานิพนธ์ที่ได้ให้คำแนะนำ แนวคิด ตลอดจนแก้ไขข้อบกพร่องต่างๆ เป็นอย่างดี จนวิทยานิพนธ์ฉบับนี้เสร็จสมบูรณ์ ซึ่งข้าพเจ้าขอขอบพระคุณเป็นอย่างสูง

ขอขอบพระคุณ บิดา มารดาผู้ซึ่งให้ความรักความเมตตาความห่วงใยและเป็นกำลังใจให้กับผู้จัดทำวิทยานิพนธ์ฉบับนี้จนสำเร็จ

ขอขอบพระคุณ คณาจารย์คณะเทคโนโลยีสารสนเทศ สถาบันพระจอมเกล้าเจ้าคุณทหารลาดกระบังทุกๆ ท่านที่ได้ประสิทธิ์ประสาทวิชาการตลอดระยะเวลาที่ได้ศึกษาในสถาบันแห่งนี้ รวมถึงคณะกรรมการที่ดำเนินการสอบที่ได้ให้คำแนะนำและคำชี้แนะในการปรับปรุงวิทยานิพนธ์ฉบับนี้ให้มีความสมบูรณ์

ขอขอบคุณพี่ๆ สมาชิกหน่วยปฏิบัติการการจัดการองค์ความรู้และวิศวกรรมความรู้ (KMAKE Lab) ทุกคนที่ได้ให้กำลังใจและช่วยเหลือในทุกด้านในการทำวิทยานิพนธ์ ขอขอบคุณคุณกฤตดาพร พัทธระสุภา ที่ได้ให้ความช่วยเหลือให้คำแนะนำองค์ความรู้ในเรื่องการประมวลผลภาษาธรรมชาติเป็นอย่างดีในการทำวิทยานิพนธ์ฉบับนี้

ขอบคุณสาขาวิชาคอมพิวเตอร์ มหาวิทยาลัยราชภัฏเชียงใหม่ ที่ได้ให้เวลาในการทำวิทยานิพนธ์ รวมถึงเพื่อนร่วมงานทุกท่านที่ได้ให้กำลังใจเป็นอย่างดี

ความดีของการศึกษาค้นคว้าครั้งนี้ขอมอบเป็นเครื่องบูชาบิดา มารดาและบูรพาจารย์ทุกท่าน ผู้เขียนมีความซาบซึ้งในความกรุณาอันดีจากทุกท่านที่ได้กล่าวนามมาและขอขอบพระคุณมา ณ โอกาสนี้

ประธาน คำจិនะ

สารบัญ

| | หน้า |
|---|------|
| บทคัดย่อภาษาไทย..... | I |
| บทคัดย่อภาษาอังกฤษ..... | II |
| กิตติกรรมประกาศ..... | III |
| สารบัญ..... | IV |
| สารบัญตาราง..... | VI |
| สารบัญรูป..... | VII |
| บทที่ 1 บทนำ | |
| 1.1 ความเป็นมาและความสำคัญของปัญหา..... | 1 |
| 1.2 วัตถุประสงค์ของการศึกษา..... | 2 |
| 1.3 ทฤษฎีหรือแนวคิดที่ใช้ในงานวิจัย..... | 2 |
| 1.4 ประโยชน์ที่คาดว่าจะได้รับจากงานวิจัย..... | 2 |
| 1.5 ขอบเขตของงานวิจัย..... | 2 |
| 1.6 ขั้นตอนของการศึกษา..... | 3 |
| บทที่ 2 ทฤษฎีและงานวิจัยที่เกี่ยวข้อง | |
| 2.1 พื้นฐานของภาษาไทยล้านนา..... | 4 |
| 2.1.1 รูปพยัญชนะ..... | 4 |
| 2.1.2 รูปสระ..... | 5 |
| 2.1.3 รูปวรรณยุกต์และสัญลักษณ์พิเศษ..... | 6 |
| 2.1.4 พยัญชนะตัวสะกด..... | 6 |
| 2.1.5 หลักการเขียนภาษาไทยล้านนา..... | 7 |
| 2.2 การแปลภาษา..... | 8 |
| 2.2.1 การแปลในระดับคำ..... | 8 |
| 2.2.2 การแปลในระดับวลี..... | 12 |
| 2.2.3 การแปลในระดับประโยค..... | 16 |

สารบัญ(ต่อ)

| | หน้า |
|---|------|
| บทที่ 3 วิธีดำเนินงานวิจัย | |
| 3.1 สถาปัตยกรรมของการแปลภาษาไทยเป็นภาษาไทยล้านนา..... | 21 |
| 3.2 ไวยากรณ์ไม่พึงบริบทใช้เทคนิคทางสถิติ..... | 23 |
| 3.3 โครงสร้างพจนานุกรมไทย-ล้านนา..... | 25 |
| 3.4 การแปลงรูปภาษาไทยล้านนา..... | 27 |
| 3.5 ไวยากรณ์ปริวรรติเพิ่มพูน..... | 29 |
| 3.6 แบบจำลองฮิดเดินมาร์คอฟและการลดความกำกวม โดยวิเทอร์บีอัลกอริทึม..... | 30 |
| 3.7 การทดลองและการวัดประเมินผลการทดลอง..... | 32 |
| บทที่ 4 วิธีดำเนินงานวิจัย | |
| 4.1 ผลการวิจัย..... | 33 |
| 4.2 การอภิปรายผล..... | 34 |
| 4.3 อุปสรรคในการทำงาน..... | 35 |
| บทที่ 5 วิธีดำเนินงานวิจัย | |
| 5.1 สรุปผลการวิจัย..... | 37 |
| 5.2 ข้อเสนอแนะ..... | 38 |
| บรรณานุกรม..... | 39 |
| ภาคผนวก..... | 40 |
| ภาคผนวก ก. ตัวอย่างผลลัพธ์จากการทดลอง..... | 41 |
| ภาคผนวก ข. บทความและผลงานวิจัยที่ได้รับการตีพิมพ์..... | 44 |
| ประวัติผู้เขียน..... | 56 |

สารบัญตาราง

| ตารางที่ | หน้า |
|---|------|
| 3.1 ตารางโครงสร้างนามวลี..... | 23 |
| 3.2 ตารางโครงสร้างกริยาวลี..... | 24 |
| 3.3 ตารางโครงสร้างพิเศษวลี..... | 24 |
| 3.4 ตารางโครงสร้างกาลวิเศษณ์วลี..... | 24 |
| 3.5 โครงสร้างสถานวิเศษณ์วลี..... | 25 |
| 3.6 ประเภทของคำ..... | 26 |
| 3.7 ตารางตัวอย่าง ของกฎของกฎการแปลรูปภาษาไทยเป็นภาษาล้านนา..... | 28 |
| 4.1 ตารางประสิทธิภาพของการแปลภาษาไทยเป็นภาษาล้านนา..... | 34 |



สารบัญรูป

| รูปที่ | หน้า |
|--------|---|
| 2.1 | สถาปัตยกรรมในการแปลภาษาไทยด้านนาเป็นภาษาไทย.....10 |
| 2.2 | สถาปัตยกรรมของการแปลอักษรธรรมอีสานให้เป็นภาษาไทย.....11 |
| 2.3 | โครงสร้างวลีของนามวลี.....12 |
| 2.4 | โครงสร้างวลีของกริยาวลี.....13 |
| 2.5 | โครงสร้างวลีของสถานวิเศษณ์วลี.....13 |
| 2.6 | โครงสร้างวลีของสถานกาลวิเศษวลี.....13 |
| 2.7 | โครงสร้างวลีของสถานพิเศษวิเศษณ์วลี.....14 |
| 2.8 | โครงสร้างวลีของอาลปนวลี.....14 |
| 2.9 | รายงานความสัมพันธ์ของส่วนมูลฐานในประโยค.....14 |
| 2.10 | รายงานความสัมพันธ์ของส่วนเสริมในประโยค.....15 |
| 2.11 | รายงานความสัมพันธ์ของความสัมพันธ์ของส่วนมูลฐานและส่วนเสริม.....15 |
| 2.12 | การแจงประโยคจากบนลงล่าง (Top-down parsing).....16 |
| 2.13 | การแจงประโยคจากล่างขึ้นบน (Bottom-up parsing).....17 |
| 2.14 | ขั้นตอนการแปลภาษาอาหรับเป็นภาษาอังกฤษ.....17 |
| 2.15 | ฟังก์ชันวิเทอร์บี สำหรับหาโครงสร้างประโยคภาษาอังกฤษ.....19 |
| 3.1 | สถาปัตยกรรมของการแปลภาษาไทยเป็นภาษาไทยด้านนา.....22 |
| 3.2 | ตัวอย่างโครงสร้างข้อมูลของพจนานุกรมไทย-ไทยด้านนา.....26 |
| 3.3 | วิเทอร์บีอัลกอริทึมสำหรับการลดคำกำกวม ในการแปลภาษาไทยเป็นภาษาไทยด้านนา.....31 |

บทที่ 1

บทนำ

1.1 ความเป็นมาและความสำคัญของปัญหา

ภาษาเป็นวิธีการในการสื่อสารของมนุษย์อีกวิธีการหนึ่ง ซึ่งสามารถสื่อสารกันโดยการเรียนรู้ ความหมายและความต้องการของผู้สื่อสารจาก เสียงและสัญลักษณ์ที่แสดงหรือตัวอักษรนั่นเอง ในแต่ละชนชาติบนโลกนี้ก็มีวิธีสื่อสารกันซึ่งมีความแตกต่างกันตามสังคม วัฒนธรรม รวมถึงภูมิประเทศ ประเทศไทยก็เป็นประเทศหนึ่งที่มีความหลากหลายทางวัฒนธรรม แต่ละภูมิภาคก็มีภาษาถิ่นที่มีเอกลักษณ์ นอกเหนือจากภาษาไทยกลาง ที่เป็นภาษาราชการของประเทศไทย นอกจากนี้ภาษาถิ่นแต่ละภาคก็มีความแตกต่างทางด้าน โทณเสียงและสัญลักษณ์ ภาษาถิ่นเหล่านั้นเป็นภาษาที่แสดงออกถึงภูมิปัญญาของบรรพบุรุษ ที่จะใช้ตัวอักษรเหล่านั้น จารึกหรือบันทึก ข้อมูลที่สำคัญๆ ของอดีตจนถึงปัจจุบันเอาไว้

ภาษาไทยล้านนาเป็นภาษาถิ่นที่ใช้กันอย่างแพร่หลายในภาคเหนือตอนบนของประเทศไทย ภาษาไทยล้านนาเป็นภาษาที่มีลักษณะพิเศษเฉพาะตัว มีต้นกำเนิดมาจากภาษาขอม ภาษาไทยล้านนานั้นจะมีทั้งภาษาพูดและภาษาเขียน คนเมืองหรือคนล้านนาจะพูดและเขียนภาษาคำเมือง , ตัวธรรมหรือภาษาไทยล้านนาแล้วแต่จะเรียก ซึ่งเป็นภาษาหลักในการสื่อสาร ในการบันทึกเรื่องราว เหตุการณ์ต่างๆ ตำรายาสมุนไพร หรือแม้กระทั่งคัมภีร์โบราณ กัมพูเทศ จารึก ก็ใช้อักษรไทย ล้านนาเป็นอักษรในการสื่อถึงเหตุการณ์เหล่านั้น ในปัจจุบันเมื่อมีความก้าวหน้าทางสังคมและเทคโนโลยี รวมไปถึงหลักสูตรทางการเรียนการสอนที่เป็นมาตรฐานของไทย ทำให้คนเมืองหรือคนล้านนารุ่นหลัง ได้ปรับเปลี่ยนวิถีชีวิต ภาษาคำเมืองหรือภาษาไทยล้านนา ก็เริ่มถูกกลืนหายไป เหลือเพียงบางพื้นที่ยังใช้ภาษาพูด ในการติดต่อสื่อสารกัน ส่วนการเรียนรู้ตัวอักษร หรือภาษาเขียน นั้น จะมีแต่ในวงแคบเช่น วัดหรือในบางสถานศึกษาเท่านั้นที่มีการจัดการเรียนการสอนภาษาไทย ล้านนา จึงควรค่าสืบสานไม่ให้สูญหาย

การแปลภาษาด้วยเครื่องคอมพิวเตอร์เป็นเทคโนโลยีหนึ่งที่ใช้โปรแกรมทางคอมพิวเตอร์แปล จากภาษาหนึ่งไปเป็นภาษาเป้าหมาย ซึ่งในการแปลภาษานั้นมีวิธีการแปลได้หลายแบบ เช่น การแปลแบบคำต่อคำ การแปลแบบวิเคราะห์โครงสร้างไวยากรณ์ และการแปลโดยใช้ภาษากลาง การแปลแต่ละแบบก็สามารถประยุกต์เข้ากับแต่ละข้อจำกัดของภาษา

จากเหตุผลข้างต้นจึงทำให้เกิดแนวความคิดงานวิจัย การแปลภาษาไทยเป็นภาษาไทยล้านนา โดยวิเคราะห์โครงสร้างของไวยากรณ์ของภาษาไทยล้านนา และค้นหาคำศัพท์ในพจนานุกรมไทย-ล้านนา ที่สร้างขึ้นแล้วก็จะลดความกำกวมของประโยคด้วยวิธีการของแบบจำลองฮิดเดนมาร์คอฟ ใช้วีเทอร์บีอัลกอริทึม และใช้ความน่าจะเป็นจากคลังประโยคไทยล้านนา

1.2 วัตถุประสงค์ของการศึกษา

- 1.2.1 เพื่ออนุรักษ์ภาษาไทยล้านนาที่ปัจจุบันมีผู้เชี่ยวชาญในภาษาไทยล้านนาลดน้อยลงไม่ให้สูญหาย
- 1.2.2 เพื่อศึกษาและพัฒนาแนวทางในการเลือกวิธีการแปลภาษา เพื่อใช้ในการแปลภาษาที่มีความใกล้เคียงกัน
- 1.2.3 เพื่อพัฒนาต้นแบบ โปรแกรมประยุกต์ในการแปลภาษาไทยเป็นภาษาไทยล้านนา

1.3 ทฤษฎีหรือแนวคิดที่ใช้ในการวิจัย

- 1.3.1 วิธีการปริวรรตภาษาไทยล้านนาเป็นการเขียนและวางตำแหน่งตัวอักษรไทยล้านนา ตามรูปแบบ โดยจะมีรูปแบบตามวิธีการเขียนภาษาไทยล้านนาของ จ.เชียงใหม่ ซึ่งช่วยในการแปลงรูปของคำภาษาไทยที่ค้นหาในพจนานุกรมแล้วไม่เจอคำดังกล่าว ซึ่งจะทำให้ได้รูปของคำภาษาไทยล้านนา
- 1.3.2 การวิเคราะห์ชนิดของคำภาษาไทย เพื่อนำไปวิเคราะห์โครงสร้างไวยากรณ์และเปรียบเทียบโครงสร้างกับไวยากรณ์ภาษาไทยล้านนา และใช้ในการหาคำศัพท์ไทยล้านนาที่ตรงกับประเภทของคำไทยที่ต้องการ
- 1.3.3 การวิเคราะห์โครงสร้างไวยากรณ์ภาษาไทยด้วย PCFG
- 1.3.4 แบบจำลองฮิดเดนมาร์คอฟโมเดลและเวกเตอร์บี้อัลกอริทึม จะใช้ในการคำนวณลำดับของประโยคที่มีความน่าจะเป็นสูงสุด

1.4 ประโยชน์ที่คาดว่าจะได้รับจากงานวิจัย

- 1.4.1 ได้ระบบการแปลภาษาไทยเป็นภาษาไทยล้านนาเพื่อสนับสนุนการเรียนรู้อาษาไทยล้านนาสำหรับผู้สนใจ
- 1.4.2 สามารถเผยแพร่ภูมิปัญญาท้องถิ่นและอนุรักษ์วัฒนธรรมไทย ออกไปสู่กว้างซึ่งจะช่วยให้องค์ความรู้ทางด้านภาษาไทยล้านนาถูกสืบสานต่อไป
- 1.4.3 ได้แนวทางในการแปลภาษาไทยเป็นภาษาอื่นที่มีลักษณะ ข้อจำกัดของภาษาที่มีความใกล้เคียงกันกับภาษาไทยล้านนา

1.5 ขอบเขตการวิจัย

- 1.5.1 พจนานุกรมที่ใช้ในงานวิจัยคือ พจนานุกรมล้านนา-ไทย ฉบับเฉลิมพระเกียรติพระบาทสมเด็จพระเจ้าอยู่หัว ในวโรกาสเจริญพระชนมายุ 80 พรรษา ของสถาบันภาษาและศิลปวัฒนธรรม มหาวิทยาลัยราชภัฏเชียงใหม่

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

1.5.2 การแปลภาษาไทยล้านนาจะอยู่ในระดับประโยค

1.5.3 ในขั้นตอนการลดความกำกวม โดยใช้วีเทอร์บีอัลกอริทึม จะใช้คลังข้อมูลภาษาที่เป็นภาษาไทยล้านนา จากประโยคภาษาไทยล้านนาที่รวบรวมได้

1.6 ขั้นตอนของการศึกษา

ขั้นตอนในการศึกษาวิทยานิพนธ์มีขั้นตอนดังต่อไปนี้

1.6.1 ศึกษาทฤษฎีและงานวิจัยที่เกี่ยวข้อง

1.6.2 ออกแบบขั้นตอนการแปลภาษาไทยเป็นภาษาไทยล้านนา

1.6.3 สร้างพจนานุกรมไทย-ล้านนา จากพจนานุกรมล้านนา-ไทย ฉบับเฉลิมพระเกียรติพระบาท สมเด็จพระเจ้าอยู่หัว ในวโรกาสเจริญพระชนมายุ 80 พรรษา ของสถาบันภาษาและศิลปวัฒนธรรม มหาวิทยาลัยราชภัฏเชียงใหม่

1.6.4 เตรียมข้อมูลคลังภาษาไทยล้านนา เพื่อใช้ในการเรียนรู้ในการลดความกำกวมโดยใช้วีเทอร์บีอัลกอริทึม

1.6.5 สร้างเครื่องมือต้นแบบ การแปลภาษาไทยเป็นภาษาไทยล้านนา

1.6.6 ทดลองการแปล เก็บผลลัพธ์ของการแปล และให้ผู้เชี่ยวชาญช่วยตรวจสอบผลลัพธ์ที่ถูกต้องของคำแปลไทยล้านนาที่ได้จากการแปล

1.6.7 วิเคราะห์ผลลัพธ์และสรุปผลการทดลอง

บทที่ 2

ทฤษฎีและงานวิจัยที่เกี่ยวข้อง

ในบทนี้จะกล่าวถึงทฤษฎีพื้นฐานที่ใช้ในการวิจัย และงานวิจัยที่เกี่ยวข้องซึ่งจะเป็นพื้นฐานของงานวิจัยการแปลภาษาไทยเป็นภาษาไทยล้านนา ทฤษฎีที่เกี่ยวข้องกับงานวิจัย ได้แก่ ภาษาไทย ล้านนา งานวิจัยที่เกี่ยวข้อง

2.1 พื้นฐานของภาษาไทยล้านนา [1]

ในการศึกษาภาษาไทยล้านนานั้น จะต้องทราบถึงรูปของตัวอักษร ก่อนเพื่อนที่จะเป็นพื้นฐานในการเขียนและการอ่าน รูปของภาษาไทยล้านนาสามารถแบ่งออกได้ดังนี้

2.1.1 รูปพยัญชนะ อักษรไทยล้านนา หรือตัวเมือง จัดตามกลุ่มภาษาบาลีแบ่งออกเป็น 5วรรค วรรคละ 5 ตัวอักษร ซึ่งเรียกว่า พยัญชนะวรรค หรือ พยัญชนะในวรรค และมีอีก 8 ตัวอักษรที่ไม่จัดอยู่ในวรรค เรียกว่า พยัญชนะอวรรค หรือ พยัญชนะนอกวรรค

พยัญชนะในวรรค มีทั้งหมด 5 วรรค 25 คือ

| | | | | | |
|------------|-------|-------|-----|-------|-------|
| กะ วรรค | กะ | ขะ | ก๊ะ | ฆะ | งะ |
| | ก | ข | ค | ฆ | ง |
| จะ วรรค | จะ | ฉะ | จ๊ะ | ฉะ | ญะ |
| | จ | ฉ | ค | ช | ญ |
| ระ ฐะ วรรค | ระ ฐะ | ระ ฐะ | คะ | ระ ฒะ | ระ ณะ |
| | ร | ฐ | ค | ฆ | ณ |
| ตะ วรรค | ตะ | ปะ | ติ้ | ธะ | นะ |
| | ต | ด | ด | ด | ด |
| ปะ วรรค | ปะ | ผะ | ป๊ะ | ภะ | มะ |
| | ป | ผ | บ | บ | ม |

พยัญชนะนอกวรรค มีทั้งหมด 8 ตัวอักษร ได้แก่

| | | | | | | | |
|----|----|----|----|----|----|----|-----|
| ยะ | ระ | ละ | วะ | สะ | หะ | พะ | อ้ง |
| ย | ร | ล | ว | ส | ฮ | ป | ฮ |

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

พยัญชนะเพิ่ม พยัญชนะยังไม่ครบตามสำเนียงภาษาพูด คนโบราณเลยเพิ่มพยัญชนะอีก 10 ตัว
คือ

| | | | | |
|-----|-----|----|---------------|----|
| คะ | ชะ | บะ | ยะ (อยะ, หยะ) | ฝะ |
| โคะ | โชะ | โบ | โย | ฝ |
| พะ | ศะ | ษะ | สะสองห้อง | นา |
| โคะ | โศ | โษ | งว | ณ |

พยัญชนะพิเศษ เป็นพยัญชนะที่มีใช้บ้างแต่ก็เป็นบางคำเท่านั้น ไม่ค่อยได้ใช้บ่อยนัก

- ๑. ตัวปะหลวง มีค่าเท่ากับตัว พ และ พพ เช่น สพพ์(๑๑๑)
- ๒. ตัวรือ, ฤ เช่น ฤกษ์(๑๑๑)
- ๓. ะโอง เท่ากับตัว ร ควบกล้ำ เช่น พระ(๑๑) พรหมา(๑๑๑)
- ๔. ไม้กะ (สระอะ) บางครั้งเท่ากับตัว ก สะกด เรียกว่า กะปุญาต เช่น (๑๑)

2.1.2 รูปสระ สระในภาษาสันสกฤตแบ่งออกเป็นหลายประเภท คือ สระตามภาษาบาลี สระ
เดี่ยว และสระผสม ซึ่งต่างก็มีลักษณะการใช้งานและหน้าที่แตกต่างกันไป
สระที่มาจากภาษาบาลี

| | | | | | | | |
|----|----|----|----|----|----|----|---|
| อะ | อา | อิ | อี | อุ | อู | เอ | เ |
| อ | อา | อิ | อี | อุ | อู | เอ | เ |

สระเดี่ยว เป็นสระแท้ (ไม่สามารถออกเสียงได้ด้วยตัวเอง) เรียกว่า “ไม้” ใช้ผสมกับพยัญชนะมี
14 ตัว

| | | | | | |
|----|-------|-------|--------|-----------------------|-------|
| -ะ | -า | ไม้กะ | -า | -อ | ไม้กั |
| -ิ | -ี | ไม้กิ | -า | -อ | ไม้กั |
| -ุ | -ู | ไม้กู | -า | -อ | ไม้กั |
| -ึ | -ุ | ไม้กึ | -า | -อ | ไม้กั |
| -เ | -เ | ไม้เก | -แ | -เ | ไม้เก |
| -เ | -เ/-เ | ไม้เก | -เ | -เ | ไม้เก |
| -อ | -อ | ไม้กั | นิคหิต | ไม้กั, ไม้กัมน (บาลี) | |

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สระผสม เป็นการนำเอาสระแท้ หรือสระเดี่ยวมาผสมรวมกัน เรียกว่าการผสมไม้ ตัวอย่างเช่น

| | | | |
|------|---------|------|----------------|
| ๑๓๙ | ไม้เาะ | ๑๓๙๙ | ไม้เาะ |
| ๑๔๐ | ไม้โอะ | ๑๓๖ | ไม้โถ้(บาลี) |
| ๑๓๖ | ไม้เก๋า | ๑๓๘ | ไม้เาะ |
| ๑๓๗ | ไม้ก้อ | ๑๓๙ | ไม้กัวะ |
| ๑๓๘ | ไม้กัว | ๑๓๙๙ | ไม้เก็ยะ |
| ๑๓๙ | ไม้เก็ย | ๑๓๙๙ | ไม้เก็ยะ |
| ๑๓๙๙ | ไม้เก็ย | ๑๓๙๙ | ไม้เาะ |
| ๑๓๙ | ไม้เก็ย | ๑๓๙๙ | ไม้เก็ย |
| ๑๓๙๙ | ไม้เก็ย | ๑๓๙๙ | ไม้เก็ย (บาลี) |

2.1.3 รูปวรรณยุกต์/สัญลักษณ์พิเศษ อดีตไม่มีวรรณยุกต์เหล่านี้ปรากฏซึ่งในโบราณเก่า จะพบบ่อยเมื่อสูญปลายโบราณ

| | | |
|---|---------------------------------|------------------|
| ๑ | ไม้เอก | ไม้หยัก, ไม้เหาะ |
| ๒ | ไม้โท | ไม้ขอช้าง |
| ๓ | ไม้หันอากาศ | ไม้ซัด |
| ๔ | การันต์ | ระห้าม |
| ๕ | ไม้ซำคำ | ไม้สองน้อย |
| ๖ | ไม้ไตคู่ | ไม้ไตคู่ |
| ๗ | ไม้กั้ง, ตั้วข่ม, ไม้โถ้(ลครูป) | ไม้กั้งไหล |
| ๘ | | ไม้เก๋าจู้ |
| ๙ | | ไม้เก๋าห่อหนึ่ง |

2.1.4 พยัญชนะตัวสะกด

หลักในการใช้พยัญชนะสะกดของภาษาไทยล้านนา มีวิธีการใช้ตามลักษณะกฎเกณฑ์ดังนี้

- ตัวสะกดอยู่ด้านล่าง เช่น ๑๓๖ ๑๓๗ ๑๓๘
- เมื่อ ไม้หรือสระหรือพยัญชนะควบกล้ำอยู่ข้างข้างล่างตัวพยัญชนะ มักนิยมเขียนตัวสะกดอยู่ถัดไปแนวเดียวกันเช่น ๑๓๖ ๑๓๗ ๑๓๘

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

- ถ้าพยัญชนะตัวสะกดเป็น พยัญชนะเหล่านี้ให้ใช้หางของตัวสะกดแทนและอยู่ข้างล่าง หรือห้อยด้านข้างพยัญชนะเสมอตามลักษณะของแต่ละตัว ยกเว้นมีไม้หรือพยัญชนะควบกล้ำอยู่ได้ตัวอักษร

ตัว น สะกด ๔ ให้ใช้ ๔ สะกด

ตัว บ สะกด ๖ ให้ใช้ ๖ สะกด

ตัว ป สะกด ๒ ให้ใช้ ๒ สะกด

ตัว พ สะกด ๓ ให้ใช้ ๓ สะกด

ตัว ผ สะกด ๗ ให้ใช้ ๗ สะกด

ตัว ม สะกด ๕ ให้ใช้ ๕ สะกด

ตัว ย สะกด ๘ ให้ใช้ ๘ สะกด

2.1.5 หลักการเขียนภาษาไทยล้านนา [2]

หลักการเขียนพยัญชนะไทยล้านนาจะมีความแตกต่างจากหลักภาษาไทย เพราะพยัญชนะไทยล้านนาเมื่อเป็นตัวสะกด บางที่จะห้อยอยู่ด้านล่าง ด้านข้างหรือเปลี่ยนรูป ซึ่งมีวิธีใช้ตามลักษณะกฎเกณฑ์ดังนี้

- ตัวสะกดโดยปกติจะอยู่ด้านล่าง
- เมื่อมีไม้หรือสระ หรือพยัญชนะควบกล้ำอยู่ด้านล่างตัวพยัญชนะ มักนิยมเขียนตัวสะกดอยู่ถัดไปในแนวเดียวกัน

- ถ้าพยัญชนะตัวสะกดเป็น ๔ ๖ ๒ ๓ ๗ ๕ ๘ ให้ใช้หางของตัวเองสะกดแทนและอยู่ข้างล่างหรือห้อยด้านข้างพยัญชนะเสมอ(ตามลักษณะของแต่ละตัว) ยกเว้นมีไม้หรือพยัญชนะควบกล้ำอยู่ได้ตัวอักษร

วิธีใช้สระและไม้ต่างๆ

- ๓๔ ไม้กะ (สระอะ) คำที่ออกเสียงสระอะ จะไม่เติมไม้กะลงไป เช่น จะไป ๑ไป ๓เน ๑๑๑ ๓วัน ๑๑ เมื่ออยู่ท้ายพยัญชนะ จะเท่ากับตัวไม้ไต่คู้ คือบังคับให้ออกเสียงสั้น เร็ว เช่น เห็ด ๑๑๑

- ๓๕ ไม้กุก (สระอุ) ไม้กุก โดยทั่วไปก็จะใช้เหมือนปกติ แต่จะมีข้อยกเว้นที่เวลาที่มีพยัญชนะอยู่ด้านล่าง จะเขียนอยู่ด้านล่างหรือด้านข้างก็ได้

- ๓๖ ไม้โก (สระโ) ใช้เหมือนสระโ แต่ถ้าเป็นคำที่มาจากภาษาบาลี หรือ คำบาลีจะใส่ ๑-๑ แทนไม้ไต่คู้ เช่น คำว่า พุทฺโธ ธมฺโม สงฺโฆ โทติ ๑๑๑ ๑๑๑ ๑๑๑ ๑๑๑

- c-ɯ ไม้เก็ย การใช้ไม้เก็ย มีที่ต้องสังเกตคือ หากไม่มีตัวสะกด เช่น เมีย เสีย ก็จะเขียนว่า ฉย ฉย เมื่อมีตัวสะกด ก็จะละ c คงไว้แต่ ɯ ต่อด้วยตัวสะกด เช่น เสียง เวียง ฉยฉ ฉยฉ

- ɯ-ɯ ไม้ไถย ไม้ไถย นี้ มีค่าเท่ากับตัว ใ-ย และ -ย เช่น ชัย ฉย วิทยาลัย ฉยฉย ไทย ฉย

- ๓๓ ซึ่งจะเป็นการใช้เมื่อมีพยัญชนะสองตัวซ้อนกันเสมอ ใช้เขียนบนพยัญชนะเพื่อให้รู้ว่าเป็นตัวสะกด เช่น ลูก ฉยฉ กอด ฉยฉ เสียม ฉยฉ

- ๔ ไม้สองน้อย ใช้เหมือนไม้จ้ำคำ แต่เขียนด้านบนพยัญชนะ เช่น ไป ๆ มา ๆ ฉย ฉย ใช้เขียนเพื่อให้รู้ว่าต้องอ่านแยก เช่น สนาม ฉยฉ สมัย ฉยฉ

- ๕ ไม้ก่ง, ไม้งาม, ตัวخم, ไม้โกศครูป ใช้เขียนบนพยัญชนะเพื่อบอกให้รู้ว่าเป็นแม่ กก กน กบ เช่น ฉย ฉย ฉย ฉย แต่ถ้าตัวสะกดเป็นตัว ว ก็จะออกเสียงเป็น อว แทน เช่น ฉย ฉย

- ๖ ไม้เก๊าจู้ เป็นรูปพิเศษ เท่ากับ -เา เช่น เา ฉย

- ๗ ไม้เก๊าห่อนึ่ง เป็นรูปพิเศษ เท่ากับ -เา เช่น เา ฉย

2.2 การเปลี่ยภาษา

การเปลี่ยภาษาเป็นศาสตร์หนึ่งที่ทำให้ซอฟต์แวร์คอมพิวเตอร์ทำการเปลี่ยข้อความหรือคำพูดในภาษาธรรมชาติ จากภาษาหนึ่งทีเรียกว่าภาษาต้นทาง ไปยังอีกภาษาหนึ่งทีเรียกว่าภาษาเป้าหมาย โดยทีสามารถจำแนกการเปลี่ยภาษาจากองค์ประกอบในประโยค ได้ดังต่อไปนี้

2.2.1 การเปลี่ยภาษาระดับคำ

“คำ”[3] ประกอบขึ้นด้วยส่วนย่อยทีเรียกว่าหน่วยคำ(morpheme) คำหนึ่งอาจประกอบขึ้นด้วยหน่วยคำเพียง 1 หน่วยหรือมากกว่าก็ได้ ซึ่งในภาษาไทยภาษาไทยเป็นคำโดด การวิเคราะห์หน่วยคำในลักษณะแยกพยางค์จึงไม่เหมาะสมกับภาษาไทย ต่างจากภาษาอังกฤษทีมีการเติมพยัญชนะลงไปที่ท้ายคำจะทำให้หน่วยคำดังกล่าวทำหน้าทีแตกต่างกันไป เช่น เมื่อเติม s ทีจะทำให้คำนั้นกลายเป็นพหูพจน์

การเปลี่ยภาษาจากภาษาต้นทาง ไปยังภาษาเป้าหมายวิธีการทีง่ายทีสุดในระดับคำคือ ใช้วิธีการเปลี่ยแบบตรงไปตรงมา โดยเทียบความหมายแบบคำต่อคำ อย่างไรก็ตามในการเปลี่ยภาษาในระดับคำแบบตรงไปตรงมานั้น จะเหมาะสมกับบางภาษาทีมีโครงสร้างของไวยากรณ์ของประโยคทีมีความใกล้เคียงกัน ทีจะสามารถใช้วิธีการดังกล่าวได้ แต่ทีจะเกิดปัญหาในกรณีทีโครงสร้างของไวยากรณ์ของภาษาต้นทางและภาษาเป้าหมายมีความแตกต่างกัน ซึ่งโดยส่วนใหญ่แล้วภาษาในแต่ละชาติแต่ละท้องถิ่น จะมีความแตกต่างกันพอสมควร

ในการศึกษางานวิจัยการเปลี่ยภาษาล้านนาเป็นภาษาไทย[4] พบว่างานวิจัยนี้มีเป้าหมายในการเปลี่ยภาษาล้านนาให้เป็นภาษาไทยซึ่งเป็นคู่ภาษาเดียวกันกับงานวิจัยการเปลี่ยภาษาไทยเป็น

ภาษาไทยล้านนา งานวิจัยดังกล่าวได้นิยามว่าในการแปลภาษาไทยล้านนาเป็นภาษาไทยนั้น โครงสร้างไวยากรณ์ของประโยคของภาษาไทยล้านนาที่เป็นภาษาต้นทาง เป็นโครงสร้างเดียวกันกับสร้างไวยากรณ์ของประโยคของภาษาไทยที่เป็นภาษาเป้าหมายในการแปล ซึ่งจากเหตุผลข้างต้นในงานวิจัยดังกล่าวจึงได้เลือกใช้วิธีการแปลในระดับคำด้วยวิธีการแบบตรงไปตรงมา โดยการใช้วิธีการค้นหาคำที่ได้จากประโยคต้นทางที่ผ่านการตัดคำแล้ว นำคำที่ต้องการแปลมาค้นหาคำแปลจากพจนานุกรมภาษาไทยล้านนา-ภาษาไทยที่ได้เตรียมไว้ ส่วนคำไหนที่ไม่พบในพจนานุกรมก็จะใช้ช่างานเพิ่มขยายในการสร้างหน่วยคำ งานวิจัยการแปลภาษาไทยล้านนาเป็นภาษาไทยเป็นงานที่มีความใกล้เคียงกับงานวิจัยการแปลภาษาไทยเป็นภาษาไทยล้านนาเพราะใช้คู่ภาษาเดียวกันคือภาษาไทยล้านนาและภาษาไทย แต่จะสลับภาษาต้นทางและภาษาเป้าหมายกัน นอกจากนี้ยังใช้วิธีการแปลโดยอ้างอิงพจนานุกรมคู่ภาษาเดียวกันอีกด้วย

สถาปัตยกรรมในการแปลภาษาไทยล้านนาเป็นภาษาไทย ในงานวิจัยดังกล่าวมีขั้นตอนที่เกี่ยวข้องกับโครงสร้างระดับคำ ดังต่อไปนี้ พจนานุกรมล้านนา-ไทย ที่ใช้ในงานวิจัยการแปลภาษาไทยล้านนาเป็นภาษาไทย เป็นพจนานุกรมภาษาไทยล้านนาเป็นภาษาไทย มีจำนวนคำทั้งหมด 7,497 คำ มีองค์ประกอบ 3 ส่วน คือ

- คำล้านนา

- เครื่องหมายกำกับ

- ส่วนวิเคราะห์ความหมาย ประกอบด้วย ประเภทของคำ คำที่ถอดออกมาเป็นภาษาไทย ส่วนกำกับหมวดคำ ลักษณะประจำคำ ความหมายภาษาอังกฤษ และศัพท์อักษร การจัดเก็บ จัดเก็บเป็นโครงสร้างแบบทรี มีลักษณะคล้ายโครงสร้างข้อมูลแบบต้นไม้ โดยมีเส้นเชื่อมของโหนด โดยที่แต่ละโหนดจะแสดงตำแหน่งของตัวอักษรถัดไป การค้นคืนของโครงสร้างแบบทรีจะทำให้รวดเร็วแต่ต้องใช้หน่วยความจำในการจัดเก็บมาก จุดเด่นคือ สะดวกในการควบคุมข้อมูลที่มีความยาวแตกต่างกัน มีความรวดเร็วในการเข้าถึงข้อมูลเพราะเข้าถึงข้อมูลโหนดหนึ่งไปยังโหนดถัดไป

วิธีการในการแปลเมื่อรับข้อมูลต้นทางที่เป็นประโยคภาษาไทยมาล้านนา ก็จะมีขั้นตอนประกอบไปด้วย การตัดคำ การเลือกคำแปลในพจนานุกรมและการปริวรรตคำไทยจากคำล้านนาที่ไม่พบในพจนานุกรม

- ในขั้นตอนของการตัดคำจะใช้วิธีการตัดคำล้านนาที่นำเข้ามาแบบเลือกคำยาวที่สุด (Longest matching) โดยจะนำประโยคที่ต้องการตัดคำ ตัดตัวอักษรทางขวาออกทีละ 1 ตัวอักษร จากนั้นจะนำอักษรที่เหลือทางด้านซ้ายมือไปค้นหาในพจนานุกรม เพื่อหารูปแบบของคำที่ตรงกับคำที่เหลือทางซ้ายมือ เมื่อกันหาพบก็จะสามารถตัดพยางค์และได้คำที่ตรงกับพจนานุกรม ผลลัพธ์จะได้คำที่ผ่านการตัดคำจากประโยคนำเข้า

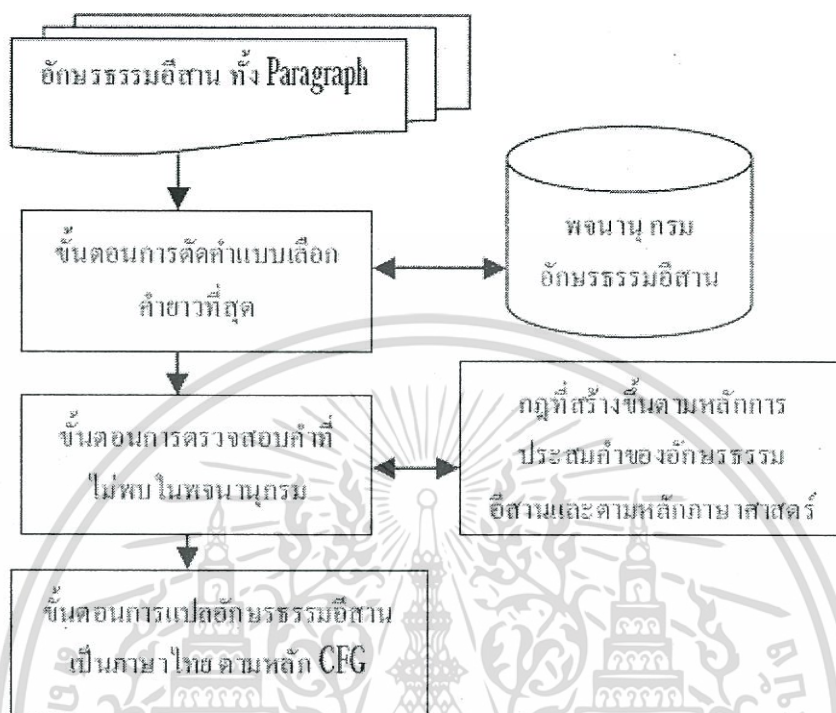
- การปริวรรตคำที่ไม่พบในพจนานุกรม ในงานวิจัยดังกล่าวจะใช้ข่ายงานเพิ่มขยาย (Nondeterministic Finite Automata NFA) สร้างกฎการปริวรรตล้านนาเป็นอักษรไทยเพื่อทำการแปลงคำไทยล้านนาที่ตัดคำและนำไปค้นหาในพจนานุกรมแล้งปรากฏว่าไม่พบข้อมูลคำไทยล้านนาคำดังกล่าวในพจนานุกรม มาแปลงรูปจากโครงสร้างของคำภาษาไทยล้านนาให้เป็นโครงสร้างคำภาษาไทย



ภาพที่ 2.1 : สถาปัตยกรรมในการแปลภาษาไทยล้านนาเป็นภาษาไทย

นอกจากนั้นในการศึกษาวิจัยการแปลภาษาอักษรธรรมอีสานให้เป็นภาษาไทย[5] งานวิจัยนี้เป็นงานวิจัยที่ใกล้เคียงกับงานวิจัยที่สนใจตรงที่เป็นงานวิจัยที่เป็นภาษาถิ่น นอกจากนั้นอักษรธรรมอีสานยังมีลักษณะและรูปของตัวอักษรที่ใกล้เคียงกับอักษรไทยล้านนาเนื่องจากในอดีต อาณาจักรล้านช้างที่เป็นอาณาจักรของชนชาติที่อยู่ติดกับลุ่มน้ำโขง กับอาณาจักรล้านนาที่เป็นอาณาจักรของหลายจังหวัดภาคเหนือตอนบนของประเทศไทย นั้นได้รับอิทธิพลวัฒนธรรมจากอาณาจักรขอมซึ่งรวมไปถึงตัวอักษรที่ใช้จดบันทึกด้วย งานวิจัยดังกล่าวในระดับคำนั้นได้ใช้วิธีการคล้ายกับงานวิจัยในการแปลภาษาไทยล้านนาเป็นภาษาไทย โดย จะนำคำที่ผ่านการตัดคำมาค้นหาในพจนานุกรมที่ได้เตรียมไว้ หากไม่พบคำดังกล่าวก็จะทำการแปลงอักษรธรรมอีสานแปลงรูปให้เป็นภาษาไทยด้วยกฎที่สร้างขึ้นตามหลักการประสมคำของอักษรธรรมอีสาน ตามหลักภาษาศาสตร์

ซึ่งมีความคล้ายคลึงกับงานวิจัยการแปลภาษาไทยล้านนาเป็นภาษาไทย ในงานวิจัยดังกล่าวมีขั้นตอนดังต่อไปนี้



ภาพที่ 2.2 : สถาปัตยกรรมของการแปลอักษรธรรมอีสานให้เป็นภาษาไทย

วิธีการในการแปลอักษรธรรมอีสานให้เป็นภาษาไทยเมื่อรับข้อมูลต้นทางที่เป็นอักษรธรรมอีสาน ก็จะมีขั้นตอนประกอบไปด้วย การตัดคำ การเลือกคำแปลในพจนานุกรมและการสร้างคำที่ไม่พบในพจนานุกรม

- ในขั้นตอนของการตัดคำจะใช้วิธีการตัดแบบเลือกคำยาวที่สุด(Longest matching) วิธีการตัดคำแบบเลือกคำยาวมีวิธีคล้ายกับงานวิจัยการแปลภาษาไทยล้านนาเป็นภาษาไทย ผลลัพธ์ก็จะได้อักษรธรรมอีสานในแต่ละคำจากนั้นก็ให้นำคำที่ได้ไปค้นหาคำแปลในพจนานุกรมอักษรธรรมอีสานผลลัพธ์จะได้คำแปลที่เป็นภาษาไทย

- พจนานุกรมอักษรธรรมอีสานในการแปลอักษรธรรมอีสานให้เป็นภาษาไทย ก็จะทำให้การสร้างพจนานุกรมอักษรธรรมอีสาน มีการจัดเก็บเป็นโครงสร้างแบบทรี เพื่อใช้ในการเลือกคำแปลอักษรธรรมอีสานให้เป็นภาษาไทย

- เมื่อนำคำที่ต้องการ ไปค้นในพจนานุกรมแล้ว หากไม่พบคำศัพท์ดังกล่าวในพจนานุกรมอักษรธรรมอีสาน ก็จะนำคำอักษรธรรมอีสานที่ไม่พบในพจนานุกรมนั้นเข้าสู่ขั้นตอน

การประสมคำจากหลักการประสมคำของอักษรธรรมอีสานให้เป็นภาษาไทยตามหลักภาษาศาสตร์ โดยใช้กฎที่สร้างขึ้นตามหลักการประสมคำอักษรธรรมอีสาน

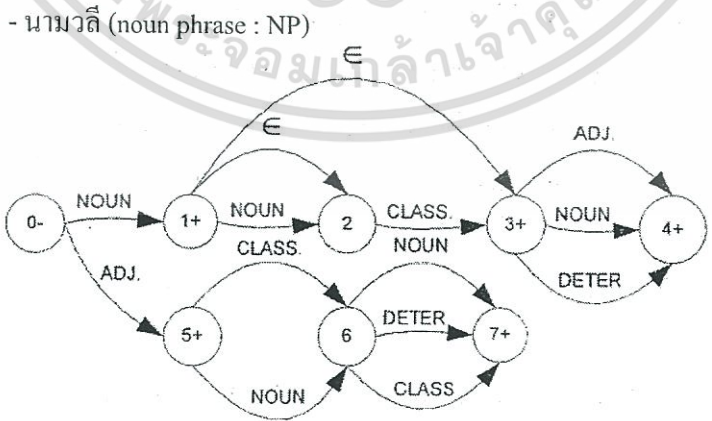
- ผลลัพธ์จากขั้นตอนการค้นหาคำศัพท์และการประสมคำจากหลักการประสมคำของอักษรธรรมอีสานให้เป็นภาษาไทย ก็จะประโยคภาษาไทยที่ยังไม่ใช่โครงสร้างไวยากรณ์ภาษาไทย ต้องผ่านขั้นตอนของการแปลโดยใช้หลักวิเคราะห์โครงสร้างไวยากรณ์ไม่พึ่งบริบท (Context-free grammar: CFG)

2.2.2 การแปลภาษาระดับวลี

วลี (Phrases) [3] คือกลุ่มคำที่ประกอบด้วยคำต่างๆ ที่นำมาเรียงกันอย่างมีความหมายและทำหน้าที่ใดหน้าที่หนึ่งในประโยค เช่น เป็นประธาน กริยา กรรม ส่วนเสริม ส่วนกริยาวิเศษณ์ วลีไม่ได้ประกอบด้วยภาคประธานและภาคแสดง นักภาษาศาสตร์อย่าง Chomsky ได้กล่าวว่า คำเพียงหนึ่งคำก็เป็นวลี แต่ในที่นี้จะถือว่า วลีคือคำตั้งแต่ 2 คำขึ้นไปที่ทำหน้าที่ใดหน้าที่หนึ่งในประโยค จากงานวิจัยการแปลภาษาด้านนาเป็นภาษาไทย[4] ในขั้นตอนของการวิเคราะห์โครงสร้างวลีและโครงสร้างประโยค จะใช้ข่ายงานเพิ่มขยาย Augmented Transition Network (ATNs) ซึ่งเป็นกราฟชนิดหนึ่งประกอบด้วยสถานะซึ่งแทนด้วยรูปวงกลมและการส่งผ่านคำซึ่งแทนด้วยเส้นโค้ง สถานะสุดท้ายจะกำกับด้วย “/1” ในแต่ละเส้นโค้งจะมีตัวเก็บค่าและหมายเลขกำกับเพื่อส่งผ่านคำในประโยคเข้าสู่ข่ายงาน ตัวเก็บค่าจะทำหน้าที่วิเคราะห์และบันทึกบทบาทของคำตามเงื่อนไขที่สร้างขึ้น

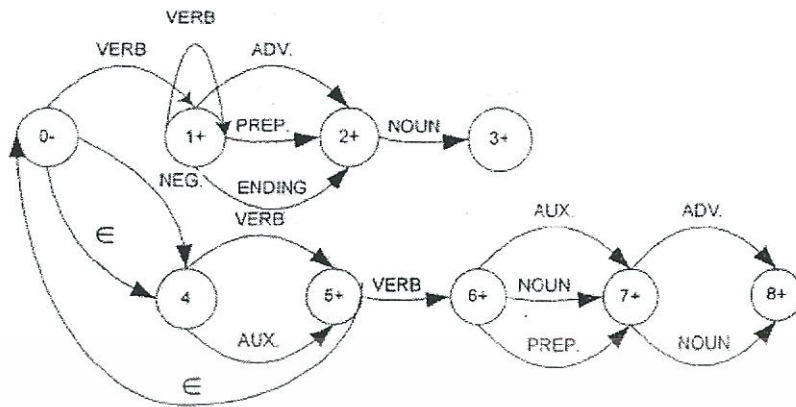
เมื่อนำประโยคมาผ่าน โครงสร้างเพื่อหาความสัมพันธ์ในประโยคแล้วก็จะได้ความสัมพันธ์ของวลีในประโยคภาษาตั้งต้น โครงสร้างวลีของภาษาไทย จะประกอบด้วย 2 ส่วนคือส่วนมูลฐาน ทำหน้าที่เป็นภาคประธาน ภาคแสดงและส่วนเสริม ทำหน้าที่เป็นส่วนขยาย

โครงสร้างวลีของภาษาไทยที่เป็นส่วนมูลฐาน



ภาพที่ 2.3 : โครงสร้างวลีของนามวลี

- กริยาวลี (Verb phrase : VP)



ภาพที่ 2.4 : โครงสร้างวลีของกริยาวลี

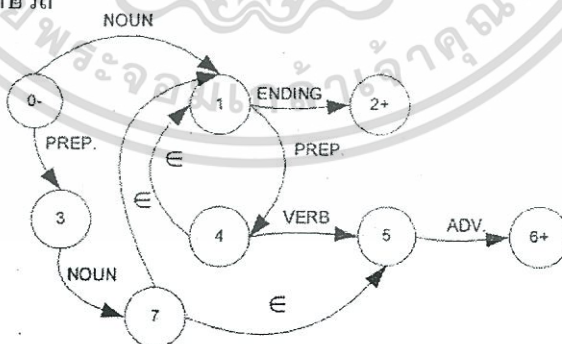
โครงสร้างวลีภาษาไทยที่เป็นส่วนเสริม

- สถานวิเศษณ์วลี



ภาพที่ 2.5 : โครงสร้างวลีของสถานวิเศษณ์วลี

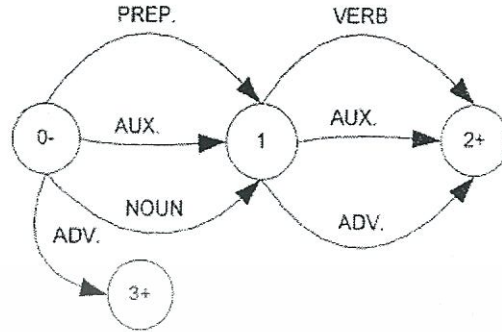
- กาลวิเศษวลี



ภาพที่ 2.6 : โครงสร้างวลีของกาลวิเศษวลี

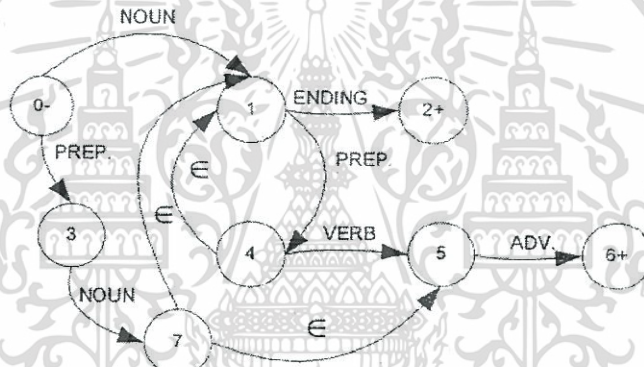
เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

- พิเศษวิเศษณ์วลี



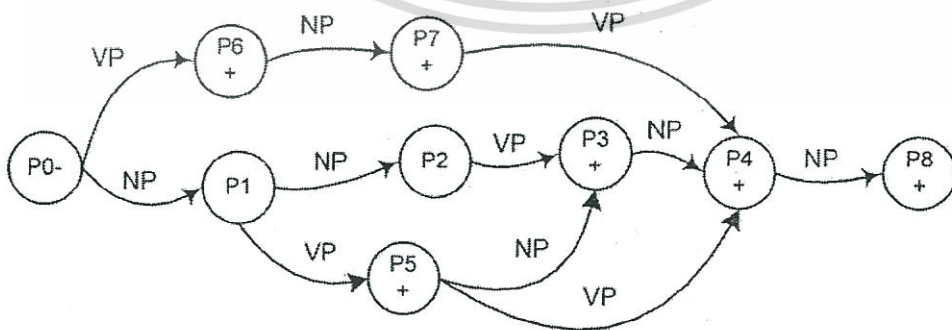
ภาพที่ 2.7 : โครงสร้างวลีของพิเศษวิเศษณ์วลี

- อาลปนะวลี



ภาพที่ 2.8 : โครงสร้างวลีของอาลปนะวลี

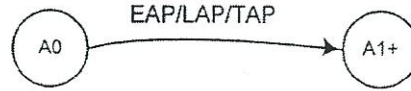
โครงสร้างประโยค ส่วนมูลฐานทำหน้าที่เป็นหน่วยประธาน หน่วยกรรม หน่วยนามเดี่ยวน (NP) และทำหน้าที่เป็นหน่วยกริยา (VP) ข่ายงานความสัมพันธ์ของส่วนมูลฐานแสดงดังภาพที่ 2.9



ภาพที่ 2.9 : ข่ายงานความสัมพันธ์ของส่วนมูลฐานในประโยค

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ส่วนเสริมสามารถเรียงลำดับ โดยเริ่มต้นจากวลีใดๆ และสลับตำแหน่งโดยไม่ทำให้ความหมายของประโยคเปลี่ยนแปลงจากภาพที่ 2.10 เป็นข่ายงานการเรียงลำดับของส่วนเสริม เมื่อ EAP คือพิเศษวิเศษณ์วลี , LAP คือ สถานวิเศษณ์วลี และ TAP คือ กาลวิเศษณ์วลี



ภาพที่ 2.10 : ข่ายงานความสัมพันธ์ของส่วนเสริมในประโยค

จากความสัมพันธ์ของส่วนมูลฐานและการเรียงลำดับของส่วนเสริม สามารถสร้างโครงสร้างประโยคภาษาไทยด้วย ATNs ได้ทั้งหมด 4 โครงสร้าง โดยให้ P เป็นส่วนมูลฐานและให้ A เป็นส่วนเสริมดังข่ายงานที่แสดงในภาพที่ 2.11



ภาพที่ 2.11 : ข่ายงานความสัมพันธ์ของความสัมพันธ์ของส่วนมูลฐานและส่วนเสริม

ขั้นตอนการแปลประโยคโดยใช้ข่ายงานเพิ่มขยาย ประกอบไปด้วย การวิเคราะห์โครงสร้างวลีและโครงสร้างประโยค และการเลือกความหมายโดยพิจารณาจากบริบทที่เกิดขึ้น

- การวิเคราะห์โครงสร้างวลีและโครงสร้างประโยค โดยใช้ข่ายงานเพิ่มขยาย Augmented Transition Network (ATNs)

- การเลือกคำแปลที่เหมาะสมพิจารณาจากบริบท และการเกิดร่วมของคำ จะใช้คุณสมบัติ relation ซึ่งเป็นตัวแปลคลาสของแต่ละคำ เช่นคำว่า “ม่วน” มีหลายคำแปล แต่ละคำตัวแปลคลาส relation จะเก็บคำที่เกิดร่วมเช่น

คำว่าม่วน แปลว่า สนุก relation เล่น-เที่ยว

คำว่าม่วน แปลว่า เพราะ relation ฟัง-เพลง-พูด

ขั้นตอนการแปลเป็นภาษาไทยด้วยโครงสร้างไวยากรณ์ ตามหลัก CFG โดยมีกฎโครงสร้างวลี 5 กฎและกฎย่อยภายใต้โครงสร้างอีก 40 กฎ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2.2.3 การแปลภาษาระดับประโยค

ประโยค [3] เกิดจากกลุ่มคำหรือวลีประกอบกันเข้าเป็นประโยค โดยกลุ่มคำที่นำมาประกอบกันนั้นจะมีเนื้อความสมบูรณ์ บอกการกระทำความเป็นอยู่หรือความเป็นไปของสิ่งหนึ่งสิ่งใด การประกอบเข้ากันเป็นประโยคจะต้องประกอบด้วยภาคแสดงอย่างน้อยหนึ่งหน่วย เช่น “สุดา อ่านหนังสือ เมื่อวานนี้” ประโยคประกอบด้วย ภาคประธานและภาคแสดงโดย สุดา คือภาคประธานของประโยค และ อ่านหนังสือ เมื่อวานนี้ คือภาคแสดงในประโยค

ในการประมวลผลภาษารวมชาติ เพื่อให้คอมพิวเตอร์สามารถรับรู้และเข้าใจภาษามนุษย์ จำเป็นต้องนำหลักการทางภาษาศาสตร์ไปใช้ โครงสร้างประโยคของภาษานักภาษาศาสตร์จะมีการกำหนดกฎเกณฑ์มาช่วยตรวจวิเคราะห์และสร้างประโยค ตัวอย่างเช่น ไวยากรณ์ไม่พึ่งบริบท (Context Free Grammar) , ข่ายงานเพิ่มขยาย Augmented Transition Network (ATNs) เป็นต้น

ไวยากรณ์ไม่พึ่งบริบท (Context Free Grammar)[9] การแจงประโยคสำหรับไวยากรณ์ CFG อาจจะแจกแจงประโยคจากบนลงล่าง (Top-down parsing) และการแจงประโยคจากล่างขึ้นบน (Bottom-up parsing) การแจงประโยคจากบนลงล่าง จะเริ่มที่สัญลักษณ์ประโยค S เขียนแทนด้วยสัญลักษณ์ นามวลี (NP) และตามด้วย กริยาวลี (VP) สัญลักษณ์ NP และ VP จะเขียนแทนต่อไปจนถึงสัญลักษณ์ จบท้ายที่แทนด้วยคำศัพท์ต่างๆ ลักษณะของการแจงประโยคแบบ บนลงล่างจะเป็นการใช้กฎไวยากรณ์ของสัญลักษณ์ ทางด้านขวาไปเขียนใหม่แทนสัญลักษณ์ทางด้านซ้าย ดังภาพที่ 2.12

| | |
|---|----------------------------|
| S | → NP + VP |
| | → Art + N + VP |
| | → The + N + VP |
| | → The birds + VP |
| | → The birds + V + NP |
| | → The birds eat + NP |
| | → The birds eat + Art + N |
| | → The birds eat the + N |
| | → The birds eat the worms. |

ภาพที่ 2.12 : การแจงประโยคจากบนลงล่าง (Top-down parsing)

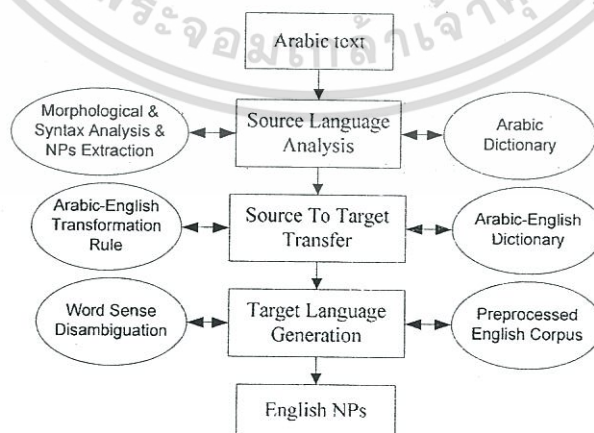
การแจงประโยคจากล่างขึ้นบน จะเริ่มต้นจากประโยค แล้วแทนด้วยคำศัพท์และชนิดของคำแทนสัญลักษณ์ของกฎทางขวาด้วยสัญลักษณ์ทางซ้ายไปจนถึงสัญลักษณ์ S ดังภาพที่ 2.13

- The birds eat the worms.
- Art + birds eat the worms.
- Art + N + eat the worms.
- Art + N + V + the worms.
- Art + N + V + Art + worms.
- Art+N + V + Art + N
- NP + V + Art+N
- NP + V+ NP
- NP + VP
- S

ภาพที่ 2.13 : การแจงประโยคจากล่างขึ้นบน (Botton-up parsing)

ในงานวิจัยการแปลภาษาอาหรับเป็นภาษาอังกฤษ[6] งานวิจัยนี้ใช้วิเทอริบีอัลกอริทึมในการเลือกตำแหน่งของโครงสร้างประโยคที่มีความเป็นไปได้สูงสุด จากคลังข้อมูลภาษาอังกฤษโดยสถาปัตยกรรมของการแปลภาษาอาหรับเป็นภาษาอังกฤษ ซึ่งแปลเฉพาะ Noun Phases ของประโยคอาหรับ

เป้าหมายหลักของระบบคือสร้างระบบการแปลโดยไม่ใช่ คลังคู่ขนานและสร้างกฎที่ครอบคลุม ภาพที่ 2.14 จะแสดงสถาปัตยกรรมของระบบแปลภาษาอาหรับเป็นภาษาอังกฤษ และทรัพยากรที่เกี่ยวข้องในแต่ละขั้นตอนการแปล ระบบแบ่งออกเป็น 3 ส่วนประกอบ คือ วิเคราะห์ประโยคต้นทาง (Source Language Analysis) , การถ่ายทอดจากภาษาต้นทางไปยังภาษาเป้าหมาย (Source to Target Transfer) และการสร้างประโยคภาษาเป้าหมาย (Target Language Generation)



ภาพที่ 2.14 : ขั้นตอนการแปลภาษาอาหรับเป็นภาษาอังกฤษ

ขั้นตอนที่ 1 วิเคราะห์ประโยคต้นทาง (Source Language Analysis) จะวิเคราะห์และแยกประโยคภาษาอาหรับ ในงานวิจัยดังกล่าวจะครอบคลุมเอกสารทางด้านการเกษตร หลังจากขั้นตอนการพาสซิงของประโยคนำเข้าภาษาอาหรับ นามวลีของภาษาอาหรับ (NP) จะถูกแยกโดยโครงสร้างไวยากรณ์ของภาษาอาหรับ

ขั้นตอนที่ 2 การถ่ายทอดจากภาษาต้นทางไปยังภาษาเป้าหมาย (Source to Target Transfer) จะทำการเปรียบเทียบ โครงสร้างของประโยค ทำการเลือกโครงสร้างต้นไม้อิงภาษา วิเคราะห์ต้นไม้อิงภาษาต้นทางเพื่อนำมาวิเคราะห์โครงสร้างต้นไม้อิงภาษาเป้าหมาย ด้านหนึ่งของต้นไม้อิงภาษาอาหรับ จะมีกฎของการถ่ายโอนโครงสร้างต้นไม้อิงภาษา ในการจับคู่กับโครงสร้างข้อมูลนำเข้า (ภาษาอาหรับ) ผลลัพธ์จะได้ต้นไม้อิงภาษาทางด้านขวามือ

ในการแปล NP ที่นำเสนอ การแปลจะเกิดขึ้นในขั้นตอน transfer phase ในการแปลภาษาอาหรับเป็นภาษาอังกฤษ transfer phase มี 2 ขั้นตอน

- การถ่ายโอนคำศัพท์ (Lexical transfer) นอกจากจะเปรียบเทียบคำศัพท์ภาษาอาหรับเป็นภาษาอังกฤษแล้วยัง รวมถึงคุณสมบัติทางสัณฐานวิทยา (morphological) ที่ตรงกับคุณสมบัติภาษาอังกฤษ

- การถ่ายโอนโครงสร้าง (Structural transfer) จะเปรียบเทียบโครงสร้างต้นไม้อิงภาษาอาหรับ เทียบเท่ากับโครงสร้างของภาษาอังกฤษ

ขั้นตอนที่ 3 การสร้างประโยคภาษาเป้าหมาย (Target Language Generation) เป้าหมายคือสร้างภาษาเป้าหมาย โดยใช้คุณลักษณะของ สัณฐานวิทยา (morphological) แม้ในโมดูลก่อนหน้าจะลดความคลุมเครือในการแปลแต่ก็ยังมีความคลุมเครือในระดับคำอยู่

ในโมดูลนี้ใช้วิธีการ dictionary-graph based WSD เป็นการตรวจสอบการทำงานร่วมกันของวิธีการ dictionary-based WSD และ graph-based WSD ทำให้การใช้คำที่ได้รับการแปลเป็นบริบท, บรรลวดัตถุประสงค์ขั้นพื้นฐาน WSD ในขั้นตอนของการแปล จะใช้พจนานุกรมอาหรับ-อังกฤษ ที่ได้อธิบายในข้างต้นแปลคำในแต่ละคำจาก ประโยคนำเข้าภาษาอาหรับ NP เพื่อแก้ปัญหาคความคลุมเครือในการแปลภาษาอาหรับ NP ระบบการแปลระบุความคลุมเครือของคำและความสัมพันธ์ระหว่างกัน และใช้การค้นหาโดยอัลกอริทึมวิเทอร์บี (viterbi search) หากการแปลที่เหมาะสม ของคำภาษาอาหรับในการสร้างภาษาอังกฤษ NP

```

function VITERBI(observations of len T, state-graph) returns best-path
    num-states ← NUM-OF-STATES(state-graph)
    Create a path probability matrix viterbi[num-states+2, T+2]
    viterbi[0,0] ← 1.0
    for each time step t from 1 to T do
        for each state s from 1 to num-states do
            viterbi[s,t] ← max1 ≤ s' ≤ num-states [viterbi[s',t-1] * as',s] * bs(ot)
            back-pointer[s,t] ← argmax1 ≤ s' ≤ num-states [viterbi[s',t-1] * as',s]
    Backtrace from highest probability state in final column of viterbi[] and return path

```

ภาพที่ 2.15 : ฟังก์ชันวิเทอร์บี สำหรับหาโครงสร้างประโยคภาษาอังกฤษ [6]

แนวคิดพื้นฐาน หากเรามีความคลุมเครือของคำภาษาอาหรับ S ซึ่งมี 2 ประโยคคือ S1 และ S2 โดยที่ S1 จะสามารถแปลได้เป็น T1 และ S2 สามารถแปลได้เป็น T2 เพื่อที่จะทำให้เกิดความไม่คลุมเครือในการเกิด S ในภาษาอาหรับ . เราจะระบุนวลีที่เกิดขึ้นและใช้วิเทอร์บีเพื่อที่จะหาคำแปลที่เหมาะสมของคำภาษาอาหรับ S ขึ้นอยู่กับโครงสร้างภาษาเป้าหมายภาษาอังกฤษ NP

วิเทอร์บีอัลกอริทึม เป็นเทคนิคที่มีประสิทธิภาพ จะคำนวณลำดับโครงสร้างของประโยคภาษาเป้าหมายที่เป็นไปได้มากที่สุด จากงานวิจัยนี้ ทำให้ได้แนวคิดในการนำทฤษฎีของแบบจำลองฮิดเดินมาร์คอฟ โดยใช้วิเทอร์บีอัลกอริทึม มาทำการลดความกำกวมของคำศัพท์ในการแปลภาษาไทยเป็นภาษาไทยล้านนา

ในงานวิจัยการแปลภาษาไทยล้านนาเป็นภาษาไทยนั้น[4] ได้ใช้แนวคิดที่ว่าโครงสร้างของประโยคภาษาไทยล้านนาแปลเป็นภาษาไทยนั้นเป็นโครงสร้างประโยคเดียวกัน ซึ่งเมื่อแปลแล้วตำแหน่งของคำแต่ละคำในประโยคไทยล้านนา ก็จะอยู่ในตำแหน่งเดิมของประโยค ภาษาไทย แต่ในส่วนของโครงสร้างประโยคในการแปลจากภาษาไทยเป็นภาษาไทยล้านนานั้น พบว่ายังมีโครงสร้างประโยคบางประโยค ที่มีโครงสร้างที่ไม่ตรงตำแหน่งเดิมอยู่

ตัวอย่างของโครงสร้างประโยคไทยและโครงสร้างประโยคไทยล้านนา ที่โครงสร้างของประโยคมีความแตกต่างกัน

เขาวิ่งเร็วมาก เมื่อ แปลเป็นภาษาไทยล้านนาจะได้ “เขาล่นซาดเวย” ซึ่งคำใดๆ ของประโยคภาษาไทย ที่ทำหน้าที่เป็นกริยา เมื่อตามด้วยคำว่า มาก จะพบว่าสามารถแปลเป็นภาษาไทยล้านนาแล้วมีโครงสร้างของประโยคที่แตกต่างจากเดิม

ซึ่งในงานวิจัยนี้จะใช้วิธีการปรับปรุงตำแหน่งโครงสร้างประโยคด้วย ไวยากรณ์ปริวรรติเพิ่มพูน(Transformation Rule) ในการทำให้ในประโยคมีโครงสร้างคำที่มีความใกล้เคียงกับการที่มนุษย์มากยิ่งขึ้น

ตัวอย่างกฎไวยากรณ์ปริวรรติเพิ่มพูน ของภาษาไทยเป็นภาษาไทยล้านนา ประโยค “ส้มตำอร่อยมาก” ในภาษาไทยสามารถแปลเป็นภาษาไทยล้านนาได้มากกว่า 1 โครงสร้าง คือ “ตำส้มจ้าดตำ” และ “ตำส้มตำขนาด”

ส้มตำ/N อร่อย/V มาก/Adv

ตำส้ม/N จ้าด/Adv ตำ/V ตำส้ม/N ตำ/V ขนาด/Adv

กฎการไวยากรณ์ปริวรรติของกริยาช่วยของคำว่า มาก (Adv) เมื่อไวยากรณ์ที่ผ่านเข้ามายังไวยากรณ์ปริวรรติเมื่อพบคำว่า “มาก” ที่มีคู่ภาษาล้านนาเป็น “จ้าด” ในภาษาไทยล้านนา จะปรับตำแหน่งประเภทของคำที่เป็นกริยาช่วย ไปอยู่ด้านหน้าของกริยา ดังต่อไปนี้

SD : NP V Adv

[มาก]

SC : NP Adv V

[จ้าด]

ผลลัพธ์ : จะได้โครงสร้างใหม่เพิ่มและโครงสร้างเดิม คือ

- ตำส้ม/N จ้าด/Adv ตำ/V

- ตำส้ม/N ตำ/V ขนาด/Adv

ดังนั้น เมื่อผ่านไวยากรณ์ปริวรรติเพิ่มพูน ประโยคไทยล้านนาดังกล่าวจะทำให้ผลลัพธ์ของการแปลภาษาไทยเป็นภาษาไทยล้านนามีโครงสร้างประโยคเมื่อปรับปรุงโครงสร้างแล้วมีความใกล้เคียงกับการใช้งานของมนุษย์มากขึ้น

บทที่ 3

วิธีการดำเนินงานวิจัย

งานวิจัยนี้มีวัตถุประสงค์ในการแปลภาษาไทยเป็นภาษาล้านนา โดยเลือกความหมายจากพจนานุกรมไทยล้านนา วิเคราะห์โครงสร้างประโยคภาษาไทยด้วยไวยากรณ์ไม่พึ่งบริบทประยุกต์เทคนิคทางสถิติ (Probabilistic Context free grammar : PCFG) ปรับโครงสร้างไวยากรณ์ให้เป็นโครงสร้างไวยากรณ์ไทยล้านนาด้วย ไวยากรณ์ปริวรรต (Transformation Grammar) และลดความกำกวมด้วยวิธีการของแบบจำลองฮิดเดินมาร์คอฟและวิเทอร์บีอัลกอริทึม (Viterbi algorithm) ผลลัพธ์จะได้ประโยคภาษาล้านนาซึ่งเป็นภาษาเป้าหมาย

3.1 สถาปัตยกรรมของการแปลภาษาไทยเป็นภาษาล้านนา

สถาปัตยกรรมของการแปลภาษาไทยเป็นภาษาล้านนาได้แบ่งขั้นตอนการทำงานออกเป็นขั้นตอนดังภาพที่ 3.1 ซึ่งประกอบด้วยขั้นตอนดังต่อไปนี้

ขั้นตอนที่ 1 หลังจากที่ใช้ผู้ใช้ใส่ประโยคภาษาไทยที่ต้องการแปล จะเข้าสู่ขั้นตอนก่อนการประมวลผล (Preprocessing) ซึ่งในขั้นตอนนี้จะเป็นขั้นตอนของการเตรียมข้อมูลก่อนการแปล จะเป็นการตัดคำภาษาไทยโดยใช้เครื่องมือตัดคำ (KU-CUT) ซึ่งเป็นเครื่องมือที่ตัดคำด้วยเทคนิคการเรียนรู้จากข้อมูลแบบไม่ใช้ตัวอย่างร่วมกับการใช้พจนานุกรมและคำ จากนั้นจะนำข้อมูลที่ผ่านการตัดคำนั้นๆ ไปเข้าสู่การวิเคราะห์โครงสร้าง ผ่านเครื่องจักรวิเคราะห์ประโยคโดยใช้ไวยากรณ์ไม่พึ่งบริบทประยุกต์เทคนิคทางสถิติ (Probabilistic Context free grammar : PCFG) เพื่อวิเคราะห์และจัดรูปแบบตำแหน่งของประเภทของคำ ให้ตรงกับโครงสร้างประโยคไทย ผลลัพธ์จะได้โครงสร้างประโยคภาษาไทยที่ถูกต้องและมีการกำกับ ชนิดของคำ (Part of speech : POS) จากประโยคตั้งต้น

ขั้นตอนที่ 2 การแปลประโยคไทยล้านนา (Mapping Thai-Lanna) ในขั้นตอนนี้จะทำการนำคำศัพท์ภาษาไทยในขั้นตอนที่ 1 ไปค้นหาคำแปลภาษาไทยล้านนาจากพจนานุกรมไทย-ล้านนาซึ่งมีการจัดเก็บข้อมูลในโครงสร้างแบบทรี (Trie) หากในกรณีที่ไม่มีพบคำแปลในพจนานุกรมจะใช้วิธีการแปลงรูปคำไทยล้านนาตามกฎของโครงสร้างของการเขียนภาษาไทยล้านนาและระบุชนิดของคำนั้นเป็นคำนามชนิดวิสามานยนาม ผลลัพธ์ในขั้นตอนนี้ก็จะได้คำแปลล้านนากำกับอยู่ในแต่ละคำไทย

ขั้นตอนที่ 3 การสร้างประโยคล้านนา (Lanna Syntactic Analysis) คำภาษาไทย ที่มีการแจกแจงโครงสร้างประโยค และกำกับชนิดของคำไทย จะนำเข้าสู่ไวยากรณ์ปริวรรตโดยจะมีกฎของการปรับให้เป็นโครงสร้างประโยคไทยล้านนา จากนั้นจะลดความกำกวมของโครงสร้างประโยคและลำดับของคำล้านนา โดยใช้แบบจำลองฮิดเดินมาร์คอฟ วิธีการวิเทอร์บีอัลกอริทึม โดยเอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่นิยามให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

คำนวณสถิติจากคลังภาษาไทยล้านนา โดย ผลลัพธ์ก็จะได้ประโยคไทยล้านนาที่โครงสร้างและลำดับคำที่ถูกต้อง



ภาพที่ 3.1 : สถาปัตยกรรมของการแปลภาษาไทยเป็นภาษาไทยล้านนา

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

3.2 ไวยากรณ์ไม่พึ่งบริบทใช้เทคนิคทางสถิติ (Probabilistic Context Free Grammar : PCFG)

ไวยากรณ์ไม่พึ่งบริบท [9] เป็นส่วนหนึ่งของไวยากรณ์โครงสร้างวลี เพื่อใช้อธิบายโครงสร้างของภาษาและความสัมพันธ์ระหว่างโครงสร้างต่างๆ ที่ประกอบในประโยคนั้นๆ ถูกนิยามโดย

$$G = (V, T, S, P)$$

โดยที่ V คือ เซตของตัวแปร (Variable)

T คือ เซตของสัญลักษณ์เทอร์มินอล (Terminal symbol)

S คือ สัญลักษณ์เริ่มต้น (Start symbol)

P คือ เซตของโปรดักชัน (Production) ซึ่งมีรูปแบบดังนี้ $A \rightarrow a$

โดยที่ A คือ ตัวแปร ($A \in V$) และ a คือสตริงของ ($V \cup T$)

ในงานวิจัยนี้กฎโครงสร้างวลีและโครงสร้างประโยคของภาษาไทย อ้างอิงจากบทความ การวิเคราะห์ประโยคภาษาไทย ของเรืองเดช ปิ่นเขื่อนขันธ์ [7] ซึ่งสามารถสร้างกฎโครงสร้างวลีจำนวน 5 กฎ และกฎย่อยภายใต้โครงสร้างวลีอีก 79 กฎ สามารถแยกตามกฎโครงสร้างวลีเพื่อนำมาสร้างเป็นกฎโครงสร้างประโยคภาษาไทยดังตาราง

| โครงสร้างนามวลี (NP) แบ่งออกเป็น 13 กฎ | |
|--|--------------------------|
| 1) NOUN | 8) NOUN + CLASS + NOUN |
| 2) ADJ | 9) NOUN + CLASS + CLASS |
| 3) NOUN + CLASS | 10) NOUN + CLASS + DETER |
| 4) NOUN + NOUN + CLASS | 11) ADJ + NOUN + NOUN |
| 5) NOUN + ADJ | 12) ADJ + NOUN + CLASS |
| 6) NOUN + NOUN | 13) ADJ + NOUN + DETER |
| 7) NOUN + DETER | |

ตารางที่ 3.1 : ตารางโครงสร้างนามวลี

| โครงสร้างกริยาวลี (VP) แบ่งออกเป็น 41 กฎ | |
|--|-------------------------------------|
| 1) VERB | 22) AUX + VERB + NOUN |
| 2) VERB + VERB | 23) AUX + VERB + PREP |
| 3) VERB + ADV | 24) NEG + VERB + VERB + AUX + ADV |
| 4) VERB + PREP | 25) NEG + VERB + VERB + AUX + NOUN |
| 5) VERB + ENDING | 26) NEG + VERB + VERB + NOUN + ADV |
| 6) VERB + ADV + NOUN | 27) NEG + VERB + VERB + NOUN + NOUN |

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

| | |
|------------------------------|-------------------------------------|
| 7) VERB + PREP + NOUN | 28) NEG + VERB + VERB + PREP + ADV |
| 8) VERB + ENDING + NOUN | 29) NEG + VERB + VERB + PREP + NOUN |
| 9) NEG + VERB | 30) NEG + AUX + VERB + AUX + ADV |
| 10) NEG + AUX | 31) NEG + AUX + VERB + AUX + NOUN |
| 11) AUX | 32) NEG + AUX + VERB + NOUN + ADV |
| 12) NEG + VERB + VERB | 33) NEG + AUX + VERB + NOUN + NOUN |
| 13) NEG + AUX + VERB | 34) NEG + AUX + VERB + PREP + ADV |
| 14) AUX + VERB | 35) NEG + AUX + VERB + PREP + NOUN |
| 15) NEG + VERB + VERB + AUX | 36) AUX + VERB + AUX + ADV |
| 16) NEG + VERB + VERB + NOUN | 37) AUX + VERB + AUX + NOUN |
| 17) NEG + VERB + VERB + PREP | 38) AUX + VERB + NOUN + ADV |
| 18) NEG + AUX + VERB + AUX | 39) AUX + VERB + NOUN + NOUN |
| 19) NEG + AUX + VERB + NOUN | 40) AUX + VERB + PREP + ADV |
| 20) NEG + AUX + VERB + PREP | 41) AUX + VERB + PREP + NOUN |
| 21) AUX + VERB + AUX | |

ตารางที่ 3.2 :ตาราง โครงสร้างกริยาวิ

| | |
|--|----------------|
| โครงสร้างพิเศษวิ (EAP) แบ่งออกเป็น 10 กฎ | |
| 1) PREP + VERB | 6) AUX + ADV |
| 2) PREP + AUX | 7) NOUN + VERB |
| 3) PREP + ADV | 8) NOUN + AUX |
| 4) AUX + VERB | 9) NOUN + ADV |
| 5) AUX + AUX | 10) ADV |

ตารางที่ 3.3 :ตาราง โครงสร้างพิเศษวิ

| | |
|--|------------------------------------|
| โครงสร้างกาลวิเศษณ์วิ (TAP) แบ่งออกเป็น 5 กฎ | |
| 1) NOUN + ENDING | 4) PREP + NOUN + ENDING |
| 2) NOUN + PREP + VERB + ADV | 5) PREP + NOUN + PREP + VERB + ADV |
| 3) PREP + NOUN + ADV | |

ตารางที่ 3.4 :ตาราง โครงสร้างกาลวิเศษณ์วิ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

| โครงสร้างสถานวิเศษณ์วลี (TAP) แบ่งออกเป็น 10 กฎ | |
|---|----------------|
| 1) PREP + VERB | 6) AUX + ADV |
| 2) PREP + AUX | 7) NOUN + VERB |
| 3) PREP + ADV | 8) NOUN + AUX |
| 4) AUX + VERB | 9) NOUN + ADV |
| 5) AUX + AUX | 10) ADV |

ตารางที่ 3.5 : โครงสร้างสถานวิเศษณ์วลี

จากแจกประโยคโดยไวยากรณ์ไม่พึงบริบทจะได้ผลลัพธ์ออกมาเป็นต้นไม้ไวยากรณ์ที่มีได้หลายโครงสร้าง ดังนั้นในงานวิจัยนี้จึงใช้เทคนิคทางสถิติ (Probabilistic Context Free Grammar : PCFG) มาช่วยในการวิเคราะห์โครงสร้างไวยากรณ์ จะมีการกำหนดค่าความน่าจะเป็นสำหรับแต่ละกฎไวยากรณ์ ซึ่งมีสมมุติฐานว่า แต่ละไวยากรณ์มีการเกิดขึ้นแบบอิสระต่อกัน โดยโอกาสเกิดของกฎแต่ละวลีนั้นมีโอกาสเกิดขึ้นเป็น 1

ค่าสถิติของแต่ละกฎสามารถหาได้จากการนับการเกิดขึ้นของกฎนั้นๆ เทียบกับกฎอื่นๆ ที่มีวลีนั้นๆ เป็นชนิดเดียวกันจากคลังประโยคของ Orchid Corpus ซึ่งเป็นคลังข้อมูลประโยคภาษาไทยของศูนย์เทคโนโลยีอิเล็กทรอนิกส์และคอมพิวเตอร์แห่งชาติ มีการกำกับชนิดของคำ และมีโครงสร้างของประโยค ผลลัพธ์จะทำให้ได้โครงสร้างไวยากรณ์ภาษาไทยที่มีค่าความน่าจะเป็นสูงสุด 1 โครงสร้าง

3.3 โครงสร้างพจนานุกรมไทย-ล้านนา

พจนานุกรมคู่ภาษาไทย-ไทยล้านนาสร้างขึ้นจาก พจนานุกรมล้านนา-ไทย ฉบับเฉลิมพระเกียรติพระบาท สมเด็จพระเจ้าอยู่หัว ในวโรกาสเจริญพระชนมายุ 80 พรรษา ของสถาบันภาษาและศิลปวัฒนธรรม มหาวิทยาลัยราชภัฏเชียงใหม่ รวบรวมข้อมูลคำศัพท์ทั้งหมด 19,650 คำ โดยจะจัดเก็บโครงสร้างข้อมูลแบบทรี (Trie)

โครงสร้างข้อมูลแบบทรี (Trie) มีลักษณะคล้ายกับโครงสร้างข้อมูลแบบต้นไม้แต่วิธีการจัดเก็บข้อมูลจะแตกต่างกัน โดยที่โครงสร้างข้อมูลแบบทรีนี้จะจัดเก็บตัวอักษรของคำศัพท์ซึ่งโครงสร้างข้อมูลแบบต้นไม้จะจัดเก็บข้อมูลทั้งคำ โดยโครงสร้างของทรีจะประกอบไปด้วยโหนดต่างๆ จะประกอบไปด้วยดัชนีที่ชี้ไปยังโหนดของตัวอักษรถัดไปซึ่งมีจำนวนดัชนีเท่ากับค่าของจำนวนตัวอักษรที่จะอนุญาตให้มีได้ในพจนานุกรมบวกกับอักษรที่ใช้ระบุเป็นตัวจบคำศัพท์ (Terminator)

สำหรับการสืบค้นใน โครงสร้างข้อมูลแบบทรีนี้จะทำโดยเริ่มต้นที่โหนดศูนย์ถ้าต้องการค้นหาคำศัพท์ก็ให้นำอักษรที่ละตัวจากคำศัพท์ที่ต้องการมาดูว่าภายในโหนด 0 นั้นมีดัชนีของตัวอักษรที่ต้องการไปชี้โหนดอื่นหรือไม่ ถ้าไม่มีแสดงว่าคำนั้น ไม่มีอยู่ในพจนานุกรม แต่ถ้ามีดัชนีเอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ที่ชี้ไปไหนลดก็ให้เดินที่ไหนที่ดัชนีนั้นชี้ไปแล้วนำตัวอักษรตัวถัดไปมาทำตามขั้นตอนแบบเดิมจนหมด เมื่อนำตัวอักษรทั้งหมดจากดัชนีมาเดินในทรีแล้วให้สลับอักษร ถ้าดัชนีมีค่าเท่ากับค่าว่าง (Null) แสดงว่าไม่มีคำศัพท์นั้นในพจนานุกรม แต่ถ้าไม่เท่ากับค่าว่างแสดงว่ามีคำศัพท์นั้นอยู่ในพจนานุกรม โดยดัชนีเป็นตัวชี้ตำแหน่งของข้อมูลของคำนั้น ส่วนความเร็วในการค้นหาภายในโครงสร้างข้อมูลแบบทรีนั้นจะไม่ขึ้นอยู่กับจำนวนคำที่มีในพจนานุกรม แต่จะขึ้นอยู่กับจำนวนดัชนีของตัวอักษร

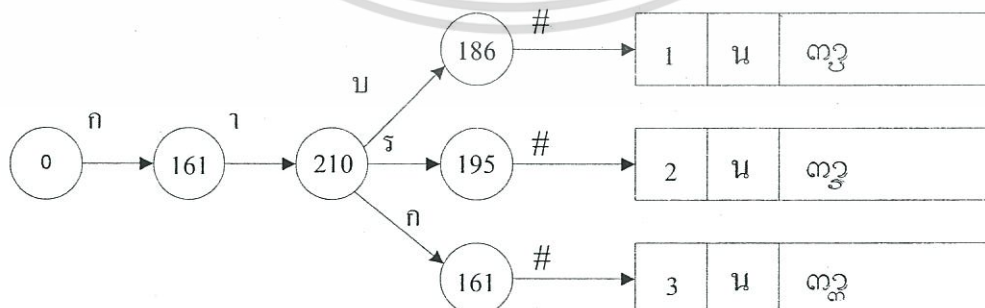
ประเภทของคำในพจนานุกรม ประเภทของคำแบ่งตามหลักภาษาศาสตร์แบ่งออกเป็น 11

ประเภทดังตารางที่ 3.6

| ประเภทคำ | ความหมาย |
|---------------|-----------------|
| Noun | คำนาม |
| Pronoun | คำสรรพนาม |
| Verb | คำกริยา |
| Auxillary | คำวิเศษณ์ |
| Preposition | คำบุพบท |
| Conjunction | คำสันธาน |
| Ending | คำลงท้าย |
| Interjunction | คำอุทาน |
| Prefix | หน่วยคำเติมหน้า |
| Negation | ปฏิเสธ |

ตารางที่ 3.6 : ประเภทของคำ

จากภาพที่ 3.2 เป็นตัวอย่างการจัดเก็บข้อมูลในพจนานุกรมไทย-ไทยล้านนา ด้วยโครงสร้างข้อมูลแบบทรี



ภาพที่ 3.2 : ตัวอย่างโครงสร้างข้อมูลของพจนานุกรมไทย-ไทยล้านนา

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

3.4 การแปลงรูปภาษาไทยล้านนา

การแปลงรูปคำไทยเป็นไทยล้านนา ในขั้นตอนการค้นหาคำศัพท์ในพจนานุกรมจะนำคำภาษาไทยที่ได้จากการตัดคำไปค้นหาคำแปลที่เป็นคู่คำศัพท์ จากพจนานุกรมคู่ภาษาไทย-ภาษาไทยล้านนา หากคำภาษาไทยไม่สามารถจับคู่ได้ จะแปลงคำไทยนั้นเป็นคำไทยล้านนาด้วยกฎการเรียงอักษรไทยล้านนาที่ได้มาจากผู้เชี่ยวชาญและการแปลงรูปภาษาไทยล้านนา

การปรัวรรตอักษรธรรมล้านนาเป็นอักษรไทยกลางดังกล่าว อาจทำได้ 3 ระดับ คือ

1. การปรัวรรตโดยการถ่ายเทียบตัวอักษรและเครื่องหมายทุกชนิด ที่ปรากฏในต้นฉบับ โดยเคร่งครัด
2. การปรัวรรตโดยใช้อักษรและวิธีผสมอักษรแบบไทยกลาง
3. การปรัวรรตโดยใช้อักษรไทยกลางบันทึกคำอ่านในภาษาถิ่นล้านนา

ที่ประชุมข้อตกลงในการปรัวรรตอักษรธรรมล้านนาเป็นอักษรไทยกลาง มติที่ประชุมเชิงปฏิบัติการเรื่องมาตรฐาน การปรัวรรตอักษรพื้นเมืองล้านนาโดย โครงการศูนย์ส่งเสริมศิลปวัฒนธรรม มหาวิทยาลัยเชียงใหม่ ได้สรุปว่า วิธีการปรัวรรตตามแบบแรกนั้น เหมาะสำหรับการปรัวรรตเอกสารที่ต้องการความถี่ถ้วนเป็นพิเศษ เช่น จารึกชนิดต่างๆ การปรัวรรตตามแบบที่สองเหมาะสำหรับการเผยแพร่แก่สาธารณชนผู้สนใจวรรณกรรมล้านนาโดยทั่วไป ส่วนการปรัวรรตตามแบบสุดท้ายนั้น เหมาะที่จะเผยแพร่แก่ประชาชนในท้องถิ่นที่พูดภาษาล้านนา โดยเฉพาะ

กฎที่สร้างขึ้นพิจารณาลำดับของการประสมคำไทยล้านนาด้วยสระต่างๆ เป็นหลักโดยกำหนดสัญลักษณ์ต่างๆ ที่ใช้ในการสร้างกฎดังนี้

- C แทนพยัญชนะต้นหรือพยัญชนะท้าย ในกรณีที่ไม่ทราบแน่นอนว่าพยัญชนะนั้นเป็นตัวสะกดหรือพยัญชนะท้าย
- T แทนวรรณยุกต์
- V แทนสระ
- S แทนพยัญชนะท้ายหรือตัวสะกดในกรณีที่ทราบแน่นอนว่าพยัญชนะนั้นเป็นตัวสะกด หรือพยัญชนะต้น

$[a_1, a_2, a_3, \dots, a_n]$ แทนการเลือกเอาตัวใดตัวหนึ่งระหว่าง $a_1, a_2, a_3, \dots, a_n$

เนื่องจากภาษาไทยและภาษาล้านนามีตัวควบกล้ำ ซึ่งหากใช้ CC แทนพยัญชนะควบกล้ำแล้ว จะทำให้กำกวมได้เพราะในภาษาล้านนา พยัญชนะบางตัวทำหน้าที่เป็นสระด้วย ดังนั้นจึงเพิ่มสัญลักษณ์ควบกล้ำเพื่อลดความกำกวมดังกล่าว

- CR แทนพยัญชนะควบกล้ำ ร
- CL แทนพยัญชนะควบกล้ำ ล
- CW แทนพยัญชนะควบกล้ำ ว
- HC แทนพยัญชนะ ห นำพยัญชนะต้น

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

- AI แทนพยัญชนะ อย
 O แทนสระ ออ ที่มีพยัญชนะท้าย
 Y แทนสระเอีย

| กฎที่ | กฎ | โครงสร้าง | กฎที่ | กฎ | โครงสร้าง |
|-------|--------|-----------|-------|---------|-----------|
| 1 | [CV]ะ | CTVS | 17 | [CT]ัวะ | CTV |
| 2 | [CV]า | CTVS | 18 | [CT]ัว | CTV |
| 3 | [CV]ิ | CTVS | 19 | [CT]ียะ | VCTV |
| 4 | [CV]ี | CTVS | 20 | [CT]ีย | VCTV |
| 5 | [CV]ึ | CTVS | 21 | [CT]ือะ | VCTV |
| 6 | [CV]ื | CTVS | 22 | [CT]ือ | VCTV |
| 7 | [CV]ุ | CTVS | 23 | [CT]ูะ | VCTV |
| 8 | [CV]ู | CTVS | 24 | [CT]ู | VCTV |
| 9 | [CT]ะ | VCTV | 25 | [CT] | VCTS |
| 10 | [CT] | VCT | 26 | [CT] | VCTS |
| 11 | [CT]ะ | VCTV | 27 | [CT]ย | VCTV |
| 12 | [CT] | VCT | 28 | [CT]ยย | VCTV |
| 13 | [CT]ะ | VCTV | 29 | [CT]า | VCTV |
| 14 | [CT] | VCT | 30 | [CT]า | CTV |
| 15 | [CT]าะ | VCTV | 31 | [CT]ั | CTV |
| 16 | [CT]อ | CTV | | | |

ตารางที่ 3.7 : กฎการแปลงรูปภาษาไทยเป็นภาษาล้านนาที่ได้สร้างขึ้น

ในการพิมพ์อักษรภาษาไทยล้านนาในคอมพิวเตอร์ ในอดีตนั้นจะใช้วิธีการพิมพ์ภาษาอังกฤษเข้าไปแทนที่อักขระพิเศษเช่น คำว่า “กิน” ซึ่งมี “น” เป็นตัวสะกด ในการพิมพ์ภาษาไทยล้านนาตัวสะกดจะต้องอยู่ได้คำ ดังนั้นจะต้องเขียนเป็น “กือ” ทำให้ผู้ใช้ต้องสลับโหมดระหว่างภาษาไทยกับภาษาอังกฤษบ่อยครั้ง การพิมพ์จึงล่าช้า และไม่สะดวกเท่าที่ควร ด้วยเหตุนี้ในปี พ.ศ.2549 สำนักส่งเสริมศิลปวัฒนธรรม มหาวิทยาลัยเชียงใหม่ จึงได้ทำโครงการ “พัฒนาระบบการพิมพ์อักษรธรรมล้านนา” และโครงการ “พัฒนาแม่แบบชุดอักษร *Lanna OTF Template*” ขึ้น เพื่อแก้ปัญหาการพิมพ์ดังกล่าว โดยพัฒนาฟอนต์ Ln-tilok และระบบการพิมพ์ให้ใช้ง่าย และมีผลกระทบต่อความเคยชิน ในระบบพิมพ์สัมผัสของผู้ใช้ให้น้อยที่สุด ทั้งนี้

ได้ศึกษาค้นคว้าหาวิธีที่จะทำให้การพิมพ์อักษรธรรมล้านนานั้น อยู่ภายในเป็นพิมพ์เดียว พัฒนา มาจนถึงรุ่น Ln-tilok 6.00 ในปี พ.ศ.2552

ในงานวิจัยการแปลภาษาไทยเป็นภาษาไทยล้านนาในการเขียนภาษาไทยล้านนา ใน คอมพิวเตอร์นั้นได้อ้างอิงวิธีการพิมพ์ตัวอักษรภาษาไทยล้านนาโดยใช้ฟอนต์ Ln-tilok ในการ พิมพ์รูปแบบตัวอักษรไทยล้านนา ซึ่งในการพิมพ์จะใช้อักขระ “๑” และ “๒” มาช่วยในการพิมพ์ เพราะตัวอักษรสองตัวนี้ไม่พบในภาษาไทยล้านนาโดยเมื่อใดก็ตามที่ต้องการพิมพ์ตัวสะกดเช่น

| | | |
|------|---------------------------|-------|
| กิน | เขียนเป็นภาษาไทยล้านนาได้ | กิน๑ |
| ทาน | เขียนเป็นภาษาไทยล้านนาได้ | ทาน๑ |
| การ | เขียนเป็นภาษาไทยล้านนาได้ | การ๑ |
| วัด | เขียนเป็นภาษาไทยล้านนาได้ | วัด๑ |
| ชาติ | เขียนเป็นภาษาไทยล้านนาได้ | ชาติ๑ |

ส่วนสัญลักษณ์ “๒” จะใช้สำหรับสระ โอะลดรูป ตัวอย่างเช่น

| | | |
|----|---------------------------|-----|
| มด | เขียนเป็นภาษาไทยล้านนาได้ | ม๒ด |
| ลด | เขียนเป็นภาษาไทยล้านนาได้ | ล๒ด |
| จด | เขียนเป็นภาษาไทยล้านนาได้ | จ๒ด |

3.5 ไวยากรณ์ปริวรรติเพิ่มพูน (Transformation generative grammar)

ชอมสกี [9] ได้กล่าวถึงไวยากรณ์โครงสร้างวลีและแสดงให้เห็นว่าหลักไวยากรณ์นี้ยังไม่ เพียงพอในการอธิบายภาษาธรรมชาติ จึงได้พัฒนาไวยากรณ์ปริวรรติเพิ่มพูน (Transformation generative grammar) ประโยคที่ผ่านไวยากรณ์ CFG แล้วจะให้ประโยคที่เรียกว่าโครงสร้างลึก โดยปกติได้อยู่ในรูปโครงสร้างต้นไม้ แต่โดยปกติประโยคที่เราใช้ในการพูดเรียกว่าโครงสร้างผิว

กฎการเปลี่ยนโครงสร้างวลีให้เป็นโครงสร้างผิวโดยการเพิ่มหรือลดตำแหน่งของคำซึ่ง อาจจะแบ่งได้เป็นกฎบังคับที่ใช้ในการเปลี่ยนรูปโครงสร้างเพราะถ้าไม่เปลี่ยนจะทำให้ประโยค ผิดหลักไวยากรณ์ อีกประเภทเป็นกฎที่ใช้หรือไม่ใช้ก็ได้ ประโยคจะไม่เปลี่ยนแปลงตามกฎก็ได้

กฎการปริวรรตจะแบ่งออกเป็น 3 ส่วนคือ

- ส่วนอธิบายโครงสร้าง (Structural description : SD) เพื่อแสดงโครงสร้างของประโยค ว่าควรเป็นอย่างไรถึงใช้กฎนี้

- ส่วนแสดงการเปลี่ยนแปลง (Structural Change :SC) เพื่อแสดงโครงสร้างของประโยค ที่เปลี่ยนแปลงไปหลังจากที่ปริวรรตแล้ว

- ส่วนของชุดเงื่อนไขที่ต้องทำการเปรียบเทียบให้ตรงกัน

ตัวอย่างกฎไวยากรณ์ปริวรรตเพิ่มพูน ของภาษาไทยเป็นภาษาไทยล้านนา ประโยค “ส้มตำอร่อยมาก” ในภาษาไทยสามารถแปลเป็นภาษาไทยล้านนาได้มากกว่า 1 โครงสร้าง คือ “ตำส้มจ้ำดล้า” และ “ตำส้มล้าขนาด”

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สั้นคำ/N อร่อย/V มาก/Adv



คำสั้น/N จ้าค/Adv คำ/V

คำสั้น/N คำ/V ขนาด/Adv

กฎการไวยกรณ์ปริวรรตของกริยาช่วยของคำว่า มาก (Adv) เมื่อไวยกรณ์ที่ผ่านเข้ามายังไวยกรณ์ปริวรรตเมื่อพบคำว่า “มาก” ที่มีคู่ภาษาล้านนาเป็น “จ้าค” ในภาษาไทยล้านนา จะปรับตำแหน่งประเภทของคำที่เป็นกริยาช่วย ไปอยู่ด้านหน้าของกริยา ดังต่อไปนี้

SD : NP V Adv

[มาก]

SC : NP Adv V

[จ้าค]

ผลลัพธ์ : จะได้โครงสร้างใหม่เพิ่มและโครงสร้างเดิม คือ

- คำสั้น/N จ้าค/Adv คำ/V

- คำสั้น/N คำ/V ขนาด/Adv

3.6 แบบจำลองฮิดเดนมาร์คอฟและการลดความกำกวมด้วยวิเทอร์บีอัลกอริทึม

แบบจำลองฮิดเดนมาร์คอฟ [10] เป็นเทคนิคการเรียนรู้ทางสถิติ (Statistical Machine Learning) ที่ใช้จำแนกรูปแบบของลำดับเหตุการณ์หนึ่งๆ ซึ่งลำดับของเหตุการณ์จะถูกจำลองให้อยู่ในรูปแบบของลำดับการเปลี่ยนแปลงของสถานะ (State) จะมีค่าความน่าจะเป็นอยู่ 2 ชนิด คือ ค่าความน่าจะเป็นของการเปลี่ยนสถานะ และค่าความน่าจะเป็นของการทำให้เกิดผลลัพธ์ เมื่อมีการเปลี่ยนสถานะเกิดขึ้นและผลลัพธ์นี้จะถูกนำไปใช้ในการตัดสินใจว่าข้อมูลที่เข้ามามีโอกาสเป็นอะไรมากที่สุด

วิเทอร์บีอัลกอริทึมเป็นอัลกอริทึมหนึ่งในแบบจำลองฮิดเดนมาร์คอฟ ที่ใช้ในการจัดตำแหน่งลำดับของสถานะต่างๆ โดยลำดับของแต่ละสถานะจะขึ้นอยู่กับความน่าจะเป็นสูงสุด ที่จะเป็นไปได้ทั้งหมด

การทำงานของวิเทอร์บีอัลกอริทึม มี 3 ขั้นตอน คือ

ขั้นตอนที่ 1 คำนวณค่าความน่าจะเป็นของสถานะเริ่มต้นของข้อมูลนำเข้า

ขั้นตอนที่ 2 คำนวณค่าความน่าจะเป็นของสถานะถัดไป ทุกสถานะ เพื่อจะเก็บเส้นทางที่ดีที่สุด สำหรับการเลือกเส้นทางในขั้นตอนถัดไป

ขั้นตอนที่ 3 การย้อนเส้นทาง (Backtracking) เพื่อค้นหาลำดับความน่าจะเป็นสูงสุด จากสถานะที่ได้คำนวณไว้แล้ว ซึ่งจะได้อันดับเส้นทางสถานะที่เป็นไปได้สูงสุดด้วย

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ในการลดค่ากำกวมในการแปลภาษาไทยเป็นภาษาไทยล้านนา ซึ่งจะลดความกำกวมของโครงสร้างประโยค และลำดับของคำศัพท์โดยใช้แบบจำลองฮิดเดินมาร์คอฟ และวิเทอร์บีอัลกอริทึม(Viterbi Algorithm) ซึ่งวิเทอร์บีอัลกอริทึมจะมี ขั้นตอนดังภาพที่ 3.3

```

For i=1 to N do
  SEQSCORE(i,1)=  $PROB(w_1 | L_i) * PROB(L_i | \emptyset)$ 
  BACKPTR(i,1)=0

For t=2 to T
  For i=1 to N
    SEQSCORE(i,t)=  $Max_{j=1,N} (SEQSCORE(j,t-1) * PROB(L_t | L_j)) * PROB(w_t | L_i)$ 
    BACKPTR(i,t)=index of j that gave the max above

C(T)=i that maximizes SEQSCORE(i,T)
For i=T-1 to 1 do
  C(i)=BACKPTR(C(i+1),i+1)
  
```

ภาพที่ 3.3 : วิเทอร์บีอัลกอริทึม(Viterbi Algorithm) สำหรับการลดค่ากำกวมในการแปลภาษาไทยเป็นภาษาไทยล้านนา [11]

จากอัลกอริทึมมีตัวแปรที่เกี่ยวข้องดังนี้

- T คือ จำนวนคำในโครงสร้างประโยคภาษาไทยล้านนาที่มีการระบุชนิดของคำ ที่ต้องการเลือกโครงสร้าง
- N คือ จำนวนชนิดของคำภาษาไทยล้านนาที่พบในประโยคที่ต้องการแปล
- SEQSCORE คือ ระเบียบขนาด $N * T$ ที่เก็บค่าความน่าจะเป็นชุดหนึ่งที่มีการเปลี่ยนสถานะ(คำ) ตั้งแต่สถานะเริ่มต้นไปจนถึงสถานะ(คำ) ปัจจุบัน
- BACKPTR คือ ระเบียบขนาด $N * T$ ที่เก็บเส้นทางที่มีค่าความน่าจะเป็นสะสมสูงสุด
- w คือ คำในโครงสร้างประโยคภาษาไทยล้านนา มีการระบุชนิดของคำไทยล้านนาที่ต้องการแปล
- L คือ ชนิดของคำภาษาไทยล้านนาในประโยคที่ต้องการแปล
- $PROB(w_i | L_j)$ คือ ความน่าจะเป็นของคำภาษาไทยล้านนา w_i ที่มีชนิดของคำภาษาไทยล้านนาเป็น L_j ซึ่งคำนวณจากจากสูตร

$$\text{PROB}(w_i|L_j) \cong \frac{\text{Count}(w \text{ at position } i \text{ and } L \text{ at position } j)}{\text{Count}(L \text{ at position } j)}$$

- $\text{PROB}(L_i|L_j)$ คือ ความน่าจะเป็นที่เมื่อเกิดชนิดของคำภาษาไทยสั้นนา L_j แล้วจะเกิดชนิดของคำภาษาไทยสั้นนา L_i ตามหลังซึ่งคำนวณจากสูตร

$$\text{PROB}(L_i|L_j) \cong \frac{\text{Count}(L \text{ at position } i \text{ and } L \text{ at position } j)}{\text{Count}(L \text{ at position } j)}$$

- C คือ ลิสต์ที่มีขนาดเท่ากับจำนวนคำในประโยคภาษาไทยสั้นนาที่ต้องการเลือกโครงสร้าง ซึ่งจะเก็บลำดับของคำภาษาไทยสั้นนาที่มีค่าความน่าจะเป็นสูงสุด คลังข้อมูลที่นำมาคำนวณค่าสถิติ จะเตรียมคลังข้อมูลภาษาไทยสั้นนาที่มีการกำกับชนิดของคำ ในประโยคไทยสั้นนา วิเทอร์บีอัลกอริทึมจะคำนวณเพื่อหาโครงสร้างของประโยคสั้นนาทั้งหมดทุกโครงสร้าง จากขั้นตอนของการปรับโครงสร้างสั้นนาในไวยากรณ์ปริวรรตก่อนหน้า ผลลัพธ์ของขั้นตอนนี้จะได้ลำดับของคำในแต่ละโครงสร้างประโยคสั้นนา และจะนำค่าสถิติที่คำนวณได้ มาเลือกโครงสร้างที่มีความน่าจะเป็นสูงสุด จะทำให้ได้โครงสร้างที่เป็นผลลัพธ์ของการแปลภาษาไทยเป็นภาษาไทยสั้นนา

3.7 การทดลองและการวัดประเมินผลการทดลอง

เนื่องจากภาษาไทยและภาษาไทยสั้นนามีความใกล้เคียงกันในลักษณะโครงสร้างของภาษาและคำที่ใช้งานในประโยค เพราะภาษาไทยสั้นนาเป็นภาษาถิ่น ดังนั้นในการออกแบบการทดลองเพื่อวัดประสิทธิภาพของการแปลจึงมีการออกแบบการทดลอง วิธีการเพื่อตรวจสอบผลลัพธ์ของการแปล ได้แก่

การแปลโดยใช้ไวยากรณ์ปริวรรต วิเคราะห์และจัดรูปแบบโครงสร้างของประโยคสั้นนาก่อนเข้าสู่กระบวนการเลือกโครงสร้างและลดความกำกวมของคำศัพท์โดย วิเทอร์บีอัลกอริทึม

การวัดประเมินผลการทดลอง ผลจากการทดลองจะนำมาวัดประสิทธิภาพการแปล โดยการคำนวณค่าร้อยละ จากประโยคที่แปลถูกต้องต่อประโยคที่นำมาทดสอบทั้งหมด ดังสมการ

$$\text{ประสิทธิภาพการแปล (\%)} = \frac{X}{N} \times 100$$

เมื่อกำหนดให้

ตัวแปร X คือ จำนวนของประโยคที่แปลถูกต้อง

ตัวแปร N คือ จำนวนของประโยคที่นำมาทดสอบทั้งหมด

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 4

การประเมินผลความถูกต้องของการแปล

จุดประสงค์ของงานวิจัยนี้คือการแปลภาษาไทยเป็นภาษาไทยล้านนา โดยใช้วิธีการเลือกความหมายจากพจนานุกรมไทยล้านนา วิเคราะห์โครงสร้างประโยคภาษาไทยด้วยไวยากรณ์ไม่พึ่งบริบท ประยุกต์เทคนิคทางสถิติ (Probabilistic Context free grammar : PCFG) ปรับโครงสร้างไวยากรณ์ให้เป็นโครงสร้างไวยากรณ์ไทยล้านนาด้วย ไวยากรณ์ปริวรรต (Transformation Grammar) และลดความกำกวมด้วยวิธีการของแบบจำลองฮิดเดินมาร์คอฟและวิเทอร์บีอัลกอริทึม (Viterbi algorithm)

ข้อมูลที่น่ามาทดสอบเป็นข้อมูลที่รวบรวมจากข้อมูลที่มีการบันทึกไว้ในเอกสาร โดยผู้เชี่ยวชาญภาษาไทยล้านนา ซึ่งได้นำมารวบรวม ซึ่งการที่นำข้อมูลชุดนี้มาทำการทดสอบเนื่องจากเป็นข้อมูลที่น่าสนใจและสามารถตรวจสอบผลลัพธ์ความถูกต้องของประโยคภาษาไทยล้านนาที่เป็นผลลัพธ์

ตัวอย่างข้อมูลที่น่ามาทำการทดลองประกอบไปด้วย นิทานพื้นบ้าน 3 เรื่อง จำนวน 121 ประโยค คำประโยคทดสอบที่ถูกต้องตามไวยากรณ์ภาษาไทยจำนวน 85 ประโยค

4.1 ผลการวิจัย

การทดสอบโดยการนำประโยคทดสอบเข้าสู่กระบวนการตัดคำโดยใช้เครื่องมือตัดคำ (KU-CUT) ซึ่งเป็นเครื่องมือที่ตัดคำตัดคำด้วยเทคนิคการเรียนรู้จากข้อมูลแบบไม่ใช้ตัวอย่างร่วมกับการใช้พจนานุกรมและคำ จากนั้นทำการแปลโดยค้นหาคำแปลภาษาไทยล้านนาจากพจนานุกรมไทยล้านนา หากในกรณีที่ไม่มีพบคำแปลในพจนานุกรมจะใช้วิธีการปริวรรตคำไทยล้านนาตามกฎของโครงสร้างของการเขียนภาษาไทยล้านนา จากนั้นนำเข้าสู่ไวยากรณ์ปริวรรตปรับให้เป็นโครงสร้างประโยคไทยล้านนา จากนั้นจะลดความกำกวมของโครงสร้างประโยคและลำดับของคำล้านนา โดยวิธีการวิเทอร์บีอัลกอริทึม ผลลัพธ์ก็จะได้ประโยคไทยล้านนาที่โครงสร้างและลำดับคำที่ถูกต้อง

การวัดประสิทธิภาพของการแปล โดยการคำนวณค่าร้อยละ จากประโยคที่แปลถูกต้องต่อประโยคที่น่ามาทดสอบทั้งหมด ดังสมการ

$$\text{ประสิทธิภาพการแปล (\%)} = \frac{X}{N} \times 100$$

เมื่อกำหนดให้

ตัวแปร X คือ จำนวนของประโยคที่แปลถูกต้อง

ตัวแปร N คือ จำนวนของประโยคที่น่ามาทดสอบทั้งหมด

ผลลัพธ์จากการวัดประสิทธิภาพของการแปลภาษาไทยเป็นภาษาไทยล้านนาของงานวิจัยนี้ ได้แสดงผลลัพธ์ดังตารางที่ 4.1

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

| ชื่อเอกสาร | จำนวน ประโยค | แปลถูกต้อง | อัตราส่วน ร้อยละ |
|---|-----------------|------------|---------------------|
| นิทานพื้นบ้าน เรื่อง ปู่กับนกยาง | 29 | 21 | 72.41 |
| นิทานพื้นบ้าน เรื่อง คนไม่รู้บุญคุณ | 27 | 20 | 74.07 |
| นิทานพื้นบ้าน เรื่อง กินข้าวอร่อยทุกมือ | 65 | 51 | 78.46 |
| ประโยคทดสอบ | 85 | 70 | 87.5 |
| รวมทั้งสิ้น | 206 | 162 | 78.64 |

ตารางที่ 4.1 : ประสิทธิภาพของการแปลภาษาไทยเป็นภาษาไทยล้านนา

4.2 การอภิปรายผล

จากการแปลภาษาไทยเป็นภาษาไทยล้านนาด้วยเลือกความหมายจากพจนานุกรมไทยล้านนา วิเคราะห์โครงสร้างประโยคภาษาไทยด้วยไวยากรณ์ไม่พึ่งบริบท ประยุกต์เทคนิคทางสถิติ (Probabilistic Context free grammar : PCFG) ปรับโครงสร้างไวยากรณ์ให้เป็นโครงสร้างไวยากรณ์ไทยล้านนาด้วย ไวยากรณ์ปริวรรต (Transformation Grammar) และลดความกำกวมด้วยวิธีการของแบบจำลองฮิดเดินมาร์คอฟและวิเทอร์บีอัลกอริทึม (Viterbi algorithm) ผลลัพธ์จะได้ประโยคภาษาไทยล้านนาซึ่งเป็นภาษาเป้าหมาย จากการทดลองพบว่า จากประโยคทั้งสิ้น 206 ประโยคที่นำมาทดลอง มีประสิทธิภาพในการแปล 78.64%

จากการทดลองพบว่า การแปลภาษาไทยเป็นภาษาไทยล้านนาด้วยวิธีเลือกคำศัพท์จากพจนานุกรมนั้น คำไทย 1 คำมักจะแปลเป็นคำไทยล้านนาได้หลายคำ ทำให้ส่งผลกระทบต่อที่จะนำคำศัพท์มาสร้างประโยคที่เป็นไปได้ทั้งหมด เมื่อเข้าสู่ขั้นตอนการลดความกำกวมด้วยวิเทอร์บีอัลกอริทึมก็จะใช้เวลาการประมวลผลค่อนข้างนานเพราะจะต้องคำนวณหาค่าทางสถิติทุกประโยคเพื่อเรียงประโยคผลลัพธ์ในการแปลด้วย

อย่างไรก็ตามในการปรับโครงสร้างไวยากรณ์ไทยให้เป็นโครงสร้างไวยากรณ์ไทยล้านนาด้วยไวยากรณ์ปริวรรตเพิ่มพูน(Thai-Thai lanna transformation rule) นั้นจะทำให้ประโยคบางประโยคมีความสมบูรณ์มากยิ่งขึ้น ตัวอย่างเช่น

ตัวอย่างผลลัพธ์ของการแปลที่ผ่านไวยากรณ์ปริวรรต

วลี “อร่อยมาก” จะประกอบด้วยคำว่า “อร่อย” ที่มีคำแปลภาษาไทยล้านนาคือคำว่า ลำ(๗๖) และคำว่า “มาก” มีคำแปลภาษาไทยล้านนา หลายคำแปล คือ ขนาด(๕๒๕) , นัก(๕๓) , ซาด(๕๔) ซึ่งเมื่อแปลออกมาแล้วจะได้ผลลัพธ์คือ

ลำขนาด(๗๖ ๕๒๕)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ล้านัก(๗๖ ๙๓)

ลำชาด(๗๖ ๙๔)

ซึ่งเมื่อปรับโครงสร้างด้วย ไวยากรณ์ปริวรรต ที่ตรงกับกฎของไวยากรณ์ปริวรรตแล้ว จะทำให้คำว่า ลำชาด(๗๖ ๙๔) มีการสลับตำแหน่งของคำเป็น ชาดลำ(๙๔ ๗๖) ซึ่งเป็นผลลัพธ์คำแปลที่ถูกต่อนั่นเอง

สาเหตุของความผิดพลาดของการแปลสามารถจำแนกสาเหตุได้ดังนี้

1. ในการตัดคำและกำกับชนิดของคำ มีผลทำให้การแปลผิดพลาดเช่น ประโยค “ป่าแก้ว ทำกับข้าวเลี้ยงแขก” เมื่อคำว่า “แก้ว” คือคำนาม เมื่อค้นหาในพจนานุกรมแล้วได้ คำแปลของคำว่า แก้วคือ รัตน จึงได้ประโยคนี้เป็น “ป่า ๑๓๑ ๑๑๑๓ ๓๖๖๐ ๗๖๑ ๙๙๓” อ่านว่า ป่ารัตน เนรมิตรกับข้าวเลี้ยงแขก

2. ในขั้นตอนการแปลงรูปเป็นคำไทยล้านนาหากไม่พบในพจนานุกรม มีความผิดพลาดเนื่องจากคำไทยคำนั้นมีโครงสร้างที่ซับซ้อนเช่นคำว่า “เนิบนาบ” จะพบว่าคำว่าเนิบมีตัวสะกดคือ “บ” แต่คำถัดไปคือ “น” จากคำว่านาบ ซึ่งทำให้เกิดความผิดพลาดในการแปลงรูปเพราะเจตนาของการแปลคำว่า “บน”

3. ในขั้นตอนของการคำนวณด้วยวิเทอบัณฑิตอริทิม หากคำศัพท์และชนิดของคำดังกล่าว ไม่มีใน คลังข้อมูลของประโยคภาษาไทยล้านนาที่ใช้ในการทดสอบ จะทำให้ส่งผลต่อผลลัพธ์ของการแปลที่ควรจะเป็น และทำให้การแปลความหมายมีความคลาดเคลื่อนกับประโยคที่ควรจะเป็นได้

4.3 อุปสรรคในการทำงาน

อุปสรรคบางประการส่งผลให้การแปลผิดพลาด อุปสรรคดังกล่าวได้แก่

1. ความไม่ครอบคลุมของคำศัพท์ในพจนานุกรม การค้นหาคำศัพท์ในพจนานุกรม คำไทยล้านนาบางคำศัพท์เป็นคำศัพท์เฉพาะ ทำให้เมื่อนำไปแทนที่คำไทยมีความหมายที่ผิดเพี้ยนไป คำไทยบางคำก็มีคำศัพท์ไทยล้านนาเพียง 1 คำที่เป็นคำเฉพาะ แต่เมื่อแทนความหมายของคำศัพท์นั้นในประโยค ประโยคดังกล่าวก็มีความผิดเพี้ยนในความหมาย

2. คำไทยบางคำแปลเป็นคำไทยล้านนาได้มากเกินไป คำศัพท์ไทยล้านนาที่มีมากเกินไปทำให้ส่งผลต่อการคำนวณที่จะใช้เวลาในการแปลประโยคนั้นนาน เช่น ประโยค “พ่อกินข้าว” พบว่า

คำว่า พ่อ มีคำศัพท์ 1 คำในพจนานุกรมคือ พ่อ (๕)

คำว่า กิน มีคำศัพท์ 7 คำในพจนานุกรมคือ กิน (๙), ญูชะ (๖๖๖=), ย้า (๖๖๖), รูดเนา (๑๕๑๕), เสพ (๘๘), เสวย (๘๘๖), ฉั่น (๖)

คำว่า คำว่า ข้าว มีคำศัพท์ 2 คำในพจนานุกรมคือ ภัตตะ (๖๖๖), เข้า (๘๘๖)

เมื่อนำมาสร้างประโยค ก็จะมีประโยคทั้งหมด 14 ประโยคได้แก่

- ດູ່ ທັງ
- ດູ່ ທຽນ = ທັງ
- ດູ່ ພວ ທັງ
- ດູ່ ອຸ ຊຸ ດູ່ ທັງ
- ດູ່ ລູ ທັງ
- ດູ່ ລຽງ ທັງ
- ດູ່ ອຸ ທັງ
- ດູ່ ທັງ ເຊັ່ນ
- ດູ່ ທຽນ = ເຊັ່ນ
- ດູ່ ພວ ເຊັ່ນ
- ດູ່ ອຸ ຊຸ ດູ່ ເຊັ່ນ
- ດູ່ ລູ ເຊັ່ນ
- ດູ່ ລຽງ ເຊັ່ນ
- ດູ່ ອຸ ເຊັ່ນ

ດັ່ງນັ້ນຈຳນວນປະໂຫຍດ ທີ່ຄຳນວນທັງໝົດຈະເປັນຜົນຄູນຂອງ ຄຳໄທຍ່າງທີ່ພົບໃນແຕ່ລະຄຳໄທ ຈຶ່ງໃຫ້ການຄຳນວນຕ້ອງທຳທຸກປະໂຫຍດສ່ວນໃຫຍ່ໃຫ້ເປັນ

3. ການຂາດແຄບຄຳຂໍ້ມູນພາສາໄທຍ່າງສ້າງໄດ້ຍາກ ການລວມຄຳກວມດ້ວຍວິທີການຂອງວິເທອຣ໌ບີອັດຄອຣ໌ທິມ ຈຳເປັນຕ້ອງພິຈາລະນາ ຄຳຂໍ້ມູນພາສາ ທີ່ມີການສຶກສາແລະເຂົ້າໝາຍກັນເພາະກຸ່ມ ຈຶ່ງໃຫ້ເຂົ້າສ່ວນໃນການສ້າງຄຳຂໍ້ມູນພາສາໄທຍ່າງເປັນໄປດ້ວຍຄວາມລຳບາກ ຂໍ້ຈຳກັດດັ່ງກ່າວໃຫ້ປະໂຫຍດທີ່ຕ້ອງການຂໍ້ມູນຈາກຄຳຂໍ້ມູນພາສາໄທຍ່າງ ໃນການລວມຄຳກວມນັ້ນ ໄດ້ຮັບຄຸນຄ່າໃນປະໂຫຍດບາງປະໂຫຍດໄດ້ໃຫ້ສ່ວນຕໍ່ຜົນສຳເລັດຂອງການປ່ຽນພາສາໄທຍ່າງ

ເອກສານນີ້ເປັນເອກສານທີ່ສ່ວນໄວ້ສຳຮັບການໃຊ້ງານເພື່ອການສຶກສາເທົ່ານັ້ນ ໄດ້ອະນຸຍາດໃຫ້ນຳໄປໃຊ້ປະໂຫຍດດ້ານການຄ້າ ໄດ້ຮັບຄຸນຄ່າ ທັງໝົດ ອີກທັງຫ້າມມີໃຫ້ແກ້ໄຂ ແລະ ຕ້ອງອ້າງອິງເຊິ່ງເຈົ້າຂອງເອກສານທຸກຄັ້ງທີ່ມີການນຳໄປໃຊ້

บทที่ 5

สรุปผลการวิจัยและข้อเสนอแนะ

5.1 สรุปผลการวิจัย

วิทยานิพนธ์ฉบับนี้ได้นำเสนอการแปลภาษาไทยเป็นภาษาไทยล้านนา โดยใช้ความหมายจากพจนานุกรมไทยล้านนา วิเคราะห์โครงสร้างประโยคภาษาไทยด้วยไวยากรณ์ไม่พึ่งบริบท ประยุกต์เทคนิคทางสถิติ (PCFG) ปรับโครงสร้างไวยากรณ์ให้เป็นโครงสร้างไวยากรณ์ไทยล้านนาด้วยไวยากรณ์ปริวรรต (Transformation Grammar) และลดความกำกวมด้วยวิธีการของแบบจำลองฮิดเดินมาร์คอฟและวิเทอบีอัลกอริทึม (Viterbi algorithm)

ภาษาไทยและภาษาไทยล้านนามีโครงสร้างไวยากรณ์ที่ใกล้เคียงกัน แต่ก็มีบางโครงสร้างที่มีความแตกต่าง ไวยากรณ์ปริวรรต (Transformation Grammar) จะช่วยทำให้การแปลภาษาไทยเป็นภาษาไทยล้านนามีความถูกต้องและใกล้เคียงกับคำแปลภาษาไทยล้านนามากยิ่งขึ้น

ข้อมูลที่ใช้ในการวัดประสิทธิภาพของการแปลภาษาไทยเป็นภาษาไทยล้านนาประกอบด้วยนิทานจำนวน 3 เรื่อง 121 ประโยค, ประโยคทดสอบจำนวน 85 ประโยค พบว่ามีประสิทธิภาพของการแปล 78.64%

โดยมีสาเหตุที่ทำให้การแปลมีความผิดพลาดคือ การตัดคำและกำกับชนิดของคำมีผลทำให้การแปลผิดพลาด, ขั้นตอนการแปลงรูปเป็นคำไทยล้านนา มีความผิดพลาดเนื่องจากคำไทยคำนั้นมีโครงสร้างที่ซับซ้อน, การคำนวณด้วยวิเทอบีอัลกอริทึม หากคำศัพท์และชนิดของคำไม่มีในคลังข้อมูลของประโยคภาษาไทยล้านนา จะทำให้ส่งผลต่อผลลัพธ์ของการแปล

5.2 ข้อเสนอแนะ

จากปัญหาที่พบในงานวิจัย เนื่องจากภาษาเป้าหมายนั้นคือภาษาไทยล้านนา เป็นภาษาถิ่นเหนือ ซึ่งแม้ว่าจะมีการเผยแพร่องค์ความรู้ภาษาไทยล้านนามากยิ่งขึ้นในปัจจุบัน แต่ก็ยังถือว่าเป็นวงแคบ และยังมีผู้ที่เชี่ยวชาญค่อนข้างจำกัด ส่งผลถึงเอกสาร และข้อมูลที่เกี่ยวข้องกับ ภาษาไทยล้านนาที่ออกสู่สาธารณะค่อนข้างที่จะมีน้อย ทำให้ในบางกระบวนการที่ต้องการข้อมูลเหล่านั้นเช่น การสร้างคลังข้อมูลภาษาทำได้ค่อนข้างยากดังนั้นควรส่งเสริมให้มีการเรียนรู้ภาษาไทยล้านนาในวงกว้างเพื่อจะได้มีข้อมูล เอกสารต่างๆ ออกสู่สาธารณะต่อไป

ส่วนของโปรแกรมในการสร้างประโยคที่เป็นไปได้ทั้งหมด หากพบว่าคำศัพท์ที่ได้จากการค้นหาในพจนานุกรมปริมาณมากจะทำให้ใช้เวลาในการประมวลผลค่อนข้างนาน เนื่องจากต้องประมวลผลคำนวณข้อมูลทางสถิติของทุกประโยคตามไปด้วย

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ด้วยทรัพยากรคลังภาษาไทยและภาษาไทยล้านนา ก่อนข้างมีจำกัดและไม่ได้เจาะจงเฉพาะขอบเขตเฉพาะด้าน ทำให้เกิดปัญหาในข้อมูลบางอย่าง ทำให้เมื่อนำมาใช้ในการวิเคราะห์โครงสร้างประโยคภาษาไทยด้วยไวยากรณ์ไม่พึ่งบริบท ประยุกต์เทคนิคทางสถิติ และการลดความกำกวมด้วยวิธีการของแบบจำลองฮิดเดินมาร์คอฟและวิเทอร์บีอัลกอริทึม ไม่ครอบคลุมดังเช่นสภาพแวดล้อมของการใช้งานจริง ดังนั้นหากทรัพยากรคลังข้อมูลทั้งภาษาไทยและภาษาไทยล้านนามีการรวบรวมที่เฉพาะเจาะจง และมีปริมาณที่ครอบคลุมทุกบริบท คาดว่าจะทำให้ผลลัพธ์ของการแปลภาษาไทยเป็นภาษาไทยล้านนามีความถูกต้องมากยิ่งขึ้น



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บรรณานุกรม

- [1] ศรีเนียน สวัสดิ์. ตำราเรียนหนังสือ ภาษาไทยล้านนา เชียงใหม่ : ร้านประเทืองวิทยา, 2540
- [2] พระมหามิลินท์ อนุคาริโก. คู่มือเรียนภาษาไทยล้านนา 24 ชั่วโมงด้วยตนเอง. ค้นเมื่อ 24 กุมภาพันธ์ 2554 , จาก <http://www.bmwit.net/word-pdf/Test-003.pdf>
- [3] อัสนีย์ ก่อตระกูล. การประมวลภาษามนุษย์ด้วยคอมพิวเตอร์. กรุงเทพฯ: หน่วยปฏิบัติการวิจัยเชี่ยวชาญเฉพาะการประมวลผลภาษาธรรมชาติและเทคโนโลยีสารสนเทศอัจฉริยะ ภาควิชาวิศวกรรมศาสตร์ คณะวิศวกรรมศาสตร์ มหาวิทยาลัยเกษตรศาสตร์, 2549.
- [4] พุชยดี ศิริแสงตระกูล และ นิกร ยาพรหม, “การแปลภาษาล้านนาเป็นภาษาไทยกลาง”, วารสารวิจัย มหาวิทยาลัยขอนแก่น, มีนาคม 2552 หน้า 264-274.
- [5] โสภัน พรหมโสดา และ พุชยดี ศิริแสงตระกูล, “การแปลงอักษรธรรมอีสานเป็นภาษาไทย”, The 7th National Conference on Computing and Information Technology (NCCIT2011), พฤษภาคม 2554 หน้า 503-509.
- [6] Ola Mohammad Ali, Mahmoud GadAlla, and Mohammad Said Abdelwahab, “Improving Machine Translation using Hybrid Dictionary-Graph Based Word Sense Disambiguation with Semantic and Statistical Methods” *International Journal of Computer and Electrical Engineering*, Vol.1, No.5, 2009, pp.1793 –8163.
- [7] เรื่องเดช ปันเขื่อนขัตติย์. การวิเคราะห์ประโยคภาษาไทย. สารภาษาไทยและวัฒนธรรมไทย ภาควิชาภาษาไทย คณะมนุษยศาสตร์ มหาวิทยาลัยเกษตรศาสตร์, 2550
- [8] ทรงศักดิ์ ปรางค์วัฒนากุล, บรรณาธิการ. ข้อตกลงในการปริวรรตอักษรธรรมล้านนาเป็นอักษรไทยกลาง. เชียงใหม่: โรงพิมพ์ช้างเผือก; 2529.
- [9] ยืน ภู่วรรณ และ ชัยยงค์ วงศ์ชัยสุวัฒน์. การประมวลผลภาษาธรรมชาติ(Natural Language Processing). กรุงเทพฯ: ศูนย์เทคโนโลยีอิเล็กทรอนิกส์ และคอมพิวเตอร์แห่งชาติ กระทรวงวิทยาศาสตร์และเทคโนโลยีและสิ่งแวดลอม; 2535. หน้า 76-84
- [10] ยศวิธน์ แสนสิงห์ และ อัครา ประโยชน์, “การรู้จำชื่อเฉพาะภาษาไทยโดยใช้แบบจำลองฮิดเดน มาร์คอฟ”, *Proceedings of the Second Conference on Knowledge and Smart Technologies 2010 (KST-2010)*, Chonburi Thailand, 24-25 July 2010.

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บรรณานุกรม(ต่อ)

- [11] Allen James, **Natural language understanding**, 2nd ed, 1994 , pp.202



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



ภาคผนวก

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



ภาคผนวก ก
ตัวอย่างผลลัพธ์ของการทดสอบการแปล

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



THE 7TH NATIONAL CONFERENCE ON COMPUTING AND INFORMATION TECHNOLOGY

PROCEEDINGS OF NCCIT 2011

THE 7TH NATIONAL CONFERENCE ON COMPUTING AND INFORMATION TECHNOLOGY

11-12 MAY 2011

WWW.NCCIT.NET

FACULTY OF INFORMATION TECHNOLOGY
KING MONGKUT'S UNIVERSITY OF TECHNOLOGY NORTH BANGKOK

VOLUME 2



บทความวิจัย

การประชุมทางวิชาการระดับชาติด้านคอมพิวเตอร์และเทคโนโลยีสารสนเทศ
ครั้งที่ 7

11-12 พฤษภาคม 2554



คณะเทคโนโลยีสารสนเทศ
มหาวิทยาลัยเทคโนโลยีพระจอมเกล้าพระนครเหนือ



THE 7TH NATIONAL CONFERENCE
ON COMPUTING AND INFORMATION
TECHNOLOGY

PROCEEDINGS OF NCCIT 2011

THE 7TH NATIONAL CONFERENCE ON COMPUTING AND INFORMATION TECHNOLOGY

11-12 MAY 2011

FACULTY OF INFORMATION TECHNOLOGY
KING MONGKUT'S UNIVERSITY OF TECHNOLOGY NORTH BANGKOK (KMUTNB)
BANGKOK, THAILAND

WWW.NCCIT.NET



บทความวิจัย

การประชุมทางวิชาการระดับชาติด้านคอมพิวเตอร์และเทคโนโลยีสารสนเทศ
ครั้งที่ 7

11-12 พฤษภาคม 2554



คณะเทคโนโลยีสารสนเทศ

มหาวิทยาลัยเทคโนโลยีพระจอมเกล้าพระนครเหนือ

เอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

โปรดกรอกชื่อหนังสือพิมพ์ให้ชัดเจน และต้องแจ้งถึงผู้จัดพิมพ์เอกสารทุกครั้งที่จะกรอกไปให้

| Room 7 | | |
|--|--|-------------|
| Session 6 : Image Processing and Pattern Recognition | | |
| Time/Paper-ID | Title/Author | Page |
| 13.00-13.20 NCCIT2011-158 | Comparison of vehicle license plate localization algorithm: EP & WEP <i>Suchitra Audulkasem, Jitdumrong Preechasuk, Piticha Paewchompoo, and Saowaluk Jirakulkanok</i> | 784 |
| 13.20-13.40 NCCIT2011-194 | Automatic Detection and Counting Scalp Hairs <i>Sumaree Arayasombat and Nongluk Covavisaruch</i> | 790 |
| 13.40-14.00 NCCIT2011-187 | Measuring of a Vitiligo Lesion Area on a Curved Surface Using Digital Image Processing <i>Bharima Clangphukhiew and Nongluk Covavisaruch</i> | 795 |
| Session 7: Web Application, Web Service, Information Retrieval, Natural Language Processing, and Ontology | | |
| Time/Paper-ID | Title/Author | Page |
| 14.00-14.20 NCCIT2011-191 | Developing an Ontology Knowledge Based for Automatic Online News Analysis <i>Wichuda Chotirat, Pudsadee Boonrawd, and Sageemas Na Wichian</i> | 800 |
| 14.20-14.40 NCCIT2011-33 | A Development of Ontology Prototyping for Thai Massage Therapy Searching <i>Wiraiwan Sanchana, Orasa Tetiwat, Nattavadee Hongboonmee, and Marut Buranarach</i> | 806 |
| 14.40-15.00 | <i>Coffee Break</i> | |
| 15.00-15.20 NCCIT2011-161 | Reducing Ambiguity in Thai Text to Thai Lanna Text Translation System using Viterbi Algorithm <i>Prathan Comejina and Ponrudee Netisopakul</i> | 812 |
| 15.20-15.40 NCCIT2011-18 | A Web-based Single Sign-on (SSO) using SAML 2.0 <i>Tatchai Russameroj, Pornchai Mongkolnam, and Kriengkrai Porkaew</i> | 818 |
| 15.40-16.00 NCCIT2011-24 | The System Judging Criteria of the Web Design Competition Case Study Development of Skill Department <i>J.Jamsai, U.Rodrueng, and Somboon Supattarakulchai</i> | 824 |
| 16.00-16.20 NCCIT2011-40 | (Server Management System On CentOS Linux Server Case Study: Website Skill Competition (Web Design), Development of Skill Department) <i>Nut Pornsophon, Arisa Getgosol, and Anumas Sangsawang</i> | 830 |
| 16.20-16.40 NCCIT2011-182 | Web Service System for Social Network, Facebook <i>Prapatsorn Sripadet and Cholatif Yawut</i> | 836 |

การลดความกำกวม ในการแปลภาษาไทยเป็นภาษาไทยล้านนาด้วยวิเทอร์บีอัลกอริทึม

Reducing Ambiguity in Thai Text to Thai Lanna Text Translation System using Viterbi Algorithm

ประธาน คำจិនะ (Prathan Comejina)¹ และ พรฤดี เนติโสภาคกุล (Ponrudee Netisopakul)²

^{1,2}คณะเทคโนโลยีสารสนเทศ สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง
pcomejina@gmail.com, ponrudee@it.kmitl.ac.th

บทคัดย่อ

การแปลภาษาไทยเป็นภาษาไทยล้านนา จะพบปัญหาการจับคู่ระหว่างคำไทยหนึ่งคำ กับ คำไทยล้านนาได้หลายคำ ดังนั้นงานวิจัยนี้จึงนำเสนอ วิธีการลดความกำกวมในการเลือกคำไทยล้านนา เพื่อการแปลภาษาไทยเป็นภาษาไทยล้านนา ด้วยวิธีการของแบบจำลองฮิดเดนมาร์คอฟและวิเทอร์บีอัลกอริทึม ร่วมกับคลังคู่ประโยคไทย-ไทยล้านนา จากการทดลองเบื้องต้นโดยใช้วิเทอร์บีอัลกอริทึม พบว่ามีความถูกต้องของการแปลเพิ่มขึ้นร้อยละ 60

คำสำคัญ: แบบจำลองฮิดเดนมาร์คอฟ วิเทอร์บีอัลกอริทึม การแปลภาษาด้วยคอมพิวเตอร์

To translate Thai language to Thai Lanna language, the result shows one Thai word can be translated into a few Thai Lanna words. To select the correct translation, this paper proposes to use hidden markov model and viterbi algorithm with Thai – Thai Lanna corpus. The accuracy of reducing ambiguity in translation using viterbi algorithm increases 60 percent.

Keyword: Hidden Markov Model, Viterbi Algorithm, Machine Translation

1. บทนำ

ภาษาไทยล้านนา เป็นภาษาถิ่นของคนไทยในภาคเหนือตอนบน พัฒนามาจากอักษรขอม เมื่อประมาณ พ.ศ.1919 ซึ่งเป็นช่วงเดียวกันกับการเกิดอักษรภาษาไทย ภาษาไทยล้านนาจึงเป็นภาษาเก่าแก่ที่มีคุณค่าทางวัฒนธรรม และเป็นเอกลักษณ์เฉพาะตัวของคนไทยในภาคเหนือ ปัจจุบันผู้เชี่ยวชาญการใช้

ภาษาไทยล้านนามีน้อยลง จึงควรที่จะอนุรักษ์และสืบทอดภาษาไทยล้านนาไม่ให้สูญหาย

งานด้านการแปลภาษาด้วยคอมพิวเตอร์ โดยทั่วไปจะพบปัญหาความกำกวม ในระดับคำและระดับโครงสร้างประโยค เช่น งานวิจัยการแปลล้านนาเป็นภาษาไทย [1] เป็นงานวิจัยที่มีความคล้ายคลึงกับงานวิจัยของเรา ตรงการใช้คู่ภาษาสำหรับการแปลเดียวกัน งานวิจัย [1] พบความกำกวมในระดับโครงสร้างประโยค และใช้วิธีช่วยงานเพิ่มขยายช่วยเลือกโครงสร้างประโยคที่ถูกต้อง ส่วนงานวิจัย [2] ก็พบปัญหาการแปลในระดับโครงสร้างและแก้ปัญหานี้ โดยใช้วิเทอร์บีอัลกอริทึมช่วยเลือกโครงสร้างที่เหมาะสมสำหรับการแปล ซึ่งสอดคล้องกับแนวคิดของงานวิจัยนี้ ในการนำวิเทอร์บีอัลกอริทึมมาช่วยในการเลือกคำแปลและลำดับคำที่ถูกต้อง

ดังนั้น งานวิจัยนี้จึงขอนำเสนอ การลดความกำกวมในการแปลภาษาไทยเป็นภาษาไทยล้านนาโดยอิงพจนานุกรม ด้วยวิธีการของแบบจำลองฮิดเดนมาร์คอฟและวิเทอร์บีอัลกอริทึม ร่วมกับคลังคู่ประโยคไทย-ไทยล้านนา

ตัวอย่างความกำกวมในการแปลภาษาไทยเป็นภาษาไทยล้านนา เช่น ในภาษาไทยคำว่า “บอก” จะพบคำแปลเป็นภาษาไทยล้านนาได้ 3 คำ คือ ກ່ວາ (ก่า), ດູນ (ดูน) และ ບອກ (บอก)

ซึ่งแต่ละคำ จะถูกใช้ในสถานการณ์ที่แตกต่างกัน ขึ้นอยู่กับบริบทที่เกิดร่วมในประโยค คือ

“ວ່າຍຸດ ກ່ວາ ປູ້ ດູງ” (พอเสียดก่าปิดงาน) ก่า ใช้กับการพูดอย่างเป็นพิธีการ

“ຮີ ດູນ ດ້ວງ” (พี่ดูนเจ้า) ดูน ใช้ในการพูดกับบุคคลระดับเจ้านาย(ฝ่ายเหนือ)

“ບອກ ບອກ ບອກ” (หละอ่อนบอกแม่) บอก ใช้ในการพูดในระดับสามัญชนทั่วไป

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

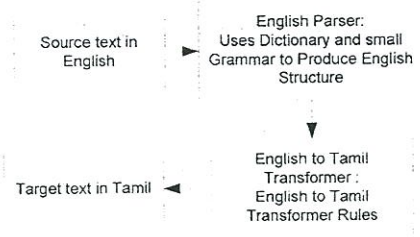
2. งานวิจัยที่เกี่ยวข้อง

งานวิจัยที่เกี่ยวกับการแปลภาษาไทยล้านนาในปัจจุบัน พบมีงานวิจัยทางการแปลจากภาษาไทยล้านนาเป็นภาษาไทย [1] โดยใช้พจนานุกรมร่วมกับช่วยงานเพิ่มขยาย ซึ่งเป็นงานที่มีความใกล้เคียงกับงานวิจัยของเรา มีขั้นตอนดังภาพที่ 1



ภาพที่ 1: การแปลภาษาไทยล้านนาเป็นภาษาไทยกลาง [1]

- ขั้นตอนที่ 1 การตัดคำโดยใช้อัลกอริทึมแบบเลือกคำยาวสุด
- ขั้นตอนที่ 2 ค้นหาคำแปลในพจนานุกรม ซึ่งมีกรเก็บข้อมูลแบบทรี
- ขั้นตอนที่ 3 หากคำไหนไม่พบในพจนานุกรม จะใช้การปริวรรตภาษาไทยล้านนา แทนโครงสร้างพยางค์ไทยล้านนาด้วยอโตมาต้าเชิงไม่กำหนด และกำหนดให้ชนิดของคำเป็นคำนามชนิดวิสามานนาม
- ขั้นตอนที่ 4 เลือกโครงสร้างประโยคไทยล้านนาและแจกแจงประโยคด้วยกฎโครงสร้างวลีและกฎโครงสร้างประโยคโดยใช้ช่วยงานเพิ่มขยาย
- ขั้นตอนที่ 5 เลือกคำแปล โดยพิจารณาจากบริบทและการเกิดรวมของคำ ในการเลือกคำแปลหลายๆ คำ



ภาพที่ 2: การแปลภาษาอังกฤษเป็นภาษาทมิฬ [2]

งานวิจัยแปลภาษาอังกฤษเป็นภาษาทมิฬ [2] ได้นำวิเทอร์บีอัลกอริทึมมาประยุกต์ใช้ในการลดความกำกวมระดับโครงสร้างประโยค

งานวิจัยนี้มีกระบวนการทำงานหลักๆ อยู่ 2 กระบวนการ ดังภาพที่ 2 คือ

- กระบวนการเรียนรู้(Training phases) จะสร้างโมเดลทางสถิติสำหรับการแปล โดยใช้คลังข้อมูลภาษาอังกฤษภาษาทมิฬ
 - กระบวนการแปล(Translation phases) จะใช้วิธีการค้นหาแบบฮิวริสติก เพื่อหาคำแปลที่ดีที่สุดของข้อความลักษณะของภาษาทมิฬ เนื่องจากเป็นภาษาที่คำในภาษาทมิฬเกิดจากการรวมกันของหลายๆ หน่วยคำมีข้อมูลเพศและกริยาาลงท้ายที่แตกต่างกัน ทำให้ต้องพิจารณาโดยบริบท ดังนั้นในกระบวนการแปลของงานวิจัยดังกล่าว จึงใช้แบบจำลองฮิดเดนมาร์คอฟและวิเทอร์บีอัลกอริทึม ในการค้นหาโครงสร้างประโยคที่เหมาะสมสำหรับการแปลที่สุด มีขั้นตอนการทำงานดังนี้
 - ขั้นตอนที่ 1 คำนวณหาค่าความน่าจะเป็นของสถานะเริ่มต้น
 - ขั้นตอนที่ 2 คำนวณความน่าจะเป็นของสถานะถัดมาไปจนถึงสถานะสุดท้าย โดยที่สถานะข้างหน้า จะเป็นตัวกำหนดสถานะถัดมาที่ว่าจะเป็นไปได้
 - ขั้นตอนที่ 3 ค้นหาเส้นทางที่มีความน่าจะเป็นสะสมสูงที่สุด
- จากงานวิจัย [2] ทำให้ผู้วิจัยเห็นแนวทางการแก้ปัญหาความกำกวม ที่อาจเกิดขึ้นในการแปลภาษาไทยเป็นภาษาไทยล้านนาด้วยการประยุกต์ใช้แบบจำลองของฮิดเดนมาร์คอฟและวิเทอร์บีอัลกอริทึม โดยอิงพจนานุกรม ในการลดความกำกวมสำหรับกรณีที่พบคำแปลภาษาไทยล้านนามากกว่า 1 คำ ซึ่งคำแปลที่เป็นผลลัพธ์จะได้อาจมาจากคำแปลที่มีความน่าจะเป็นสูงสุด ทำให้การ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

แปลภาษามีความใกล้เคียงกับธรรมชาติของภาษาล้านนามากยิ่งขึ้น

3. การลดความกำกวมในการแปลภาษาไทย-ไทยล้านนา

3.1 ภาพรวมของการแปลภาษาไทยเป็นภาษาไทยล้านนา

ในการแปลภาษาไทยเป็นภาษาไทยล้านนา โดยอิงพจนานุกรม มีขั้นตอนดังนี้



ภาพที่ 3: ภาพรวมของการแปลภาษาไทยเป็นภาษาไทยล้านนา

ขั้นตอนที่ 1 ตัดคำในประโยคภาษาไทยที่รับเข้ามา ด้วยโปรแกรม KU CUT พัฒนาโดยมหาวิทยาลัยเกษตรศาสตร์ ใช้เทคนิคการเรียนรู้จากข้อมูลแบบไม่ใช้ตัวอย่างร่วมกับการใช้พจนานุกรมและคำ [3] มีความถูกต้องของการตัดคำร้อยละ 79

ขั้นตอนที่ 2 การจับคู่คำไทยกับคำไทยล้านนา ขั้นตอนนี้จะนำคำไทยที่ตัดได้ไปค้นหาคำแปลในพจนานุกรมคำไทย-ไทยล้านนา ซึ่งคำไทยคำหนึ่งอาจจับคู่กับคำแปลไทยล้านนาได้หลายคำ จึงเป็นสาเหตุทำให้เกิดความกำกวมขึ้น

ขั้นตอนที่ 3 แปลงคำไทยล้านนาที่ไม่พบในพจนานุกรมและเขียนเรียงเป็นประโยคไทยล้านนา ขั้นตอนนี้จะนำคำไทยที่ไม่พบคำแปลในพจนานุกรมมาแปลงเป็นคำไทยล้านนาตามกฎการเรียงอักษรไทยล้านนาที่ได้มาจากผู้เชี่ยวชาญและการปริวรรตภาษาไทยล้านนา จากนั้นจะเขียนเรียงคำไทยล้านนาที่ได้จากขั้นตอนที่ 2 และ 3 เป็นประโยคไทยล้านนาที่เป็นไปได้ทั้งหมด

ขั้นตอนที่ 4 ค้นหาประโยคที่มีความน่าจะเป็นสูงสุด ด้วยวิเทอร์บีอัลกอริทึม เป็นขั้นตอนสำหรับลดความกำกวมของคำศัพท์ จะถูกใช้เมื่อประโยคไทยล้านนาที่เรียบเรียงได้ มีมากกว่า 1 ประโยค จะนำประโยคไปคำนวณหาประโยคที่มีความน่าจะเป็นสูงสุด ค่าความน่าจะเป็นที่จะนำมาคำนวณของแต่ละคำ ได้มาจากการนับความถี่ของคำในคลังคู่ประโยคไทย-ไทยล้านนา และผลลัพธ์สุดท้าย จะได้ประโยคไทยล้านนาที่มีความน่าจะเป็นสะสมสูงสุด ซึ่งจะเป็นประโยคที่มีความเหมาะสมที่จะใช้ในการแปลที่สุด

3.2 การประยุกต์ใช้วิเทอร์บีอัลกอริทึมสำหรับการแปลภาษาไทยเป็นภาษาไทยล้านนา

แบบจำลองฮิดเดนมาร์คอฟ [4] เป็นเทคนิคการเรียนรู้ทางสถิติ (Statistical Machine Learning) ที่ใช้จำแนกรูปแบบของลำดับเหตุการณ์หนึ่งๆ ซึ่งลำดับของเหตุการณ์จะถูกจำลองให้อยู่ในรูปแบบของลำดับการเปลี่ยนแปลงของสถานะ (State) จะมีค่าความน่าจะเป็นอยู่ 2 ชนิด คือ ค่าความน่าจะเป็นของการเปลี่ยนสถานะ และค่าความน่าจะเป็นของการทำให้เกิดผลลัพธ์ เมื่อมีการเปลี่ยนสถานะเกิดขึ้น และผลลัพธ์นี้จะถูกนำไปใช้ในการตัดสินใจว่าข้อมูลที่เข้ามามีโอกาสเป็นอะไรมากที่สุด

วิเทอร์บีอัลกอริทึม เป็นอัลกอริทึมหนึ่งในแบบจำลองฮิดเดนมาร์คอฟ ที่ใช้ในการจัดตำแหน่งลำดับของสถานะต่างๆ โดยลำดับของแต่ละสถานะ จะขึ้นอยู่กับความน่าจะเป็นสูงสุด ที่จะเป็นไปได้ทั้งหมด

การทำงานของวิเทอร์บีอัลกอริทึม [5] มี 3 ขั้นตอน คือ
ขั้นตอนที่ 1 คำนวณค่าความน่าจะเป็นของสถานะเริ่มต้นของข้อมูลนำเข้า

ขั้นตอนที่ 2 คำนวณความน่าจะเป็นของสถานะถัดไป ทุกสถานะ เพื่อจะเก็บเส้นทางที่ดีที่สุด สำหรับการเลือกเส้นทางในขั้นตอนถัดไป

ขั้นตอนที่ 3 การย้อนเส้นทาง เพื่อค้นหาลำดับความน่าจะเป็นสูงสุด จากสถานะที่ได้คำนวณไว้แล้ว ซึ่งจะได้ลำดับเส้นทางสถานะที่เป็นไปได้สูงสุดด้วย

การลดความกำกวมในการแปลภาษาไทยเป็นภาษาไทยล้านนาด้วยวิเทอร์บีอัลกอริทึม มีสมการและตัวแปรที่ต้องใช้ในอัลกอริทึม ดังต่อไปนี้

$$- \text{PROB}(w_i|L_i) \cong \frac{\text{จำนวนการเกิดที่ } w_i \text{ มี } L_i \text{ เป็นคู่}}{\text{จำนวน } L_i \text{ ทั้งหมด}} \quad (1)$$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สมการที่ (1) ใช้ในการคำนวณหาความน่าจะเป็นที่ L_i จะเป็นคำคู่ของ w_i โดยที่ w_i คือ คำไทยที่ต้องการแปล และ L_i คือ คำแปลไทยล้านนาที่พบในพจนานุกรม

$$- \text{PROB}(L_i|L_{i-1}) \cong \frac{\text{จำนวนการเกิดของ } L_i \text{ ตามหลัง } L_{i-1}}{\text{จำนวน } L_{i-1} \text{ ทั้งหมด}} \quad (2)$$

สมการที่ (2) ใช้ในการคำนวณหาความน่าจะเป็น เมื่อคำภาษาไทยล้านนา L_{i-1} จะตามมาด้วยคำไทยล้านนา L_i

- T คือ จำนวนคำในประโยคภาษาไทยที่ต้องการแปล
- N คือ จำนวนคำแปลภาษาไทยล้านนาที่พบในพจนานุกรมสำหรับคำไทยที่ต้องการแปล
- SEQSCORE คือ อาร์เรย์ขนาด $N \times T$ ใช้เก็บผลลัพธ์การคำนวณความน่าจะเป็นสูงสุด ของลำดับคำที่อาจจะแปลได้ทั้งหมด ตั้งแต่ลำดับแรกจนถึงคำในลำดับปัจจุบัน
- BACKPTR คือ อาร์เรย์ขนาด $N \times T$ ใช้เก็บเส้นทางที่มีค่าความน่าจะเป็นสะสมสูงสุด
- C คือ ลิสต์ที่มีขนาด T ใช้เก็บชุดลำดับของคำไทยล้านนาที่มีค่าความน่าจะเป็นสูงสุด
- Con คือ ลิสต์ที่ใช้เก็บข้อมูลคำไทยที่จับคู่กับคำไทยล้านนาได้เพียง 1 คำ ลิสต์นี้มีประโยชน์จะช่วยลดจำนวนรอบในการคำนวณหาความน่าจะเป็น โดยไม่ต้องคำนวณทุกๆ คำไทยล้านนา

```

For i=1 to N do
    if w1 in Con do
        SEQSCORE(i,1)= PROB(w1|L1)*PROB(L1|c)
        BACKPTR(i,1)=0
        Break;
    else
        SEQSCORE(i,1)= PROB(w1|L1)*PROB(L1|c)
        BACKPTR(i,1)=0
For t=2 to T
    For i=1 to N
        if i in Con do
            SEQSCORE(i,t)= Max_{j,1..N}(SEQSCORE(j,t-1)
                *PROB(Lj|L)) * PROB(wt|L)
            BACKPTR(i,t)=index of j that gave the max above
            Break;
        else
            SEQSCORE(i,t)= Max_{j,1..N}(SEQSCORE(j,t-1)
                *PROB(Lj|L)) * PROB(wt|L)
            BACKPTR(i,t)=index of j that gave the max above
C(T)=i that maximizes SEQSCORE(i,T)
For i=T-1 to 1 do: C(i)=BACKPTR(C(i+1),i+1)
    
```

ภาพที่ 4: วิเทอร์บีอัลกอริทึมในการลดความกำกวมในการแปลภาษาไทยเป็นภาษาไทยล้านนา

จากภาพที่ 4 วิเทอร์บีอัลกอริทึมที่ใช้ลดความกำกวมในการแปลภาษาไทยเป็นภาษาไทยล้านนา มีขั้นตอนดังนี้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ขั้นตอนที่ 1 คำนวณค่าความน่าจะเป็นของคำไทยล้านนาที่จะเป็นคำแปลของคำไทยคำแรก

ขั้นตอนที่ 2 คำนวณค่าความน่าจะเป็นของคำไทยล้านนาที่จะเป็นคำแปลของคำไทยคำถัดๆ ไปจนถึงคำสุดท้าย โดยจะเก็บตำแหน่งของคำไทยล้านนาที่มีความน่าจะเป็นสะสมสูงสุด ณ ตำแหน่งที่กำลังพิจารณาไว้ในลิสต์ BACKPTR เพื่อจดจำเส้นทาง

ขั้นตอนที่ 3 ย้อนเส้นทางในตัวแปร BACKPTR เพื่อหาลำดับคำไทยล้านนา ที่มีค่าความน่าจะเป็นสะสมสูงสุด จากการย้อนเส้นทางนี้จะทำให้ประโยคภาษาไทยล้านนา ที่มีค่าความน่าจะเป็นสะสมสูงสุดอยู่ในลิสต์ C

3.3 ตัวอย่างการลดความกำกวมในการแปลภาษาไทยเป็นภาษาไทยล้านนา

การแปลภาษาไทยเป็นภาษาไทยล้านนา ในที่นี้จะยกตัวอย่างโดยใช้คลังคู่ประโยคไทย-ไทยล้านนาจากนิทานพื้นบ้านเรื่อง “ปู่กับนกยาง” ซึ่งมีจำนวนคู่ประโยค 29 ประโยค โดยมีประโยคตั้งต้นคือ “นก ยาง บอก” เมื่อผ่านขั้นตอนการตัดคำ และจับคู่คำศัพท์ สามารถเรียงเป็นประโยคไทยล้านนาได้ 6 ประโยคดังนี้

- ประโยคที่ 1 “ลื้อ ฆว ฆัว” (นก ฆาง ก่าว)
- ประโยคที่ 2 “ลื้อ ฆว ฐว” (นก ฆาง ฐน)
- ประโยคที่ 3 “ลื้อ ฆว ฆวต” (นก ฆาง บอก)
- ประโยคที่ 4 “ฆวต ฆว ฆัว” (สะกุนะ ฆาง ก่าว)
- ประโยคที่ 5 “ฆวต ฆว ฐว” (สะกุนะ ฆาง ฐน)
- ประโยคที่ 6 “ฆวต ฆว ฆวต” (สะกุนะ ฆาง บอก)

ตารางที่ 1: แสดงค่าความถี่ของคำไทยคู่ล้านนาจาก คลังคู่ประโยคไทย-ไทยล้านนา

| คำไทยคู่ล้านนา | ความถี่ | คำไทยคู่ล้านนา | ความถี่ |
|-----------------|---------|----------------|---------|
| นก/ลื้อ (นก) | 11 | บอก/ฆัว (ก่าว) | 0 |
| นก/ฆวต (สะกุนะ) | 0 | บอก/ฐว (ฐน) | 0 |
| ยาง/ฆวต (ฆาง) | 11 | บอก/ฆวต (บอก) | 2 |

ตารางที่ 2: แสดงค่าความถี่แบบยูนิแกรมของคำล้านนาจากคลังคู่ประโยคไทย-ไทยล้านนา

| คำล้านนา | ความถี่ | คู่คำล้านนา | ความถี่ |
|--------------|---------|-------------|---------|
| ลื้อ (นก) | 11 | ฆัว (ก่าว) | 0 |
| ฆวต (สะกุนะ) | 0 | ฐว (ฐน) | 0 |
| ฆว (ฆาง) | 11 | ฆวต (บอก) | 2 |

ตารางที่ 3: แสดงค่าความถี่แบบไบนารีของคำล้านนาจากคลังคำ

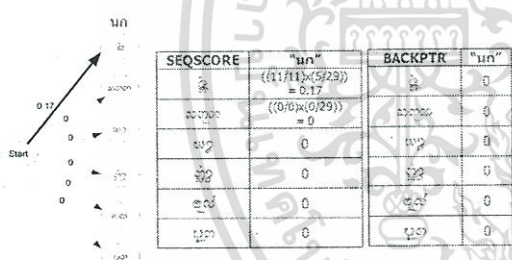
ประโยคไทย-ไทยล้านนา

| คู่คำล้านนาที่ตามกันมา | ความถี่ | คู่คำล้านนาที่ตามกันมา | ความถี่ |
|--------------------------|---------|-------------------------------|---------|
| ขจฺ ฐึ (หญิงนก) | 11 | ขจฺ ขวค (หญิงขาง) | 0 |
| ขจฺ ขจฺขจ (หญิงสะกุนะ) | 0 | ขจฺ ขวค (บอกขาง) | 0 |
| ขจฺ ขจฺ (เก่าขาง) | 0 | ฐึ Start (นกเริ่มต้นประโยค) | 5 |

ข้อสังเกต ค่าความถี่ของคำไทยคู่คำล้านนา และ ค่าความถี่แบบยูนิแกรมของคำล้านนา ของคำอื่นๆ ที่ไม่ได้ปรากฏในตารางที่ 1 ถึง ตารางที่ 3 ให้กำหนดมีค่าเท่ากับ 0 ส่วนค่าความถี่แบบไบนารีของคำล้านนาของคำอื่นๆ ที่ไม่ได้ปรากฏในตารางที่ 1 ถึง ตารางที่ 3 ให้กำหนดมีค่าเท่ากับ 0.0001 [5]

จากนั้นนำประโยคทั้ง 6 มาผ่านกระบวนการลดความกำกวมด้วยวิเทอรบีอัลกอริทึม ดังนี้

ขั้นตอนที่ 1 คำนวณค่าความน่าจะเป็นของคำไทยล้านนาที่จะเป็นค่าแปลของคำไทยคำแรก ในที่นี้คือ คำว่า "นก" ตามการทำงานในรอบที่ (1) ของภาพที่ 4

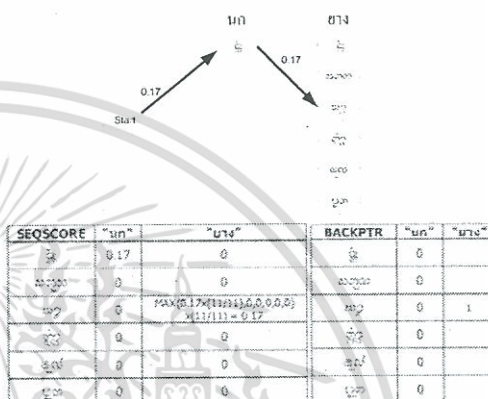


ภาพที่ 5: แสดงการคำนวณหาค่าความน่าจะเป็นของคำแปลของคำว่า "นก" ในตัวแปร SEQSCORE และ การเริ่มบันทึกเส้นทางในตัวแปร BACKPTR

จากภาพที่ 5 ในตัวแปร SEQSCORE จะได้ค่าความน่าจะเป็นของคำว่า "ขจฺ" (นก) ที่จะเป็ค่าแปลของคำไทย "นก" สูงที่สุดคือ 0.17 และ เริ่มการบันทึกเส้นทางลงในตัวแปร BACKPTR ในตำแหน่งคอลัมน์ของคำว่า "นก" กำหนดให้เป็น 0 เนื่องจากเป็นค่าแรกของประโยค

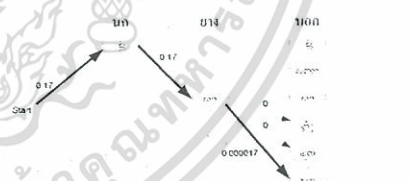
ขั้นตอนที่ 2 คำนวณค่าความน่าจะเป็นของคำไทยล้านนาที่จะเป็นค่าแปลของคำไทยคำถัดมา คือ คำว่า "ขาง" ตามการทำงานในรอบที่ (2) ของภาพที่ 4 เนื่องจากในขั้นตอนของการจับคู่คำไทยกับคำไทยล้านนา คำว่า "ขาง" สามารถจับคู่กับคำล้านนาได้เพียง 1 คำ จึงมีการบันทึกตำแหน่งของคำไว้ในตัวแปร Con ซึ่งจะช่วยให้ไม่ต้องไปคำนวณค่าความน่าจะเป็นค่าแปล

ของคำไทยล้านนาที่เหลืออื่นๆ ที่มีในลิสต์ L ดังภาพที่ 6 แสดงการคำนวณหาค่าความน่าจะเป็นของคำว่า "ขจฺ" (หญิง) ที่จะเป็ค่าแปลของคำว่า "ขาง" ในตัวแปร SEQSCORE โดยที่ความน่าจะเป็นที่ได้จะเป็นค่าความน่าจะเป็นที่สะสม มาตั้งแต่คำเริ่มต้นมาจนถึงคำที่กำลังพิจารณา ซึ่งก็คือ คำว่า "ขาง" มีความน่าจะเป็นที่คำว่า "ขจฺ" (หญิง) จะเป็ค่าแปลเท่ากับ 0.17 และทำการบันทึกเส้นทางเท่ากับ 1 ลงในตัวแปร BACKPTR



ภาพที่ 6: แสดงการคำนวณหาค่าความน่าจะเป็นของคำแปลของคำว่า "ขาง" ในตัวแปร SEQSCORE และ การบันทึกเส้นทางในตัวแปร BACKPTR

ต่อมาคำนวณหาค่าความน่าจะเป็นของคำไทยล้านนา ที่จะเป็ค่าแปลของคำว่า "บอก"



ภาพที่ 7: แสดงการคำนวณหาค่าความน่าจะเป็นของคำแปลของคำว่า "บอก" ในตัวแปร SEQSCORE และ การบันทึกเส้นทางในตัวแปร BACKPTR

จากภาพที่ 7 ในตัวแปร SEQSCORE จะได้ค่าความน่าจะเป็นของคำว่า "ขจฺ" (บอก) ที่จะเป็ค่าแปลของคำไทย "บอก" สูงที่สุดคือ 0.000017 และ และทำการบันทึกเส้นทางเท่ากับ 2 ลงในตัวแปร BACKPTR

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ขั้นตอนที่ 3 ย้อนเส้นทางเพื่อหาลำดับคำแปลไทยล้านนา ที่มีค่าความน่าจะเป็นสูงสุด จากลิสต์ BACKPTR ตามการทำงานในกรอบที่ (3) ของภาพที่ 4 จะทำให้ได้ประโยคภาษาไทยล้านนาที่มีค่าความน่าจะเป็นสะสมสูงสุดเป็นผลลัพธ์สุดท้ายอยู่ในตัวแปร C คือ ประโยค “ลื้อ ๗๖ ๒๓” (นก ฃาง บอก)

4. การดำเนินงานในปัจจุบันและอนาคต

4.1 การดำเนินงานในปัจจุบัน

สร้างคลังคู่ประโยคไทย-ไทยล้านนาต้นแบบ จากนิทานเรื่อง “ปู่กับนกกระยาง” ทำให้ได้คู่ประโยคไทย-ไทยล้านนา 29 คู่ ประกอบด้วยคำจำนวน 201 คำ โดยเป็นคำที่ไม่ซ้ำกันจำนวน 93 คำ และทดสอบการแปลภาษาไทยเป็นภาษาไทยล้านนา ด้วยประโยคทดสอบจำนวน 10 ประโยคที่ผู้วิจัยสร้างขึ้นเอง

วัดประสิทธิภาพการแปล โดยการคำนวณจากประโยคที่แปลถูกต้องต่อประโยคที่นำมาทดสอบทั้งหมด ดังสมการด้านล่าง

$$\text{ประสิทธิภาพ(\%)} = \frac{X}{N} \times 100 \quad (3)$$

ตัวแปร X คือ จำนวนของประโยคที่แปลถูกต้อง
ตัวแปร N คือ จำนวนของประโยคที่นำมาทดสอบทั้งหมด

4.1.1 ผลการทดลอง

แบ่งการทดลองเป็น 2 วิธี คือ วิธีที่ 1 ใช้วีเทอร์บี อัลกอริทึมช่วยลดความกำกวม ในกรณีที่คำไทยคำหนึ่งมีคำแปลไทยล้านนาได้หลายคำ เปรียบเทียบกับ วิธีที่ 2 ที่จะแก้ปัญหานี้ โดยการเลือกคำแปลไทยล้านนาคำแรก ที่พบในพจนานุกรม ได้ผลการทดลองดังตารางที่ 4

ตารางที่ 4: ตารางเปรียบเทียบประสิทธิภาพของการแปล

| วิธีการทดลอง | ผลลัพธ์ที่แปลถูกต้อง |
|--------------|----------------------|
| วิธีที่ 1 | 90% |
| วิธีที่ 2 | 30% |

จากการทดลอง วิธีที่ 1 ให้ผลการแปลที่มีความถูกต้องมากกว่าวิธีที่ 2 ถึง 60% ดังนั้นจึงสรุปได้ว่าวีเทอร์บีอัลกอริทึมสามารถช่วยลดความกำกวม ในกรณีที่คำไทยคำหนึ่งมีคำแปล

ไทยล้านนาหลายคำได้ และการแปลจะยังมีประสิทธิภาพมากขึ้นเมื่อจำนวนคู่ประโยคในคลังคู่ประโยคไทย-ไทยล้านนามีมากขึ้น

4.2 การดำเนินงานในอนาคต

จะจัดสร้างคลังคู่ประโยคไทย-ไทยล้านนา 500 ประโยค โดยนำมาจากนิทานและวรรณกรรมพื้นบ้าน และใช้ประโยคทดสอบ 100 ประโยค โดยจะวัดประสิทธิภาพเปรียบเทียบกับความเห็นของผู้เชี่ยวชาญ 4 คน

5. สรุป

จากการทดลองเบื้องต้นวีเทอร์บีอัลกอริทึมสามารถช่วยลดความกำกวมในการแปลภาษาไทยเป็นภาษาไทยล้านนาได้ โดยเมื่อทดสอบกับประโยคที่สร้างจากคำไทยที่มีปรากฏในคลังคู่ประโยคไทย-ไทยล้านนา พบว่ามีประสิทธิภาพเพิ่มขึ้นร้อยละ 60 ในอนาคตผู้วิจัยจะรวบรวมประโยคในคลังคู่ประโยคไทย-ไทยล้านนาให้ได้มากขึ้น เพื่อรองรับการแปลภาษาไทยเป็นภาษาไทยล้านนาให้ดียิ่งขึ้น

เอกสารอ้างอิง

[1] พุทธิดี ศิริแสงตระกูล และ นิกร ยาทรม, “การแปลภาษาล้านนาเป็นภาษาไทยกลาง”, *วารสารวิจัย มหาวิทยาลัยขอนแก่น*, มีนาคม 2552 หน้า 264-274.

[2] S.Vetrivel and Diana Baby, “English to Tamil Statistical Machine Translation and Alignment Using HMM”, *Recent Advances in Networking, VLSI and Signal Processing* University of Cambridge, UK, 20-22 February 2010, pp.182-186.

[3] สุทธิ สดประเสริฐ และ อัสนีย์ ก่อตระกูล, “การตัดคำภาษาไทยด้วยเทคนิคการเรียนรู้แบบไม่ใช้ตัวอย่าง”, *The 7th National Computer Science and Engineering Conference (NCSEC'2003)*, Chonburi Thailand, 2003.

[4] ยศวัฒน์ แสนสิงห์ และ อัครา ประโยชน์, “การรู้จำชื่อเฉพาะภาษาไทย โดยใช้แบบจำลองฮิดเดนมาร์คอฟ”, *Proceedings of the Second Conference on Knowledge and Smart Technologies 2010 (KST-2010)*, Chonburi Thailand, 24-25 July 2010.

[5] James Allen, *Natural Language Understanding*, The Benjamin/Cummings Publishing Company, Inc.:Redwood City, CA, 1995.

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



THE 8TH NATIONAL CONFERENCE
ON COMPUTING AND INFORMATION
TECHNOLOGY

11-12 MAY 2012

FACULTY OF INFORMATION TECHNOLOGY
KING MONGKUT'S UNIVERSITY OF TECHNOLOGY NORTH BANGKOK (KMUTNB)

WWW.NCCIT.NET

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ประวัติผู้เขียน

ชื่อ-นามสกุล นายประธาน คำจិនะ

วัน เดือน ปีเกิด วันที่ 20 เมษายน 2526

ประวัติการศึกษา

- ระดับประถมศึกษา โรงเรียนวัดท่าโป่ง
- ระดับมัธยมศึกษาตอนต้น โรงเรียนสันป่าตองวิทยาคม
- ระดับมัธยมศึกษาตอนปลาย โรงเรียนสันป่าตองวิทยาคม
- ระดับอุดมศึกษา วิทยาศาสตร์บัณฑิต สาขาวิชาการคอมพิวเตอร์ มหาวิทยาลัยราชภัฏ

เชียงใหม่

ประวัติการทำงาน

- พ.ศ. 2548 - 2550 พนักงานมหาวิทยาลัย ตำแหน่ง นักวิชาการคอมพิวเตอร์ สำนักส่งเสริมวิชาการ มหาวิทยาลัยราชภัฏเชียงใหม่
- พ.ศ. 2550 - ปัจจุบัน พนักงานมหาวิทยาลัยสายวิชาการ สาขาวิชาคอมพิวเตอร์ คณะวิทยาศาสตร์และเทคโนโลยี มหาวิทยาลัยราชภัฏเชียงใหม่