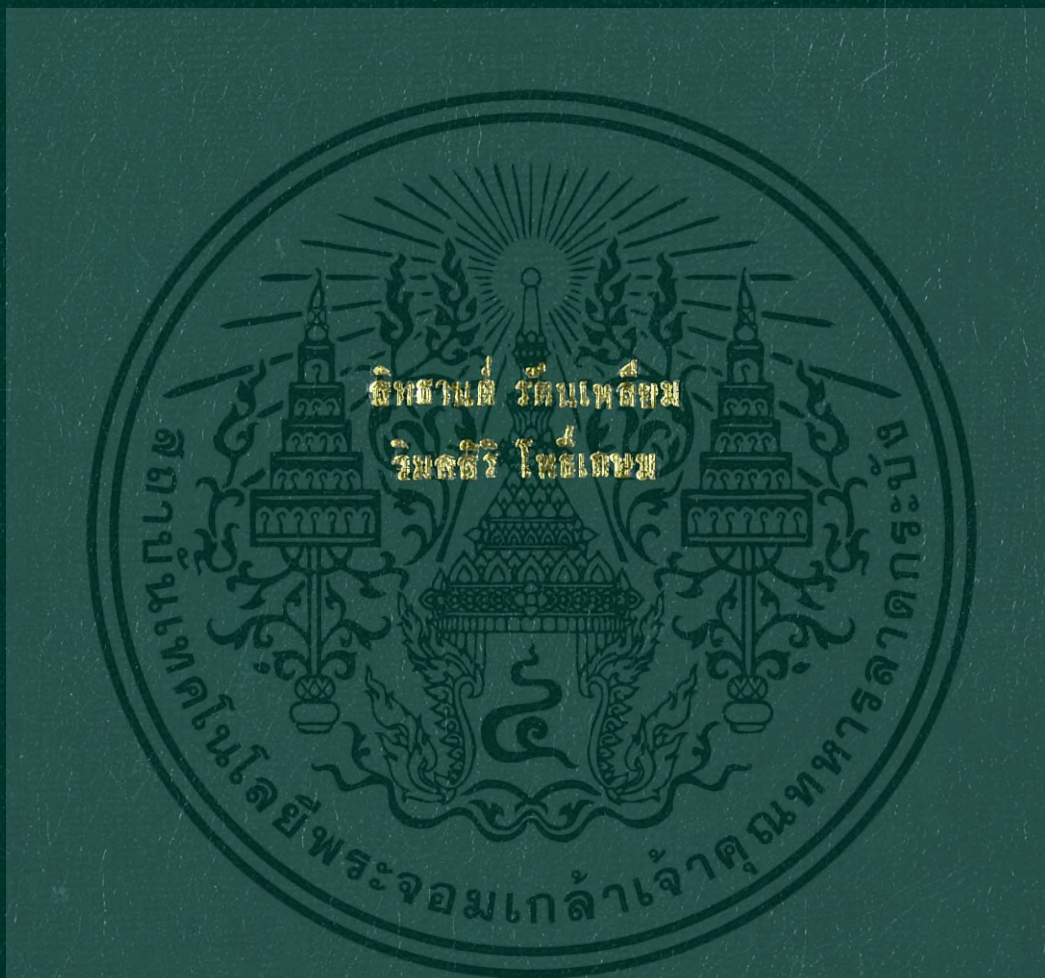


ระบบบันทึกข้อมูลภาษาไทยด้วยตัวอักษรโรมัน
ROMANIZED THAI INPUT METHOD EDITOR



ปริญญานิพนธ์นี้เป็นส่วนหนึ่งของกรณีศึกษาทางเทคโนโลยีสารสนเทศ วิทยาลัยเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

สาขาวิชาเทคโนโลยีสารสนเทศ

คอมพิวเตอร์และเทคโนโลยีสารสนเทศ

สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

ภาคเรียนที่ ๒ ปีการศึกษา ๒๕๕๘

สำนักหอสมุดกลาง พระจอมเกล้าลาดกระบัง

ระบบป้อนข้อมูลภาษาไทยด้วยตัวอักษรโรมัน
ROMANIZED THAI INPUT METHOD EDITOR

โดย



สิทธานต์ รัตนเหลียม

SITTAN RATTANALIAM

วิมลสิริ โปธิเกษม

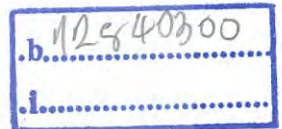
WIMONSIRI POKASAME

อาจารย์ที่ปรึกษา

ดร.สุภวรรณ อันนันทน์

อาจารย์ที่ปรึกษาร่วม

ดร.นล เปรมัชเชียร



มหาวิทยาลัย.....
เลขทะเบียน..... 146195
วันเดือนปี..... 25 ใสอ. 2560

ปริญญานิพนธ์นี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญาวิทยาศาสตรบัณฑิต

สาขาวิชาเทคโนโลยีสารสนเทศ

คณะเทคโนโลยีสารสนเทศ

สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับภาคเรียนที่ 2 ปีการศึกษา 2558
กรุณาให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ระบบป้อนข้อมูลภาษาไทยด้วยตัวอักษรโรมัน
ROMANIZED THAI INPUT METHOD EDITOR

โดย



ปริญญานิพนธ์นี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญาวิทยาศาสตรบัณฑิต
สาขาวิชาเทคโนโลยีสารสนเทศ
คณะเทคโนโลยีสารสนเทศ
สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับภาคเรียนที่ 2 ปีการศึกษา 2558 อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ROMANIZED THAI INPUT METHOD EDITOR



A PROJECT SUBMITTED IN PARTIAL FULFILLMENT
OF THE REQUIREMENT FOR THE DEGREE OF
BACHELOR OF SCIENCE PROGRAM IN INFORMATION TECHNOLOGY
FACULTY OF INFORMATION TECHNOLOGY
KING MONGKUT'S INSTITUTE OF TECHNOLOGY LADKRABANG

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาปี 2015 เท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



COPYRIGHT 2016

FACULTY OF INFORMATION TECHNOLOGY

KING MONGKUT'S INSTITUTE OF TECHNOLOGY LADKRABANG

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น เมื่ออนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ใบรับรองปริญญาโท ประจำปีการศึกษา 2558
คณะเทคโนโลยีสารสนเทศ
สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

เรื่อง

ระบบป้อนข้อมูลภาษาไทยด้วยตัวอักษรโรมัน

ROMANIZED THAI INPUT METHOD EDITOR

ผู้จัดทำ

1. นายสิทธิชานต์ รัตนเหลี่ยม รหัสนักศึกษา 55070130
2. นางสาววิมลสิริ โพธิ์เกษม รหัสนักศึกษา 55070109

สุวรรณ อัม

..... อาจารย์ที่ปรึกษา

(ดร.สุภวรรณ อัมพันธ์)

นล เปรมชัย

..... อาจารย์ที่ปรึกษาร่วม

(ดร.นล เปรมชัย)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

หัวข้อโครงการ	ระบบป้อนข้อมูลภาษาไทยด้วยตัวอักษรโรมัน	
นักศึกษา	นายสิทธิานต์ รัตนเหลี่ยม	รหัสนักศึกษา 55070130
	นางสาววิมลสิริ โพธิ์เกษม	รหัสนักศึกษา 55070109
ปริญญา	วิทยาศาสตร์บัณฑิต	
สาขาวิชา	เทคโนโลยีสารสนเทศ	
ปีการศึกษา	2558	
อาจารย์ที่ปรึกษา	ดร.สุภวรรณ อันนันทน์	
อาจารย์ที่ปรึกษาร่วม	ดร.นล เปรมชัยเชิฐ	

บทคัดย่อ

ระบบป้อนข้อมูลภาษาไทยด้วยตัวอักษรโรมันเป็นวิธีการหนึ่งที่จะทำให้สามารถพิมพ์ภาษาไทยด้วยการใช้ตัวอักษรโรมันได้ เนื่องจากปัจจุบันชาวต่างชาติได้เข้ามามีบทบาทในประเทศไทยมากขึ้น มีการติดต่อสื่อสารกับชาวต่างชาติเพิ่มขึ้น ราชบัณฑิตยสถานจึงได้ประกาศกฎเกณฑ์การถอดอักษรไทยด้วยตัวอักษรโรมันขึ้น เพื่อใช้เป็นมาตรฐานในการเขียนข้อความหรือคำภาษาไทยด้วยการใช้ตัวอักษรโรมันเพื่ออำนวยความสะดวกให้กับชาวต่างชาติมากขึ้น ดังที่เห็นในแผนที่หรือป้ายบอกทางต่างๆ ในวิทยานิพนธ์นี้ผู้จัดทำจึงได้พัฒนาระบบป้อนข้อมูลภาษาไทยด้วยตัวอักษรโรมันขึ้น โดยยึดหลักเกณฑ์การถอดอักษรไทยด้วยตัวอักษรโรมันแบบถ่ายเสียงของราชบัณฑิตยสถานมาใช้ประกอบการพัฒนา วิทยานิพนธ์ฉบับนี้ได้นำเสนอขั้นตอนวิธีการถอดอักษรโดยใช้วิธีการเข้ารหัสผ่านตัวกลาง โดยใช้ชวาร์ตส์ค็อกซ์ ซึ่งเป็นวิธีที่จะอนุญาตให้ใช้คำภาษาไทยโดยใช้ตัวอักษรโรมันได้เพื่อค้นคืนคำไทยที่หลักการเขียนที่ตรงกันในอีกภาษาหนึ่ง โดยจะมีตัวกลางที่ตรงกันสำหรับสองภาษา ขั้นตอนวิธีที่นำเสนอแบ่งออกเป็นสองส่วนคือ ขั้นตอนวิธีการเข้ารหัสตัวอักษรภาษาไทยและขั้นตอนวิธีการเข้ารหัสตัวอักษรภาษาอังกฤษ ผลการทดลองจะแสดงให้เห็นว่าขั้นตอนวิธีการเข้ารหัสภาษาไทยด้วยตัวอักษรโรมันโดยใช้วิธีชวาร์ตส์ค็อกซ์สามารถแสดงผลออกมาได้สูงถึง 90 เปอร์เซ็นต์ นอกจากนั้นจะมีปัญหาในเรื่องความกำกวมในการออกเสียงของคำนั้นๆ และการสะกดคำที่แตกต่างกันในแต่ละบุคคลก็ส่งผลให้ระบบไม่สามารถค้นคืนคำไทยออกมาได้

Project Title	Thai Romanized Input Method Editor	
Student	Mr. Sittan Rattanaliam	Student ID 55070130
	Miss Wimonsiri Pokasame	Student ID 55070109
Degree	Bachelor of Science	
Program	Information Technology	
Academic Year	2015	
Advisor	Dr. Supawan Annannub	
Co-Advisor	Dr. Nol Premasathian	

ABSTRACT

Thai Romanized Input Method Editor is the way to write Thai language using Roman alphabets. Nowadays, the foreigners has come to Thailand a lots and they have to communicate with Thai people. The Royal Institute has established the standard by proposing the principle of Romanization[5] on the basis of transcription. In this study, we aim to develop such a system and we base on the rules of The Royal Institute. In this paper, we present an algorithm of transcription called 'Soundex' that allow the users input Roman alphabet and transform to Numeric codes. The algorithm has divided to two parts for Thai alphabets and Roman alphabet. The results has shown that the algorithm can perform to 90%, but some missing is a result of ambiguous of language and each people are spelling difference.

กิตติกรรมประกาศ

ผู้วิจัยขอขอบพระคุณ ดร.สุภวรรณ อันันันหนับ อาจารย์ที่ปรึกษาและ ดร.นล เปรมชัยเชียร อาจารย์ที่ปรึกษาร่วมเป็น อย่างสูง ปริญญาณิพนธ์เล่มนี้มีอาจสำเร็จลุล่วงได้ หากไม่ได้รับความช่วยเหลือและคำปรึกษาต่างๆจากอาจารย์ทั้งสองท่าน

นอกจากนี้ผู้วิจัยขอขอบคุณประยูทธ สุวรรณวิสารท ที่ให้ความร่วมมือในการให้ข้อมูลการทำชาวน้เด็กซ์ภาษาไทย และได้ให้คำแนะนำต่างๆ เพื่อใช้เป็นแนวทางในการดำเนินโครงการ ให้เป็นไปในแนวทางที่ถูกต้องและสำเร็จลุล่วงไปได้ด้วยดี

สิทธิชานต์ รัตน์เหลี่ยม
วิมลศิริ โพธิ์เกษม



สารบัญ

	หน้า
บทคัดย่อภาษาไทย.....	I
บทคัดย่อภาษาอังกฤษ.....	II
กิตติกรรมประกาศ.....	III
สารบัญ.....	IV
สารบัญตาราง.....	VI
สารบัญรูป.....	VII

บทที่ 1. บทนำ	1
1.1 ความเป็นมาและความสำคัญของปัญหา	1
1.2 ความมุ่งหมายและวัตถุประสงค์ของการศึกษา	1
1.3 สมมติฐานของการศึกษา.....	2
1.4 ทฤษฎีหรือแนวคิดที่ใช้ใน โครงการงาน.....	2
1.5 ขอบเขตโครงการ	2
1.6 วิธีการดำเนิน โครงการงาน	2
1.7 โครงสร้างของปริญญาานิพนธ์	3
บทที่ 2. เอกสารและงานวิจัยที่เกี่ยวข้อง	4
2.1 การถอดอักษร	4
2.2 ประวัติความเป็นมาของการถอดอักษรในภาษาไทย.....	5
2.3 การถอดอักษรไทยเป็น อักษร โรมันแบบถ่ายเสียงของราชบัณฑิตยสถาน	6
2.4 วิธีเชิงสัทลักษณ์ (Phonetic algorithm).....	8
2.5 ขั้นตอน วิธีชวณเด็กซ์ภาษาอังกฤษ	9
2.6 ขั้นตอน วิธีชวณเด็กซ์ภาษาไทย	10
2.6.1 งานวิจัยการใช้หลักคำพ้องเสียงเพื่อค้นหาชุดอักขระภาษาไทยที่ออกเสียง เหมือนกัน	10
2.6.2 งานวิจัย A Thai Soundex System for Spelling Correction	13
2.6.3 งานวิจัยการเข้ารหัสคำทับศัพท์เพื่อการค้น ค้นข้ามภาษาไทย-อังกฤษ	15

สารบัญ (ต่อ)

	หน้า
บทที่ 3. การวิเคราะห์และออกแบบระบบ	18
3.1 เก็บรวบรวมข้อมูล	19
3.2 การเขียนโปรแกรมเข้ารหัสชาวน์เด็กซ์สำหรับข้อมูลภาษาไทย	21
3.3 การเขียน โปรแกรมเข้ารหัสชาวน์เด็กซ์สำหรับภาษาอังกฤษ	26
3.4 การกำหนด Soundex สำหรับสระในภาษาไทย	29
บทที่ 4. ผลการดำเนินงาน	31
4.1 การเปรียบเทียบความสอดคล้องของรหัสชาวน์เด็กซ์จากอักษรภาษาอังกฤษและจาก อักษรภาษาไทย	31
4.2 การประเมินผลความถูกต้องแม่นยำของ โปรแกรมการถอดอักษร	32
4.3 การแก้ไขอัลกอริทึมเพื่อเพิ่มความแม่นยำในการทำงาน	34
4.4 กำกับความหมายและประเภทของคำศัพท์ในคลังข้อมูล	34
4.5 สรุปผลจากการศึกษาแนวทางที่เหมาะสมในการพัฒนาระบบการถอดอักษรอังกฤษ เป็นไทยแบบถ่ายเสียง.....	35
บทที่ 5. สรุปผล โครงการงาน	37
5.1 สรุปผลการทำโครงการงาน	37
5.2 ปัญหาที่พบระหว่างดำเนิน โครงการงาน	37
5.3 ประโยชน์ที่ได้รับจากการดำเนิน โครงการงาน.....	38
5.4 แนวทางการพัฒนาในอนาคต	38
บรรณานุกรม	39
ภาคผนวก	41

สารบัญตาราง

หน้า

ตารางที่

2.1 ตารางเทียบเสียงพยัญชนะและสระ ตามหลักเกณฑ์โดยราชบัณฑิตยสถาน	7
2.2 ตารางการเข้ารหัส ซาวน์เด็กซ์ของ Odell และ Russell	10
2.3 การกำหนดรหัสตัวอักษรของรหัสซาวน์เด็กซ์ภาษาไทยสำหรับอักษรตัวแรก	11
2.4 การกำหนดรหัสตัวอักษรของรหัสซาวน์เด็กซ์ภาษาไทยสำหรับตัวอักษรที่เหลือ	12
2.5 ตัวอย่างการเข้ารหัสซาวน์เด็กซ์ภาษาไทย	12
2.6 การเข้ารหัสสำหรับพยัญชนะต้น	13
2.7 การเข้ารหัสสำหรับพยัญชนะสระ	14
2.8 การเข้ารหัสสำหรับพยัญชนะสะกด	14
2.9 ตัวอย่างการเข้ารหัสซาวน์เด็กซ์ภาษาไทย	15
2.10 การกำหนดรหัสซาวน์เด็กซ์สำหรับอักษรไทยและอักษรอังกฤษ โดยประยุกต์ สุวรรณวิสารท	17
3.1 ตัวอย่างข้อมูลที่เก็บมาจาก โครงการคลังข้อมูลภาษาไทยแห่งชาติ	20
3.2 ตารางแสดงกลุ่มพยัญชนะในภาษาไทยตามการออกเสียง	21
3.3 ตารางกลุ่มตัวอักษรไทยในวิธีการเข้ารหัส ซาวน์เด็กซ์ ของ Odell และ Russell	22
3.4 ตารางแสดงการแทนรหัสตัวอักษรตัวแรกสำหรับการเข้ารหัสตัวอักษรไทย	22
3.5 แสดงการเข้ารหัสสำหรับตัวอักษรตัวที่ไม่ใช่ตัวแรกของคำ	23
3.6 แสดงการเข้ารหัสสำหรับสระ	24
3.7 ตารางแสดงตัวอักษรที่ทำให้เกิดข้อผิดพลาดของโปรแกรม	25
3.8 ตัวอย่างข้อมูลที่ผ่านการเข้ารหัสซาวน์เด็กซ์	26
3.9 ตัวอย่างการเข้ารหัสคำภาษาอังกฤษ	26
3.10 แสดงการแทนรหัสตัวอักษรตัวแรกในขั้นตอนการแปลงรหัสสำหรับ ตัวอักษรโรมัน	28
3.11 แสดงการแทนรหัสสำหรับตัวอักษรตัวที่ไม่ใช่ตัวแรกของคำในขั้นตอนการเข้า รหัสตัวอักษรโรมัน	29
4.1 แสดงรายการคำที่มีรหัสซาวน์เด็กซ์ตรงกันเมื่อค้นหาคำว่า	31
4.2 แสดงการเปรียบเทียบรหัสซาวน์เด็กซ์คำภาษาอังกฤษที่ได้กับซาวน์เด็กซ์ คำภาษาไทย	33

สารบัญรูป

หน้า

รูปที่

3.1 กระบวนการพัฒนาโปรแกรมการถอดอักษรภาษาอังกฤษเป็นภาษาไทยโดยใช้ตัวกลาง ในการถอดอักษรแบบถ่ายเสียง	19
3.2 แผนภาพแสดงการทำงาน โดยรวมของระบบ	27
3.3 แสดงตัวอย่างขั้นตอน การเข้ารหัสชวาว์เด็กซ์ของคำว่า “และ”	30
4.1 ตัวอย่างการแสดงผลคำไทยที่พบเจอโดยการเทียบรหัสชวาว์เด็กซ์ที่ตรงกัน	31
4.2 ตัวอย่างการแสดงผลข้อมูลด้วยการใช้ตัวอักษรที่ไม่ได้ระบุในกฎการถอดอักษรโดย ราชบัณฑิตยสถาน	34
4.3 ตัวอย่างการแสดงความหมายประกอบของคำที่ค้นหา	35



บทที่ 1

บทนำ

1.1 ความเป็นมาและความสำคัญของปัญหา

ภาษาแต่ละภาษามีลักษณะที่แตกต่างกันออกไป แต่การออกเสียงยังคงมีลักษณะที่คล้ายกัน เช่น การใช้พยัญชนะ สระ ตัวสะกด รวมกันเป็นคำ ในรูปของตัวอักษรที่ต่างกัน หรืออาจจะคล้ายกันแล้วแต่ภูมิภาคที่ใกล้เคียงกัน มนุษย์เราใช้อักษรเหล่านี้ในการบันทึก หรือศิลปะวรรณกรรม แต่ส่วนที่สำคัญที่สุดของตัวอักษรเหล่านี้คือ การสื่อสาร ตัวอักษรสามารถที่จะเป็นเครื่องมือในการสื่อสารได้เป็นอย่างดี

การพิมพ์เป็นขั้นตอนในการสื่อสารผ่านตัวอักษรในโลกเทคโนโลยี การที่เราจะพิมพ์บนแป้นพิมพ์ในภาษาต่างๆ ก็ต้องอาศัยความรู้และความเข้าใจในไวยากรณ์และคำศัพท์ของภาษานั้นๆ เช่นเดียวกับกับภาษาไทย ที่มีหลักภาษาและตัวอักษรที่มีลักษณะเฉพาะ

เนื่องจากมีชาวต่างชาติที่สนใจที่จะศึกษาภาษาไทย ที่ต้องการพิมพ์ภาษาไทยได้อย่างถูกต้อง แต่ยังมีปัญหาเรื่องของการเขียนภาษาไทยหรือการสะกดคำต่างๆ อยากพิมพ์ภาษาไทยให้ได้ตรงตามความหมายที่จะสื่อสารได้อย่างถูกต้อง หรือต้องการค้นคว้าวิจัยแหล่งข้อมูลต่างๆ ที่เป็นภาษาไทย ได้อย่างมีประสิทธิภาพ สำหรับคนไทยที่ต้องใช้เป็นพิมพ์ภาษาต่างประเทศและยังจำตำแหน่งต่างๆ ของแป้นพิมพ์ภาษาไทยไม่ได้ ก็สามารถที่จะพิมพ์ได้ถูกต้องเช่นกัน

สำหรับปริญญาโทปีนี้จะมุ่งเน้นในการพัฒนา Input Method Editor (IME) ที่สามารถแปลงตัวอักษรโรมันที่ผสมกันแล้วอ่านออกเสียงคล้ายคำที่มีความหมายในภาษาไทย แล้วแปลงคำเป็นคำนั้นๆ ในภาษาไทย โดยสามารถนำไปใช้ได้กับระบบปฏิบัติการ Windows

1.2 ความมุ่งหมายและวัตถุประสงค์ของการศึกษา

1.2.1 เพื่อศึกษาและพัฒนาตัวแก้ไขวิธีการป้อนข้อมูล (IME) สำหรับการป้อนข้อมูลลงในคอมพิวเตอร์ด้วยการพิมพ์ตัวอักษรบน US – Standard Keyboard และแสดงผลเป็นตัวอักษรภาษาไทย

1.2.2 เพื่อจะช่วยอำนวยความสะดวกให้กับชาวต่างชาติที่ต้องการพิมพ์ภาษาไทย เพื่อการเรียนรู้ และเขียนได้อย่างถูกต้อง

1.2.3 เพื่อจะช่วยอำนวยความสะดวกแก่คนไทยที่ใช้แป้นพิมพ์ภาษาอังกฤษ (ที่ไม่มีภาษาไทย) เพื่อพิมพ์ได้อย่างถูกต้อง

1.3 สมมติฐานของการศึกษา

1.3.1 คำทุกคำในภาษาไทยสามารถนำมาพิมพ์ทับศัพท์ด้วยอักษรโรมันได้ทุกคำ (อาจจะมีการเขียนเหมือนกันหรือคล้ายกัน)

1.3.2 คำบางคำในภาษาไทยนำเอามาจากภาษาอังกฤษ คำเหล่านี้ต้องใช้กฎการถอดอักษรจากงานวิจัยอื่นๆ หรือสร้างขึ้นเองโดยอ้างอิงหลักการของราชบัณฑิตยสถาน

1.3.3 โปรแกรมหรือ IME ที่สร้างขึ้นสามารถนำไปใช้เป็นเครื่องมือในการช่วยพิมพ์ได้จริง และมีความถูกต้องไม่ต่ำกว่า 90%

1.4 ทฤษฎีหรือแนวคิดที่ใช้ในโครงการ

การ Romanization คือการแปลงจากการเขียนด้วยอักษรภาษาดั้งเดิม เป็นการเขียนด้วยอักษรโรมัน อาศัยการออกเสียงของตัวอักษรโรมันมาประสมกันให้มีเสียงอ่านที่เหมือนหรือคล้ายกันกับภาษาดั้งเดิม โดยในโครงการนี้จะยึดหลักการการแปลงอักษรตามแบบฉบับของราชบัณฑิตยสถาน ตัวระบบเราจะใช้ภาษา C# ในการพัฒนาเป็นหลัก การหาคำต่างๆเราใช้อัลกอริทึม Metaphone เป็นตัวในการจำแนกคำโดยการสะกดคำ

1.5 ขอบเขตโครงการ

1.5.1 คำศัพท์ในตัวอักษรภาษาอังกฤษไม่รวมถึงคำย่อหรือรหัสดิจิทัล (Acronym)

1.5.2 คำที่ใช้ในการทดสอบขั้นตอนวิธี จะใช้หลักเกณฑ์การถอดอักษรของราชบัณฑิตยสถาน

1.6 วิธีการดำเนินงาน

1.6.1 ศึกษาขั้นตอนวิธีการค้นคืนสารสนเทศข้ามภาษา

1.6.2 ศึกษาหลักภาษาในการถอดอักษร และหลักเกณฑ์การทับศัพท์จากภาษาอังกฤษเป็นภาษาไทย และจากภาษาไทยไปเป็นภาษาอังกฤษ

1.6.3 ศึกษาขั้นตอนวิธีการใช้อัลกอริทึมชานว์เด็กซ์ ภาษาอังกฤษและภาษาไทย

1.6.4 ออกแบบและพัฒนาขั้นตอนวิธีการเข้ารหัสคำทับศัพท์และคำไทยเพื่อการค้นคืนข้ามภาษา

1.6.5 ออกแบบวิธีการทดสอบขั้นตอนวิธีในข้อ 4

1.6.6 ทดสอบและปรับปรุงคุณภาพของขั้นตอนวิธี

1.6.7 สรุปผลการวิจัย

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

1.7 โครงสร้างของปฏิญญานิพนธ์

ปฏิญญานิพนธ์ฉบับนี้ประกอบไปด้วยเนื้อหา 5 บท โดยแต่ละบทจะมีเนื้อหา ดังนี้

บทที่ 1 บทนำ บอกถึงความเป็นมาและความสำคัญของปัญหา วัตถุประสงค์ของการศึกษา ประโยชน์ ขอบเขตโครงการ และวิธีการดำเนินงาน

บทที่ 2 ทฤษฎีและหลักการ บอกถึงทฤษฎีและหลักการต่างๆที่ใช้ในการดำเนินโครงการ

บทที่ 3 วิเคราะห์และ การออกแบบ แสดงถึงการวิเคราะห์ห้ต่างเพื่อที่จะมาใช้ในการพัฒนาโครงการ และการออกแบบโครงสร้าง ลักษณะหน้าตาหรือรูปแบบของโครงการ

บทที่ 4 ระบบต้นแบบ นำเสนอระบบที่ได้จากการพัฒนาลักษณะการใช้งานต่างๆ

บทที่ 5 สรุปและประเมินผลโครงการ วิเคราะห์ผลที่ได้จากการพัฒนาโครงการ และ

นำเสนอ



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 2

เอกสารและงานวิจัยที่เกี่ยวข้อง

ในบทนี้จะกล่าวถึงทฤษฎีต่าง ๆ ที่เกี่ยวกับหลักภาษาไทย และ หลักการถอดอักษร ได้แก่ การถอดอักษร ประวัติความเป็นมาของการถอดอักษรในภาษาไทย การถ่ายเสียงด้วยตัวอักษร โรมัน ตามหลักเกณฑ์ของราชบัณฑิตยสถาน การเข้ารหัสชาวน์เด็กซ์ ขึ้นตอนวิชาวน์เด็กซ์ภาษาอังกฤษ ขึ้นตอนวิชาวน์เด็กซ์ภาษาไทย งานวิจัยการใช้หลักคำพ้องเสียง เพื่อค้นหาชุดอักษรภาษาไทยที่ ออกเสียงเหมือนกัน [4] งานวิจัย A Thai Soundex System for Spelling Correction [3] งานวิจัยการเข้ารหัสคำทับศัพท์เพื่อการค้นคืนข้ามภาษาไทย-อังกฤษ [1]

2.1 การถอดอักษร

การถอดอักษร (Transliteration) หมายถึง การนำคำในภาษาหนึ่งมาเขียนด้วยตัวอักษรอีก ภาษาหนึ่งแบบอักษรต่ออักษร โดยพยายามให้หน่วยเสียงของอักษรทั้งสองภาษาใกล้เคียงกันมากที่สุด [1] เช่นคำว่า “CAMBRIDGE” ในภาษาอังกฤษถอดอักษรเป็น “เคมบริดจ์” ในภาษาไทย เป็นต้น การถอดอักษรระหว่างภาษาเป็น กระบวนการหนึ่งที่เกิดขึ้น เมื่อเราต้องการนำคำจากภาษาหนึ่ง มาใช้ในอีกภาษาหนึ่ง ซึ่งในแต่ละภาษาก็จะมีวิธีการเขียนที่แตกต่างกัน เนื่องจากระบบเสียงในแต่ละภาษานั้นแตกต่างกัน จึงต้องใช้คำจากภาษาต่างประเทศมาช่วยในการถอดอักษรระหว่างภาษา นี้ขึ้น และนอกจากนั้นยังช่วยให้ชาวต่างชาตินั้นอ่านและใช้คำของเราได้ การถอดอักษรระหว่างภาษา นี้ทำได้หลายวิธีและมีคำศัพท์ที่ใช้แตกต่างกันอยู่ 3 คำหลักๆ ได้แก่ Transliteration, Transcription และ Romanization

Transliteration หมายถึง การเขียนตัวเขียน ในภาษาหนึ่งด้วยตัวเขียนที่สอดคล้องกันในอีก ภาษาหนึ่งคำว่า Transliteration นี้มักจะหมายถึง “การถอดอักษร” แต่บางคนก็ใช้คำอื่น เช่น “การ ถ่ายถอดอักษร” “การถ่ายถอดตัวอักษร” “การถ่ายรูปร่างอักษร” เวลาที่ชาวตะวันตกศึกษาเกี่ยวกับเรื่อง เกี่ยวกับประเทศอื่น ๆ ที่มีระบบตัวเขียนต่างออกไป ก็มักใช้วิธีการ Transliteration นี้ คือถอดอักษร ภาษานั้นๆออกมาด้วยตัวอักษรที่ตนใช้ [4]

Transcription หมายถึง การบันทึกข้อมูลเสียงของภาษาโดยใช้ระบบการเขียนที่กำหนด มัก หมายถึง “การถ่ายเสียง” (ธีระพันธ์ 2526), การถ่ายถอดเสียง (อุไรศรี และ อรวรรณ 2545) ในที่นี้จะ ขอใช้คำว่า “การถ่ายเสียง” เพื่อให้แตกต่างจากคำ “การถอดอักษร” อย่างชัดเจน นักภาษาศาสตร์ได้ กำหนดสัญลักษณ์สากลที่ใช้ในการถ่ายเสียงภาษาต่างๆ อักษรที่ใช้ในการถ่ายเสียงนั้นเรียกว่า สัท อักษร (Phonetic) และในการถ่ายเสียงนั้น อาจใช้ตัวอักษร โรมันมาแทนคำในภาษาหนึ่งด้วยก็ได้

เอกสารดังนั้นก็คาบเกี่ยวกับความหมายของ Romanization ด้วย [4] ไม่น่าจะพูดให้ไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ดังนั้น Romanization จึงหมายถึงการแทนคำในภาษาหนึ่งด้วยตัวอักษรโรมัน ซึ่งอาจใช้วิธีการถอดอักษร (Transliteration) หรือ วิธีการถ่ายเสียง (Transcription) ก็ได้ แต่ความจริงมีคำไทยอีกประเภทหนึ่งซึ่งมีความหมายทำนองเดียวกับ Romanization คือคำว่า “การทับศัพท์” เพียงแต่เรามักใช้ในความหมายของการเขียนคำภาษาต่างประเทศด้วยอักษรไทย ในขณะที่ Romanization เป็นการเขียนคำภาษาไทยหรือภาษาอื่นๆด้วยอักษรโรมัน แต่จริงๆแล้วหลักการของการทับศัพท์อาจจะใช้วิธีการถอดอักษรหรือการถ่ายเสียงก็ได้ (ความหมายตามพจนานุกรมไทยฉบับราชบัณฑิตยสถาน พ.ศ. ๒๕๔๒ คือ การที่รับเอาคำของภาษาหนึ่งมาใช้ในอีกภาษาหนึ่ง โดยวิธีการถ่ายเสียงและถอดอักษร)

2.2 ประวัติความเป็นมาของการถอดอักษรในภาษาไทย

คริส โวลด์ได้กล่าวถึงความเป็นมาของการถอดอักษรไทยไว้ว่า มีมาเนิ่นนานแล้ว ตัวอย่างเช่น ชาวโปรตุเกสในศตวรรษที่ 16 เรียกเมืองไทยว่า SIAO หรือ MUANTAI มีเมืองหลวงชื่อ HUDIA (อยุธยา) ซึ่งคริสโวลด์เห็นว่าเป็นการถอดอักษรไทยที่ไม่เหมาะสม การถอดอักษรไทยโดยชาวต่างประเทศนั้นมักเป็นไปตามอำเภอใจของแต่ละคน ในสมัยรัชกาลสมเด็จพระนารายณ์มหาราช (พ.ศ. 2200-2231) กลุ่มมิชชันนารีชาวฝรั่งเศสได้คิดวิธีการถอดอักษรไทยเพื่อให้ได้เสียงใกล้เคียงกับเสียงภาษาไทย เช่น “เมื่อท่านมาเราได้กินสำเร็จแล้ว” จะเขียนว่า “MEÛÀ TÂN MÃ RÃO DÁI KIN SAM-RED LÊOU” ซึ่งปรากฏอยู่ในหนังสือของลาตูแบร์ เอกอัครราชทูตของพระเจ้าหลุยส์ที่ 14 แต่ดูเหมือนว่าชาวต่างประเทศต่าง ๆ ก็จะมีวิธีการเขียนคำไทยโดยอิงกับระบบภาษาของตน ไม่ได้ใช้ระบบที่บาทหลวงฝรั่งเศสเสนอขึ้น จนในสมัยรัชกาลที่ 5 เมื่อเริ่มมีปัญหาเรื่องการเขียนชื่อทางภูมิศาสตร์ในแผนที่ซึ่งไทยได้ทำร่วมกับอังกฤษและฝรั่งเศส จึงมีความพยายามร่วมกันระหว่างไทยและฝรั่งเศสว่าจะสร้างระบบการเขียนคำไทยด้วยตัวอักษรโรมัน แต่ระบบนี้ก็ไม่ได้เป็นที่แพร่หลาย เพราะปัญหาเรื่องการถอดอักษรไทยเป็น โรมันนี้เป็นที่ถกเถียงกันมากขึ้น ในหมู่นักวิชาการในเวลาต่อมา[9] จะเห็นได้จากบทความของพระบาทสมเด็จพระมงกุฎเกล้าเจ้าอยู่หัว[10] ที่ทรงวิตกกังวลว่า หากปล่อยให้ต่างคนต่างเขียนก็จะทำให้เกิดความสับสน จึงได้เสนอให้ใช้ระบบที่ยึดตามตัวสะกดของคำศัพท์เดิมในภาษาบาลีสันสกฤตหรือที่รู้จักกันในชื่อของระบบกราฟิก (Graphic System หรือ Hunterian System) พระองค์ทรงไม่เห็นด้วยกับการถอดอักษรโดยวิธีการถ่ายเสียง จนกระทั่งปี พ.ศ. 2474 กระทรวงกรมการ (กระทรวงศึกษาธิการในปัจจุบัน) ก็ได้ตั้งคณะกรรมการเพื่อพิจารณาเรื่องหลักการถอดอักษรไทยเป็นโรมันที่จะใช้เป็นมาตรฐานร่วมกันและได้นำเสนอผลเพื่อรับฟังความคิดเห็นจากกระทรวงกลาโหม ราชบัณฑิตยสถาน สยามสมาคม และศาสตราจารย์ออร์ซ เซเคส์ แห่งวิทยาลัยฝรั่งเศสของตะวันออกไกลและต่อมาราชบัณฑิตยสถานได้

สานงานต่อจนเป็นที่มาของประกาศเกณฑ์การถอดอักษรไทยเป็นโรมันฉบับแรกในปี พ.ศ. 2482 เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการไขว่คว้าเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า โดยกำหนดไว้สองระบบคือ ระบบทั่วไป (General System) และ ระบบพิศดาร (Precise System) ไม่วางกรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมีเหตุดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ทั้งสองระบบจะใช้วิธีการต่างกัน โดยระบบทั่วไปจะใช้วิธีการถ่ายเสียง เช่นคำว่า กษัตริย์ ถ่ายเสียงได้เป็น KASAT แต่ระบบพิสดารจะใช้วิธีการถอดอักษร ก็จะถอดอักษรเป็น KASATRIY แต่ประกาศที่ออกมานั้นก็ไม่เป็นที่ใช้งานกันอย่างแพร่หลายนัก มีเพียงบางกลุ่ม บางหน่วยงานที่ใช้เกณฑ์ดังกล่าว[6] เช่น กรมทางหลวง ใช้เขียนชื่อถนน กรมแผนที่ทหาร การรถไฟ ใช้เขียนชื่อสถานีรถไฟ กรุงเทพมหานคร ใช้เขียนชื่อถนนและซอยต่างๆ หลังจากนั้นจึงมีการพัฒนาต่อมาเรื่อยๆ กระทั่งปี พ.ศ. 2541 มาตรฐาน ISO11940:1998 ได้รับการรับรองให้ใช้เป็นแบบ Transliteration คือแต่ละตัวอักษรจะถอดแบบหนึ่งต่อหนึ่ง และได้รับการพัฒนามาเรื่อยๆ จนกระทั่งปีพ.ศ. 2542 ราชบัณฑิตยสถานได้ยกเลิกประกาศเกณฑ์การถอดอักษรไทยเป็นอักษรโรมันเดิมและยกเลิกแถลงการณ์เรื่องการเขียนชื่อจังหวัด เขตอำเภอ และกิ่งอำเภอ โดยใช้หลักการเขียนแบบทั่วไป และให้ใช้หลักเกณฑ์การถอดอักษรไทยเป็นโรมันแบบถ่ายเสียงที่ได้ปรับปรุงใหม่ให้เหมาะสมยิ่งขึ้น[5]

2.3 การถอดอักษรไทยเป็นอักษรโรมันแบบถ่ายเสียงของราชบัณฑิตยสถาน

การถอดอักษรไทยเป็นอักษรโรมันมี 2 แบบ ได้แก่ การถอดอักษรตามวิธีเขียน (Transliteration) และการถอดอักษรตามวิธีอ่าน หรือการถอดแบบถ่ายเสียง (Transcription) การถอดอักษรตามวิธีเขียน เป็นการถอดตามตัวอักษรทุกตัวที่มีอยู่ในคำนั้นๆ[6] มีข้อดีคือทราบที่มาของคำสามารถสะกดคำได้ และสามารถถอดอักษรกลับไปยังคำเดิม (retransliteration) ได้ ข้อดีคืออ่านออกเสียงได้ยาก แม้จะใส่เครื่องหมายออกเสียงช่วยไว้ด้วยก็ตาม และคำที่ถอดออกมาดูยุ่งยาก รุงรัง ส่วนการถอดอักษรแบบถ่ายเสียงเป็นการถอดอักษรที่จำเป็น ต้องถอดครบทุกตัวอักษร เน้นการออกเสียงเป็นหลัก ข้อดีคือทำให้อ่านคำไทยที่ถอดเป็นอักษรโรมันได้เสียงใกล้เคียง คำที่ถอดออกมา มีความกะทัดรัด เข้าใจง่าย ข้อคือคือ ไม่เอื้อให้ถอดอักษรกลับคืนสู่คำเดิมได้

ตัวอย่างการถอดอักษรตามวิธีเขียน มะม่วง → Mamwng

ตัวอย่างการถอดอักษรแบบถ่ายเสียง มะม่วง → Mamuang

ระบบการถอดอักษรด้วยวิธีการถ่ายเสียง นั้นคือการแทนตัวอักษรในภาษาหนึ่งด้วยรูปแทนเสียงของภาษานั้น เพื่อให้อ่านคำภาษาไทยที่เขียนด้วยตัวอักษรโรมันให้ได้เสียงใกล้เคียง โดยไม่คำนึงถึงการสะกดตัวการ์นต์และวรรณยุกต์[6] เช่น จันทร์ = chan, พระ = phra, แก้ว = kao การถอดอักษรไทยเป็นอักษรโรมันแบบถ่ายเสียง เป็นหลักการถอดตัวอักษรไทยเป็นอักษรโรมันอย่างเป็นทางการ โดยราชบัณฑิตยสถาน ใช้สำหรับหนังสือและสิ่งพิมพ์ของรัฐบาล และป้ายชื่อถนนต่างๆ ในประเทศไทย โดยใช้หลักเกณฑ์การถอดอักษรไทยเป็นอักษรโรมันตั้งแต่พ.ศ. 2482 และในพ.ศ. 2542 ได้ปรับปรุงให้เหมาะสมยิ่งขึ้น โดยยกเลิกประกาศเดิมและให้ใช้หลักเกณฑ์การถอดอักษรไทยเป็นอักษรโรมันแบบถ่ายเสียง[5]

การถอดอักษรไทยเป็นโรมันตามแบบราชบัณฑิตยสถานนี้ เป็นที่ยอมรับและใช้กันอย่าง
เอกสารนี้เป็นเอกสารทูลเกล้าฯ ถวายเพื่อใช้ในการเรียนการสอนเพื่อการศึกษาเท่านั้น เมื่ออนุญาตให้นำไปใช้ประโยชน์ตามการ
แพร่หลาย โดยเฉพาะตามหน่วยงานราชการต่าง ๆ รวมทั้งได้รับการรองรับจากองค์การ
ไม่หวังกำไรใดๆ ทั้งสิ้น อีกทั้งห้ามมีเหตุดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งหากมีการนำไปใช้

สหประชาชาติให้ใช้เป็นมาตรฐานในการถอดอักษรไทยเป็นโรมันสำหรับชื่อต่างๆ อย่างไรก็ตาม ก็มีผู้วิจารณ์ว่าระบบที่ใช้นี้ไม่เหมาะสมสำหรับใช้สอนภาษาไทยให้ชาวต่างชาติ เนื่องจากไม่ได้มีการกำหนดวรรณยุกต์ ความสั้นหรือยาวของสระ รวมทั้งไม่แยกความแตกต่างของเสียงสระ เช่น โอะ – โอ ฯลฯ และไม่ได้แยกความแตกต่างของเสียงระหว่าง จ กับ ช ด้วย จึงมีผลทำให้ชาวต่างประเทศออกเสียงผิดได้ แต่ทั้งนี้คงต้องเข้าใจว่า จุดมุ่งหมายหลักของระบบแบบราชบัณฑิตยสถานนั้น ไม่ได้มุ่งเน้นที่การเรียนการสอนภาษาไทยให้ชาวต่างประเทศ แต่ต้องการให้ชาวต่างประเทศโดยทั่วไปที่จำเป็นต้องติดต่อสื่อสาร สามารถอ่านและเขียนคำไทยด้วยตัวอักษรโรมันได้

องค์การระหว่างประเทศว่าด้วยการมาตรฐาน (International Organization for Standardization : ISO) ได้กำหนดมาตรฐานการถอดอักษรไทยเป็นอักษรโรมันขึ้นในปี 1998 เพื่อจัดให้มีเครื่องมือสำหรับการสื่อสารสากลในการเขียนที่ชัดเจน ไม่คลุมเครือ เพื่อการถ่ายทอดอัตโนมัติที่สามารถถอดกลับได้ ถูกต้องสมบูรณ์ โดยอาจจะเป็นคนหรือเครื่องจักร โดยได้กำหนดสัญลักษณ์สำหรับการออกเสียงสระ สัญลักษณ์เพื่อแยกความแตกต่างระหว่างพยัญชนะที่ออกเสียงเหมือนกัน ด้วยวิธีนี้ทำให้อักษรแต่ละตัวถูกลดเป็นอักษรโรมันและสัญลักษณ์กำกับที่ไม่ซ้ำกันสามารถถอดเป็นอักษรโรมันและถอดกลับเป็นอักษรไทยได้อย่างถูกต้อง

ตารางที่ 2.1 ตารางเทียบเสียงพยัญชนะ และ สระ ตามหลักเกณฑ์ที่กำหนดโดยราชบัณฑิตยสถาน

พยัญชนะ			สระ	
พยัญชนะไทย	อักษรโรมัน		สระไทย	อักษรโรมัน
	ตัวต้น	ตัวสะกด		
ก	k	k	อะ, ะ (อะ ลดรูป), รร (มีตัวสะกด), อา	a
ข ฃ ค ฅ ฌ	kh	k	รร (ไม่มีตัวสะกด)	an
ง	ng	ng	อ่า	am
จ ฉ ช จ	ch	t	อิ, อี้	i
ซ ฌร (เสียง ซ)	s	t	อึ, อี้	ue
ศ ษ ส	s	t	อุ, อู	u
ญ	y	n	เอะ, ะ (เอะ ลดรูป), เอ	e
ฎ ฏ (เสียง ด) ด	d	t	แอะ, แอ	ae
ฏ ต	t	t	โอะ, - (โอะ ลดรูป), โอ, เอาะ, ออ	o
ฐ ฑ ฒ ถ ฑ ฐ	th	t	เออะ, ะ (เออะ ลดรูป), เออ	oe
ณ น	n	n	เอียะ, เอีย	ia

เอกสารนี้เป็นเอกสารที่สงวนลิขสิทธิ์ไว้สำหรับนรใช้งานเพื่อเอื้อประโยชน์เท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 2.1 (ต่อ) ตารางเทียบเสียงพยัญชนะและสระตามหลักเกณฑ์ที่กำหนดโดย

ราชบัณฑิตยสถาน

พยัญชนะ			สระ	
พยัญชนะไทย	อักษรโรมัน		สระไทย	อักษรโรมัน
	ตัวต้น	ตัวสะกด		
บ	b	p	เอื้อะ, เอื้อ	uea
ป	p	p	อ้าวะ, อ้าว, -ว- (อ้าว ทรูปร)	ua
พ พ ภ	ph	p	ไอ, ไอ, อัย, ไอย, อาย	ai
ฝ ฟ	f	p	เอา, อาว	ao
ม	m	m	อูย	ui
ย	y	-	โอย, ออย	oi
ร	r	n	เอย	oei
ล พ	l	n	เอื้อย	ueai
ว	w	-	อวย	uai
ห ส	h	-	อิ้ว	io
			เอื้อว, เอว	eo
			แเอ้ว, แเอว	aeo
			เอื้อยว	iao
			ฤ (เสียง รึ) ฤ๑	rua
			ฤ (เสียง ริ)	ri
			ฤ (เสียง เรอ)	roe
			ฤ, ฤ๑	lue

2.4 วิธีเชิงสัทลักษณะ (Phonetic algorithm)[12]

ขั้นตอนวิธีเชิงสัทลักษณะ (Phonetic algorithm) คือขั้นตอนวิธีสำหรับการกำหนดดัชนีของคำต่างๆ โดยใช้การออกเสียงเป็นเกณฑ์ ขั้นตอนวิธีเชิงสัทลักษณะส่วนใหญ่พัฒนาขึ้นเพื่อใช้กับเอกสารภาษาอังกฤษ ดังนั้นการใช้กฎเกณฑ์ดังกล่าวกับคำในภาษาอื่นอาจไม่ให้ผลลัพธ์ที่มีความหมายไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ขั้นตอนวิธีเหล่านี้มีความซับซ้อน ด้วยกฎและข้อยกเว้นหลายประการในภาษานั้นๆ เนื่องจากการสะกดคำ และการออกเสียงในภาษาอังกฤษ[12]

ขั้นตอนวิธีเชิงสัญลักษณ์ที่เป็นที่รู้จักนั้นมีดังนี้

- ซาวนด์เด็กซ์ ถูกพัฒนาขึ้น โดย Robert C. Russell และ Margaret K. Odell พัฒนาขึ้นเพื่อใช้เข้ารหัสสำหรับชื่อคนในภาษาอังกฤษในการทำสำมะโนประชากร ใช้แก้ปัญหาเนื่องจากบางคำนั้นสามารถสะกดได้หลายแบบอ่านออกเสียงคล้ายกัน อาจจะแตกต่างกันเพราะสำเนียง หรือวัฒนธรรมซึ่งจะเป็นปัญหาหลักในการเก็บข้อมูลประวัติของชื่อบุคคล
- Metaphone และ Double Metaphone ใช้สำหรับคำในภาษาอังกฤษทั่วไป ไม่เพียงชื่อเฉพาะบุคคล ขั้นตอนวิธีนี้เป็นรากฐานของซอฟต์แวร์ตรวจสอบการสะกดคำในหลายโปรแกรม (Spell Checking)
- New York State Identification and Intelligence System (NYSIIS) เป็นการจับคู่หน่วยเสียง (phoneme) ที่คล้ายกันให้เข้ากับอักษรตัวเดียวกัน ให้ผลลัพธ์เป็นสายอักขระซึ่งสามารถอ่านได้โดยไม่ต้องถอดรหัส
- การเข้าหาด้วยการจัดระดับคู่ (Match Rating Approach) ถูกพัฒนาขึ้นโดยสายการบินเวสเทอร์นแอร์ไลน์เมื่อปี พ.ศ. 2520 ใช้เทคนิคการเข้ารหัสและเปรียบเทียบและจัดระดับความคล้ายคลึงให้กับชื่อภาษาอังกฤษที่มีลักษณะพ้องเสียงกัน ประโยชน์ของขั้นตอนวิธีนี้คือเพิ่มความแม่นยำและรวดเร็วในการค้นหาชื่อภาษาอังกฤษไม่เพียงแต่ตรวจหาคำที่ตรงที่สุดเพียงอันเดียวแต่ครอบคลุมไปถึงชื่อที่มีลักษณะพ้องเสียงกันเมื่ออาจจะมีการสะกดที่แตกต่างกัน

2.5 ขั้นตอนวิธีซาวนด์เด็กซ์ภาษาอังกฤษ

M.K. Odell และ R. C. Russell ได้ออกแบบขั้นตอนวิธีเข้ารหัสชื่อในภาษาอังกฤษโดยยึดหลักของการอ่านออกเสียง เพื่อให้ชื่อที่อ่านออกเสียงคล้ายกันก็จะได้รับรหัสเหมือนกัน หรือที่เรียกว่า “ซาวนด์เด็กซ์” ขั้นตอนวิธีดังกล่าวได้ใช้แนวคิดทางภาษาศาสตร์และตัวเลขที่ว่าชื่อในภาษาอังกฤษสามารถจำแนกความแตกต่างได้โดยพิจารณาเพียงพยัญชนะเท่านั้น [13]

ขั้นตอนวิธีเข้ารหัสของซาวนด์เด็กซ์ โดยเริ่มจากการนำตัวอักษรตัวแรกของคำไปเป็นรหัส ส่วนตัวอักษรที่เหลือจะแปลงเป็นตัวเลขโดยใช้ตารางการกำหนดรหัสซาวนด์เด็กซ์ ดังแสดงในตารางที่ 2.2 จากนั้นจะตัดรหัสตัวเลขศูนย์ออกไป แล้วถักรหัสตัวเลขที่อยู่ตำแหน่งติดกันมีค่าเท่ากัน

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 2.4 การกำหนดรหัสตัวอักษรของรหัสชาวน์เด็กซ์ภาษาไทยสำหรับตัวอักษรที่เหลือ

รหัสตัวเลข	ตัวอักษร
0	ม ว ำ
1	ก ข ฃ ค ฅ ฆ
2	ง ย
3	ญ ฌ ฎ
4	ฎ ฏ ค ต ศ ษ ส
5	บ ป พ ภ
6	ฝ ฟ พ ห อ ฮ
7	จ ฉ ช ซ ฌ ฎ
8	ฐ ฑ ฒ ถ ท ธ
9	ร ล พ ฤ ฎ

ส่วนข้อกำหนดเพื่อให้เหมาะสมกับภาษาไทยแบ่งเป็นกรณีต่าง ๆ ดังนี้

- กรณีพบ ใ- ใ- ใ-ย และ ัย จะเปลี่ยนให้อยู่ในรูปแบบเดียวกันคือ ัย ก่อนทำการเข้ารหัส เพื่อให้ได้รหัสชาวน์เด็กซ์ที่เหมือนกัน เนื่องจากสระดังกล่าวอ่านออกเสียงเหมือนกัน เช่น ไท ไทไทย และ ทัย จึงจะเปลี่ยนเป็น ทัย
- กรณีพบ รร จะเปลี่ยนเป็น ริน เมื่อไม่มีตัวสะกดตามหลัง และจะเปลี่ยนเป็น ร์ เมื่อมีตัวสะกดตามหลัง
- กรณีพบการันต์ จะตัดการันต์และพยัญชนะที่มีตัวการันต์กำกับรวมทั้งสระและอักษรควบการันต์ทิ้ง เช่น คำว่า จันทร ศักดิ์ และ พันธุ์ เปลี่ยนเป็น จัน ศัก และ พัน ตามลำดับ

ตารางที่ 2.5 ตัวอย่างการเข้ารหัสชาวน์เด็กซ์ภาษาไทย

ชื่อภาษาไทย	รหัสชาวน์เด็กซ์ภาษาไทย
อัมพร	อ059000
อำภรณ์	อ059000
พรรณศักดิ์	พ341000
พันธุ์ศักดิ์	พ341000
เนืองนิจ	น6234000
เนืองนิตย์	น6234000

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2.6.2 งานวิจัย A Thai Soundex System for Spelling Correction[3]

งานวิจัย A Thai Soundex System for Spelling Correction ได้นำเสนอวิธีในการเข้ารหัสสำหรับชาวเน็ตเด็กซ์ภาษาไทย โดยมีแนวคิดว่าจะระบบหน่วยคำในภาษาไทยไม่มีเครื่องหมายแบ่งพยางค์ จึงทำให้บางคำนั้น สามารถอ่านออกเสียงได้หลายแบบ เช่น คำว่า ตากลม สามารถอ่านเป็น ตา-กลม หรือ ตาก-ลม ก็ได้ ซึ่งการเข้ารหัสชาวเน็ตเด็กซ์เป็นการเข้ารหัสแบบไม่มีเครื่องหมายแบ่งพยางค์ จึงทำให้รหัสที่ได้เกิดความผิดพลาดไปด้วย[3] ดังนั้นเพื่อลดความผิดพลาดในการอ่านออกเสียง จึงได้ใช้เทคนิค Nondeterministic Finite Automaton with Output เพื่อสร้างรหัสชาวเน็ตเด็กซ์สำหรับการอ่านออกเสียงในทุกกรณีของคำ และในการค้นคืน ก็จะนำรหัสที่ได้ไปเปรียบเทียบกับรหัสทุกตัวของคำที่ต้องการค้นคืน ถ้าพบรหัสที่เหมือนกันเพียงหนึ่ง คู่ก็จะถือว่าทั้งคู่เป็นคู่คำที่อ่านออกเสียงคล้ายกัน[3]

การกำหนดรหัสชาวเน็ตเด็กซ์สำหรับอักขระไทยดังนี้

1. กำหนดรหัสชาวเน็ตเด็กซ์สำหรับพยัญชนะต้น จะแบ่งตามเสียงของพยัญชนะ 21 เสียงของไทย และได้รวม ร กับ ล เป็นกลุ่มเดียวกัน โดยจะใช้พยัญชนะไทย 1 ตัวแทนกลุ่มเสียง ดังแสดงในตารางที่ 2.6
2. กำหนดรหัสชาวเน็ตเด็กซ์สำหรับสระ จะแบ่งตามเสียงของสระและรวมสระเสียงสั้นและยาวเข้าด้วยกัน โดยจะใช้อักขระโรมันมาแทนกลุ่มเสียงสระ ดังตารางที่ 2.7
3. กำหนดรหัสชาวเน็ตเด็กซ์สำหรับตัวสะกด จะแบ่งกลุ่มเสียงตามมาตราต่างๆ ของอักษรไทย และได้เพิ่มกลุ่ม ฮ สำหรับคำที่ไม่มีตัวสะกดดังแสดงในตารางที่ 2.8

ตารางที่ 2.6 การเข้ารหัสสำหรับพยัญชนะต้น

รหัสตัวอักษร	ตัวอักษร	รหัสตัวอักษร	ตัวอักษร
ก	ก	บ	บ
ค	ข ฃ ค ฅ ฌ	ป	ป
ง	ง	พ	พ ภ ผ
จ	จ	ฟ	ฝ ฟ
ช	ช ฉ ฌ	ม	ม
ฌ	ฌ ศ ษ ส	ย	ณ ย
ด	ด ฎ ฑ*	ร	ร ล พ
ต	ต ฏ	ว	ว
ท	ฐ ฑ* ฒ ถ ท ฑ	อ	อ
น	ณ น	ฮ	ห ฮ

* ฑ สามารถอ่านออกเสียงได้สองแบบ คือ /ค/ และ /ท/ จึงสร้างรหัสทั้ง 2 แบบ

เอกสารนี้เป็นเอกสารที่จัดทำขึ้นเพื่อใช้ในการศึกษาวิจัยเท่านั้น ไม่สามารถนำข้อมูลไปใช้

- ใช้รหัสตัวเลขเป็นอักษรตัวแรกของรหัส แทนการใช้ตัวอักษรของคำเป็นตัวแรก เนื่องจากพบว่ามีตัวอักษรอังกฤษหลายตัวถอดเป็นอักษรไทยเหมือนกัน เช่น V และ W ถอดเป็นตัว ว และได้มีการเสนอรหัสเพิ่มอีก 3 ตัว คือ 7 8 และ 9 ซึ่งจะใช้ในกรณีที่เป็นตัวอักษรแรกของคำเท่านั้น โดยมีรายละเอียดดังนี้
 - รหัส 7 สำหรับคำที่ขึ้นต้นด้วย AEIOU หรือ อ เนื่องจากพบว่าคำภาษาอังกฤษที่ขึ้นต้นด้วยสระมักจะใช้ อ เป็นพยัญชนะต้นในการอ่านออกเสียง เช่น ABRAHAM (อับราฮัม) EDWARD (เอ็ดเวิร์ด) OHM (โอห์ม) เป็นต้น
 - รหัส 8 สำหรับคำที่ขึ้นต้นด้วย H เนื่องจาก H ที่เป็นตัวอักษรแรกของคำ มักจะเป็นพยัญชนะและอ่านออกเสียงเป็น ฮ เช่น HOPSKINS (ฮอปกินส์) ส่วน H ที่อยู่ส่วนอื่นของคำมักจะเป็นอักษรควบและไม่อ่านออกเสียง เช่น WHITE (ไวต์) JOHN (จอห์น) SHOW (โชว์) เป็นต้น
 - รหัส 9 สำหรับตัวอักษร Y เนื่องจากพบว่า Y ที่เป็นตัวอักษรแรกของคำ มักจะเป็นพยัญชนะและอ่านออกเสียงเป็น ย เช่น YAHOO (ยาฮู) ส่วน Y ที่อยู่ส่วนอื่นของคำมักจะเป็นสระ เช่น ONYX (โอนิกซ์) PHYSICS (ฟิสิกส์) เป็นต้น
- เพิ่มรหัส 52 สำหรับตัวอักษร ง เนื่องจากอักษร ง ถอดมาจากตัวอักษรอังกฤษ NG โดยรหัส 5 สำหรับอักษร N และรหัส 2 สำหรับ G
- ขยายความยาวของรหัสคำที่ได้จากเดิมคือ 4 หลัก (ตามหลักของ Odell และ Russell) เป็น ไม่จำกัดความยาวของรหัสคำที่ได้ เนื่องจากพบว่าคำที่มีความยาวมาก อาจได้รหัสตรงกันทั้งๆ ที่ออกเสียงไม่คล้ายกันเลย เช่น คริสต์เตียน ได้รหัส 262395 และ คริสต์โตฟเฟล ได้รหัส 262314 ซึ่งถ้าพิจารณาเพียง 4 หลัก คือ 2623 จะถือว่าคำทั้งสองออกเสียงคล้ายกัน
- เพิ่มพารามิเตอร์ K ซึ่งเป็นความยาวน้อยสุดของรหัสคำ โดยพิจารณาเฉพาะคำที่มีความยาวน้อยสุดของรหัสคำที่มากกว่า K หลัก
- วรรณยุกต์ และสระในภาษาไทยจะไม่นำมาพิจารณาในการเข้ารหัส

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 3

การวิเคราะห์และออกแบบระบบ

หลังจากที่ได้ทบทวนงานวิจัยเรื่องการถอดอักษรในภาษาต่างๆแล้วนั้น พบว่าในการถอดอักษรนั้นมีหลายขั้นตอนมีหลายแนวทาง เพื่อให้ได้มาซึ่ง โปรแกรมการถอดอักษรที่สมบูรณ์สำหรับวิทยานิพนธ์เล่มนี้จะเน้นความสำคัญไปที่การพัฒนาโปรแกรมการถอดอักษรภาษาอังกฤษไปเป็นภาษาไทยโดยใช้เกณฑ์การถอดอักษรไทยเป็นโรมันแบบถ่ายเสียงของราชบัณฑิตยสถาน ผู้วิจัยต้องการ จะใช้วิธีการถอดอักษร โดยใช้ตัวกลางในการถอดอักษร หรือรูปแทนเสียง เพราะเห็นว่าหลักเกณฑ์ของราชบัณฑิตยสถานนั้นยังไม่ครอบคลุมมากพอ ยังคงมีปัญหาเรื่องการกำหนดความสัมพันธ์ของการถ่ายเสียง ซึ่งไม่เป็นแบบหนึ่งต่อหนึ่ง ตัวอักษร กล่าวคือ ตัวอักษรโรมันหนึ่งตัวสามารถถ่ายเสียงออกมาได้เป็นหลายตัวอักษรของภาษาไทย นอกจากนี้ในการถอดอักษรภาษาอังกฤษเป็นภาษาไทยโดยไม่ใช้ตัวกลางนั้น มีความจำเป็นที่จะต้องมามีข้อมูลคำไทยและคำที่ถอดเป็นอักษรโรมันคู่กัน ซึ่งไม่มีหน่วยงานใดเผยแพร่ซึ่งงานนี้ ทำให้ไม่สามารถวิจัยด้วยการถอดอักษรโดยตรงโดยไม่ใช้ตัวกลางได้ ดังนั้นผู้วิจัยจึงเห็นว่า การถอดอักษรภาษาอังกฤษเป็นภาษาไทยโดยใช้ตัวกลางในการถ่ายเสียงนั้นเหมาะสมที่สุด โดยผู้วิจัยได้เก็บรวบรวมข้อมูลที่เผยแพร่โดยศูนย์เทคโนโลยีอิเล็กทรอนิกส์และคอมพิวเตอร์แห่งชาติ และ คณะอักษรศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย เป็นคลังข้อมูลที่มีขนาด 7,235 คำ แล้วนำข้อมูลเหล่านี้มาประมวลผลโดยการเข้ารหัสชวาน์เด็กซ์ โดยการเข้ารหัสจะใช้ตารางที่ได้ปรับปรุงจากงานวิจัยของประยูทธ สุวรรณวิสารเรื่องการเข้ารหัสคำทับศัพท์เพื่อการค้นคืนข้ามภาษาไทย – อังกฤษ[1] และงานวิจัยของเทพพิทักษ์ การบุญญานันท์เรื่อง A Thai Soundex System for Spelling Correction[3]

ในรายงานฉบับนี้จะแบ่งออกเป็นสองส่วนคือการเข้ารหัสชวาน์เด็กซ์ภาษาไทย และรหัสชวาน์เด็กซ์สำหรับคำทับศัพท์ ในส่วนของกระบวนการพัฒนาโปรแกรมนั้นมีขั้นตอนแรกเป็นการเก็บรวบรวมข้อมูลคำไทยเพื่อสร้างเป็นคลังคำศัพท์ของราชบัณฑิตยสถานขนาด 7,235 คำ ส่วนขั้นตอนต่อไปจะเป็นการเข้ารหัสชวาน์เด็กซ์เพื่อนำข้อมูลคำไทยที่ผ่านการเข้ารหัสมาใส่ลงในฐานข้อมูลเพื่อใช้เป็นแบบจำลองการเปรียบเทียบ จากนั้นจะนำรหัสที่ได้มาเปรียบเทียบกับข้อมูลตัวอักษรภาษาอังกฤษที่ได้ป้อนข้อมูลเข้ามาว่ามีข้อมูลรหัสชวาน์เด็กซ์ตัวใดที่สอดคล้องกับคำภาษาอังกฤษบ้าง จากนั้นจึงจะนำข้อมูลที่สอดคล้องมาแสดงผลใน โปรแกรม ขั้นตอนสุดท้ายจะเป็น

เอกสารนี้เป็นเอกสารที่จัดทำขึ้นเพื่อใช้ในการเรียนการสอนและเป็นเอกสารที่พัฒนาขึ้น
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

กระบวนการพัฒนาโปรแกรมการถอดอักษรภาษาอังกฤษเป็นภาษาไทยโดยใช้ตัวกลางในการถอดอักษรแบบถ่ายเสียงสามารถแสดงได้ดังรูปที่ 3.1



รูปที่ 3.1 กระบวนการพัฒนาโปรแกรมการถอดอักษรภาษาอังกฤษเป็นภาษาไทยโดยใช้ตัวกลางในการถอดอักษรแบบถ่ายเสียง

3.1 การเก็บรวบรวมข้อมูล

ขั้นตอนแรกของการพัฒนา คือ การเก็บรวบรวมข้อมูล ข้อมูลที่ต้องการคือคำศัพท์ไทยที่พิมพ์เผยแพร่โดยศูนย์วิจัยอิเล็กทรอนิกส์และคอมพิวเตอร์แห่งชาติ และข้อมูลคำไทยใช้บ่อยที่เผยแพร่โดยโครงการคลังข้อมูลภาษาไทยแห่งชาติ ซึ่งจัดรวบรวมโดยคณะอักษรศาสตร์จุฬาลงกรณ์มหาวิทยาลัยทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

มหาวิทยาลัย โดยรวบรวมจากงานประเภทงานเขียน เช่น งานวิชาการ กึ่งวิชาการ เรื่องแต่ง นิยาย โคลงกลอน เรื่องสั้น อื่นๆ รวมกันเป็นจำนวนประมาณ 7,235 คำ

ขั้นตอนการเตรียมข้อมูลเพื่อใช้สร้างฐานข้อมูล ในการดำเนินโครงการ

1. รวบรวมข้อมูลประกอบด้วยคำไทยได้มาทั้งหมด 7,235 คำ
2. นำข้อมูลที่เตรียมไว้มาจัดเก็บลงในฐานข้อมูล Microsoft Access แบ่งเป็น 2 สดมภ์ โดย
 - สดมภ์ที่ 1 จัดเก็บคำไทยที่รวบรวมได้
 - สดมภ์ที่ 2 จัดเก็บ รหัสชวาร์เด็กซ์ของคำไทยที่ได้
3. เขียนโปรแกรมการสลับที่ตัวอักษร ให้ตรงกับการถอดอักษรโดยใช้ตัวอักษรโรมัน โดยจะกล่าวในขั้นตอนต่อไป
4. นำข้อมูลที่ได้จากข้อ 3. ไปเข้ารหัสชวาร์เด็กซ์ทีละคำ โดยจะกล่าวในขั้นตอนต่อไป
5. จัดเก็บรหัสที่ได้ลงในฐานข้อมูล สดมภ์ที่สอง แสดงเป็นตัวอย่าง ได้ดังนี้

ตารางที่ 3.1 ตัวอย่างข้อมูลที่เก็บรวบรวมโดยโครงการคลังข้อมูลภาษาไทยแห่งชาติ

คำศัพท์ภาษาไทย	รหัสชวาร์เด็กซ์ที่ได้
การ	K97500
เป็น	P10150
ใน	N97105
จะ	C97000
ของ	K11152
มี	M41050
และ	L97101
ไม่	M49710
ได้	D97105

ข้อมูลที่รวบรวมได้มีทั้งหมดจำนวนประมาณ 7,235 คำ จัดเก็บอยู่ในตาราง (Table) ทั้งหมด

1 ตาราง โดยจำแนกตามที่มาของการเผยแพร่ หลังจากนั้นจะนำข้อมูลทั้งหมดที่เก็บได้มาเข้ารหัส

ชวาร์เด็กซ์ในขั้นตอนต่อไป การใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

3.2 การเขียนโปรแกรมเข้ารหัสชาวน์เด็กซ์สำหรับข้อมูลคำภาษาไทย

หลังจากที่ได้ข้อมูลมาอยู่ในลักษณะของฐานข้อมูลแล้ว ก็จะเริ่มเรียงลำดับข้อมูลใหม่อยู่ในรูปของการสลับที่ตัวอักษรสระประสมให้สอดคล้องกับลำดับการเขียนในภาษาอังกฤษเนื่องจากชาวน์เด็กซ์จะไม่สามารถแยกคำไทยได้ เพราะคำไทยมีความกำกวมหลายประการ แตกต่างจากภาษาต้นฉบับ โดยเริ่มเขียน โปรแกรมแปลงอัตโนมัติเพื่อความสะดวกและรวดเร็ว เช่น ‘เมื่อ’ ระบบจะเรียงลำดับตัวอักษรที่เป็นสระประสมเป็น ‘มื่อ’ จะตรงกับการเขียนทับศัพท์คือ ‘muea’ แล้วก็จะนำคำที่ได้เก็บไว้ลงตารางรอเข้าโปรแกรมแปลงรหัสอีกครั้ง หลังจากนั้น จะเขียน โปรแกรมเข้ารหัสชาวน์เด็กซ์ เพื่อถอดอักษรคำไทยให้อยู่ในรูปข้อมูลของรูปแทนเสียง เพื่อใช้ในการเปรียบเทียบคำภาษาอังกฤษ ด้วยวิธีการเข้ารหัสแบบอัตโนมัติ นั่นคือการนำคำภาษาไทยมาเข้ารหัส โดยใช้โปรแกรมการเข้ารหัส เหตุผลที่ใช้โปรแกรมเข้ารหัสเพราะเป็นวิธีที่สะดวกและรวดเร็ว

ในการทำโครงการนี้ ชาวน์เด็กซ์ที่นำมาใช้เป็นตัวที่ถูกพัฒนาเพื่อใช้สำหรับภาษาไทย โดยหลักการคือ แบ่งกลุ่มพยัญชนะ ในภาษาไทยออกตามการออกเสียงที่เหมือนกันเป็น 21 กลุ่มด้วยกัน ดังตารางที่ 3.2

ตารางที่ 3.2 ตารางแสดงกลุ่มพยัญชนะในภาษาไทยตามการออกเสียง

ก	ฎ ด	ฝ ฟ
ข ฃ ค ฅ ฌ	ฏ ต	ม
ง	ฐ ฑ ฒ ถ ท ธ	ร
จ	ณ น	ล พ
ฉ ช จ	บ	ว
ซ ศ ษ ส	ป	ห ฮ
ญ ย	ผ พ ภ	อ

จากนั้นนำ 21 กลุ่มมาจัดให้เหลือ 7 กลุ่ม อ้างอิงจากตาราง อัลกอริทึมชาวน์เด็กซ์ ของ Odell and Russell (ในตาราง 2.2) โดยอาศัยการออกเสียงที่คล้ายกัน จะได้ผลลัพธ์ตามตาราง 3.3

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 3.3 ตารางกลุ่มตัวอักษรไทยในวิธีการเข้ารหัสชาวเน็ตเวิร์กของ Odell and Russell

ตัวอักษรอังกฤษ	ตัวอักษรไทย
A E I O U H W Y	อ ห อ ว ญ ย
B F P V	บ ภ ฟ ป ผ พ ภ ว
C G J K Q S X Z	ข ข ค ค ฅ ฌ ช ฌ ก จ ฅ ศ ษ ส
D T	ฎ ก ฎ ต ฐ ฑ ฒ ถ ฑ ฐ
L	ล ฬ
M N	ม ฌ น
R	ร

สำหรับการเขียน โปรแกรมการเข้ารหัสจะยึดหลักเกณฑ์งานวิจัยการเข้ารหัสคำทับศัพท์เพื่อค้นคืนข้ามภาษาไทย – อังกฤษ โดยประยูทธ สุวรรณวิสารท ซึ่งผู้วิจัยได้มีการปรับปรุงอัลกอริทึมเพื่อให้สามารถใช้กับคำไทยได้ดังนี้

- ใช้ตัวอักษรภาษาอังกฤษที่ตรงกับเสียงนั้นเป็นตัวแรกของรหัสชาวเน็ตเวิร์ก หากตัวอักษรตัวแรกของคำเป็นเสียงสระ (เสียง อ) จะใช้ตัวอักษร 'AEIOU' แทนสระนั้น (แสดงในตารางที่ 3.4)
- ใช้รหัส 52 สำหรับตัวอักษร 'ง' เมื่ออยู่ส่วนอื่นของคำ และใช้ตัวอักษร 'NG' เมื่อเป็นตัวแรกของคำ (แสดงในตารางที่ 3.4 และ ตารางที่ 3.5)

ตารางที่ 3.4 แสดงการแทนรหัสตัวอักษรตัวแรกสำหรับการเข้ารหัสตัวอักษรไทย

Thai alphabet	Code
อ	A E I O U
บ	B
ฝ ฟ	F
ป ผ พ ภ	P
ก ข ฅ ค ฅ ฌ	K
จ ฌ ช ฌ	C
ซ ศ ษ ส	S
ฎ ฑ	D
ฎ ฑ ฐ ฑ ฒ ถ ฑ ฐ	T
ล ฬ	L
ม	M

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

- ใช้รหัส ASCII เป็นรหัสสำหรับสระ (แสดงในตารางที่ 3.6)
- ไม่มีการเข้ารหัสสำหรับวรรณยุกต์
- แปลงสระในตัวอักษรไทย ให้อยู่ในรูปของการเขียนในภาษาอังกฤษ เช่น ‘-า’ -> ‘ao’, ‘ไ’ -> ‘ai’. (แสดงในตารางที่ 3.6)

ตารางที่ 3.6 แสดงการเข้ารหัสสำหรับสระ

ตัวอักษรไทย	ตัวอักษรโรมัน	รหัสที่ได้
อะ, อา, อึ	a	97
อัวะ, อัว	ua	11797
อำ	am	9764
อิ, อึ	i	105
อุ, อู	u	117
อึ, อึ	ue	11769
เอย	oei	111101105
เอะ, เอ็, เอ	e	101
เออะ, เออ, เอ็	oe	111101
เอา	ao	97111
เอ็ยะ, เอ็ย	ia	10597
เอ็อะ, เอ็อ	uea	11710197
แอะ, แ็, แอ	ae	97101
โอะ, โอ	o	111
เอาะ, อ็, ออ	o	111
ไอ, ไอ, ไอย	ai	97105
เอ็ย	ueai	11710197105
-วย	uai	11797105
ิว	io	105111
เ็ว, เ-ว	eo	101111
เอ็ยว	iao	10597111

- ตัวสะกด ในภาษาไทย เช่น จ ฌ ฌ ฌ ศ ษ ส ญ ร ล พ เมื่อนำตัวอักษรเหล่านี้มาเขียนแบบ

เอกสารนี้เป็นเอกสารลิขสิทธิ์ของมหาวิทยาลัยเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง
 ถ่ายเสียงภาษาไทยโดยใช้ตัวอักษรโรมันตามเกณฑ์การถอดอักษรของราชบัณฑิตยสถาน
 ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

จะพบว่าตัวสะกดที่ได้ จะถูกแปลงเป็นคนละตัวกับอักษรไทย เช่น จ.จาน เมื่อใช้เป็นพยัญชนะต้น จะใช้ 'ch' แต่เมื่อเป็นตัวสะกดจะใช้ตัวอักษร 't' ดังนั้น รหัส ชวาน์เด็กซ์ ที่แปลงได้ในฐานข้อมูลก็จะไม่ตรงกับคำไทยที่ใช้ตัวอักษรโรมัน ที่แปลงได้จากโปรแกรม

ตารางที่ 3.7 ตารางแสดงตัวอักษรที่ทำให้เกิดข้อผิดพลาดของโปรแกรม

รูปพยัญชนะไทย	พยัญชนะต้น	ตัวสะกด	รหัสที่ถอดได้	
			พยัญชนะต้น	ตัวสะกด
จ ฉ ช ฌ	ch	t	2	3
ซ ศ ษ ส ทร (ออกเสียง ซ)	s	t	2	3
ญ	y	n	9	5
ร	r	n	6	5
ล พ	l	n	4	5

- กรณีที่พบการอ่านออกเสียงเชื่อมระหว่างพยางค์ เช่น พัฒนา มีการออกเสียง /tha/ เชื่อมระหว่างพยางค์ โปรแกรมจะต้องเพิ่มตัวอักษรเสียงเดียวกันเพิ่มขึ้นมาในคำ เพื่อให้ตรงกับเสียงที่อ่าน
- คำที่มีตัวการันต์ประกอบ ให้นำไปตัดตัวการันต์ออกก่อนที่จะนำไปเข้ารหัสใช้แก้ไขปัญหากรณีที่พบตัวการันต์ ตัวอักษรก่อนหน้าจะต้องไม่ออกเสียงด้วย และในกรณีที่ มีตัวการันต์ไม่ออกเสียง 2 ตัว เช่น จันทร์ ก็จะต้องไม่นำตำแหน่ง ทร ไปเข้ารหัสด้วย
- หากมีอักษรนำเสียงสนิท เช่น อย-हन-หล-หม-หญ-หง-หฺร-หฺว-หฺย จะต้องไม่มีการเข้ารหัสตัวอักษร ห หรือ อ จะตัดตัวอักษร อ หรือ ห ที่ แล้วให้อักษรตัวถัดไปมาเป็นรหัสแทน กรณีคำที่พบเป็นอักษรนำเสียงสนิท
- กรณีที่มีสระ ็ (อำ) อยู่ในคำ ให้เปลี่ยนสระ ็ (อำ) เป็น ั (อัม) โดยแทนรหัสเป็น 5
- กรณีที่พบ รร (ร หัน) เช่น สรรค์ สรรพ มหัสจรรย์ จะแก้ไขโดยแยกเป็นอีก 2 กรณีคือ
 1. ร หัน แต่ไม่มีตัวสะกด จะเปลี่ยนให้เป็น ัน (อัน)
 2. ร หัน มีตัวสะกด จะเปลี่ยน ร หัน ให้เป็น ั (สระ อะ ลรูป)

เมื่อเขียน โปรแกรมการเข้ารหัสแล้ว จากนั้นจะดึงคำไทยจากฐานข้อมูลมาเข้ารหัสทีละตัวแล้ว

จัดเก็บลงใน Microsoft Access ในสัปดาห์ที่ 2 แสดงเป็นตัวอย่าง ได้ดังนี้
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 3.8 ตัวอย่างข้อมูลที่ผ่านการเข้ารหัสชาวน์เด็กซ์

ที่	T10500
การ	K97500
เป็น	P10150
ใน	N97105
จะ	C97000
ของ	K11152

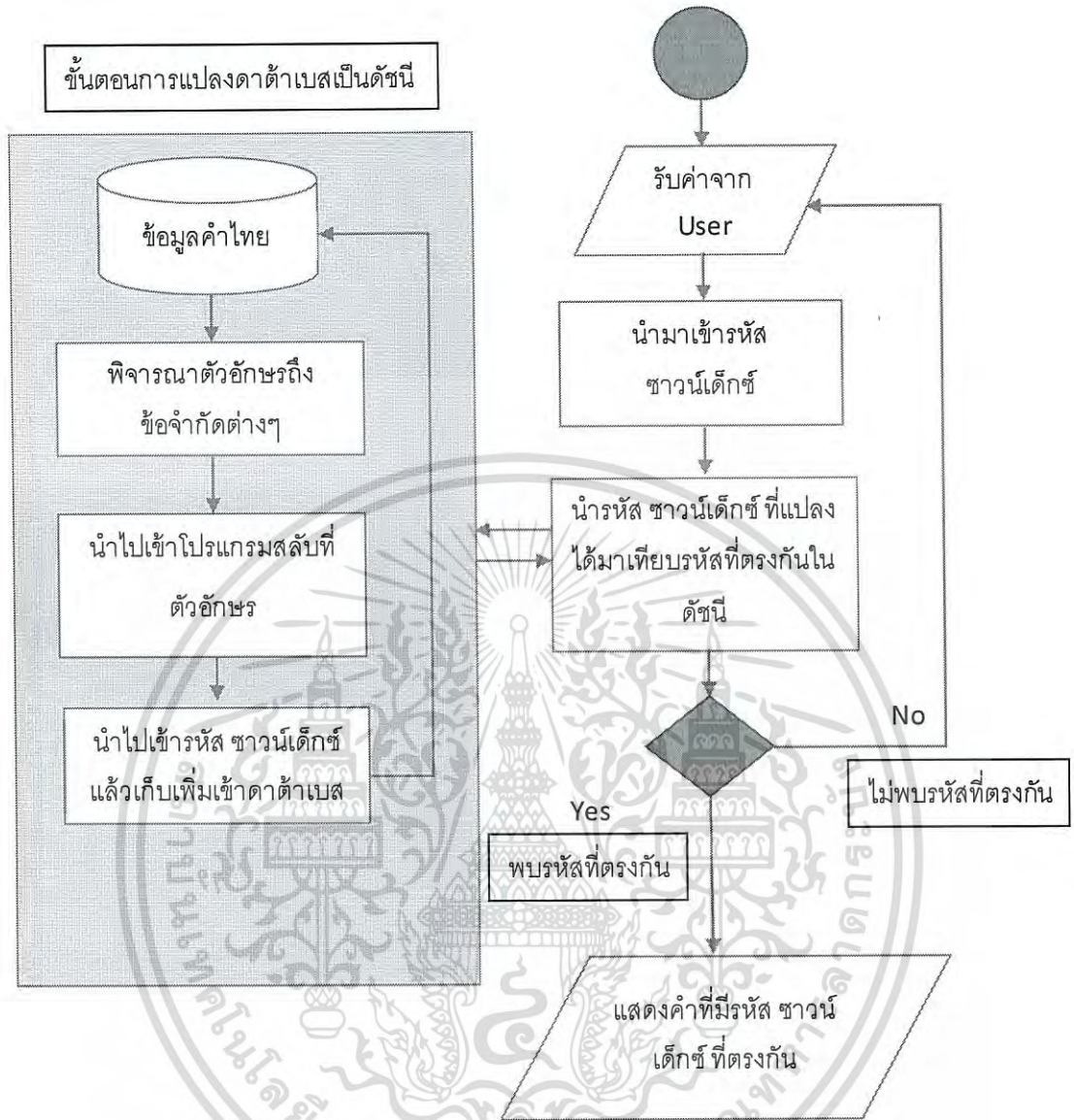
3.3 การเขียนโปรแกรมเข้ารหัสชาวน์เด็กซ์สำหรับคำภาษาอังกฤษ

หลังจากเตรียมข้อมูลรหัสชาวน์เด็กซ์คำไทยเพื่อใช้สำหรับการเปรียบเทียบแล้ว จากนั้นจะเป็นขั้นตอนการเขียนโปรแกรมการเข้ารหัสชาวน์เด็กซ์สำหรับคำภาษาอังกฤษ โดยหลักๆแล้วจะใช้เกณฑ์การเข้ารหัสโดยใช้ตารางเดียวกันกับการเข้ารหัสชาวน์เด็กซ์ภาษาไทย โดยขั้นแรกจะให้คำที่ป้อนเข้ามาในระบบเปลี่ยนให้เป็นตัวพิมพ์ใหญ่ทั้งหมด จากนั้นจึงนำมาเข้ารหัสโดยใช้ตารางที่ 3.7 ตัวอย่างคำภาษาอังกฤษที่ผ่านการเข้ารหัสดังนี้

ตารางที่ 3.9 ตัวอย่างการเข้ารหัสคำภาษาอังกฤษ

ที่	THI	T10500
การ	KAN	K97500
เป็น	PEN	P10150
ใน	NAI	N97105
จะ	CHA	C97000
ของ	KHONG	K11152
มี	MI	M41050
และ	LAE	L97101
ไม่	MAI	M49710
ได้	DAI	D97105

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 3.2 แผนภาพแสดงการทำงาน โดยรวมของระบบ

ซึ่งผู้วิจัยได้มีการปรับปรุงขั้นตอนวิธีการเข้ารหัสฆาวันเด็กซ์คำภาษาอังกฤษดังนี้

- หากตัวอักษรแรกของคำเป็นสระ AEIOU จะแทนรหัสด้วย 'AEIOU' เช่นเดียวกับขั้นตอนการเข้ารหัสในอักษรไทย แต่ตัวอักษรนั้นจะขึ้นอยู่กับตัวอักษรตัวแรกของสระที่ได้จากการแปลงเป็นตัวอักษรภาษาอังกฤษ (แสดงในตารางที่ 3.6)
- ไม่มีการเข้ารหัสสำหรับตัวอักษร 'H' ที่อยู่หลังตัวอักษร C, K, T, P
- คำภาษาอังกฤษที่ขึ้นต้นด้วย NG จะให้รหัสฆาวันเด็กซ์ตัวแรกเป็น NG เพื่อเป็นการกำหนดว่าคำที่ต้องค้นนั้นเป็น ตัวอักษร ง ไม่ใช่ตัวอักษร น หรือ ณ และเมื่ออยู่ส่วนอื่นของคำให้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับบุคลากรใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 แทนรหัสเป็น 52 โดยที่ 5 แทนตัวอักษร N และ 2 แทนตัวอักษร G (แสดงในตารางที่ 3.11)
 ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมีเหตุดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

- สำหรับสระ AEIOU ที่ไม่ใช่ตัวแรกของคำ จะแทนด้วยรหัส ASCII เช่นเดียวกับขั้นตอนการเข้ารหัสในอักษรไทย (ดังแสดงในตารางที่ 3.6)
- รหัสที่ได้จะมีความยาวเท่ากับ 9 ตัวอักษร โดยที่ตัวแรกจะเป็นตัวอักษร และอีก 6 ตัวที่เหลือจะเป็นตัวเลขเช่นเดียวกับขั้นตอนการเข้ารหัสในอักษรไทย

จากนั้นเมื่อได้รหัสแล้วจะนำรหัสที่ได้มาเทียบกับฐานข้อมูลดังกล่าวในขั้นตอนต่อไป

ตารางที่ 3.10 แสดงการแทนรหัสตัวอักษรตัวแรกในขั้นตอนการแปลงรหัสสำหรับตัวอักษรโรมัน

Roman alphabet	Code
A E I O U	A E I O U
B	B
F V	F
PH P	P
K KH CG	K
C CH J	C
S	S
D	D
T	T
L	L
M	M
N	N
R	R
NG	NG
H	H
Y	Y
W	W

หมายเหตุ: A E I O U จะขึ้นอยู่กับตัวแรกของ สระ
ที่ได้จากการแปลงเป็นตัวอักษรภาษาอังกฤษ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 3.11 แสดงการแทนรหัสสำหรับตัวอักษรตัวที่ไม่ใช่ตัวแรกของคำในขั้นตอนการเข้ารหัสตัวอักษรโรมัน

Roman alphabet	Thai Sound	Code	Note
B F P V PH	บ ฝ ฟ ฝ ผ พ ภ	1	
K C G J KH CH S	ก ข ฃ ค ฅ ฉ จ ฉ ช ฌ ซ ศ ษ ส	2	
	จ ฉ ช ฌ ซ ศ ษ ส	3	#1
D T	ฎ ฏ ฏ ฐ ฑ ฒ ถ ฑ ฐ	3	
L	ล พ	4	
M	ม	@	
N	ณ น	5	
R	ร	6	
	ร ฤ ล พ	5	#1
H	ห ฮ	7	
W	ว	8	
Y	ญ ย	9	
NG	ง	52	

#1 : for final consonant

3.4 การกำหนดรหัสขานี้เด็กซ์ สำหรับสระ ในภาษาไทย

ในการศึกษานี้ผู้วิจัยได้มีวิธีการเข้ารหัสสระเพิ่มขึ้นมาแตกต่างจากต้นฉบับเนื่องจากผลลัพธ์ที่ค้นหาได้มีจำนวนที่เยอะเกินไป และคำศัพท์ส่วนใหญ่มีความเกี่ยวข้องน้อย หรือไม่ตรงกับที่ผู้ใช้ต้องการ เป็นผลมาจากเราเข้ารหัสขานี้เด็กซ์เฉพาะตัวอักษรเท่านั้น เราจึงเพิ่มการเข้ารหัสให้กับสระด้วย เพื่อให้การค้นหาคำมีขอบเขตที่แคบและ ใกล้เคียงกับคำที่ต้องการมากขึ้น โดยจะใช้รหัส ASCII ของตัวอักษรสระในภาษาอังกฤษ มาเป็นรหัสขานี้เด็กซ์แทน

ยกตัวอย่างเช่น เข้ารหัสขานี้เด็กซ์คำว่า “และ”

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ขั้นตอนแรก ก่อนนำมาเข้ารหัสตัวเด็กซ์จะต้องทำการสลับตำแหน่งเมื่อมีสระประสม หรือสระที่มีตำแหน่งนำข้างหน้าพยัญชนะต้น จากคำว่า ”และ” จะถูกสลับเป็น “ลแะ”

ขั้นตอนที่สอง แปลงสระนั้นๆให้อยู่ในรูปของตัวอักษรภาษาอังกฤษ ซึ่งในตัวอย่างนี้สระ คือ “แะ” (แอะ) จะถูกแปลงเป็น “ae” (ดูได้จากตารางที่3.6) เมื่อแปลงเสร็จจะได้เป็น “lae”

ขั้นตอนสุดท้าย นำผลลัพธ์ที่ได้จากขั้นตอนที่สอง มาแปลงเป็นรหัสตัวเด็กซ์ จาก ตัวอย่างในขั้นตอนที่สอง นำ “lae” มาแปลงจะได้เป็น L97101000 (L คือตัว ‘ล’ จากตารางที่ 3.4 เลขรหัส 97101 คือ “ae” จากตารางที่ 3.6 ส่วนเลขศูนย์อีกสามตัวท้ายถูกเพิ่มให้รหัสครบ 9 ตัว)



รูปที่ 3.3 แสดงตัวอย่างขั้นตอนการเข้ารหัสตัวเด็กซ์ของคำว่า “และ”

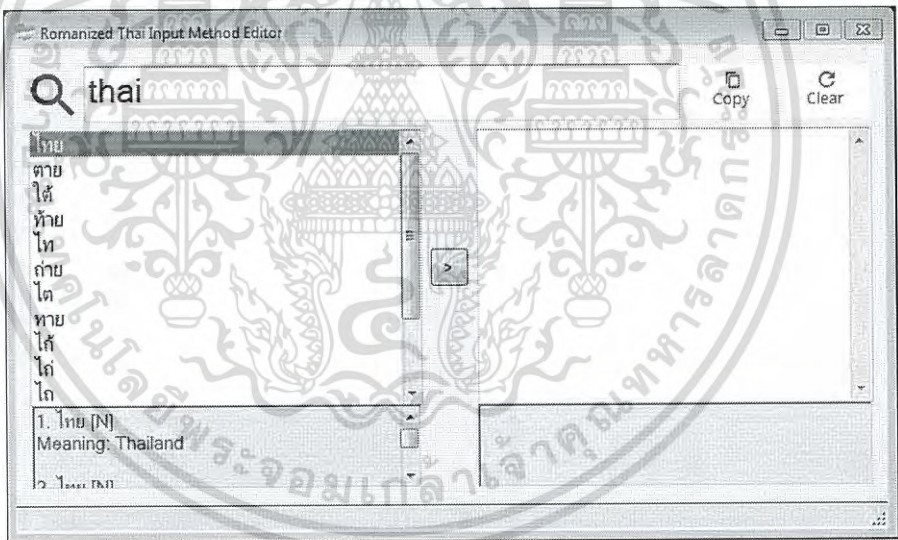
บทที่ 4

ผลการดำเนินงาน

4.1 การเปรียบเทียบความสอดคล้องของรหัสชาวน์เด็กซ์จากอักษรภาษาอังกฤษและจากอักษรภาษาไทย

เมื่อได้ข้อมูลรหัสชาวน์เด็กซ์ที่ใช้เป็นตัวกลางในการถอดอักษรแล้ว ก็จะนำรหัสชาวน์เด็กซ์ภาษาอังกฤษที่ได้มาเทียบกับข้อมูลชาวน์เด็กซ์ภาษาไทยที่ได้จัดเก็บไว้ในขั้นตอนที่ 3.2 วิธีการเปรียบเทียบข้อมูลจะใช้คำสั่ง SQL ในการเปรียบเทียบว่า ให้เลือกคำไทยในสดมภ์ที่ 1 เมื่อรหัสชาวน์เด็กซ์ภาษาอังกฤษตรงกับรหัสชาวน์เด็กซ์ในสดมภ์ที่ 2 จากนั้นจะนำไปแสดงผลลงในโปรแกรมออกมาเป็นรายการคำที่พบเจอ ดังรูปที่ 3.3

เมื่อพิมพ์คำว่า “thai” ลงในช่อง INPUT ซึ่งคำที่ต้องการคือคำว่า “ไทย” ระบบจะคืนค่าคำไทยที่มีรหัสชาวน์เด็กซ์ตรงกับคำว่า “thai” ขึ้นมา



รูปที่ 4.1 ตัวอย่างการแสดงผลคำไทยที่พบเจอ โดยการเทียบรหัสชาวน์เด็กซ์ที่ตรงกัน

ตารางที่ 4.1 แสดงรายการคำที่มีรหัสตรงกัน เมื่อค้นหาคำว่า thai

คำไทยที่ได้	รหัสชาวน์เด็กซ์
ไทย	T9710500
ตาย	T9710500
ใต้	T9710500
ท้าย	T9710500
ไท	T9710500

ตารางที่ 4.1 (ต่อ) แสดงรายการคำที่มีรหัสตรงกัน เมื่อค้นหาคำว่า thai

คำไทยที่ได้	รหัสชวาร์นเด็กซ์
ถ่าย	T9710500
ไต่	T9710500
ทาย	T9710500
ไต	T9710500
ไต่	T9710500
ไต่	T9710500
ไต่	T9710500
ไต่	T9710500
ไต่	T9710500
ไต่	T9710500
ไต่	T9710500
ไต่	T9710500

จะเห็นได้ว่าข้อมูลทั้งหมดที่พบนั้น มีรูปแบบการเขียนที่สอดคล้องกับคำว่า Thai ทั้งสิ้น ส่วนคำที่มีหลายพยางค์ระบบก็นำข้อมูลนั้นออกมาให้เพื่อช่วยให้ผู้ใช้สามารถใช้งานได้ง่ายขึ้น โดยไม่จำเป็นต้องพิมพ์ให้ครบทั้งคำ

4.2 การประเมินผลความถูกต้องแม่นยำของโปรแกรมการถอดอักษร

ในการทดสอบว่าโปรแกรมการถอดอักษรที่พัฒนาขึ้นมาทำงานได้ดีเพียงใด ก็จะต้องมีการตรวจสอบความถูกต้องของคำที่ถอดอักษรได้จาก โปรแกรมเปรียบเทียบกับคำที่ถอดอักษรโดยราชบัณฑิตยสถาน ซึ่งผู้วิจัยได้มีการเตรียมข้อมูลเพื่อใช้ทดสอบนี้เป็นข้อมูลคำไทยที่ถอดอักษรเป็นโรมัน โดยใช้เกณฑ์การถอดอักษรของราชบัณฑิตยสถาน[3] เป็นข้อมูลจำนวน 4,808 คู่คำ โดยจะเปรียบเทียบระหว่างสดมภ์ของ SOUNDEX_CODE (รหัสที่ได้เมื่อผ่านกระบวนการเข้ารหัส) กับ COMPARE_CODE (รหัสที่ควรได้) แล้วนำค่า Boolean ที่ได้มาจัดเก็บในสดมภ์ RESULT ดังแสดงตารางใน Microsoft Excel ดังนี้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 4.2 แสดงการเปรียบเทียบรหัสชวาว์เด็กซ์คำภาษาอังกฤษที่ได้กับรหัสชวาว์เด็กซ์คำ
ภาษาไทย

THAI_WORD	SOUNDEX_CODE	COMPARE_CODE	RESULT
กักขัง	K9729752	K9729752	TRUE
กั้ง	K9752000	K9752000	TRUE
กั้งวด	K9752811	K9752811	TRUE
กัก	K9730000	K9730000	TRUE
กั้น	K9750000	K9750000	TRUE
กั้น	K9750000	K9750000	TRUE
กั้นเอง	K9751015	K9751015	TRUE
กั้นและกั้น	K9754971	K9754971	TRUE
กั้นยายน	K9759979	K9759979	TRUE
กั๊	K9710000	K9710000	TRUE
กั๊ข้าว	K9712971	K9712971	TRUE
กั๊พูชา	K9764111	K9764111	TRUE
กา	K9700000	K9700000	TRUE

ขั้นตอนทั้งหมดที่กล่าวไปเป็นวิธีการดำเนินการวิจัยโดยละเอียด มีการทดสอบและ
เปรียบเทียบข้อมูลหลายๆ ครั้ง เนื่องจากข้อมูลมีปริมาณมาก บางครั้ง โปรแกรมก็คำนวณออกมา
ผิดพลาด ดังนั้นจึงต้องทดสอบให้ได้ค่าที่คงที่ที่สุด

จากนั้นจะทำการนับข้อมูลที่มีค่า RESULT เป็น TRUE นั้น หมายถึงข้อมูลที่เปรียบเทียบกัน
มีความถูกต้อง ได้รหัสที่ตรงกัน จากนั้นก็จะนำมาหาค่าเฉลี่ยของข้อมูลทั้งหมดที่ได้จัดเก็บว่ามีความ
ถูกต้องคิดเป็นกี่เปอร์เซ็นต์ และผลความถูกต้องแม่นยำของ โปรแกรมที่พัฒนาในบทต่อไปจะ

กล่าวถึงผลการถอดอักษร โดยใช้โปรแกรมคอมพิวเตอร์กับตัวกลางในการถอดอักษร

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

4.3 การแก้ไขอัลกอริทึมเพื่อเพิ่มความแม่นยำในการทำงาน

หลังจากที่เราได้ทดสอบการทำงานของโปรแกรม จะพบว่าผลจากการค้นคำมีจำนวนมาก และคำบางคำที่ได้ค่อนข้างที่จะไม่มีความสอดคล้องกับคำที่ป้อนลงไป ในโปรแกรมเช่น คำที่มีตัวสะกดเป็น ม จะแสดงพร้อมกับคำที่มีตัวสะกดเป็นตว น เนื่องจากพบว่ารหัสชาวน์เด็กซ์ที่นำมาใช้ กำหนดให้ พยัญชนะเสียง ม และ พยัญชนะเสียง น มีการเข้ารหัสเป็นรหัสตัวเดียวกันคือ 5 ดังนั้นจึงกำหนดรหัสใหม่ให้พยัญชนะเสียง ม มีรหัสเป็น @ เพื่อลดความผิดพลาดในการค้นคำ

หรือกรณีที่ผู้ใช้ไม่ได้ใช้หลักการถอดอักษรในการพิมพ์คำไทยด้วยตัวอักษรภาษาอังกฤษ เช่น ใช้ตัวอักษร J แทน จ ซึ่งตามมาตรฐานจะต้องใช้ CH ผู้วิจัยได้มีการแก้ไขอัลกอริทึมให้รองรับความเคยชินของผู้ใช้งานมากขึ้น ผู้ใช้สามารถใช้ตัวอักษรที่ไม่ตรงตามมาตรฐานได้ แต่ทำได้เพียงบางตัวเท่านั้น เพราะยังมีข้อจำกัดอยู่บางประการ

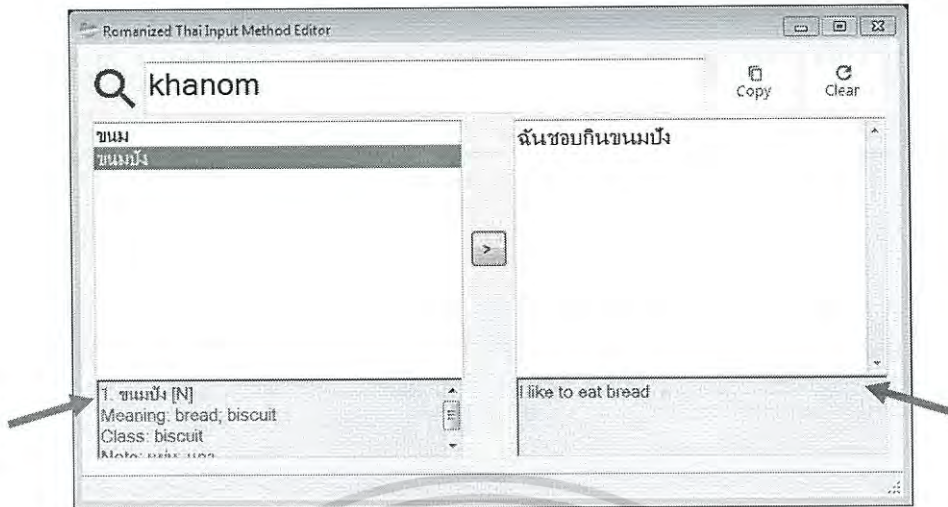


รูปที่ 4.2 ตัวอย่างการแสดงผลข้อมูลด้วยการใช้ตัวอักษรที่ไม่ได้ระบุในกฎการถอดอักษรโดยราชบัณฑิตยสถาน

4.4 กำกับความหมายและประเภทของคำศัพท์ในคลังข้อมูล

เพื่อให้ผู้ใช้สามารถเลือกใช้คำได้อย่างถูกต้อง เนื่องจากคำในภาษาไทยมีคำพ้องรูปและคำพ้องเสียงมากมาย จึงต้องมีการกำกับคำศัพท์ด้วยความหมายและประเภทของคำเพิ่มเติม โดยได้ใช้ Google api มาเป็นตัวช่วยในการแปลความหมายของคำ ความหมายที่ได้จะใช้เป็นภาษาอังกฤษซึ่งสามารถอ่านเข้าใจได้ในระดับสากล ทั้งชาวไทยและชาวต่างชาติ

เอกสารนี้เป็นเอกสารที่สงวนลิขสิทธิ์สำหรับการใช้งานเพื่อการศึกษาเท่านั้น เมื่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4.3 ตัวอย่างการแสดงความหมายประกอบของคำที่ค้นหา

4.5 สรุปผลจากการศึกษาแนวทางที่เหมาะสมในการพัฒนาระบบการถอดอักษรอังกฤษเป็นไทยแบบถ่ายเสียง

ผู้จัดทำได้ทำการทดสอบโดยใช้คำไทย – อังกฤษ ทั้งหมด 4,807 คู่คำหลังจากที่ได้มีการปรับปรุงอัลกอริทึม ทำให้ค่าความถูกต้องของอัลกอริทึมรหัสชวอน์เด็กซ์ของคำไทยและอังกฤษตรงกันเพิ่มขึ้น

จากผลการทดสอบพบว่า มีบางคำที่ระบบไม่สามารถดึงข้อมูลออกมาแสดงผลได้ เนื่องจากบางคำอัลกอริทึมไม่สามารถแยกเสียงของตัวอักษรนั้นออกมาได้อย่างถูกต้อง โดยตัวอักษรที่เป็นปัญหามากที่สุดคือ อักษร ข และ อ เนื่องจากความหลากหลายในการอ่านออกเสียงสองตัวนี้ในคำต่างๆ ในขั้นตอนของการแปลงสระให้เป็นตัวอักษรโรมัน เช่น คำว่า ทายาท ระบบแยกคำออกมาได้เป็น ทาย-าท ซึ่งไม่ถูกต้อง เมื่อนำคำไปสืบค้นจึงไม่ปรากฏคำนี้ขึ้น

ประเด็นสำคัญที่ทำให้คำระบบไม่สามารถค้นคำที่ต้องการออกมาแสดงผลได้คือ ขั้นตอนของการเรียงลำดับตัวอักษรใหม่ให้สอดคล้องกับลำดับของการเขียนในรูปของคำทับศัพท์ เพราะถ้าหากมีการเรียงลำดับตัวอักษรผิด รหัสที่ได้ก็จะผิดไปด้วยเพราะระบบจะต้องนำเอาลำดับของตัวอักษรไปพิจารณาว่าสอดคล้องกับพยัญชนะ หรือสระตัวใดในภาษาไทย และอีกประเด็นหนึ่งที่สำคัญไม่แพ้กัน คือความหลากหลายของการออกเสียง โดยปกติแล้วชาวต่างชาติมักจะถอดเสียงคำ

ไทยโดยใช้เสียงของภาษาดั้งเดิมของตนเอง ซึ่งจะแตกต่างกับภาษาไทย เช่น คำว่า เป็น ชาวต่างชาติ เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ส่วนใหญ่มักถอดเสียงออกมาเป็น ben , bin หรือเสียง C , G ชาวต่างชาติส่วนใหญ่มักถอดเสียงเป็นเสียง ก หรือ จ ทำให้ระบบไม่สามารถค้นเจอคำที่การออกอย่างหลีกเลี่ยงไม่ได้

จากปัญหาที่กล่าวมาข้างต้น ผู้วิจัยจึงได้เพิ่มส่วนของการแสดงความหมายของคำขึ้นมา หากผู้ใช้ถอดเสียงผิดไปจากกฎการถอดอักษรหรือใช้ความคุ้นชินของเสียงมาเป็นหลัก เพื่อช่วยให้เกิดความแม่นยำ และคำที่แสดงผลออกมานั้นได้เรียงลำดับตามความถี่ในการใช้งานในปัจจุบัน



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 5

สรุปผลโครงการ

5.1 สรุปผลการทำโครงการและข้อเสนอแนะ

โครงการนี้เกี่ยวกับการสร้างวิธีขั้นตอนในการพิมพ์ภาษาไทยด้วยตัวอักษรภาษาอังกฤษเพื่อใช้เป็นเครื่องมือในการศึกษาภาษาไทยและเป็นการช่วยในการพิมพ์ภาษาไทยด้วยตัวอักษรภาษาอังกฤษโดยคำนึงจากการอ่านออกเสียงของตัวอักษรและการประสมกันเป็นพยางค์หรือคำนั้นๆ โดยได้ใช้หลักการจับคู่ตัวอักษรในการแยกสระกับพยัญชนะ และใช้ชวัญเด็กซ์อัลกอริทึมในการเข้ารหัสเพื่อจำแนกคำด้วยคำอ่านของคำศัพท์

และด้วยข้อจำกัดในการใช้ฐานข้อมูลคำศัพท์ ซึ่งเป็นวิธีที่ไม่เหมาะสมสำหรับการทำ Romanization เพราะอาจจะมีคำใหม่ๆเกิดขึ้นได้ในอนาคตอยู่ตลอดเวลา จึงไม่เหมาะสมที่จะเผยแพร่โปรแกรมที่พัฒนาเพื่อใช้เป็นเครื่องมือการเขียนคำทับศัพท์สำหรับผู้ใช้โดยทั่วไป ผู้วิจัยจึงคาดหวังไว้ว่าต้องการให้รายงานฉบับนี้เป็นข้อมูลอ้างอิง หรือเป็น โปรแกรมต้นแบบเพื่อนำไปศึกษาและพัฒนาให้มีประสิทธิภาพต่อไปในอนาคต

5.2 ปัญหาที่พบระหว่างดำเนินโครงการ

1. หลักการถอดอักษรไทยเป็นอักษรโรมันแบบถ่ายเสียงของราชบัณฑิตยสถานยังมีความไม่ชัดเจนในเรื่องของรายละเอียดในภาษาไทย โดยที่พบคือ
 - สระเสียงสั้นเสียงยาว ใช้ตัวอักษรเดียวกัน เช่น อะ กับ อา ใช้ตัว 'A' เหมือนกัน
 - ตัวอักษรภาษาไทย 1 เสียงบางตัวอักษรมีหลายรูป อีกทั้งมีตัวที่เสียงไม่เหมือนกันใช้ตัวเดียวกัน เช่น 'ส ศ ษ' กับ 'ซ' แทนด้วย 'S'โดยปัญหาเหล่านี้ทำให้การพิมพ์ไม่สามารถเจาะจงคำ หรือทำให้ผลลัพธ์ที่ได้มีขอบเขตกว้างเกินไป ส่งผลให้ผู้ผู้ใช้
2. พยัญชนะต้น หรือคำควบกล้ำ ควบกล้ำทั้งแท้และควบกล้ำไม่แท้ในภาษาไทยทำให้ระบบถอดตัวอักษรทำงานได้ไม่ดันทัก เช่น จรวด จรด ทราบ
3. ตัวสะกดบางตัวทำให้ระบบไม่สามารถถอดออกมาให้ตรงกับคำนั้นได้ เช่น มาตรฐาน สามารถ การใช้ตัวสะกด ตร หรือ รด เป็นการยืมคำมาจากภาษาอื่น
4. อัลกอริทึมชวัญเด็กซ์ที่นำมาใช้ยังมีข้อบกพร่องอยู่มาก จึงต้องมีการปรับเปลี่ยนเพื่อให้เกิดความถูกต้องแม่นยำมากขึ้น

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

5.3 ประโยชน์ที่ได้รับจากการดำเนินโครงการ

จากการจัดทำโครงการ 1 ผู้จัดทำได้รับประโยชน์ดังนี้

1. ได้ศึกษาทฤษฎีความรู้ถึงความคล้ายคลึงกันของภาษาต่างๆ
2. ฝึกการแก้ไขปัญหาการลดความบกพร่องของระบบ

5.4 แนวทางการพัฒนาในอนาคต

- พัฒนาให้ระบบสามารถถอดอักษรได้แม่นยำมากขึ้นที่ระดับ 100%
- พัฒนาขั้นตอนวิธีที่ครอบคลุมในส่วนของคำอังกฤษทับศัพท์ภาษาไทย และภาษาไทยทับศัพท์อังกฤษให้เป็นขั้นตอนวิธีเดียว เพื่อให้เกิดความสะดวกในการใช้งาน
- มีระบบแนะนำความน่าจะเป็นของคำต่อไป โดยใช้หลัก n-gram มาพิจารณาความเป็นไปได้ของคำถัดไป
- พัฒนาให้รองรับการใช้งานบนมือถือ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บรรณานุกรม

- [1] ประยูท สุวรรณวิสารท. การเข้ารหัสคำทับศัพท์เพื่อการค้นคืนข้ามภาษาไทย-อังกฤษ. วิทยานิพนธ์ปริญญาโทมหาบัณฑิต ภาควิชาวิศวกรรมคอมพิวเตอร์ คณะวิทยาศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย, 2541.
- [2] วิโรจน์ อรุณมานะกุล. อักษรวิธีไทยและการถอดอักษรระหว่างภาษาไทยและภาษาอังกฤษ. พิมพ์ครั้งที่ 1 กรุงเทพมหานคร : โครงการเผยแพร่ผลงานวิชาการ คณะอักษรศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย, 2551
- [3] T. Karoonboonyanan, V. Sornlertlamvanich and S. Meknavin. **A Thai Soundex System for Spelling Correction. Proceedings of the Natural Language Processing Pacific Rim Symposium (1997) : 633-636.**
- [4] วรณี อุดมพานิชย์. การใช้หลักคำพ้องเสียง เพื่อค้นหาชุดอักษรภาษาไทยที่ออกเสียงเหมือนกัน. วิทยานิพนธ์ปริญญาโทมหาบัณฑิต ภาควิชาวิศวกรรมคอมพิวเตอร์ บัณฑิตวิทยาลัย จุฬาลงกรณ์มหาวิทยาลัย, 2526.
- [5] ราชบัณฑิตยสถาน. ประกาศราชบัณฑิตยสถาน เรื่อง การถอดอักษรโรมัน, 2482.
- [6] วลัยวรา ไชยฤกษ์. การพัฒนาโปรแกรมการถอดอักษรภาษาอังกฤษเป็นไทยโดยใช้คลังคำทับศัพท์ของราชบัณฑิตยสถาน. วิทยานิพนธ์ปริญญาโทมหาบัณฑิต ภาควิชาภาษาศาสตร์คณะอักษรศาสตร์ บัณฑิตวิทยาลัย จุฬาลงกรณ์มหาวิทยาลัย, 2547.
- [7] Suwanvisat, Prayut and Prasitjutrakul, Somchai. 1998. **Thai-English Cross-Language Transliterated Word Retrieval using Soundex Technique.** The National Computer Science and Engineering Conference, Kasetsart University, Bangkok, Thailand.
- [8] Suwanvisat, Prayut and Prasitjutrakul, Somchai, 1999. **Transliterated Word Encoding and Retrieval Algorithms for Thai-English Cross-Language Retrieval.** The National Computer Science and Engineering Conference, Bangkok, Thailand.

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บรรณานุกรม (ต่อ)

- [9] A.B. Grisworld. **Thought on the Romanization of Siamese.** Printed for distribution at the obsequies of Colonel Kasem Nandhakij, February 25, 1969
- [10] Vajiravudh, King. 1912. **The Romanization of Siamese Words.** In **Journal of the Siam Society**, Vol.9, Part 4, 1-10.
- [11] Travis Leithead, Takayoshi Kochi, Kenji Baheux, Hironori Bono. **“Input Method Editor PI”**[Online]. Available: <http://www.w3.org/TR/ime-api>. 2015.
- [12] “ขั้นตอนวิธีเชิงสัญลักษณ์” [Online]. Available: <https://th.wikipedia.org/wiki/ขั้นตอนวิธีเชิงสัญลักษณ์>.
- [13] "Soundex" [online]. Available: <https://en.wikipedia.org/wiki/Soundex>



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

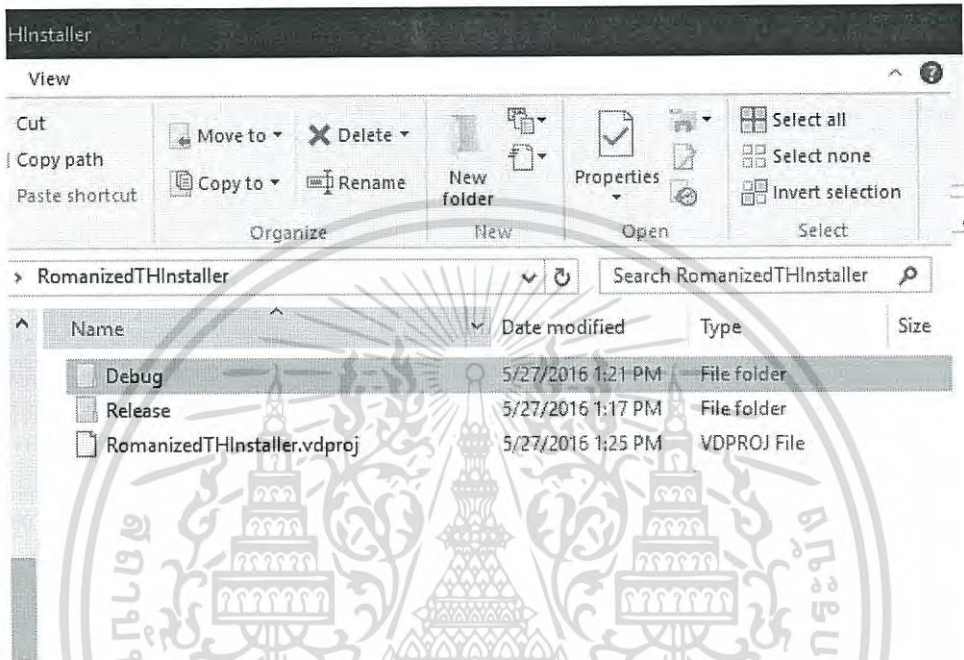


เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

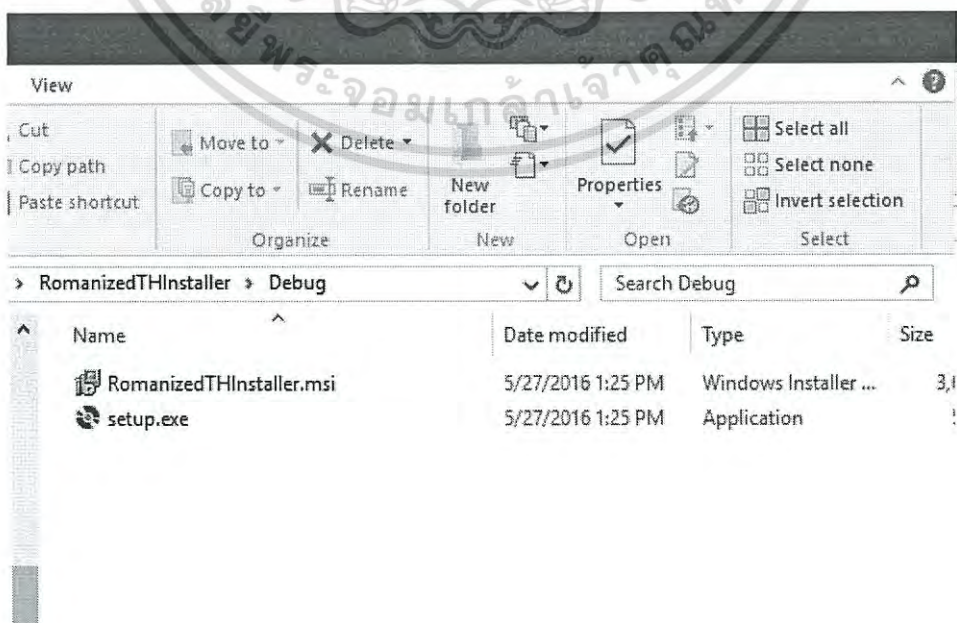
คู่มือการติดตั้งโปรแกรม

Romanized Thai Input Method Editor

1. เข้าไปใน folder ของโปรแกรมจากนั้นเข้าไปที่ folder ชื่อ Debug

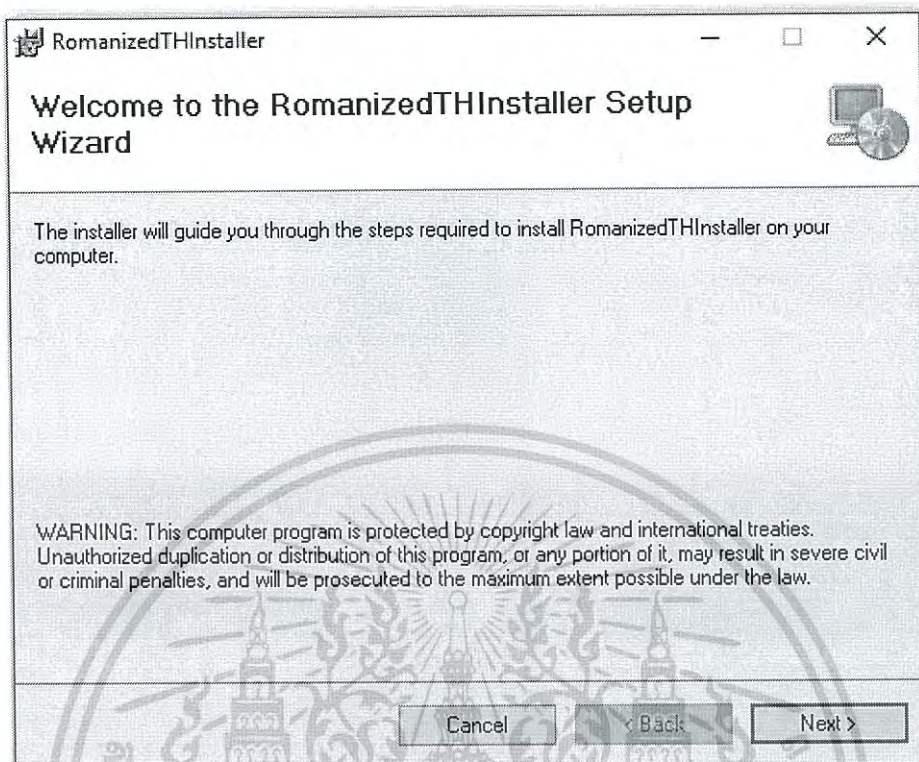


2. ดับเบิลคลิกที่ RomanizedTHInstaller.msi

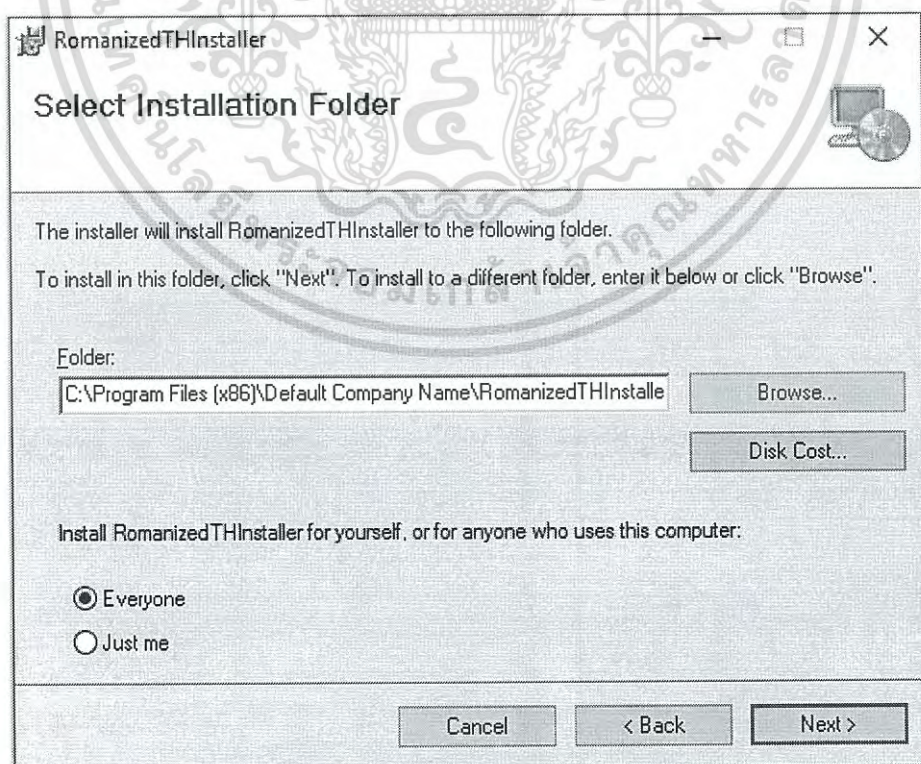


เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

3. คลิก Next

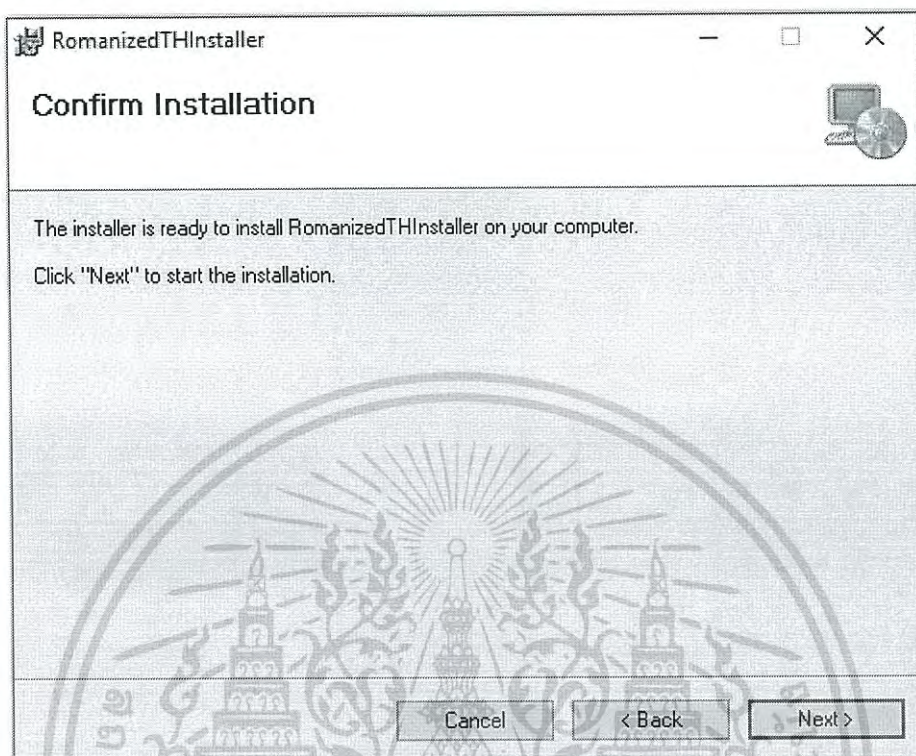


4. คลิก Browse เพื่อเลือกพื้นที่ ที่ต้องการติดตั้ง จากนั้น คลิก Next

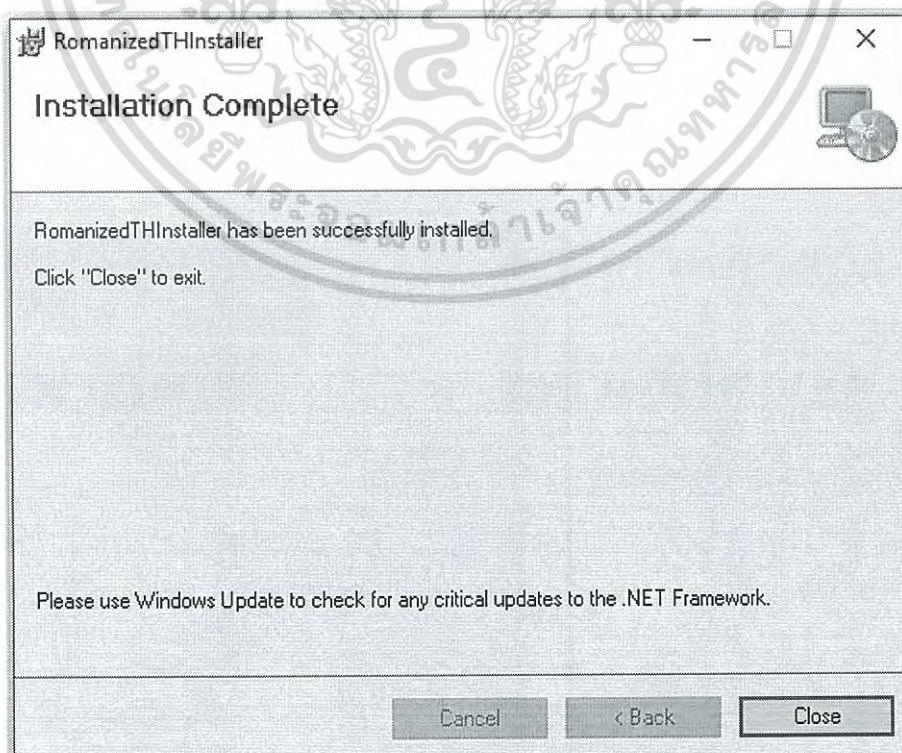


เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

5. หลังจากขั้นตอน Select Installation Folder เสร็จแล้วให้คลิก Next อีกครั้ง เพื่อเป็นการยืนยันการ Install จากนั้นรอนจนกว่าจะ ดำเนินการเสร็จ



6. หลังจากตัว Install ดำเนินการเสร็จแล้วให้คลิก Close



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

คู่มือการใช้งานโปรแกรม



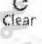
Romanized Thai Input Method Editor

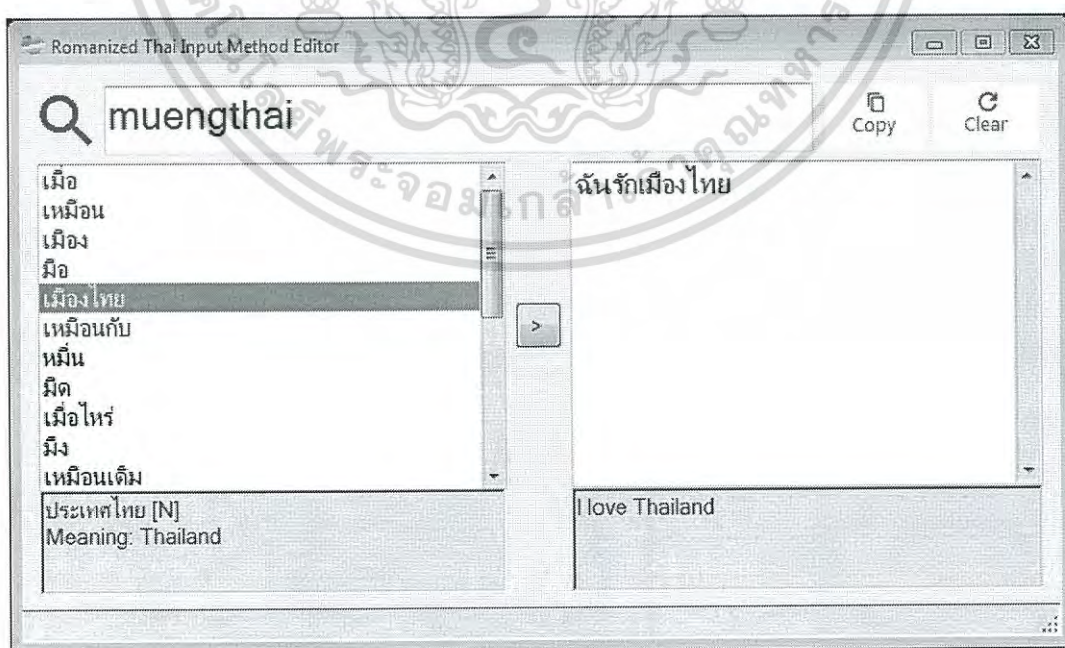


โปรแกรม Romanized Thai Input Method Editor เป็น โปรแกรมที่ช่วยแปลงคำไทยที่พิมพ์ด้วยตัวอักษรโรมัน ให้เป็นคำไทยที่พิมพ์ด้วยตัวอักษรไทย จุดมุ่งหมายของโปรแกรมคือต้องการช่วยเหลือชาวต่างชาติที่ต้องการศึกษาภาษาไทย หรือต้องการพิมพ์ภาษาไทยให้ได้อย่างถูกต้อง แหล่งที่มาของคำไทยในโปรแกรมนั้น ได้มาจากคณะอักษรศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย และศูนย์เทคโนโลยีอิเล็กทรอนิกส์และคอมพิวเตอร์แห่งชาติ (NECTEC) ซึ่งได้จากการรวบรวมข่าวสาร บทความ นวนิยายต่างๆมารวมกัน เพื่อให้มีประโยชน์ในการใช้งานในชีวิตประจำวันมากที่สุด โดยคำศัพท์ที่ค้นเจอจะเป็นคำศัพท์ที่ถูกเรียงตามความถี่ของการใช้งานในปัจจุบัน นอกจากนี้ยังช่วยเหลือให้ชาวต่างชาติที่ไม่แน่ใจในคำศัพท์ไทยด้วยการแสดงความหมายของรายการคำไทยที่ค้นเจอเป็นภาษาอังกฤษ และแสดงความหมายรวมทั้งประโยคเมื่อผู้ใช้ได้เลือกรายการคำเหล่านั้น แล้วนำมาสร้างเป็นประโยคใหม่ขึ้น

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ขั้นตอนการใช้งาน

1. เปิดโปรแกรม Romanzied Thai Input Method Editor ขึ้นมา
2. ช่อง  ใช้พิมพ์คำไทยที่ต้องการค้นหาด้วยตัวอักษรโรมัน
3. เมื่อพิมพ์คำไทยในข้อที่ 2 แล้วระบบจะทำการค้นหาคำไทยที่ตรงกันในฐานข้อมูล มาแสดงในช่องแสดงรายการคำศัพท์ พร้อมทั้งแสดงความหมายของแต่ละรายการคำศัพท์ที่ค้นเจอ
4. กดปุ่ม  เพื่อเลือกรายการคำที่ต้องการ จากช่องแสดงรายการคำศัพท์ที่ค้นเจอ หรือผู้ใช้สามารถใช้ปุ่ม enter บนแป้นพิมพ์เพื่อเลือกคำศัพท์ก็ได้ ระบบจะนำมาแสดงในช่องรายการคำศัพท์ที่เลือกทางด้านขวา และพร้อมแสดงความหมายเมื่อนำคำศัพท์ที่เลือกมารวมกันเป็นประโยค
5. ปุ่ม  จะทำการคัดลอกรายการคำที่ช่องแสดงรายการคำศัพท์ที่เลือกทางด้านขวา
6. ปุ่ม  จะทำการล้างข้อความทั้งหมดที่อยู่ในช่องป้อนข้อความ และช่องแสดงรายการคำศัพท์ที่เลือก



ตัวอย่างการใช้งาน เมื่อพิมพ์ chan rak muengthai ตามลำดับ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ระบบป้อนข้อมูลภาษาไทยด้วยตัวอักษรโรมัน

วิมลสิริ โพธิ์เกษม และ สิทธานต์ รัตนเหลียม

คณะเทคโนโลยีสารสนเทศ สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

Emails: wimonsiri.pks@gmail.com, sittan.rattanaliam@gmail.com

บทคัดย่อ

ระบบป้อนข้อมูลภาษาไทยด้วยตัวอักษร โรมันเป็นวิธีการหนึ่งที่จะช่วยให้พิมพ์ภาษาไทยด้วยการใช้ตัวอักษร โรมันได้ เนื่องจากปัจจุบันชาวต่างชาติได้เข้ามามีบทบาทในประเทศไทยมากขึ้น มีการติดต่อสื่อสารกับชาวต่างชาติเพิ่มขึ้น ราชบัณฑิตยสถานจึงได้ประกาศกฎเกณฑ์การถอดอักษรไทยด้วยตัวอักษร โรมันขึ้น เพื่อใช้เป็นมาตรฐานในการเขียนข้อความหรือคำภาษาไทยด้วยการใช้ตัวอักษรโรมันเพื่ออำนวยความสะดวกให้กับชาวต่างชาติมากขึ้น ดังที่เห็น ในแผนที่ หรือป้ายบอกทางต่างๆ ในวิทยานิพนธ์นี้ผู้จัดทำจึงได้พัฒนาระบบป้อนข้อมูลภาษาไทยด้วยตัวอักษร โรมันขึ้น โดยยึดหลักเกณฑ์การถอดอักษรไทยด้วยตัวอักษร โรมันแบบถ่ายเสียงของราชบัณฑิตยสถานมาใช้ประกอบการพัฒนา วิทยานิพนธ์ฉบับนี้ได้นำเสนอขั้นตอนวิธีการถอดอักษร โดยใช้วิธีการเข้ารหัสผ่านตัวกลางโดยใช้ ชวาว์เค็ชซ์ ซึ่งเป็นวิธีที่จะอนุญาตให้ใช้คำภาษาไทย โดยใช้ตัวอักษรโรมันได้เพื่อคั่นคั่นคำไทยที่หลักการเขียนที่ตรงกัน ในอีกภาษาหนึ่ง โดยจะมีตัวกลางที่ตรงกันสำหรับสองภาษา ขั้นตอนวิธีที่นำเสนอแบ่งออกเป็นสองส่วนคือ ขั้นตอนวิธีการเข้ารหัสตัวอักษรภาษาไทยและขั้นตอนวิธีการเข้ารหัสตัวอักษรภาษาอังกฤษ ผลการทดลองจะแสดงให้เห็นว่าขั้นตอน วิธีการเข้ารหัสภาษาไทยด้วยตัวอักษร โรมัน โดยใช้วิธีชวาว์เค็ชซ์สามารถแสดงผลออกมาได้สูงถึง 90 เปอร์เซ็นต์ นอกจากนั้น จะมีปัญหาในเรื่องความกำกวมในการออกเสียงของคำนั้นๆ และการสะกดคำที่แตกต่างกันในแต่ละบุคคลก็ส่งผลให้ระบบไม่สามารถคั่นคั่นคำไทยออกมาได้

คำสำคัญ – Romanization; Thai; Automatic Thai Romanization; Input method editor

1. บทนำ

การถอดอักษรไทย-โรมันเริ่มใช้ในช่วงทศวรรษที่ 17 คิดค้นขึ้นโดยกลุ่มมิชชันนารี [4] อย่างไรก็ตามยังไม่มีระบบที่สามารถทำให้ชาวต่างชาติออกเสียงได้อย่างถูกต้อง ในปัจจุบันได้มีการทำระบบการถอดอักษรโรมันอยู่มากรวมถึงการใช้คำทับศัพท์ ต่อมา

ISO-11940 และถูกแก้ไขพัฒนาเพิ่มเติมในปี ค.ศ. 2003 โดยระหว่างการพัฒนา นี้ทางราชบัณฑิตยสถาน ได้ตีพิมพ์หลักเกณฑ์การถอดอักษรไทยเป็น โรมันแบบถ่ายเสียงในปี ค.ศ. 1999 [2] แต่มาตรฐานนี้ยังคงต้องการการตรวจสอบความถูกต้อง ซึ่งปัจจุบันบุคคลทั่วไปจึงนิยมถอดอักษรตามแบบของตัวเองมากกว่าใช้หลักของราชบัณฑิตยสถาน โดยใช้อัลกอริทึม

เอกสารนี้เป็นเอกสารต้นฉบับเพื่อการศึกษาค้นคว้าเท่านั้น เมื่ออนุญาตให้เผยแพร่โดยไม่เสียค่าใช้จ่าย
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ชวาว์เด็กซ์(เราจะอธิบายในตอน ที่ 3) คลังข้อมูลคำศัพท์ที่ใช้ถูกรวบรวมมาจาก ศูนย์เทคโนโลยีอิเล็กทรอนิกส์และคอมพิวเตอร์แห่งชาติ (NECTEC) ก่อนทำการถอดอักษรไทย-โรมัน คำศัพท์ภาษาไทยจะถูกนำมาเรียงใหม่ตามตัวอักษรที่ใช้แทนในเสียงภาษาอังกฤษ และรูปวรรณยุกต์จะถูกตัดออก โดยจะถูกเรียง โดยขั้นตอนวิธีด้วยคอมพิวเตอร์ หลังจากนั้นจะตรวจสอบและแก้ไขให้อีกถูกต้องโดยบุคคลอีกรอบ เป็นไปได้ที่จะมีตัวอักษรภาษาอังกฤษมากกว่าหนึ่งตัวเพื่อใช้ในการแทนตัวอักษรภาษาไทย เช่น ‘th’ แทน ‘ท’, ‘ia’ แทน ‘- ีย’ ยกตัวอย่าง ขั้นตอนการเรียงตัวอักษรใหม่ได้ดังต่อไปนี้

เสียด	→	ส ี ย	→	Sia
ร่างกาย	→	ร าง กาย	→	Rangkai
และ	→	ล แะ	→	Lae
ทำให้	→	ท ำ ห ใ	→	Tamhai
เสียง	→	ส ี ย ิง	→	Siang

โดยในบทความนี้เราจะแบ่งออกเป็น 5 ส่วน ส่วนที่ 2 คือ การสลับตำแหน่งสระประสม ส่วนที่ 3 คือ อัลกอริทึมชวาว์เด็กซ์ ส่วนที่ 4 คือ การเข้ารหัสตัวอักษร และส่วนที่ 5 คือ สรุปผล

การอ่านออกเสียงซึ่งจะทำให้ชื่อที่อ่านออกเสียงเหมือน หรือคล้ายกัน มีรหัสเสียงที่เหมือนกัน อัลกอริทึมนี้จึงมีพื้นฐานมาจากชื่อที่เป็นภาษาอังกฤษและยังเข้ารหัสเพียงพยัญชนะเท่านั้น การเข้ารหัสเสียงนี้จะแปลงรหัสตาม ตารางที่ 1. จากนั้นรหัสที่เป็นเลข 0 จะถูกตัดออกทั้งหมด และรหัสที่ซ้ำกันและ

อยู่ติดกันจะถูกลดให้เหลือเพียงตัวเดียว ขั้นตอนสุดท้ายคือให้เรียงโดยใช้ตัวอักษรแรกเป็นรหัสตามด้วยรหัสตัวเลขสามตัว เลขที่แปลงได้ ตัวอย่างเช่น ALEXANDER จะ ถูกแปลงเป็น A425

2. การสลับตำแหน่งสระประสม

ก่อนการถอดอักษรไทย-โรมัน เราต้องเตรียมข้อมูลเพื่อสร้างเป็นคลังคำศัพท์ที่รวบรวมการจับคู่กัน ระหว่างคำศัพท์ภาษาไทยและคำอ่านภาษาอังกฤษ วิธีการนี้เราเรียกว่าการจับคู่ข้อมูลวิธีนี้จะ จับคู่คำศัพท์ภาษาไทยกับคำภาษาอังกฤษที่มีคำอ่านเหมือนหรือคล้ายกับคำศัพท์ภาษาไทยนั้นๆ ในขั้นตอนนี้ตัวอักษรจะถูกแยกและนำมาเรียงเป็นพยางค์โดยอัลกอริทึมบนคอมพิวเตอร์ [5] เพื่อการนำมาเรียงได้อย่างถูกต้อง เราต้องกำหนด คู่ตัวอักษร ของ ภาษา ไทย และ ภาษาอังกฤษก่อน แต่ตัวอักษรบางกลุ่มยังไม่สามารถหาคู่ที่ถูกต้องได้เนื่องจากการผันเสียงหรือการงออกเสียงของคำ เช่น ‘ทร’สามารถออกเสียงได้ 2 แบบ ทั้งเสียง ‘ทร’ และเสียง ‘ซ’ ในขั้นตอนนี้จึงมีความสำคัญมากต่อการนำไปเข้ารหัสในขั้นตอนต่อไป

3. ชวาว์เด็กซ์อัลกอริทึม

ชวาว์เด็กซ์ถูกใช้ในการตรวจสอบการสะกดคำ [1]และการสืบค้นข้ามภาษา[3] ในขั้นตอนนี้เราพบว่า การจับคู่เสียงของตัวอักษรระหว่างภาษาไทยกับภาษาอังกฤษไม่เป็นความสัมพันธ์แบบหนึ่งต่อหนึ่ง เนื่องจากตัวอักษรไทยหนึ่งเสียงสามารถเขียนได้ในหลายรูป เพื่อการแก้ปัญหาที่เราได้ใช้อัลกอริทึมชวาว์เด็กซ์เริ่มแรกถูกพัฒนาขึ้น โดย M.K. Odell and R.C.

Russell ถูกใช้ในการเข้ารหัสเสียงของชื่อตามเสียงพูดของพยัญชนะ (ไม่สนใจตัวอักษรทางซ้ายสุด) เป็นตัวเลขตาม

หลักการเข้ารหัสเหมือนกับของ M.K. Odell and R.C. Russell คือ ใช้ตัวแรกเป็นตัวอักษร และตามด้วยตัวเลขโดยไม่สนใจวรรณยุกต์และสระ

ตารางที่ 1. การจับคู่รหัสชวามันเด็กซ์กับตัวอักษรภาษาอังกฤษของ M.K. Odell and R.C. Russell

ตัวอักษร	รหัสตัวเลข
A E I O U H W Y	0
B F P V	1
C G J K Q S X Z	2
D T	3
L	4
M N	5
R	6

ชวามันเด็กซ์อัลกอริทึมเป็นวิธีที่เร็วและค่อนข้างมีความแม่นยำ แต่ก็อาจจะพบปัญหาได้หากตัวอักษรนั้นออกเสียงไม่ตรงกับที่กำหนดไว้

4. การเข้ารหัสตัวอักษร

สำหรับการเข้ารหัสชวามันเด็กซ์ในภาษาไทย โดยส่วนใหญ่จะมีต้นแบบมาจากการเข้ารหัสชวามันเด็กซ์ของ M.K. Odell and R.C. Russell เช่น ของคุณประยูทธ สุวรรณวิสารท ได้ใช้

- การใช้ตัวอักษรตัวแรกเป็นรหัส ถ้าเป็นสระให้ใช้อักษร 'A' เป็นตัวแทนสำหรับสระนั้นๆ

- ขยายความยาวของรหัสชวามันเด็กซ์ที่ได้เป็น 9 ตัว จาก 4 ตัว เพื่อรหัสที่มีความแจ่มจ่มมากขึ้น

- ใช้รหัส ASCII แทนสระ (ตามตารางที่ 4) และแปลงสระเป็นรูปของอักษรภาษาอังกฤษ เช่น 'เอ' เป็น 'ao', 'ไอ' เป็น 'ai' (แสดงในตารางที่ 4)

- ไม่คำนึงถึงเสียงสูงต่ำในขั้นตอนการเข้ารหัสชวามันเด็กซ์

ในโครงการนี้ได้นำอัลกอริทึมชวามันเด็กซ์ภาษาไทยที่พัฒนาโดยคุณประยูทธ สุวรรณวิสารท มาดัดแปลงเพิ่มเติม โดยจะเข้ารหัสตัวเลขทั้งพยัญชนะและสระ

ตารางที่ 2. แสดงการจับคู่รหัสสำหรับตัวอักษรตัวหลังจากตัวแรก(รหัสตัวเลข)

ตัวอักษร อังกฤษ	ตัวอักษรไทย	รหัส	หมายเหตุ
B F P V Ph	บ ฟ พ ฝ พ ฝ	1	
K C G Kh Ch S	ก ข ค ก ข จ ฉ ช ฉ ช ศ ส ส	2	
	จ ฉ ช ฉ ศ ส ส	3	#1
D T	ฎ ฏ ฐ ฑ ฒ ถ ฑ ฐ	3	
L	ล พ	4	
M	ม	@	
N	ณ น	5	
R	ร	6	
	ร ฤ ล พ	5	#1
H	ห ฮ	7	
W	ว	8	
Y	ญ ย	9	
Ng	ง	52	

#1 : สำหรับตัวอักษรตัวสุดท้าย

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น เมื่อนำไปใช้โดยไม่ได้รับอนุญาตให้เผยแพร่โดยไม่ได้รับอนุญาต หรือต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

เอกสารอ้างอิง

- [1] T. Karoonboonyanan, V. Sornlertlamvanich and S. Meknavin. **A Thai Soundex System for Spelling Correction. Proceedings of the Natural Language Processing Pacific Rim Symposium (1997): 633-636.**
- [2] The Royal Institute. The principle of Romanization on the basis of transcription, 2482.
- [3] Suwanvisat, Prayut and Prasitjutrakul, Somchai. 1998. **Thai-English Cross-Language Transliterated Word Retrieval using Soundex Technique.** The National Computer Science and Engineering Conference, Kasetsart University, Bangkok, Thailand.
- [4] A.B. Grisworld. **Thought on the Romanization of Siamese.** Printed for distribution at the obsequies of Colonel Kasem Nandhakij, February 25, 1969
- [5] Wanwara Chaiyapruerk, The Development of an English Thai Transliteration Program Base on Transliterated-word Corpus of the Royal Institute, Department of Linguistics, Faculty of Arts Chulalongkorn University, 2004.

Thai alphabet (อักษรไทย)

Created by Simon Ager, Omniglot.com – the guide to writing systems and languages

Consonants (พยัญชนะ)

ก	ข	ฃ	ค	ฅ	ฉ	ง	จ
ก ไก่	ข ไข่	ฃ ขวด	ค ควาย	ฅ คน	ฉ ระฆัง	ง งู	จ จาน
ko khai	kho khai	kho khuat	kho khwai	kho khon	kho rakhang	ngo ngu	cho chan
k/k	kh/k	kh/k	kh/k	kh/k	kh/k	ng/ng	ch/t
mid	high	high	low	low	low	low	mid
ฌ	ช	ฌ	ฎ	ญ	ฎ	ฏ	ฐ
ฌ ฌิ่ง	ช ช้าง	ฌ โข่	ฎ เถอ	ญ หยิ่ง	ฎ ฎา	ฏ ปฏัก	ฐ ฐาน
cho ching	cho chang	so so	cho choe	yo ying	do cha-da	to pa-tak	tho than
ch/-	ch/t	s/t	ch/-	y/n	d/t	t/t	th/t
high	low	low	low	low	mid	mid	high
ฑ	ฒ	ณ	ด	ต	ถ	ท	ธ
ฑ มณโฑ	ฒ ผู้เต่า	ณ เณร	ด เด็ก	ต เต่า	ถ ถุง	ท พหาร	ธ ธง
tho montho	tho phuthao	no nen	do dek	to tao	tho thung	tho thahan	tho thong
th/t	th/t	n/n	d/t	t/t	th/t	th/t	th/t
low	low	low	mid	mid	high	low	low
น	บ	ป	ฝ	ฝ	พ	ฟ	ภ
น นู	บ ใบไม้	ป ปลา	ฝ ฝิ่ง	ฝ ฝ่า	พ พาน	ฟ ฟัน	ภ สำภา
no nu	bo baimai	po pla	pho phueng	fo fa	pho phan	fo fan	pho samphao
n/n	b/p	p/p	ph/-	f/-	ph/p	f/p	ph/p
low	mid	mid	high	high	low	low	low
ม	ย	ร	ล	ว	ศ	ษ	ส
ม ม้า	ย ยักษ์	ร เรือ	ล ลิง	ว แหวน	ศ ศาลา	ษ ษายี่	ส เสือ
mo ma	yo yak	ro ruea	lo ling	wo waen	so sala	so rue-si	so suea
m/m	y/y	r/n	l/n	w/w	s/t	s/t	s/t
low	low	low	low	low	high	high	high

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Thai alphabet (อักษรไทย)

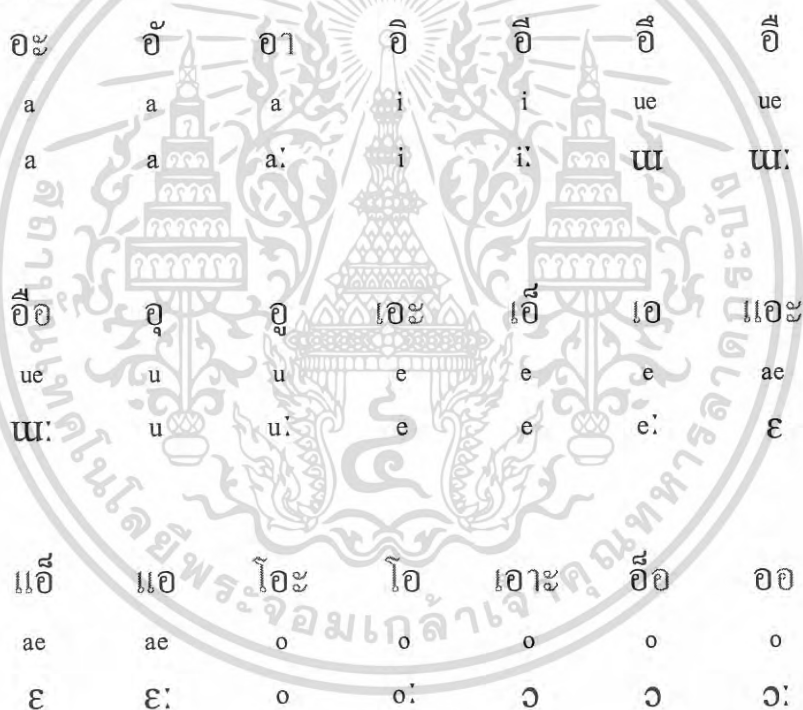
Created by Simon Ager, Omniglot.com – the guide to writing systems and languages

ห	ฬ	อ	ฮ
ห หีบ	ฬ จุฬา	อ อ่าง	ฮ นกฮูก
ho hip	lo chu-la	o ang	ho nok-huk
h/-	l/n	-	h/-
high	low	mid	low

Notes

- The sounds represented by some consonants change when they are used at the end of a syllable (indicated by the letters on the right of the slash). Some consonants can only be used at the beginning of a syllable.
- Duplicate consonants represent different Sanskrit and Pali consonants sounds which are pronounced identically in Thai.

Vowel diacritics (สระ)



อะ	อิ	อา	อึ	อึ	อึ	อึ
a	a	a	i	i	ue	ue
a	a	a:	i	i:	ue	ue:
อุ	อู	อู	เอะ	เอ	เอ	เอะ
ue	u	u	e	e	e	ae
ue:	u	u:	e	e	e:	ε
แอ	แอ	โอะ	โอ	เออะ	อึ	ออ
ae	ae	o	o	o	o	o
ε	ε:	o	o:	ว	ว	ว:

อึ	เออะ	เออ	เอ็
o	oe	oe	oe
ว:	ว?	ว:	ว:

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Thai alphabet (อักษรไทย)

Created by Simon Ager, Omniglot.com – the guide to writing systems and languages

Diphthongs

เียะ	เีย	เือะ	เือ	ัวะ	ัว	ว
ia	ia	uea	uea	ua	ua	ua
iaʔ	ia	๓aʔ	๓a	uaʔ	ua	ua

ิว	เือ	เิว	เือว	เือ	เือว	เียว
io	eo	eo	aeo	ao	ao	iao
iu/iw	eu/ew	e.u/e.w	ε.u/ε.w	au/aw	a.u	iau/iaw

อัย	ไอ	ไอ	ไย	อาย	อัย	อัย
ai	ai	ai	ai	ai	oi	oi
ai/aj	ai/aj	ai/aj	ai/aj	a.i/a.j	ว/ว	ว:i/ว:j

ไย	อุย	เอย	อวย	เือย	อำ	ฤ
oi	ui	oei	uai	ueai	am	rue
o:i/o:j	ui/uj	๕:i/๕:j	uai/uaj	๓ai/๓aj	am	r๓/ri

ฤ	ฤ	ฤ	อ	อ
rue	lue	lue	๓	silences
r๓:	๓	๓:	~	consonant

The vowel diacritics are shown with the letter o ang (อ), which acts as a silent vowel carrier at the beginning of words that start with a vowel.

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Thai alphabet (อักษรไทย)

Created by Simon Ager, Omniglot.com – the guide to writing systems and languages

Punctuation

๐	คาไก่	<u>dtaa gai</u>	๐	โคมูตร	<u>khoh muut</u>
		beginning of a new paragraph			marks the end of a story
					Tone markers

Tone markers

ก๋	ไม้เอก	low tone	กั	ไม้โท	falling tone
	mai ehk			mai toh	
กั	ไม้ตรี	high tone	กั+	ไม้จัตวา	rising tone
	maai tree			maai jat dta	
				waa	

Tone indication

	Open syllables			Closed syllables *	
	unmarked	อ๋	อั	short vowel	long vowel
Class 1	mid	low	falling	low	low
Class 2	rising	low	falling	low	low
Class 3	mid	falling	high	high	falling

* Closed syllables are those ending with p, t or k

Numerals (เลขไทย)

๐	๑	๒	๓	๔	๕	๖	๗	๘	๙	๑๐
ศูนย์	หนึ่ง	สอง	สาม	สี่	ห้า	หก	เจ็ด	แปด	เก้า	สิบ
sun	neung	song	sam	si	ha	hok	chet	paet	kao	sip
0	1	2	3	4	5	6	7	8	9	10

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้