

การเข้ารหัสคำทับศัพท์เพื่อการค้นคืนข้ามภาษาไทย-อังกฤษ



นาย ประยุทธ์ สุวรรณวิสารท

สถาบันวิทยบริการ

วิทยานิพนธ์นี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญาวิทยาศาสตรมหาบัณฑิต
สาขาวิชาวิทยาศาสตร์คอมพิวเตอร์ ภาควิชาวิศวกรรมคอมพิวเตอร์

บัณฑิตวิทยาลัย จุฬาลงกรณ์มหาวิทยาลัย

ปีการศึกษา 2541

ISBN 974-332-123-3

ลิขสิทธิ์ของบัณฑิตวิทยาลัย จุฬาลงกรณ์มหาวิทยาลัย

11 ก.ย. 2545

I1940329X

**TRANSLITERATED WORD ENCODING FOR THAI-ENGLISH
CROSS-LANGUAGE RETRIEVAL**



Mr. Prayut Suwanvisat

สถาบันวิทยบริการ
จุฬาลงกรณ์มหาวิทยาลัย

A Thesis Submitted in Partial Fulfillment of the Requirements

for the Degree of Master of Science in Computer Science

Department of Computer Engineering

Graduate School

Chulalongkorn University

Academic Year 1998

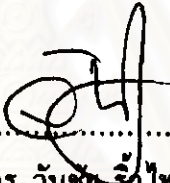
ISBN 974-332-123-3

หัวข้อวิทยานิพนธ์ การเข้ารหัสคำทับศัพท์เพื่อการค้นคืนข้ามภาษาไทย-อังกฤษ
โดย นาย ประยุทธ์ สุวรรณวิสารท
ภาควิชา วิศวกรรมคอมพิวเตอร์
อาจารย์ที่ปรึกษา ผู้ช่วยศาสตราจารย์ ดร. สมชาย ประสิทธิ์ภูตระกูล


บัณฑิตวิทยาลัย จุฬาลงกรณ์มหาวิทยาลัย อนุมัติให้บัณฑิตวิทยาลัยฉบับนี้เป็นส่วนหนึ่ง
ของการศึกษาตามหลักสูตรปริญญาโทบัณฑิต


..... คณบดีบัณฑิตวิทยาลัย
(ศาสตราจารย์ นายแพทย์ สุภวัฒน์ ชูติวงศ์)

คณะกรรมการสอบวิทยานิพนธ์


..... ประธานกรรมการ
(รองศาสตราจารย์ ดร. วันชัย ใจไพบูลย์)


..... อาจารย์ที่ปรึกษา
(ผู้ช่วยศาสตราจารย์ ดร. สมชาย ประสิทธิ์ภูตระกูล)


..... กรรมการ
(ผู้ช่วยศาสตราจารย์ ดร. ประภาส จงสถิตย์วัฒนา)


..... กรรมการ
(อาจารย์ ดร. บุญเสริม กิจศิริกุล)

ฉบับนี้ได้รับแจ้งให้เผยแพร่โดยสำนักพิมพ์มหาวิทยาลัยเทคโนโลยีพระจอมเกล้าพระนครเหนือเพียงเล่มนี้เท่านั้น

ประยูทธ สุวรรณวิสารท : การเข้ารหัสคำทับศัพท์เพื่อการค้นคืนข้ามภาษาไทย-อังกฤษ
(TRANSLITERATED WORD ENCODING FOR THAI-ENGLISH CROSS-LANGUAGE
RETRIEVAL) อ.ที่ปรึกษา : ผศ. ดร. สมชาย ประสิทธิ์จูตระกูล. 85 หน้า. ISBN 974-332-123-3.

วิทยานิพนธ์ฉบับนี้นำเสนอขั้นตอนวิธีการเข้ารหัสคำทับศัพท์เพื่อการค้นคืนข้ามภาษาไทย-อังกฤษ ซึ่งอนุญาตให้ใช้ข้อความที่เป็นคำทับศัพท์ภาษาอังกฤษหรือภาษาไทยในการค้นคืนเอกสารที่มีคำหลักตรงกันในอีกภาษา โดยมีข้อสมมุติฐานว่าสามารถทำการค้นคืนข้ามภาษาไทย-อังกฤษได้โดยไม่ต้องอาศัยพจนานุกรม ขั้นตอนวิธีที่นำเสนอแบ่งออกเป็นสองส่วนคือ (1) ขั้นตอนวิธีการเข้ารหัสคำทับศัพท์เพื่อการค้นคืนข้ามแบบภาษาไทยทับศัพท์ภาษาอังกฤษ และ (2) ขั้นตอนวิธีการเข้ารหัสคำทับศัพท์เพื่อการค้นคืนข้ามแบบภาษาอังกฤษทับศัพท์ภาษาไทย ขั้นตอนวิธีการค้นคืนข้ามภาษานี้จะทำงานโดยการเข้ารหัสคำในข้อความแล้วนำรหัสคำที่ได้ไปเปรียบเทียบกับรหัสคำในดัชนีคำหลัก การเปรียบเทียบรหัสคำในการข้ามภาษาแบบภาษาไทยทับศัพท์ภาษาอังกฤษจะอาศัยวิธีการเปรียบเทียบแบบเหมือนกันทุกประการ ส่วนการเปรียบเทียบรหัสคำในการข้ามภาษาแบบอังกฤษทับศัพท์ภาษาไทยจะอาศัยวิธีการเปรียบเทียบเชิงประมาณและแยกเปรียบเทียบส่วนพยัญชนะและสระออกจากกัน โดยใช้เทคนิคกำหนดการพลวัต ผลการทดลองแสดงให้เห็นว่าขั้นตอนวิธีการเข้ารหัสคำทับศัพท์เพื่อการค้นคืนข้ามภาษาไทย-อังกฤษแบบภาษาไทยทับศัพท์ภาษาอังกฤษมีค่าเรียกคืนสูงถึง 90 เปอร์เซ็นต์ และค่าแม่นยำสูงถึง 78 เปอร์เซ็นต์ เมื่อคำทับศัพท์มีความยาวมากกว่า 7 ตัวอักษร และแบบภาษาอังกฤษทับศัพท์ภาษาไทยมีค่าเรียกคืนสูงถึง 73 เปอร์เซ็นต์ และค่าแม่นยำสูงถึง 69 เปอร์เซ็นต์

สถาบันวิทยบริการ
จุฬาลงกรณ์มหาวิทยาลัย

ภาควิชา วิศวกรรมคอมพิวเตอร์
สาขาวิชา วิศวกรรมศาสตร์คอมพิวเตอร์
ปีการศึกษา 2541

ลายมือชื่อนิติ ประยูทธ สุวรรณวิสารท
ลายมือชื่ออาจารย์ที่ปรึกษา ประสิทธิ์จูตระกูล
ลายมือชื่ออาจารย์ที่ปรึกษาร่วม

4070330721 : MAJOR COMPUTER SCIENCE

KEY WORD: TRANSLITERATED WORD / ENCODING / CROSS-LANGUAGE / SOUNDEX

PRAYUT SUWANVISAT : TRANSLITERATED WORD ENCODING FOR THAI-ENGLISH CROSS-LANGUAGE RETRIEVAL. THESIS ADVISOR : ASSIST. PROF. SOMCHAI PRASITJUTRAKUL, Ph.D. 85 pp. ISBN 974-332-123-3.

This thesis presents two algorithms for transliterated word encoding for Thai-English cross-language retrieval. The algorithms enable retrieval of documents containing either the English-to-Thai or Thai-to-English transliterated keywords. We have a hypothesis that cross-language retrieval does not use a dictionary. The proposed algorithms are (1) English-to-Thai transliterated word encoding for cross-language retrieval algorithm and (2) Thai-to-English transliterated word encoding for cross-language retrieval algorithm. This cross-language retrieval is done by encoding each word in the query terms and then matching the query code with codes of keywords in the index. The English-to-Thai cross-language retrieval uses exact code matching. On the other hand, the Thai-to-English uses approximate code matching (separately done for consonant and vowel parts) by using dynamic programming technique. Experimental results showed that for keywords of length longer than seven characters the recall and precision of the English-to-Thai transliterated word cross-language retrieval are 90% and 78%, respectively. The recall and precision of the Thai-to-English transliterated word are around 73% and 69%, respectively.

สถาบันวิทยบริการ

จุฬาลงกรณ์มหาวิทยาลัย

ภาควิชา.....วิศวกรรมคอมพิวเตอร์

สาขาวิชา.....วิทยาศาสตร์คอมพิวเตอร์

ปีการศึกษา..... 2541

ลายมือชื่อผู้ผลิต..... วิเชษณ์ สุวรรณวิสารท

ลายมือชื่ออาจารย์ที่ปรึกษา..... วิเชษณ์ สุวรรณวิสารท

ลายมือชื่ออาจารย์ที่ปรึกษาร่วม.....

กิตติกรรมประกาศ

วิทยานิพนธ์ฉบับนี้สำเร็จรุด่วงไปได้ด้วยความช่วยเหลืออย่างดียิ่งของผู้ช่วยศาสตราจารย์ ดร. สมชาย ประสิทธิ์ฐิตระกูล อาจารย์ที่ปรึกษาวิทยานิพนธ์ ซึ่งท่านได้ให้คำแนะนำและข้อคิดเห็นต่าง ๆ ในการวิจัยมาด้วยดีตลอด รวมทั้งตรวจแก้วิทยานิพนธ์ฉบับนี้อีกหลายครั้ง ผู้เขียนขอขอบพระคุณผู้ช่วยศาสตราจารย์ ดร. ประกาศ จงสถิตย์วัฒนา ซึ่งท่านได้ให้แนวคิดในการทำงานวิจัยในภาควิชาวิศวกรรมคอมพิวเตอร์

ขอขอบพระคุณมูลนิธิเพื่อการศึกษาคอมพิวเตอร์และการสื่อสาร (Computer and Communication Education Foundation) และมูลนิธิเพื่อการศึกษาและวิจัยวิทยาศาสตร์คอมพิวเตอร์ จุฬาลงกรณ์มหาวิทยาลัย ที่สนับสนุนทุนการศึกษาและทุนในการทำงานวิจัยครั้งนี้

ท้ายนี้ ผู้วิจัยขอกราบขอบพระคุณมารดา ซึ่งสนับสนุนและให้กำลังใจแก่ผู้วิจัยเสมอจนสำเร็จการศึกษา



สถาบันวิทยบริการ
จุฬาลงกรณ์มหาวิทยาลัย

สารบัญ

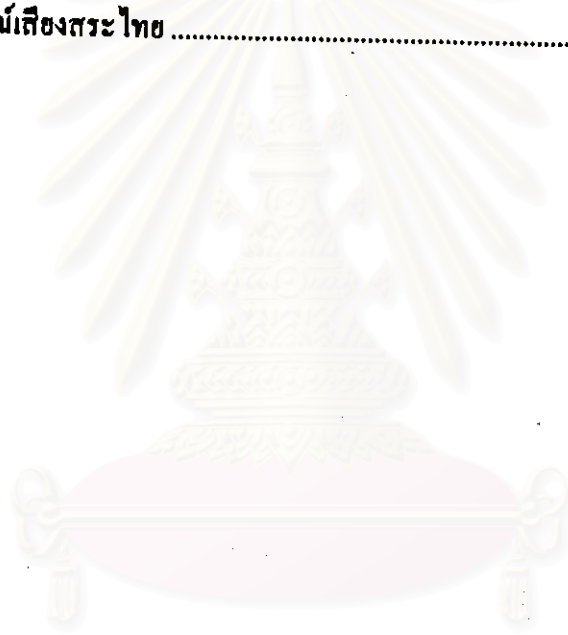
	หน้า
บทคัดย่อภาษาไทย	ง
บทคัดย่อภาษาอังกฤษ	จ
กิตติกรรมประกาศ.....	ฉ
สารบัญ	ช
สารบัญตาราง	ณ
สารบัญภาพ	ญ
1. บทนำ	1
1.1 ความเป็นมาและความสำคัญของปัญหา.....	1
1.2 วัตถุประสงค์ของการวิจัย	3
1.3 ขอบเขตของการวิจัย.....	3
1.4 ขั้นตอนและวิธีดำเนินการวิจัย.....	3
1.5 ประโยชน์ที่คาดว่าจะได้รับ	3
1.6 ผลงานที่ตีพิมพ์จากงานวิจัย.....	4
1.7 โครงสร้างของวิทยานิพนธ์	4
2. ทฤษฎีและงานวิจัยที่เกี่ยวข้อง	5
2.1 การถอดอักษร	5
2.2 การใช้ตัวอักษรโรมันเพื่อการถ่ายเสียง	7
2.3 หลักเกณฑ์การทับศัพท์	8
2.4 การค้นคืนข้ามภาษา	8
2.5 ขั้นตอนวิธีชาวค้เด็กซ์ภาษาอังกฤษ	9
2.6 ขั้นตอนวิธีชาวค้เด็กซ์ภาษาไทย	11
2.7 ขั้นตอนวิธีระยะแก้ไขสั้นที่สุด (Minimum Edit Distance).....	22
2.8 สรุป.....	24
3. ขั้นตอนวิธีการค้นคืนข้ามภาษาแบบภาษาไทยทับศัพท์ภาษาอังกฤษ	25
3.1 โครงสร้างของระบบค้นคืนข้ามภาษาไทยทับศัพท์ภาษาอังกฤษ.....	25
3.2 ขั้นตอนวิธีการเข้ารหัสคำ.....	25
3.3 วิธีการทดลอง.....	28
3.4 ผลการทดลอง.....	29
3.5 สรุป.....	32

4. ขั้นตอนวิธีการค้นคืนข้ามภาษาแบบภาษาอังกฤษทับศัพท์ภาษาไทย	33
4.1 โครงสร้างของระบบค้นคืนข้ามภาษาอังกฤษทับศัพท์ภาษาไทย.....	33
4.2 ขั้นตอนวิธีการเข้ารหัสคำ.....	34
4.2.1 การประมวลผลตัวอักษรเบื้องต้น.....	34
4.2.1.1 การลดรูปและตัดวรรณยุกต์.....	35
4.2.1.2 การแทนที่สระประสมด้วยสัญลักษณ์เสียงสากล.....	39
4.2.1.3 การถอดอักษร.....	40
4.2.2 การย้ายตำแหน่งสระ.....	44
4.3 ขั้นตอนวิธีเปรียบเทียบรหัสคำ.....	44
4.3.1 การคำนวณหาค่าความแตกต่างของรหัสคำ.....	45
4.3.1.1 การกำหนดต้นทุนในการแก้ไขอักขระ.....	47
4.3.2 เงื่อนไขในการเปรียบเทียบรหัสคำ.....	55
4.4 วิธีการทดลอง.....	57
4.5 ผลการทดลอง.....	58
4.6 สรุป.....	60
5. สรุปผลการวิจัยและข้อเสนอแนะ.....	61
5.1 สรุปผลการวิจัย.....	61
5.2 ข้อดีและข้อเสียของขั้นตอนวิธี.....	62
5.3 ข้อเสนอแนะ.....	63
รายการอ้างอิง.....	64
ภาคผนวก.....	66
ภาคผนวก ก.....	67
ภาคผนวก ข.....	73
ภาคผนวก ค.....	81
ประวัติผู้เขียน.....	85

สารบัญตาราง

	หน้า
ตารางที่ 2.1 ตัวอย่างการถอดอักษรจากภาษาต้นแบบ ไปยังภาษาเป้าหมาย	6
ตารางที่ 2.2 การกำหนดรหัสตัวอักษรภาษาอังกฤษ	11
ตารางที่ 2.3 การกำหนดรหัสตัวอักษรของรหัสตัวอักษรภาษาไทย	12
ตารางที่ 2.4 การกำหนดรหัสตัวเลขของรหัสตัวอักษรภาษาไทย	13
ตารางที่ 2.5 ตัวอย่างการเข้ารหัสตัวอักษรภาษาไทย	14
ตารางที่ 2.6 การกำหนดรหัสอักษรสำหรับอักขระตัวแรกของรหัสตัวอักษรภาษาไทย	15
ตารางที่ 2.7 การกำหนดรหัสอักษรสำหรับอักขระตัวถัดไปของรหัสตัวอักษรภาษาไทย	15
ตารางที่ 2.8 ตัวอย่างการเข้ารหัสตัวอักษรภาษาไทย	16
ตารางที่ 2.9 การเข้ารหัสสำหรับพยัญชนะต้น	17
ตารางที่ 2.10 การเข้ารหัสสำหรับสระ	18
ตารางที่ 2.11 การเข้ารหัสสำหรับตัวสะกด	18
ตารางที่ 2.12 ตัวอย่างการเข้ารหัสตัวอักษรภาษาไทย	21
ตารางที่ 3.1 กลุ่มเสียงของพยัญชนะไทย 21 กลุ่ม	26
ตารางที่ 3.2 กลุ่มอักษรไทยและกลุ่มอักษรอังกฤษที่ออกเสียงคล้ายกันในรหัสตัวอักษร	26
ตารางที่ 3.3 การกำหนดรหัสตัวอักษรสำหรับอักษรไทยและอักษรอังกฤษที่น่าเสนอ	28
ตารางที่ 3.4 รายละเอียดจำนวนคำศัพท์ที่ใช้ในการทดลอง	29
ตารางที่ 3.5 ความสัมพันธ์ระหว่างความยาวน้อยสุดของรหัสคำกับความยาวคำเฉลี่ย	30
ตารางที่ 3.6 ความสัมพันธ์ระหว่างความยาวน้อยสุดของรหัสคำกับประสิทธิภาพ	30
ตารางที่ 4.1 การใช้อักษรโรมันแทนพยัญชนะไทยของ ISO	41
ตารางที่ 4.2 การถอดอักษรอังกฤษเป็นอักษรไทยในส่วนของพยัญชนะที่น่าเสนอ	42
ตารางที่ 4.3 การถอดอักษรอังกฤษเป็นอักษรไทยในส่วนของสระที่น่าเสนอ	43
ตารางที่ 4.4 กลุ่มเสียงของพยัญชนะไทย	48
ตารางที่ 4.5 อักขระที่ไม่สามารถแยกความแตกต่างในการถอดอักษร	48
ตารางที่ 4.6 อักขระแทนมาตราต่าง ๆ	49
ตารางที่ 4.7 อักขระควบไม้แท้	50
ตารางที่ 4.8 อักขระนำเสียงสนิท	50
ตารางที่ 4.9 อักขระควบที่เป็นตัวสะกด	50
ตารางที่ 4.10 ตัวอย่างตารางการกำหนดต้นทุนในการแทนที่อักขระที่เป็นพยัญชนะ	51
ตารางที่ 4.11 ตารางการกำหนดต้นทุนในการแทนที่อักขระที่เป็นสระ	52

ตารางที่ 4.12 ตารางการกำหนดต้นทุนในการเพิ่มหรือการลบอักษร.....	53
ตารางที่ 4.13 ความสัมพันธ์ระหว่างค่าแอฟกับประสิทธิภาพของระบบคั่นคืน	58
ตารางที่ ก.1 ตารางกำหนดต้นทุนในการแก้ไขอักษรพยัญชนะ.....	68
ตารางที่ ข.1 การใช้อักษรโรมันแทนพยัญชนะ ไทยของ ISO	73
ตารางที่ ข.2 การใช้อักษรโรมันแทนอักษรสระไทยของ ISO	74
ตารางที่ ข.3 การใช้อักษรโรมันแทนพยัญชนะ ไทยของราชบัณฑิตยสถาน	75
ตารางที่ ข.4 การใช้อักษรโรมันแทนสระไทยของราชบัณฑิตยสถาน	76
ตารางที่ ข.5 การใช้อักษรโรมันแทนพยัญชนะ ไทยของจันทร์เพ็ญ ไหวหารสุนทร	77
ตารางที่ ข.6 การใช้อักษรโรมันแทนสระไทยของจันทร์เพ็ญ ไหวหารสุนทร	78
ตารางที่ ข.7 กลุ่มเสียงพยัญชนะไทยและสัญลักษณ์เสียง	79
ตารางที่ ข.8 สัญลักษณ์เสียงสระไทย	80



สถาบันวิทยบริการ
จุฬาลงกรณ์มหาวิทยาลัย

สารบัญภาพ

	หน้า
รูปที่ 2.1 โปรแกรมการเข้ารหัสชาวคเด็กซ์ภาษาอังกฤษ.....	10
รูปที่ 2.2 โครงร่างเครื่องเข้ารหัสชาวคเด็กซ์.....	19
รูปที่ 2.3 เครื่องเข้ารหัสชาวคเด็กซ์.....	20
รูปที่ 2.4 โปรแกรมระยะแก้ไขสั้นที่สุด โดยใช้เทคนิคกำหนดการพลวัต	23
รูปที่ 3.1 ความสัมพันธ์ระหว่างความยาวน้อยสุดของรหัสคำกับประสิทธิผล	31
รูปที่ 4.1 ลำดับการทำงานของระบบคั่นคินข้ามภาษา	33
รูปที่ 4.2 การประมวลผลตัวอักษรเบื้องต้น.....	35
รูปที่ 4.3 การเปรียบเทียบรหัสคำที่น่าสนใจ.....	45
รูปที่ 4.4 ความสัมพันธ์ระหว่างค่าแอดฟากับประสิทธิผลของระบบคั่นคิน	59

สถาบันวิทยบริการ
จุฬาลงกรณ์มหาวิทยาลัย