

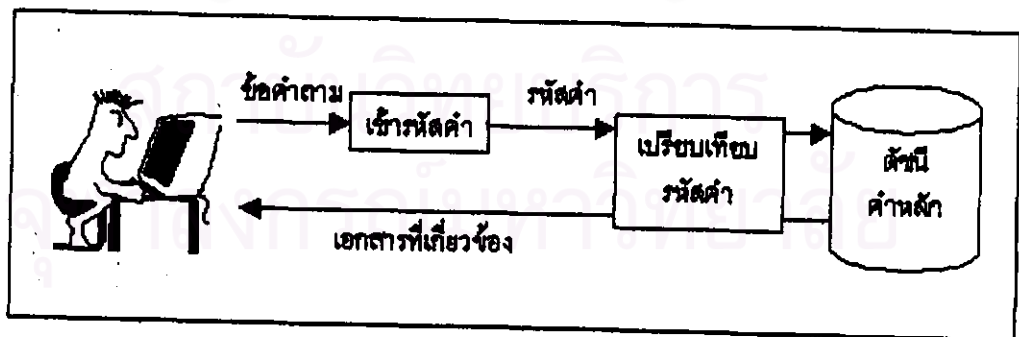


4. ขั้นตอนวิธีการค้นคืนข้ามภาษาแบบภาษาอังกฤษทับศัพท์ภาษาไทย

จากบทที่แล้วได้กล่าวถึงขั้นตอนวิธีสำหรับการเข้ารหัสคำทับศัพท์เพื่อการค้นคืนข้ามภาษาแบบภาษาไทยทับศัพท์ภาษาอังกฤษ ในบทนี้จะกล่าวถึงการค้นคืนในส่วนภาษาอังกฤษทับศัพท์ภาษาไทย โดยจะแบ่งการทำงานของขั้นตอนวิธีออกเป็น 2 ขั้นตอน คือ ขั้นตอนวิธีการเข้ารหัสคำ และขั้นตอนวิธีการเปรียบเทียบรหัสคำ สำหรับการค้นคืน

4.1 โครงสร้างของระบบค้นคืนข้ามภาษาอังกฤษทับศัพท์ภาษาไทย

ขั้นตอนวิธีที่น่าเสนอสำหรับการเข้ารหัสคำทับศัพท์เพื่อการค้นคืนข้ามภาษาไทย-อังกฤษ ในส่วนภาษาอังกฤษทับศัพท์ภาษาไทยนั้นมีลำดับการทำงานดังนี้ เมื่อผู้ใช้งานได้ระบุข้อความที่เป็นภาษาอังกฤษทับศัพท์ภาษาไทยให้กับระบบค้นคืนข้ามภาษาแล้ว ระบบจะทำการเข้ารหัสคำในข้อความ เมื่อได้รหัสคำแล้วจะนำไปเปรียบเทียบกับรหัสคำในดัชนีคำหลักของเอกสารที่ได้เข้ารหัสไว้แล้วในขั้นตอนการทำงานนี้ คำหลักใดที่ผ่านเงื่อนไขการเปรียบเทียบจะถือว่าคำหลักนั้นเป็นคำหลักที่ตรงกันในอีกภาษาหนึ่งดังแสดงในรูปที่ 4.1 ขั้นตอนวิธีที่น่าเสนอนี้เปิดโอกาสให้ผู้ใช้สามารถระบุข้อความที่เป็นภาษาอังกฤษทับศัพท์ภาษาไทยเพื่อค้นคืนเอกสารที่เป็นภาษาไทย และข้อความที่เป็นภาษาไทยเพื่อค้นคืนเอกสารที่เป็นภาษาอังกฤษทับศัพท์ภาษาไทย โดยขั้นตอนวิธีนี้สามารถแบ่งออกเป็น 2 ส่วนคือ ขั้นตอนวิธีการเข้ารหัสคำ และขั้นตอนวิธีการเปรียบเทียบรหัสคำ



รูปที่ 4.1 ลำดับการทำงานของระบบค้นคืนข้ามภาษา

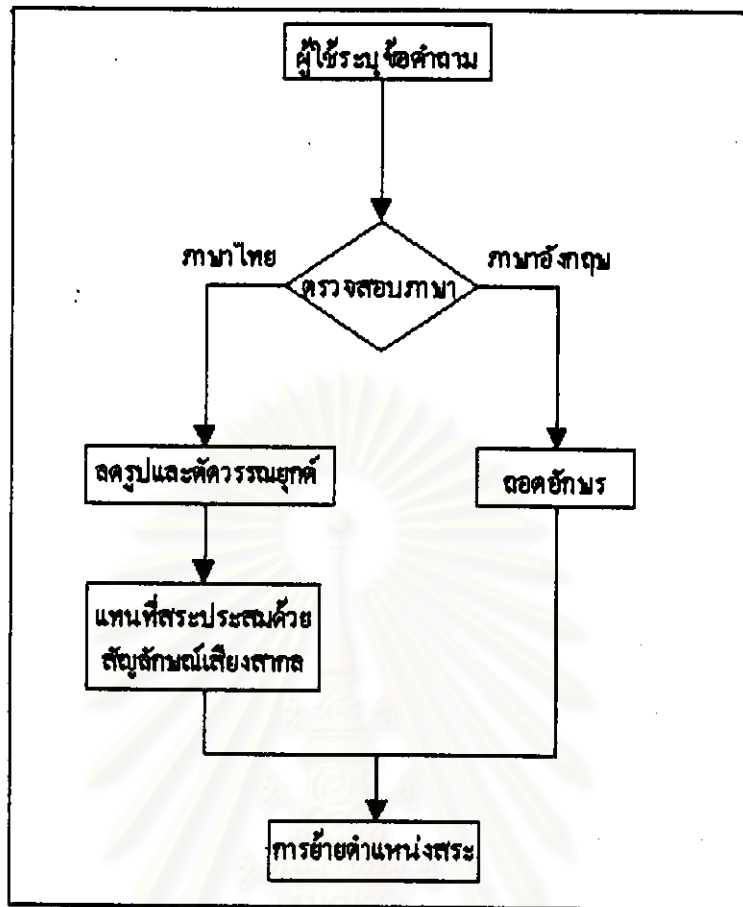
4.2 ขั้นตอนวิธีการเข้ารหัสคำ

ขั้นตอนวิธีการเข้ารหัสคำที่น่าเสนอนั้นแบ่งออกเป็นสองขั้นตอนย่อย คือ ขั้นตอนการประมวลผลตัวอักษรเบื้องต้นและการย้ายตำแหน่งสระ เพื่อให้สอดคล้องกับโครงสร้างของระบบค้นคืนข้ามภาษาอังกฤษทับศัพท์ภาษาไทยที่กล่าวไว้ในหัวข้อ 4.1 ซึ่งขั้นตอนการประมวลผลตัวอักษรเบื้องต้นจะตรวจสอบข้อความที่ผู้ใช้ระบุว่าเป็นภาษาอังกฤษทับศัพท์ภาษาไทยหรือเป็นภาษาไทยโดยจะเปลี่ยนให้ข้อความทั้งสองแบบนี้อยู่ในรูปแบบเดียวกันก่อนจึงจะทำการย้ายตำแหน่งสระต่อไป

4.2.1 การประมวลผลตัวอักษรเบื้องต้น

เพื่อให้การเข้ารหัสคำไทยหรือคำอังกฤษทับศัพท์ได้รหัสตรงกันนั้น งานวิจัยนี้ได้อาศัยหลักการอ่านออกเสียงที่เหมือนกัน ซึ่งการอาศัยแค่หน่วยอักษรบางครั้งไม่สามารถบอกได้ว่ามีหน่วยเสียงที่ถูกต้องเป็นอย่างไร เช่น หน่วยอักษร “ร” จะถอดเป็นหน่วยเสียง /r/ แต่ถ้ามี “ร” สองตัวติดกันอาจจะเป็น ร หัน ซึ่งมีหน่วยเสียงเป็น / ฺร / ในกรณีที่มีตัวสะกดตามหลัง หรือมีหน่วยเสียงเป็น / ฺน / ในกรณีที่ไม่เป็นตัวสะกดตามหลัง เป็นต้น ดังนั้นผู้วิจัยจึงได้นำเสนอขั้นตอนการประมวลผลตัวอักษรเบื้องต้นเพื่อแก้ปัญหาในเรื่องของการอ่านออกเสียงของคำในภาษาไทย และทำให้ข้อความที่เป็นภาษาไทยหรือภาษาอังกฤษทับศัพท์ภาษาไทยอยู่ในรูปแบบเดียวกัน เพื่อให้ง่ายต่อการประมวลผลต่อไป

การประมวลผลตัวอักษรเบื้องต้น ประกอบด้วย การลดรูปและตัดวรรณยุกต์ การแทนที่สระประสมด้วยสัญลักษณ์เสียงสากล และการถอดอักษร โดยจะพิจารณาจากข้อความที่ผู้ใช้ระบุให้กับระบบ ถ้าข้อความเป็นภาษาไทยจะต้องทำการการประมวลผลตัวอักษรเบื้องต้นในส่วนของการลดรูปและตัดวรรณยุกต์ และการแทนที่สระประสมด้วยสัญลักษณ์เสียงสากล แต่ถ้าข้อความเป็นภาษาอังกฤษทับศัพท์ภาษาไทยจะต้องทำการประมวลผลตัวอักษรเบื้องต้นในส่วนของการถอดอักษร ก่อนดังแสดงในรูปที่ 4.2



รูปที่ 4.2 การประมวลผลตัวอักษรเบื้องต้น

4.2.1.1 การลดรูปและตัดวรรณยุกต์

การลดรูปและตัดวรรณยุกต์เป็นการประมวลผลตัวอักษรเบื้องต้นส่วนหนึ่งที่พยายามเปลี่ยนรูปแบบการสะกดของคำไทยให้อยู่ในรูปแบบที่ง่ายต่อการประมวลผลในขั้นต่อไป โดยพยายามทำให้ผลลัพธ์ที่ได้จากขั้นตอนนี้มีรูปแบบเหมือนกันผลลัพธ์ที่ได้จากการถอดอักษรคำภาษาอังกฤษทับศัพท์ภาษาไทย ผู้วิจัยได้นำเสนอขั้นตอนวิธีการอ่านออกเสียงคำไทยโดยแบ่งเป็นกรณีต่าง ๆ ดังนี้

- วรรณยุกต์ และไม้ไต่คู้ ทำการตัดวรรณยุกต์และไม้ไต่คู้ทิ้ง เนื่องจากในการถอดอักษรไทยเป็นอักษรอังกฤษ จะไม่พิจารณาวรรณยุกต์และไม้ไต่คู้ เช่น ช้าง ถอดอักษรเป็น CHANG
- ร ทำการเปลี่ยน ร เป็น ~น ในกรณีที่ไม่มีตัวสะกดตามหลัง เช่น จรรยา บรรจบ ครรภ์ เป็นต้น และเป็น ~ ในกรณีที่มิด้ตัวสะกดตามหลัง เช่น ชรรม พรรณ กรรม

เป็นต้น จากการวิเคราะห์รูปแบบของคำผู้วิจัยได้สร้างรูปแบบของ รร โดยยึดหลักทางภาษาศาสตร์¹ และข้อมูลทางสถิติ² ซึ่งมีรูปแบบดังนี้

กำหนดให้

... อักษรระใด	<V> สระ
<E> จุดสิ้นสุดสายอักษร	<IV> สระหน้า คือ เ-แ- ใ- ใ- และ โ-
<C> พยัญชนะ	

รูปแบบ รร ที่เปลี่ยนเป็น ัน

...รร<E>	เช่น สรร
...รร<C><C><E>	เช่น บรรทม
...รร<C>' ...	เช่น สรรค์
...รร<IV>...	เช่น บรรแดง
...รร<C><V>...	เช่น สรรช้อย

รูปแบบ รร ที่เปลี่ยนเป็น ัย

...รร<C><E>	เช่น ชรรม
...รร<C><IV>...	เช่น บรรณโลก

- สระ ใ- ใ- และ ใ-ย ทำการเปลี่ยนสระดังกล่าวทั้งหมดเป็น ัย รูปแบบเดียวเพื่อเพิ่มความถูกต้องยิ่งขึ้นในการเข้ารหัสคำ เนื่องจากสระดังกล่าวมีการอ่านออกเสียงที่คล้ายกันมาก และในการถอดอักษรมักจะถอดเป็น -AI เหมือนกัน เช่น ไซ ไซ ไซย และ ัย ต่างก็ถอดอักษรเป็น CHAI เป็นต้น การเปลี่ยนสระ ใ- ใ- และ ใ-ย มีหลักเกณฑ์ต่าง ๆ ดังนี้

1. การเปลี่ยนสระ ใ- เป็น ัย เนื่องจากสระ ใ- ที่ใช้กับอักษรควบซึ่งมีทั้งหมด 5 คำ³ คือ ไก้ ไหม้ โหด โหญ่ และ โคร(โคร) โดยจะเปลี่ยน ใ เป็น ัย หลัง

¹ กาชัย ทองหล่อ, หลักภาษาไทย, พิมพ์ครั้งที่ 10 (กรุงเทพมหานคร : อมรการพิมพ์, 2540).

² S. Sethaputra, Thaisoft So Sethaputra Dictionary Version 1.5 [Computer Software], Bangkok : Thaisoft Co.,Ltd, 1996.

³ กาชัย ทองหล่อ, หลักภาษาไทย, พิมพ์ครั้งที่ 10 (กรุงเทพมหานคร : อมรการพิมพ์, 2540).

พยัญชนะควบ (หลังจากทำการตัดวรรณยุกต์แล้ว) เช่น โกล์ ใหม่ และใหญ่ จะเปลี่ยนเป็น กถ์ หม้อย และหญ้อย ตามลำดับ นอกเหนือจาก 5 คำนี้แล้วจะเปลี่ยน ใ เป็น ้อย หลังพยัญชนะต้นของสระ ใ- เช่น ใน ไจ และใส จะเปลี่ยนเป็น น้อย จ้อย และถ้อย ตามลำดับ

2. การเปลี่ยนสระ ใ-ย เป็น ้อย เนื่องจากสระ ใ-ย ที่มีใช้ส่วนใหญ่เป็นคำที่มาจากบาลี-สันสกฤต ซึ่งมีใช้อยู่ 7 คำ⁴ คือ ไชย ไทย ไคย หิมพิลาไสย ภาคินัย อภินัย อสงไสย โดยจะเปลี่ยน ใ- เป็น ้อย หลังพยัญชนะต้นตัวอักษร ย มีอยู่แล้ว
3. การเปลี่ยนสระ ใ- ที่ใช้กับอักษรควบจะเปลี่ยน ใ เป็น ้อย หลังพยัญชนะควบ (หลังจากทำการตัดวรรณยุกต์แล้ว) เช่น ไกร ไคร และไทร จะเปลี่ยนเป็น กร้อย คร้อย และทร้อย ตามลำดับ ซึ่งรูปแบบที่กำหนดดังนี้

...ใ<C><C>

เช่น ไกร

...ใ<C><C><IV>...

เช่น ไกถ์แก่ถ้อย

นอกจากนี้แล้วจะถือว่าเป็นสระ ใ- ที่ใช้กับพยัญชนะธรรมดาซึ่งจะเปลี่ยนเป็น ้อย หลังพยัญชนะต้น

- สระ ำ และ ำ้ม ทำการเปลี่ยนสระ ำ เป็น ำ้ม รูปแบบเดียวเพื่อเพิ่มความถูกต้องยิ่งขึ้นในการเข้ารหัสคำ เนื่องจากสระดังกล่าวมีการอ่านออกเสียงที่เหมือนกันและในการถอดอักษรมักจะถูกถอดเป็น -AM เหมือนกัน เช่น คำ และ ำ้ม ต่างก็ถอดอักษรเป็น KAM เป็นต้น
- การันต์และอักษรควบการันต์ ทำการตัดอักษรและอักษรควบที่มีตัวการันต์กำกับออก เนื่องจากการถอดอักษรโดยปกติจะไม่ถอดอักษรที่มีการันต์กำกับ เช่น สิทธิ ถอดเป็น SITH พันธุ์ ถอดเป็น PHAN เป็นต้น การตัดอักษรและอักษรควบการันต์ออกมี 7 กรณี⁵ คือ

⁴ S. Sethaputra, Thaisoft So Sethaputra Dictionary Version 1.5 [Computer Software], Bangkok : Thaisoft Co.,Ltd, 1996.

⁵ คำขวัญ ทองหล่อ, หลักภาษาไทย, พิมพ์ครั้งที่ 10 (กรุงเทพมหานคร : อมรการพิมพ์, 2540).

1. สระเสียงสั้นที่ไม่มีตัวสะกดควบگارันต์ คือ พยัญชนะตามด้วย สระ อี หรือ สระ อุ และگارันต์ โดยจะตัดตัวสะกด สระ และگارันต์ทิ้ง ตัวอย่างเช่น คำว่า สิทธิ และพันธุ์ จะตัดเป็น สิทธิ และ พัน ตามลำดับ
2. พยัญชนะตัวเดียว คือ มีพยัญชนะหนึ่งตัวและگارันต์ โดยจะตัดพยัญชนะหนึ่งตัวและگارันต์ ตัวอย่างเช่น คำว่า กรรม ศักดิ์ บุรุษ และทรัพย์ จะถูกตัดเป็น กร ศัก บุร และทรัพย์ ตามลำดับ
3. พยัญชนะสองตัว คือ -จน์ โดยจะตัด จน์ ทิ้ง ตัวอย่างเช่น คำว่า กาญจน์ สิญจน์ และวิญจน์ จะตัดเป็น กาญ สิญ และ วิญ ตามลำดับ
4. พยัญชนะสามตัว คือ -มณ โดยจะตัด มณ ทิ้ง ตัวอย่างเช่น ถักมณ จะตัดเป็น ถัก
5. อักษรควบแท้ จะตัดอักษรควบทั้งสองตัวและگارันต์ ตัวอย่างเช่น พักตร์ และ อินทร์ จะตัดเป็น พัก และอิน ตามลำดับ
6. อักษรควบไม่แท้ จะตัดพยัญชนะหนึ่งตัวและگارันต์ ตัวอย่างเช่น คำว่า ชรรม์ สรรพ์ และหรรษ์ จะถูกตัดเป็น ชรร สรร และหรร ตามลำดับ
7. อักษรนำ จะตัดพยัญชนะทั้งสองตัวและگارันต์ ตัวอย่างเช่น คำว่า ถักขณ จะถูกตัดเป็น ถัก

- สระ ฤ และ ฌ ทำการเปลี่ยน ฤ เป็น รี และ ฌ เป็น รือ เช่น พฤษา เปลี่ยนเป็น พริกษา ฤษี เปลี่ยนเป็น รือษี เป็นต้น ส่วน ฤ และ ฌ ไม่นำมาพิจารณา เนื่องจากปัจจุบันไม่มีที่ใช้แล้ว^๑

หมายเหตุ ทุกกรณีที่กำลังมานั้นจะต้องทำงานตามลำดับก่อนหลัง เพื่อลดการตรวจสอบซ้ำซ้อน เช่น จะทำการวรรณยุกต์ และไม่ได้คู่ก่อน และในส่วนของ รร จะไม่พิจารณาว่ามีวรรณยุกต์อยู่ เป็นต้น

หลังจากที่นำคำไทยผ่านขั้นตอนวิธีการลดรูปและตัดวรรณยุกต์เพื่อเปลี่ยนคำไทยให้อยู่ในรูปแบบที่ง่ายต่อการประมวลผลแล้ว จะทำการแทนที่สระสระประสมด้วยสัญลักษณ์เสียงสากล ซึ่งจะกล่าวในหัวข้อถัดไป

^๑ กาชัย ทองหล่อ, หลักภาษาไทย, พิมพ์ครั้งที่ 10 (กรุงเทพมหานคร : อมรการพิมพ์, 2540).

4.2.1.2 การแทนที่สระประสมด้วยสัญลักษณ์เสียงสากล

การแทนที่สระประสมด้วยสัญลักษณ์เสียงสากลเป็นการประมวลผลตัวอักษรเบื้องต้นส่วนหนึ่ง เพื่อให้คำในภาษาไทยที่มีการใช้สระประสมและสระเดี่ยวที่ใช้อักขระตั้งแต่สองตัวขึ้นไป อยู่รูปแบบที่ง่ายต่อการประมวลผล และลดจำนวนอักขระที่ต้องเปรียบเทียบ ซึ่งขั้นตอนนี้จะใช้สัญลักษณ์เสียงสากลแทนเสียงสระดังกล่าวโดยจะวางสัญลักษณ์เสียงสากลไว้หลังพยัญชนะต้นเพื่อให้มีรูปแบบเหมือนกับผลลัพธ์ที่ได้จากการถอดอักษรคำ การใช้สัญลักษณ์เสียงสากลแทนที่สระมีดังนี้

- ใช้สัญลักษณ์เสียง e แทนสระ เ-ะ เช่น เกะ เปลี่ยนเป็น ge
- ใช้สัญลักษณ์เสียง x แทนสระ แ-ะ แ-ะ เช่น แกะ แหทะ เปลี่ยนเป็น gx หลx ตามลำดับ
- ใช้สัญลักษณ์เสียง q แทนสระ เ-ิ เ-ิ เช่น เกิด เจริญ เปลี่ยนเป็น gqd ผจqu ตามลำดับ
- ใช้สัญลักษณ์เสียง I (คัดแปลงจาก ih) แทนสระ เ-ิอะ เ-ิอะ เ-ิย และ เ-ิย เช่น เสียง เกวียน เปลี่ยนเป็น sig กวเิน ตามลำดับ
- ใช้สัญลักษณ์เสียง U (คัดแปลงจาก uh) แทนสระ เ-ือะ เ-ือ เ-ือ ัวะ และ ัว เช่น เรือง เกถือ ัวะ และ ัว เปลี่ยนเป็น rbg glU ผU และ wU ตามลำดับ
- ใช้สัญลักษณ์เสียง @ แทนสระ เ-า เ-า เ-าะ และ เ-าะ เช่น เา เา เาะ และ เาะ เปลี่ยนเป็น g@ พร@ ง@ และ ฉพ@ ตามลำดับ

หมายเหตุ สามารถดูรายละเอียดการใช้สัญลักษณ์เสียงสากลสำหรับอักขระไทยได้ใน

ภาคผนวก ข

4.2.1.3 การถอดอักษร

ในกรณีที่ข้อความเป็นภาษาอังกฤษทับศัพท์ภาษาไทยจะทำการถอดอักษรอังกฤษเป็นภาษาไทย การถอดอักษรที่นำเสนอในที่นี้ จะเป็นการเปลี่ยนพยัญชนะอังกฤษเป็นพยัญชนะไทย ส่วนสระอังกฤษจะใช้หลักเกณฑ์เหมือนกับการตรวจหาสระประสม คือจะถอดสระอังกฤษเป็นสระไทย แต่ถ้าสระไทยนั้นเป็นสระประสมหรือสระเดี่ยวที่ใช้อักษรตั้งแต่สองตัวขึ้นไปจะใช้สัญลักษณ์เสียงสากล (หนึ่งตัว) แทนเสียงสระดังกล่าว

หลักเกณฑ์ในการถอดในส่วนของพยัญชนะจะใช้หลักการเทียบตัวอักษรโรมัน-ไทย ของ ISO⁷ ดังตารางที่ 4.1 ยกเว้นอักษรบางตัวที่ ISO ไม่ได้เทียบให้ ได้แก่อักษร G Q และ X หรือบางตัวที่เทียบให้แล้วแต่ไม่นิยมใช้ เช่น จ เท่ากับ C ซึ่งส่วนใหญ่นิยมใช้ CH หรือ J สำหรับ จ⁸ ดังนั้นผู้วิจัยจึงนำเสนอการถอดอักษรอังกฤษเป็นอักษรไทยในส่วนของพยัญชนะ ดังที่แสดงในตารางที่ 4.2 ซึ่งได้ทำการปรับเปลี่ยนและเพิ่มเติมบางตัวจากต้นแบบของ ISO (ที่มีเครื่องหมาย * ในตาราง) เพื่อให้สมบูรณ์ยิ่งขึ้น เช่น ถอด DH เป็น ท และ ถอด BH เป็น พ เนื่องจากเป็นการเขียนแบบบาลีสันสกฤต ซึ่งยังมีใช้กันอยู่มากในปัจจุบัน⁹

⁷ International Organization for Standardization, "Information and Documentation – Transliteration of Thai", Draft International Standard ISO/DIS 11940, 1996.

⁸ จันทรเพ็ญ ไหวหารสุนทร, "การศึกษาการใช้อักษรโรมันแทนอักษรไทย" (ปริญาณานิพนธ์ คณะอักษรศาสตร์ มหาวิทยาลัยศรีนครินทรวิโรฒ ประสานมิตร, 2530).

⁹ พจนต์ ทิมเจริญ, "การเขียนชื่อภาษาไทยด้วยอักษรโรมัน," วารสารแผนที่ ปีที่ 27 ฉบับที่ 2 (ตุลาคม-ธันวาคม 2527) : 61-74.

พยัญชนะไทย	อักษรโรมัน	พยัญชนะไทย	อักษรโรมัน
ก	K	ท	TH
ข	KH	น	N
ฃ	KH	บ	B
ค	KH	ป	P
ก	KH	ฝ	PH
ฅ	KH	ฝ	F
ฉ	KH	พ	PH
ง	NG	ฟ	F
ด	C	ภ	PH
ด	CH	ม	M
ช	CH	ย	Y
ซ	S	ร	R
ฌ	CH	ร	R
ญ	Y	ล	L
ฎ	D	ฬ	L
ฏ	T	ว	W
ฐ	TH	ศ	S
ฑ	TH	ษ	S
ฒ	TH	ส	S
ณ	N	ห	H
ด	D	ฬ	L
ต	T	อ	X
ถ	TH	ฮ	H
ท	TH		

ตารางที่ 4.1 การใช้อักษรโรมันแทนพยัญชนะไทยของ ISO

อักษรอังกฤษ	อักษรไทย	หมายเหตุ	อักษรอังกฤษ	อักษรไทย	หมายเหตุ
B	บ		N	น (ณ)	
BH	พ	*	NG	ง	
C	ช	*	P	ป	
CH	ช (ณ ณ)		PH	พ (ผ ก)	
CK	ก	*	Q	ค	*
D	ด (ฎ)		R	ร (ฤ)	
DH	ท	*	S	ส (ซ ศ ษ)	
F	ฟ (ฝ)		T	ต (ฏ)	
G	ก	*	TH	ท (ฐ ฑ ฒ ถ ฑ)	
H	ห (ฮ)		V	ว	
J	จ	*	W	ว	
K	ก		X	ก	*
KH	ข (ข ค ฅ ฌ)		Y	ย (ญ)	
L	ล (ฬ พ)		Z	ซ	*
M	ม				

* ส่วนที่ปรับเปลี่ยนและเพิ่มเติมจากแบบของ ISO

ตารางที่ 4.2 การถอดอักษรอังกฤษเป็นอักษรไทยในส่วนของพยัญชนะที่นำสนอ

ส่วนหลักเกณฑ์ในการถอดอักษรในส่วนสระอังกฤษเป็นสระไทยนั้น ผู้วิจัยพบว่ามีปัญหาอย่างมากในการถอดอักษรคือ หนึ่งหน่วยอักษรโรมันสามารถถอดได้เป็นหลายหน่วยอักษรไทย เช่น A ถอดอักษรเป็น -เ- -แ- -ะ -า และ O ถอดอักษรเป็น -อ -เ- เป็นต้น และหลายหน่วยอักษรโรมันสามารถถอดเป็นหนึ่งหน่วยอักษรไทย เช่น U หรือ OO ถอดเป็น - , เป็นต้น และจากการศึกษาของผู้วิจัยพบว่ามีความหลากหลายในการใช้อักษรโรมันแทนอักษรไทย เช่น คำว่า พร มีการเขียนเป็น Phorn Phon Porn Pon เป็นต้น ซึ่งการใช้ om เป็นความนิยมใช้ซึ่งไม่ถูกต้องตามหลักภาษาที่ราชบัณฑิตยสถานกำหนด ดังนั้นในงานวิจัยนี้จึงได้พยายามชี้หลักการถอดอักษรของราชบัณฑิตยสถาน¹⁰ และ การใช้อักษรโรมันแทนอักษรไทยของจันทร์เพ็ญ ไหวหารสุนทร¹¹ (ได้จาก

¹⁰ ราชบัณฑิตยสถาน, "ประกาศราชบัณฑิตยสถาน เรื่อง การถอดอักษรไทยเป็นโรมัน", 2482.

¹¹ จันทร์เพ็ญ ไหวหารสุนทร, "การศึกษาการใช้อักษรโรมันแทนอักษรไทย" (ปริญญาโท คณะอักษรศาสตร์ มหาวิทยาลัยศรีนครินทรวิโรฒ ประสานมิตร, 2530).

การสำรวจความนิยมในการใช้อักษรโรมันแทนอักษรไทย) เป็นต้นแบบและเพิ่มเติมบางส่วนจากการศึกษาของผู้วิจัยที่มีผู้นิยมใช้เข้าไป (ที่มีเครื่องหมาย * ในตาราง) เพื่อความถูกต้องในการใช้งาน ดังที่แสดงในตารางที่ 4.3

ตัวอักษร	ถอดอักษรเป็น	หมายเหตุ
-A	= ะ	
-AA	= ำ	*
-AE	= x (เ-ะ เ-)	
-AI	= ัย	
-AO	= @ (เ-)	
-AIU	= I (เ-ัย)	*
-ARN	= ำน	*
-ART	= ำท	*
-E	= ี (เ-ะ เ-)	
-EE	= ี	*
-EO	= ีว	
-ER	= q (เ-อ เ-) หลัง R ต้องไม่เป็นสระ	
-EU	= ี	
-I	= ิ	
-IA	= I (เ-ยะ เ-ย)	
-IE	= I (เ-ยะ เ-ย)	*
-O	= อ (-อ)	
-OE	= q (เ-อ เ-)	
-OI	= อย	
-OO	= ุ	
-ORN	= ร (-อน)	*
-U	= ู (เ-อ เ-อ เ-อ เ-อ)	
-UA	= U (เ-อะ เ-อ เ-วะ เ-ว)	
-UE	= ุ	

* ส่วนที่เพิ่มเติมจากแบบของราชบัณฑิตยสถาน

ตารางที่ 4.3 การถอดอักษรอังกฤษเป็นอักษรไทยในส่วนของสระที่นำสนอ

หมายเหตุ ถ้าตัวอักษรแรกของคำเป็นสระได้แก่ AEIO และ U ให้ถอด A เป็น อ ถอด E เป็น เอ (สระ เ- กับ อ) ถอด I เป็น อิ (อ กับสระ -) ถอด O เป็น โอ (สระ เ- กับ อ) และ ถอด U เป็น อุ (อ กับสระ -) และถ้าตัวอักษรถัดไปเป็นสระอีกให้นำอักษรตัวแรกไปรวมด้วยในการถอดอักษรโดยเทียบตามตารางที่ 4.3

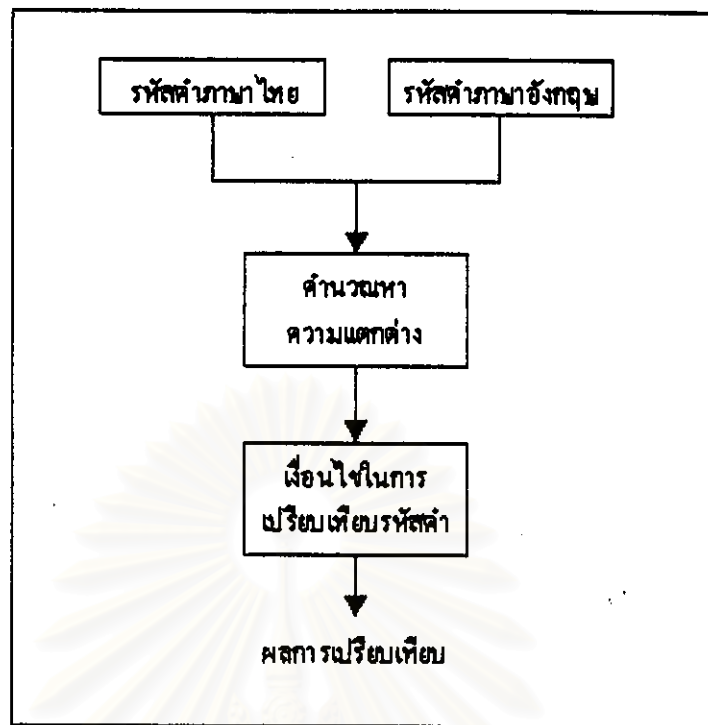
4.2.2 การย้ายตำแหน่งสระ

หลังจากคำภาษาอังกฤษทับศัพท์ภาษาไทยและคำภาษาไทยที่ได้ผ่านขั้นตอนประมวลผลตัวอักษรเบื้องต้นแล้ว คำศัพท์ที่ได้จะอยู่ในรูปแบบเดียวกันคือ อักษรไทยผสมกับสัญลักษณ์เสียงสากล (แทนเสียงสระ) ขั้นตอนสุดท้ายของการเข้ารหัสคำคือ การย้ายตำแหน่งสระ โดยจะทำการย้ายสระไปข้างหลังของสายอักขระตามลำดับ เช่น อุดมศักดิ์ เข้ารหัสเป็น อุดมศก.~ เป็นต้น

4.3 ขั้นตอนวิธีเปรียบเทียบรหัสคำ

ขั้นตอนวิธีในการเปรียบเทียบรหัสคำที่นำเสนอในที่นี้มีขั้นตอนดังนี้ เริ่มจากนำรหัสคำทั้งสองภาษามาทำการคำนวณหาค่าความแตกต่าง (Distance) ของรหัสคำด้วยเทคนิคระยะแก้ไขขั้นที่น้อย (Minimum Edit Distance) และนำค่าความแตกต่างของคู่รหัสคำ มาทดสอบกับเงื่อนไขในการเปรียบเทียบ ถ้าผ่านการทดสอบจะถือว่ารหัสคำทั้งสองรหัสนี้เป็นรหัสที่มาจากคำหลักที่ตรงกันในอีกภาษา ดังรูปที่ 4.3





รูปที่ 4.3 การเปรียบเทียบรหัสคำที่น่าเสนอ

การคำนวณหาค่าความแตกต่าง จะพิจารณาความแตกต่างกันที่เสียงของอักขระไม่ใช่ความแตกต่างที่รูปของอักขระ ส่วนเฟื่อนใจในการเปรียบเทียบรหัสคำนั้นจะนำเอาความยาวของรหัสคำมาร่วมพิจารณากับค่าความแตกต่างของรหัสคำด้วย ซึ่งรายละเอียดของขั้นตอนนั้นจะกล่าวในหัวข้อถัดไป

4.3.1 การคำนวณหาค่าความแตกต่างของรหัสคำ

ผู้วิจัยได้นำเสนอการคำนวณหาค่าความแตกต่างของรหัสคำด้วยเทคนิคระยะแก้ไขเสียงอ่านสั้นที่สุด (Minimum Phonetic Edit Distance) ของ J. Zobel¹² ซึ่งมีขั้นตอนวิธีการคำนวณคล้ายกับเทคนิคระยะแก้ไขสั้นที่สุด คือการคำนวณหาค่าความแตกต่างของคำจะได้มาจากการคำนวณหาต้นทุนน้อยที่สุดในการแก้ไขอักขระให้คำทั้งสองเหมือนกัน แต่เทคนิคระยะแก้ไขเสียงอ่านสั้นที่สุดจะพิจารณาความแตกต่างกันทางเสียงของอักขระแทนรูปของอักขระ โดยใช้กลุ่มอักขระของ

¹² J. Zobel and P. Dart, *Phonetic String Matching: Lessons from Information Retrieval*, Proceedings of the 19th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, Zurich, Switzerland, pp. 166-172, 1996.

ชาวดีเด็กซ์ช่วยในการกำหนดกลุ่มเสียงที่คล้ายกัน การคำนวณค่าความแตกต่างสามารถแสดงในรูปของสมการเวียนเกิด Edit (P_j, W_k) ได้ดังนี้

$$\begin{aligned} \text{Edit}(P_0, W_0) &= 0 \\ \text{Edit}(P_j, W_0) &= \text{Edit}(P_{j-1}, W_0) + D(p_{j-1}, p_j) \\ \text{Edit}(P_0, W_k) &= \text{Edit}(P_0, W_{k-1}) + D(w_{k-1}, w_k) \\ \text{Edit}(P_j, W_k) &= \min[\text{Edit}(P_{j-1}, W_k) + D(p_{j-1}, w_j), \\ &\quad \text{Edit}(P_j, W_{k-1}) + D(p_{j-1}, w_k), \\ &\quad \text{Edit}(P_{j-1}, W_{k-1}) + R(p_j, w_k)] \end{aligned}$$

โดยที่

$P_j = p_1 p_2 p_3 \dots p_j$ เป็นสายอักขระต้นแบบ มีความยาว j ตัวอักษร

$W_k = w_1 w_2 w_3 \dots w_k$ เป็นสายอักขระเป้าหมาย มีความยาว k ตัวอักษร

$R(p_j, w_k)$ เป็นฟังก์ชันที่กำหนดต้นทุนในการแทนที่อักขระ โดยแบ่งออกเป็น 3 กรณี คือ (1) ต้นทุนเท่ากับ 0 ถ้า p_j และ w_k เป็นอักขระที่เหมือนกัน (2) ต้นทุนเท่ากับ 1 ถ้า p_j และ w_k เป็นอักขระที่อยู่กลุ่มเสียงเดียวกันในชาวดีเด็กซ์ (ดูรายละเอียดได้ในตารางที่ 2.2 การกำหนดรหัสชาวดีเด็กซ์) และ (3) นอกจากนี้แล้ว ต้นทุนจะเท่ากับ 2

$D(p_j, w_k)$ เป็นฟังก์ชันที่กำหนดต้นทุนในการเพิ่มหรือลบอักขระ โดยแบ่งออกเป็น 3 กรณี คือ (1) ต้นทุนเท่ากับ 0 ถ้า p_j และ w_k เป็นอักขระที่เหมือนกัน (2) ต้นทุนเท่ากับ 1 ถ้า p_j และ w_k เป็นอักขระที่อยู่กลุ่มที่ไม่ออกเสียง ได้แก่ อักขระ H และ W หรือเป็นอักขระที่อยู่กลุ่มเสียงเดียวกันในชาวดีเด็กซ์ (ตารางที่ 2.2 การกำหนดรหัสชาวดีเด็กซ์) และ (3) นอกจากนี้แล้ว ต้นทุนจะเท่ากับ 2

จากการศึกษาของผู้วิจัยพบว่าระยะแก้ไขเสียงอ่านถิ่นที่สุดของ J. Zobel ยังขาดความเหมาะสมเมื่อนำมาใช้กับภาษาไทย ผู้วิจัยจึงได้ทำการปรับเปลี่ยนฟังก์ชันการกำหนดต้นทุนในการแก้ไขอักขระเพื่อความเหมาะสมเมื่อใช้กับคำในภาษาไทย โดยมีแนวคิดดังนี้

- แบ่งการคำนวณหาค่าความแตกต่างของรหัสคำออกเป็น 2 ส่วนคือ (1) ค่าความแตกต่างของรหัสคำในส่วนที่เป็นพยัญชนะและ (2) ค่าความแตกต่างของรหัสคำกับในส่วนที่เป็นสระ เนื่องจากในการถอดอักษรด้วยกฎที่ใช้ในงานวิจัยนี้พบว่าตำแหน่ง

ของสระที่ได้จากการถอดมัทจะคลาดเคลื่อนไปจากตำแหน่งที่ควรจะเป็น ซึ่งจะทำให้การคำนวณหาค่าความแตกต่างของรหัสคำผิดพลาดด้วย

- กำหนดต้นทุนต่อการแก้ไขอักษรให้มี 4 ระดับ คือ C_1, C_2, C_3 และ C_4 โดยที่ C_1 จะหมายถึงการใช้ต้นทุนน้อยที่สุดในการแก้ไขอักษร และ C_4 จะหมายถึงการใช้ต้นทุนมากที่สุดในการแก้ไขอักษรตามลำดับ ในที่นี้จะกำหนดให้ $C_1=0, C_2=1, C_3=4$ และ $C_4=7$ ซึ่งจะเพิ่มความละเอียดในการคำนวณหาค่าความแตกต่างของรหัสคำ
- กำหนดให้ต้นทุนในการแก้ไขแต่ละอักษรไม่เท่ากัน จากเดิมกำหนดให้การเพิ่ม การลบ และการแทนที่อักษรใด ๆ จะกำหนดต้นทุนในการแก้ไขเท่ากับหนึ่งหรือสองแล้วแต่กรณี ในขณะที่วิธีการที่นำเสนอจะกำหนดต้นทุนในการแก้ไขโดยจะพิจารณาจากความสำคัญของอักษรที่แก้ไข เช่น ต้นทุนในการลบสระ -ะ ควรจะน้อยกว่าต้นทุนการลบพยัญชนะ เนื่องจากภาษาไทยมีการลดรูปสระ -ะ แต่ภาษาอังกฤษทับศัพท์ภาษาไทยไม่มีการลดรูปเสียงสระ -ะ ต้นทุนในการแทนที่อักษรที่อยู่ในมาตราเดียวกันควรจะน้อยกว่าต้นทุนการแทนที่อักษรที่อยู่ต่างมาตรา เป็นต้น
- อนุญาตให้มีการแทนที่อักษรแบบหนึ่งอักษรต่อสองอักษร จากเดิมหนึ่งอักษรต่อหนึ่งอักษร เนื่องจากในภาษาไทยมีลักษณะการใช้อักษรควบ โดยอักษรควบที่ส่งผลต่อการอ่านออกเสียง เช่น การเปลี่ยน “จก” ให้เป็น “จกร” ต้นทุนในการแทนที่ ก ด้วย กร ควรจะน้อยกว่าต้นทุนในการเพิ่ม ร

รายละเอียดต่าง ๆ ในการออกแบบฟังก์ชันกำหนดต้นทุนในการแก้ไขอักษรจะกล่าวในหัวข้อถัดไป

4.3.1.1 การกำหนดต้นทุนในการแก้ไขอักษร

จากการศึกษาของผู้วิจัยพบว่าการใช้กลุ่มเสียงภาษาไทยของชาวเด็กซ์ร่วมกับการกำหนดต้นทุนในการแก้ไขอักษร ทำให้การคำนวณค่าความแตกต่างคลาดเคลื่อนไปจากความเป็นจริงมาก เช่น ไม่มีการจัดการเรื่องอักษรควบ

ผู้วิจัยได้ออกแบบตารางการกำหนดต้นทุนในการแทนที่อักษรโดยกำหนดให้แต่ละอักษรมีต้นทุนที่ต่างกันอย่างแล้วแต่ตัวที่จะไปแทนที่แบบอักษรต่ออักษร โดยมีแนวคิดดังนี้

- พยัญชนะไทยมี 44 รูป สามารถแยกความแตกต่างทางเสียงได้ 21 เสียง ดังนั้นพยัญชนะที่มีเสียงเหมือนกันตามหลักภาษาไทยจะให้ค้ำต้นในการแทนที่อักขระเท่ากับ C_1 ดังตารางที่ 4.4

ก	ฎ ค	ฝ ฟ
ข ขค คข	ฏ ต	ม
ง	ฐ ฑ ฒ ก ท ษ	ร
จ	ณ น	ล ฬ
ฉ ชฌ	บ	ว
ซ ศษ ส	ป	ห ฮ
ญ ย	ผ พ ภ	อ

ตารางที่ 4.4 กลุ่มเสียงของพยัญชนะไทย

- อักขระที่ไม่สามารถแยกความแตกต่างในการถอดอักษร โดยอักษรไทยเหล่านี้ไม่ได้ อยู่กลุ่มเสียงเดียวกันตามตารางที่ 4.4 ซึ่งจะกำหนดค้ำต้นในการแทนที่อักขระเหล่านี้เท่ากับ C_2 ดังตารางที่ 4.5

อักษรโรมัน	อักษรไทย
K	ก ค
T	ต ท
P	ป พ
CH	จ ช
A	-ะ -า
E	เ - เอ
U	ุ - อู

ตารางที่ 4.5 อักขระที่ไม่สามารถแยกความแตกต่างในการถอดอักษร

- อักขระที่ไม่สามารถแยกความแตกต่างทางเสียงในภาษาอังกฤษ เนื่องจากการใช้อักขระอังกฤษอย่างเดียวไม่สามารถแทนเสียงในภาษาไทยได้สมบูรณ์ จะให้ต้นทุนในการแทนที่อักขระเท่ากับ C_1 ตัวอย่างเช่น สระเสียงสั้นกับสระเสียงยาว (A เท่ากับ $-e$ หรือ $-a$)
- มาตรฐานของอักขระไทย จะเลือกบางอักขระมาเป็นตัวแทนมาตรา ดังตาราง 4.6 โดยจะกำหนดต้นทุนในการแทนที่อักขระตัวแทนมาตรากับอักขระที่อยู่ในมาตราเดียวกันเท่ากับ C_2 ซึ่งวิธีนี้จะทำให้ความแตกต่างของเสียงได้มากกว่าวิธีจัดกลุ่มแบบชาวค็เด็กซ์ภาษาไทยที่รวมอักขระต่าง ๆ ที่อยู่มาตราเดียวกันเข้าเป็นกลุ่มเสียงเดียวกัน¹³ เช่น คำว่า “สาร” กับ “ชาญ” วิธีที่นำเสนอสามารถบอกได้ว่า ส-ช และ ร-ญ ต่างกันสิ้นเชิง เนื่องจากทั้ง ส และ ช มีความสัมพันธ์กับ ค หรือ ค (มาตรา กค) แต่ ส และ ช ไม่มีความสัมพันธ์กันเอง แต่ถ้าเป็นวิธีชาวค็เด็กซ์จะบอกว่าอักขร ส และ ช มีความสัมพันธ์กันเนื่องจากอยู่ในมาตราค และ อักขร ร และ ญ มีความสัมพันธ์กันเนื่องจากอยู่ในมาตรากน

มาตรา	ตัวอักษร	ตัวแทนมาตรา
กค	ก ข ค ฅ	ก
กค	ค ด ต ท ษ ฎ ฏ ฐ ฒ อ ฌ ศ ษ ส ซ	ค ค
กง	ง	ง
กน	น ฌ ญ ฎ ฬ ร	น
กบ	บ ป ฝ ภ ฝ	บ
กม	ม	ม

ตารางที่ 4.6 อักขระแทนมาตราต่าง ๆ

- การแทนที่อักขระแบบหนึ่งอักขระต่อสองอักขระ กำหนดให้ต้นทุนในการแทนที่อักขระเท่ากับ C_1 โดยอักขรควบที่ส่งผลต่อการอ่านออกเสียงที่พิจารณาในงานวิจัยนี้คือ¹⁴

¹³ นิลเนตร ทรูวงค์ ฅ ฤชยา, “การเปลี่ยนอักขระของคำในภาษาไทย โดยใช้หลักการของชาวค็เด็กซ์” (ปริญาญานิพนธ์ คณะวิศวกรรมศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย, 2534).

¹⁴ กำชัย ทองหล่อ, หลักภาษาไทย, พิมพ์ครั้งที่ 10 (กรุงเทพมหานคร : อมรการพิมพ์, 2540).

■ อักษรควบไม่แท้ มี 5 ตัว

อักษร	ออกเสียงเป็น	ตัวอย่างเช่น
ทร	ช	ทราบ ทรง
จร	จ	จริง
สว	ส	สรวง เสริม
ศร	ศ	ศรี
ชร	ช	ไจรี

ตารางที่ 4.7 อักษรควบไม่แท้

■ อักษรนำเสียงสนิท มี 9 ตัว

อักษร	ตัวอย่าง	อักษร	ตัวอย่าง
อช	อยู่อยาก	หช	หยก ห่า
หง	เหงือก	หร	หรั่ง หูหระ
หญ	ใหญ่ หยั้ง	หถ	หลวง หล่อ
หน	หนู หน้อย	หว	หวาน หวัง
หม	หมู หมอ		

ตารางที่ 4.8 อักษรนำเสียงสนิท

■ อักษรควบที่เป็นตัวสะกด เช่น

อักษร	ตัวอย่าง	อักษร	ตัวอย่าง
กร	จักร	คร	สมัคร
ดร	จิตร	ทร	สมุทร
ปร	กอปร	รค	ชลมารค
รท	สามารถ	หม	พรหม

ตารางที่ 4.9 อักษรควบที่เป็นตัวสะกด

- ส่วนอื่น ๆ นอกเหนือจากนี้กำหนดคตินิยมในการแทนที่อักษรเท่ากับ C_4 (มากที่สุด)

จากแนวคิดดังกล่าวสามารถสร้างตารางการกำหนดต้นทุนในการแทนที่อักขระในส่วน
ของพยัญชนะและสระได้ดังแสดงเพียงบางส่วนในตารางที่ 4.10 และตารางที่ 4.11 ตามลำดับ
และได้นิยาม REPLACE[a, b] หมายถึงตารางค้นหาต้นทุนในการแทนที่อักขระ b ด้วยอักขระ a
โดยต้นทุนจะคำนวณจากตารางที่ 4.10 และตารางที่ 4.11

	ก	ข	ฃ	ค	...	ท	ธ	น	บ	ป	ผ	ฝ	พ	ฟ	...	พ	อ	ย	กร	ทร	...	หล	หว	อย	
ก	C ₁	C ₂	C ₂	C ₂	...	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	C ₁	C ₁	...	C ₄	C ₄	C ₄	
ข	C ₂	C ₁	C ₁	C ₁	...	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	
ฃ	C ₂	C ₁	C ₁	C ₁	...	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	
ค	C ₂	C ₁	C ₁	C ₁	...	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	
...
ท	C ₄	C ₄	C ₄	C ₄	...	C ₁	C ₁	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	
ธ	C ₄	C ₄	C ₄	C ₄	...	C ₁	C ₁	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	
น	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₁	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₃	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	
บ	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	C ₁	C ₃	C ₄	C ₄	C ₃	C ₃	...	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	
ป	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	C ₃	C ₁	C ₂	C ₄	C ₂	C ₄	...	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	
ผ	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	C ₂	C ₁	C ₄	C ₁	C ₄	C ₄	...	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	
ฝ	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	C ₄	C ₄	C ₁	C ₄	C ₁	C ₄	...	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	
พ	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	C ₃	C ₂	C ₁	C ₄	C ₁	C ₄	...	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	
ฟ	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	C ₃	C ₄	C ₄	C ₁	C ₄	C ₁	...	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	
...
พ	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₃	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₁	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	
อ	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₁	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	
ย	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₁	C ₄	C ₄	...	C ₄	C ₄	C ₄	
กร	C ₁	C ₂	C ₄	C ₄	...	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	C ₁	C ₄	...	C ₄	C ₄	C ₄	
ทร	C ₁	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	C ₄	C ₁	...	C ₄	C ₄	C ₄	
...
หล	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₁	C ₄	C ₄	
หว	C ₂	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₁	C ₄	
อย	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₄	C ₄	C ₄	...	C ₄	C ₄	C ₁	

ตารางที่ 4.10 ตัวอย่างตารางการกำหนดต้นทุนในการแทนที่อักขระที่เป็นพยัญชนะ

หมายเหตุ สามารถดูรายละเอียดตารางการกำหนดต้นทุนในการแทนที่อักขระที่เป็น
พยัญชนะทุกอักขระได้ในภาคผนวก ก

- พัญชนะที่เป็นส่วนประกอบของสระประสม ได้แก่ อ ย และ ว เช่น เ-อ เ-ย และ -ว- เป็นต้น ซึ่งมักจะเกิดความผิดพลาดได้ง่ายในการถอดอักษร เช่น คำว่า Chuan (ชวาน) การถอดอักษรจะได้ เซียน ซึ่งไม่ได้อักษร ว ดังนั้นจะกำหนดต้นทุนในการเพิ่มหรือลบอักษรกลุ่มนี้เท่ากับ C_2
- การทับศัพท์ที่ไม่ตรงตามประกาศของราชบัณฑิตยสถาน เช่น คำว่า ธรรม-Thama จิตร-Chitara และ รัตน-Ratana เป็นต้น จากการทับศัพท์แบบนี้ทำให้เวลาถอดอักษรแล้ว จะได้สระ -ะ หรือ -า เกิน ดังนั้นจะกำหนดต้นทุนในการเพิ่มหรือลบอักษรกลุ่มนี้เท่ากับ C_2

สรุปการกำหนดต้นทุนในการเพิ่มหรือลบอักษระดังตารางที่ 4.12 และได้นิยาม DELADD[a] หมายถึง ตารางค้นหาต้นทุนในการเพิ่มหรือลบอักษระ a โดยคำนวณต้นทุนจากตารางที่ 4.12

อักษระ	ต้นทุน
-ะ	C_2
เ-	C_2
อ	C_2
ย	C_2
ว	C_2
า	C_2

ตารางที่ 4.12 ตารางการกำหนดต้นทุนในการเพิ่มหรือการลบอักษระ

ส่วนฟังก์ชันกำหนดต้นทุนในการแก้ไขอักษระสามารถแบ่งออกเป็น 2 ฟังก์ชันคือ (1) การแทนที่อักษระ และ (2) การเพิ่มหรือการลบอักษระ โดยมีรายละเอียดดังนี้

การแทนที่อักษระ

การกำหนดต้นทุนในการแทนที่อักษระ โดยเปลี่ยนฟังก์ชัน $R(p_i, w_i)$ ในหัวข้อ 4.3.1 เป็น $R(p_i, p_j, w_i)$ เพื่อสนับสนุนการเปรียบเทียบแบบหนึ่งอักษระต่อสองอักษระ (อักษระควบ) ซึ่งฟังก์ชันจะเลือกต้นทุนที่น้อยสุดระหว่างการเปรียบเทียบแบบหนึ่งอักษระต่อหนึ่งอักษระกับหนึ่งอักษระต่อสองอักษระเป็นค่าตอบของฟังก์ชัน โดยมีรายละเอียดดังนี้

$$R(p_{j-1}p_j, w_k) = \min(\text{REPLACE}[p_j, w_k], \text{REPLACE}[p_{j-1}p_j, w_k])$$

โดยที่	$p_{j-1}p_j$	เป็นอักขระสองตัวใด ๆ ที่มีตำแหน่งติดกัน
	w_k	เป็นอักขระใด ๆ
	$\text{REPLACE}[p_j, w_k]$	ตารางค้นหาต้นทุนในการแทนที่อักขระ w_k ด้วย p_j
	$\text{REPLACE}[p_{j-1}p_j, w_k]$	ตารางค้นหาต้นทุนในการแทนที่อักขระ w_k ด้วย $p_{j-1}p_j$

การเพิ่มหรือการลบอักขระ

การกำหนดต้นทุนในการเพิ่มหรือการลบอักขระ (ฟังก์ชัน $D(p_j, w_k)$ ในหัวข้อ 4.3.1) จะเลือกต้นทุนที่น้อยสุดระหว่างต้นทุนในการเพิ่มหรือลบอักขระ และต้นทุนในการแทนที่อักขระ (แทนที่สองอักขระในสายอักขระ P ด้วยหนึ่งอักขระในสายอักขระ W) เป็นค่าตอบของฟังก์ชัน โดยมีรายละเอียดดังนี้

$$D(p_j, w_k) = \min(\text{DELADD}[p_j], \text{REPLACE}[p_{j-1}p_j, w_k])$$

โดยที่	$p_{j-1}p_j$	เป็นอักขระสองตัวใด ๆ ที่มีตำแหน่งติดกัน
	w_k	เป็นอักขระใด ๆ
	$\text{DELADD}[p_j]$	ตารางค้นหาต้นทุนในการเพิ่มหรือลบอักขระ p_j
	$\text{REPLACE}[p_{j-1}p_j, w_k]$	ตารางค้นหาต้นทุนในการแทนที่อักขระ $p_{j-1}p_j$ ด้วย w_k

สรุปการคำนวณค่าความแตกต่างสามารถแสดงในรูปของสมการเวียนเกิด Edit (P_j, W_k) แทนสมการในหัวข้อ 4.3.1 ดังนี้

$$\text{Edit}(P_0, W_0) = 0$$

$$\text{Edit}(P_j, W_0) = \text{Edit}(P_{j-1}, W_0) + D(p_{j-1}p_j)$$

$$\text{Edit}(P_0, W_k) = \text{Edit}(P_0, W_{k-1}) + D(w_{k-1}, w_k)$$

$$\text{Edit}(P_j, W_k) = \min[\text{Edit}(P_{j-1}, W_k) + D(p_{j-1}, w_j),$$

$$\text{Edit}(P_j, W_{k-1}) + D(p_{k-1}, w_k),$$

$$\text{Edit}(P_{j-1}, W_{k-1}) + R(p_{j-1}p_j, w_k)]$$

โดยที่

$P_j = p_1 p_2 p_3 \dots p_j$ เป็นสายอักขระต้นแบบ มีความยาว j ตัวอักษร

$W_k = w_1 w_2 w_3 \dots w_k$ เป็นสายอักขระเป้าหมาย มีความยาว k ตัวอักษร

$p_{j-1} p_j$ เป็นสองอักขระใด ๆ ที่มีตำแหน่งติดกัน

$D(p_j, w_k)$ มีค่าเท่ากับ $\min(\text{DELADD}[p_j], \text{REPLACE}[p_{j-1} p_j, w_k])$

$R(p_{j-1} p_j, w_k)$ มีค่าเท่ากับ $\min(\text{REPLACE}[p_j, w_k], \text{REPLACE}[p_{j-1} p_j, w_k])$

4.3.2 เงื่อนไขในการเปรียบเทียบรหัสคำ

เงื่อนไขที่ใช้ในการทดสอบว่ารหัสคำที่นำมาเปรียบเทียบเป็นรหัสคำที่ได้มากจากคำหลักที่ตรงกันระหว่างภาษาไทยกับภาษาอังกฤษทับศัพท์ภาษาไทยหรือไม่ นั้น มีดังนี้

$$\text{Edit}(P_m^c, W_n^c) \leq \alpha \times \text{Max}(\text{Len}(P_m^c), \text{Len}(W_n^c)) \times C_4 \text{ และ}$$

$$\text{Edit}(P_m^v, W_n^v) \leq \alpha \times \text{Max}(\text{Len}(P_m^v), \text{Len}(W_n^v)) \times C_4$$

โดยที่

- P_m^c คือรหัสคำ P เฉพาะในส่วนของพยัญชนะ มีความยาว m ตัวอักษร
- P_m^v คือรหัสคำ P เฉพาะในส่วนของสระ มีความยาว m ตัวอักษร
- $\text{Edit}(P, W)$ ค่าความแตกต่างของรหัสคำ P กับ รหัสคำ W
- α (แอลฟา) คือพารามิเตอร์ที่ปรับประสิทธิภาพ มีค่าระหว่าง 0 – 1
- $\text{Max}(j, k)$ คือฟังก์ชันที่จะเลือกค่าที่มากที่สุดระหว่าง j กับ k
- $\text{Len}(P)$ ความยาวของรหัสคำ P
- C_4 คือต้นทุนในการแก้ไขอักขระที่มากที่สุด

จากเงื่อนไขดังกล่าวในส่วนของ $\text{Max}(\text{Len}(P), \text{Len}(W)) \times C_4$ หมายถึงต้นทุนมากที่สุดที่ใช้ในการแก้ไขอักขระคือแก้ไขทุกอักขระของ W ให้เท่ากับ P ส่วนค่าของแอลฟานี้เป็นพารามิเตอร์ของระบบที่ผู้ใช้สามารถกำหนดได้ ซึ่งจะส่งผลต่อค่าแม่นยำและค่าเรียกคืนของระบบ ค่าแอลฟามีค่าระหว่าง 0 ถึง 1 ซึ่งจะเป็นตัวกำหนดเกณฑ์การยอมรับความเหมือนกันของรหัสคำแบบ

ประมาณ โดยที่ 0 หมายถึงรหัสคำทั้งสองต้องเหมือนกันทุกประการจึงจะยอมรับ ในขณะที่ 1 หมายถึงการยอมรับทุก ๆ คู่รหัสคำไม่ว่าจะแตกต่างกันเท่าไรก็ตาม

รหัสคำที่ได้มาจากคำหลักภาษาไทยกับภาษาอังกฤษทับศัพท์ภาษาไทยนั้นจะตรงกันก็ต่อเมื่อ สมการเงื่อนไขต้องเป็นจริงทั้งส่วนของพจน์ชนะ และส่วนของชนะ

ตัวอย่าง การทดสอบคำว่า "CHULERTTIYAWONG" และ "ชุลลิตติยวงษ์" เป็นคำทับศัพท์ที่ตรงกันในภาษาไทย-อังกฤษหรือไม่ โดยจะเริ่มค้นจากการเข้ารหัสคำ การคำนวณหาค่าความแตกต่าง และนำค่าความแตกต่างที่ได้ไปทดสอบกับเงื่อนไขเพื่อสรุปว่าเป็นคำทับศัพท์ที่ตรงกันในภาษาไทย-อังกฤษหรือไม่

การเข้ารหัสคำ

- CHULERTTIYAWONG → ชุลลิตติยวงษ์ → ชลคคชวงษ์, q^๖
- ชุลลิตติยวงษ์ → ชุลลิตติยวง → ชุลลิตติยวง → ชลคคชวงษ์, q^๖

การคำนวณหาความแตกต่าง

- ส่วนพจน์ชนะของรหัสคำ $Edit(ชลคคชวงษ์, ชลคคชวงษ์) = 5$
- ส่วนชนะของรหัสคำ $Edit(q^๖, q^๖) = 1$

การทดสอบเงื่อนไขในการเปรียบเทียบรหัสคำ

กำหนดให้ $\alpha = 0.15$

$$C_1 = 7$$

$$Edit(P_m^c, W_n^c) \leq \alpha \times \text{Max}(\text{Len}(P_m^c), \text{Len}(W_n^c)) \times C_1 \text{ และ}$$

$$Edit(P_m^v, W_n^v) \leq \alpha \times \text{Max}(\text{Len}(P_m^v), \text{Len}(W_n^v)) \times C_1$$

$$Edit(ชลคคชวงษ์, ชลคคชวงษ์) \leq 0.15 \times \text{Max}(\text{Len}(ชลคคชวงษ์), \text{Len}(ชลคคชวงษ์)) \times 7 \text{ และ}$$

$$Edit(q^๖, q^๖) \leq 0.15 \times \text{Max}(\text{Len}(q^๖), \text{Len}(q^๖)) \times 7$$

$$4 \leq 0.15 (\text{Max}(8, 7) \times 7) \text{ และ}$$

$$1 \leq 0.15 (\text{Max}(4, 3) \times 7)$$

$$4 \leq 0.15 (8 \times 7) \text{ และ}$$

$$1 \leq 0.15 (4 \times 7)$$

$$4 \leq 8.4 \text{ และ}$$

$$1 \leq 4.2$$

จากตัวอย่างพบว่า $4 \leq 8.4$ และ $1 \leq 4.2$ เป็นจริง เพราะฉะนั้นขั้นตอนวิธีสรุปคำศัพท์ทั้งสองคำเป็นคำทับศัพท์ที่ตรงกันในภาษาไทย-อังกฤษ

4.4 วิธีการทดลอง

ผู้วิจัยได้ทำการทดลองขั้นตอนวิธีที่ได้นำเสนอโดยใช้ชื่อและชื่อสกุลทั้งภาษาไทยและภาษาอังกฤษที่ตรงกันของนิสิตปัจจุบัน ระดับปริญญาโท สาขาวิชาภาษาอังกฤษ จำนวน 5,000 คู่ เป็นตัวอย่างข้อมูลในการทดลอง โดยนำคำศัพท์ทั้งหมดไปทำการเข้ารหัสด้วยขั้นตอนวิธีที่นำเสนอในหัวข้อ 4.2 และจัดเก็บคำศัพท์และรหัสคำในฐานข้อมูล หลังจากนั้นนำคำศัพท์ทั้งหมด 10,000 คำ ไปค้นคืนทีละคำศัพท์กับฐานข้อมูลด้วยขั้นตอนวิธีการเปรียบเทียบรหัสคำที่ได้นำเสนอในหัวข้อ 4.3 เพื่อทำการคำนวณค่าแม่นยำ และค่าเรียกคืนสำหรับการค้นคืนใช้สูตรดังนี้

$$\text{ค่าแม่นยำ} = \frac{\text{จำนวนคำศัพท์ที่เกี่ยวข้องที่คืนกลับมา}}{\text{จำนวนคำศัพท์ที่คืนกลับมา}} \times 100$$

$$\text{ค่าเรียกคืน} = \frac{\text{จำนวนคำศัพท์ที่เกี่ยวข้องที่คืนกลับมา}}{\text{จำนวนคำศัพท์ที่เกี่ยวข้องทั้งหมด}} \times 100$$

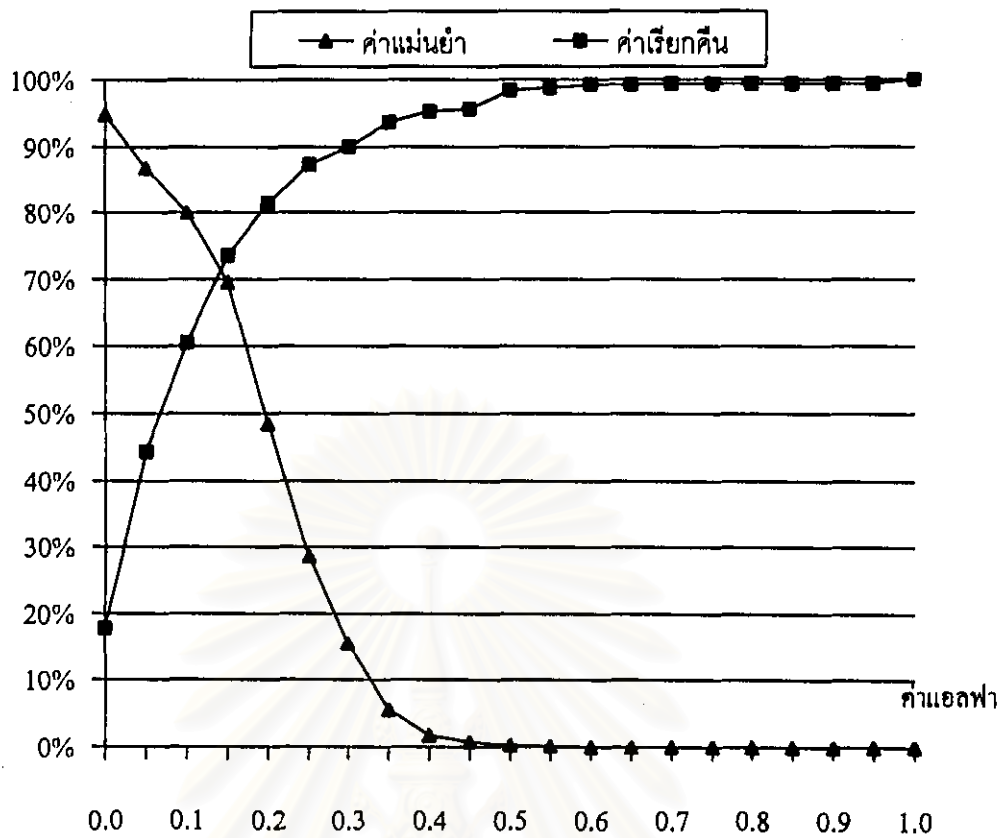
ผู้วิจัยได้ทำการทดลองโดยการใช้ค่าแอกฟาในระดับต่าง ๆ เพื่อหาความสัมพันธ์ระหว่างค่าแอกฟากับประสิทธิภาพของระบบค้นคืน

4.5 ผลการทดลอง

จากการทดลองในหัวข้อ 4.4 เพื่อหาความสัมพันธ์ระหว่างค่าแอลฟากับประสิทธิภาพของระบบคั่นคืนโดยการเปลี่ยนแปลงค่าแอลฟาทุก ๆ 0.05 จาก 0.00 ถึง 1.00 ได้ผลแสดงในตารางที่ 4.13 และรูปที่ 4.4

ค่าแอลฟา	ค่าแม่นยำ	ค่าเรียกคืน
0.00	0.947752	0.177200
0.05	0.867973	0.442200
0.10	0.800750	0.605800
0.15	0.694093	0.734800
0.20	0.484050	0.814800
0.25	0.286261	0.872800
0.30	0.154865	0.899000
0.35	0.056202	0.936200
0.40	0.018329	0.952800
0.45	0.007337	0.955800
0.50	0.002760	0.984000
0.55	0.001368	0.986800
0.60	0.000750	0.992200
0.65	0.000517	0.992200
0.70	0.000348	0.993600
0.75	0.000271	0.993600
0.80	0.000235	0.993600
0.85	0.000222	0.993600
0.90	0.000214	0.993600
0.95	0.000211	0.993600
1.00	0.000200	1.000000

ตารางที่ 4.13 ความสัมพันธ์ระหว่างค่าแอลฟากับประสิทธิภาพของระบบคั่นคืน



รูปที่ 4.4 ความสัมพันธ์ระหว่างค่าแอลฟากับประสิทธิผลของระบบคั่นคืน

จากรูปที่ 4.4 แสดงให้เห็นว่าค่าแม่นยำของระบบคั่นคืนสูงประมาณ 94 เปอร์เซ็นต์ และจะลดต่ำลงอย่างต่อเนื่อง เมื่อค่าของแอลฟาเพิ่มขึ้น ส่วนพฤติกรรมของค่าเรียกคืนจะเริ่มต้นประมาณ 17 เปอร์เซ็นต์และจะเพิ่มค่าขึ้นเมื่อค่าแอลฟาเพิ่มขึ้น ประสิทธิภาพของระบบที่ได้จากการทดลองพบว่า สอดคล้องกับพฤติกรรมโดยปกติของค่าแม่นยำและค่าเรียกคืนในระบบคั่นคืนใด ๆ คือแนวโน้มของค่าแม่นยำและค่าเรียกคืนจะตรงข้ามกัน

จากผลการทดลองพบว่าขั้นตอนวิธีการเข้ารหัสคำทับศัพท์เพื่อการคั่นคืนข้ามภาษาไทย-อังกฤษจะเกิดประสิทธิผลสูงสุดคือค่าแม่นยำเท่ากับ 69 เปอร์เซ็นต์และค่าเรียกคืนเท่ากับ 73 เปอร์เซ็นต์ เมื่อกำหนดค่าแอลฟาเท่ากับ 0.15

4.6 สรุป

ในบทนี้ได้กล่าวถึงขั้นตอนวิธีการเข้ารหัสคำทับศัพท์เพื่อการค้นคืนข้ามภาษาไทย-อังกฤษแบบภาษาอังกฤษทับศัพท์ภาษาไทย การค้นคืนข้ามภาษานั้นสามารถทำได้โดยนำข้อคำถามไปเข้ารหัสคำแล้วนำรหัสคำที่ได้ไปค้นหาในดัชนีคำหลักของเอกสารที่ได้เข้ารหัสไว้แล้วในขั้นตอนการทำดัชนี

ขั้นตอนวิธีที่ได้นำเสนอนี้ ผู้วิจัยได้ใช้หลักภาษาไทยมาช่วยในการออกแบบการประมวลผลตัวอักษรเบื้องต้นเพื่อให้คำศัพท์ทั้งสองภาษามีรูปแบบที่เหมือนกันก่อนที่ไปนำไปเข้ารหัสคำ หลังจากนั้นได้ใช้เทคนิคระยะแก้ไขการอ่านออกเสียงสั้นที่สุดมาดัดแปลงเพื่อให้เหมาะสมกับการอ่านออกเสียงของภาษาไทยเพื่อใช้ในการเปรียบเทียบกับคู่รหัสคำที่ได้ และได้โดยการออกแบบเงื่อนไขที่ใช้ทดสอบว่ารหัสคำทั้งสองว่าเป็นรหัสคำที่ได้มาจากคำทับศัพท์ที่ตรงกันในภาษาไทย-อังกฤษ

ผลการทดลองขั้นตอนวิธีที่นำเสนอ พบว่าค่าแม่นยำสูงถึง 69 เปอร์เซ็นต์ และค่าเรียกคืนสูงถึง 73 เปอร์เซ็นต์ เมื่อกำหนดให้ค่าแอลฟาเท่ากับ 0.15

สถาบันวิทยบริการ
จุฬาลงกรณ์มหาวิทยาลัย