

2. ทฤษฎีและงานวิจัยที่เกี่ยวข้อง

ในบทนี้จะกล่าวถึงทฤษฎีต่าง ๆ ที่เกี่ยวกับหลักภาษาไทย ได้แก่ การถอดอักษร การใช้ อักษรโรมันเพื่อการถ่ายเสียง และหลักเกณฑ์การทับศัพท์ เพื่อใช้เป็นฐานความรู้ในการออกแบบ ขั้นตอนวิธีการเข้ารหัสคำทับศัพท์เพื่อการค้นคืนข้ามภาษาไทย-อังกฤษให้ทำงานได้ถูกต้องตาม หลักภาษา

ส่วนงานวิจัยต่าง ๆ ที่ผ่านมาที่เกี่ยวข้องและมีอิทธิพลต่องานวิจัยนี้ได้แก่ การค้นคืนสารสนเทศข้ามภาษา ขั้นตอนวิธีชาวคัลเดิร์สภาษาอังกฤษ ขั้นตอนวิธีชาวคัลเดิร์สภาษาอังกฤษ และขั้นตอนวิธีระยะแก้ไขสั้นที่สุด

2.1 การถอดอักษร

การถอดอักษร (Transliteration) หมายถึง การนำคำในภาษาหนึ่งมาเขียนด้วยตัวอักษรอีก ภาษาหนึ่งแบบอักษรต่ออักษร โดยพยายามใช้หน่วยเสียงของอักษรทั้งสองภาษาใกล้เคียงกันมากที่สุด¹ ตัวอย่างเช่น คำว่า “HARVARD” ในภาษาอังกฤษถอดอักษรเป็น “ฮาร์วาร์ด” ในภาษาไทย เป็นต้น การถอดอักษรแบ่งเป็น 3 ขั้นตอนหลัก ๆ ดังนี้

1. ถอดหน่วยอักษรในภาษาดั้งเดิม (Source Language) เป็นหน่วยเสียงในภาษาดั้งเดิม เช่น ถอดหน่วยอักษร “B” ในภาษาอังกฤษเป็นหน่วยเสียง /b/ ในภาษาอังกฤษ เป็นต้น
2. แทนหน่วยเสียงในภาษาดั้งเดิม ด้วยหน่วยเสียงในภาษาเป้าหมาย (Target Language) โดยพยายามใช้หน่วยเสียงที่ใกล้เคียงกันมากที่สุด เช่น แทนหน่วยเสียง /b/ ในภาษาอังกฤษเป็นหน่วยเสียง /บ/ ในภาษาไทย เป็นต้น
3. ถอดหน่วยเสียงในภาษาเป้าหมาย เป็นหน่วยอักษรในภาษาเป้าหมาย เช่น ถอดหน่วยเสียง /b/ ในภาษาไทยเป็นหน่วยอักษร “บ” ในภาษาไทย เป็นต้น

¹ อุไรรัตน์ บุญधानนท์, “การถอดอักษรภาษาอังกฤษเป็นไทยโดยใช้หลักวิชาภาษาศาสตร์” (วิทยานิพนธ์ปริญญาโทบัณฑิต ภาควิชาภาษาศาสตร์ บัณฑิตวิทยาลัย จุฬาลงกรณ์มหาวิทยาลัย, 2528), หน้า 6-18.

จากขั้นตอนดังกล่าวสามารถแสดงตัวอย่างการถอดอักษร ได้ดังตารางที่ 2.1

ภาษาดั้งเดิม		ภาษาเป้าหมาย	
หน่วยอักษร	หน่วยเสียง	หน่วยเสียง	หน่วยอักษร
B	/b/	/b/	บ
M	/m/	/m/	ม
T	/t/	/tʰ/	ท

ตารางที่ 2.1 ตัวอย่างการถอดอักษรจากภาษาดั้งเดิมไปยังภาษาเป้าหมาย

ปัญหาต่าง ๆ ในการถอดอักษร

- ความสัมพันธ์ของหน่วยอักษรและหน่วยเสียง มีความสัมพันธ์แบบหนึ่งตัวอักษร แทนหลายหน่วยเสียง เช่น ในภาษาอังกฤษ “C” แทนด้วย /k/ หรือ /s/ เป็นต้น และมีความสัมพันธ์แบบหลายหน่วยอักษรแทนหนึ่งหน่วยเสียง เช่น ในภาษาอังกฤษ “N, TN, GN, PN” แทนด้วย /n/ ในภาษาไทย “ร, ฤ, ทร” แทนด้วย /r/ และ “ฉ, ช, ฌ” แทนด้วย /ch/ เป็นต้น
- การแบ่งพยางค์ในภาษาดั้งเดิม เมื่อมีพยางค์เดียวอยู่ระหว่างสระ เช่น คำว่า money ในภาษาอังกฤษ จะแบ่งพยางค์อย่างไร จะถอดพยางค์ที่สองตัวเพื่อให้อ่านได้สะดวกเป็น มัน-นีย์ หรือจะถอดอักษรเพียงตัวเดียวตามที่ปรากฏในภาษาอังกฤษเป็น มะ-นีย์ หรือ มัน-อีย์
- ปัญหาอันเนื่องมาจากช่วงเวลาของการอิมพอร์ตศัพท์ คำทับศัพท์บางคำอิมมาเป็นเวลานาน ซึ่งในอดีตมีหลักการการทับศัพท์ไม่ตรงกับหลักการในปัจจุบัน เช่น “C” ที่แทน /k/ ในอดีตนิยมถอดเป็นอักษร “ก” เช่น ก๊ก (Cook) กัปตัน (Captain) กระรัต (Carat) แก๊ป (Cap) เป็นต้น แต่ปัจจุบัน “C” ที่แทน /k/ มักจะถอดเป็น “ค” ในตำแหน่งพยางค์ต้น เช่น คอนโดมิเนียม (Condominium) แคปซูล (Capsule) แครีอต (Carrot) เป็นต้น

2.2 การใช้ตัวอักษรโรมันเพื่อการถ่ายเสียง

การใช้ตัวอักษรโรมันเพื่อการถ่ายเสียง (Romanization)² คือการถ่ายเสียงตัวอักษรของภาษาอื่น ๆ ที่ไม่ใช่อักษรโรมัน เช่น ไทย จีน ญี่ปุ่น ฯลฯ ให้เป็นตัวอักษรโรมัน เพื่อให้ผู้ที่ไม่รู้อักขรณานัน ๆ สามารถอ่านออกเสียงได้ เช่น คำว่า “อยุธยา” สามารถถ่ายเสียงเป็น “AYUTTHAYA” เป็นต้น การเขียนชื่อเฉพาะต่าง ๆ จากภาษาไทยด้วยอักษรโรมันนั้นมีปัญหามาก เช่น คำว่า “ไทยท努” เขียนเป็น “Thai Danu” ซึ่งการใช้ตัวอักษร “D” แทนตัว “ท” นั้นเป็นการเขียนโดยอาศัยหลักเกณฑ์ของภาษาบาลี³

จากปัญหาต่าง ๆ ที่เกิดขึ้น ทางราชบัณฑิตยสถานจึงได้พยายามสร้างระบบมาตรฐานในการใช้ตัวอักษรโรมันเพื่อการถ่ายเสียงสำหรับตัวอักษรไทย โดยมีแนวคิดว่าจะยึดหลักภาษาไทยเป็นที่ตั้งและหาวิธีการถ่ายเสียงตัวอักษรตามวิธีเขียนและการออกเสียงในภาษาไทย เพราะถ้ายึดหลักการถ่ายเสียงตัวอักษรตามศัพท์ในภาษาเดิมซึ่งเป็นที่มาของคำไทยแล้ว อักษรไทยตัวเดียวกันอาจถ่ายเสียงเป็นอักษรโรมันต่างกันได้ เช่น ทธ จะถ่ายเสียงเป็น Thon, S, Tr., Tara แต่ถ้าจะถ่ายเสียงตัวอักษรตามเสียงอย่างเดิวนั้น เมื่อเขียนกลับเป็นภาษาไทยก็เขียนได้หลายอย่าง เช่น Ban Ma เขียนเป็น บ้านมา, บ้านม่า, บ้านม้า หรือ บ้านหมาได้ เป็นต้น

ในที่สุด ทางราชบัณฑิตยสถานจึงได้กำหนดระบบการใช้ตัวอักษรโรมันเพื่อการถ่ายเสียงสำหรับตัวอักษรไทยออกเป็น 2 ระบบ⁴ คือระบบทั่วไปและระบบพิศดาร โดยระบบทั่วไปจะใช้สำหรับกรณีที่ต้องการออกเสียงสำคัญกว่าการเขียนตัวสะกด ซึ่งจะอาศัยหลักการออกเสียงเป็นสำคัญ ต้องสอดคล้องกับไวยากรณ์ของไทย และสามารถขยายเป็นระบบเฉพาะได้ เช่น คำว่า “กษัตริย์” ถ่ายเสียงเป็น Kasat ส่วนระบบเฉพาะจะใช้ในกรณีที่แสดงตัวอักษรให้ละเอียดแม่นยำ เพื่อให้คงความหมายของคำนั้นไว้ เช่น คำว่า “กษัตริย์” ถ่ายเสียงเป็น Kasatriy

² อุไรรัตน์ บุญยานนท์, “การถอดอักษรภาษาอังกฤษเป็นไทยโดยใช้หลักวิชาภาษาศาสตร์” (วิทยานิพนธ์ปริญญาโทบัณฑิต ภาควิชาภาษาศาสตร์ บัณฑิตวิทยาลัย จุฬาลงกรณ์มหาวิทยาลัย, 2528), หน้า 6-18.

³ พยงค์ ทิมเจริญ, “การเขียนชื่อภาษาไทยด้วยอักษรโรมัน,” วารสารแพทย ปีที่ 27 ฉบับที่ 2 (ตุลาคม-ธันวาคม 2527) : 61-74.

⁴ อุไรรัตน์ บุญยานนท์, “การถอดอักษรภาษาอังกฤษเป็นไทยโดยใช้หลักวิชาภาษาศาสตร์” (วิทยานิพนธ์ปริญญาโทบัณฑิต ภาควิชาภาษาศาสตร์ บัณฑิตวิทยาลัย จุฬาลงกรณ์มหาวิทยาลัย, 2528), หน้า 6-18.

2.3 หลักเกณฑ์การทับศัพท์

ในปัจจุบันได้มีการบัญญัติศัพท์วิชาการขึ้นใช้กันอย่างแพร่หลาย และในการบัญญัติศัพท์นั้น บางครั้งไม่สามารถหาคำไทยมาใช้ได้ตรงความหมายที่ต้องการ หรือกรณีที่เป็นคำที่ไม่เคยมีการบัญญัติศัพท์มาก่อนจึงต้องใช้วิธีทับศัพท์ นอกจากนี้การเขียนคำวิสามานยนามต่าง ๆ เช่น ชื่อคน ชื่อสถานที่ ก็ต้องใช้วิธีทับศัพท์เช่นเดียวกัน เพื่อให้เป็นมาตรฐานเดียวกันในการทับศัพท์⁵ ทางราชบัณฑิตยสถานจึงได้กำหนดหลักเกณฑ์การทับศัพท์ไว้ดังนี้

1. การทับศัพท์ให้ถอดอักษรในภาษาเดิมพอควรแก่การแสดงที่มาของรูปศัพท์ และให้เขียนในรูปที่อ่านได้สะดวกในภาษาไทย
2. การวางหลักเกณฑ์ได้แยกกำหนดหลักเกณฑ์การทับศัพท์ภาษาต่าง ๆ แต่ละภาษาไป เช่น หลักเกณฑ์การทับศัพท์สำหรับภาษาอังกฤษ หลักเกณฑ์การทับศัพท์สำหรับภาษาญี่ปุ่น เป็นต้น
3. คำทับศัพท์ที่ใช้กันมานานจนถือเป็นคำไทย และปรากฏในพจนานุกรมฉบับราชบัณฑิตยสถานแล้ว ให้ใช้ต่อไปตามเดิม เช่น ช็อกโกแลต ช็อกโกแลต เช็ค แก๊ส ก๊าซ
4. ศัพท์วิชาการซึ่งใช้เฉพาะกลุ่มไม่ใช้ศัพท์ทั่วไปอาจเพิ่มหลักเกณฑ์ขึ้นตามความจำเป็น

2.4 การค้นคืนข้ามภาษา

การค้นคืนสารสนเทศข้ามภาษา หมายถึง การค้นคืนสารสนเทศ โดยภาษาที่ใช้ในข้อคำถามแตกต่างจากภาษาที่ใช้ในการจัดเก็บเอกสารแล้วระบบจะค้นคืนสารสนเทศข้ามภาษาให้โดยอัตโนมัติ โดยทั่วไปแล้วสามารถแบ่งการค้นคืนสารสนเทศข้ามภาษาออกเป็น 2 วิธี⁶ คือ

1. การแปลข้อคำถาม คือ การแปลภาษาของข้อคำถามที่ผู้ใช้ระบุขณะที่ทำการสืบค้น ให้เป็นภาษาเดียวกับที่จัดเก็บในเอกสารก่อน แล้วจึงทำการสืบค้น
2. การแปลเอกสาร คือ การแปลภาษาของเอกสารทั้งหมดให้เป็นภาษาเดียวกับข้อคำถามก่อนแล้วจึงทำการสร้างดัชนี

⁵ ราชบัณฑิตยสถาน, หลักเกณฑ์การทับศัพท์ ฉบับราชบัณฑิตยสถาน (2535).

⁶ D. Oard and B. Dorr, A Survey of Multilingual Text Retrieval. Technical Report UMIACS-TR-96-19 CD-TR-3615, University of Maryland, College Park, April 1996.

งานวิจัยทางการค้นคืนสารสนเทศข้ามภาษาที่ผ่านมานั้น สามารถแบ่งออกเป็น 3 แนวคิดหลัก ๆ⁷ คือ

1. การแปลข้อความโดยใช้เทคนิคการแปลภาษาด้วยเครื่อง (Machine Translation)⁸ คือ การแปลภาษาที่เอกสารหรือข้อความด้วยเครื่องแปลภาษา ซึ่งข้อดีของวิธีการนี้ขึ้นอยู่กับคุณภาพของการแปลภาษา

2. การใช้ฐานพจนานุกรมหลายภาษา (Multilingual Dictionary)⁹ คือ การใช้พจนานุกรมหลายภาษาในลักษณะของบัญชีคำอรรถาภิธาน หรือการสร้างบัญชีคำพ้องความหมายแต่เป็นการพ้องความหมายข้ามภาษา เช่น คำว่า "Food" จะพ้องความหมายกับคำว่า "อาหาร" เป็นต้น

3. การใช้ฐานคำแบบขนาน (Parallel Corpus-Based)¹⁰ เป็นอีกทางเลือกของการใช้บัญชีคำอรรถาภิธาน คือการใช้ประโยชน์จากฐานข้อมูลที่มีอยู่แล้ว ซึ่งจัดเก็บในลักษณะสองภาษาควบคู่กันไป ตัวอย่างเช่น ข้อมูลหนังสือที่มีการจัดเก็บชื่อเรื่องที่เป็นภาษาไทยและภาษาอังกฤษควบคู่กันไป เป็นต้น โดยจะนำข้อมูลแต่ละคู่มาจัดเก็บเป็นฐานคำแบบขนาน

2.5 ขั้นตอนวิธีชาวเด็กรหัสภาษาอังกฤษ

M. K. Odell และ R. C. Russell ได้ออกแบบขั้นตอนวิธีการเข้ารหัสชื่อในภาษาอังกฤษโดยยึดหลักของการอ่านออกเสียง เพื่อให้ชื่อที่อ่านออกเสียงคล้ายกันได้รับรหัสเหมือนกัน หรือที่เรียกว่า "ชาวเด็กรหัส" (Soundex) ขั้นตอนวิธีดังกล่าวได้ใช้แนวคิดทางภาษาศาสตร์และชวเลขที่ว่าชื่อในภาษาอังกฤษสามารถจำแนกความแตกต่างได้โดยพิจารณาเพียงพยัญชนะเท่านั้น¹¹

⁷ D. Oard and B. Dorr, A Survey of Multilingual Text Retrieval, Technical Report UMIACS-TR-96-19 CID-TR-3615, University of Maryland, College Park, April 1996.

⁸ B. R. Pevzner, Automatic Translation of English Text to the Language of the Pusto-Nepusto-2 System, *Automatic Documentation and Mathematical Linguistics*, 3(40):40-48, 1969.

⁹ F. Semturs, STAIRS/TLS-A System for Free Text and Descriptor Searching, *Proceedings of the ASIS Annual Meeting*, Volume 15, pp. 295-298, November 1978.

¹⁰ G. Salton, Experiments in Multi-Lingual Information Retrieval, *Information Processing Letters*, 2(1):6-11, 1973.

¹¹ A. Binstock and J. Rex, Practical Algorithms for Programmers (New York: Addison Wesley, 1995), pp. 157-172.

```

char *SOUNDEX(char *Name)
{
    /*ABCDEFGHIJKLMNPOQRSTUVWXYZ*/
    char Table[] = "01230120022455012623010202";
    char Code[] = "0000";
    int Count = 0;
    char Ch;
    /*----- For the First Character -----*/
    Code[Count++] = Name[0];
    /*----- For the Rest Character -----*/
    for (i=2; i < strlen (Name); i++) {
        Ch = Table[Name[i] - 'A'];
        if (Ch != PrevCode && Ch != '0') {
            Code[Count++] := Ch;
            if (Count = 5)
                return(Code);
        }
        PrevCode := Ch;
    }
    return(Code);
}

```

รูปที่ 2.1 โปรแกรมการเข้ารหัสชาวเด็กซ์ภาษาอังกฤษ

จากรูปที่ 2.1 แสดงขั้นตอนการเข้ารหัสชาวเด็กซ์ โดยเริ่มจากการนำตัวอักษรตัวแรกของคำไปเป็นรหัส ส่วนตัวอักษรที่เหลือจะแปลงเป็นตัวเลขโดยใช้ตารางการกำหนดรหัสชาวเด็กซ์ ดังแสดงในตารางที่ 2.2 จากนั้นจะตัดรหัสตัวเลขศูนย์ออกไป และถ้ารหัสตัวเลขที่อยู่ตำแหน่งติดกันมีค่าเท่ากันจะเก็บเพียงหนึ่งรหัสเท่านั้น ชุดทำรหัสนี้ที่ได้คือตัวอักษรตัวแรกของชื่อตามด้วยรหัสตัวเลขตามตัวแรกที่ได้จากการแปลง ถ้ารหัสที่ได้มีความยาวไม่ถึงสี่หลักจะเติมตัวเลขศูนย์จนครบสี่หลัก ตัวอย่างเช่น คำว่า "ALEXANDER" มีรหัสชาวเด็กซ์เท่ากับ A425 เป็นต้น

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

ตารางที่ 2.2 การกำหนดรหัสชาวคเด็กซ์ภาษาอังกฤษ

ขั้นตอนวิธีชาวคเด็กซ์ดังกล่าวเป็นขั้นตอนวิธีที่ง่ายและทำงานได้รวดเร็ว ส่วนในเรื่องของประสิทธิภาพในการค้นคืนพบว่าเปอร์เซ็นต์ของคำเรียกคืนสูงมาก แต่เปอร์เซ็นต์ของคำแม่นยำต่ำ

2.6 ขั้นตอนวิธีชาวคเด็กซ์ภาษาไทย

ปัจจุบันมีงานวิจัยทางด้านชาวคเด็กซ์ภาษาไทย เพื่อแก้ไขปัญหาต่าง ๆ เช่น การสืบค้นชื่อและนามสกุลที่อ่านออกเสียงเหมือนกันในทะเบียนรายชื่อ การสืบค้นคำไทยที่มักสะกดผิด และการตรวจสอบตัวสะกด เป็นต้น การออกแบบและพัฒนาขั้นตอนวิธีเพื่อแก้ปัญหาดังกล่าวจึงแตกต่างกันไป ในที่นี้จะขอกกล่าวถึงงานวิจัยของวรรณิ อุดมพานิชย์¹² งานวิจัยของนิลเนตร อรุณวงศ์ ณ อยุธา¹³ และงานวิจัยของ Theppitak Karoonboonyanan¹⁴

¹² วรรณิ อุดมพานิชย์, "การใช้หลักคำห้องเสียง เพื่อค้นหาชุดอักษรภาษาไทยที่ออกเสียงเหมือนกัน" (วิทยานิพนธ์ปริญญาโทบริหาร ภาควิชาวิศวกรรมคอมพิวเตอร์ บัณฑิตวิทยาลัย จุฬาลงกรณ์ มหาวิทยาลัย, 2526), หน้า 8-19.

¹³ นิลเนตร อรุณวงศ์ ณ อยุธา, "การเปลี่ยนอักขระของคำในภาษาไทย โดยใช้หลักการของชาวคเด็กซ์" (ปริญญาโทบริหาร คณะวิศวกรรมศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย, 2534).

¹⁴ T. Karoonboonyanan, V. Somtertlamvanich and S. Meknavin, A 'Thai Soundex System for Spelling Correction, Proc. Of the Natural Language Processing Pacific Rim Symposium 1997, Phuket, Thailand, December 2-4, pp.633-636.

งานวิจัยของวรรณิ อุดมพานิชย์ มีการออกแบบขั้นตอนวิธีการเข้ารหัสชาวคเด็กซ์ภาษาไทยเพื่อให้ชื่อที่อ่านออกเสียงเหมือนกันได้รหัสตรงกัน เพื่อแก้ปัญหาในการสืบค้นชื่อและนามสกุลที่อ่านออกเสียงเหมือนกันในทะเบียนรายชื่อ หลักเกณฑ์และข้อกำหนดของขั้นตอนวิธีดัดแปลงมาจากชาวคเด็กซ์ภาษาอังกฤษของ Odell และ Russell ซึ่งได้เพิ่มเติมข้อกำหนดบางส่วนเพื่อเหมาะสมกับภาษาไทย หลักเกณฑ์การสร้างรหัสชาวคเด็กซ์มีดังนี้

- กำหนดให้ขนาดของรหัสมีความยาวเท่ากับ 7 หลัก ประกอบด้วยตัวอักษร 1 หลักและตัวเลข 6 หลัก โดยมีรูปแบบของรหัสชาวคเด็กซ์ดังนี้

$$(\text{รหัสชาวคเด็กซ์}) = (\text{ตัวอักษร})(\text{ตัวเลข})(\text{ตัวเลข})(\text{ตัวเลข})(\text{ตัวเลข})(\text{ตัวเลข})(\text{ตัวเลข})$$
- เปลี่ยนพยัญชนะตัวแรกของคำเป็นรหัสตัวอักษรตามตารางที่ 2.3
- เปลี่ยนอักษรตั้งแต่ตัวที่สองขึ้นไปของคำเป็นรหัสตัวเลขตามตารางที่ 2.4
- โดยปกติจะไม่นำสระ วรรณยุกต์ และไม้ได้คู่ มาสร้างรหัส ยกเว้นสระที่ทำให้เสียงสะกดเป็นพยัญชนะ ได้แก่ สระ -า ใ ใ
- ถ้ารหัสชาวคเด็กซ์ที่ได้มีความยาวน้อยกว่า 7 หลักให้เติม 0 ต่อไปจนรหัสมีความยาวครบ 7 หลัก

รหัสตัวอักษร	ตัวอักษร
ก	ก
ข	ข ช ค คม
ง	ง
จ	จ
ช	ช ฉ ฉ
ต	ต ศ ส
ด	ด ฎ
ค	ค ฎ
ท	ท ฑ ฒ ก ท ษ
น	ณ น

รหัสตัวอักษร	ตัวอักษร
บ	บ
ป	ป
พ	พ ภ ผ
ฟ	ฝ ฟ
ม	ม
ย	ย อย
ร	ร ล ฬ ฤ ฦ
ว	ว
อ	อ
ฮ	ฮ อ

ตารางที่ 2.3 การกำหนดรหัสตัวอักษรของรหัสชาวคเด็กซ์ภาษาไทย

รหัสตัวเลข	ตัวอักษร
0	ม ว ำ
1	ก ข ฉ ค ด ช
2	ง อ
3	ญ ณ
4	ฎ ฏ ค ต ศ ช ส
5	บ ป พ ก
6	ฝ ฟ ฟ ห อ ฮ
7	จ ฉ ช ซ ฌ
8	ฐ ฑ ฒ ณท ษ
9	ร ฤ ล ฬ

ตารางที่ 2.4 การกำหนดรหัสตัวเลขของรหัสชาวค้เด็กซ์ภาษาไทย

ส่วนข้อกำหนดเพื่อให้เหมาะสมกับภาษาไทยแบ่งดังนี้กรณีต่าง ๆ ดังนี้

- กรณีพบ ใ- ใ- ใ-ย และ ัย จะเปลี่ยนให้อยู่ในรูปแบบเดียวกันคือ ัย ก่อนทำการเข้ารหัสชาวค้เด็กซ์ เพื่อให้ได้รหัสชาวค้เด็กซ์เหมือนกัน เนื่องจากสระดังกล่าวอ่านออกเสียงเหมือนกัน เช่น ไท ไท ไทย และ ทัย จะเปลี่ยนเป็น ทัย
- กรณีพบ รร จะเปลี่ยนเป็น ัน ในกรณีที่ไม่มีตัวสะกดตามหลัง และเปลี่ยนเป็น ัน กรณีที่มีตัวสะกดตามหลัง เช่น เปลี่ยนคำว่า สรรเพชร รังสรรค์ พรรณทิ และ ชรรมรัตน์ เป็น สันเพชร รังสัน ัพณทิ และ ชัมรัตน์ ตามลำดับ
- กรณีพบการันต์ จะตัดการันต์และพยัญชนะที่มีตัวการันต์กำกับรวมทั้งสระและอักษรควบการันต์ทั้ง เช่น คำว่า จันทร์ ศักดิ์ และ พันธุ์ เปลี่ยนเป็น จัน ศัก และ พัน ตามลำดับ

ตัวอย่างการเข้ารหัสชาวเด็กซ์ภาษาไทย ดังตารางที่ 2.5

ชื่อภาษาไทย	รหัสชาวเด็กซ์ภาษาไทย
อัมพร	๐059000
อำภรณ์	๐059000
พรรณศักดิ์	พ341000
พันธุ์ศักดิ์	พ341000
เนืองนิจ	น623400
เนืองนิตย์	น623400
ทับประด	ท559400
ทับปรท	ท594000

ตารางที่ 2.5 ตัวอย่างการเข้ารหัสชาวเด็กซ์ภาษาไทย

งานวิจัยของนิลเนตร อรรถวงศ์ ณ อยุธยา มีการออกแบบขั้นตอนวิธีการเข้ารหัสชาวเด็กซ์ภาษาไทยเพื่อให้รหัสที่ได้ตรงกันระหว่างคำที่สะกดถูกและสะกดผิด เพื่อแก้ปัญหาการสับสนคำไทยที่มักสะกดผิด โดยมีแนวคิดที่ว่าพยายามพิจารณาอักษรในทุกตัวในการเข้ารหัส หลักเกณฑ์ต่าง ๆ ในการเข้ารหัสชาวเด็กซ์มีดังนี้

- ใช้อักษรไทยทั้งหมดในการเข้ารหัส เพื่อสะดวกในการอ่าน
- ถ้าตัวอักษรที่อยู่ติดกันเป็นตัวควบกล้ำให้เปลี่ยนเป็นรหัสตัวเดียวกันคือเป็นตัวหน้าของตัวควบกล้ำ เช่น คำว่า ปรับ คลอง และ จริง ให้เข้ารหัสเฉพาะ ปีบ คอง และ จึง ตามลำดับ
- ถ้าตัวอักษรที่อยู่ติดกันเหมือนกันให้เปลี่ยนเป็นรหัสเพียงหนึ่งตัว
- ถ้าตัวอักษรที่อยู่ติดกันเป็นอักษรนำเสียงสนธิมี 9 ตัว ได้แก่ ออ หง หญ หน หม หย หร หล และ หว เช่นคำว่า อยู่ เหงือก หญิง หนุ หมอ หยก หริง หลวง และ หวัง เป็นต้น โดยจะตัดอักษร อ หรือ ห แล้วเปลี่ยนอักษรที่เหลือเป็นรหัสเพียงหนึ่งตัว
- ไม่นำวรรณยุกต์และไม่ได้นำมาเข้ารหัส
- เปลี่ยนพยัญชนะตัวแรกของคำเป็นรหัสตัวอักษรตามตารางที่ 2.6
- เปลี่ยนอักษระตั้งแต่ตัวที่สองขึ้นไปเป็นรหัสตามตารางที่ 2.7
- การเปลี่ยนตำแหน่งสระหน้า ได้แก่ เ- แ- โ- ใ- ไ- ไปไว้หลังสุดของคำ เพื่อให้คำที่สะกดผิดได้รหัสเหมือนกันเช่น ขโมย กับ โขมย ถไล กับ ไถล เป็นต้น

ตัวอย่างการเข้ารหัสชาวคเด็กซ์ภาษาไทย ดังตารางที่ 2.8

ชื่อภาษาไทย	รหัสชาวคเด็กซ์ภาษาไทย
พหูสุต	พหูสุต
พหูสุต	พหูสุต
กิตติพรรณ	กตบน
กิติพันธุ์	กตบน
เมืองนิง	นีองนค
เมืองนิตย์	นีองนค
สับประด	ชบนค
สัปรด	ชบนค
เถลไถล	ทนคนไ
ถลเถล	ทนคนไ

ตารางที่ 2.8 ตัวอย่างการเข้ารหัสชาวคเด็กซ์ภาษาไทย

งานวิจัยของ Theppitak Karoonboonyanan ได้นำเสนอวิธีใหม่ในการเข้ารหัสสำหรับชาวคเด็กซ์ภาษาไทย โดยมีแนวคิดที่ว่าระบบหน่วยคำในคำในภาษาไทยไม่มีเครื่องหมายแบ่งพยางค์ ทำให้คำบางคำสามารถอ่านออกเสียงได้หลายรูปแบบ เช่น คำว่า “เสนา” สามารถอ่านออกเสียงได้ทั้ง เสนา และ เส-นา เป็นต้น ซึ่งการเข้ารหัสชาวคเด็กซ์ได้ใช้หลักการอ่านออกเสียงของคำในการเข้ารหัส ถ้าหากว่าแบ่งพยางค์ผิดพลาดทำให้เสียงที่ได้ในการเข้ารหัสก็ผิดพลาด ดังนั้นเพื่อลดความผิดพลาดในการอ่านออกเสียงจึงได้ใช้เทคนิค ออโตมาตันสถานะจำกัดแบบไม่แน่นอน มีข้อมูลออก (Nondeterministic Finite Automaton with Output) เพื่อสร้างรหัสชาวคเด็กซ์สำหรับการอ่านออกเสียงทุกกรณีของคำ ทำให้การเข้ารหัสคำหนึ่งคำอาจจะได้หลายรหัสชาวคเด็กซ์ซึ่งจะต่างจากวิธีการเข้ารหัสชาวคเด็กซ์แบบอื่นที่ได้กล่าวมาแล้วคือการเข้ารหัสคำหนึ่งคำจะได้รหัสชาวคเด็กซ์เพียงรหัสเดียว ส่วนในการค้นคืนจะนำรหัสทุกตัวของคำที่เป็นข้อความไปเปรียบเทียบกับรหัสทุกตัวของคำที่ต้องการค้นคืน ถ้าพบว่ามีรหัสเหมือนกันเพียงคู่หนึ่งก็จะถือว่าคำทั้งคู่เป็นคำที่อ่านออกเสียงคล้ายกัน การกำหนดรหัสชาวคเด็กซ์ และขั้นตอนวิธีการเข้ารหัสชาวคเด็กซ์มีดังนี้

การกำหนดรหัสชาวค์เด็กซ์สำหรับอักษรไทยมีดังนี้

1. กำหนดรหัสชาวค์เด็กซ์สำหรับพยัญชนะต้น จะแบ่งตามเสียงของพยัญชนะ 21 กลุ่มของไทย และได้รวม ร กับ ล เป็นกลุ่มเดียวกัน โดยจะใช้พยัญชนะไทยหนึ่งตัวแทนกลุ่มเสียง ดังแสดงในตารางที่ 2.9
2. กำหนดรหัสชาวค์เด็กซ์สำหรับสระ จะแบ่งตามเสียงของสระและรวมสระเสียงสั้นและยาวเข้าด้วยกัน โดยจะใช้อักษรโรมันกลุ่มเสียงสระดังแสดงในตารางที่ 2.10
3. กำหนดรหัสชาวค์เด็กซ์สำหรับตัวสะกด จะแบ่งกลุ่มเสียงตามมาตราต่าง ๆ ของอักษรไทย และได้เพิ่มกลุ่ม ฮ สำหรับคำที่ไม่มีตัวสะกดดังแสดงในตารางที่ 2.11

รหัสชาวค์เด็กซ์	ตัวอักษร
ก	ก
ค	ข ฅ ค ฅ ฅ
ง	ง
จ	จ
ช	ช ฅ ฅ
ซ	ซ ฅ ฅ ฅ
ด	ด ฅ* ฅ
ต	ต ฅ
ท	ท ฅ* ฅ ฅ ฅ ฅ
น	ณ น

รหัสชาวค์เด็กซ์	ตัวอักษร
บ	บ
ป	ป
พ	พ ฅ ฅ
ฟ	ฝ ฝ
ม	ม
ย	ฅ ฅ
ร	ร ฅ ฅ
ว	ว
อ	อ
ฮ	ห ฅ

* ฅ สามารถอ่านออกเสียงได้ 2 แบบคือ /ค/ และ /ก/ จึงสร้างรหัสทั้ง 2 แบบ

ตารางที่ 2.9 การเข้ารหัสสำหรับพยัญชนะต้น

จุฬาลงกรณ์มหาวิทยาลัย

รหัสขาค์เด็กซ์	รูปแบบสระ
A	-ะ ั -รร- -อ ำ ใ- ใ- ใ-อ ใ-า -า
I	ิ- -ฤ- (ร)
V	ึ- -ฤ- (ร)
U	ู-
E	เ-ะ เ-็ เ-
X	แ-ะ แ-็ แ-
O	โ-ะ อ+ตัวสะกด โ-
C	เ-าะ -อ- -อ- -ร (อร)
D	เ-อะ เ-อ เ-็ เ-ย
J	เ-็ยะ เ-็ช
W	เ-็อะ เ-็อ
S	-ัวะ -ัว -ัว-

ตารางที่ 2.10 การเข้ารหัสสำหรับสระ

รหัสขาค์เด็กซ์	ตัวอักษร
ก	ก ข ค ฉ
ง	ง
ค	จ ฉ ช ซ ฌ ญ ฎ ฏ ฐ ฑ ฒ ค,คด ท ษ ศ ส ษ
น	ญ ฒ น ร ก พ
บ	บ ป ฟ ฟ ก ผ ฝ
ม	ม ำ
ย	ย
ว	ว
ย	-

ตารางที่ 2.11 การเข้ารหัสสำหรับตัวสะกด

ขั้นตอนวิธีการเข้ารหัสชาวคเด็กซ์มีดังนี้

1. สร้างโครงร่างเครื่องเข้ารหัสชาวคเด็กซ์ จากการรวบรวมกฎเกณฑ์ต่าง ๆ ในการสร้างคำในภาษาไทย
2. ขยายเครื่องเข้ารหัสชาวคเด็กซ์ โดยกำหนดข้อมูลนำออกให้สำหรับแต่ละกฎเกณฑ์
3. นำค่าที่ต้องการเข้ารหัสผ่านเครื่องเข้ารหัสชาวคเด็กซ์ ผลลัพธ์ที่สมบูรณ์แต่ละตัวที่ได้คือรหัสชาวคเด็กซ์

ตัวอย่างขั้นตอนการสร้างเครื่องเข้ารหัสชาวคเด็กซ์

สมมติว่ากฎเกณฑ์การสร้างคำที่ถูกต้องมีรูปแบบดังนี้

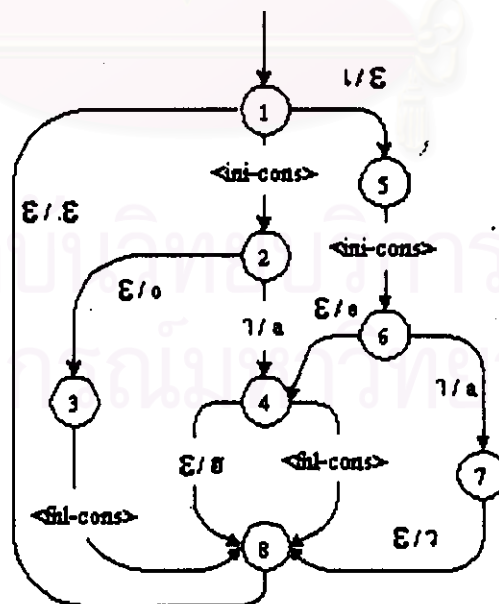
1. $\langle \text{int-cons} \rangle \langle \text{fnl-cons} \rangle$
2. $\langle \text{int-cons} \rangle 1 [\langle \text{fnl-cons} \rangle]$
3. $1 \langle \text{int-cons} \rangle 1$
4. $1 \langle \text{int-cons} \rangle [\langle \text{fnl-cons} \rangle]$

โดยที่ $\langle \text{ini-cons} \rangle ::= \langle \text{cons} \rangle \mid \langle \text{cons} \rangle \langle \text{cons} \rangle$

$\langle \text{cons} \rangle ::= \text{ส} \mid \text{น}$

$\langle \text{fnl-cons} \rangle ::= \text{ส} \mid \text{น}$

จากกฎเกณฑ์ดังกล่าวสามารถแสดงดังในรูปที่ 2.2



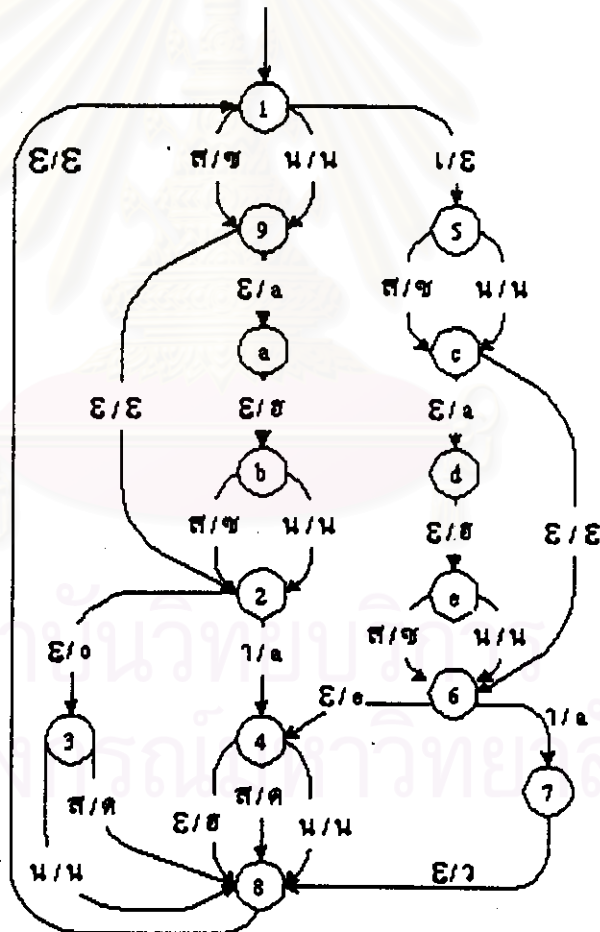
รูปที่ 2.2 โครงร่างเครื่องเข้ารหัสชาวคเด็กซ์

กำหนดข้อมูลนำออกให้แต่ละกฎเกณฑ์ดังนี้

1. $\langle \text{int-cons} \rangle (\epsilon / o) \langle \text{fml-cons} \rangle$
2. $\langle \text{int-cons} \rangle (1 / a) (\epsilon / \theta)$
3. $\langle \text{int-cons} \rangle (1 / a) \langle \text{fml-cons} \rangle$
4. $(1 / \epsilon) \langle \text{int-cons} \rangle (1 / a) (\epsilon / 1)$
5. $(1 / \epsilon) \langle \text{int-cons} \rangle (\epsilon / e) (\epsilon / \theta)$
6. $(1 / \epsilon) \langle \text{int-cons} \rangle (\epsilon / e) \langle \text{fml-cons} \rangle$

โดยที่ $\langle \text{ini-cons} \rangle ::= \langle \text{cons} \rangle \mid \langle \text{cons} \rangle (\epsilon / a) (\epsilon / \theta) \langle \text{cons} \rangle$
 $\langle \text{cons} \rangle ::= (\text{ส} / \text{ช}) \mid (\text{น} / \text{น})$
 $\langle \text{fml-cons} \rangle ::= (\text{ส} / \text{ค}) \mid (\text{น} / \text{น})$

จากการกำหนดข้อมูลนำออกสามารถแสดงดังในรูปที่ 2.3



รูปที่ 2.3 เครื่องเข่านักซาวด์เด็กซ์

ตัวอย่างนำคำว่า “เสนา” ผ่านเครื่องเข้ารหัสชาวค้เด็กซ์ ได้ผลลัพธ์ที่สมบูนดตองค่าคือ

1. (1)—1/ε →(5)—π/ซ →(c)—ε/a →(d)—ε/ธ →(e)—น/น →(6)→ 1/a
 →(7)—ε/v →(8)*
2. (1)—1/ε →(5)—π/ซ →(c)—ε/ε →(6)—ε/c →(4)—ε/ธ →(8)→ ε/ε
 →(1)→ น/น →(9)—ε/ε →(2)—1/a →(4)—ε/ธ →(8)*

สรปรหัสชาวค้เด็กซ์ที่ได้คือ “ซซสนาว” และ “ซซนธธ”

ตัวอย่างการเข้ารหัสชาวค้เด็กซ์ภาษาไทย ดังตารางที่ 2.12

ชื่อภาษาไทย	รหัสชาวค้เด็กซ์ภาษาไทย
บุญญา	บนนยลธ, บนธปอนยลธ, บนนยลธยลธ, บนธคยลธ, บนธยลธยลธ, บนนยลธ
บุณยา	บนนยลธ, บนนณลธยลธ, บนนณลธยลธ, บนนณลธยลธ
ต้บประด	ซลบปลธลค, ซลบบลธปลธลคคลธ, ซลบปลธลคคลธ, ซลบบลธปลธลคยลธ, ซลบบลธปลธลคยลธ, ซลบบลธลคยลธ, ซลบบลธลคยลธ, ซลบบลธลคยลธ, ซลบบลธปลธลคยลธ, ซลบบลธลคยลธ, ซลบบลธลคยลธ, ซลบบลธปลธลค
ต้ปรถ	ซลบปลธลค, ซลบบลธลคซลธ, ซลบลคซลธ, ซลบบลธซลธ, ซลบซลธ, ซลบบลธลคซลธ, ซลบบลธลคซลธ, ซลบลคซลธ, ซลบลคซลธ, ซลบบลธซลธ, ซลบซลธ, ซลบบลธลคซลธ, ซลบบลธลคซลธ, ซลบลคซลธ, ซลบลคซลธ, ซลบลค
ธรรมะ	ทลมมลธ, ทลนมลธ, ทคณคณมมลธ, ทคณคยลธ, ทคณลธยลธ, ทคยลคยลคณมมลธ, ทลคยลคยลคณมมลธ, ทคยลคยลคณมมลธ, ทลคยลคยลคณมมลธ, ทคยลคยลคยลธ, ทลคยลคยลคยลธ, ทคยลคยลคยลธ, ทคยลคยลคยลธ, ทลคยลคยลคยลธ, ทลคยลคยลคยลธ, ทคยลคยลคยลธ, ทลคยลคยลคยลธ
ธัมมะ	ทลมมลธ, ทลมมลธยลธ

ตารางที่ 2.12 ตัวอย่างการเข้ารหัสชาวค้เด็กซ์ภาษาไทย

2.7 ขั้นตอนวิธีระยะแก้ไขสั้นที่สุด (Minimum Edit Distance)

ระยะแก้ไขสั้นที่สุด¹⁵ เป็นเทคนิคหนึ่งในการวัดความคล้ายคลึงกันระหว่าง 2 สายอักขระ ซึ่งจะทำการคำนวณหาจำนวนคำสั่งน้อยที่สุดที่ใช้ในการเพิ่ม การลบ และการแทนที่แต่ละตัวอักขระ เพื่อให้สายอักขระทั้งสองสายเหมือนกัน ตัวอย่างเช่น ระยะห่างของการแก้ไขให้ EXSAMBL เป็น EXAMPLE เท่ากับ 3 ซึ่งมีวิธีการคำนวณดังนี้

- | | | | |
|-------------------------------|---------|----|---------|
| 1. การลบตัวอักษร S | EXSAMBL | => | EXAMBL |
| 2. การแทนที่ตัวอักษร B ด้วย P | EXAMBL | => | EXAMPL |
| 3. การเพิ่มตัวอักษร E | EXAMPL | => | EXAMPLE |

จากวิธีการคำนวณข้างต้นสามารถเขียนในอยู่ในรูปการคำนวณด้วยความสัมพันธ์เวียนเกิด Edit (P_j, W_k) ได้ดังนี้

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

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
 $r(p_j, w_k) = 0$ ถ้า p_j เท่ากับ w_k
 1 ถ้า p_j ไม่เท่ากับ w_k

¹⁵ 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.

ขั้นตอนวิธีการคำนวณหาค่าระยะการแก้ไขสั้นที่สุดใช้เวลาทำงานเป็น $O(mn)$ โดย m และ n คือความยาวของสายอักขระที่ 1 และ 2 ตามลำดับ เมื่อใช้เทคนิคกำหนดการพลวัต (Dynamic Programming)¹⁶ ดังแสดงในรูปที่ 2.4

```

int EditDist(char *P; char *W) {
    int F [STRLEN+1][STRLEN+1];
    int i, j, l1, l2;

    l1 = strlen (P);
    l2 = strlen (W);
    for (i=0; i <= STRLEN; i++)
        F[0][i] = F[i][0] = i;
    for (i=0; i <= STRLEN; i++)
        for (j=0; j <= STRLEN; j++)
            F[i][j] := Min( F[i-1][j] + 1,
                           F[i][j-1] + 1,
                           F[i-1][j-1] + Equal(P[i-1], W[j-1]) );

    return (F[l1][l2]);
}

int Equal(char a, char b) {
    return ( a == b ? 0 : 1 );
}

```

รูปที่ 2.4 โปรแกรมระยะแก้ไขสั้นที่สุด โดยใช้เทคนิคกำหนดการพลวัต

¹⁶ 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.

2.8 สรุป

ในบทนี้ได้กล่าวถึงทฤษฎีและงานวิจัยต่าง ๆ ที่เกี่ยวข้องในการพัฒนาขั้นตอนวิธีการเข้ารหัสคำทับศัพท์เพื่อการค้นคืนข้ามภาษาไทย-อังกฤษ โดยใช้แนวคิดในการแปลข้อความที่เป็นคำทับศัพท์ด้วยทฤษฎีของการถอดอักษร และใช้เทคนิคชาวค้เด็กซ์ช่วยในการเข้ารหัสคำ โดยรายละเอียดต่าง ๆ จะกล่าวไว้ในบทถัดไป



สถาบันวิทยบริการ
จุฬาลงกรณ์มหาวิทยาลัย