

การระบุค่าไทยและค่าทับศัพท์ด้วยแบบจำลองเอ็นแกรม

นายอัศวพล เอกวงศ์อนันต์

วิทยานิพนธ์นี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญาอักษรศาสตรมหาบัณฑิต

สาขาวิชาภาษาศาสตร์ ภาควิชาภาษาศาสตร์

คณะอักษรศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย

ปีการศึกษา 2548

ISBN 974-53-2360-8

ลิขสิทธิ์ของจุฬาลงกรณ์มหาวิทยาลัย

IDENTIFICATION OF THAI AND TRANSLITERATED WORDS
BY N-GRAM MODELS

Mr. Akarapol Ekwonganan

A Thesis Submitted in Partial Fulfillment of the Requirements
for the Degree of Master of Arts in Linguistics

Department of Linguistics

Faculty of Arts

Chulalongkorn University

Academic Year 2005

ISBN 974-53-2360-8

หัวข้อวิทยานิพนธ์ การระบุค่าไทยและค่าทับศัพท์ด้วยแบบจำลองเอ็นแกรม
โดย นายอัศวพล เอกวงศ์อนันต์
สาขาวิชา ภาษาศาสตร์
อาจารย์ที่ปรึกษา ผู้ช่วยศาสตราจารย์ ดร.วิโรจน์ อรุณมานะกุล

คณะอักษรศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย อนุมัติให้หัวข้อวิทยานิพนธ์ฉบับนี้เป็น
ส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญาโท

.....คณบดีคณะอักษรศาสตร์
(ศาสตราจารย์ ดร.ธีระพันธ์ เหลืองทองคำ)

คณะกรรมการสอบวิทยานิพนธ์

.....ประธานกรรมการ
(ผู้ช่วยศาสตราจารย์ ดร.สุตาพร ลักษณะียนาวิน)

.....อาจารย์ที่ปรึกษา
(ผู้ช่วยศาสตราจารย์ ดร.วิโรจน์ อรุณมานะกุล)

.....กรรมการ
(อาจารย์ ดร. ณัฐกร ทับทอง)

อัครพล เอกวงศ์อนันต์ : การระบุคำไทยและคำทับศัพท์ด้วยแบบจำลองเอ็นแกรม.
(IDENTIFICATION OF THAI AND TRANSLITERATED WORDS BY N-GRAM
MODELS) อาจารย์ที่ปรึกษา : ผศ. ดร. วิโรจน์ อรุณมานะกุล, จำนวน 270 หน้า. ISBN
974-53-2360-8.

วัตถุประสงค์ของการวิจัยครั้งนี้ เพื่อค้นหาสายอักขระเฉพาะสำหรับใช้ในการระบุภาษา
ของคำโดยใช้คลังข้อมูลคำไทย คำทับศัพท์ภาษาอังกฤษ ภาษาญี่ปุ่นและภาษาฝรั่งเศส และ
พัฒนาระบบการระบุภาษาของคำไทยและคำทับศัพท์ภาษาต่างประเทศโดยใช้สายอักขระ
เฉพาะและใช้แบบจำลองเอ็นแกรมขนาด 1-5 แกรม

คลังข้อมูลที่ใช้ในงานวิจัยนี้ คือ คลังข้อมูลคำไทย คำทับศัพท์ภาษาอังกฤษ ภาษาญี่ปุ่น
ภาษาละ 10,000 คำ และคำทับศัพท์ภาษาฝรั่งเศส 1,000 คำ โดยเก็บจากข้อมูลที่พบใน
ภาษารวมชาติซึ่งอาจจะไม่ได้ทับศัพท์ถูกต้องตามเกณฑ์ของราชบัณฑิตยสถานก็ได้ 80%
ของคลังข้อมูลถูกนำมาใช้เพื่อหาสายอักขระเฉพาะและสร้างแบบจำลองเอ็นแกรมของแต่ละ
ภาษา ในขณะที่อีก 20% ถูกใช้เพื่อการทดสอบระบบแบบต่าง ๆ

สายอักขระเฉพาะที่พบสะท้อนให้เห็นถึงลักษณะเฉพาะของแต่ละภาษาได้ในระดับหนึ่ง
จึงมีผลให้ระบบที่ใช้สายอักขระเฉพาะในการระบุภาษาสามารถตัดสินภาษาได้ถูกต้อง 50.58%
48.71% 54.09% และ 20.40% สำหรับคำไทย คำทับศัพท์ภาษาอังกฤษ ภาษาญี่ปุ่น และ
ฝรั่งเศส ตามลำดับ

เมื่อใช้แบบจำลองเอ็นแกรมในการระบุภาษา ระบบสามารถระบุภาษาของคำไทย คำทับ
ศัพท์ภาษาอังกฤษ และญี่ปุ่นได้ถูกต้องกว่า 90% แต่ได้เพียงประมาณ 60% สำหรับคำทับศัพท์
ฝรั่งเศส ผลที่ได้ยืนยันว่าขนาดของข้อมูลการฝึกมีผลต่อการทำงานของระบบการระบุภาษาทั้ง
สองระบบ นอกจากนี้ จากผลที่พบว่าระบบที่ใช้แบบจำลอง 3-แกรมให้ผลดีกว่าระบบที่ใช้ขนาด
แกรมอื่นๆ ทำให้สรุปได้ว่า ขนาดของเอ็นแกรมมีผลต่อการทำงานของระบบการระบุภาษา

ภาควิชา	ภาษาศาสตร์	ลายมือชื่อนิติ
สาขาวิชา	ภาษาศาสตร์	ลายมือชื่ออาจารย์ที่ปรึกษา..... ..
ปีการศึกษา	2548	

##4580259022 : MAJOR LINGUISTICS

KEY WORD: LANGUAGE IDENTIFICATION

AKARAPOL EKWONGANAN : IDENTIFICATION OF THAI AND TRANSLITERATED WORDS BY N-GRAM MODELS. THESIS ADVISOR : ASST. PROF. WIROTE AROONMANAKUN, Ph.D., 270 pp. ISBN 974-53-2360-8.

This research aims to find the unique character sequences of Thai and transliterated words (English, Japanese, and French), and implement language identification systems using unique character sequences and n-gram models (1-5 gram).

The corpora in this research consist of 10,000 Thai words, 10,000 English transliterated words, 10,000 Japanese transliterated words, and 1,000 French transliterated words. Transliterated words are collected from naturally occurring texts, even some of them are not conformed to the Royal Institute guidelines of transliteration. 80% of the corpus is used to extract unique character sequences and to build an n-gram language model of each language, while the other 20% is used for testing the systems.

The unique character sequences reflect some characteristics of the languages. As a result, the system using unique character sequence can identify languages correctly 50.58%, 48.71%, 54.09%, and 20.40% for Thai words, English, Japanese, and French transliterated words respectively.

When an n-gram language model is used, the system can identify languages correctly more than 90% for Thai, English and Japanese transliterated word, but only about 60% for French transliterated words. This confirms that the size of training corpus affects the performances of both systems. The results also show that the system using 3-gram model performs better than other n-gram models. Therefore, we can conclude that the size of n-gram does affect the performance of the language identification system.

Department	LINGUISTICS	Student's signature.....
Field of study	LINGUISTICS	Advisor's signature.....
Academic year	2005	

กิตติกรรมประกาศ

วิทยานิพนธ์ฉบับนี้สำเร็จลงได้ด้วยความช่วยเหลืออย่างดียิ่งของผู้ช่วยศาสตราจารย์ ดร. วิโรจน์ อรุณมานะกุล อาจารย์ที่ปรึกษาวิทยานิพนธ์ ซึ่งท่านได้ให้คำปรึกษา แนะนำ ช่วยตรวจทาน แก้ไข งานวิจัยด้วยความเอาใจใส่เป็นอย่างดี และจุดประกายให้ผู้วิจัยทำงานวิจัยนี้

ผู้วิจัยขอกราบขอบพระคุณผู้ช่วยศาสตราจารย์ ดร. สุตาพร ลักษณะียนาวิน และ ดร. ณิชกร ทับทอง ที่ได้กรุณาให้คำแนะนำและข้อคิดเห็นในการทำวิทยานิพนธ์ฉบับนี้ด้วยความเอาใจใส่ ทั้งยังช่วยตรวจแก้ไขให้สำเร็จสมบูรณ์ยิ่งขึ้น

ผู้วิจัยขอขอบพระคุณคณาจารย์ภาควิชาภาษาศาสตร์ คณะอักษรศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย ที่ให้ความรู้ความเมตตาสั่งสอนทั้งทางด้านวิชาการและการทำงานวิจัย และผู้วิจัยขอกราบขอบพระคุณบิดา คุณแม่และคุณป้าที่ให้ความรัก ความเข้าใจ ความเอาใจใส่ ให้โอกาส และสนับสนุนทุกด้าน รวมทั้งเป็นกำลังใจให้แก่ผู้วิจัยตลอดมา สุดท้ายขอขอบคุณพี่ชายและเพื่อน ๆ ที่คอยให้คำแนะนำ ปรึกษา ในการทำวิจัย รวมทั้งเป็นกำลังใจให้งานวิจัยนี้สำเร็จด้วยดี

สารบัญ

	หน้า
บทคัดย่อภาษาไทย	ง
บทคัดย่อภาษาอังกฤษ	จ
กิตติกรรมประกาศ	ฉ
สารบัญ	ช
สารบัญตาราง	ฅ
สารบัญภาพและแผนภูมิ	ญ
บทที่	
1 บทนำ.....	1
1.1 ความเป็นมาและความสำคัญของปัญหา.....	1
1.2 วัตถุประสงค์ของการวิจัย	3
1.3 สมมติฐานการวิจัย	4
1.4 ขอบเขตของการวิจัย.....	4
1.5 ประโยชน์ที่คาดว่าจะได้รับ	4
1.6 โครงสร้างวิทยานิพนธ์.....	5
1.7 ศัพท์เฉพาะที่ใช้ในงานวิจัย.....	6
2 การระบุภาษาของคำและข้อความ	7
2.1 แนวทางในการสร้างแบบจำลองภาษา.....	10
2.2 แนวทางในการสร้างวิธีรวมสายอักขระเฉพาะ.....	17
2.3 แนวทางในการจัดทำระบบการระบุภาษาของคำ	19
2.4 หลักเกณฑ์การทับศัพท์ภาษาต่างประเทศ	22
3 คลังข้อมูลคำทับศัพท์.....	24
3.1 การสร้างคลังข้อมูล	24
3.2 เกณฑ์ในการคัดเลือกข้อมูล	25
3.3 การสร้างข้อมูลฝึกสอนและข้อมูลทดสอบสำหรับการทดลอง	29
4 การระบุภาษาของคำด้วยสายอักขระเฉพาะ	31
4.1 ภาพโดยรวมของระบบ.....	31

สารบัญ (ต่อ)

บทที่	หน้า
4.2 โปรแกรมระบุภาษาของคำ	32
4.3 ผลการทดลอง	34
4.4 วิเคราะห์ผลการทดสอบ	37
4.5 ปัญหาการตัดสินค้าผิดพลาดที่เกิดขึ้นจากการใช้สายอักขระเฉพาะ 1-5 ตัว	38
4.6 สายอักขระเฉพาะที่ใช้จริงในระบบ	41
4.7 การวิเคราะห์สายอักขระเฉพาะ	51
4.8 สรุปผลการใช้วิธีสายอักขระเฉพาะ	59
5 การระบุภาษาของคำด้วยแบบจำลองภาษา	61
5.1 การสร้างแบบจำลองภาษา (Language model)	61
5.2 โปรแกรมระบุภาษาของคำด้วยแบบจำลองภาษา	62
5.3 สรุปผลการทดลอง	64
5.4 ปัญหาการตัดสินค้าผิดพลาดที่เกิดขึ้นจากแบบจำลองภาษา	67
5.5 สรุปการใช้แบบจำลองภาษาเอ็นแกรม	69
6 โปรแกรมระบุภาษาของคำ	70
6.1 ภาพโดยรวมของระบบ	70
6.2 โปรแกรมระบุภาษาของคำ	71
6.3 สรุปผลการทดลอง	73
7 สรุปผลการวิจัย อภิปรายผลและข้อเสนอแนะ	76
7.1 สรุปผลการวิจัย	76
7.2 อภิปรายผลการวิจัย	77
7.3 ข้อเสนอแนะ	79
รายการอ้างอิง	81
ภาคผนวก	84
ภาคผนวก ก สายอักขระเฉพาะ	85
ภาคผนวก ข สายอักขระเฉพาะที่ใช้จริงในระบบและจำนวนครั้งที่ใช้	230
ประวัติผู้เขียนวิทยานิพนธ์	270

สารบัญตาราง

ตารางที่	หน้า
2.1	ตัวอย่างการจำแนกภาษาของข้อความ..... 7
2.2	ตัวอย่างการระบุภาษาของคำ 8
2.3	ตัวอย่างการเปรียบเทียบความน่าจะเป็นของการระบุประเภท ของข้อความ 16
2.4	ตัวอย่างการใช้วิธีรวมสายอักขระเฉพาะ (Dunning, 1994)..... 17
3.1	การเขียนทับศัพท์ที่ต่างจากเกณฑ์ทางราชบัณฑิตยสถาน 28
4.1	การแยกสายอักขระตามเอ็นแกรม 1-5 แกรม 33
4.2	ผลการตัดสีนภาษาได้ถูกต้อง ผลการตัดสีนภาษาผิดพลาดและผลการตัดสีน ภาษาไม่ได้ ด้วยสายอักขระเฉพาะ 1-5 ตัว 35
4.3	ผลการตัดสีนภาษาของคำผิดพลาดจากการใช้สายอักขระเฉพาะ 1-5 ตัว 39
4.4	การเปรียบเทียบจำนวนสายอักขระเฉพาะที่พบในคลังข้อมูลการฝึกและจำนวน ของสายอักขระเฉพาะที่ใช้จริงในแต่ละภาษา 42
5.1	การแยกสายอักขระตามเอ็นแกรม 2-5 แกรม 63
5.2	ตัวอย่างผลการทดลองที่ได้แบบจำลองภาษา 64
5.3	ผลการตัดสีนภาษาของคำได้ถูกต้องด้วยแบบจำลอง 2-5 แกรม..... 65
5.4	ผลการตัดสีนภาษาของคำผิดพลาดจากการใช้แบบจำลองภาษา 2-5 แกรม..... 67
6.1	ผลการตัดสีนภาษาของคำได้ถูกต้องด้วยโปรแกรมระบุคำ (สายอักขระเฉพาะ 1-5 ตัว และแบบจำลอง 2-5 แกรม) 74
6.2	ผลการเปรียบเทียบค่าความถูกต้องของระบบสายอักขระเฉพาะ ระบบแบบจำลอง ภาษา และระบบที่ประยุกต์ทั้ง 2 ระบบ 75

สารบัญภาพและแผนภูมิ

ภาพและแผนภูมิ	หน้า
ภาพที่ 3.1 การแบ่งส่วนคลังข้อมูลในแต่ละภาษา	30
ภาพที่ 4.1 ลักษณะของการสร้างสายอักขระเฉพาะ 1-5 แกรม.....	32
ภาพที่ 4.2 ขั้นตอนการทำงานของโปรแกรมระบุภาษาของคำด้วยสายอักขระเฉพาะ	33
ภาพที่ 5.1 ลักษณะของการแบบจำลองภาษาเอ็นแกรม.....	61
ภาพที่ 5.2 ขั้นตอนการทำงานของโปรแกรมระบุภาษาของคำด้วยแบบจำลองภาษา.....	62
แผนภูมิที่ 5.3 การเปรียบเทียบการใช้แบบจำลองของแต่ละภาษา	66
ภาพที่ 6.1 ขั้นตอนการทำงานของโปรแกรมระบุภาษาของคำ	71

บทที่ 1

บทนำ

1.1 ความเป็นมาและความสำคัญของปัญหา

เนื่องจากในปัจจุบัน ประเทศไทยมีการติดต่อกับต่างประเทศมากขึ้น จึงได้รับอิทธิพลทั้งทางวัฒนธรรมและเทคโนโลยีใหม่ๆ ภาษาที่ใช้ในปัจจุบัน จึงไม่เพียงพอต่อความต้องการ โดยเฉพาะศัพท์วิชาการและศัพท์เทคนิค มีการแพร่หลายอย่างรวดเร็ว จึงมีการสร้างคำหรือศัพท์ใหม่ๆ ขึ้นมาใช้ในภาษาไทย เพื่อเรียกชื่อสิ่งที่เกิดขึ้นใหม่ๆ จากต่างประเทศ ในรูปของการยืมคำ การแปล การบัญญัติศัพท์ใหม่ ฯลฯ ในการบัญญัติศัพท์ใหม่นั้น ทางราชบัณฑิตยสถานได้วางหลักเกณฑ์ไว้ 2 ข้อ คือ 1. สร้างคำใหม่จากคำศัพท์ที่เป็นคำไทยแท้หรือจากบาลีสันสกฤต 2. หากไม่สามารถบัญญัติศัพท์ใหม่ได้ จะด้วยเหตุผลใดๆ ก็ตามก็ให้ใช้การทับศัพท์แทน การทับศัพท์จึงเป็นวิธีหนึ่งของการสร้างคำใหม่

คำทับศัพท์ หมายถึง คำภาษาต่างประเทศที่เขียนด้วยตัวอักษรไทย อาจเป็นคำสามัญนาม เช่น คอมพิวเตอร์ เทคโนโลยี ฟุตบอล หรือคำวิสามานยนาม เช่น แปซิฟิก เมดิเตอร์เรเนียน ไอบีเอ็ม ยูเนสโก ฯลฯ (ราชบัณฑิตยสถาน, 2538)

คำทับศัพท์เริ่มใช้กันอย่างแพร่หลาย และสามารถทับศัพท์จากภาษาต่างประเทศต่าง ๆ ได้มากมาย เห็นได้จากเกณฑ์ทางราชบัณฑิตยสถาน (2535) มีการกำหนดเกณฑ์การทับศัพท์ไว้ทั้งหมด 9 ภาษา คือ ภาษาอังกฤษ ภาษาฝรั่งเศส ภาษาญี่ปุ่น ภาษาเยอรมัน ภาษอิตาลี ภาษาสเปน ภาษารัสเซีย ภาษาอาหรับ และภาษามลายู ตัวอย่างคำทับศัพท์จากราชบัณฑิตยสถานได้แก่

<u>คำศัพท์ภาษาอังกฤษ</u>	<u>ราชบัณฑิตยสถานบัญญัติว่าเป็น</u>	<u>ทับศัพท์ว่าเป็น</u>
computer	เครื่องจักรสมองกล	คอมพิวเตอร์
plastics	-	พลาสติก
bacteria	-	แบคทีเรีย

คำทับศัพท์เหล่านี้เขียนโดยใช้ตัวอักษรภาษาไทยและเขียนปะปนกับคำไทยทั่วไป ดังจะเห็นได้จากคำที่ขีดเส้นใต้ในตัวอย่างข้อความจากหนังสือเทคโนโลยีคอมพิวเตอร์และสารสนเทศ (ศรีไพร, 2544) ซึ่งจะแสดงให้เห็นการใช้ศัพท์เฉพาะทางดังนี้

หน่วยความจำหลัก (MAIN MEMORY UNIT) เป็นวงจรรวมหรือชิปที่ใช้บันทึกโปรแกรมและข้อมูล หน่วยความจำหลักจะบรรจุอยู่บนเมนบอร์ดหรือแผงวงจรหลัก หน่วยความจำบางประเภทก็ถูกออกแบบให้อยู่ในชิปพีซีอยู่แล้ว หน่วยของข้อมูลที่จัดเก็บในหน่วยความจำเรียกว่า ไบต์ (byte) ซึ่ง 1 ไบต์ จะประกอบไปด้วย 8 บิต นอกจากนี้ยังมีหน่วยเป็น กิโลไบต์ (kilobyte หรือ KB)

การแยกคำใดเป็นคำทับศัพท์มาหรือเป็นคำไทยอาจเป็นเรื่องง่ายสำหรับมนุษย์ แต่สำหรับคอมพิวเตอร์แล้ว การสามารถระบุแยกว่าสายอักขระที่พบเป็นคำไทยหรือคำทับศัพท์มา และทับศัพท์จากภาษาใด ถือเป็นงานหนึ่งที่ทำทลายความสามารถของคอมพิวเตอร์ และเป็นพื้นฐานที่จำเป็นในการพัฒนาระบบประมวลผลภาษาไทย (Natural Language Processing) ต่าง ๆ

ตัวอย่างเช่น ในการพัฒนาโปรแกรมทางด้านการค้นคืนสารสนเทศข้ามภาษา (Cross-Language Information Retrieval) จะมีโปรแกรมส่วนหนึ่งที่ทำหน้าที่ถอดอักขระคำทับศัพท์กลับจากภาษาไทยเป็นภาษาต้นฉบับ (backward transliteration) เพื่อจะได้รู้ว่าคำเดิมในภาษานั้นคืออะไร การจะทำเช่นนี้ได้ จำเป็นต้องมีการระบุภาษาของคำในข้อความ ก่อนที่จะนำไปใช้กฎการถอดอักขระ (transliteration) ในแต่ละภาษา

จากตัวอย่างข้างบน จะเห็นได้ว่าไม่มีคำหรือสัญลักษณ์ใดๆ ชี้ให้เห็นว่าคำที่ขีดเส้นใต้เป็นคำทับศัพท์ภาษาอังกฤษ เพราะในภาษาไทย มีการใช้คำไทยแท้และคำทับศัพท์จากภาษาอื่นๆ กลมกลืนกันภายใต้การใช้ตัวเขียนภาษาไทย ถึงแม้ว่าคนทั่วไปอาจจะรู้ว่าคำใดเป็นคำไทย หรือ เป็นคำทับศัพท์จากภาษาอังกฤษ ฝรั่งเศส หรือ ญี่ปุ่น แต่คอมพิวเตอร์ไม่สามารถแยกความแตกต่างของภาษาเหล่านี้ได้เอง หากเราจะนำข้อความนี้ไปใช้กับโปรแกรมถอดตัวอักษรจากภาษาไทยเป็นภาษาอังกฤษ เราจะต้องมีการระบุว่าคำใดเป็นคำไทย หรือคำทับศัพท์ภาษาอังกฤษก่อน เพื่อให้คอมพิวเตอร์สามารถเลือกถอดอักขระเฉพาะคำที่ไม่ใช่คำทับศัพท์ มิฉะนั้นข้อความทั้งหมดจะถูกนำไปถอดตัวอักษร ผู้วิจัยจึงต้องการจะพัฒนาโปรแกรมที่ใช้ในการระบุภาษาของคำ เพื่อนำไปใช้ในการตัดสินใจประมวลผลให้ถูกต้อง กล่าวคือ เมื่อให้คำที่ต้องการมาโปรแกรมจะระบุว่าทับศัพท์มาจากภาษาใด เช่น

แบดมินตัน = คำทับศัพท์ภาษาอังกฤษ \Longrightarrow นำไปถอดตัวอักษรเป็น Badminton
 ฮาตารี = คำทับศัพท์ภาษาญี่ปุ่น \Longrightarrow นำไปถอดตัวอักษรเป็นคำว่า Hatari
 หนังสือ = คำไทย \Longrightarrow นำไปแปลด้วยข้อมูลด้วย พจนานุกรม เป็นคำ
 ว่า Book

การพัฒนาโปรแกรมที่ใช้ในการระบุภาษาของคำ อาศัยการรวบรวมคลังข้อมูลคำทับศัพท์ภาษาต่างๆ และคำไทย ซึ่งคำไทยในที่นี้รวมไปถึงคำที่มาจากบาลีสันสกฤตด้วย เพราะผู้วิจัยเห็นว่ากระบวนการระบุที่มาของคำ ไม่ว่าจะเป็นคำไทยแท้หรือเป็นบาลีสันสกฤตนั้น ไม่มีความจำเป็นสำหรับการประมวลผลภาษาไทย เช่น ในการถอดตัวอักษรกลับจากภาษาไทยเป็นภาษาอังกฤษ ทั้งคำไทยและคำบาลีสันสกฤตจะไม่ถูกนำไปถอดอักษร ในการแปลจากภาษาไทยเป็นภาษาอื่นๆ ทั้งคำไทยและคำบาลีสันสกฤตจะถูกนำไปหาคำแปลจากพจนานุกรมสองภาษาเหมือนกัน ในขณะที่ถ้าเป็นคำทับศัพท์ภาษาอื่นจะถูกถอดกลับเป็นภาษาเดิม หลังจากที่ได้คลังข้อมูลคำทับศัพท์แล้ว จึงให้คอมพิวเตอร์ได้เรียนรู้ข้อมูลทางภาษาและสายอักขระเฉพาะ (unique character sequence) ของแต่ละภาษาในเชิงสถิติ จากคลังข้อมูลนี้และในคลังข้อมูลคำทับศัพท์ที่ใช้ในการทดลองนี้ ผู้วิจัยเลือกคำทับศัพท์ภาษาอังกฤษ ญี่ปุ่น และฝรั่งเศส เนื่องจากคำทับศัพท์ภาษาอังกฤษและภาษาญี่ปุ่น มีขนาดของข้อมูลมากและเป็นตัวอย่างเพื่อทดลองเปรียบเทียบระหว่างภาษาแถบตะวันตก และภาษาแถบเอเชีย อีกทั้งยังเลือกทดลองกับคำทับศัพท์ภาษาฝรั่งเศสซึ่งเป็นภาษาที่มีตัวอย่างคำทับศัพท์น้อย เพราะต้องการทดสอบประสิทธิภาพของระบบว่าสามารถระบุที่มาของภาษาเมื่อเรียนรู้จากคลังข้อมูลขนาดเล็กได้หรือไม่

1.2 วัตถุประสงค์ของการวิจัย

1.2.1 หาสายอักขระเฉพาะ (unique character sequence) ในการระบุภาษาของคำโดยใช้คลังข้อมูลคำไทย คำทับศัพท์ภาษาอังกฤษ ญี่ปุ่น และภาษาฝรั่งเศส

1.2.2 พัฒนาระบบการระบุภาษาของคำไทยและคำทับศัพท์ภาษาต่างประเทศโดยใช้แบบจำลองเอ็นแกรม

1.2.3 ประเมินประสิทธิภาพของแบบจำลองเอ็นแกรมขนาดต่างๆ ในการระบุภาษาของคำ

1.3 สมมติฐานการวิจัย

1.3.1 คำไทยและคำทับศัพท์จากภาษาต่าง ๆ จะมีสายอักขระเฉพาะที่ใช้ระบุภาษาของคำนั้นได้

1.3.2 ระบบการระบุภาษาของคำสามารถใช้ระบุภาษาของคำได้ถูกต้องมากกว่า 90%

1.3.3 ขนาดของแบบจำลองเอ็นแกรมมีความสำคัญต่อประสิทธิภาพระบบการระบุภาษา

1.4 ขอบเขตของการวิจัย

1.4.1 คำทับศัพท์ที่ใช้เป็นคำทับศัพท์จากภาษาอังกฤษ ภาษาญี่ปุ่น และภาษาฝรั่งเศส เท่านั้น

1.4.2 คำไทย คำทับศัพท์ภาษาอังกฤษและคำทับศัพท์ภาษาญี่ปุ่น จะรวบรวมจากพจนานุกรมราชบัณฑิตยสถาน และหนังสือต่าง ๆ จำนวนไม่ต่ำกว่าภาษาละ 10,000 คำ ส่วนคำทับศัพท์ภาษาฝรั่งเศสจะรวบรวมจากหนังสือต่างศัพท์เฉพาะ ๆ จำนวนไม่ต่ำกว่า 1,000 คำ อย่างไรก็ตาม คำทับศัพท์ที่พบในภาษาธรรมชาติมีทั้งที่ทับศัพท์ตามเกณฑ์ของราชบัณฑิตยสถาน และที่ไม่ใช้เกณฑ์ตามราชบัณฑิตยสถาน ในที่นี้ผู้วิจัยรวบรวมข้อมูลคำทับศัพท์โดยไม่สนใจว่าเป็นการทับศัพท์ตามเกณฑ์ของราชบัณฑิตยสถานหรือไม่ เพราะในการทำงานจริงโปรแกรมถอดอักษรกลับ (backward transliteration) ของแต่ละภาษาจะต้องถอดคำทับศัพท์ต่าง ๆ นั้นกลับเป็นภาษาเดิมให้ได้ โดยไม่จำเป็นต้องรู้ว่าข้อมูลนั้นทับศัพท์ตามเกณฑ์หรือไม่

1.4.3 แบบจำลองเอ็นแกรมที่ใช้จะทดลองใช้ 1-5 แกรม เพื่อทดสอบหาประสิทธิภาพของระบบเอ็นแกรมที่ดีที่สุด และเปรียบเทียบผลการทดลองเมื่อใช้เอ็นแกรมขนาดต่าง ๆ กัน

1.5 ประโยชน์ที่คาดว่าจะได้รับ

1.5.1 ได้ระบบการระบุภาษาของคำไทยและคำทับศัพท์ภาษาต่าง ๆ

1.5.2 ระบบนี้สามารถนำไปประยุกต์ใช้ในการค้นคืนสารสนเทศข้ามภาษา (cross-language information retrieval) และการถอดอักษรคำยืมทับศัพท์กลับเป็นภาษาต้นฉบับ (backward transliteration)

1.6 โครงสร้างวิทยานิพนธ์

ผู้วิจัยได้จัดทำทั้งหมด 7 บท มีดังนี้

บทที่ 1 บทนำ โดยนำเสนอภาพโดยรวมของความเป็นมาและความสำคัญของปัญหา การระบุภาษาของคำ จุดประสงค์ของการวิจัย สมมติฐานการวิจัย ขอบเขตของงานวิจัย ประโยชน์ที่คาดว่าจะได้รับ และศัพท์เฉพาะที่ใช้ในงานวิจัย

บทที่ 2 การระบุภาษาของคำและข้อความ โดยนำเสนอภาพรวมของงานการระบุภาษาของคำที่ต่างกับการระบุภาษาของข้อความ แนวทางการใช้แบบจำลองภาษาต่างๆ เพื่อแก้ปัญหาวิธีการและแบบจำลองที่เลือกใช้แก้ปัญหา รวมไปถึงเหตุผลที่เลือกใช้

บทที่ 3 คลังข้อมูลคำทับศัพท์ โดยนำเสนอวิธีการสร้างคลังข้อมูล กรณีในการคัดเลือกข้อมูลและ ปัญหาในการสร้างคลังข้อมูล สำหรับทำข้อมูลการฝึกและข้อมูลทดสอบ เพื่อให้ระบบมีประสิทธิภาพมากที่สุด

บทที่ 4 การระบุภาษาของคำด้วยสายอักขระเฉพาะ โดยนำเสนอภาพรวมของระบบระบุภาษาของคำด้วยสายอักขระเฉพาะ 1-5 ตัว ผลการทดสอบระบบ การอภิปรายผล รวมทั้งนำสายอักขระเฉพาะที่ได้ไปวิเคราะห์ลักษณะเฉพาะของภาษา โดยหลังจากรวบรวมคลังข้อมูลและสร้างชุดอักขระเฉพาะในแต่ละภาษาแล้ว เราจะทราบถึงลักษณะการรวมกันของพยัญชนะเด่นๆ ในแต่ละภาษา ซึ่งจะมีรูปแบบตายตัวไม่ซ้ำซ้อนกัน เป็นลักษณะเฉพาะของแต่ละภาษา ซึ่งสายอักขระเฉพาะเหล่านี้จะสะท้อนให้เห็นถึงลักษณะเฉพาะทางด้านเสียงที่เกิดขึ้นของแต่ละภาษาด้วย ดังนั้นผู้วิจัยจึงสรุปลักษณะของเสียงเด่นๆ ที่เกิดขึ้นจากชุดอักขระเฉพาะเหล่านั้น และการปรากฏใช้ของอักขระเด่นๆ ภายในภาษานั้นๆ

บทที่ 5 การระบุภาษาของคำด้วยแบบจำลองภาษา โดยนำเสนอภาพรวมของระบบระบุภาษาของคำด้วยแบบจำลองภาษาตั้งแต่ 2-5 แกรม ผลการทดสอบระบบ การอภิปรายผล และผลการเปรียบเทียบการใช้ขนาดของเอ็นแกรมต่างๆกัน

บทที่ 6 โปรแกรมระบุภาษาของคำ โดยนำเสนอภาพรวมของระบบระบุภาษาของคำที่ใช้ทั้ง 2 วิธีร่วมกัน ผลการทดสอบระบบโดยรวมและการอภิปรายผล

บทที่ 7 สรุป อภิปรายผลและข้อเสนอแนะ โดยสรุปผลการทดสอบระบบทั้งหมดและสรุปอภิปรายผล รวมทั้งเสนอแนวทางในการแก้ไขระบบ และแนวทางในการพัฒนาระบบต่อไป

1.7 ศัพท์เฉพาะที่ใช้ในงานวิจัย

สายคำ (word sequence) คือ คำที่เขียนเรียงกันมากกว่า 1 คำ เช่น สายคำว่า “...to be positive” จะเห็นได้ว่าสายคำนี้ประกอบจากคำ 3 คำคือ คำว่า “to” “be” และ “positive”

สายอักขระหรือชุดอักขระ (string) คือ ตัวอักษรหรืออักขระที่เขียนประกอบเรียงกันเป็นคำหรือเป็นสายอักขระ ในงานวิจัยนี้หมายถึงตัวอักษรภาษาไทย เช่น คำว่า โมดูล โดยจะนับตัวอักษร 1 ตัวเท่ากับ 1 อักขระ ดังนี้ “โ” “ม” “ด” “ู” และ “ล”

ชุดอักขระเฉพาะ (unique character sequence) คือ ตัวอักษรที่เกิดขึ้นเรียงกันในภาษาใดภาษาหนึ่งเท่านั้น โดยไม่ปรากฏในภาษาอื่น ในงานวิจัยนี้กำหนดให้ใช้สายอักขระเฉพาะจากแบบจำลองตั้งแต่ 1-5 แกรม เท่านั้น

การระบุภาษาของข้อความ (Language Identification) คือ งานสำหรับระบุข้อความว่าเป็นภาษาใด ซึ่งส่วนใหญ่ภาษาที่ต้องการจะระบุมักเป็นภาษาที่ใช้อักษรโรมันเหมือนกัน เช่น ภาษาอังกฤษ ภาษาเยอรมัน ภาษาฝรั่งเศส

แบบจำลองภาษา (language model) คือ แบบจำลองที่ใช้ข้อมูลภาษาเชิงสถิติที่ได้จากข้อมูลการฝึกเป็นตัวกำหนดว่าแต่ละข้อความน่าจะเป็นภาษานั้นมากน้อยเพียงใด โดยเก็บข้อมูลหน่วยภาษา ตั้งแต่ 1-5 แกรม

ตัวจำแนกภาษา (language classifier) คือ โมดูลสำหรับใช้ระบุภาษาของข้อความหรือของคำ

ข้อมูลการฝึก (training data) คือ ข้อมูลภาษาที่ใช้สำหรับสร้างแบบจำลองทางสถิติ

ชุดของคำที่ใช้ทดสอบระบบ (test set) คือ คำสำหรับใช้ในการทดลองในแต่ละโมดูล โดยในขั้นตอนทดสอบของงานวิจัยนี้ จะแบ่งเป็น 2 ส่วน คือคำที่ใช้เป็นข้อมูลการฝึก 80% และคำที่ใช้เป็นข้อมูลทดสอบอีก 20%

บทที่ 2

การระบุภาษาของคำและข้อความ

การระบุภาษาของคำในวิทยานิพนธ์นี้มีความคล้ายคลึงกับการระบุภาษาของข้อความ (Language Identification) ในงานวิจัยอื่นๆ คือ มุ่งไปที่จะระบุคำหรือข้อความนั้นว่าเป็นภาษาใด โดยทำให้คอมพิวเตอร์ได้เรียนรู้ลักษณะสำคัญของภาษาจนสามารถระบุภาษาของคำหรือข้อความที่ให้ได้ ปัญหาที่พบเหมือนกัน คือจำเป็นต้องสร้างระบบที่มีประสิทธิภาพ ทำให้สามารถระบุภาษาได้ครอบคลุมทุกๆ ภาษาในโลก แม้การระบุภาษาของคำทับศัพท์จะไม่มีทางเสนอแนวทางแก้ปัญหาไว้อย่างชัดเจน แต่ในเรื่องการระบุภาษาของข้อความได้มีผู้ศึกษาไว้เป็นจำนวนมาก ผู้วิจัยจึงได้ทบทวนวรรณกรรมงานดังกล่าวไว้ในวิทยานิพนธ์นี้ เพื่อนำไปประยุกต์ใช้ต่อไป

การระบุภาษาของข้อความ หรือ Language Identification ถือว่าเป็นงานวิจัยที่สำคัญอย่างหนึ่ง และมีการพัฒนาจนเป็นระบบใช้กันเรื่อยมา เช่น ระบุว่าข้อความที่พบเป็นข้อความภาษาอังกฤษ ภาษาเยอรมัน ภาษาฝรั่งเศส ฯลฯ ซึ่งภาษาที่จะระบุส่วนมากเป็นภาษาที่มีความกำกวมของสายอักขระหรือตัวอักษร (string) เนื่องจากใช้ตัวอักษรชุดเดียวกัน เช่น ตัวอักษรโรมัน จึงต้องหาลักษณะที่สำคัญของสายอักขระ ที่จะบ่งบอกว่าสายอักขระหรือข้อความนั้นเป็นของภาษาใด โดยในการทดลองส่วนใหญ่นี้ใช้วิธีรับข้อมูลเข้า เป็นข้อความที่เขียนขึ้นด้วยภาษาใดภาษาหนึ่งเท่านั้น แล้วใช้ แบบจำลองภาษา (language model) ที่หาจากสายคำ (word sequence) เป็นระบบจำแนก (classifier) ภาษา ผลที่ได้ ดังตัวอย่างตารางที่ 2.1

ตารางที่ 2.1 ตัวอย่างการจำแนกภาษาของข้อความ

TEXT	LANGUAGE
InterPlas Conference will be the meeting place for delegates to be amidst top-notchd.....	English
Il faut faire demi-tour. Traversez le carrefour et lorsque vous voyez la station	French

ข้อความ (text) ในตารางใช้สายอักขระหรือตัวอักษรชนิดเดียวกัน คือตัวอักษรโรมัน แต่จริงๆ แล้วเป็นคนละภาษากัน คือภาษาอังกฤษและภาษาฝรั่งเศส

การระบุภาษาของข้อความนี้มีส่วนคล้ายคลึงกับงานที่ผู้วิจัยมุ่งจะทำ คือ เป็นการระบุที่มาของภาษาเหมือนกัน ทั้งสองงานนี้เป็นการจำแนกประเภทของภาษาจากสายอักขระที่ประกอบขึ้นจากตัวอักษรชุดเดียวกัน แต่ต่างกันตรงที่การระบุภาษาของข้อความเป็นการระบุว่าข้อความทั้งหมดนั้นเป็นภาษาอะไร ส่วนการระบุที่มาของภาษาของคำเป็นการระบุว่าคำที่พิจารณาเป็นคำไทยหรือทับศัพท์มาจากภาษาอะไร ดังแสดงในตัวอย่างตารางที่ 2.2

ตารางที่ 2.2 ตัวอย่างการระบุภาษาของคำ

คำ	ที่มา
ชานา	คำไทย
ปาปัว	คำทับศัพท์ภาษาฝรั่งเศส
ฟิวส์	คำทับศัพท์ภาษาอังกฤษ

การระบุภาษาของข้อความสามารถทำได้โดยการหาลักษณะที่สำคัญของภาษานั้นก่อน เช่น ลักษณะของคำศัพท์ (vocabulary term), ความยาวโดยเฉลี่ยของคำ (word average length), คุณสมบัติทางโครงสร้างและทางความหมายของคำ (syntactic และ semantic properties) ลักษณะเหล่านี้ของแต่ละภาษาจะไม่เหมือนกัน ทำให้สามารถระบุได้ว่าข้อความนั้นๆ เป็นภาษาอะไร ในทำนองเดียวกันหากเป็นคำไทยหรือคำทับศัพท์ ซึ่งมีที่มาจากภาษาต่างๆ กัน แม้ว่าจะอยู่ในรูปภาษาไทยเหมือนกัน แต่ก็น่าจะมีลักษณะสำคัญของภาษาที่ต่างกันด้วย ดังนั้นวิธีที่ใช้ในการระบุภาษาของข้อความก็น่าจะสามารถนำมาประยุกต์ใช้กับการระบุที่มาของภาษาในที่นี้ด้วย

จากปัญหาของการระบุภาษาของข้อความที่ผ่านมา แม้ว่ามีงานวิจัยที่ใช้วิธีทางภาษาศาสตร์ เช่น การระบุภาษาโดยใช้สายอักขระเฉพาะ (Unique character string identification) เช่น ในงานของ Churcher (1993) และ Churcher et al. (1994) และ การระบุภาษาของข้อความโดยใช้คำที่เกิดขึ้นบ่อยๆ ในแต่ละภาษาเป็นตัวตัดสิน (Frequent word recognition) เช่น ในงานของ Hayes (1993), Johnson (1993) และ Churcher et al. (1994) แต่ผู้วิจัยเห็นว่าวิธีการใช้สายอักขระเฉพาะมีข้อจำกัดคือไม่สามารถระบุภาษาของข้อความได้ทุกๆ ภาษา และเห็นว่าวิธีที่ถูกนำมาใช้มากที่สุดคือ วิธีทางสถิติ เนื่องจากเป็นวิธีที่ง่าย ไม่จำเป็นต้องใช้การวิเคราะห์ทางภาษาศาสตร์อย่างลึกซึ้ง เพียงแค่ใช้ข้อมูลการฝึกเพียงอย่างเดียว ทำให้ใช้ต้นทุนในการวิเคราะห์ต่ำกว่า มีประสิทธิภาพสูง และยังแก้ปัญหาข้อจำกัดของวิธีทางภาษาศาสตร์ได้ คือ ทำให้สามารถระบุภาษาได้ทุกภาษา เพียงแค่เตรียมข้อมูลสำหรับใช้ในการฝึกสำหรับภาษานั้นเท่านั้น

ระบบวิธีทางสถิติที่ใช้ในการระบุภาษาของข้อความส่วนใหญ่ เป็นการสร้างแบบจำลองทางภาษา (language model) ที่เก็บข้อมูลทางสถิติสำหรับภาษาต่างๆ จากคลังข้อมูลการฝึก (training corpus) เพื่อสร้างแบบจำลองภาษา (language model) หรือ แบบลักษณะ (profile) ของแต่ละภาษา แล้วนำแบบจำลองเหล่านี้มาใช้ทดสอบกับข้อความที่ต้องการจะระบุ หลังจากนั้นจึงจะประเมินผลเพื่อดูว่าข้อความที่ทดสอบนั้นใกล้เคียงกับแบบจำลองของภาษาใดมากที่สุด ดังนั้นระบบ (system) จะมี 2 ระบบ คือ ระบบสำหรับสร้างแบบจำลองภาษา เพื่อใช้เป็นตัวจำแนกภาษา (language classifier) และระบบสำหรับหาความน่าจะเป็นของข้อความ ที่จะเป็นภาษาใดภาษาหนึ่ง(เมื่อเทียบกับแบบจำลองภาษาที่มี)

เป้าหมายของการสร้างแบบจำลองภาษาก็เพื่อใช้ในการคาดเดาความเป็นไปได้ของสายคำ ถ้าสายคำเกิดขึ้นบ่อยๆ ในภาษา A ก็จะมีความเป็นไปได้สูงว่าเป็นข้อความในภาษา A นั้น

ลักษณะของแบบจำลองภาษา คือ

1. ยิ่งใช้ข้อมูลในการฝึกในแต่ละภาษามีจำนวนมาก แบบจำลองภาษาก็ยิ่งให้ผลดีขึ้นด้วย
2. คำที่ใช้คือ ค่าความถี่ และค่าความน่าจะเป็นของการเกิดคำร่วมกัน
3. เนื่องจากว่าคำในชุดของคำที่ใช้ทดสอบระบบ (test set) อาจไม่มีอยู่ในข้อมูลการฝึก ดังนั้นค่าความน่าจะเป็นของภาษาที่ต้องการจะระบุที่ออกมาจะเท่ากับ 0 จึงต้องมีการปรับแต่งค่า (smoothing) ระบบในการทดลอง

อย่างไรก็ตามแบบจำลองภาษา ที่ใช้ในการระบุภาษาของข้อความนั้นหาได้จากค่าทางสถิติของสายคำ (word sequence) คือดูคำที่รวมกันจนเป็นข้อความนั้น แต่แบบจำลองที่ใช้ในการระบุที่มาภาษาของคำนั้นจะหาจากค่าทางสถิติของสายอักขระ (character sequence) แทน คือ คาดเดาความเป็นไปได้ของสายอักขระหรือตัวอักษรที่เกิดขึ้นว่าตรงกับภาษาใด โดยในงานวิจัยนี้ทดลองด้วยการรับข้อมูลเข้าทีละคำ

เพื่อที่จะให้เข้าใจวิธีทางสถิติสำหรับระบุภาษาของคำ ผู้วิจัยจะทบทวนวรรณกรรมที่เกี่ยวข้องกับการระบุภาษาของข้อความก่อนว่ามีการทำอะไรบ้าง และวิธีใดเหมาะสำหรับนำมาประยุกต์ใช้กับงานการระบุภาษาของคำ

2.1 แนวทางในการสร้างแบบจำลองภาษา (Language model) หรือ ตัวจำแนกภาษาของข้อความ (text classifiers)

ในงานการระบุภาษาของข้อความนี้ สามารถทำได้หลายวิธี แต่ในที่นี้จะกล่าวถึง วิธีที่ใช้แบบจำลองมาร์คอฟ (Markov Model) หรือแบบจำลองเอ็นแกรม (N-Gram model) โดยละเอียด เพราะเป็นวิธีการที่ใช้ในการวิจัยนี้

2.1.1 แบบจำลองมาร์คอฟ หรือ แบบจำลองเอ็นแกรม

แบบจำลองเอ็นแกรม คือ แบบจำลองที่ใช้คำนวณค่าความน่าจะเป็นของสายอักขระ (character sequence) ที่เกิดขึ้นร่วมกันเป็นคำ หรือค่าความน่าจะเป็นของสายคำ (word sequence) ที่เกิดขึ้นร่วมกันเป็นประโยค โดยค่าความน่าจะเป็นของสายอักขระหรือคำ ประเมินได้จากคลังข้อมูลที่สร้างไว้

แกรม คือ หน่วยที่ใช้ในการสร้างแบบจำลอง อาจจะเป็นเสียง คำ หรือ อักขระก็ได้ และแกรมมีได้หลายขนาดแล้วแต่จะกำหนด ตั้งแต่ 1 จนถึง n

ในแบบจำลองเอ็นแกรมนี้ใช้ความยาวของสายอักขระและสายคำแตกต่างกัน ได้แก่ 2-แกรม 3-แกรม 4-แกรม ฯลฯ ถ้าจะประมาณค่าความน่าจะเป็นของสายคำหรือสายอักขระจากคลังข้อมูลโดยใช้วิธีเอ็นแกรม ผลที่ได้มีดังนี้

การประมาณค่าด้วย 2-แกรม (Probability bigram) คือ การประมาณค่าความน่าจะเป็นของสายอักขระ (สายคำ) ที่เกิดขึ้นร่วมกันว่ามีค่าเท่ากับผลคูณของความน่าจะเป็นที่จะพบอักขระ (คำ) ทีละ 2 ตัว (คำ) ติดกันในสายอักขระ (สายคำ) นั้น

การประมาณค่าด้วย 3-แกรม (Probability trigram) คือ การประมาณค่าความน่าจะเป็นของสายอักขระ (สายคำ) ที่เกิดขึ้นร่วมกันว่ามีค่าเท่ากับผลคูณของความน่าจะเป็นที่จะพบอักขระ (คำ) ทีละ 3 ตัว (คำ) ติดกันในสายอักขระ (สายคำ) นั้น

การประมาณค่าด้วย 4-แกรม (Probability quadigram) คือ การประมาณค่าความน่าจะเป็นของสายอักขระ (สายคำ) ที่เกิดขึ้นร่วมกันว่ามีค่าเท่ากับผลคูณของความน่าจะเป็นที่จะพบอักขระ (คำ) ทีละ 4 ตัว (คำ) ติดกันในสายอักขระ (สายคำ) นั้น

หรืออาจประมาณค่าความน่าจะเป็นจากความยาวของเอ็นแกรมมากกว่า 4-แกรม ก็ได้ ขึ้นอยู่กับความจำเป็นในการทดลอง แต่ระบบของเอ็นแกรมก็ยิ่งซับซ้อนมากขึ้นตามลำดับ

การประมาณค่าความน่าจะเป็นของสายอักขระ โดยการใช้อินแกรมดังที่กล่าวมา คือ การใช้สมมติฐานของมาร์คอฟ (Markov assumption) ว่า การปรากฏของตัวอักษรตัวหนึ่งขึ้นกับตัวอักษรก่อนหน้าเพียง $n-1$ ตัว เราจึงเรียกแบบจำลองลักษณะนี้ว่าเป็นแบบจำลองอินแกรม (n-gram model) หรือแบบจำลองมาร์คอฟ ซึ่งวิธีนี้มักนิยมใช้ในงานระบุภาษาของข้อความกันมาก เห็นได้จากงานของ Cavnar และ Trenkle (1994), Dunning (1994), Combrick และ Botha (1995), Sibun และ Reynar (1996), Piotrowski (2000), Peng et al. (2003) เนื่องจากสามารถใช้เพื่อระบุภาษาได้อย่างมีประสิทธิภาพและเรียบง่ายกว่า โดยสามารถประมาณได้ดังนี้

$$\text{ไบแกรม } P(c_1c_2c_3\dots c_n) = P(c_1) P(c_2|c_1) P(c_3|c_2) \dots P(c_n|c_{n-1})$$

$$\text{ไตรแกรม } P(c_1c_2c_3\dots c_n) = P(c_1) P(c_2|c_1) P(c_3|c_1c_2) \dots P(c_n|c_{n-2}c_{n-1})$$

$$\text{ควอดริแกรม } P(c_1c_2c_3\dots c_n) = P(c_1) P(c_2|c_1) P(c_3|c_1c_2) P(c_4|c_1c_2c_3) \dots P(c_n|c_{n-3}c_{n-2}c_{n-1})$$

หมายเหตุ: P แทน ค่าความน่าจะเป็น, c แทน อักขระหรือตัวอักษร และ $(c_1c_2c_3\dots c_n)$ แทน สายอักขระที่ประกอบด้วยอักขระตั้งแต่ 3 ตัวขึ้นไปจนถึง n ตัว

ส่วนความน่าจะเป็นของสายคำที่รวมกันเป็นประโยค $w_1w_2w_3\dots w_n$ หากประมาณค่าด้วยอินแกรมต่างๆ ผลที่ได้ดังนี้ โดย w แทน คำ n แทน จำนวนนับต่อไป P แทน ค่าความน่าจะเป็น และ $(w_1w_2w_3\dots w_n)$ แทน สายคำที่ประกอบด้วยคำมากกว่า 3 คำขึ้นไป

- ความน่าจะเป็นของประโยคโดยใช้วิธี 2-แกรม คือ

$$P(w_1w_2w_3\dots w_n) = P(w_1) P(w_2|w_1) P(w_3|w_2) \dots P(w_n|w_{n-1})$$

- ความน่าจะเป็นของประโยคโดยใช้วิธี คือ 3-แกรม คือ

$$P(w_1w_2w_3\dots w_n) = P(w_1) P(w_2|w_1) P(w_3|w_1w_2) \dots P(w_n|w_{n-2}w_{n-1})$$

- ความน่าจะเป็นของประโยคโดยใช้วิธี 4-แกรม คือ

$$P(w_1w_2w_3 \dots w_n) = P(w_1) P(w_2|w_1) P(w_3|w_1w_2) P(w_4|w_1w_2w_3) \dots P(w_n|w_{n-3}w_{n-2}w_{n-1})$$

ตัวอย่าง การประมาณค่าความน่าจะเป็นของประโยค "He like to eat"

ผู้วิจัยใช้สมการประมาณค่าความน่าจะเป็นพื้นฐานดังนี้

$$P(w_1w_2 \dots w_T) = \prod_{i=1}^T P(w_i|w_1 \dots w_{i-1})$$

โดยที่ P คือ ค่าความน่าจะเป็น (Probability) ซึ่งประมาณได้จากคลังข้อมูล

T คือ จำนวนของคำ

i คือ ลำดับของคำโดยเริ่มต้นที่ลำดับที่ 1

$P(w_i | w_1 \dots w_{i-1})$ คือ ความน่าจะเป็นของคำ w_i หลังจากเกิดคำ $w_1 w_2 \dots w_{i-1}$ ก่อนหน้านี้แล้ว

ความน่าจะเป็นของประโยคนี้ $P(w_1 w_2 \dots w_T)$ สามารถประมาณได้ โดยถือว่าการปรากฏของคำ w_i นั้นขึ้นอยู่กับจำนวนคำข้างหน้า $n-1$ ตัวเท่านั้นหรือขึ้นอยู่กับขนาดของเอ็นแกรม ดังนั้นถ้าหากประมาณค่าความน่าจะเป็นของประโยคนี้ โดยใช้ 2-แกรม จะปรับเปลี่ยนสมการดังนี้

$$P(w_1 w_2 \dots w_T) = P(w_1 | < s >) P(w_2 | w_1) \dots P(w_i | w_{i-1})$$

$P(w_1 | < s >)$ หมายถึง ความน่าจะเป็นของคำที่หนึ่งเมื่อเกิดเป็นคำแรกของประโยค ซึ่งในที่นี้คือช่องว่าง

$P(w_2 | w_1)$ หมายถึง ความน่าจะเป็นของคำ w_2 หลังจากเกิดคำ w_1

$P(w_i | w_{i-1})$ หมายถึง ความน่าจะเป็นของคำ w_i หลังจากเกิดคำ w_{i-1}

ดังนั้น จากสูตรประมาณค่าความน่าจะเป็นด้วย 2-แกรม หากต้องการหาความน่าจะเป็นของประโยค "He like to eat" โดยใช้เอ็นแกรมในระดับของคำ ผลออกมาคือ $P(\text{He} | < s >) P(\text{like} | \text{He}) P(\text{to} | \text{like}) P(\text{eat} | \text{to})$

ส่วนค่าความน่าจะเป็นจะหาได้จากคลังข้อมูล เช่น

$$P(\text{like} | \text{He}) = \frac{c(\text{He-like})}{c(\text{He})} \text{ โดยที่ } c \text{ คือจำนวนนับ}$$

จากสูตรนี้หมายถึง ค่าความน่าจะเป็นของ like เมื่อเกิดร่วมกับคำว่า He คำนวณได้จากนำจำนวนนับของ He ที่เกิดร่วมกับคำว่า likeหารด้วยจำนวนนับของการเกิด He เดี่ยวๆ

วิธีการใช้แบบจำลองเอ็นแกรมนี้เป็นวิธีทางสถิติที่นิยมใช้กันมากที่สุด เพราะเป็นวิธีที่เรียบง่าย มีประสิทธิภาพสูง และเหมาะสำหรับวิเคราะห์ภาษา สามารถใช้ระบุภาษาได้ดีกว่าวิธีอื่นๆ โดยเมื่อดูจากทบทวนวรรณกรรมต่อไปจะแสดงให้เห็นประสิทธิภาพของแบบจำลองเอ็นแกรม วิธีการใช้เอ็นแกรมในงานระบุภาษาและผลการทดลอง

2.1.2 ตัวอย่างงานวิจัยของการระบุภาษาของข้อความที่ใช้แบบจำลองเอ็นแกรม

แบบจำลองเอ็นแกรมถูกนำมาใช้ในงานต่างๆ มากมาย เช่น

Dunning (1994) ใช้ แบบจำลองมาร์คอฟ (Markov model) ในระดับตัวอักษรเพื่อระบุภาษาของข้อความที่เขียนด้วยอักษรโรมัน เหตุผลที่เขาเลือกแบบจำลองนี้เพราะ คิดว่าการระบุ

ภาษาของข้อความเป็นเรื่องง่าย ไม่จำเป็นต้องใช้แบบจำลองอื่นๆ ที่ยุ่งยากกว่านี้ เขากล่าวไว้ว่า สามารถระบุภาษาได้เลยแม้ใช้ข้อมูลทดสอบเพียงแค่ 10 ตัวอักษรเท่านั้นแต่จะให้ผลดีมากที่สุดหากใช้ตัวอักษร 50 ตัวหรือมากกว่านั้นก็ได้ ส่วนจำนวนข้อมูลการฝึกก็ใช้น้อยเพียงแค่ว่าข้อความที่ประกอบด้วย 2-3 พันคำก็เพียงพอ ในการทดลองเขาใช้ข้อมูลการฝึกเป็นเอกสารข้อความที่มีขนาดตั้งแต่ 1-50 กิโลไบต์ และใช้ข้อมูลทดสอบเป็นข้อความจำนวน 50 ข้อความ โดยขนาดตั้งแต่ 10-500 ไบต์ ลักษณะของแบบจำลองของเขา คือ การดูความน่าจะเป็นของตัวอักษรหรืออักขระที่เกิดขึ้นภายในข้อความ โดยจะดูต่อเนื่องกันทั้งข้อความ ในผลการทดลองเขาก็เปรียบเทียบให้ดูว่า หากใช้ข้อมูลการฝึก และข้อมูลทดสอบยิ่งมาก ค่าความแม่นยำก็ยิ่งสูงตามด้วย โดยอาจสูงถึง 99.9% แต่ก็พบปัญหาคือ เมื่อทดสอบกับภาษาอังกฤษ ฝรั่งเศส เยอรมัน และภาษาสเปน อาจระบุภาษาผิดพลาดได้เนื่องจากแต่ละภาษามีรากฐานคล้ายกัน

Sibun and Reynar (1996) ใช้วิธีแบบจำลองเอ็นแกรม เพื่อระบุภาษาที่เขียนด้วยอักษรโรมันซึ่งมีทั้งหมด 27 ภาษา เช่น ภาษาดัตช์ อังกฤษ อิตาลี ลาติน สเปน เป็นต้น เขาใช้แบบจำลองเอ็นแกรมระดับตัวอักษรตั้งแต่ 1-3 แกรม เพื่อสร้าง probability distribution หรือแบบจำลองภาษาจากข้อมูลการฝึก วิธีการสร้างคลังข้อมูลการฝึกของเขาจะมีลักษณะเหมือนกับงานของ Dunning คือเป็นการดูทีละเอกสารข้อความต่อเนื่องกัน โดยในเอกสารจะต้องมีการตัดสัญลักษณ์และเครื่องหมายต่างๆ (markup, accented character and symbol) ออกก่อนเพื่อรับเฉพาะตัวอักษรเท่านั้น และข้อมูลทดสอบก็เป็นเอกสารข้อความเช่นเดียวกัน ในการทดลองเขาใช้วิธีเอ็นโทรปีสัมพัทธ์ (relative entropy) สำหรับวัดค่าความสัมพันธ์ระหว่างข้อมูลทดสอบ (test set) กับชุดข้อมูลการฝึก (training) หรืออาจกล่าวได้ว่า เป็นการวัดค่าความน่าจะเป็นของชุดข้อมูลทดสอบว่าตรงกับชุดข้อมูลการฝึกของภาษาใดภาษาหนึ่ง ผลของการทดลองระบุภาษาทั้งหมด 27 ภาษาโดยใช้ 1-แกรม และ 2-แกรม ค่าความแม่นยำโดยเฉลี่ยมากกว่า 90%

Piotrowski (2000) ระบุภาษาของข้อความที่เขียนด้วยตัวโรมันเพื่อการค้นคืนสารสนเทศว่าเป็นภาษาเยอรมัน อังกฤษ ฝรั่งเศส หรืออิตาลี เขาดูความน่าจะเป็นในระดับของสายอักขระหรือตัวอักษรที่ประกอบขึ้นเป็นข้อความ โดยวิธีที่เขาใช้คือ 2-แกรม ร่วมกับ มาร์คอฟ เซน (Markov Chain) เพื่อสร้างแบบจำลองภาษากับทุกภาษาที่ต้องการระบุ ในแบบจำลองภาษาจะเป็นที่รวมข้อมูลความน่าจะเป็นของสายอักขระที่เกิดขึ้นร่วมกัน ผลการทดลองเขาได้เปรียบเทียบการใช้ขนาดของข้อมูลให้ดูว่า หากใช้ขนาดของข้อมูลการฝึกมาก ค่าความแม่นยำที่ได้ยิ่งสูงตามด้วย เช่น ขนาดของข้อมูลการฝึกประมาณ 20 ไบต์ (bytes) ให้ผลการทดลองค่าความแม่นยำเฉลี่ย 98.73% แต่ถ้าใช้ขนาดของข้อมูลการฝึกประมาณ 50 ไบต์ ให้ผลการทดลองค่าความแม่นยำสูงถึง 99.69%

Combrinck และ Botha (1995) เสนอวิธีสำหรับระบุภาษาของข้อความ โดยทดสอบกับภาษาทั้งหมด 12 ภาษาทั้งภาษาแถบอัฟริกันและยุโรปเียน ระบบของเขาจะเป็นการสร้างแบบจำลองภาษาจากข้อมูลการฝึกด้วยเอ็นแกรม แต่เขาเรียกวิธีนี้ว่าทรานซิชัน-เวกเตอร์ (transition vector) โดยหลักการคือใช้เพื่อดูสายอักขระ หรือดูการรวมของสายอักขระตั้งแต่ 2-3 ตัวในข้อความแต่ละภาษาแล้วจึงสร้างแบบจำลองของภาษานั้นขึ้นมา ระบบของพวกเขาประกอบไปด้วย 3 ขั้นตอน คือ

1. pre-processor เพื่อตัดส่วนที่ไม่ใช่ตัวอักษรออกไป เช่น สัญลักษณ์
2. training module เพื่อสร้างแบบจำลองในแต่ละชุดข้อมูลการฝึก โดยในแบบจำลองจะอยู่ในรูปของทรานซิชัน-เวกเตอร์ ซึ่งจะให้ภาพของการรวมกันของตัวอักษร 2 และ 3 ตัว เช่นเดียวกับเอ็นแกรม 1-2 แกรม
3. recognition module เป็นส่วนที่ใช้ระบุภาษาของข้อความ จาก ทรานซิชัน-เวกเตอร์ที่เด่นๆ ในแต่ละภาษา ซึ่งในที่นี้คือ ทรานซิชัน-เวกเตอร์ ที่มีความถี่ของการรวมกันของตัวอักษรสูงกว่าภาษาอื่นๆ

Combrinck และ Botha (1995) ยังเสนอว่า ระบบของเขาสร้างขึ้นเพื่อรองรับความยาวของสายอักขระมากกว่า 3- แกรม แต่ทรานซิชัน เวกเตอร์ ที่ให้ผลที่สุด คือ 3-แกรม เขาอ้างถึงผลการทดลองว่าตัวจำแนกภาษาของเขาอาจสามารถระบุภาษาได้ถูกต้อง 100% โดยผลที่ดีที่สุดคือ เมื่อใช้กับข้อมูลทดสอบประมาณ 200-300 ตัวอักษรเท่านั้น แต่ระบบของเขายังมีปัญหาคือขาดเสถียรภาพในการใช้งานจริง เขายังเปรียบเทียบให้เห็นว่าหากใช้ข้อมูลขนาดใหญ่เกินไปทั้งข้อมูลการฝึกและข้อมูลทดสอบ ค่าการวัดผลทดลอง (Relative Performance Measure) ที่ออกมาจะต่ำลงตามลำดับ

Peng et al. (2003) เสนอวิธีสำหรับระบุประเภทของข้อความ (classification or categorization) ทั้งหมด 24 ประเภทในภาษาจีนจาก TREC และภาษาญี่ปุ่นจาก NTCIR ในการทดลองเขาใช้ข้อมูลการฝึกทั้งหมด 310,355 ข้อความ แต่ละประเภทมีประมาณ 1,747 - 83,668 ข้อความ และข้อมูลทดสอบมีทั้งหมด 10,000 ข้อความ ซึ่งข้อความในแต่ละประเภทจะมีประมาณ 56 - 2,656 ข้อความ

ลักษณะของแบบจำลองภาษาของเขา คือ การคาดเดาความน่าจะเป็นของสายอักขระที่เกิดขึ้นร่วมกันในข้อความแต่ละประเภท โดยดูที่ละข้อความต่อเนื่องกัน ในแบบจำลองเอ็นแกรมที่ใช้มีตั้งแต่ 1-8 แกรม หลังจากที่ได้แบบจำลองภาษาแต่ละประเภทแล้ว เขาใช้แบบจำลองภาษาสำหรับเป็นตัวจำแนกภาษาและประเภทของข้อความ (text classifier) ในตัวจำแนกภาษาและประเภทของข้อความนี้จะดูลักษณะที่สำคัญของภาษาเช่น คำศัพท์ (vocabulary term),

ความยาวเฉลี่ยของคำ (word average length), ลักษณะทางวากยสัมพันธ์และอรรถศาสตร์ (syntactic and semantic properties) สำหรับใช้เพิ่มการตรวจสอบที่ระบุประเภท (categories) ของข้อความ เนื่องจากข้อความในแต่ละประเภทจะมีลักษณะทางภาษาที่ต่างกัน หากใช้แบบจำลองไบแกรมจะช่วยให้คอมพิวเตอร์สามารถพบลักษณะทางภาษาที่ต่างกันได้และใช้ระบุประเภทของข้อความเหล่านั้นได้ ตัวอย่างสมมติ เช่น เมื่อใช้ 2-แกรม เราอาจจะพบสายอักขระ “ic” และ “er” ซึ่งปรากฏเฉพาะในข้อความประเภทข่าวธุรกิจเป็นส่วนใหญ่ ดังนั้น หากนำสายอักขระเหล่านี้มาตรวจสอบ ก็จะสามารถระบุประเภทได้ถูกต้องมากยิ่งขึ้น ในการทดลองเขาทดสอบกับข้อความที่ต้องการระบุประเภทด้วยเอ็นแกรมเหมือนกับที่ใช้ในการฝึก ผลการทดลองประมาณสูงกว่า 80% และเขายังเปรียบเทียบผลการใช้ขนาดเอ็นแกรมต่างๆ ด้วย โดยขนาดของเอ็นแกรม 3-แกรมให้ผลดีที่สุด นอกจากนี้ Pang et al. ยังวิจารณ์ไว้ว่าการสร้างแบบจำลองภาษาที่ใช้สายอักขระ (character level) นี้จะสามารถใช้กับงานการระบุประเภทของข้อความได้ดีกับภาษาที่ไม่มีการเขียนแยกคำ โดยเฉพาะในภาษาแถบเอเชีย มากกว่าแบบจำลองที่ใช้สายคำ (word level)

ในงานทบทวนวรรณกรรมเหล่านี้ใช้วิธีเอ็นแกรมในระดับตัวอักษร เพื่อระบุภาษาของข้อความ แต่ยังมีบางงานที่เสนอการใช้เอ็นแกรมในระดับของคำ เช่น

งานของ Cavnar และ Trenkle (1994) ซึ่งเป็นการระบุหาประเภท (categories) ของข้อความในแต่ละภาษา เช่น ระบุว่าเอกสารนั้นเป็นเอกสารประเภทใด เช่น ด้าน security ด้าน AI ด้าน compilers ด้าน compression และด้านกราฟฟิค งานของเขาจะคล้ายคลึงกับงานที่ทบทวนมาก่อนหน้านี้คือ เป็นการดูทีละข้อความต่อเนื่องกัน ทั้งในข้อมูลการฝึกและข้อมูลทดสอบ แบบจำลองที่เขาใช้คือ แบบจำลองเอ็นแกรมของสายคำ โดยเขาเลือกใช้ตั้งแต่ 1-5 แกรม เพื่อสร้างแบบลักษณะหรือแบบจำลองภาษาของข้อความแต่ละประเภทภายในแบบจำลองภาษาจะแสดงค่าความถี่และค่าการเกิดร่วมกันของเอ็นแกรมทั้งหมด หลังจากได้ข้อมูลการฝึกของแต่ละประเภทแล้ว เขาก็จะเปรียบเทียบข้อมูลทดสอบกับข้อมูลการฝึกแต่ละประเภทว่า มีค่าความน่าจะเป็นที่จะตรงกันมากน้อยเพียงใด โดยใช้วิธี out-of-place measure หรือค่าความห่าง (distance measure) แล้วเลือกประเภทที่มีค่าความห่างกับข้อมูลทดสอบน้อยที่สุดเป็นคำตอบ ผลจากการทดลองระบุภาษาได้ถูกต้องมากกว่า 90% โดยผลจะต่างกันขึ้นอยู่กับขนาดของข้อมูลการฝึกในแต่ละประเภทและข้อมูลทดสอบ ยิ่งข้อมูลการฝึกและข้อมูลทดสอบมาก ค่าความแม่นยำก็สูงขึ้นด้วย ในงานของเขาจะเป็นการใช้เอ็นแกรมต่างกับงานที่ผู้วิจัยจะทำคือ เป็นการระบุหาประเภทแต่ผู้วิจัยเสนองานการระบุภาษา

จากบททวนวรรณกรรมแนวทางแก้ปัญหาการระบุภาษาและประเภทของข้อความที่ผ่าน มา แบบจำลองที่ผู้วิจัยเลือกใช้คือ แบบจำลองเอ็นแกรม (n-gram model) เพราะ จากบททวน วรรณกรรมที่ผ่านมาเห็นว่า ผลจากการทดลองที่ได้จากงานการระบุประเภทของข้อความด้วย แบบจำลองนี้สูงมากกว่า 90% โดยใช้เพียงแค่ 2-แกรมเท่านั้น และเมื่อเปรียบเทียบกับแบบจำลอง ภาษาเอ็นแกรม (LM) กับวิธีการอื่นๆ เช่น Naïve Bays classifier (NB), adhoc n-gram (OOP)¹ และ support vector machine classifier (SVM) พบว่า แบบจำลองภาษา ให้ผลดีที่สุด เห็นได้จากตารางที่ 2.3 ซึ่ง Peng ได้แสดงผลการทดลองเปรียบเทียบค่าความน่าจะเป็นของการ ระบุประเภทของข้อความประเภทหนึ่งด้วยการใช้แบบจำลองภาษานาอีฟ เบย์ และวิธีการอื่นๆ โดยอ้างอิงผลการทดลองระบุประเภทของข้อความของ Aizawa (2001) ด้วย แสดงให้เห็นว่าเมื่อ ระบุประเภทในข้อความภาษาจีนด้วยแบบจำลองภาษา จะให้ค่าความถูกต้องสูงที่สุด

ตารางที่ 2.3 ตัวอย่างการเปรียบเทียบความน่าจะเป็นของการระบุประเภทของข้อความ (Peng, et al. 2003 อ้างถึงใน Aizawa (2001))

LM	NB	OOP	SVM
Chinese Character Level (ระบุข้อความภาษาจีนในระดับตัวอักษร)			
0.868	0.856	0.8087	0.817
Japanese Byte Level (ระบุข้อความภาษาญี่ปุ่นในระดับตัวอักษร)			
0.84	0.66	0.4990	85%

นอกจากนี้ ผู้วิจัยเลือกใช้แบบจำลองเอ็นแกรมเพื่อพัฒนาระบบการระบุภาษาของคำ โดยมีเหตุผลดังนี้

1. เอ็นแกรมเป็นวิธีทางสถิติที่เป็นพื้นฐานและเรียบง่ายมากที่สุด ซึ่งจะช่วยลดต้นทุน ในการทำระบบการระบุภาษา
2. สามารถนำวิธีเอ็นแกรมมาใช้ร่วมกับวิธีอื่นๆได้ เนื่องจากจะช่วยเพิ่มประสิทธิภาพ ของระบบให้ดีขึ้น

¹ เป็นวิธีประมาณค่าความน่าจะเป็นโดยใช้แบบจำลองเอ็นแกรมเหมือนกับ LM แต่ต่างกับ LM ตรงที่ วิธี OOP จะเลือกข้อความที่ทดสอบว่าตรงกับประเภทใดจากค่าความห่าง (distance measure) โดยที่เลือกใช้ คือ out-out-place (OOP) หรือ distance (d,c) ส่วน วิธี LM เลือกข้อความที่ทดสอบว่าตรงกับประเภทใดจาก ค่า likelihood หรือ $\Pr(d|c)$

3. วิธีเอ็นแกรมสามารถใช้ระบุภาษาได้ทุกภาษา ในกรณีระบุภาษาที่ใช้สายอักขระที่เป็นตัวอักษรประเภทเดียวกัน แต่ต่างภาษากัน เช่น อักษรโรมันที่ใช้เขียนภาษาอังกฤษ เยอรมัน อิตาลี หรืออักษรไทย ที่ใช้เขียนคำทับศัพท์ภาษาอังกฤษ ญี่ปุ่นและฝรั่งเศส เป็นต้น
4. เนื่องจากรูปแบบของข้อมูลนำเข้าเพื่อทดสอบและข้อมูลการฝึกควรจะเรียบง่ายมากที่สุด ไม่มีการประมวลผลเบื้องต้น (pre-process) เช่น การ encoding ข้อมูลหรือการใส่ tag ลงบนข้อมูล ดังนั้นระบบเอ็นแกรมจึงเป็นวิธีที่ดีที่สุดที่จะลดความยุ่งยากในการพัฒนาระบบ
5. ในระบบเอ็นแกรม เพียงใช้คลังข้อมูลจำนวนน้อยก็สามารถระบุภาษาของคำได้ถูกต้อง

จากการทบทวนวรรณกรรม ผู้วิจัยพบว่าวิธีรวมสายอักขระซึ่งจะกล่าวต่อไปนี้เป็นอีกวิธีหนึ่งที่มีผู้วิจัยเคยทำไว้และใช้ในงานระบุภาษาของข้อความ จึงเป็นที่น่าสนใจว่าจะสามารถใช้ในงานระบุภาษาของคำในภาษาไทยได้หรือไม่ ผู้วิจัยจึงต้องการทบทวนวรรณกรรมแนวทางในการสร้างวิธีรวมสายอักขระในส่วนต่อไป เพื่อศึกษาวิธีการและการนำวิธีรวมสายอักขระมาประยุกต์ใช้เพื่อเพิ่มประสิทธิภาพของงานระบุภาษาของคำ

2.2 แนวทางในการสร้างวิธีรวมสายอักขระเฉพาะ (Unique Character Combinations)

วิธีรวมสายอักขระคือ การดูความเป็นไปได้ของจำนวนสายอักขระ (character sequence) ที่เกิดขึ้นตั้งแต่ 3 ตัวขึ้นไป ซึ่งจะแสดงลักษณะเฉพาะของแต่ละภาษาให้เห็นชัดเจน คือเป็นสายอักขระสั้นๆ ที่เกิดขึ้นในภาษาใดภาษาหนึ่งเท่านั้น โดยไม่เกิดขึ้นซ้ำกับภาษาอื่นๆ เลย สายอักขระเฉพาะเหล่านี้พบได้จากคลังข้อมูลภาษา โดยเราจะได้ข้อมูลการเกิดร่วมกันของตัวอักษร พยัญชนะนำ-ตาม สระหรือตัวสะกดที่เป็นลักษณะเด่นของภาษา เห็นได้จากตารางที่ 2.4

ตารางที่ 2.4 ตัวอย่างการใช้วิธีรวมสายอักขระเฉพาะ (Dunning, 1994)

Language	Unique Character
Dutch	“vnd”
English	“ery”
French	“eux”
Italian	“cchi”
Serbo-croat	“lj”

ในอักขระเฉพาะที่กล่าวไว้ในตาราง เป็นสายอักขระเฉพาะที่ประกอบขึ้นจาก 3 ตัวอักษรที่เกิดขึ้นในท้ายคำของภาษานั้นๆ โดยจะไม่ปรากฏซ้ำกับภาษาอื่น ๆ เลย นอกจากนี้ในแต่ละภาษาสามารถมีสายอักขระเฉพาะมากกว่าหนึ่งชุด และมีสายอักขระมากกว่า 3 ตัวติดกันก็ได้

Dunning ยังแสดงให้เห็นปัญหาของวิธีนี้ว่า ถ้าในข้อความนั้นมีคำภาษาอังกฤษอย่าง Pinocchio แล้วจะใช้สายอักขระเฉพาะ cchi สรุปว่าเป็นภาษาอิตาลีก็อาจผิดได้ กรณีแบบนี้เป็นเพราะภาษาเหล่านี้ใช้สายอักขระชุดเดียวกัน และยืมคำข้ามภาษามาโดยไม่ต้องทับศัพท์ด้วยตัวเขียนอีกภาษา เช่นจากตัวอย่าง เมื่อทดลองด้วยข้อมูลการฝึกในภาษาอิตาลี อาจทำให้ได้ cchi เป็นสายอักขระเฉพาะ เพราะบังเอิญทดลองฝึกในข้อมูลภาษาอังกฤษแล้วไม่มีคำว่า Pinocchio อยู่ แต่พอเอาไปใช้จริง ในข้อความภาษาอังกฤษมีคำว่า pinocchio อยู่ ระบบจึงผิดพลาดได้ เพราะมีข้อมูลว่า cchi เป็นสายอักขระเฉพาะของภาษาอิตาลีซึ่งก็จริง แต่บังเอิญข้อความที่ทดสอบยืมคำจากภาษาอิตาลีมาใช้และทั้งสองภาษาใช้ตัวอักษรโรมันเหมือนกัน จึงเป็นปัญหาที่เกิดขึ้น แต่ในภาษาไทย เราทับศัพท์ภาษาต่างประเทศ คำยืมจึงเขียนด้วยตัวอักษรไทย จึงไม่มีปัญหาที่เกิดจากการยืมคำและคงรูปเดิมไว้

ลักษณะของงานแบบนี้มีมากมาย เช่น จากงานของ Churcher et al. (1994) พวกเขาใช้วิธีสังเกตจากสายอักขระเฉพาะที่ประกอบขึ้นจากตัวอักษรเพื่อระบุภาษาของข้อความในภาษาเขียนด้วยอักษรโรมันถึง 9 ภาษา ในการทดลองเขาสามารถหาสายอักขระเฉพาะ 2 ตัวอักษรติดกันในแต่ละภาษาได้ แต่พวกเขาวิจารณ์ถึงผลการทดลองของวิธีนี้ว่า ไม่มีประสิทธิภาพเพียงพอ เพราะสายอักขระเฉพาะที่ใช้เพียงแค่ 2-แกรม อาจยาวไม่เพียงพอสำหรับกำหนดให้เป็นสายอักขระเฉพาะ และหากใช้ทดสอบกับคำโดยใช้แค่ 2-แกรม ในคำที่ทดสอบคำๆ หนึ่งจะมีสายอักขระที่ตรงกับสายอักขระเฉพาะในหลายๆ ภาษาได้

จากการศึกษาวิธีการใช้สายอักขระเฉพาะ ผู้วิจัยสามารถสรุปปัญหาของวิธีนี้ได้ดังนี้

1. สายอักขระเฉพาะส่วนใหญ่ประกอบด้วยอักขระแค่ 2 หรือ 3 ตัวกันเท่านั้น สายอักขระเฉพาะยาวๆ มีแนวโน้มที่จะเกิดขึ้นยาก เพราะยิ่งต้องการสายอักขระเฉพาะยาวๆ ก็จำเป็นต้องใช้คลังข้อมูลขนาดใหญ่มากขึ้น
2. สายอักขระเฉพาะมีจำกัดคือ ไม่สามารถใช้ระบุภาษาได้ทุกๆ ข้อความหรือทุกคำ เพราะบางข้อความหรือบางคำอาจไม่มีคำที่ตรงกับสายอักขระเฉพาะเลย จึงต้องนำวิธีการระบุภาษาอื่นๆ เข้ามาช่วย ในที่นี้ผู้วิจัยแก้ปัญหาโดยใช้วิธีทางสถิติร่วมกับการสังเกตจากสายอักขระเฉพาะ

ดังนั้นจากการทบทวนวรรณกรรมทั้งแบบจำลองเอ็นแกรมและวิธีรวมสายอักขระที่ผ่านมา ผู้วิจัยจึงเลือกใช้วิธีเอ็นแกรมร่วมกับวิธีการสังเกตจากสายอักขระเฉพาะ ดังนั้นผู้วิจัยจึงมุ่งที่จะ

นำเสนอแนวทางในการใช้สายอักขระเฉพาะและการสร้างแบบจำลองภาษา เพื่อใช้ในงานการพัฒนา
ระบบการระบุภาษาของคำนี้

2.3 แนวทางในการจัดทำระบบการระบุภาษาของคำ

จากการศึกษาแนวทางต่างๆของการระบุภาษาของข้อความ ในการพัฒนาระบบการระบุ
ภาษาของคำในวิทยานิพนธ์นี้ ผู้วิจัยจึงประยุกต์วิธีที่เป็นไปได้และเหมาะสมสำหรับภาษาไทย
เพื่อกำหนดแนวทางในการพัฒนาระบบการระบุภาษาของคำไทยและคำทับศัพท์ ดังนี้

1. ใช้วิธีการสังเกตจากสายอักขระเฉพาะ (Unique character sequence) เนื่องจาก
เห็นว่าคำทับศัพท์หรือคำไทยมักมีรูปแบบการเกิดที่ตายตัว เป็นไปได้ที่จะมีลักษณะเฉพาะและมี
ตัวอักษรหรือพยัญชนะที่ให้ลักษณะเด่นของภาษาได้ชัดเจน เราจึงสามารถระบุภาษาของคำ
ทับศัพท์จากสายอักขระ (character sequence) ที่เป็นตัวอักษรภาษาไทย ตัวอย่างเช่น ในคำ
ทับศัพท์ภาษาอังกฤษ อาจจะมี “ตอร์” เป็นชุดอักขระเฉพาะ ในคำว่า คอมพิวเตอร์ มิเตอร์
ฟิลเตอร์ “ล์ม” ในคำว่า ฟิล์ม และ “อล” เป็นชุดอักขระเฉพาะ ในคำว่า ดิจิตอล ฟุตบอล โดยใน
สายอักขระเหล่านี้จะปรากฏเฉพาะในคำทับศัพท์ภาษาใดภาษาหนึ่งเท่านั้น จุดประสงค์หลักใน
การทำสายอักขระเฉพาะ คือต้องการลดขั้นตอนของระบบ โดยไม่จำเป็นต้องนำไปเทียบกับ
แบบจำลองภาษา และเป็นการเพิ่มประสิทธิภาพของระบบให้ถูกต้องมากยิ่งขึ้น

แต่ทว่าในการใช้สายอักขระเฉพาะไม่สามารถแก้ปัญหาคำบางคำอาจไม่มีสายอักขระตรงกับสายอักขระเฉพาะของภาษาใดๆ เลย
นอกจากนี้ คำ 1 คำอาจมีสายอักขระที่ไปตรงกับสายอักขระเฉพาะมากกว่าหนึ่งภาษา ตัวอย่าง
สมมติ คำว่า คอมพายเลอร์ ตรงกับสายอักขระเฉพาะของคำทับศัพท์ภาษาอังกฤษว่า “คอม”
และตรงกับสายอักขระเฉพาะของคำทับศัพท์ภาษาฝรั่งเศส “อร์” ด้วย จึงทำให้ตัดสินใจจากการ
เทียบสายอักขระเฉพาะไม่ได้ว่าควรเป็นภาษาใด ในกรณีนี้ระบบที่พัฒนาจะนำคำไปทดสอบ
ต่อไปในแบบจำลองภาษาของคำทับศัพท์ที่เป็นไปได้ของทั้ง 2 ภาษา ในที่นี้จึงนำไปทดสอบ
กับคำทับศัพท์ภาษาอังกฤษและคำทับศัพท์ภาษาฝรั่งเศส

ในงานวิจัยนี้ หากสายอักขระเฉพาะด้วย 1-5 แกรม คือพิจารณาสายอักขระที่เป็นอักษร
1-5 ตัวติดกัน เพราะวิธีเหล่านี้น่าจะให้ลักษณะเด่นของอักขระ (character) ในแต่ละภาษาที่
ชัดเจนมากกว่าเพียงแค่ 2-แกรม ซึ่งเป็นการพิจารณาอักษรเพียงสองตัวติดกัน ส่วนข้อมูล
การฝึก (training data) สำหรับสร้างสายอักขระเฉพาะ ใช้ข้อมูลเดียวกันกับสำหรับฝึก (training)
ทำแบบจำลองภาษา

2. ใช้แบบจำลองทางภาษา คือ แบบจำลองเอ็นแกรม ในที่นี้จะทดลองโดยใช้ 2-5 แกรม นอกจากนี้วิธีเอ็นแกรมที่เลือกใช้จะต่างกับงานทั่วไปคือ เปลี่ยนจากการดูทีละคำ (w) เป็นการดูเอ็นแกรมทีละอักขระหรือตัวอักษร (c) เพราะในจุดประสงค์ของการทดลองผู้วิจัยมุ่งจะระบุภาษาของคำทับศัพท์ในภาษาไทยเป็นคำๆ ดังนั้นในข้อมูลทั้งข้อมูลการฝึกและข้อมูลสำหรับทดสอบจึงใช้ทีละคำ และหากจะใช้วิธีเอ็นแกรมสำหรับหาลักษณะเด่นของคำในแต่ละภาษา จะต้องใช้เอ็นแกรมทีละเอียดยกกว่าระดับคำ คือ ระดับสายอักขระเท่านั้น จึงจะได้ผลดีที่สุด โดยกำหนดให้ $w = c_1 c_2 c_3 \dots c_n$ และ c คือ สายอักขระหรือตัวอักษร ในสายอักขระผู้วิจัยจะนับทั้งพยัญชนะวรรณยุกต์ และสระ เป็นอักขระตัวหนึ่ง เช่น คำว่า “แอกเซส” จะแยกเป็น แ อ ก เซ และ ส แล้วจึงนับความถี่และเก็บค่าความน่าจะเป็นตามเอ็นแกรม โดยเก็บเป็น 3-แกรม ตัวอย่างสมมติดังนี้

ข้อมูลไตรแกรม

ตัวอักษร	ความถี่
แอก	180
อกเ	954
กเซ	123
เซส	1180

โดยผู้วิจัยจะเก็บข้อมูลลักษณะแบบนี้จากทุกๆ คำในคลังข้อมูล เป็นค่าทางสถิติ

ในการประมาณค่าความน่าจะเป็นของคำ ผู้วิจัยจะประมาณโดยใช้เอ็นแกรมด้วย เช่นในกรณีที่ใช้ 3-แกรม จากคำว่า “แอกเซส” สามารถประมาณค่าความน่าจะเป็นของคำด้วย 3-แกรมได้ดังนี้

$$P(\text{แ}) P(\text{อ|แ}) P(\text{ก|แ อ}) P(\text{เ|อ ก}) P(\text{ซ|ก เ}) P(\text{ส|เซ})$$

$$\text{โดยที่ } P(c_3|c_1 c_2) \text{ คำนวณได้จาก } = \frac{C(c_1-c_2-c_3)}{C(c_1-c_2)}$$

ดังนั้น จากคำว่า แอกเซส เมื่อคำนวณโดยใช้สมการนี้ จะได้ดังนี้ โดยที่ ANS คือ คำตอบ

$$P(\text{แ}) = \frac{\text{จำนวนสายอักขระ (ก)}}{\text{จำนวนสายอักขระทั้งหมด}}$$

จำนวนสายอักขระทั้งหมด

$$P(a|a) = \frac{\text{จำนวนสายอักขระ (a a)}}{\text{จำนวนสายอักขระ (a)}} = \text{ANS2}$$

$$P(a|a a) = \frac{\text{จำนวนสายอักขระ (a a a)}}{\text{จำนวนสายอักขระ (a a)}} = \text{ANS3}$$

$$P(a|a a) = \frac{\text{จำนวนสายอักขระ (a a a)}}{\text{จำนวนสายอักขระ (a a)}} = \text{ANS4}$$

$$P(a|a a) = \frac{\text{จำนวนสายอักขระ (a a a)}}{\text{จำนวนสายอักขระ (a a)}} = \text{ANS5}$$

$$P(a|a a) = \frac{\text{จำนวนสายอักขระ (a a a)}}{\text{จำนวนสายอักขระ (a a)}} = \text{ANS6}$$

หลังจากที่คำนวณได้แล้ว จึงนำแต่ละคำตอบมาคูณกัน ผลที่ได้คือ ค่าความน่าจะเป็นของคำว่า แอ กเซส ทั้งคำเมื่อประมาณโดยใช้ 3-แกรม โดยค่าความน่าจะเป็นยิ่งใกล้ ค่า 1 แสดงว่าเป็นค่าที่พบบ่อยในคลังข้อมูล

อย่างไรก็ตาม จากการคำนวณนี้ อาจเกิดกรณีที่คำนวณแล้วค่าที่ได้เป็นศูนย์ขึ้นได้ เนื่องจากสายอักขระบางสายไม่พบในคลังข้อมูล จำนวนสายอักขระที่พบจึงเป็นศูนย์ หากนำค่าศูนย์ไปคำนวณในสมการนี้ จะไม่สามารถคำนวณผลออกมาได้ หรือ ผลที่ออกมาจะเท่ากับศูนย์ เช่น จากการคำนวณคำว่า แอ กเซส

$$P(a|a a) P(a|a a) P(a|a a) P(a|a a)$$

$$0.052 * 0.032 * 0 * 0.045 * 0.000007 = 0$$

จึงต้องมีการปรับค่าไม่ให้เป็นศูนย์ ในที่นี้เลือกใช้วิธีแบบง่าย คือ เพิ่มค่า 0.00000001^2 ในการคำนวณความน่าจะเป็นของสายอักขระ $P(a|a a)$ ดังนี้

$$P(a|a a) = \frac{\text{จำนวนสายอักขระ(a a a)} + 0.00000001}{\text{จำนวนสายอักขระ(a a)} + 0.00000001}$$

โดยกำหนดเงื่อนไขไว้ว่า ในกรณีที่พบว่า จำนวนสายอักขระ (a a) หรือ จำนวนสายอักขระ (a a) มีค่าเท่ากับศูนย์หรือไม่พบค่าในคลังข้อมูลการฝึกเท่านั้น

² ในงานนี้ไม่ได้ใช้วิธีการปรับค่าที่ใช้กันทั่วไป เช่น Witten-Bell, Good Turing, deleted Interpolation, etc. เพราะคิดว่าวิธีที่ใช้นี้ก็เพียงพอกับงานแล้ว

ก่อนจะนำเสนอบทต่อไป ผู้วิจัยจะทบทวนวรรณกรรมหลักเกณฑ์การยืมคำจากภาษาต่างประเทศด้วยการทับศัพท์ก่อน โดยยึดหลักเกณฑ์ทางราชบัณฑิตยสถาน เพื่อแสดงเหตุผลว่าทำไมลักษณะเฉพาะของแต่ละภาษาสามารถใช้บอกที่มาของภาษาได้ และแสดงเหตุผลที่ผู้วิจัยนำสายอักขระเฉพาะมาวิเคราะห์เพื่อหาลักษณะเฉพาะของแต่ละภาษา

2.4 หลักเกณฑ์การทับศัพท์ภาษาต่างประเทศ

จากนิยามคำว่า “คำทับศัพท์” ในพจนานุกรมฉบับราชบัณฑิตยสถาน พ.ศ. 2525 ทำให้สรุปได้ว่าการทับศัพท์ คือ การถ่ายเสียงศัพท์ภาษาต่างประเทศด้วยตัวอักษรไทยทีละตัว หรือ การถ่ายคำศัพท์ภาษาต่างประเทศให้เป็นอักษรไทย ตัวอย่างการทับศัพท์ตามเกณฑ์ทางราชบัณฑิตยสถาน เช่น

คำทับศัพท์ภาษาอังกฤษ คำว่า	Calculus เขียนทับศัพท์ว่า แคลคูลัส
	Technique เขียนทับศัพท์ว่า เทคนิค
คำทับศัพท์ภาษาฝรั่งเศส คำว่า	Bastia เขียนทับศัพท์ว่า บาสตียา
	Technique เขียนทับศัพท์ว่า เตกนีก

เมื่อดูจากเกณฑ์การทับศัพท์ตามราชบัณฑิตยสถาน จะเห็นได้ว่าการทับศัพท์ภาษาต่างประเทศไม่ได้ใช้ตัวอักษรไทยทุกตัว เช่น ชม ฎ ฏ ฐ ศ ษ ฒ เป็นต้น เพราะฉะนั้นจึงมีสายอักขระเฉพาะอย่างน้อยที่สุดก็ในคำไทย และการคงรูปตัวอักษรแทนเสียงควบในคำภาษาต่างประเทศ อาจทำให้มีสายอักขระที่ไม่พบในไทย เช่น จากตัวอย่างในคำทับศัพท์ภาษาอังกฤษ เราอาจพบสายอักขระเฉพาะ “เทค” เพราะในคำไทยไม่มีการใช้อักษร ค เป็นตัวสะกดคำ ส่วนในคำทับศัพท์ภาษาฝรั่งเศสเราอาจพบสายอักขระเฉพาะ “เตก” เพราะในคำไทยรูปสระ อะ + ตัวสะกด จำเป็นต้องมีไม้ไต่คู้กำกับเพื่อให้ออกเสียงสระ เอ้- ได้ เช่นคำว่า เล็ก เป็นต้น จากเหตุผลนี้จึงอธิบายได้ว่าสายอักขระเฉพาะสามารถแสดงลักษณะเฉพาะของภาษาได้ ดังนั้นหากเรานำสายอักขระเฉพาะมาประยุกต์ใช้ในโปรแกรมระบุภาษาของคำ เราก็สามารถใช้สายอักขระเฉพาะบอกที่มาของภาษาได้

แม้ว่าในการทับศัพท์จริงจะแตกต่างจากเกณฑ์ทางราชบัณฑิตยสถานได้ เช่น คำว่า lock หากทับศัพท์ตามเกณฑ์ทางราชบัณฑิตยสถานจะใช้คำว่า “ล็อก” เพราะเสียงพยัญชนะท้าย /k/ ให้แทนด้วยตัวอักษร ก และทางราชบัณฑิตยสถานกำหนดให้ใช้ ไม้ไต่คู้เพื่อแสดงความ

แตกต่างจากคำไทย (ลอก) ด้วย แต่ในการใช้จริงอาจปรากฏการทับศัพท์แบบอื่นได้ เช่น ล็อค ล็อค เพราะผู้ใช้อาจออกเสียงคำทับศัพท์ตามความนิยม ความถนัด หรือออกเสียงให้เข้ากับระบบเสียงของภาษาต่างประเทศที่ใช้ แต่ว่าหากเราพิจารณาเปรียบเทียบหลักเกณฑ์การทับศัพท์ของราชบัณฑิตยสถาน ที่ใช้สำหรับภาษาต่างๆ ก็แสดงในเบื้องต้นได้ว่า มีความแตกต่างระหว่างคำทับศัพท์ภาษาต่างๆ อยู่ เช่น ตัวอักษรที่ใช้ในสำหรับแต่ละภาษาจะไม่เหมือนกัน โดยในคำไทยมีการใช้ตัวอักษร ค ฃ เพื่อออกเสียง ค ใช้ ส ษ ศ เพื่อออกเสียง ส แต่ในภาษาอังกฤษใช้ตัวอักษร c และ s เท่านั้น และเสียงที่ถ่ายออกมาเวลาทับศัพท์แต่ละภาษาก็แตกต่างกัน ได้แก่

คำทับศัพท์ภาษาอังกฤษ คำว่า Europe เขียนทับศัพท์ว่า ยูโรป

คำทับศัพท์ภาษาฝรั่งเศส คำว่า Europe เขียนทับศัพท์ว่า เออโรป

จะเห็นได้ว่าคำๆ เดียวกันแต่ออกเสียงต่างกัน ในการถ่ายเสียงเวลาทับศัพท์จึงออกมาต่างกัน ซึ่งสิ่งเหล่านี้ทำให้สะท้อนความแตกต่างออกมาในรูปของสายอักขระที่ปรากฏได้ เช่น จากตัวอย่างคำทับศัพท์ภาษาอังกฤษ เราอาจได้สายอักขระเฉพาะ “โรป” และคำทับศัพท์ภาษาฝรั่งเศส เราอาจได้สายอักขระเฉพาะ “รอป” ดังนั้น เราจึงสามารถนำสายอักขระเฉพาะเหล่านี้มาใช้เพื่อตัดสินภาษาของคำทับศัพท์ได้

บทที่ 3

คลังข้อมูลคำทับศัพท์

ในบทนี้ผู้วิจัยจะนำเสนอวิธีการสร้างคลังข้อมูลและเกณฑ์ในการสร้างคลังข้อมูลคำทับศัพท์ เพื่อให้เข้าใจภาพรวมของคลังข้อมูลที่ใช้ในการทดลอง และเหตุผลที่เลือกใช้คลังข้อมูลคำทับศัพท์ภาษาอังกฤษ ญี่ปุ่นและฝรั่งเศส

3.1 การสร้างคลังข้อมูล

คลังข้อมูลที่ใช้ในงานวิจัยนี้ประกอบด้วย คลังข้อมูลคำไทย คำทับศัพท์ภาษาอังกฤษ คำทับศัพท์ภาษาญี่ปุ่น และคำทับศัพท์ภาษาฝรั่งเศส สาเหตุที่เลือกศึกษาคำทับศัพท์ภาษาอังกฤษและภาษาญี่ปุ่นเพราะเป็นคำทับศัพท์ที่พบมากในภาษาไทยและเป็นภาษาที่แตกต่างกันมาก จึงเลือกมาเป็นตัวอย่างของภาษาที่ใช้ทดสอบว่าแบบจำลองเอ็นแกรมในที่นี้จะสามารถระบุที่มาภาษาที่แตกต่างกันนี้ได้หรือไม่ นอกจากนี้ที่ผู้วิจัยเลือกทดลองกับภาษาฝรั่งเศสด้วย เพราะว่าคำทับศัพท์ภาษาฝรั่งเศสนี้จะพบน้อยกว่าภาษาอังกฤษและภาษาญี่ปุ่น จึงถือเป็นการทดสอบระบบว่าจะยังคงมีความทนทานของระบบ (robustness) ว่าจะยังคงดีหรือไม่หากใช้คลังข้อมูลการฝึกจำนวนน้อย

คลังข้อมูลการฝึก (training corpus) ทั้งสำหรับทำแบบจำลองภาษา และชุดอักขระเฉพาะ นำมาจาก พจนานุกรมฉบับราชบัณฑิตยสถาน พ.ศ. 2525 (2538) หลักเกณฑ์การทับศัพท์ภาษาอังกฤษ (2538) หลักเกณฑ์การทับศัพท์ภาษาฝรั่งเศส เยอรมัน อิตาลี สเปน รัสเซีย ญี่ปุ่น อาหรับ และมลายู (2535) ศัพท์เทคนิควิศวกรรมไฟฟ้าสื่อสาร (2541) รวมศัพท์ญี่ปุ่น (2544) และหนังสืออื่นๆ เช่น หนังสือการท่องเที่ยวในต่างประเทศต่างๆ และหนังสือที่แสดงศัพท์เฉพาะทาง (technical term) ต่างๆ นอกจากนี้ยังรวบรวมจากบทความในอินเทอร์เน็ตด้วย ซึ่งจะมีทั้งคำไทย และคำทับศัพท์ภาษาอังกฤษ ญี่ปุ่นและฝรั่งเศส เหตุผลที่เลือกเก็บข้อมูลคำทับศัพท์จากแหล่งข้อมูลหลายแหล่ง เพราะต้องการความหลากหลายของคำทับศัพท์ โดยคำที่ปรากฏในแหล่งข้อมูล ผู้วิจัยเป็นผู้ตัดสินใจเองว่าเป็นคำไทยหรือคำทับศัพท์ภาษาใด โดยใช้ความรู้เรื่องภาษาและคำทับศัพท์ของผู้วิจัยเองและใช้การสังเกตจากคำบางคำที่ปรากฏทั้งรูปเดิมและรูปทับศัพท์ด้วย หากมีข้อสงสัยก็จะตรวจสอบจากพจนานุกรมภาษานั้นๆ ว่ามีคำนั้นใช้ในความหมายที่พบหรือไม่ จากคลังข้อมูลที่รวบรวมมาเหล่านี้ เราสามารถนำมาใช้ในการฝึก โดย

การนำแต่ละคำแยกเป็นคำไทย และคำทับศัพท์ของแต่ละภาษาก่อน ขนาดของคลังข้อมูลมี ทั้ง 4 ประเภท เป็นดังนี้

- คำไทย รวบรวมทั้งหมด 10,000 คำ คำไทยในที่นี่รวมไปถึงคำบาลีและสันสกฤตด้วย เพราะระบบประมวลผลภาษาโดยทั่วไปไม่จำเป็นต้องแยกความต่างระหว่างคำไทยแท้ กับคำไทยที่ยืมมาจากภาษาบาลีหรือสันสกฤต
- คำทับศัพท์ภาษาอังกฤษ รวบรวมทั้งหมด 10,000 คำ คำทับศัพท์ภาษาอังกฤษเป็นคำทับศัพท์ที่แพร่หลายมากที่สุดทั้งในหนังสือวิชาการ บทความ นิตยสาร หนังสือพิมพ์และศัพท์วิชาการส่วนใหญ่ใช้คำทับศัพท์ภาษาอังกฤษ เช่น ศัพท์ทางวิทยาศาสตร์ คอมพิวเตอร์ วิศวกรรมศาสตร์ เป็นต้น
- คำทับศัพท์ภาษาญี่ปุ่น รวบรวมทั้งหมด 10,000 คำ คำทับศัพท์ภาษาญี่ปุ่นส่วนใหญ่รวบรวมจากหนังสือเกี่ยวกับวัฒนธรรม ประเพณี การแต่งกาย อาหารและประวัติศาสตร์ในประเทศญี่ปุ่น
- คำทับศัพท์ภาษาฝรั่งเศส จะรวบรวมเพียงแค่ 1,000 คำ คำทับศัพท์ภาษาฝรั่งเศสจะรวบรวมจากหนังสือเฉพาะทางที่เกี่ยวกับประเทศฝรั่งเศส วัฒนธรรม สถาปัตยกรรมเป็นส่วนใหญ่

3.2 เกณฑ์ในการคัดเลือกข้อมูล

เนื่องจากในปัจจุบันมีการใช้คำทับศัพท์กันอย่างแพร่หลาย และการเขียนทับศัพท์เป็นภาษาไทยของแต่ละภาษาไม่มีการกำหนดกฎเกณฑ์ที่แน่นอน จึงทำให้เกิดคำทับศัพท์หลายรูปแบบ เช่น ตัวอย่างจากข้อความที่พบ

“นายเรจินัลด์ จี แม็คคาร์ตี ประกอบอาชีพการมเครดิตการ์ดระบบออนไลน์ ณ เมืองเลควิลล์ ขณะที่ลูกค้าขอเปิดแอลซีไปยังประเทศฝรั่งเศส”

“อย่าเวิร์กฮาร์ดนั้เลย ไปเอ็นจอยไลฟ์ชะบ้าง”

“ผู้รู้ภาษาอย่างดีกับผู้ไม่รู้จะออกเสียงต่างกัน เช่น เซนทรัล เซ็นทรัล เซินทรัล หรือ สเปเชียล สเปเชียล ต่างก็มาจากคำทับศัพท์คำเดียวกัน”

จะพบว่าภายในข้อความเหล่านี้มีการใช้คำทับศัพท์ภาษาอังกฤษมากมายและต่างกัน เช่น คำทับศัพท์ที่เขียนเรียงกันทั้งๆที่เป็นคนละคำ (เวิร์กฮาร์ด เครดิตการ์ด เลควิลล์) ชื่อเฉพาะ (เรจินัลด์ จี แม็คคาร์ตี) คำย่อ (แอลซี) และคำทับศัพท์ที่เขียนต่างกัน ทั้งๆที่ทับศัพท์จากคำๆ เดียวกัน (เซนทรัล เซ็นทรัล เซินทรัล) เป็นต้น หากไม่ได้กำหนดเกณฑ์ใดๆ ในการเก็บรวบรวมข้อมูลคำทับศัพท์เหล่านี้ เพียงแค่พบคำทับศัพท์ก็เก็บแล้ว อาจทำให้ข้อมูลที่รวบรวมไม่สม่ำเสมอ

และเกิดความผิดพลาดของระบบระบุภาษาได้ ในที่นี้ จึงได้กำหนดเกณฑ์ในการเก็บรวบรวมข้อมูลคำทับศัพท์ ดังนี้

1. กรณีที่เป็นคำทับศัพท์ที่มาจากรสรพจน์ (Acronym) และคำย่อ (Abbreviation) เห็นได้จากตัวอักษรพิมพ์ใหญ่ภายในวงเล็บต่อท้ายจากคำทับศัพท์ต่างๆ โดยคำทับศัพท์ทั้ง 2 ประเภทนี้ เมื่อเขียนทับศัพท์เป็นไทยแล้วจะอ่านต่างกัน คือ คำรสรพจน์จะอ่านเป็นคำ ส่วนคำย่อจะอ่านแบบคำย่อ ตัวอย่างเช่น

คำย่อ television มีคำย่อว่า TV อ่านทับศัพท์ ทีวี

Non-performing loan มีคำย่อว่า NPL อ่านทับศัพท์ว่า เอ็นพีแอล

The school of Germological Sciences มีคำย่อว่า SGS มีคำย่อว่า เอสจีเอส

และอีกมากมาย เช่น ซีบีเอส ซีเอดีเอฟ ซีเอฟเออาร์ ซีไอเอ ทีเอ็นที บีทียู บีทีเอส

บีทูเอส พีซีโอ พีอาร์ ยูพีเอส ยูวี ยูอาร์แอล ยูเอสทีอาร์ วีดีโอ วีเอชเอส อาร์ซี

รสรพจน์ signal plus noise plus distortion to noise plus distortion ratio (SINAD) อ่านทับศัพท์ ไชนันด์ ไมไซ เอสไอเอ็นเอดี

UNESCO อ่านทับศัพท์ว่า ยูเนสโก

ASEAN อ่านทับศัพท์ว่า อาเซียน

และอีกมากมาย เช่น เอเปก โอเปค เอแบค เฮฟต้า ออฟต้า

ในที่นี้ผู้วิจัยเลือกเก็บข้อมูลคำทับศัพท์ภาษาอังกฤษที่เป็นอ่านเป็นคำ หรือ คำรสรพจน์เท่านั้น เพราะคำย่อเป็นเพียงตัวคำอ่านจากอักษรเท่านั้น ไม่ได้ให้ลักษณะสำคัญของภาษามากเท่ากับสายอักขระที่เกิดร่วมกันเป็นคำ และคำย่อไม่เป็นปัญหาในระบบการระบุภาษาของคำทับศัพท์ เนื่องจากคำย่อมีจำนวนจำกัด เช่น เอ บี ซี ดี จนถึง แซด เท่านั้นจึงไม่นำมารวมอยู่ในปัญหาการระบุภาษาของคำของงานนี้

ส่วนในกรณีที่มีการใช้ผสมระหว่างคำทับศัพท์จากรสรพจน์และคำย่อ เช่น เอเพ็กซ์ (Apex) เอ็มลิงค์ (M-link) เอ็มเว็บ (M-web) ไอทีซิตี (IT City) อีเมลล์ (E-mail) เอ็นเกจ (N-gage) อีคอมเมิร์ซ (E-commerce) เอ็นเกจ (N-gage) ในที่นี้ผู้วิจัยเลือกเก็บทั้งคำโดยไม่สนใจว่าเป็นคำย่อหรือไม่ เพราะเขียนติดกันจนถือว่าเป็นคำๆ เดียวกัน ยกเว้นในกรณีที่เขียนทับศัพท์แล้วมีเครื่องหมายจุดคั่นหรือมีการเว้นช่องไฟ เช่น จี.เอช.ดี. โคลด์ (G.H.D. Cold) ผู้วิจัยจะเลือกเก็บแต่เฉพาะคำทับศัพท์คำว่า โคลด์ เท่านั้น

2. ในกรณีที่เป็นคำทับศัพท์ที่มาจากชื่อเฉพาะ ผู้วิจัยก็จะเลือกเก็บทั้งหมด โดยไม่สนใจว่าเป็นคำทับศัพท์ที่มาจากชื่อเฉพาะประเภทใด เช่น คำว่า พอลลา ซีเมนต์ โมโตโรล่า พานาโซนิค แต่ชื่อเฉพาะที่ปรากฏในภาษาไทย มีทั้งชื่อเฉพาะที่เขียนติดกัน และชื่อเฉพาะที่มีการเขียนทับศัพท์ผสมระหว่างคำไทยและคำทับศัพท์ ดังนั้นจึงมีเกณฑ์ ดังนี้

2.1 ชื่อเฉพาะที่มีการเขียนทับศัพท์ติดกันไป ไม่มีการแยกคำตามภาษาเดิมหรือมีเครื่องหมายใดๆ มากั้น กรณีนี้ให้เก็บทั้งคำโดยไม่แยก เพราะถือว่าเป็นคำๆ หนึ่ง³ เช่น ยูไนเต็ทฟลาวมิลล์ รัอกเวิช ทิบโก้ฟู้ดส์ เดทตอล เอปสัน การ์เนชั่น เป็นต้น

2.2 ชื่อเฉพาะที่มีการเขียนทับศัพท์ผสมระหว่างคำไทยและคำทับศัพท์ กรณีนี้ให้เก็บเฉพาะคำทับศัพท์อย่างเดียวแม้จะเขียนทับศัพท์ติดกันจนเป็นคำๆ เดียวกัน เช่น

“อมรินทร์พรีนติ้ง แอนด์ พับลิชชิ่ง” เก็บคำว่า พรีนติ้ง แอนด์ และพับลิชชิ่ง

“ห้องเย็นเอเซีย ชิฟู้ด” เก็บคำว่า เอเซีย และชิฟู้ด

“คอนโดวิภาวดีสวีท” เก็บคำว่า คอนโด และสวีท

จะเห็นว่าคำที่ขีดเส้นใต้ไม่ใช่ทับศัพท์จึงต้องตัดทิ้ง และแยกเก็บเฉพาะคำที่เป็นคำทับศัพท์ภาษาอังกฤษเป็นคำๆ แยกกัน เนื่องจากผู้วิจัยต้องการให้ระบบสามารถระบุภาษาของคำทับศัพท์ เฉพาะในส่วนของคำที่มีที่มาจากคำทับศัพท์ภาษาต่างประเทศจริงๆ เท่านั้น

3. เก็บรวบรวมรูปแบบคำทับศัพท์ทั้งหมดแม้ว่าจะเป็นคำทับศัพท์ภาษาต่างประเทศที่เขียนจากคำศัพท์คำๆ เดียวกันแต่เขียนในรูปต่างกัน เนื่องจากผู้ใช้คำทับศัพท์อาจเขียนคำทับศัพท์ผิดแปลกไปจากเกณฑ์ของราชบัณฑิตยสถาน อาจมีเหตุผลคือไม่ค่อยเป็นที่นิยม ไม่มีความรู้เพียงพอในการเขียนคำทับศัพท์ให้ถูกต้อง หรือผิดพลาดในการตรวจแก้ไขคำทับศัพท์ ดังนั้นคำทับศัพท์ที่ปรากฏตามหนังสือหรือนิตยสารจึงแตกต่างกัน แม้ว่าจะเขียนทับศัพท์จากคำศัพท์ภาษาต่างประเทศตัวเดียวกัน ตัวอย่างเห็นได้จากตารางที่ 3.1 ดังนี้

³ เหตุผลที่ใช้เกณฑ์เช่นนี้ เพราะระบบจะไม่รู้ว่าคำทับศัพท์ที่เห็นเขียนติดกันอยู่นั้นประกอบด้วยคำต่างประเทศกี่คำ โปรแกรมระบุภาษาของคำที่พัฒนานี้ทำได้เพียงตัดสินว่า สายอักขระนั้นมาจากภาษาใด ซึ่งเมื่อรู้คำตอบแล้ว จึงส่งให้โปรแกรมถอดอักษรกลับ (backward transliteration) หาเอาเองว่าเป็นคำต่างประเทศคำใด ซึ่งอาจประกอบด้วยคำมากกว่าหนึ่งคำก็ได้ ดังนั้น ข้อมูลการฝึกจึงเก็บทั้งสายอักขระโดยไม่สนใจว่าประกอบจากคำต่างประเทศมากกว่าหนึ่งคำหรือไม่

ตารางที่ 3.1 การเขียนทับศัพท์ที่ต่างจากเกณฑ์ทางราชบัณฑิตยสถาน

ศัพท์	ทับศัพท์ตามหลักราชบัณฑิตยสถาน	ทับศัพท์แบบอื่น
Lock	ล็อก	ล๊อค ล็อค ล็อค ล็อค
Tent	เต็นท์	เต้นท์
Knot	นอต	น็อต น็อต
Pyramid	พีระมิด	ปिरามิด ปิรามิด พีรามิด

ในที่นี้ผู้วิจัยเลือกเก็บทั้งหมดโดยไม่สนใจว่าทับศัพท์ตามเกณฑ์ทางราชบัณฑิตยสถานหรือไม่ เนื่องจากผู้วิจัยคาดหวังว่าระบบมีประสิทธิภาพเพียงพอที่จะสามารถระบุภาษาของคำที่รับเข้ามาทดสอบได้ทั้งหมดโดยคำที่รับเข้ามาอาจจะตรงหรือไม่ตรงตามเกณฑ์ของราชบัณฑิตยสถานก็ได้ ดังนั้นข้อมูลการฝึกจึงจำเป็นต้องมีทั้งคำที่ตรงและไม่ตรงตามเกณฑ์ของราชบัณฑิตยสถาน เพื่อให้ตรงกับข้อมูลจริงที่พบใช้กันทั่วไป

4. ในกรณีคำทับศัพท์ที่สร้างจากคำภาษาต่างประเทศที่มีเครื่องหมายยัติภังค์ (hyphen) คั่น จะเก็บรวบรวมทั้งคำ เนื่องจากในคำทับศัพท์ภาษาไทยจะเขียนทับศัพท์ติดต่อกันไปโดยไม่มีเครื่องหมายยัติภังค์แยกคำตามภาษาเดิม เช่น คำทับศัพท์ภาษาอังกฤษ

ครอสสติตช์ (cross-stitch) เก็บคำว่า ครอสสติตช์

แคลเซียมคาร์บอเนต (calcium-carbonate) เก็บคำว่า แคลเซียมคาร์บอเนต

ไวร์คัท (wire-cut) เก็บคำว่า ไวร์คัท

แอนตี้สลิป (Anti-Slip) เก็บคำว่า แอนตี้สลิป

แต่ในกรณีคำทับศัพท์ที่เขียนจากคำศัพท์ที่มีเครื่องหมายยัติภังค์คั่น โดยเขียนทับศัพท์คั่นด้วยเครื่องหมายยัติภังค์เช่นเดียวกับภาษาเดิม ผู้วิจัยจะแยกเก็บเป็นคำๆ ตามภาษาเดิมด้วย เช่น คำว่า McGraw-Hill เขียนทับศัพท์ว่า แมกกรอว์-ฮิลล์ ในกรณีนี้ผู้วิจัยเลือกยกเว้นคำทับศัพท์คำนี้ แต่แยกเก็บเป็นคำๆ แทนคือเก็บคำว่า แมกกรอว์ และฮิลล์

5. เก็บรวบรวมคำทับศัพท์ที่สร้างจากคำภาษาต่างประเทศที่เขียนแยกกัน แม้ว่าภาษาเดิมจะมีการเขียนแยกกันแต่เมื่อทับศัพท์แล้วจะเขียนติดกันไป ไม่ต้องแยกตามภาษาเดิม เช่น คำทับศัพท์ภาษาอังกฤษ

ฟังก์ชันไฮเพอร์โบลิก (Hyperbolic function) เก็บคำว่า ฟังก์ชันไฮเพอร์โบลิก เป็น 1 คำ

อาร์กติกเซอร์เคิล (Arctic Circle) เก็บคำว่า อาร์กติกเซอร์เคิล เป็น 1 คำ และมีอีกมากมาย เช่น กรีนพีซ (green piece) กรีนไวลด์ (green wild) โกลด์โคสต์ (gold coast) คอมมอนเซนส์ (common sense) เป็นต้น

ลักษณะของคำทับศัพท์แบบนี้จะเกิดขึ้นบ่อยที่สุด จนบางครั้งอาจใช้จนผู้ใช้ไม่รู้ว่คำศัพท์คำนั้นประกอบไปด้วยคำทับศัพท์ภาษาอังกฤษกี่คำ บางครั้งอาจคิดว่าเป็นคำทับศัพท์ที่มาจากภาษาอังกฤษคำเดียวกัน ดังนั้นผู้วิจัยจึงเลือกเก็บทั้งคำ โดยไม่ต้องแยกเป็นคำๆ ตามคำศัพท์ภาษาเดิม

อย่างไรก็ตาม มีข้อยกเว้นบางคำที่ภาษาเดิมมีการเขียนแยกกัน เมื่อมีการเขียนทับศัพท์โดยเขียนแยกกันตามภาษาเดิม ในกรณีนี้ให้เก็บแยกเป็นคำๆ เช่น คำทับศัพท์ภาษาอังกฤษ

“แฟนซีวูด อินดัสตรีส์” เก็บคำว่า แฟนซีวูด และ อินดัสตรีส์

“ไดอาน่า ดีพาร์ทเมนท์สโตร” เก็บคำว่า ไดอาน่า และดีพาร์ทเมนท์สโตร

“ลา구나 รีสอร์ท” เก็บคำว่า ลา구나 และรีสอร์ท

นอกจากนี้ในกรณีคำทับศัพท์ที่เขียนติดกันแต่มีคำคุณศัพท์ประกอบกับคำทับศัพท์คำนั้น เช่น คำว่า “เป็นของ” “แบบ” “อย่าง” “ชนิด” “ระบบ” ฯลฯ ในกรณีนี้ผู้วิจัยจะแยกเก็บเป็นที่ละคำ เพราะคำเหล่านี้ไม่ใช่มีที่มาจากคำภาษาต่างประเทศ จะต้องเก็บคำที่เป็นคำภาษาต่างประเทศจริงๆ เท่านั้นเพื่อไม่ให้เกิดปัญหาในการใช้งานระบุภาษาของคำทับศัพท์จริงๆ ตัวอย่างคำทับศัพท์ภาษาอังกฤษที่เป็นปัญหาเช่น

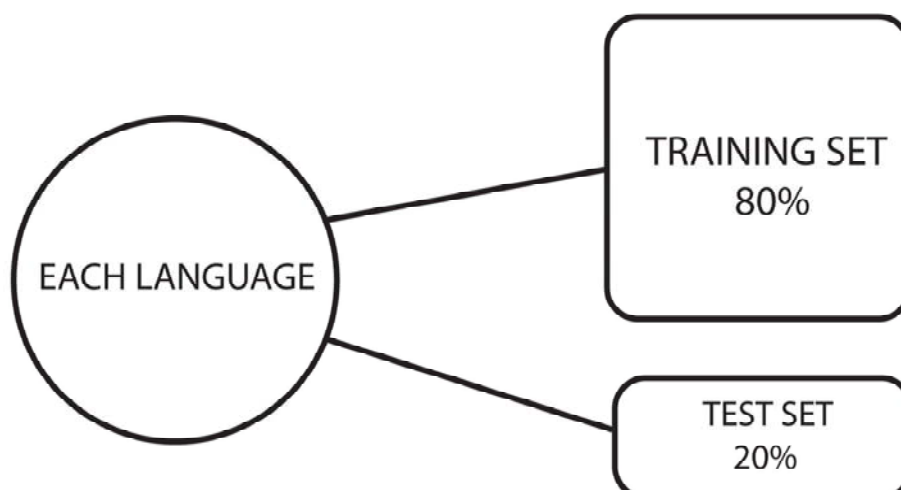
Eulerian function จะเขียนทับศัพท์ว่า ฟังก์ชันแบบออยเลอร์

Thermosetting plastic จะเขียนทับศัพท์ว่า พลาสติกชนิดเทอร์โมเซตติง
จากตัวอย่างนี้ ผู้วิจัยจะเก็บที่ละคำคือ ฟังก์ชัน ออยเลอร์ พลาสติก และเทอร์โมเซตติง

3.3 การสร้างข้อมูลการฝึกและข้อมูลทดสอบสำหรับใช้ในการทดลอง

เนื่องจากภายในระบบระบุภาษาจากคำจำเป็นต้องใช้ข้อมูลการฝึกสำหรับสร้างแบบจำลองภาษาและใช้ข้อมูลสำหรับการทดลอง ดังนั้นจากการรวบรวมคลังข้อมูลคำไทยและคำทับศัพท์ภาษาต่างประเทศ ซึ่งมีคำไทย 10,000 คำ คำทับศัพท์ภาษาอังกฤษ 10,000 คำ คำทับศัพท์ภาษาญี่ปุ่น 10,000 คำ และคำทับศัพท์ภาษาฝรั่งเศส 1,000 คำ ผู้วิจัยได้แบ่งส่วนข้อมูลเหล่านี้

โดยใช้หลักการ 5-fold cross validation approach ซึ่งแยกเป็นข้อมูลการฝึก (training set) ภาษาละ 80% จากคลังข้อมูลทั้งหมดในแต่ละภาษา ในการทดลองนี้ ใช้คำไทย 8,000 คำ คำทับศัพท์ภาษาอังกฤษ 8,000 คำ และคำทับศัพท์ภาษาญี่ปุ่น 8,000 คำ และคำทับศัพท์ภาษาฝรั่งเศส 800 คำ และข้อมูลทดสอบ (test set) ภาษาละ 20% จากคลังข้อมูลทั้งหมดในแต่ละภาษา ในการทดลองนี้ ใช้คำไทย 2,000 คำ คำทับศัพท์ภาษาอังกฤษ 2,000 คำ และคำทับศัพท์ภาษาญี่ปุ่น 2,000 คำ และคำทับศัพท์ภาษาฝรั่งเศส 200 คำ โดยได้แบ่งด้วยวิธีการสุ่มเลือกจากคลังข้อมูลโดยใช้โปรแกรมภาษาเพิร์ลเป็นเครื่องมือในการแบ่งส่วน เห็นได้จากภาพที่ 3.1



ภาพที่ 3.1 การแบ่งส่วนคลังข้อมูลในแต่ละภาษา

และเพื่อให้ได้ผลที่ไม่บังเอิญขึ้นกับข้อมูลที่ใช้ฝึกและทดสอบ จะทำการทดลอง 5 ครั้ง โดยแต่ละครั้งจะใช้ข้อมูลทดสอบ 20% ที่ต่างกัน ดังนั้น ผลที่ได้ในแต่ละภาษาจะมีชุดข้อมูล 5 ชุด แล้วจึงนำผลแต่ละชุดมาสรุปรวมกันเพื่อแสดงผลการทดสอบการระบุค่าของระบบ ในคำไทย คำทับศัพท์ภาษาอังกฤษ และญี่ปุ่น แสดงผลค่าทดสอบภาษาละ 10,000 คำ ส่วนคำทับศัพท์ภาษาฝรั่งเศส 1,000 คำ

บทที่ 4

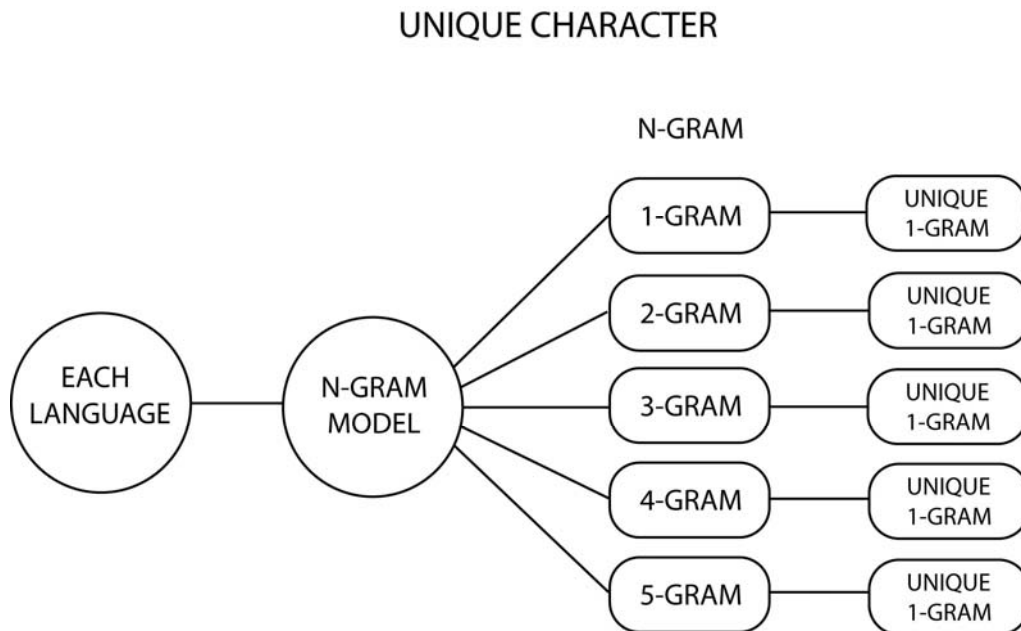
การระบุภาษาของคำด้วยสายอักขระเฉพาะ

ในบทนี้ผู้วิจัยจะกล่าวถึงระบบที่ผู้วิจัยพัฒนา คือ การระบุภาษาของคำด้วยสายอักขระเฉพาะก่อน โดยแสดงถึงภาพรวมของระบบการระบุภาษาโดยใช้สายอักขระเฉพาะ ผลการทดลอง สายอักขระเฉพาะที่พบในคลังข้อมูล สายอักขระเฉพาะที่ใช้จริงในการทดสอบระบบ รวมทั้งนำสายอักขระเฉพาะที่ได้จากคลังข้อมูลไปวิเคราะห์เพื่อหาลักษณะเฉพาะของแต่ละภาษา

4.1 ภาพโดยรวมของวิธีการใช้สายอักขระเฉพาะ (unique characters)

วิธีการใช้สายอักขระเฉพาะในงานวิจัยนี้ เป็นการนำคลังข้อมูลแต่ละภาษาทั้งหมดมาหาสายอักขระเฉพาะของแต่ละภาษาโดยใช้วิธีเอ็นแกรมของสายอักขระตั้งแต่ 1-5 แกรม เช่น ในการสร้างสายอักขระเฉพาะด้วย 3-แกรม เราอาจพบสายอักขระเฉพาะ “ทร” ในคำไทย และอาจพบสายอักขระเฉพาะ “ตร์” ในคำทับศัพท์ภาษาอังกฤษ เป็นต้น ในการสร้างสายอักขระเฉพาะของแต่ละภาษา จะใช้วิธีนำคลังข้อมูลที่รวบรวมเป็นคำมาแยกเป็นสายอักขระตั้งแต่ 1-5 แกรม แล้วนำมาเปรียบเทียบกันในแต่ละแกรม โดยไม่ขึ้นอยู่กับค่าความถี่ของสายอักขระ แต่ใช้วิธีสังเกตจากการปรากฏของสายอักขระ โดยหากสายอักขระใดที่ปรากฏแล้วไม่ซ้ำกับสายอักขระในภาษาอื่น ๆ ที่เหลือ จะถือว่าเป็นสายอักขระเฉพาะของภาษานั้น

ผลลัพธ์ที่ได้จากการสร้างสายอักขระเฉพาะจะได้ สายอักขระเฉพาะทั้งหมด 4 ภาษาคือ คำไทย คำทับศัพท์ภาษาอังกฤษ ภาษาญี่ปุ่น และภาษาฝรั่งเศส โดยในแต่ละภาษาจะประกอบไปด้วย สายอักขระเฉพาะ 1-5 ตัว เห็นได้จากภาพที่ 4.1

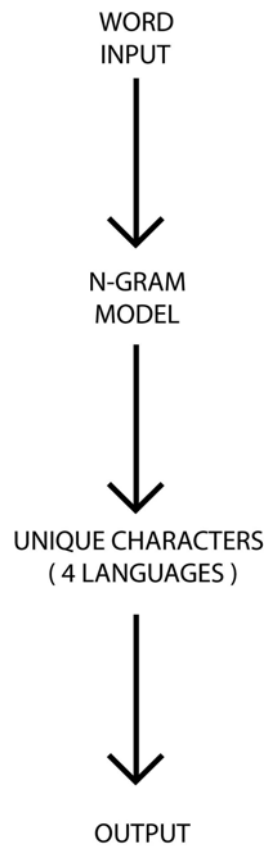


ภาพที่ 4.1 ลักษณะการสร้างสายอักขระเฉพาะ 1-5 แกรม

4.2 โปรแกรมระบุภาษาของคำด้วยสายอักขระเฉพาะ

โปรแกรมระบุภาษาของคำด้วยสายอักขระเฉพาะ แบ่งออกเป็น 2 ขั้นตอน ได้แก่

1. ขั้นตอนรับค่าเข้า
2. ขั้นตอนเทียบกับสายอักขระเฉพาะ เห็นได้จากภาพที่ 4.2



ภาพที่ 4.2 ขั้นตอนการทำงานของโปรแกรมระบุภาษาของคำด้วยสายอักขระเฉพาะ

1. ขั้นตอนรับคำเข้า เป็นขั้นตอนที่รับคำเข้ามาแล้วแยกคำๆ นั้น ออกเป็นสายอักขระด้วยเอ็นแกรมทั้ง 1-5 แกรม เช่น คำที่รับเข้ามาคำว่า “อักขระ” เมื่อผ่านระบบเอ็นแกรม จะแยกเป็นแต่ละเอ็นแกรม เห็นได้จากตารางที่ 4.1 ดังนี้

ตารางที่ 4.1 การแยกสายอักขระตามเอ็นแกรม 1-5 แกรม

1-แกรม	2-แกรม	3-แกรม	4-แกรม	5-แกรม
อ	อ้	อ้ก	อ้กข	อ้กขร
-	-ก	-กข	-กขร	-กขระ
ก	กข	กขร	กขระ	
ข	ขร	ขระ		
ร	ระ			
ะ				

2. ขั้นตอนเทียบกับสายอักขระเฉพาะ หลังจากที่แยกสายอักขระเป็นเอ็นแกรมต่างๆ แล้ว ขั้นตอนต่อไปเป็นการนำสายอักขระที่แยกด้วยเอ็นแกรม มาเทียบกับสายอักขระเฉพาะทั้ง 4 ภาษา โดยแยกเทียบกันทั้ง 1-5-แกรม ตามแต่ละเอ็นแกรมว่าตรงกันหรือไม่ หากสายอักขระที่นำมาทดสอบตรงกับสายอักขระเฉพาะของภาษาใดภาษาหนึ่งเท่านั้น ก็ถือว่าคำนั้นเป็นภาษาตามสายอักขระเฉพาะนั้นทันที แต่หากไม่มีสายอักขระที่ตรงกับสายอักขระเฉพาะของภาษาใดๆ เลย หรือ หากพบว่าคำนั้นมีสายอักขระตรงกับสายอักขระเฉพาะของภาษาต่างๆ มากกว่า 1 ภาษา ทำให้ยังตัดสินใจไม่ได้ว่าเป็นคำจากภาษาใด ตัวอย่างเช่น ในคำว่า “อักขระ” เมื่อทดลองด้วย 4-แกรม อาจะพบ สายอักขระ “กขระ” ตรงกับสายอักขระเฉพาะภาษาญี่ปุ่น และพบสายอักขระ “อักข” ตรงกับสายอักขระเฉพาะคำไทย ทำให้ตัดสินใจไม่ได้ว่าเป็นคำไทยหรือมาจากภาษาญี่ปุ่น ระบบจึงไม่ระบุว่าเป็นภาษาอะไร และข้ามไประบุคำถัดไป ดังนั้น ผลการทดลองจะได้เฉพาะคำที่สามารถระบุภาษาได้ด้วยสายอักขระเฉพาะเท่านั้น ซึ่งในที่นี้ผู้วิจัยเป็นผู้ตรวจสอบผลการระบุภาษาของคำเองว่า ถูกต้องหรือไม่

ในขั้นตอนการเปรียบเทียบกับสายอักขระเฉพาะนี้ ระบบที่ผู้วิจัยพัฒนาจะทดสอบเพียงครั้งเดียว คือ หลังจากที่รับคำที่ทดสอบแล้วแยกเป็นสายอักขระด้วยเอ็นแกรม 1-5-แกรม แล้วระบบจะนำสายอักขระที่ทดสอบ 1-แกรม มาเทียบกับสายอักขระเฉพาะ 1-แกรมในแต่ละภาษาก่อน ถ้ายังไม่พบว่าตรงกัน ระบบก็จะนำสายอักขระที่ทดสอบ 2-แกรม มาเทียบกับสายอักขระเฉพาะ 2-แกรมในแต่ละภาษา หากเปรียบเทียบกันแล้ว ยังไม่ตรงกันอีกก็จะทำงานแบบนี้ต่อไปใน 3 4 และ 5-แกรม ตามลำดับ และหากว่าระบบนำสายอักขระที่ทดสอบเปรียบเทียบกับสายอักขระเฉพาะจนถึง 5-แกรมแล้วยังไม่พบว่าตรงกัน ระบบก็จะข้ามไปทดสอบคำต่อไปเลย

4.3 ผลการทดลอง

จากการทดสอบโปรแกรมระบุภาษาของคำไทยและคำทับศัพท์แสดงให้เห็นผลของระบบระบุภาษาของคำด้วยสายอักขระเฉพาะ 1-5 ตัว ผลที่ได้มี 3 ผลการทดลองคือ ผลการตัดสินใจภาษาได้ถูกต้อง ผลการตัดสินใจภาษาผิดพลาด และผลการตัดสินใจภาษาไม่ได้ โดยเห็นได้จาก ตารางที่ 4.2

ผลการตัดสินใจภาษาได้ถูกต้อง หมายถึงระบบสามารถระบุภาษาของคำได้ถูกต้อง

ผลการตัดสินใจภาษาผิดพลาด หมายถึงระบบระบุภาษาของคำผิดพลาด

ผลการตัดสินใจภาษาไม่ได้ หมายถึงระบบไม่สามารถตัดสินใจภาษาของคำนั้นได้

ตารางที่ 4.2 ผลการตัดสีนภาษาได้ถูกต้อง ผลการตัดสีนภาษาผิดพลาด และผลการตัดสีนภาษาไม่ได้ ด้วยสายอักขระเฉพาะ 1-5 ตัว โดยผลในตาราง คือ จำนวนคำ ส่วนที่อยู่ในเครื่องหมายวงเล็บคือจำนวนเปอร์เซ็นต์

จำนวนการตัดสีนภาษาด้วยสายอักขระเฉพาะ 1-5 ตัวอักษร คำ (เปอร์เซ็นต์)				
ภาษา		ตัดสีนภาษาถูกต้อง	ตัดสีนภาษาผิดพลาด	ตัดสีนภาษาไม่ได้
คำไทย	1	1414 (14.14%)	0	
	2	756 (7.56%)	2 (0.02%)	
	3	2104 (21.04%)	274 (2.74%)	
	4	651 (6.51%)	164 (1.64%)	
	5	133 (1.33%)	23 (0.23%)	
รวม		5058 (50.58%)	463 (4.63%)	4479 (44.79%)
คำทับศัพท์ ภาษาอังกฤษ	1	0	2 (0.02%)	
	2	122 (1.22%)	14 (0.14%)	
	3	2018 (20.18%)	264 (2.64%)	
	4	1769 (17.69%)	277 (2.77%)	
	5	962 (9.62%)	81 (0.81%)	
รวม		4871 (48.71%)	638 (6.38%)	4491 (44.91%)
คำทับศัพท์ ภาษาญี่ปุ่น	1	0	0	
	2	160 (1.6%)	7 (0.07%)	
	3	3392 (33.92%)	186 (1.86%)	
	4	1683 (16.83%)	168 (1.68%)	
	5	174 (1.74%)	38 (0.38%)	
รวม		5409 (54.09%)	399 (3.99%)	4192 (41.92%)
คำทับศัพท์ภาษา ฝรั่งเศส	1	0	0	
	2	0	1 (0.10%)	
	3	57 (5.70%)	122 (12.20%)	
	4	101 (10.10%)	99 (9.90%)	
	5	46 (4.60%)	34 (3.40%)	
รวม		204 (20.40%)	256 (25.60%)	540 (54%)

ตารางที่ 4.2 แสดงผลดังนี้ คือ ผลการตัดสีนภาษาได้ถูกต้องด้วยสายอักขระเฉพาะ 1-5 ตัว ผลการตัดสีนภาษาผิดพลาด และผลการตัดสีนภาษาไม่ได้

จากการดูผลการตัดสินภาษาได้ถูกต้องด้วยสายอักขระเฉพาะในแต่ละภาษา เห็นว่า เมื่อระบุภาษาของคำไทย คำทับศัพท์ภาษาอังกฤษ และญี่ปุ่น ด้วยสายอักขระเฉพาะ 3-ตัว ให้ผลการระบุภาษาที่ถูกต้องสูงที่สุด คือ คำไทย 21.04% คำทับศัพท์ภาษาอังกฤษ 20.18% และคำทับศัพท์ภาษาญี่ปุ่น 33.92% แต่ในการระบุภาษาของคำทับศัพท์ภาษาฝรั่งเศส ด้วยสายอักขระเฉพาะ 4-ตัว ให้ผลดีที่สุด คือ 10.1% ส่วนการระบุภาษาด้วยสายอักขระเฉพาะ 1-ตัว สามารถใช้ระบุภาษาได้ถูกต้องได้เฉพาะในคำไทยเท่านั้น

นอกจากนี้ จากตารางที่ 4.2 ยังแสดงให้เห็นผลโดยรวมของระบบการใช้สายอักขระได้ โดยจะต้องดูในช่อง "รวม" เนื่องจากในการทำงานของระบบการใช้สายอักขระเฉพาะจะเป็นการทำงานทีละตัว โดยเริ่มจากอักขระเฉพาะ 1 ตัว ก่อน พอไม่ผ่านจึงมาที่ 2 3 4 และ 5-ตัว ตามลำดับ ดังนั้นหากคิดผลการทดลองของทั้งระบบการใช้สายอักขระเฉพาะ จึงต้องคิดรวมผลทั้งหมด 1-5 ตัว เช่น จากตารางที่ 4.2 เมื่อต้องการดูผลการระบุของคำไทยจะแสดงผลทีละตัว และนำมารวมกันดังนี้ $14.14 + 7.56 + 21.04 + 6.51 + 1.33 = 50.58\%$ ผลลัพธ์ที่ได้ก็คือ ค่าความถูกต้องจากการตัดสินภาษาของคำไทยด้วยการใช้ระบบสายอักขระเฉพาะตั้งแต่ 1-5 ตัว โดยผลในช่องรวม แสดงให้เห็นว่า วิธีการใช้สายอักขระเฉพาะมีผลต่อโปรแกรมระบุคำประมาณ 50% โดยสายอักขระเฉพาะมีผลต่อการระบุคำทับศัพท์ภาษาญี่ปุ่นมากที่สุด คือ 54.09% รองลงมาคือ คำไทย 50.58% และคำทับศัพท์ภาษาอังกฤษ 48.71% โดยอาจเป็นเพราะภาษาญี่ปุ่นเป็นภาษาที่พบสายอักขระเฉพาะมากที่สุด หรือ เป็นภาษาใช้สายอักขระซ้ำๆ กันมาก ไม่มีการผสมสระมากเท่าภาษาอื่นๆ ซึ่งจะอภิปรายเหตุผลที่แต่ละภาษาได้ผลการทดสอบแตกต่างกันต่อไป ในส่วนวิเคราะห์ผลการทดสอบในหัวข้อ 4.4

จากการดูผลการตัดสินภาษาผิดพลาด พบว่าคำทับศัพท์ภาษาอังกฤษผิดพลาดมากที่สุด ถึง 6.38% ซึ่งในส่วนของผลการตัดสินภาษาของคำที่ผิดพลาดนั้นจะกล่าวอย่างละเอียดในหัวข้อ 4.5

และเมื่อดูผลการตัดสินภาษาของคำไม่ได้ พบว่า หลังจากใช้สายอักขระเฉพาะแล้ว ยังเหลือคำที่ระบบยังตัดสินภาษาไม่ได้อีกประมาณ 41-45% คำไทยเหลือ 44.79% คำทับศัพท์ภาษาอังกฤษเหลือ 44.91% คำทับศัพท์ภาษาญี่ปุ่นเหลือ 41.92% และคำทับศัพท์ภาษาฝรั่งเศสเหลือ 54% คำที่ตัดสินภาษาไม่ได้เหล่านี้ เป็นคำที่ไม่มีสายอักขระใดๆ ตรงกับสายอักขระเฉพาะ 1-5 ตัวเลย ดังนั้น คำที่เหลือเหล่านี้จึงต้องนำไปตัดสินภาษาด้วยวิธีอื่นต่อไป โดยในที่นี้ นำไปตัดสินภาษาต่อไปด้วยแบบจำลองภาษา 2-5 แกรม (ซึ่งกล่าวถึงในบทที่ 6 การระบุภาษาของคำด้วยสายอักขระเฉพาะ 1-5 ตัวร่วมกับแบบจำลองภาษา 2-5 แกรม)

4.4 วิเคราะห์ผลการทดสอบ

จากการทดสอบด้วยสายอักขระเฉพาะแต่ละภาษา ผู้วิจัยพบว่าทุกภาษาที่ทดสอบ เมื่อระบุด้วยสายอักขระเฉพาะ 3-ตัว จะให้ผลดีที่สุด แต่เราไม่สามารถตัดสินได้ว่าการใช้สายอักขระเฉพาะตัวใดมีประสิทธิภาพดีที่สุด เพราะระบบการทำงานของสายอักขระเฉพาะจะเริ่มจาก 1-ตัว ก่อนเสมอและทำงานต่อไป 2 3 4 และ 5 ตามลำดับ

4.4.1 ผลการตัดสินภาษาจากสายอักขระเฉพาะคำไทย

ผลจากการระบุคำไทยด้วยสายอักขระเฉพาะ เมื่อระบุคำไทยด้วยสายอักขระเฉพาะ 3-ตัวให้ผลถูกต้องดีที่สุดคือ 21.04% และผลจะต่ำลงเมื่อใช้สายอักขระ 4-ตัว (6.51%) และ 5-ตัว (1.33%) ตามลำดับ และเมื่อเปรียบเทียบผลระหว่างภาษาพบว่า สายอักขระเฉพาะ 1 ตัว มีผลเฉพาะระบุคำไทยเท่านั้น เพราะ คำไทยมีการใช้พยัญชนะที่ภาษาอื่นไม่ใช้ คือพยัญชนะที่มาจากภาษาบาลี สันสกฤต ได้แก่ ข ฃ ฉ ฌ ฎ ฏ ฌ ฎ ฏ ฐ ฑ ฒ ณ ด ต ซึ่งในภาษาอื่นใช้พยัญชนะตัวอื่นแทนเสียงเหล่านี้ ทำให้พยัญชนะเหล่านี้ไม่ปรากฏในภาษาอื่น และเมื่อระบุด้วยสายอักขระเฉพาะ 2-ตัว ให้ผลถูกต้องมากกว่าภาษาอื่น ๆ คือ 7.56%

4.4.2 ผลการตัดสินภาษาจากสายอักขระเฉพาะคำทับศัพท์ภาษาอังกฤษ

ผลจากการระบุคำไทยด้วยสายอักขระเฉพาะ เมื่อระบุคำทับศัพท์ภาษาอังกฤษด้วยสายอักขระเฉพาะ 3 ตัว ให้ผลถูกต้องดีที่สุดคือ 20.18% และผลจะต่ำลงเมื่อระบุด้วยสายอักขระเฉพาะ 4-ตัว และ 5-ตัว ตามลำดับ และเมื่อเปรียบเทียบผลระหว่างภาษาพบว่า เมื่อระบุภาษาด้วยสายอักขระเฉพาะ 4 และ 5-ตัว จะให้ผลถูกต้องมากกว่าภาษาอื่น คือ 17.69% และ 9.62% และเมื่อระบุด้วยสายอักขระเฉพาะ 2 และ 3-ตัว ให้ผลต่ำกว่าภาษาอื่นมาก คือ 1.22% และ 20.18% อาจเป็นเพราะว่า ภาษาอังกฤษเป็นภาษาที่เขียนขึ้นด้วยอักขระยาวๆ เห็นได้จากมีการใช้สายอักขระเฉพาะยาวๆ ในระบบมาก เช่นสายอักขระเฉพาะ เตอร์ เบอร์ เลอร์ เป็นต้น

4.4.3 ผลการตัดสินภาษาจากสายอักขระเฉพาะคำทับศัพท์ภาษาญี่ปุ่น

ผลจากการระบุคำไทยด้วยสายอักขระเฉพาะ เมื่อระบุคำทับศัพท์ภาษาญี่ปุ่นด้วยสายอักขระเฉพาะ 3-ตัว ให้ผลถูกต้องดีที่สุดคือ 33.92% และผลจะต่ำลงเมื่อระบุด้วยสายอักขระเฉพาะ 4-ตัว และ 5-ตัว ตามลำดับ และเมื่อเปรียบเทียบผลระหว่างภาษาพบว่า เมื่อระบุภาษาด้วยสายอักขระเฉพาะ 3-ตัว ให้ผลถูกต้องมากกว่าภาษาอื่น อาจเป็นเพราะภาษาญี่ปุ่นเป็นภาษาที่เขียนด้วยคำเดี่ยวๆรวมกัน ได้แก่ อะ อิ อุ เอะ โอะ เช่น คำว่า อิคิยะ อียุ คิโยะ เป็นต้น

4.4.4 ผลการตัดสินภาษาจากสายอักขระเฉพาะคำทับศัพท์ภาษาฝรั่งเศส

ผลจากการระบุคำไทยด้วยสายอักขระเฉพาะ พบว่าสายอักขระเฉพาะ 2-ตัว ไม่มีผลต่อการระบุคำทับศัพท์ภาษาฝรั่งเศส และเมื่อระบุคำทับศัพท์ภาษาฝรั่งเศสด้วยสายอักขระเฉพาะ 4-ตัว ให้ผลถูกต้องดีที่สุดคือ คือ 10.1% โดยผลการทดลองระบุคำทับศัพท์ภาษาฝรั่งเศสอาจไม่ให้ผลที่แน่นอนเนื่องจากใช้คลังข้อมูลน้อย แต่อาจกล่าวได้ว่า หากใช้ขนาดของคลังข้อมูลน้อย จำนวนของสายอักขระเฉพาะที่พบและผลการใช้สายอักขระเฉพาะในการระบุภาษาของคำก็จะน้อยด้วย

4.5 ปัญหาการตัดสินคำผิดพลาดที่เกิดขึ้นจากการใช้สายอักขระเฉพาะ 1-5 ตัว

จากผลการทดสอบระบบ พบว่ามีการระบุภาษาผิดพลาดจากการใช้สายอักขระเฉพาะ 1-5 ตัว โดยที่คำภาษาหนึ่งถูกตัดสินผิดพลาดให้เป็นอีกภาษาหนึ่ง เช่น คำไทย “ตะนอย” ถูกตัดสินผิดพลาดเป็นคำทับศัพท์ภาษาอังกฤษ ข้อผิดพลาดในส่วนนี้แม้จะคิดเป็นเปอร์เซ็นต์ไม่มาก แต่ก็น่าสนใจศึกษาว่า โปรแกรมระบุภาษาคำโดยวิธีการใช้สายอักขระเฉพาะนี้ได้ตัดสินผิดพลาดจากภาษาใดเป็นภาษาใดบ้าง ซึ่งแสดงไว้ในตาราง 4.3 ดังนี้ โดยแสดงผลที่ตัดสินผิดพลาด เป็น จำนวนคำ และ เปอร์เซนต์จากคำทดสอบทั้งหมด

ตารางที่ 4.3 ผลการตัดสินภาษาของคำผิดพลาดจากการใช้สายอักขระเฉพาะ 1-5 ตัว

วิธี	คำที่ทดสอบ	ตัดสินผิดเป็น				
		ไทย	อังกฤษ	ญี่ปุ่น	ฝรั่งเศส	รวม
สายอักขระ 1 ตัวอักษร	ไทย	0	0	0	0	
	อังกฤษ	2 (0.02%)	0	0	0	2 (0.02%)
	ญี่ปุ่น	0	0	0	0	
	ฝรั่งเศส	0	0	0	0	
สายอักขระ 2 ตัวอักษร	ไทย	0	2 (0.02%)	0	0	2 (0.02%)
	อังกฤษ	10 (0.1%)	0	4 (0.04%)	0	14 (0.14%)
	ญี่ปุ่น	4 (0.04%)	3 (0.03%)	0	0	7 (0.07%)
	ฝรั่งเศส	1 (0.10%)	0	0	0	1 (0.10%)
สายอักขระ 3 ตัวอักษร	ไทย	0	154 (1.54%)	106 (1.06%)	14 (0.14%)	274 (2.74%)
	อังกฤษ	139 (1.39%)	0	94 (0.94%)	31 (0.31%)	264 (2.64%)
	ญี่ปุ่น	88 (0.88%)	92 (0.92%)	0	6 (0.06%)	186 (1.86%)
	ฝรั่งเศส	22 (2.20%)	90 (9%)	10 (1%)	0	122 (12.2%)
สายอักขระ 4 ตัวอักษร	ไทย	0	87 (0.87%)	63 (0.63%)	14 (0.14%)	164 (1.64%)
	อังกฤษ	86 (0.86%)	0	116 (1.16%)	75 (0.75%)	277 (2.77%)
	ญี่ปุ่น	61 (0.61%)	102 (1.02%)	0	5 (0.05%)	168 (1.68%)
	ฝรั่งเศส	12 (1.20%)	76 (7.60%)	11 (1.10%)	0	99 (9.90%)
สายอักขระ 5 ตัวอักษร	ไทย	0	12 (0.12%)	10 (0.10%)	1 (0.01%)	23 (0.23%)
	อังกฤษ	11 (0.11%)	0	25 (0.25%)	45 (0.45%)	81 (0.81%)
	ญี่ปุ่น	10 (0.10%)	28 (0.28%)	0	0	38 (0.38%)
	ฝรั่งเศส	0	34 (3.40%)	0	0	34 (3.40%)
รวมสาย อักขระ 1-5 ตัวอักษร	ไทย	0	255 (2.55%)	179 (1.79%)	29 (0.29%)	463 (4.63%)
	อังกฤษ	248 (2.48%)	0	239 (2.39%)	151(1.51%)	638 (6.38%)
	ญี่ปุ่น	163 (1.63%)	225 (2.25%)	0	11 (0.11%)	399 (3.99%)
	ฝรั่งเศส	35 (3.50%)	200 (20%)	21 (2.10%)	0	256(25.60%)

จากตารางที่ 4.3 แสดงให้เห็นข้อผิดพลาดในการตัดสินภาษา ซึ่งผู้วิจัยได้วิเคราะห์ข้อผิดพลาดในแต่ละภาษาไว้ดังนี้

4.5.1 ข้อผิดพลาดจากการระบุคำไทย

จากข้อมูลการระบุภาษาของคำไทยที่ผิดพลาดจากการใช้สายอักขระเฉพาะ ส่วนใหญ่โปรแกรมตัดสินภาษาผิดพลาดว่าเป็นคำภาษาอังกฤษและภาษาญี่ปุ่น โดยตัดสินผิดว่าเป็นภาษาอังกฤษมากที่สุด (2.55%) เช่น คำว่า ล้อเลียน ถูกตัดสินว่าเป็นคำภาษาอังกฤษจากสายอักขระเฉพาะ “อเลี้ย” และ ตะนอย ถูกตัดสินผิดว่าเป็นคำภาษาอังกฤษจากสายอักขระเฉพาะ “นอย” เป็นเพราะบังเอิญในคลังข้อมูลการฝึกภาษาอังกฤษมีคำว่า กรอเลียร์ เอเลี่ยน อิลลินอยส์ และในคลังข้อมูลการฝึกภาษาไทยไม่มีคำที่มีอักขระนี้อยู่

4.5.2 ข้อผิดพลาดจากการระบุคำทับศัพท์ภาษาอังกฤษ

จากข้อมูลการระบุภาษาของคำไทยที่ผิดพลาดด้วยสายอักขระเฉพาะ ส่วนใหญ่โปรแกรมตัดสินภาษาผิดพลาดว่าเป็นคำไทย คำภาษาญี่ปุ่นและภาษาฝรั่งเศส โดยตัดสินผิดว่าเป็นคำไทยมากที่สุด (2.48%) เช่น คำว่า เกียร์ ถูกตัดสินว่าเป็นคำไทยจากสายอักขระเฉพาะ “กีเยร” และ คำว่า ลอการิทึม ถูกตัดสินว่าเป็นคำไทยจากสายอักขระเฉพาะ “อการ” เป็นเพราะบังเอิญในคลังข้อมูลการฝึกคำไทยมีคำว่า เกียรตฤณ เกียรตประวัตติ หอการคำและสื่อการศึกษา เป็นต้น แต่ในคลังข้อมูลการฝึกภาษาอังกฤษไม่มีสายอักขระนี้

4.5.3 ข้อผิดพลาดจากการระบุคำทับศัพท์ภาษาญี่ปุ่น

จากข้อมูลการระบุภาษาของคำไทยที่ผิดพลาดด้วยสายอักขระเฉพาะ ส่วนใหญ่โปรแกรมตัดสินภาษาผิดพลาดว่าเป็นคำภาษาอังกฤษ และคำไทย โดยตัดสินผิดว่าเป็นคำภาษาอังกฤษมากที่สุด (2.25%) เช่น คำว่า เซนทะ ถูกตัดสินว่าเป็นคำภาษาอังกฤษจากสายอักขระเฉพาะ “เซนท” เป็นเพราะบังเอิญในคลังข้อมูลการฝึกภาษาอังกฤษมีคำว่า เซนทรอยด์ และเซนท์ เป็นต้น แต่ในคลังข้อมูลการฝึกภาษาญี่ปุ่นไม่มีสายอักขระนี้

4.5.4 ข้อผิดพลาดจากการระบุคำทับศัพท์ภาษาฝรั่งเศส

จากข้อมูลการระบุภาษาของคำไทยที่ผิดพลาดด้วยสายอักขระเฉพาะ ส่วนใหญ่โปรแกรมตัดสินภาษาผิดพลาดว่าเป็นคำไทย คำภาษาอังกฤษและภาษาญี่ปุ่น โดยตัดสินผิดว่าเป็นคำภาษาอังกฤษมากที่สุด (20%) เช่นคำว่า อาเนต ถูกตัดสินว่าเป็นคำภาษาอังกฤษจากสายอักขระเฉพาะ “เนต” และคำว่า เฟรียา ถูกตัดสินว่าเป็นคำภาษาอังกฤษจากสายอักขระเฉพาะ “เฟร” เป็นเพราะบังเอิญในคลังข้อมูลการฝึกภาษาอังกฤษมีคำว่า ลามิเนต คาร์บอนเนต เดลีเฟรช และ เฟรชชี เป็นต้น แต่ในคลังข้อมูลการฝึกภาษาฝรั่งเศสไม่มีสายอักขระนี้

เหตุผลที่โปรแกรมตัดสินภาษาผิดพลาดหรือตัดสินภาษาไม่ได้ เป็นเพราะ คลังข้อมูลการฝึกอาจมีขนาดไม่มากพอที่จะทำให้พบสายอักขระเฉพาะที่ครอบคลุมทั้งภาษา ดังจะเห็นได้ชัดเจนในกรณีของคำทับศัพท์ภาษาฝรั่งเศสที่ตัดสินผิดพลาดถึง 25.6% นอกจากนี้ จากผลข้อผิดพลาดที่แสดงในตารางที่ 4.3 เราอาจกล่าวได้ว่า คำทับศัพท์ภาษาญี่ปุ่นมีลักษณะเฉพาะของสายอักขระที่ชัดเจนมากกว่าคำไทย และคำทับศัพท์ภาษาอังกฤษ เพราะมีข้อผิดพลาดเพียง 3.99% ในขณะที่คำทับศัพท์ภาษาอังกฤษจะมีลักษณะเฉพาะของสายอักขระที่น้อยกว่า ดังจะเห็นได้จากข้อผิดพลาดที่พบ 6.38% และจากการที่คำทับศัพท์ภาษาญี่ปุ่นและคำไทยก็ถูกตัดสินผิดเป็นคำทับศัพท์ภาษาอังกฤษมากที่สุด

นอกจากข้อผิดพลาดดังกล่าว ก็มีประเด็นที่น่าสนใจว่าสายอักขระเฉพาะทุกตัวที่พบในแต่ละภาษา สามารถนำมาใช้จริงได้หรือไม่ และแต่ละภาษามีการใช้สายอักขระเฉพาะใดมากเป็นพิเศษ ดังนั้นผู้วิจัยจึงได้นำเสนอสายอักขระเฉพาะที่พบว่าใช้จริงในระบบและจำนวนครั้งที่ใช้

4.6 สายอักขระเฉพาะที่ใช้จริงในระบบ

จากสายอักขระเฉพาะ 1-5 ตัวในแต่ละภาษาที่พบจากคลังข้อมูลการฝึก ผู้วิจัยพบว่าสายอักขระเฉพาะเหล่านี้ไม่ได้ถูกนำมาใช้จริงในระบบทุกตัว และจำนวนการใช้สายอักขระเฉพาะแต่ละตัวจะไม่เท่ากัน ดังนั้น ผู้วิจัยจึงเสนอตารางการเปรียบเทียบจำนวนสายอักขระเฉพาะที่พบในคลังข้อมูลการฝึกและจำนวนของสายอักขระเฉพาะที่ใช้จริงในแต่ละภาษา เห็นได้จากตารางที่ 4.4 เพื่อแสดงประสิทธิภาพของสายอักขระเฉพาะ เช่น อักขระเฉพาะที่พบจำนวนมาก แต่ถูกนำมาใช้น้อย หรือในทางกลับกัน อักขระเฉพาะที่พบจำนวนน้อย แต่สามารถใช้ระบุภาษาได้มาก

โดยในคอลัมน์ที่สองจะแสดงจำนวนสายอักขระเฉพาะที่พบในคลังข้อมูลการฝึก ตัวเลขในวงเล็บคือ จำนวนเปอร์เซ็นต์จากสายอักขระที่พบในคลังข้อมูลการฝึกทั้งหมด 1-5 ตัว และในคอลัมน์ที่สามแสดงจำนวนสายอักขระเฉพาะที่ใช้จริงในระบบ ในวงเล็บคือจำนวนเปอร์เซ็นต์คิดจากสายอักขระเฉพาะที่พบในคลังข้อมูลการฝึก

ตาราง 4.4 การเปรียบเทียบจำนวนสายอักขระเฉพาะที่พบในคลังข้อมูลการฝึกและจำนวนของสายอักขระเฉพาะที่ใช้จริงในแต่ละภาษา

สายอักขระเฉพาะ (ตัวอักษร)		จำนวนสายอักขระเฉพาะที่ พบในคลังข้อมูลการฝึก	จำนวนสายอักขระเฉพาะที่ ใช้จริงในระบบ
คำไทย	1	16 (0.13%)	13 (81.25%)
	2	542 (4.39%)	173 (31.92%)
	3	5736 (46.5%)	872 (15.13%)
	4	5522 (44.76%)	389 (7.05%)
	5	520 (4.22%)	70 (13.46%)
คำทับศัพท์ภาษาอังกฤษ	1	0	0
	2	119 (0.76%)	22 (18.49%)
	3	5247 (33.5%)	781 (14.88%)
	4	8944 (57.10%)	903 (10.10%)
	5	1353 (8.64%)	335 (24.76%)
คำทับศัพท์ภาษาญี่ปุ่น	1	0	0
	2	105 (0.95%)	32 (30.48%)
	3	3537 (31.82%)	781 (22.08%)
	4	6726 (60.51%)	864 (12.85%)
	5	747 (6.72%)	107 (14.32%)
คำทับศัพท์ภาษา ฝรั่งเศส	1	0	0
	2	14 (0.49%)	0
	3	679 (23.91%)	36 (5.3%)
	4	1785 (62.85%)	68 (3.81%)
	5	362 (12.75%)	28 (7.73%)

จากตารางที่ 4.4 พบว่าสายอักขระที่ใช้จริงในระบบทุกภาษาส่วนใหญ่ เป็นสายอักขระเฉพาะที่พบมากในคลังข้อมูลการฝึก แต่นำไปใช้ได้จริงน้อย คือไม่ถึง 50% ของสายอักขระเฉพาะในคลังข้อมูลการฝึกทั้งหมด เช่น ในสายอักขระเฉพาะคำทับศัพท์ภาษาอังกฤษ ถูกนำไปใช้จริง

ในระบบสูงสุดแค่ 24.76% ในสายอักขระเฉพาะ 5-ตัว ในสายอักขระเฉพาะคำทับศัพท์ภาษาญี่ปุ่น ถูกนำไปใช้จริงในระบบสูงสุดแค่ 30.48% ในสายอักขระเฉพาะ 2-ตัว ยกเว้นกรณีของสายอักขระเฉพาะคำไทย ซึ่งพบข้อมูลสายอักขระเฉพาะ 1-ตัวในคลังข้อมูลการฝึกเพียงแค่ 16 ตัว แต่นำไปใช้จริงในระบบถึง 13 ตัวหรือ 81.25% ได้แก่ ษ ผ ฐ ฤ ไ ฦ ๓ ภ ฎ ฏ ฃ พ ฒ ซึ่งสายอักขระเฉพาะเหล่านี้เป็นอักขระที่มาจากคำยืมภาษาบาลีและสันสกฤต อาจเป็นเพราะภาษาไทยมีการยืมคำจากบาลีและสันสกฤตมาใช้เป็นจำนวนมาก ดังนั้นคำส่วนใหญ่จึงใช้สายอักขระเฉพาะเหล่านี้ ส่วนกรณีของคำทับศัพท์ภาษาฝรั่งเศส สายอักขระเฉพาะถูกนำไปใช้จริงสูงสุดแค่ 7.73% ในสายอักขระเฉพาะ 5-ตัว ซึ่งน้อยที่สุดเมื่อเทียบกับภาษาอื่น เป็นเพราะใช้ขนาดของคลังข้อมูลน้อยที่สุด

นอกจากนี้จากการวิเคราะห์สายอักขระเฉพาะที่ใช้จริงในระบบ ผู้วิจัยพบว่า มีประเด็นที่น่าสนใจ เช่น สายอักขระใดที่ใช้บอกภาษาได้ดีที่สุด ดังนั้นผู้วิจัยจึงแสดงตัวอย่างสายอักขระที่ใช้จริงในระบบที่น่าสนใจไว้บางส่วน ส่วนรายการสายอักขระเฉพาะที่ใช้จริงทั้งหมดจะกล่าวไว้ในภาคผนวก

สายอักขระคำไทยที่ถูกนำมาใช้ในระบบการใช้สายอักขระเฉพาะ โดยที่ผลในตารางคือจำนวนครั้งที่ใช้ระบุภาษาได้ถูกต้อง ส่วนตัวเลขในวงเล็บคือจำนวนเปอร์เซ็นต์เมื่อคิดจากจำนวนรวมของสายอักขระ 1-5 ตัวที่ใช้ระบุภาษาถูกต้องทั้งหมด คือ 5058 คำ

สายอักขระเฉพาะคำไทย 1-ตัว (จำนวนครั้งที่ระบุภาษาได้ถูกต้องทั้งหมด 1414 ครั้ง)

สายอักขระ	จำนวนครั้งที่ใช้ ระบุภาษาได้ ถูกต้อง (เปอร์เซ็นต์)					
		ฤ	125 (2.47%)	ฃ	65 (1.29%)	
		ไ	118 (2.33%)	พ	32 (0.63%)	
		ฦ	110 (2.17%)	ฒ	17 (0.34%)	
		๓	104 (2.06%)			
ษ	298 (5.89%)	ภ	93 (1.84%)			
ผ	176 (3.48%)	ฎ	92 (1.82%)			
ฐ	127 (2.51%)	ฏ	57 (1.13%)			

จากสายอักขระเฉพาะ 1-ตัว เห็นได้ว่า สายอักขระเฉพาะ 1-ตัวใช้ตัดสินภาษาได้มากเมื่อเทียบกับสายอักขระเฉพาะขนาดอื่นๆ โดยเฉพาะ ษ ใช้ตัดสินภาษาได้มากที่สุดถึง 5.89%

สายอักขระเฉพาะคำไทย 2 ตัว (จำนวนครั้งที่ระบุภาษาได้ถูกต้องทั้งหมด 756 ครั้ง)

สายอักขระเฉพาะ	จำนวนครั้งที่ ใช้ระบุภาษา ได้ถูกต้อง (เปอร์เซ็นต์)	เ-ด	23 (0.45%)	ร-ภ	19 (0.38%)	มี-	12 (0.24%)
อ-ก	22 (0.43%)	ภ-	22 (0.43%)	ข-อ	19 (0.38%)	ณ-	11 (0.22%)
เ-ก	21 (0.42%)	า-ณ	21 (0.42%)	ร-ศ	16 (0.32%)	จ-ภ	11 (0.22%)
อ-	21 (0.42%)	อ-	21 (0.42%)	ศ-ร	15 (0.3%)	ถ-	10 (0.2%)
ส-	35 (0.69%)	ศ-า	21 (0.42%)	น-	13 (0.26%)	ถ-ว	10 (0.2%)
ภ-า	32 (0.63%)	ภ-	20 (0.4%)	เ-ภ	13 (0.26%)	ช-ณ	10 (0.2%)

จากสายอักขระเฉพาะ 2-ตัว สังเกตเห็นว่า ไม่มีลักษณะใดเด่นชัด แต่พบว่ามีสายอักขระเฉพาะที่ถูกนำไปใช้มากกว่าสายอักขระเฉพาะอื่นๆ คือ ส- (0.69%) ภ-า (0.63%)

สายอักขระเฉพาะคำไทย 3 ตัว (จำนวนครั้งที่ระบุภาษาได้ถูกต้องทั้งหมด 2104 ครั้ง)

สายอักขระเฉพาะ	จำนวนครั้งที่ ใช้ระบุภาษา ได้ถูกต้อง (เปอร์เซ็นต์)	น-า-ย	17 (0.34%)	ว-ห	10 (0.2%)	ม-ด-	7 (0.14%)
บ-ร-ร	39 (0.77%)	อ-ร-ร	15 (0.3%)	อ-ย	10 (0.2%)	บ-ป-ร	7 (0.14%)
'-า-ง	34 (0.67%)	'-ย-ง	15 (0.3%)	ส-ม	9 (0.18%)	น-ห-า	7 (0.14%)
จ--า	31 (0.61%)	ห-ล-ว	14 (0.28%)	ป-ระ	9 (0.18%)	ช--า	7 (0.14%)
จ-ท-ธ	27 (0.53%)	ร--อ	14 (0.28%)	เ-พ-	9 (0.18%)	ค-น-ธ	7 (0.14%)
ก-ล-	20 (0.4%)	ญ-ช	14 (0.28%)	ะ-พ-ร	8 (0.16%)	โ-ห-ม	7 (0.14%)
า-ร-ย	19 (0.38%)	ห-ล-	12 (0.24%)	ห-ว	8 (0.16%)	'-า-ว	7 (0.14%)
'-า-ม	18 (0.36%)	ว-ร-ร	10 (0.2%)	ค-ร-	8 (0.16%)	ว-ป	7 (0.14%)
ด-ร	18 (0.36%)	ล--ย	10 (0.2%)	'-า-ย	8 (0.16%)	ด-ต	7 (0.14%)
		พ-จ-น	10 (0.2%)	า-ธ-ร	7 (0.14%)	า-ย-ต	6 (0.12%)
		เ-ห-ล	10 (0.2%)	ะ-ท-	7 (0.14%)	า-ค-ม	6 (0.12%)
		จ-ม	10 (0.2%)	ล--	7 (0.14%)	ะ-ก-ร	6 (0.12%)

จากสายอักขระเฉพาะที่ใช้จริง 3-ตัว พบข้อสังเกตคือ มีการใช้สายอักขระเฉพาะที่ประกอบด้วย ตัวอักษรใดก็ตาม + ร หัน หรือ _ร เช่นสายอักขระเฉพาะ บรร อรร วรر ซึ่งลักษณะแบบนี้แสดงถึงการใช้ภาษาไทย

สายอักขระเฉพาะคำไทย 4 ตัว (จำนวนครั้งที่ระบุภาษาได้ถูกต้องทั้งหมด 651 ครั้ง)

สายอักขระเฉพาะ	จำนวนครั้งที่ใช้ระบุภาษาได้ถูกต้อง (เปอร์เซ็นต์)	เ-ต--า	9 (0.18%)	ต-ะ-ก-ร	4 (0.08%)	ว-า-ม-เ	3 (0.06%)
		ั-ม-พ-	9 (0.18%)	ง-ก-ร-า	4 (0.08%)	ว-ร-า-ช	3 (0.06%)
		อ-า-ร-	7 (0.14%)	ค-ร-อ-ง	4 (0.08%)	ระ-ะ-ติ-	3 (0.06%)
		เ-ส-น-า	5 (0.1%)	กั-น-แ	4 (0.08%)	ระ-ะ-เ-ห	3 (0.06%)
ก-ระ-ะ-โ	16 (0.32%)	เ-พ-ล-	5 (0.1%)	ก-ร-ะ-ต	4 (0.08%)	ร-ส-ุ-ม	3 (0.06%)
น-ท-ร-	12 (0.24%)	ิ-ม-พ-	5 (0.1%)	ก-ร-ะ-ช	4 (0.08%)	มู-ล-	3 (0.06%)
ก-ระ-ะ-ด	12 (0.24%)	ส-ม-ย	4 (0.08%)	า-ม-ย	3 (0.06%)	พิ-ร-	3 (0.06%)
ั-ต-น-	11 (0.22%)	ร-า-ล-	4 (0.08%)	อ-น-เ-ก	3 (0.06%)	พิ-ต-ร	3 (0.06%)
ก-ระ-ะ-ส	9 (0.18%)	พิ-น-ท	4 (0.08%)	อ-ต-อ-อ	3 (0.06%)	พ-น-ช	3 (0.06%)

จากสายอักขระเฉพาะคำไทย 4-ตัว พบข้อสังเกต คือ มีการใช้ กระ+ตัวอักษรใดก็ตาม หรือ กระ_ ในการตัดสัทภาษาได้มาก เช่น กระโ กระด กระส กระต เป็นต้น เป็นเพราะคำไทยมีการใช้ ร ควบกล้ำร่วมกับ สระอะ บ่อย

ส่วนสายอักขระเฉพาะคำไทย 5 ตัว (จำนวนครั้งที่ระบุภาษาได้ถูกต้องทั้งหมด 133 ครั้ง) ซึ่งไม่ค่อยมีอะไรเด่น เนื่องจากแต่ละสายอักขระเฉพาะมีการใช้จริงแค่ 1-7 ครั้งเท่านั้น

สายอักขระคำทับศัพท์ภาษาอังกฤษที่ถูกนำมาใช้ในระบบการใช้สายอักขระเฉพาะ โดยที่ผลในตารางคือ จำนวนครั้งที่ใช้ระบุภาษาได้ถูกต้อง ส่วนตัวเลขในวงเล็บคือจำนวนเปอร์เซ็นต์ เมื่อคิดจากจำนวนรวมของสายอักขระ 1-5 ตัวที่ใช้ระบุภาษาถูกต้องทั้งหมด คือ 4871 คำ

สายอักขระเฉพาะภาษาอังกฤษ 2 ตัว (จำนวนครั้งที่ระบุภาษาได้ถูกต้องทั้งหมด 122 ครั้ง)

สายอักขระเฉพาะ	จำนวนครั้งที่ใช้ระบุภาษาได้ถูกต้อง (เปอร์เซ็นต์)	ฟ-ต	11 (0.23%)	๓-ป	5 (0.1%)
		๓-ป	10 (0.21%)	พ - ๓	4 (0.08%)
๓-ไ	23 (0.47%)	ฮ-ฟ	7 (0.14%)	พ-พ	2 (0.04%)
๓-ค	22 (0.45%)	๓-ค	7 (0.14%)	๓-'	2 (0.04%)
		พ-ด	6 (0.12%)	๓-ส	2 (0.04%)

จากสายอักขระเฉพาะ 2-ตัว พบข้อสังเกตคือ มีการใช้ ๓- ตามด้วยตัวอักษรใดก็ตาม เช่น ๓-ค ๓-ป และ ๓- ตามด้วยตัวอักษร เช่น ๓-ค ๓-ป ๓-ส เป็นต้น

สายอักขระเฉพาะภาษาอังกฤษ 3 ตัว (จำนวนครั้งที่ระบุภาษาได้ถูกต้องทั้งหมด 2018 ครั้ง)

สายอักขระเฉพาะ	จำนวนครั้งที่ ใช้ระบุภาษา ได้ถูกต้อง (เปอร์เซ็นต์)	อี-อ-ค	18 (0.37%)	อ-ล-ต	10 (0.21%)	โ-ฟ-ล	9 (0.19%)
เนอะ		ส-เ-ป	17 (0.35%)	ส-ลิ-	10 (0.21%)	เ-อ-ล	9 (0.19%)
น-ด	113 (2.32%)	ิ-ส-ช	14 (0.27%)	ว-ท-	10 (0.21%)	-อ-ป	9 (0.19%)
อ-ล-ล	31 (0.64%)	ส-ส-	13 (0.27%)	ม-มี-	10 (0.21%)	-ค-ท	9 (0.19%)
อี-ล	28 (0.57%)	บ-ล-ล	13 (0.27%)	ิ-ส-ค	10 (0.21%)	ท-ร-	8 (0.16%)
เ-ว	28 (0.57%)	เ-จ-อ	13 (0.27%)	ิ-ล-ต	10 (0.21%)	ท-ติ-	8 (0.16%)
แ-อ-น	20 (0.41%)	ฟ-ฟ-	12 (0.25%)	-ป-า	10 (0.21%)	ท-ช-	8 (0.16%)
อ-ล-	18 (0.37%)	บ-อ-ล	12 (0.25%)	ว-ย-	9 (0.19%)	ค-อ-ร	8 (0.16%)
ล-ล-	18 (0.37%)	ค-ว-อ	12 (0.25%)	ล-ล-	9 (0.19%)	-อ-ต	8 (0.16%)
		เ-โ-ม	12 (0.25%)	ช-ช-	9 (0.19%)	ว-ส-ต	7 (0.14%)
		ู-ต	11 (0.23%)	ไ-ล-เ	9 (0.19%)	พ-ล-ย	7 (0.14%)

จากสายอักขระเฉพาะ 3-ตัว จะสังเกตเห็นรูปแบบของภาษา คือ มีการใช้ _ _ มาก เช่น นด์ อล สล ฟฟ วย ลล เป็นต้น โดยเฉพาะ นด์ ใช้ตัดสินภาษาได้ 2.32% ซึ่งสูงที่สุดเมื่อเทียบกับสายอักขระเฉพาะขนาดอื่นๆ และพบสายอักขระเฉพาะ ี-อ เช่น ี-อ-ป ี-อ-ค ี-อ-ต เป็นต้น

สายอักขระเฉพาะภาษาอังกฤษ 4 ตัว (จำนวนครั้งที่ระบุภาษาได้ถูกต้องทั้งหมด 1769 ครั้ง)

สายอักขระเฉพาะ	จำนวนครั้งที่ ใช้ระบุภาษา ได้ถูกต้อง (เปอร์เซ็นต์)	เ-ด-ช- <th>10 (0.21%)</th> <th>อ-ล-ร-ส <th>7 (0.14%)</th> <th>น-เ-ด- <th>6 (0.12%)</th> </th></th>	10 (0.21%)	อ-ล-ร-ส <th>7 (0.14%)</th> <th>น-เ-ด- <th>6 (0.12%)</th> </th>	7 (0.14%)	น-เ-ด- <th>6 (0.12%)</th>	6 (0.12%)
เ-อ-ก	33 (0.68%)	อ-ม-มี-	9 (0.18%)	พ-า-ร-	7 (0.14%)	ค-อ-ม-เ	6 (0.12%)
ท-อ-ร-	31 (0.64%)	ิ-ล-ล-	9 (0.18%)	พ-ร-เ-	7 (0.14%)	ค-ล-น	6 (0.12%)
อ-ร-ก	16 (0.33%)	ิ-ร-น	9 (0.18%)	บ-ร-อ-ด	7 (0.14%)	เ-ว-อ-ล	6 (0.12%)
เ-น-ด	16 (0.33%)	ล-ย-	8 (0.16%)	อ-ป-ป-	6 (0.12%)	เ-ว-โ-ล	6 (0.12%)
อ-ร-ส	12 (0.25%)	มี-ล-ล	8 (0.16%)	ส-โ-ต-น	6 (0.12%)	เ-ม-ติ-	6 (0.12%)
เ-ม-ต-ร	11 (0.23%)	ค-เ-ต-อ	8 (0.16%)	ล-เ-ว	6 (0.12%)	เ-ด-น	6 (0.12%)
ส-ต-า-ร	10 (0.21%)	ิ-เ-ต-อ	8 (0.16%)	มี-ย-ม	6 (0.12%)	ิ-เ-อ-ร	6 (0.12%)
		-เ-ป-	8 (0.16%)	ม-น-เ-ด	6 (0.12%)	-น-ส-	6 (0.12%)
		า-เ-น-	7 (0.14%)	บ-อ-ด-	6 (0.12%)	ฮ-า-ร-	5 (0.1%)
		อ-ร-ส	7 (0.14%)	น-น-	6 (0.12%)	อ-ร-ส	5 (0.1%)

จากสายอักขระเฉพาะ 4-ตัว ไม่มีลักษณะเด่นของภาษาให้เห็นชัดเจน แต่สังเกตได้ว่า มีสายอักขระเฉพาะบางตัวที่ใช้ได้จริงมากกว่าสายอักขระอื่นๆ อย่างเห็นได้ชัด คือ เ-อ-ก (0.68%) และ ท-อ-ร- (0.64%)

สายอักขระเฉพาะภาษาอังกฤษ 5 ตัว (จำนวนครั้งที่ระบุภาษาได้ถูกต้องทั้งหมด 962 ครั้ง)

สายอักขระเฉพาะ	จำนวนครั้งที่ ใช้ระบุภาษา ได้ถูกต้อง (เปอร์เซ็นต์)	สายอักขระเฉพาะภาษาอังกฤษ 5 ตัว (จำนวนครั้งที่ระบุภาษาได้ถูกต้องทั้งหมด 962 ครั้ง)		
		ก-ล-ต-อ-ร	ต-อ-ร-เ-	ด-ด-ิ-ง
ส-เ-ต-อ-ร	30 (0.62%)	ก-ล-ต-อ-ร 13 (0.27%)	ต-อ-ร-เ- 9 (0.18%)	ด-ด-ิ-ง 6 (0.12%)
เ-ต-อ-ร-เ	27 (0.55%)	เ-ช-อ-ร-เ 13 (0.27%)	ว-อ-ร-เ-ด 8 (0.16%)	ก-เ-ก-เ-ด 6 (0.12%)
เ-ล-อ-ร-เ	21 (0.43%)	ฟ-อ-ร-เ-ม 12 (0.25%)	ฟ-อ-ร-เ-แ 8 (0.16%)	อ-ิ-น-ด-เ 5 (0.1%)
ล-เ-ล-อ-ร	21 (0.43%)	เ-ช-อ-ร-เ 12 (0.25%)	ช-อ-ร-เ-เ 8 (0.16%)	อ-ร-เ-ริ- 5 (0.1%)
เ-ป-อ-ร-เ	20 (0.41%)	ว-อ-ร-เ-เ 11 (0.23%)	ค-า-ร-เ-เ 8 (0.16%)	อ-ร-เ-ติ- 5 (0.1%)
ฟ-อ-ร-เ-ด	16 (0.33%)	บ-อ-ร-เ-เ 11 (0.23%)	อ-ร-เ-เ-ม 7 (0.14%)	อ-ร-เ-แ-ม 5 (0.1%)
ฟ-อ-ร-เ-เ	14 (0.29%)	เ-ป-อ-ร-เ 11 (0.23%)	อ-ร-เ-เ-ช 7 (0.14%)	ว-เ-ช-อ-ร 5 (0.1%)
น-เ-ช-อ-ร	14 (0.29%)	อ-ร-เ-เ-ก 10 (0.21%)	อ-เ-ว-อ-ร 7 (0.14%)	ม-อ-ร-เ-เ 5 (0.1%)
		ว-อ-ร-เ-ช 10 (0.21%)	ล-เ-ก-ช-เ 7 (0.14%)	ม-เ-ป-อ-ร 5 (0.1%)
		ม-เ-ม-อ-ร 10 (0.21%)	ค-า-ร-เ-เ 7 (0.14%)	ป-อ-ร-เ-ม 5 (0.1%)
		ค-อ-ร-เ-เ 10 (0.21%)	เ-ริ-เ-ย-ม 7 (0.14%)	ป-อ-ร-เ-เ 5 (0.1%)
		เ-น-เ-ด-อ 10 (0.21%)	อ-อ-ร-เ-แ 6 (0.12%)	บ-อ-ร-เ-ร 5 (0.1%)

จากสายอักขระเฉพาะ 5-ตัว จะสังเกตเห็นรูปแบบของภาษา คือการใช้ เ-อ-ร เช่น เตอ-ร เบอ-ร เลอ-ร เซอ-ร เซอ-ร เป็นต้น การใช้ เ-อ-ร เช่น ฟอ-ร-ด วอ-ร-เ บอ-ร-เ วอ-ร-ช-เ ตอ-ร-เ เป็นต้น และ เ-อ-ร เช่น ส-เ-ต-อ-ร น-เ-ช-อ-ร ก-เ-ด-อ-ร น-เ-ช-อ-ร ม-เ-ป-อ-ร เป็นต้น

สายอักขระคำทับศัพท์ภาษาญี่ปุ่นที่ถูกนำมาใช้ในระบบการใช้สายอักขระเฉพาะ โดยที่ผลในตารางคือ จำนวนครั้งที่ระบุภาษาได้ถูกต้อง ส่วนตัวเลขในวงเล็บคือจำนวนเปอร์เซ็นต์ เมื่อคิดจากจำนวนของสายอักขระ 1-5 ตัวที่ใช้ระบุภาษาถูกต้องทั้งหมด คือ 5409 คำ

สายอักขระเฉพาะภาษาญี่ปุ่น 2 ตัว (จำนวนครั้งที่ระบุภาษาได้ถูกต้องทั้งหมด 160 ครั้ง)

สายอักขระเฉพาะ	จำนวนครั้งที่ ใช้ระบุภาษา ได้ถูกต้อง (เปอร์เซ็นต์)	สายอักขระเฉพาะภาษาญี่ปุ่น 2 ตัว (จำนวนครั้งที่ระบุภาษาได้ถูกต้องทั้งหมด 160 ครั้ง)		
		จ-ฉ	เ-ย	ท-ฮ
จ-ฉ	58 (1.07%)	จ-ฉ 12 (0.22%)	เ-ย 5 (0.09%)	ท-ฮ 2 (0.04%)
เ-ย	16 (0.3%)	เ-ย 8 (0.15%)	ช-ฮ 4 (0.07%)	ช-ฮ 2 (0.04%)
		จ-ฉ 6 (0.11%)	ง-ญ 3 (0.06%)	ง-ญ 2 (0.04%)
		เ-ฉ 5 (0.09%)	เ-ง 3 (0.06%)	เ-ฉ 2 (0.04%)
		ค-ะ 5 (0.09%)	เ-ฮ 3 (0.06%)	เ-อ 2 (0.04%)
		ช-จ 5 (0.09%)	เ-ช 3 (0.06%)	เ-ท 2 (0.04%)

จากสายอักขระเฉพาะ 2-ตัว จะสังเกตเห็นรูปแบบของภาษา คือการใช้ เ-ะ เช่น ฮะ-ดะ โดยที่ ฮะ ใช้ตัดสิ้นภาษา และ เ-ะ เช่น เ-ะ-จ เ-ะ-ย เ-ะ-ฮ เ-ะ-ช

สายอักขระเฉพาะภาษาญี่ปุ่น 3 ตัว (จำนวนครั้งที่ระบุภาษาได้ถูกต้องทั้งหมด 3392 ครั้ง)

สายอักขระเฉพาะ	จำนวนครั้งที่ ใช้ระบุภาษา ได้ถูกต้อง (เปอร์เซ็นต์)	อี--ว	6 (0.11%)	ึ--มิ	6 (0.11%)	ง--โ	5 (0.09%)
เฉพาะ		อ-ะ-อ	6 (0.11%)	ึ--บุ	6 (0.11%)	ง-ง-ะ	5 (0.09%)
ม-ช	7 (0.13%)	ย-จิ	6 (0.11%)	ึ--คิ	6 (0.11%)	ง-โ-ญ	5 (0.09%)
ท-สุ	7 (0.13%)	ม-ะ-ส	6 (0.11%)	ิ-ด-ะ	6 (0.11%)	ค-ย	5 (0.09%)
ช-คิ	7 (0.13%)	น-ะ-ช	6 (0.11%)	ึ--ง-ญ	6 (0.11%)	ค-ะ-อ	5 (0.09%)
ช-เ	7 (0.13%)	ท-อิ	6 (0.11%)	า-ก-	5 (0.09%)	ก-า-โ	5 (0.09%)
จิ-ช	7 (0.13%)	ญ-ค	6 (0.11%)	ะ-ฉิ	5 (0.09%)	โ-ย-ญ	5 (0.09%)
จิ-ว	7 (0.13%)	ช-ม	6 (0.11%)	อ-ร	5 (0.09%)	โ-ค-ก	5 (0.09%)
ค-ช	7 (0.13%)	ง-ญิ	6 (0.11%)	อ-ค	5 (0.09%)	เ-น-ง	5 (0.09%)
โ-ค-ย	7 (0.13%)	ง-โ-ง	6 (0.11%)	สุ-ช	5 (0.09%)	เ-ท-ะ	5 (0.09%)
เ-ระ-ะ	7 (0.13%)	คิ-จ	6 (0.11%)	สุ-โ	5 (0.09%)	เ-ช-ะ	5 (0.09%)
ุ-ท-ะ	7 (0.13%)	ค-า-อ	6 (0.11%)	สุ-ะ-ค	5 (0.09%)	เ-ก-โ	5 (0.09%)
ุ-ด-ะ	7 (0.13%)	ค-ะ-ส	6 (0.11%)	ว-โ-ด	5 (0.09%)	ุ-น-ย	5 (0.09%)
ุ-ไ-ก	7 (0.13%)	ก-า-ค	6 (0.11%)	น-ะ-ค	5 (0.09%)	ุ-คิ	5 (0.09%)
ุ-โ-ง	7 (0.13%)	ก-ะ-ม	6 (0.11%)	ท-ะ-ช	5 (0.09%)	ุ-โ-จ	5 (0.09%)
ึ-ย-ร	7 (0.13%)	ไ-ด-จ	6 (0.11%)	ท-ค-ะ	5 (0.09%)	ุ-เ-อ	5 (0.09%)
ฮิ-ง	6 (0.11%)	โ-ฮ-ก	6 (0.11%)	ด-โ-ฮ	5 (0.09%)	ิ-ญ	5 (0.09%)
		โ-ม-จ	6 (0.11%)	ช-ะ-ว	5 (0.09%)	ิ-ง-ง	5 (0.09%)
		ุ-งิ	6 (0.11%)	ช-ยุ	5 (0.09%)	ิ-โ-บ	5 (0.09%)

จากสายอักขระเฉพาะ 3-ตัว พบว่ามีการใช้ภาษาที่ส่วนใหญ่ประกอบด้วยสระเสียงสั้น
เห็นได้จาก สายอักขระเฉพาะที่ใช้เสียงอะ หรือ ะ เช่น ระ เทะ เซะ เป็นต้น สายอักขระ
เฉพาะที่ใช้เสียง อุโอ หรือ ุ โ เช่น สุ โ งุ โ ุ โ ุ โ ุ โ เป็นต้น และ สายอักขระเฉพาะที่ใช้เสียง
ออะ หรือ ุ ะ เช่น ุ ทะ ุ ตะ เป็นต้น

สายอักขระเฉพาะภาษาญี่ปุ่น 4 ตัว (จำนวนครั้งที่ระบุภาษาได้ถูกต้องทั้งหมด 1683 ครั้ง)

สายอักขระเฉพาะ	จำนวนครั้งที่ ใช้ระบุภาษา ได้ถูกต้อง (เปอร์เซ็นต์)	ก-ุ-ระ-ะ	14 (0.26%)	ช-ึ-มิ	9 (0.17%)	ร-ึ-น-โ	7 (0.13%)
เฉพาะ		อ-ะ-ริ	13 (0.24%)	ฮิ-ระ-ะ	8 (0.15%)	ร-ึ-ค-ค	7 (0.13%)
เ-ค-อิ	33 (0.61%)	ช-ึ-คิ	13 (0.24%)	ท-า-ก-า	8 (0.15%)	ระ-ะ-ค-ุ	7 (0.13%)
ช-ึ-น-โ	19 (0.35%)	ก-ุ-ช-ิ	13 (0.24%)	ช-ึ-ด-ะ	8 (0.15%)	ช-ึ-กิ	7 (0.13%)
ม-ุ-ระ-ะ	16 (0.3%)	อ-ุ-จิ	11 (0.2%)	ึ-น-จิ	8 (0.15%)	โ-อ-โ-ม	7 (0.13%)
คิ-โ-ย	14 (0.26%)	เ-ส-ง-เ	11 (0.2%)	ะ-ช-ึ-โ	7 (0.13%)	อ-ุ-ร-า	6 (0.11%)
		ฮิ-า-ร-า	9 (0.17%)	ริ-โ-น	7 (0.13%)	อิ-ว-า	6 (0.11%)

สายอักขระคำทับศัพท์ภาษาฝรั่งเศสที่ถูกนำมาใช้ในระบบการใช้สายอักขระเฉพาะ โดยที่ผลในตารางคือ จำนวนครั้งที่ใช้ระบุภาษาได้ถูกต้อง ส่วนตัวเลขในวงเล็บคือจำนวนเปอร์เซ็นต์ เมื่อคิดจากจำนวนรวมของสายอักขระ 1-5 ตัวที่ใช้ระบุภาษาถูกต้องทั้งหมด คือ 204 คำ เนื่องจากสายอักขระเฉพาะในคำทับศัพท์ภาษาฝรั่งเศสที่ถูกนำไปใช้จริงมีแค่สายอักขระเฉพาะ 3-5 ตัว และแต่ละตัวก็ถูกนำไปใช้จำนวนน้อย เป็นเพราะคลังข้อมูลน้อย ทำให้ไม่ค่อยพบสายอักขระที่ใช้จริงในระบบ ดังนั้นผู้วิจัยจึงวิเคราะห์สายอักขระเฉพาะ 3-5 ตัวเพียงคร่าวๆ เท่านั้น อาจไม่สามารถนำไปวิเคราะห์เพื่อหาลักษณะของการใช้ภาษาได้

สายอักขระเฉพาะ 3 ตัว (จำนวนครั้งที่ระบุภาษาได้ถูกต้องทั้งหมด 57 ครั้ง)

สายอักขระเฉพาะ	จำนวนครั้งที่ใช้ระบุภาษาได้ถูกต้อง (เปอร์เซ็นต์)	ร-แ-ต	4 (1.96%)	ด-ู-ช	2 (0.98%)
ง-ป-า	5 (2.45%)	ง-แ-ฟ	3 (1.47%)	ด-อ-โ	2 (0.98%)
ร-แ-ย	4 (1.96%)	ม-แ-ต	3 (1.47%)	แ-ต-เ	2 (0.98%)
		อ-แ-ต	2 (0.98%)	ม-ต-ม	2 (0.98%)
		ต-ร	2 (0.98%)	ร-ด-ง	2 (0.98%)

จากสายอักขระเฉพาะ 3-ตัว สังเกตได้ว่า สายอักขระที่ใช้จริงส่วนใหญ่ประกอบด้วย “แ” เป็นส่วนประกอบในสายอักขระเฉพาะ

สายอักขระเฉพาะ 4 ตัว (จำนวนครั้งที่ระบุภาษาได้ถูกต้องทั้งหมด 101 ครั้ง)

สายอักขระเฉพาะ	จำนวนครั้งที่ใช้ระบุภาษาได้ถูกต้อง (เปอร์เซ็นต์)	ช-ง-โ	3 (1.47%)	ม-อ-ง-ด	2 (0.98%)	ก-อ-ง	2 (0.98%)
อ-ง-ป	4 (1.96%)	ก-า-บ-า	3 (1.47%)	ม-อ-ง-ด	2 (0.98%)	ก-ร-ว	2 (0.98%)
ว-า-ล-อ	3 (1.47%)	บ-ร-อ	3 (1.47%)	ฟ-อ-ง-แ	2 (0.98%)	โ-ร-ต	2 (0.98%)
ติ-เ-ย	3 (1.47%)	า-แ-ค-เ	2 (0.98%)	ต-า-ส	2 (0.98%)	เ-ล-อ-เ	2 (0.98%)
		อ-ง-ช	2 (0.98%)	ช-เว-อ	2 (0.98%)	เ-ต-ส-ก	2 (0.98%)
		ว-า-ล	2 (0.98%)	ง-ช-า-ร	2 (0.98%)	เ-ต-ร-อ	2 (0.98%)
		ล-อ-ง-ช	2 (0.98%)	ก-า-ส-ช	2 (0.98%)	เ-ช-ล	2 (0.98%)

จากสายอักขระเฉพาะ 4-ตัว สังเกตได้ว่าการใช้ภาษาที่ประกอบด้วยเสียง ออง มาก เช่น ลองช มองต มองด ฟองแ และเสียง เออ เช่น เบรอ ซเวอ เลอ เทรอ เป็นต้น

สายอักขระเฉพาะ 5 ตัว (จำนวนครั้งที่ระบุภาษาได้ถูกต้องทั้งหมด 46 ครั้ง)

สายอักขระเฉพาะ	จำนวนครั้งที่ ใช้ระบุภาษา ได้ถูกต้อง (เปอร์เซ็นต์)	า-ร-โ-ต	2 (0.98%)	แ-ว-ร-ด	2 (0.98%)	ม-า-ล-อ	1 (0.49%)
		า-ต-า-ล-อ	2 (0.98%)	ม-อ-ด	2 (0.98%)	ป-ิ-อ-ม	1 (0.49%)
		ร-า-บ-อ	2 (0.98%)	า-ว-า-ร	1 (0.49%)	บ-ริ-ช-โ	1 (0.49%)
		ป-ิ-อ-ต	2 (0.98%)	า-ร-น-า	1 (0.49%)	ช-อ-ร-บ	1 (0.49%)
ร-เ-ล-อ	5 (2.45%)	บ-ู-ร-เ	2 (0.98%)	อ-ม-แ-บ-ร	1 (0.49%)		
ู-ร-เ-น	4 (1.96%)	บ-ร-ก-อ	2 (0.98%)	อ-ง-ท-ร-า	1 (0.49%)		
ก-า-ร-ก	3 (1.47%)	ท-ริ-อ-า	2 (0.98%)	ล-อ-ง-ก	1 (0.49%)		

จากสายอักขระเฉพาะ 5-ตัว สังเกตได้ว่า มีการใช้ภาษาที่ประกอบด้วย เสียง อา และ เออ มาก เช่น ร์เลอ การ์ก าดาลอ ราเบอ มาเลอ กาบาส เป็นต้น

จากสายอักขระเฉพาะที่ใช้จริงเหล่านี้ สามารถแสดงให้เห็นลักษณะสำคัญของแต่ละภาษาได้ เช่น สายอักขระเฉพาะคำไทย 3-ตัว ที่ใช้บ่อย คือ บรร กรร อรร สรร ทำให้วิเคราะห์ได้ว่า ภาษาไทยมีการใช้ ร หัน บ่อย จนเป็นลักษณะที่เกิดเฉพาะในภาษาไทย ซึ่งลักษณะเฉพาะของแต่ละภาษาเหล่านี้จะแสดงต่อไปใน หัวข้อ 4.7 การวิเคราะห์สายอักขระเฉพาะ

นอกจากนี้ จากการศึกษาสายอักขระเฉพาะที่ใช้จริงยังพบข้อสังเกตบางอย่าง เช่น ถ้าสังเกตในภาษาไทย อักขระตัวเดียวสามารถใช้ระบุภาษาได้ เมื่อเป็นอักขระ 2 ตัวอักษรอย่าง ษ ผ ฐ ฤ ฎ ใ จะไม่ปรากฏในสายอักขระแบบ 2 ตัวอักษร แต่บางตัวอักษรอย่าง ณ ภ กลับยังพบในสายอักขระแบบ 2 อักขระ แต่ไม่พบใน 3 4 และ 5 ตัวอักษร เพราะว่า โดยหลักแล้วเมื่อตัวอักษรอย่าง ษ ใช้ตัดสินเป็นคำไทยได้แล้ว ก็จะไม่พบใน 2 ตัวอักษรแล้ว ส่วน ณ_ หรือ ณ+ ตัวอะไรก็ตาม ใน 2 อักขระ ก็จะไม่พบใน 3-5 ตัวอักษรแล้ว ซึ่งก็เป็นจริง ทั้งหมดนี้เป็นข้อสังเกต ซึ่งจะกล่าวโดยละเอียดในหัวข้อ 4.7

4.7 การวิเคราะห์สายอักขระเฉพาะ

เนื่องจากการทดลองทำให้เห็นการปรากฏของสายอักขระเฉพาะในคำไทย คำทับศัพท์ ในภาษาอังกฤษ ภาษาญี่ปุ่นและภาษาฝรั่งเศส และสายอักขระเฉพาะสามารถใช้ในการระบุภาษาได้มากกว่า 50% ในบทนี้ผู้วิจัยจึงสนใจที่จะวิเคราะห์สายอักขระเฉพาะเหล่านี้ ทั้งตำแหน่งที่ปรากฏของสายอักขระและด้านการใช้ภาษา โดยผู้วิจัยทำการวิเคราะห์ด้านรูปและเสียง เพื่อเป็นข้อมูลสนับสนุนว่าทำไมสายอักขระเฉพาะเหล่านี้จึงใช้ในการระบุภาษาได้

หลังจากทดลองโปรแกรมใช้คลังข้อมูลแต่ละภาษาเพื่อสร้างสายอักขระเฉพาะ ผลจากการทดลองที่ได้คือ คลังข้อมูลสายอักขระเฉพาะทั้งหมด 4 ภาษา ซึ่งภายในคลังข้อมูลประกอบด้วย สายอักขระเฉพาะที่มีความยาวอักขระหรือตัวอักษรตั้งแต่ 1-5 ตัว จากการใช้อินแกรม 1-5 แกรม และเนื่องจากการทดลองต้องทดสอบเพื่อหาสายอักขระเฉพาะจากชุดข้อมูล 5 ชุดในแต่ละภาษา ดังนั้นสายอักขระเฉพาะทั้งหมดที่ได้ คือ สายอักขระเฉพาะทั้งหมด 4 ภาษาและภาษาละ 5 ชุดสายอักขระ แต่เพื่อให้เห็นสายอักขระเฉพาะจริงๆ ในแต่ละภาษาสายอักขระเฉพาะที่ผู้วิจัยจะนำมาพิจารณาคือ เอาสายอักขระเฉพาะจากข้อมูลคำรวมกันทั้งหมดมาพิจารณาทีเดียว เพราะสายอักขระเฉพาะที่ได้จากการทดลองชุดข้อมูลที่ 1 อาจไม่ใช่สายอักขระเฉพาะที่แท้จริงในชุด 2 3 4 และ 5 ก็ได้

จากการรวบรวมคำทับศัพท์และหาสายอักขระเฉพาะของแต่ละภาษา ทำให้เห็นลักษณะสำคัญของแต่ละภาษา เช่น ลักษณะทางเสียงของภาษา ที่สะท้อนให้เห็นเมื่อเขียนด้วยอักษรไทย และยังให้ข้อสังเกตของภาษาซึ่งไม่สามารถอธิบายเหตุผลได้ ผู้วิจัยจึงต้องการวิเคราะห์สายอักขระเฉพาะเพื่อแสดงลักษณะสำคัญทางคำไทย คำทับศัพท์ในภาษาอังกฤษ ภาษาญี่ปุ่นและภาษาฝรั่งเศสตามลำดับต่อไป

4.7.1 สายอักขระคำไทย

ในการวิเคราะห์สายอักขระคำไทย ซึ่งผู้วิจัยวิเคราะห์ด้วยสายอักขระเฉพาะขนาด 1-5 แกรม ทำให้พบว่าคำไทยมีการใช้ตัวอักษร อักษรวิธีการใช้คำไทย และข้อสังเกตของสายอักขระเฉพาะ ซึ่งสิ่งเหล่านี้จะสะท้อนให้เห็นลักษณะของภาษาไทยที่คำภาษาอื่นไม่ปรากฏ มีดังนี้

1. การใช้ตัวอักษร ทั้งพยัญชนะ สระ และวรรณยุกต์ การใช้ตัวอักษรก็มีส่วนในการสะท้อนให้เห็นลักษณะของคำไทยได้ โดยตัวอักษรที่ปรากฏเฉพาะในคำไทย ได้แก่ ฃ ฎ ฏ ฐ ฌ ญ ฎ ษ พ ไ ำ และ วรรณยุกต์จัตวา ซึ่งอักขระหรือตัวอักษรเหล่านี้ไม่ปรากฏในคำทับศัพท์ภาษาอังกฤษ ญี่ปุ่น และฝรั่งเศส เพราะมีการใช้อักขระหรือตัวอักษรอื่นแทนเสียงของอักขระเหล่านี้อยู่แล้ว และถือว่ามี ความหมายเดียวกัน เช่น ในการใช้พยัญชนะ อักขระ “ฃ” ในภาษาอื่นใช้ “ค”, อักขระ “ฎ” ในภาษาอื่นใช้ “ด”, อักขระ “ฏ” ในภาษาอื่นใช้ “ต”, อักขระ “ฐ ฌ ญ” ในภาษาอื่นใช้ “ท ถ”, อักขระ “ฎ” ในภาษาอื่นใช้ “พ”, อักขระ “ษ พ” ในภาษาอื่นใช้ “ร” และ อักขระ

6. สายอักขระ _อะ ซึ่งรูปแบบของสายอักขระนี้มาจากเสียง เออะ ในคำไทย เช่น เกอะ เกรอะ เป็นต้น ซึ่งในภาษาอื่นไม่มีเสียงนี้ เพราะภาษาอื่นต้องมีเสียงพยัญชนะสะกด จึงทำให้ไม่มีการใช้รูปสายอักขระ เออะ แต่ใช้รูป ออ แทนและทำให้เป็นเสียงสั้นด้วยไม้ไต่คู้ เช่น ล้อต เป็นต้น

7. สายอักขระ _ียบ ซึ่งในคำไทย สายอักขระนี้มาจากเสียง เอียบ เช่น ชียบ ทียบ บียบ พียบ รียบ เป็นต้น แต่สายอักขระนี้ก็มักมีพบบ้างในคำทับศัพท์ภาษาญี่ปุ่น เช่น เฮียบังโตะ และ เกียบเบ็ทชี

8. การใช้ ห นำ ในภาษาไทยมีการใช้ ห นำ เพื่อให้กลายเป็นเสียงจัตวา และเพื่อเพิ่มความหมาย เช่น หรัด หวัง หงิด หยิก หรี หมด เป็นต้น โดยที่การทับศัพท์ภาษาต่างๆ ในงานนี้มักไม่เขียนด้วย ห นำ เป็นเพราะในคำทับศัพท์มีการใช้เสียงจัตวาแต่ไม่นิยมใช้ ห นำ มาแทนเสียง จัตวา ส่วนใหญ่ละการใช้รูป ห นำ แต่ก็ยังออกเสียงจัตวาอยู่ เช่น คำว่า แกรนิต – แกรนิต คลินิก – คลินิก คาร์บอน – คาร์บอน เป็นต้น แต่มียกเว้นบางกรณี เช่น คำทับศัพท์ภาษาอังกฤษ มีการใช้คำว่า ไต้หวัน อารับ เป็นต้น ซึ่งในเกณฑ์ทับศัพท์ภาษาอังกฤษ ปัจจุบันไม่มีการใช้ ห นำ แต่คำเหล่านี้ทับศัพท์มาแต่เดิมจนเป็นที่ยอมรับกันแล้ว จึงอนุโลมให้ใช้ต่อไป และในกรณีของคำทับศัพท์ภาษาญี่ปุ่น ในคำว่า พุหฺยคิมะ เป็นต้น ซึ่งเกณฑ์ทับศัพท์ภาษาญี่ปุ่นก็ไม่มีการใช้ ห นำ แต่คำนี้อาจเกิดขึ้นเพราะผู้ใช้อาจเขียนทับศัพท์โดยไม่คำนึงถึงกฎเกณฑ์เช่นเดียวกับในภาษาอังกฤษ เพียงแค่ทับศัพท์ตามการออกเสียงในภาษาไทยเท่านั้น

9. ข้อสังเกตของการใช้สายอักขระคำไทย ในการสร้างคำไทยซึ่งเป็นการใช้สายอักขระร่วมกัน ทำให้พบว่าการใช้สายอักขระบางสายร่วมกันทำให้เกิดสายอักขระเฉพาะ ซึ่งแม้ว่าจะไม่ใช่สายอักขระเฉพาะของภาษาไทย แต่ก็เป็นสายอักขระที่พบมากในภาษาไทย จึงนำมาแสดงไว้ดังต่อไปนี้

- สายอักขระ _ก มาจาก เสียง อา ตามด้วยเสียง โอะ ลจรูป ตามด้วยพยัญชนะ ก เช่น ทารก จารก ชารก บาลก ลาดก วาดก ซาลก เป็นต้น ซึ่งเป็นการใช้เสียง อา ร่วมกับเสียง ออก แต่ก็พบบ้างในภาษาอื่น เช่น คำภาษาอังกฤษ (อลาสกา) คำภาษาญี่ปุ่น (อายาดกะ) และ คำภาษาฝรั่งเศส (ปาสกาล)

- สายอักขระ _นทร์ เช่น วนทร์ พนทร์ เป็นต้น ซึ่งเป็นการใช้สายอักขระที่เป็นพยัญชนะติดกันร่วมกับเครื่องหมายการันต์ สายอักขระนี้แสดงให้เห็นลักษณะเฉพาะของการใช้ตัวการันต์ที่มากกว่าหนึ่งตัวอักษรของคำไทย เช่น ทร เป็นต้น

4.7.2 สายอักขระคำทับศัพท์ภาษาอังกฤษ

ในการวิเคราะห์สายอักขระคำทับศัพท์ภาษาอังกฤษ ซึ่งผู้วิจัยวิเคราะห์ด้วยสายอักขระเฉพาะขนาด 2-5 แกรม ทำให้พบว่าคำทับศัพท์ภาษาอังกฤษ มีการใช้ตัวอักษร อักษรวิธีของภาษาซึ่งในที่นี้คือการใช้เครื่องหมายวรรคตอน และข้อสังเกตของภาษาจากสายอักขระเฉพาะ ซึ่งสิ่งเหล่านี้จะสะท้อนให้เห็นลักษณะของคำทับศัพท์ภาษาอังกฤษ ที่คำภาษาอื่นไม่ปรากฏ มีดังนี้

1. การใช้อักษรที่เป็นตัวสะกด เช่น ตัวสะกด ล ค ส ฟ เป็นตัวสะกดที่พบมากในคำทับศัพท์ภาษาอังกฤษ เนื่องจากในภาษาไทยมีการใช้เสียงพยัญชนะท้ายหรือตัวสะกดเป็นเสียงระเบิดไม่ก้องและเสียงนาสิกเท่านั้น แต่เมื่อถอดอักษรจากอักษรภาษาอังกฤษ l k s และ f จึงปรากฏอักษร ล ค ส ฟ ในพยัญชนะท้ายในคำทับศัพท์ อาจเป็นเพราะผู้ใช้ต้องการออกเสียงด้วยสำเนียงของภาษาอังกฤษ

โดยจากการใช้อักษรเหล่านี้เมื่อปรากฏร่วมกับเสียงสระอื่นๆเช่น เสียง “เอ” และ “อ” จะทำให้เกิดสายอักขระเฉพาะของภาษาอังกฤษ ได้แก่

- สายอักขระ แ_ค เช่น แมค แทค แนค แลค แอค เป็นต้น
- สายอักขระ แ_ส เช่น แจส แพส แวส แอส เป็นต้น แต่สายอักขระนี้ยังพบได้บ้างในภาษาฝรั่งเศส เช่น แตแรส นาร์บงแนส มงแตสกีเออ เป็นต้น
- สายอักขระ _อส เช่น กอส คอส ลอส เป็นต้น แต่สายอักขระนี้ยังพบได้บ้างในภาษาไทยแต่เป็นการมองข้ามพยางค์ (หอสุมุด ซอสสามสาย) และในภาษาฝรั่งเศส (กาซอส ปลูวียอส)
- สายอักขระ _อฟ เช่น คอฟ ซอฟ ทอฟ มอฟ เป็นต้น แต่สายอักขระนี้ยังพบได้บ้างในภาษาฝรั่งเศส เช่น กอฟแฟรง ออฟฟิซิเยส์ เป็นต้น

และยังพบตัวสะกดที่พบมากในภาษาอังกฤษ คือ ต เช่น สายอักขระ เ็ต เช่น เต็ต เน็ต เร็ต เล็ต เป็นต้น แต่ตัวสะกดนี้ก็ยังมีบ้างในภาษาอื่น เช่น คำไทย (ภูเก็ต) ภาษาญี่ปุ่น (เท็ตสึยะ) และในภาษาฝรั่งเศส (กัมเบ็ตตา)

2. การใช้เครื่องหมายวรรคตอน ซึ่งจะต่างกับคำในภาษาอื่นตรงที่ ในภาษาอื่นจะไม่มีการใช้เครื่องหมายวรรคตอนในลักษณะดังต่อไปนี้

- สายอักขระ _น_ เช่น เซนส์ เรนต เรนส์ เป็นต้น สายอักขระนี้เกิดขึ้นเพราะคำเหล่านี้เมื่อถอดอักษรเป็นคำทับศัพท์ภาษาอังกฤษแล้ว จะปรากฏเสียง _น_ ซึ่งในภาษาไทยออกเสียงแค่ เอน เท่านั้น ดังนั้นจึงมีการใช้เครื่องหมายวรรคตอนในพยัญชนะท้ายที่ไม่ออกเสียงในภาษาไทย เช่น ในภาษาอังกฤษ lens จะถูกทับศัพท์เป็น เลนส จึงมีการใส่การันต์บน ส เพื่อตัดเสียง ส ออกไป

- สายอักขระ _อร์ เช่น เซอร์ เคอร์ เจอร์ เบอร์ เปอร์ เลอร์ เป็นต้น สายอักขระนี้เกิดขึ้นเพราะ ในภาษาอังกฤษมีการใช้ er ทำพยยางค์ และเมื่อถอดอักษรทับศัพท์เป็นไทยแล้วจะเป็น เออร์ แต่ในภาษาไทยไม่ออกเสียง ร ดังนั้นจึงต้องใส่เครื่องหมายการันต์ เพื่อตัดเสียง ร ออกไป แต่สายอักขระนี้ยังพบได้บ้างในภาษาฝรั่งเศส เช่น เมอซีเออร์ กูร์มาเยอร์ เป็นต้น

- สายอักขระ ี ุ ู เช่น กิตส์ กิยร์ คิยร์ จิยร์ เป็นต้น คำเหล่านี้เมื่อถอดอักษรเป็นคำทับศัพท์ภาษาอังกฤษแล้วเป็นพยัญชนะตัวสุดท้ายเป็นพยัญชนะที่ไม่ออกเสียง จึงต้องใส่เครื่องหมายการันต์กำกับ โดยสายอักขระนี้พบมากใน ี ุ ู เช่น คำว่า เกียร์ (gear) เบียร์ (bear) เป็นต้น

3. ข้อสังเกตของการสร้างคำทับศัพท์ภาษาอังกฤษ ในการสร้างคำทับศัพท์ภาษาอังกฤษ เป็นการใช้อยู่ร่วมกันทั้งพยัญชนะและสระ ดังนั้น ทำให้พบว่าการใช้สายอักขระบางสายร่วมกันทำให้เกิดสายอักขระ ซึ่งแม้ว่าจะไม่ใช่สายอักขระเฉพาะของคำภาษาอังกฤษ แต่ก็ยังเป็นสายอักขระที่พบมากในคำภาษาอังกฤษในขณะที่ไม่พบในภาษาอื่นหรือพบน้อยกว่ามาก ดังนั้นจึงนำมาแสดงไว้ดังต่อไปนี้

- สายอักขระ เ เช่น คเ ไซเ ด้เ ด้เ เป็นต้น
- สายอักขระ เ เ เช่น ร์สเ ร์คเ ท์ซเ ล์ดเ เป็นต้น
- สายอักขระ เ เ เช่น ูเ ูเ ูเ ูเ ูเ เป็นต้น
- สายอักขระ แ เ เช่น แกเ แครเ แพเ แบเ แพเ แพเ แมเ เป็นต้น แต่ในสายอักขระนี้ยังพบบ้างในภาษาอื่น เช่น คำไทย (กุญแจเ) และภาษาฝรั่งเศส (แดกียเ บูแดอเ)
- สายอักขระ แ เ เช่น แมนู แคนู แวลู เป็นต้น แต่ในสายอักขระนี้ยังพบบ้างในภาษาฝรั่งเศส เช่น อาแซกูร์ เป็นต้น
- สายอักขระ า เ เช่น -าเซ -าเด เป็นต้น แต่ในสายอักขระนี้ยังพบบ้างในภาษาอื่น เช่น ภาษาฝรั่งเศส (กาเปเตียง) และภาษาญี่ปุ่น (คาเคเมะ)
- สายอักขระ า เ เช่น -าเซ -าเท -าเป -าเอ -าเว เป็นต้น แต่ในสายอักขระนี้ยังพบบ้างในคำไทย เช่น ไร่เงินปลาทอง
- สายอักขระ ิ เ เช่น -ิเค -ิเน -ิเร -ิเล -ิเอ เป็นต้น แต่ในสายอักขระนี้ยังพบบ้างในภาษาญี่ปุ่น เช่น นิเซอะ จิตซุริเสเรียวกุ เป็นต้น
- สายอักขระ ิ เ เ เช่น -ิเซน -ิเอน เป็นต้น แต่ในสายอักขระนี้ยังพบบ้างในภาษาญี่ปุ่น เช่น เคอิเอนงกิ โยบิเซนเคียะ เป็นต้น
- สายอักขระ ุ เ เช่น ูโคเ ูโทเ ูโรเ เป็นต้น แต่ในสายอักขระนี้ยังพบบ้างในภาษาญี่ปุ่น เช่น เซนซุโซเดน เป็นต้น

- สายอักขระ ิ ี เช่น ริชิม ลิชิม ริทิม เป็นต้น แต่ในสายอักขระนี้ยังพบบ้างในภาษาญี่ปุ่น เช่น คะมิชิมิยะ โคะอิชิมิ เป็นต้น
- สายอักขระ ื เ เช่น ืเอ ืลิเน ืลิเว ืดิเต เป็นต้น แต่ในสายอักขระนี้ยังพบบ้างในภาษาฝรั่งเศส เช่น ตาปีซีเยร์ เป็นต้น
- สายอักขระ อ เ เช่น ออกเตอ อนเดอ อนเนอ อมเมอ เป็นต้น แต่ในสายอักขระนี้ยังพบบ้างในภาษาอื่น เช่น ภาษาฝรั่งเศส (ซองเตอลู อองเลองส์) และภาษาญี่ปุ่น (เองเคอิ)
- สายอักขระ ั อ เช่น ังเกอ ัตเตอ ันเดอ ันเนอ ัมเปอ ัลเลอ เป็นต้น แต่ในสายอักขระนี้ยังพบบ้างอื่นๆ เช่น ภาษาญี่ปุ่น (ซังเคอิ นันเซอิ) และภาษาฝรั่งเศส (กัทเธอริน ตัสเตอแวง)

4.7.3 สายอักขระคำทับศัพท์ภาษาญี่ปุ่น

ในการวิเคราะห์สายอักขระคำทับศัพท์ภาษาญี่ปุ่น ซึ่งผู้วิจัยวิเคราะห์ด้วยสายอักขระเฉพาะขนาด 2-5 แกรม ทำให้พบว่าคำทับศัพท์ญี่ปุ่น มีการใช้อักษรวิธีของภาษาและข้อสังเกตของภาษาจากสายอักขระเฉพาะ ซึ่งสิ่งเหล่านี้จะสะท้อนให้เห็นลักษณะของคำทับศัพท์ภาษาญี่ปุ่น ดังนี้

1. สายอักขระ ทชี่ และ ทลี สายอักขระนี้สะท้อนให้เห็นการใช้พยัญชนะของภาษาญี่ปุ่นร่วมกับเสียง อี โดยทับศัพท์มาจาก ttsu และ tsu ซึ่งจากเกณฑ์ทับศัพท์ ttsu และ tsu จะถอดเสียงเป็น ตลี แต่อาจนิยมถอดเสียงเพี้ยนเป็น ทลี หรือ ทชี่ ได้

2. ข้อสังเกตของการสร้างคำทับศัพท์ภาษาญี่ปุ่น จากการที่ภาษาญี่ปุ่นมีสระเดียวน้อยกว่าคำไทย คือ มี 10 สระ ได้แก่ อะ อา เอะ เอ อี อี้ โอะ โอ อุ อู ในขณะที่คำไทยมี 18 สระ ได้แก่ อะ อา อี อี้ อี้ อ้อ อุ อู เอะ เอ แอะ แอ โอะ โอ เอาะ ออ เอาะ เออ ดังนั้น ทำให้รูปแบบภาษาญี่ปุ่นอาจมีสระเดียว 10 ตัวนี้ต่อเนื่องกันค่อนข้างมากกว่าภาษาอื่น เป็นเหตุให้พบสายอักขระเหล่านี้ในภาษาญี่ปุ่นค่อนข้างมาก ในขณะที่พบในภาษาอื่นน้อยหรืออาจไม่พบเลยโดยสายอักขระเหล่านี้ได้แก่

- สายอักขระ ะ โ เช่น ะโป ะโร ะโอ เป็นต้น แต่ในสายอักขระนี้ยังพบบ้างในภาษาอังกฤษ เช่น นโปเลียน อะโรเมติกส์ เป็นต้น
- สายอักขระ ู า ิ เช่น ูมารี ูนาซิริ ูราอิ เป็นต้น แต่ในสายอักขระนี้ยังพบบ้างในภาษาอังกฤษ เช่น ดูกาทิส ยูคาลิปตัส เป็นต้น

4.7.4 สายอักขระคำทับศัพท์ภาษาฝรั่งเศส

ในการวิเคราะห์สายอักขระคำทับศัพท์ภาษาฝรั่งเศส ซึ่งผู้วิจัยวิเคราะห์ด้วยสายอักขระ เฉพาะขนาด 2-5 แกรม ทำให้พบว่าในการเขียนทับศัพท์ภาษาฝรั่งเศส มีการใช้อักษรวิธี และ ข้อสังเกตของภาษา ซึ่งสามารถสะท้อนให้เห็นลักษณะของคำทับศัพท์ภาษาฝรั่งเศสได้ แต่ สิ่งที่สะท้อนให้เห็นในภาษานี้มักปรากฏให้เห็นน้อยกว่าภาษาอื่น ๆ เพราะผู้วิจัยใช้คลังข้อมูล คำทับศัพท์ภาษาฝรั่งเศสน้อยกว่าในภาษาอื่น ซึ่งแสดงไว้ดังนี้

1. การใช้สายอักขระที่มีเสียง ออ ซึ่งเสียง ออ (aur, am, an, aon, em) เป็นเสียงที่พบ มากในภาษานี้ แต่ก็พบมากในภาษาอื่นด้วย เช่น ในคำไทย (ขอ ฟอ) คำทับศัพท์ภาษาอังกฤษ (แมนทอล) คำทับศัพท์ภาษาญี่ปุ่น (คอนจิ) แต่ในกรณีที่มีเสียง ออ นี้ ตามหลังด้วยเสียง อู เห็นได้ จาก สายอักขระ _ออ_ เช่น ซอนู ฟอรู ลอปู เป็นต้น แต่สายอักขระนี้ยังพบบ้างในภาษาอังกฤษ เช่น เดอลุกซ์ ออกุสต์ เป็นต้น

2. การใช้สายอักขระที่มีเสียง อง นำหน้า ซึ่งเสียง อง (om) เป็นเสียงที่พบมากใน ภาษานี้แต่พบในภาษาอื่นด้วย เช่น คำไทย (กงจักร) คำทับศัพท์ภาษาญี่ปุ่น (บงโอโตริ) แต่เมื่อ เสียง อง ตามด้วยเสียง อง ทำให้เกิดสายอักขระเฉพาะ _ง_ง เช่นคำว่า องบง ดาลองซง บริยองซง ลงฟง เป็นต้น แต่สายอักขระนี้ยังพบบ้างในภาษาอื่น เช่น คำไทย (หงองแหงง) คำภาษาอังกฤษ (ฮองกง) และภาษาญี่ปุ่น (เตงบง)

3. ข้อสังเกตของการสร้างคำทับศัพท์ฝรั่งเศส จากการสร้างคำทับศัพท์ภาษาฝรั่งเศส มีการใช้สระที่มีเสียง อู ออง แอ และ โอ มากในภาษา จึงทำให้เกิดสายอักขระเฉพาะที่พบมาก ในคำทับศัพท์ภาษาฝรั่งเศส ในขณะที่ภาษาอื่นพบน้อย ได้แก่

- สายอักขระ โอ_อง เช่น โกลอง โตบอง โบมอง โวบอง เป็นต้น
- สายอักขระ - ูร์แ_ เช่น - ูร์แค - ูร์แน เป็นต้น แต่สายอักขระนี้ยังพบบ้างใน ภาษาอังกฤษ เช่น มูร์แลนด์ เป็นต้น

4.8 สรุปผลการใช้วิธีสายอักขระเฉพาะ

จากผลการตัดสินภาษาได้ถูกต้องด้วยสายอักขระเฉพาะในแต่ละภาษาเมื่อมองแต่ภาษา หลักที่ใช้จำนวนคลังข้อมูลการฝึกเท่านั้น พบว่าวิธีการใช้สายอักขระเฉพาะมีผลต่อโปรแกรมระบุ คำประมาณ 50% คือ คำทับศัพท์ภาษาญี่ปุ่น 54.09% คำไทย 50.58% และคำทับศัพท์ ภาษาอังกฤษ 48.71% โดยสายอักขระเฉพาะมีผลต่อการระบุคำทับศัพท์ภาษาญี่ปุ่นมากที่สุด

อาจเป็นเพราะภาษาญี่ปุ่นเป็นภาษาที่พบสายอักขระเฉพาะมากที่สุด และตัดสินภาษาผิดพลาดน้อยที่สุด หรือหากมองในเชิงภาษาศาสตร์ ภาษาญี่ปุ่นเป็นภาษาที่ใช้สายอักขระซ้ำๆ กันมาก ไม่มีการผสมสระมากเท่าภาษาอื่นๆ ส่วนวิธีสายอักขระเฉพาะมีผลต่อการระบุคำไทยมากถึง 50% อาจเป็นเพราะคำไทยถูกตัดสินภาษาด้วยสายอักขระเฉพาะ 1-ตัวถึง 14.14% ซึ่งในภาษาอื่นไม่มีการใช้สายอักขระเฉพาะ 1-ตัวนี้

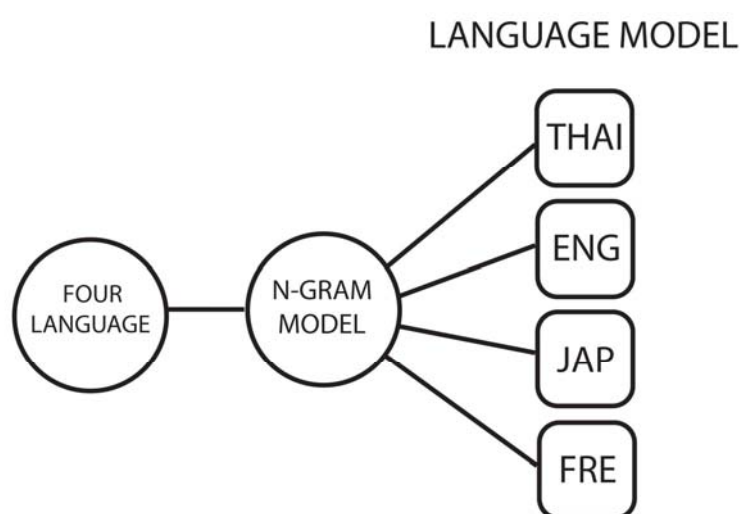
บทที่ 5

การระบุภาษาของคำด้วยแบบจำลองภาษา

ในบทนี้ผู้วิจัยจะกล่าวถึงอีกระบบหนึ่งที่ผู้วิจัยพัฒนาเพื่อใช้ร่วมกับการใช้สายอักขระเฉพาะ คือ การระบุภาษาของคำด้วยแบบจำลองภาษา 2-5 แกรม โดยแสดงถึงภาพรวมของระบบ การระบุภาษาโดยใช้แบบจำลองภาษา ผลการทดสอบระบบ ผลเปรียบเทียบการใช้ขนาดของเอ็นแกรมต่างกัน และข้อผิดพลาดที่เกิดขึ้น

5.1 การสร้างแบบจำลองภาษา (Language model)

การสร้างแบบจำลองภาษาในงานวิจัยนี้เป็นเพียงการเก็บค่าทางสถิติของชุดอักขระที่เกิดขึ้นในคลังข้อมูลด้วยวิธีเอ็นแกรมทั้ง 2-5 แกรม โดยจะใช้งานในช่วงเวลาขณะที่ต้องการจะทดสอบระบบ ผลลัพธ์ที่ได้ประกอบด้วยแบบจำลองภาษาทั้งหมด 4 ภาษา คือแบบจำลองคำไทย คำทับศัพท์ภาษาอังกฤษ ภาษาญี่ปุ่นและภาษาฝรั่งเศส เห็นได้จากภาพที่ 5.1 โดยในแบบจำลองแต่ละภาษาก็แยกเป็น 2-แกรม 3-แกรม 4-แกรม และ 5-แกรม เพื่อทดลองเปรียบเทียบว่าวิธีเอ็นแกรมขนาดใดให้ประสิทธิภาพดีกว่า

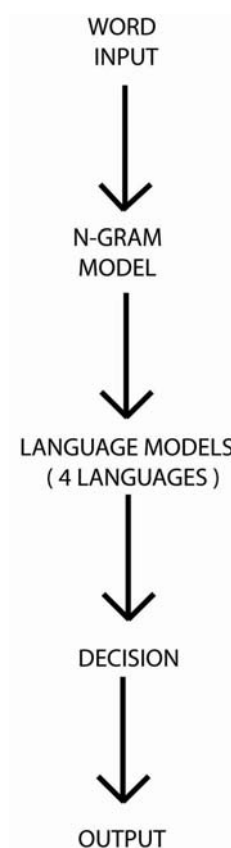


ภาพที่ 5.1 ลักษณะของแบบจำลองภาษาเอ็นแกรม

ต่อไปจะกล่าวถึงโปรแกรมระบุภาษาของคำด้วยแบบจำลองภาษา โดยจะอธิบายโครงสร้างของโปรแกรมระบุภาษาของคำด้วยแบบจำลองภาษาที่พัฒนาและขั้นตอนของโปรแกรมทั้งหมดที่ละขั้นตอน

5.2 โปรแกรมระบุภาษาของคำด้วยแบบจำลองภาษา

โปรแกรมระบุภาษาของคำด้วยแบบจำลองภาษา แบ่งออกเป็น 2 ขั้นตอน ได้แก่
 1. ขั้นตอนรับคำเข้า 2. ขั้นตอนระบุภาษาด้วยแบบจำลองภาษา เห็นได้จากภาพที่ 5.2



ภาพที่ 5.2 ขั้นตอนการทำงานของโปรแกรมระบุภาษาของคำด้วยแบบจำลองภาษา

1. ขั้นตอนรับคำเข้า เป็นขั้นตอนที่รับคำเข้ามาแล้วแยกคำๆ นั้น ออกเป็นสายอักขระ ด้วยเอ็นแกรมทั้ง 2-5 แกรม เช่น คำที่รับเข้ามา คือ คำว่า “อักขระ” เพื่อผ่านระบบเอ็นแกรม จะแยกเป็นแต่ละเอ็นแกรม เห็นได้จากตารางที่ 5.1 ดังนี้

ตารางที่ 5.1 การแยกสายอักขระตามเอ็นแกรม 2-5 แกรม

2-แกรม	3-แกรม	4-แกรม	5-แกรม
อ	อัก	อักข	อักขร
-ก	-กข	-กขร	-กขระ
กข	กขร	กขระ	
ขร	ขระ		
ระ			

2. ขั้นตอนระบุภาษาด้วยแบบจำลองภาษา เป็นขั้นตอนที่นำสายอักขระที่แยกมา คำนวณหาค่าความน่าจะเป็นที่เกิดขึ้นในคลังข้อมูล เช่น เมื่อหาความน่าจะเป็นของคำว่า อักขระ โดยใช้ 4-แกรม เราก็ประมาณค่าโดยนำความน่าจะเป็นของแต่ละสายอักขระมาคูณกันดังนี้ $P(อ) P(- | ่อ) P(ก | ่อ) P(ข | อก) P(ร | -กข) P(ะ | กขร)$ ในขั้นตอนนี้ใช้วิธีเอ็นแกรมตั้งแต่ 2-5 แกรมพร้อมกับวิธีปรับค่าไม่ให้เป็นศูนย์ (smoothing) โดยการเพิ่มค่า 0.00000001 ในกรณีที่ทดลองแล้วไม่พบสายอักขระในข้อมูลการฝึก

ในขั้นตอนระบุภาษาของคำด้วยแบบจำลองภาษา จะทดสอบระบบทีละแกรม คือ คำที่รับเข้ามาต้องตัดสินผ่านภาษาเริ่มจาก 2-แกรม พอตัดสินได้แล้วก็จะนำคำนั้นตัดสินต่อไปใน 3-แกรมอีก ระบบจะทำลักษณะแบบนี้จนถึง 5-แกรม โดยจะทำเช่นนี้กับทุกแบบจำลองภาษา (คำไทย คำทับศัพท์ภาษาอังกฤษ คำทับศัพท์ภาษาญี่ปุ่นและคำทับศัพท์ภาษาฝรั่งเศส)

ตัวอย่างผลการทดลองที่ได้ค่าความน่าจะเป็นจากการใช้แบบจำลอง 4-แกรม เห็นได้จากตารางที่ 5.2 ดังนี้

ตารางที่ 5.2 ตัวอย่างผลการทดลองที่ได้จากแบบจำลองภาษา

ค่าความน่าจะเป็นจาก 4-แกรม				
คำที่ทดสอบ จากแบบจำลอง ภาษา	คำไทย	ภาษาอังกฤษ	ภาษาญี่ปุ่น	ภาษาฝรั่งเศส
สินเชื่อ	<u>0.079547</u>	0.00000008	0.0001125	0.0004389
แจ็กเก็ต	0.0004527	<u>0.23293787</u>	0.0046872	0.035714
ทาเกจิ	0.0045454	0.00011531	<u>0.0571847</u>	0.0296052
ซองตราล	0.1150739	0.0244702	0.000025	<u>0.2577250</u>

หลังจากนั้นที่ได้ค่าความน่าจะเป็นของแบบจำลองภาษาแล้วระบบก็จะตัดสินว่าเป็นคำจากภาษาใดโดยเลือกจากแบบจำลองภาษาที่ให้ค่าความน่าจะเป็นสูงที่สุด ซึ่งจากตัวอย่างนี้ คำว่า “สินเชื่อ” จะถูกตัดสินว่าเป็นคำไทย คำว่า “แจ็กเก็ต” ตัดสินว่าเป็นคำทับศัพท์ภาษาอังกฤษ คำว่า “ทาเกจิ” ตัดสินว่าเป็นคำทับศัพท์ภาษาญี่ปุ่น และคำว่า “ซองตราล” ตัดสินว่าเป็นคำทับศัพท์ภาษาฝรั่งเศส โดยหลังจากการทดลองระบุภาษาด้วยโปรแกรมแล้ว ผู้วิจัยเป็นผู้นำค่านั้นตรวจสอบความถูกต้องเองว่าระบบระบุภาษาของคำได้ถูกต้องหรือไม่ เพื่อเป็นตรวจสอบความถูกต้องและทนทานของระบบ

5.3 สรุปผลการทดลอง

จากการทดสอบโปรแกรมระบุภาษาของคำไทยและคำทับศัพท์แสดงให้เห็นผลของระบบระบุภาษาด้วยแบบจำลอง 2-5 แกรม และเนื่องจากระบบนี้จะทดสอบทีละเอ็นแกรม คือระบบจะเริ่มตัดสินภาษาคำทุกคำที่แบบจำลอง 2- แกรม ก่อน พอระบบทำงานเสร็จ ก็นำคำทุกคำที่ตัดสินภาษาแล้ว ผ่านมายังแบบจำลองภาษา 3-แกรมเพื่อตัดสินภาษาอีกซึ่งจะตัดสินเป็นภาษาที่ต่างกันได้ โดยทำลักษณะแบบนี้จนถึง 5-แกรม ดังนั้น ผลการทดลองที่แสดงในตาราง 5.3 คือ ค่าความถูกต้อง ที่แยกกันในแต่ละขนาดของเอ็นแกรมนั้น ผลการทดสอบที่ได้จึงให้ทีละแกรมแยกจากกัน เห็นได้จาก ตารางที่ 5.3 โดยผลการทดสอบระบบที่ได้คือ ค่าความถูกต้อง ซึ่งหมายถึง ระบบสามารถระบุภาษาของคำได้ถูกต้อง

ตารางที่ 5.3 ผลการตัดสินภาษาของคำได้ถูกต้องด้วยแบบจำลอง 2-5 แกรม
โดยผลการทดลอง คือ จำนวนคำที่ระบุภาษาได้ถูกต้อง ในวงเล็บคือจำนวนเปอร์เซ็นต์

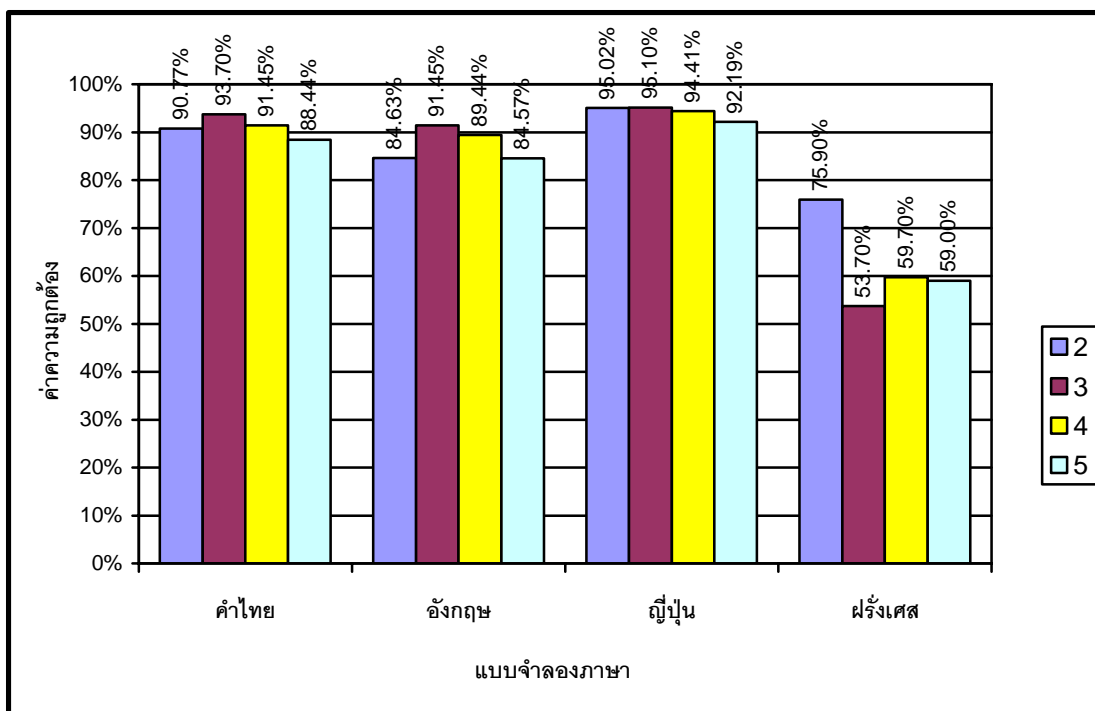
ภาษา วิธี		จำนวนการระบุภาษาของคำได้ถูกต้อง			
		จำนวนภาษาละ 10000 คำ			1000 คำ
		ไทย	อังกฤษ	ญี่ปุ่น	ฝรั่งเศส
แบบจำลองภาษา 2-5 แกรม	2	9077 (90.77%)	8463 (84.63%)	9502 (95.02%)	759 (75.90%)
	3	9370 (93.70%)	9145 (91.45%)	9510 (95.10%)	537 (53.70%)
	4	9145 (91.45%)	8944 (89.44%)	9441 (94.41%)	597 (59.70%)
	5	8844 (88.44%)	8457 (84.57%)	9219 (92.19%)	590 (59%)

จากตารางที่ 5.3 แสดงให้เห็นว่า เมื่อระบุภาษาของคำแต่ละภาษาด้วยแบบจำลองภาษา 3-แกรม ให้ค่าความถูกต้องดีที่สุด ยกเว้นเมื่อระบุภาษาของคำทับศัพท์ภาษาฝรั่งเศส ส่วนการระบุภาษาด้วย 5-แกรมให้ผลแย่งที่สุด แต่ในทางกลับกันเมื่อระบุคำทับศัพท์ภาษาฝรั่งเศสให้ผลดีที่สุด

นอกจากนี้ ผลการทดลองของแต่ละภาษาที่มีลักษณะคล้ายกัน คือ ยิ่งใช้แกรมสูงผลจะต่ำลงอาจเป็นเพราะในแต่ละภาษาไทย อังกฤษ ญี่ปุ่น จะเจอตัวอักษรที่สามได้ต้องเห็นมาก่อนสองตัวอักษร ดังนั้นแบบจำลอง 3-แกรม จึงได้ผลมากที่สุดในทุกภาษา แต่เมื่อเพิ่มเป็น 4-แกรม ผลถูกต้องลดลง อาจเป็นเพราะเมื่อขยายขนาดเอ็นแกรมมากขึ้น ทำให้จำนวนปริมาณเอ็นแกรมที่พบในข้อมูลการฝึกไม่มากพอทำให้ผลไม่ดีนัก ซึ่งยืนยันได้จากภาษาฝรั่งเศสซึ่งได้ผลไม่ดีเลยในทุกแกรมเพราะข้อมูลการฝึกมีน้อยไป ส่วนภาษาญี่ปุ่น เห็นได้ว่า ภาษาญี่ปุ่นในแต่ละแกรมให้ผลสูงกว่าภาษาอื่นทั้งหมด โดยเมื่อใช้ 2 และ 3-แกรม จะสูงกว่าภาษาอื่นมาก อาจเป็นเพราะ ลักษณะของภาษาญี่ปุ่น มีการใช้สระเดี่ยวๆ ติดกันมาก และมีเพียงสระ อะ อิ อุ เอะ โอะ เท่านั้น ทำให้เพียงแค่ใช้แบบจำลองเอ็นแกรมสั้นๆ ก็สามารถระบุภาษาได้ดี และสังเกตว่าผลจากแบบจำลอง 2-4 แกรมในภาษาญี่ปุ่นไม่ต่างกันมาก คือ 95.02% 95.10% และ 94.41% แสดงว่าเห็นตัวอักษรมาก่อน 1 ตัวก็เดาได้ง่ายกว่า ส่วนภาษาอังกฤษให้ผลต่ำกว่าภาษาอื่นทั้งหมด อาจเป็นเพราะลักษณะของภาษาอังกฤษมีการใช้ภาษาที่มีการสร้างคำทั้งคำสั้นและคำยาวๆ ไม่เหมือนกับคำไทยที่มีพื้นฐานของภาษาเป็นคำโดดหรือคำสั้นๆ หรือภาษาญี่ปุ่นที่ใช้แต่คำสั้นๆ

นอกจากนี้เรายังสามารถเปรียบเทียบผลของการใช้แบบจำลองภาษาต่างได้ ดังนี้

แผนภูมิที่ 5.3 การเปรียบเทียบการใช้แบบจำลองของแต่ละภาษา



จากแผนภูมิที่ 5.3 แสดงการเปรียบเทียบการใช้แบบจำลองของแต่ละภาษา สังเกตเห็นว่า ระบบระบุภาษาของคำใช้แบบจำลองภาษาขนาด 3-แกรม ได้ค่าความถูกต้องสูงที่สุดและค่าความถูกต้องเริ่มต่ำลงใน 4 และ 5-แกรม ตามลำดับ ยกเว้นการระบุคำทับศัพท์ภาษาฝรั่งเศสจะให้ค่าความถูกต้องสูงที่สุดเมื่อระบุภาษาด้วย 2-แกรม นอกจากนี้ จากการเปรียบเทียบการใช้ขนาดของเอ็นแกรมต่างๆกัน ในกรณีที่ใช้คลังข้อมูลเท่ากัน คือ พิจารณาเฉพาะคำไทย คำทับศัพท์ภาษาอังกฤษและญี่ปุ่น ก็อาจสรุปได้ว่า ขนาดของเอ็นแกรมมีผลต่อการตรวจสอบลักษณะเฉพาะของคำ (sparse data) โดยเมื่อขนาดของเอ็นแกรมมากกว่า 3 โอกาสที่จะเจอลักษณะเฉพาะของคำที่ทดสอบในคลังข้อมูลการฝึกก็ยิ่งน้อยลงทำให้ระบุภาษาผิดพลาด ดังนั้น ยิ่งขนาดของเอ็นแกรมมากขึ้นผลก็ยิ่งแย่งลง สังเกตได้จากแผนภูมิแสดงการเปรียบเทียบขนาดของเอ็นแกรม เมื่อขนาดเอ็นแกรมเท่ากับ 5 ค่าความถูกต้องจะน้อยลงอย่างเห็นได้ชัดในการระบุทุกภาษา และในกรณีของคำทับศัพท์ภาษาฝรั่งเศสซึ่งมีจำนวนข้อมูลการฝึกน้อยมาก ผลจากเอ็นแกรมตั้งแต่ 3-5แกรมจะต่ำกว่า 2-แกรม อย่างเห็นได้ชัด

5.4 ปัญหาการตัดสินค้าผิดพลาดที่เกิดขึ้นจากแบบจำลองภาษา

จากผลการทดสอบระบบ พบว่ามีการระบุภาษาผิดพลาดจากการใช้แบบจำลองภาษา 2-5 ตัว ซึ่งแสดงไว้ในตาราง 5.4 ดังนี้ โดยแสดงผลที่ตัดสินค้าผิดพลาด เป็นจำนวนค่าและเปอร์เซ็นต์เทียบกับข้อมูลทั้งหมด

ตารางที่ 5.4 ผลการตัดสินค้าของคำผิดพลาดจากการใช้แบบจำลองภาษา 2-5 แกรม

วิธี	คำที่ทดสอบ	ตัดสินค้าเป็น				
		ไทย	อังกฤษ	ญี่ปุ่น	ฝรั่งเศส	รวม
แบบจำลอง 2แกรม	ไทย	0	464 (4.64%)	264 (2.64%)	195 (1.95%)	923 (9.23)
	อังกฤษ	362 (3.62%)	0	239 (2.39%)	936 (9.36%)	1537 (15.37%)
	ญี่ปุ่น	171 (1.71%)	239 (2.39%)	0	88 (0.88%)	498 (4.98%)
	ฝรั่งเศส	33 (3.30%)	196(19.60%)	12 (1.20%)	0	241(24.10%)
รวม		566 (8.63%)	899 (26.63%)	515 (6.23%)	1219 (12.19%)	3199 (53.68%)
แบบจำลอง 3แกรม	ไทย	0	353 (3.53%)	220 (2.20%)	57 (0.57%)	630 (6.30%)
	อังกฤษ	289 (2.89%)	0	218 (2.18%)	348 (3.48%)	855 (8.55%)
	ญี่ปุ่น	158 (1.58%)	294 (2.94%)	0	38 (0.38%)	490 (4.90%)
	ฝรั่งเศส	58 (5.80%)	384(38.40%)	21 (2.1%)	0	463 (46.3%)
รวม		505 (10.27%)	1031 (44.77%)	459 (6.48%)	443 (4.43%)	2438 (66.05%)
แบบจำลอง 4แกรม	ไทย	0	391 (3.91%)	308 (3.08%)	156 (1.56%)	855 (8.55%)
	อังกฤษ	262 (2.62%)	0	340 (3.40%)	454 (4.54%)	1056 (10.56%)
	ญี่ปุ่น	144 (1.44%)	338 (3.38%)	0	77 (0.77%)	559 (5.59%)
	ฝรั่งเศส	46 (4.60%)	315 (31.5%)	42 (4.20%)	0	403(40.30%)
รวม		452 (8.66%)	1044 (38.79%)	690 (10.68%)	687 (6.87%)	2873 (65%)
แบบจำลอง 5แกรม	ไทย	0	473 (4.73%)	409 (4.09%)	274(2.74%)	1156 (11.56%)
	อังกฤษ	348 (3.48%)	0	517 (5.17%)	678 (6.78%)	1543 (15.43%)
	ญี่ปุ่น	218 (2.18%)	432 (4.32%)	0	131 (1.31%)	781 (7.81%)
	ฝรั่งเศส	49 (4.90%)	299(29.90%)	62 (6.20%)	0	410 (41%)
รวม		615 (10.56%)	1204 (38.95%)	988 (15.46%)	1083 (10.83%)	3890 (75.80%)

จากตารางที่ 5.4 เมื่อดูผลรวมการระบุภาษาผิดพลาดจากแบบจำลองภาษาในแต่ละแกรม (แนวตั้ง) ยกเว้นในคำทับศัพท์ภาษาฝรั่งเศสเพราะใช้ข้อมูลการฝึกไม่เท่ากัน พบว่าในการระบุคำทับศัพท์ภาษาอังกฤษด้วยแบบจำลองภาษา 2-5 แกรม ให้ข้อผิดพลาดมากที่สุด คือ 15.37% 8.55% 10.56% และ 15.43% ตามลำดับ และภาษาที่ถูกตัดสินผิดพลาดน้อยที่สุด คือ คำทับศัพท์ภาษาญี่ปุ่น 4.98% 4.90% 5.59% และ 7.81% ตามลำดับ และเมื่อดูผลรวมของภาษาที่ระบุผิดพลาดในแต่ละแกรม (แนวนอน) พบว่าแบบจำลองภาษาในแต่ละแกรมส่วนใหญ่ มักระบุผิดว่าเป็นคำทับศัพท์ภาษาอังกฤษมากที่สุด คือ 26.63% 44.77% 38.79% และ 38.95% ตามลำดับ

นอกจากนี้เมื่อเปรียบเทียบกันระหว่างแต่ละแกรม ยกเว้นคำทับศัพท์ภาษาฝรั่งเศส พบว่า เมื่อระบุภาษาทุกภาษาด้วยแบบจำลอง 5-แกรม ให้ผลผิดพลาดมากที่สุด คือ คำไทย (11.56%) คำทับศัพท์ภาษาอังกฤษ (15.43%) และคำทับศัพท์ภาษาญี่ปุ่น (7.81%) จึงสรุปได้ว่า ยิ่งใช้ขนาดเอ็นแกรมยิ่งมากผลผิดพลาดก็ยิ่งมากขึ้นด้วย ส่วนการระบุคำทับศัพท์ภาษาฝรั่งเศส ด้วย 3-แกรมให้ผลผิดพลาดมากที่สุด คือ 46.3% และเมื่อเปรียบเทียบข้อผิดพลาดกับที่ระบุในภาษาอื่น ๆ พบว่าในคำทับศัพท์ภาษาฝรั่งเศสให้ข้อผิดพลาดมากที่สุด อาจเป็นเพราะคลังข้อมูลการฝึกจำนวนน้อยกว่าภาษาอื่น

และตารางที่ 5.4 ในแต่ละภาษายังแสดงให้เห็นข้อผิดพลาดในการตัดสินภาษาผิดเป็นภาษาอื่น ซึ่งจากตารางจะเห็นว่าในสามภาษาหลัก ไทย อังกฤษ และญี่ปุ่น (ไม่พิจารณาฝรั่งเศส เพราะเป็นความต่างที่มาจากขนาดข้อมูลเป็นสำคัญ) คำไทยและภาษาอังกฤษจะตัดสินผิดเป็นภาษาญี่ปุ่นน้อยที่สุดใน 2-3 แกรม และหากเทียบกับภาษาอื่นแล้ว ภาษาญี่ปุ่นถูกตัดสินภาษาผิดน้อยที่สุด แสดงว่ามีความสับสนระหว่างภาษาไทยและภาษาอังกฤษมากกว่า จึงสรุปได้ว่า สายอักขระ 2-3 ตัวของภาษาญี่ปุ่นมีลักษณะที่ต่างจากภาษาไทยและภาษาอังกฤษ แต่พอเป็น 4-5 แกรม เริ่มไม่ชัดเจนแล้ว อาจเป็นเพราะยิ่งแกรมสูงข้อมูลการฝึกไม่เพียงพอ จึงไม่ได้ภาพความสับสนระหว่างภาษาแบบที่พบในแบบจำลอง 2-5 แกรม และจากสาเหตุข้อมูลการฝึกไม่เพียงพอ จึงเป็นเหตุผลว่าทำไมคำภาษาอังกฤษถูกตัดสินผิดพลาดว่าเป็นคำไทยใน 2-3 แกรม แต่ใน 4-5 แกรม คำอังกฤษถูกตัดสินผิดว่าเป็นคำญี่ปุ่น

ในส่วนของความสับสนระหว่างภาษาไทยและภาษาอังกฤษ คำไทยมักตัดสินภาษาผิดว่าเป็นภาษาอังกฤษใน 2-5 แกรม เช่น คำว่า ภูเก็ต คอแร้ง จามจุรี น้ำพุ และคำภาษาอังกฤษ มักตัดสินผิดว่าเป็นคำไทยใน 2-3 แกรม เช่น คำว่า วิตามิน ละติน แกรนิต เพราะลักษณะของภาษาอาจมีความคล้ายคลึงกัน เนื่องจากการเขียนทับศัพท์ภาษาอังกฤษในปัจจุบัน อาจมีการเขียนคำอังกฤษให้ตรงกับสำเนียงไทยหรือออกเสียงให้ตรงกับระบบเสียงของภาษาไทย เช่น มี

การใช้ตัวอักษร สระ วรรณยุกต์ที่เหมือนกับคำไทย ทำให้คำไทยกับอังกฤษจึงมีการเขียนที่คล้ายกัน ซึ่งจะต่างกับภาษาญี่ปุ่นที่ต่างกันชัดเจน คือมีแต่การใช้คำสั้นๆ เท่านั้น

อย่างไรก็ตาม ภาษาญี่ปุ่นถูกตัดสินผิดพลาดน้อยสุด แต่ที่ตัดสินผิดพลาดส่วนใหญ่ มักจะตัดสินผิดว่าเป็นคำภาษาอังกฤษใน 2-5 แกรม เช่นคำว่า โคะโมง ราเมง เซนเบออิ คานอิชิ อาจเป็นเพราะขนาดของคลังข้อมูลการฝึกไม่เพียงพอต่อการแสดงลักษณะของภาษาญี่ปุ่นให้ชัดเจน

5.5 สรุปการใช้แบบจำลองภาษาเอ็นแกรม

เหตุผลที่ระบบที่ใช้แบบจำลองเอ็นแกรมตัดสินผิดพลาดเป็นเพราะคลังข้อมูลอาจไม่เพียงพอต่อการแสดงลักษณะของภาษา หรือแต่ละภาษาที่มีความคล้ายคลึงกันมากจนอาจต้องใช้คลังข้อมูลมากกว่านี้ อย่างไรก็ตาม จากการทดสอบระบบระบุภาษาของคำด้วยแบบจำลองภาษาเอ็นแกรม ได้ผลการทดสอบมากกว่า 90% จึงสรุปได้ว่า เมื่อเทียบกับระบบสายอักขระเฉพาะ วิธีการใช้แบบจำลองภาษาเอ็นแกรมมีประสิทธิภาพและความทนทานของระบบดีกว่า เพราะระบบสายอักขระเฉพาะได้ผลความถูกต้องน้อยกว่า คือ คำไทย 50.58% คำทับศัพท์ภาษาอังกฤษ 48.71% คำทับศัพท์ภาษาญี่ปุ่น 54.09% และคำทับศัพท์ภาษาฝรั่งเศส 20.40% ในขณะที่ระบบที่ใช้แบบจำลอง 3-แกรม ได้ผลความถูกต้องมากกว่า คือ คำไทย 93.70% คำทับศัพท์ภาษาอังกฤษ 91.45% คำทับศัพท์ภาษาญี่ปุ่น 95.10% และคำทับศัพท์ภาษาฝรั่งเศส 53.70% แต่ระบบเอ็นแกรมนี้ มีข้อเสียคือ ใช้เวลาในการประมาณผลช้ากว่า ผู้วิจัยจึงนำทั้ง 2 ระบบผสมผสานกันเป็นระบบระบุภาษาของคำด้วยการใช้สายอักขระเฉพาะและแบบจำลองภาษาซึ่งจะกล่าวต่อไปในบทที่ 6

บทที่ 6

โปรแกรมระบุภาษาของคำ

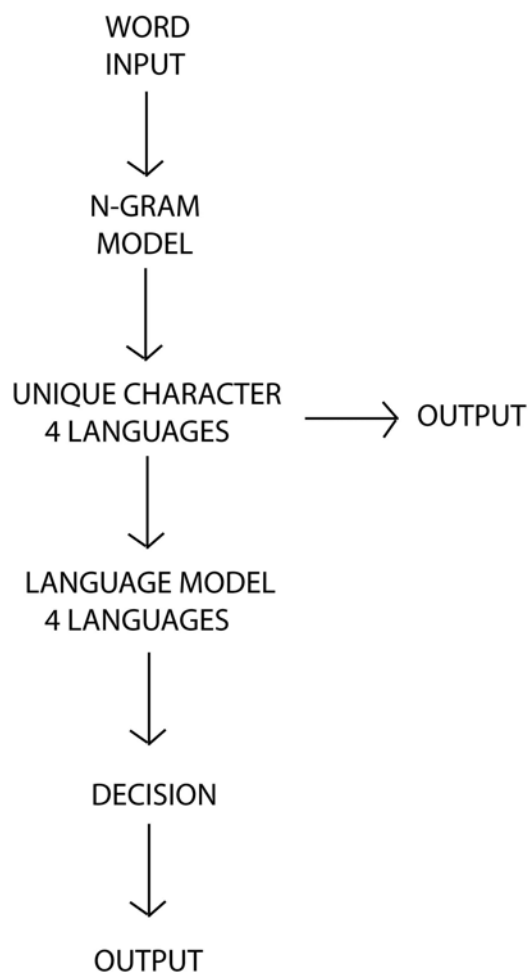
ในบทนี้ผู้วิจัยจะกล่าวถึงระบบที่ผู้วิจัยพัฒนา คือ โปรแกรมระบุภาษาของคำด้วยการใช้สายอักขระเฉพาะ 1-5 ตัวร่วมกับการใช้แบบจำลองภาษาเอ็นแกรม 2-5 แกรม รวมทั้งรายงานและวิเคราะห์ผลการทดสอบระบบที่พัฒนาร่วมกัน โดยจะอธิบายภาพโดยรวมของระบบก่อน เพื่อให้เข้าใจถึงแนวคิดและวิธีการสร้างโปรแกรมระบุภาษาของคำด้วย 2 ระบบนี้

6.1 ภาพโดยรวมของระบบ

โปรแกรมระบุภาษาของคำ คือ โปรแกรมที่ใช้ระบุที่มาภาษาของคำซึ่งพัฒนาโปรแกรมด้วยภาษาเพิร์ล(Perl) ทั้งหมด แนวคิดของการพัฒนาโปรแกรมระบุภาษาของคำนี้ คือ รับข้อมูลคำที่ต้องการตรวจสอบเข้ามาแล้วตรวจหาสายอักขระเฉพาะ โดยนำคำที่ต้องการตรวจสอบเข้ามาแยกเป็นสายอักขระก่อน แล้วจึงนำไปเทียบกับสายอักขระเฉพาะของแต่ละภาษา ถ้าเทียบแล้วตรงกับสายอักขระเฉพาะของภาษาใด ก็จะระบุว่าเป็นคำภาษานั้น แต่ถ้าเทียบแล้วไม่ตรงกัน คือ ระบุคำไม่ได้ว่าเป็นภาษาอะไร จึงนำคำที่ต้องการตรวจสอบนั้นหาความน่าจะเป็นภาษาหนึ่งๆ โดยเทียบกับแบบจำลองภาษาในแต่ละภาษา เหตุผลที่เลือกพัฒนาระบบในลักษณะนี้เพราะ การระบุภาษาด้วยสายอักขระเฉพาะจะทำงานได้เร็วกว่าการใช้แบบจำลองเอ็นแกรมมาก และสามารถระบุภาษาได้ถูกต้องแล้วประมาณ 50% โดยมีข้อผิดพลาดประมาณ 5% (ดูผลการทดลองบทที่ 4) จึงได้กำหนดให้เป็นส่วนแรกของโปรแกรมที่พัฒนา ในกรณีที่ไม่สามารถตัดสินจากสายอักขระเฉพาะได้ว่าเป็นภาษาใด จึงจะใช้แบบจำลองภาษาเป็นตัวตัดสินว่าคำที่รับเข้ามา นั้นใกล้เคียงกับภาษาใดมากที่สุด ซึ่งในการทดลองนี้ก็จะทดลองกับแบบจำลองภาษาตั้งแต่ 2-5 แกรมเพื่อยืนยันอีกครั้งว่าควรใช้แบบจำลองภาษาที่เอ็นแกรมใด และสำหรับแนวคิดและภาพรวมของวิธีการหาสายอักขระเฉพาะ และการสร้างแบบจำลองภาษาของแต่ละภาษา ผู้วิจัยได้นำเสนอและอธิบายไว้แล้วในบทที่ 4 และ 5 ดังนั้นผู้วิจัยจึงต้องการแสดงขั้นตอนของระบบต่อไป

6.2 โปรแกรมระบุภาษาของคำ

โปรแกรมระบุภาษาของคำ แบ่งออกเป็น 3 ขั้นตอน ได้แก่ 1. ขั้นตอนรับคำเข้า 2. ขั้นตอนเทียบกับสายอักขระเฉพาะและ 3. ขั้นตอนระบุภาษาด้วยแบบจำลองภาษา เห็นได้จากภาพที่ 6.1



ภาพที่ 6.1 ขั้นตอนการทำงานของโปรแกรมระบุภาษาของคำ

1. ขั้นตอนรับคำเข้า เป็นขั้นตอนที่รับคำเข้ามาแล้วแยกคำๆ นั้น ออกเป็นสายอักขระด้วยเอ็นแกรมทั้ง 1-5 แกรม เช่น คำที่รับเข้ามา คือ คำว่า “อักขระ” เพื่อผ่านระบบเอ็นแกรมจะแยกเป็นแต่ละเอ็นแกรม

2. ขั้นตอนเทียบกับสายอักขระเฉพาะ หลังจากที่แยกสายอักขระเป็นเอ็นแกรมต่างๆ แล้ว ขั้นตอนต่อไปเป็นการนำสายอักขระที่แยกด้วยเอ็นแกรมแล้ว มาเทียบกับสายอักขระเฉพาะ ทั้ง 4 ภาษา โดยแยกเทียบกันทั้ง 1-5-แกรม ตามแต่ละเอ็นแกรมว่าตรงกันหรือไม่ หากสายอักขระที่นำมาทดสอบตรงกับสายอักขระเฉพาะของภาษาใดภาษาหนึ่งเท่านั้น ก็ถือว่าคำนั้นเป็นภาษาตามสายอักขระเฉพาะนั้นทันที โดยไม่จำเป็นต้องผ่านขั้นตอนเปรียบเทียบกับแบบจำลองภาษา แต่หากไม่มีสายอักขระที่ตรงกับสายอักขระเฉพาะของภาษาใดๆ เลย ก็ให้นำคำนั้นเทียบกับแบบจำลองภาษาทั้ง 4 ภาษาในขั้นตอนต่อไป แต่หากพบว่าคำนั้นมีสายอักขระตรงกับสายอักขระเฉพาะของภาษาต่างๆ มากกว่า 1 ภาษา ทำให้ยังตัดสินไม่ได้ว่าเป็นคำจากภาษาใด ก็จะเปรียบเทียบกับแบบจำลองภาษาเหล่านั้น ตัวอย่างเช่น ในคำว่า “อักขระ” เมื่อทดลองด้วย 4-แกรม อาจะพบ สายอักขระ “กขระ” ตรงกับสายอักขระเฉพาะภาษาญี่ปุ่น และพบสายอักขระ “อักข” ตรงกับสายอักขระเฉพาะคำไทย ทำให้ตัดสินไม่ได้ว่าเป็นคำไทยหรือมาจากภาษาญี่ปุ่น ดังนั้นระบบจะเปรียบเทียบคำนั้นกับแบบจำลองภาษาญี่ปุ่นและแบบจำลองคำไทยเท่านั้น เพื่อตัดสินว่าน่าจะเป็นคำในแบบจำลองใดมากกว่ากัน

ในขั้นตอนการเปรียบเทียบกับสายอักขระเฉพาะนี้ ระบบที่ผู้วิจัยพัฒนาจะทำงานเพียงครั้งเดียว เช่นเดียวกับระบบระบุภาษาด้วยสายอักขระเฉพาะ 1-5 ตัว ในบทที่ 4 และหากว่าระบบนำสายอักขระที่ทดสอบเปรียบเทียบกับสายอักขระเฉพาะจนถึง 5-แกรมแล้วยังไม่พบว่าตรงกัน ระบบก็จะทดลองต่อไปในขั้นตอนระบุภาษาด้วยแบบจำลองภาษา

3. ขั้นตอนระบุภาษาด้วยแบบจำลองภาษา เป็นขั้นตอนสุดท้ายสำหรับระบุภาษา โดยจะใช้ในกรณีที่สายอักขระที่แยกด้วยเอ็นแกรมไม่ตรงกับสายอักขระเฉพาะใดๆเลย หรือสายอักขระที่ทดลองตรงกับสายอักขระเฉพาะมากกว่า 1 ภาษา ในขั้นตอนนี้จะทดลองใช้วิธีเอ็นแกรม 2-5 แกรม เช่นเดียวกับระบบระบุภาษาด้วยแบบจำลองภาษาในบทที่ 5

หลังจากนั้นที่ได้ค่าความน่าจะเป็นของแบบจำลองภาษาของภาษาต่างๆ แล้วระบบก็จะเลือกผลจากภาษาที่มีค่าความน่าจะเป็นสูงที่สุดเป็นคำตอบที่ถูกต้อง หลังจากการทดลองระบุภาษาด้วยโปรแกรมแล้ว ผู้วิจัยเป็นผู้นำคำนั้นตรวจสอบความถูกต้องเองว่าระบบระบุภาษาของคำได้ถูกต้องหรือไม่ เพื่อเป็นตรวจสอบประสิทธิภาพของระบบ

6.3 สรุปผลการทดลอง

จากการทดลองโปรแกรมระบุภาษาของคำทับศัพท์และคำไทย แสดงให้เห็นผลของระบบระบุภาษาด้วยการสังเกตจากสายอักขระเฉพาะ (1-5 สาย) และแบบจำลองภาษาทุกภาษาด้วยเอ็นแกรม (2-5 แกรม) โดยหลังจากที่ผ่านการระบุภาษาของคำด้วยสายอักขระเฉพาะแล้ว ซึ่งเป็นขั้นตอนที่ 2 ของระบบ พบว่าสายอักขระเฉพาะสามารถแก้ไขปัญหาระบุคำได้ประมาณ 50% ยังมีคำเหลืออีกประมาณ 50% ที่ระบบการใช้สายอักขระเฉพาะไม่สามารถตัดสินภาษาได้ จึงส่งมายังขั้นตอนต่อไป คือการระบุภาษาของคำด้วยแบบจำลองภาษา ในขั้นตอนนี้ระบบจะทำงานแยกกันในแต่ละเอ็นแกรม คือ คำที่ยังไม่ได้ตัดสินอีก 50% ก็จะผ่านมายัง แบบจำลองภาษา 2-แกรมทั้งหมด พอระบบทำงานเสร็จ ก็นำคำที่ยังไม่ได้ตัดสินอันเดียวกันนี้ 50% ผ่านมายังแบบจำลองภาษา 3-แกรมอีก โดยทำลักษณะแบบนี้จนถึง 5-แกรม ดังนั้น ผลการทดลองที่แสดงในตารางคือ ค่าความถูกต้องที่แยกกันในแต่ละขนาดของเอ็นแกรม หากต้องการจะดูผลของทั้ง 2 ระบบรวมกัน จะต้องนำผลที่ได้จากสายอักขระเฉพาะในช่อง “รวม” ซึ่งเป็นผลทดสอบที่รวมการใช้สายอักขระเฉพาะทั้งหมด 1-5 ตัว มารวมกับผลจากแบบจำลองภาษาที่ละแกรม เพราะในการทำงานของระบบการใช้สายอักขระเฉพาะจะเป็นการทำงานเพียงครั้งเดียว ดังนั้น หากคิดผลการทดลองของทั้งระบบ จะต้องนำผลจากการใช้สายอักขระทั้ง 1-5 ตัว มารวมกันก่อน

ตัวอย่างการดูผลของระบบโดยรวม เช่น หากเราต้องการวิเคราะห์ผลการระบุภาษาของคำไทย เราก็นำ ผลการทดสอบจากสายอักขระเฉพาะรวมทั้งตั้งแต่ 1-5 ตัว 50.58% + ผลทดสอบจากแบบจำลองภาษา 2-แกรม 39.83% คำตอบที่ได้ คือ 90.41% (ในช่อง “รวมถูก”) ซึ่งเป็นผลโดยรวมของระบบระบุหลังจากผ่านสายอักขระเฉพาะและแบบจำลองภาษา 2-แกรม หรือหากต้องการรู้ผลโดยรวมเมื่อใช้แบบจำลองภาษา 3-แกรม เราก็นำผลจากสายอักขระเฉพาะรวมทั้งตั้งแต่ 1-5 ตัว 50.58% + ผลทดสอบจากแบบจำลอง 3-แกรม 41.33% ผลที่ได้คือ 91.91% ซึ่งผลโดยรวมทั้งหมดแต่ละเอ็นแกรม ผู้วิจัยได้รวมไว้ในช่อง “รวมถูก” และแยกแต่ละภาษาและเอ็นแกรมไว้ เห็นได้จากตาราง 6.1

ตารางที่ 6.1 ผลการตัดสินภาษาของคำได้ถูกต้องด้วยโปรแกรมระบุคำ (สายอักขระเฉพาะ 1-5 ตัว และแบบจำลอง 2-5 แกรม) โดยผลการทดลอง คือ จำนวนคำที่ระบุภาษาได้ถูกต้อง (เปอร์เซ็นต์)

ภาษา วิธี		ผลการระบุภาษาของคำได้ถูกต้อง											
		จำนวนคำที่ทดสอบ 10,000 คำ									1,000 คำ		
		คำไทย			ภาษาอังกฤษ			ภาษาญี่ปุ่น			ภาษาฝรั่งเศส		
		ถูก	ผิด	ถูก	ผิด	ถูก	ผิด	ถูก	ผิด	ถูก	ผิด		
ผลจากสายอักขระเฉพาะ	1	1414 (14.14%)	0	0	2	0	0	0	0	0	0		
	2	756 (7.56%)	2	122 (1.22%)	14	160 (1.6%)	7	0	3				
	3	2104 (21.04%)	274	2018 (20.18%)	264	3392 (33.92%)	186	57 (5.7%)	122				
	4	651 (6.51%)	164	1769 (17.69%)	277	1683 (16.83%)	168	101 (10.1%)	99				
	5	133 (1.33%)	23	962 (9.62%)	81	174 (1.74%)	38	46 (4.60%)	34				
	รวม		5058 (50.58%)	463	4871 (48.71%)	638	5409 (54.09%)	399	204 (20.40%)	258			
ผลจากแบบจำลองภาษา	2	3983	9041 (90.41%)	496	3728	8599 (85.99%)	763	3903	9312 (93.12%)	289	410	614 (61.40%)	128
	3	4133	9191 (91.91%)	346	4036	8907 (89.07%)	455	3898	9307 (93.07%)	294	286	490 (49%)	252
	4	4012	9070 (90.7%)	467	3876	8747 (87.47%)	615	3833	9242 (92.42%)	359	319	523 (52.30%)	219
	5	3849	8907 (89.07%)	630	3569	8440 (84.4%)	922	3720	9129 (91.29%)	472	333	537 (53.70%)	205
	รวม												

จากผลการทดลองเมื่อเปรียบเทียบกันในแต่ละภาษา เห็นว่า ระบบระบุภาษาของคำภาษาญี่ปุ่นให้ผลดีที่สุดในทุกๆ แกรม คือ 93.12% 93.07% 92.42% และ 91.29% และระบบระบุภาษาของคำภาษาอังกฤษให้ผลแย่ที่สุด ในทุกๆ แกรม คือ 85.09% 89.07% 87.47% และ 84.4% และเมื่อพิจารณาในแต่ละภาษา จะพบว่าในภาษาไทย อังกฤษและญี่ปุ่น ระบบที่เลือกใช้แบบจำลองภาษา 3-แกรมจะให้ผลดีที่สุด

นอกจากนี้จากผลการทดลองซึ่งเป็นระบบที่ประยุกต์นำวิธีสายอักขระเฉพาะและแบบจำลองภาษาเข้าด้วยกัน เมื่อนำไปเปรียบเทียบกับระบบการใช้สายอักขระเฉพาะ (บทที่4) และระบบการใช้แบบจำลองภาษา (บทที่5) ไว้ในตารางที่ 6.2

ตารางที่ 6.2 ผลการเปรียบเทียบค่าความถูกต้องของระบบสายอักขระเฉพาะ ระบบแบบจำลองภาษา และระบบที่ประยุกต์ทั้ง 2 ระบบ โดยระบบสายอักขระเฉพาะ คือ UN1-5 ระบบแบบจำลองภาษา คือ LM 2-5 และ ระบบที่ประยุกต์ทั้ง 2 ระบบ คือ UN-LM

ภาษา	จำนวน 10000 คำ									จำนวน 1000 คำ		
	คำไทย			ภาษาอังกฤษ			ภาษาญี่ปุ่น			ภาษาฝรั่งเศส		
ระบบ	UN1-5	LM 2-5	UN+LM	UN1-5	LM 2-5	UN+LM	UN1-5	LM 2-5	UN+LM	UN1-5	LM 2-5	UN+LM
ผลการทดลอง	50.5	90.77%	90.41%	48.7	84.63%	85.99%	54.0	95.02%	93.12%	20.4	75.5%	61.40%
		93.70%	91.91%		91.45%	89.07%		95.10%	93.07%		53.7%	49%
	8%	91.45%	90.70%	1%	89.44%	87.47%	9%	94.41%	92.42%	0%	59.7%	52.30%
		88.44%	89.07%		84.57%	84.40%		92.19%	91.29%		59%	53.70%

จากการเปรียบเทียบพบว่า ระบบที่ประยุกต์ทั้ง 2 ระบบรวมกัน (UN+LM) ดีกว่าระบบการใช้สายอักขระเฉพาะ (UN) อย่างเห็นได้ชัด เพราะ ผลจากระบบสายอักขระเฉพาะประมาณ 50% แต่ ผลจากระบบที่ประยุกต์ทั้ง 2 ระบบ (UN+LM) ให้ผลมากกว่า 90%

ส่วนกรณีเมื่อนำระบบที่ประยุกต์ทั้ง 2 ระบบรวมกัน ไปเปรียบเทียบกับระบบที่ใช้แบบจำลองภาษาเอ็นแกรม (LM) พบว่าระบบที่ใช้แบบจำลองเอ็นแกรมได้ผลดีกว่าในทุกๆ แกรม ยกเว้นแค่ในการระบุคำทับศัพท์ภาษาอังกฤษด้วย 2-แกรม ผลจากระบบประยุกต์ดีกว่าเล็กน้อย

ดังนั้น ผู้วิจัยจึงเสนอ 2 แนวทางในการเลือกใช้ระบบคือ อาจกลับไปใช้เฉพาะวิธีเอ็นแกรมจะดีกว่าระบบที่ประยุกต์ทั้ง 2 ระบบ และปรับปรุงระบบเอ็นแกรมให้ดีขึ้น เพื่อให้ผลการทดลองใกล้เคียง 100% มากที่สุด อาจทำได้โดยการเพิ่มคลังข้อมูลการฝึกให้มากขึ้นหรือใช้ข้อมูลการฝึกที่มีการกระจายข้อมูลมากขึ้น แต่อาจมีข้อเสียคือ การประมาณผลของระบบช้ามากขึ้น

ส่วนทางเลือกที่สอง คือ อาจใช้ระบบประยุกต์ทั้ง 2 ระบบเหมือนเดิม เพียงแค่ปรับปรุงระบบให้ดีขึ้น ในการปรับปรุงส่วนของระบบสายอักขระเฉพาะ คือ สร้างสายอักขระเฉพาะให้ครอบคลุมทั้งภาษา โดยการเพิ่มคลังข้อมูลการฝึกให้มากขึ้น เพื่อให้ผลผิดพลาดที่เกิดขึ้นน้อยที่สุด ส่วนการปรับปรุงระบบแบบจำลองภาษา คือ อาจเพิ่มขนาดของข้อมูลการฝึก ซึ่งจะกล่าวไว้ในส่วนของข้อเสนอแนะ

บทที่ 7

สรุปผลการวิจัย อภิปรายผลและข้อเสนอแนะ

วิทยานิพนธ์นี้ได้นำเสนองานเกี่ยวกับการระบุภาษาของคำไทย คำทับศัพท์ภาษาอังกฤษ คำทับศัพท์ภาษาญี่ปุ่นและคำทับศัพท์ภาษาฝรั่งเศส โดยสรุปผลการวิจัยทั้งหมดไว้ดังนี้

7.1 สรุปผลการวิจัย

เนื่องจากในปัจจุบันมีการใช้คำทับศัพท์ภาษาต่างประเทศปะปนกับคำไทยจำนวนมาก หากเราต้องการพัฒนาโปรแกรมทางด้านการค้นหาสารสนเทศข้ามภาษา (Cross - Language Information Retrieval) เช่น โปรแกรมถอดอักษรคำทับศัพท์กลับจากภาษาไทยเป็นภาษาต้นฉบับ (backward transliteration) โดยก่อนที่จะถอดอักษรกลับเป็นภาษาต้นฉบับ เราจำเป็นต้องรู้ว่าคำเดิมนั้นเป็นคำในภาษาอะไร จึงต้องมีการระบุภาษาของคำในข้อความ ก่อนที่จะนำไปใช้กฎการถอดอักษร (transliteration) ในแต่ละภาษา นอกจากนี้ ยังพบว่าเคยมีผู้วิจัยทำงานในลักษณะแบบนี้ เช่น การใช้แบบจำลองภาษาเอ็นแกรม และการใช้สายอักขระเฉพาะ โดยผู้วิจัยเห็นว่าแบบจำลองภาษาเอ็นแกรมให้ผลถูกต้องสูงและเรียบง่ายมากที่สุด และวิธีการใช้สายอักขระเฉพาะสามารถประยุกต์ใช้กับงานแบบนี้ได้

ดังนั้น วัตถุประสงค์ที่ต้องการจะทำ คือ ต้องการหาสายอักขระเฉพาะสำหรับใช้ในการระบุภาษาของคำโดยใช้คลังข้อมูลคำไทย คำทับศัพท์ภาษาอังกฤษ ญี่ปุ่นและภาษาฝรั่งเศส และพัฒนาระบบการระบุภาษาของคำไทยและคำทับศัพท์ภาษาต่างประเทศโดยใช้แบบจำลองเอ็นแกรมขนาด 1-5 แกรม เพื่อทดสอบประสิทธิภาพของเอ็นแกรมที่ต่างกัน

ในการดำเนินการวิจัย ผู้วิจัยเริ่มจากเก็บคลังข้อมูลภาษาต่างๆก่อน ในคำไทย คำทับศัพท์ภาษาอังกฤษ และภาษาญี่ปุ่น เก็บภาษาละ 10,000 คำ ส่วนคำทับศัพท์ภาษาฝรั่งเศสเก็บเพียง 1,000 คำ โดยผู้วิจัยได้แบ่งส่วนของคลังข้อมูลสำหรับทำข้อมูลการฝึกภาษาละ 80% ไว้สำหรับการสร้างสายอักขระเฉพาะ 1-5 สายในแต่ละภาษา และสำหรับสร้างแบบจำลองภาษาแต่ละภาษา ส่วนอีก 20% แบ่งเป็นคลังข้อมูลทดสอบ เหตุที่แบ่งส่วนข้อมูลเพราะต้องการให้ได้ผลที่ไม่บังเอิญขึ้นกับข้อมูลที่ใช้ในการฝึกและทดสอบ

ในขั้นตอนของการทดลองระบบระบุภาษาของคำ คือ รับข้อมูล (คำ) เข้ามาแล้วตรวจหาสายอักขระเฉพาะในแต่ละภาษาถ้าตรงกันกับสายอักขระเฉพาะใดก็ตัดสินว่าเป็นภาษานั้น แต่

ถ้ายังตัดสินภาษาไม่ได้ จึงหาความน่าจะเป็นของภาษาหนึ่งๆ โดยเทียบกับแบบจำลองภาษาต่างๆ และตัดสินผลจากแบบจำลองภาษาที่ให้ค่าความน่าจะเป็นที่สูงที่สุด

หลังจากการทดลองผู้วิจัยเป็นผู้นำคำที่ถูกตัดสินภาษานั้นมาตรวจสอบความถูกต้องเองว่าระบบระบุภาษาของคำได้ถูกต้องหรือไม่ เพื่อเป็นการตรวจสอบความทนทานของระบบ

7.2 อภิปรายผลการวิจัย

ผลการทดลองของ แบ่งออกเป็น 3 ระบบ คือ การระบุภาษาของคำด้วยสายอักขระเฉพาะ 1-5 ตัว การระบุภาษาของคำด้วยแบบจำลองภาษา 2-5 แกรม และการระบุภาษาของคำด้วยสายอักขระเฉพาะร่วมกับแบบจำลองภาษา ดังนี้

7.2.1 การทดลองระบุภาษาของคำด้วยสายอักขระเฉพาะ 1-5 ตัว

สรุปได้ว่าทุกภาษาสามารถพบสายอักขระเฉพาะได้ ซึ่งสายอักขระเฉพาะ 1-ตัว พบเฉพาะในคำไทยเท่านั้น และในการทดลอง สามารถนำสายอักขระเฉพาะเหล่านี้ใช้ในการระบุภาษาได้ เห็นได้จาก ในการระบุภาษาของคำไทย คำทับศัพท์ภาษาอังกฤษ และญี่ปุ่น ด้วยสายอักขระเฉพาะ 3-ตัว ให้ผลการระบุภาษาที่ถูกต้องสูงที่สุด คือ คำไทย 21.04% คำทับศัพท์ภาษาอังกฤษ 20.18% และคำทับศัพท์ภาษาญี่ปุ่น 33.92% แต่ในการระบุภาษาของคำทับศัพท์ภาษาฝรั่งเศส ด้วยสายอักขระเฉพาะ 4-ตัว ให้ผลดีที่สุด คือ 10.10% ส่วนการระบุภาษาด้วยสายอักขระเฉพาะ 1-ตัว สามารถใช้ระบุภาษาได้ถูกต้องได้เฉพาะในคำไทยเท่านั้น

และเมื่อดูผลรวมของการระบุภาษาของคำที่ถูกต้องจากการใช้สายอักขระเฉพาะทั้ง 1-5 ตัว แสดงให้เห็นว่า วิธีการใช้สายอักขระเฉพาะมีผลต่อการระบุภาษาของคำประมาณ 50% โดยสายอักขระเฉพาะมีผลต่อการระบุคำทับศัพท์ภาษาญี่ปุ่นมากที่สุด คือ 54.09% รองลงมาคือ คำไทย 50.58% คำทับศัพท์ภาษาอังกฤษ 48.71% คำทับศัพท์ภาษาฝรั่งเศส 20.40% อาจเป็นเพราะภาษาญี่ปุ่นเป็นภาษาที่พบสายอักขระเฉพาะมากที่สุด หรือ เป็นภาษาใช้สายอักขระซ้ำๆ กันมาก ไม่มีการผสมสระมากเท่าภาษาอื่นๆ

นอกจากนี้ ยังสังเกตได้ว่า ขนาดของคลังข้อมูลมีผลต่อการใช้สายอักขระเฉพาะด้วย เห็นได้จากการระบุภาษาในคำทับศัพท์ภาษาฝรั่งเศส เมื่อใช้สายอักขระเฉพาะ 1-5 ตัว ให้ผลถูกต้องแค่ 20.40% และผลผิดพลาดถึง 25.60% ซึ่งมากกว่าภาษาอื่นมาก อาจเป็นเพราะ หาก

ใช้จำนวนคลังข้อมูลการฝึกน้อย สายอักขระเฉพาะที่ได้ก็จะน้อย ทำให้สายอักขระเฉพาะที่ได้ไม่ครอบคลุมทั้งภาษา

7.2.2 การทดลองระบุภาษาของคำด้วยแบบจำลองภาษา

เมื่อระบุภาษาของคำแต่ละภาษาด้วยแบบจำลองภาษา 3-แกรม ให้ค่าความถูกต้องดีที่สุด ยกเว้นเมื่อระบุภาษาของคำทับศัพท์ภาษาฝรั่งเศส ส่วนการระบุภาษาด้วย 5-แกรมให้ผลแย่มากที่สุด แต่ในทางกลับกันเมื่อระบุคำทับศัพท์ภาษาฝรั่งเศสให้ผลดีที่สุด เหตุผลที่ค่าความถูกต้องแตกต่างกัน อาจเป็นเพราะลักษณะของขนาดข้อมูลการฝึก เห็นได้จากการทดลองระบุคำทับศัพท์ภาษาฝรั่งเศสให้ผลการทดลองน้อยที่สุดเมื่อเทียบกับภาษาอื่น คือ สูงสุดแค่ 75.90% ใน 2-แกรม จึงสรุปได้ว่าขนาดของคลังข้อมูลมีผลต่อระบบเอ็นแกรม

และจากการเปรียบเทียบผลการทดลองจากการใช้แบบจำลองของแต่ละภาษา สังเกตเห็นว่า ระบบระบุภาษาของคำใช้แบบจำลองภาษาขนาด 3-แกรม ได้ค่าความถูกต้องสูงที่สุดและค่าความถูกต้องเริ่มต่ำลงใน 4 และ 5-แกรม ตามลำดับ ยกเว้นการระบุคำทับศัพท์ภาษาฝรั่งเศสจะให้ค่าความถูกต้องสูงที่สุดเมื่อระบุภาษาด้วย 2-แกรม

จากการเปรียบเทียบการใช้ขนาดของเอ็นแกรมต่างกัน ในกรณีที่ใช้คลังข้อมูลเท่ากัน คือ พิจารณาเฉพาะคำไทย คำทับศัพท์ภาษาอังกฤษและญี่ปุ่น จึงสรุปได้ว่า ขนาดของเอ็นแกรมมีผลต่อการตรวจสอบลักษณะเฉพาะของคำ โดยเมื่อขนาดของเอ็นแกรมมากกว่า 3 โอกาสที่จะเจอลักษณะเฉพาะของคำที่ทดสอบในคลังข้อมูลการฝึกก็ยิ่งน้อยลงทำให้ระบุภาษาผิดพลาด ดังนั้น ยิ่งขนาดของเอ็นแกรมมากขึ้นผลก็ยิ่งแย่ เห็นได้จากผลการทดลอง เมื่อใช้ขนาดเอ็นแกรมเท่ากับ 5 ค่าความถูกต้องจะน้อยลงอย่างเห็นได้ชัด และในกรณีของคำทับศัพท์ภาษาฝรั่งเศสซึ่งมีจำนวนข้อมูลการฝึกน้อยมาก ผลจากเอ็นแกรมตั้งแต่ 3-5 แกรมจะต่ำกว่า 2-แกรมอย่างเห็นได้ชัด

นอกจากนี้ เมื่อนำระบบเอ็นแกรมไปเทียบกับระบบสายอักขระเฉพาะ พบว่าวิธีการใช้แบบจำลองภาษาเอ็นแกรมมีประสิทธิภาพและความทนทานของระบบดีกว่า เพราะระบบเอ็นแกรมให้ผลความถูกต้องประมาณ 90% แต่ระบบสายอักขระเฉพาะได้ผลความถูกต้องประมาณ 50% คือ คำไทย 50.58% คำทับศัพท์ภาษาอังกฤษ 48.71% คำทับศัพท์ภาษาญี่ปุ่น 54.09% และคำทับศัพท์ภาษาฝรั่งเศส 20.40% แต่ระบบเอ็นแกรมนี้ มีข้อเสียคือ ใช้เวลาในการประมาณผลช้ากว่า

7.2.3 การทดลองระบบระบุภาษาของคำด้วยสายอักขระเฉพาะร่วมกับแบบจำลองภาษา

จากการทดลองทำให้สรุปได้ว่า ระบบที่พัฒนาเป็นไปตามเป้าหมาย คือให้ผลค่าความถูกต้องมากกว่า 90% โดยเฉพาะใน การระบุภาษาของคำไทยด้วย 3-แกรม ได้ถูกต้อง 91.91% และในการระบุคำทับศัพท์ภาษาญี่ปุ่นด้วย 3-แกรม ได้ถูกต้อง 93.07% แต่เมื่อเทียบระบบนี้กับระบบที่ใช้สายอักขระเฉพาะ พบว่า ผลจากระบบนี้ดีกว่าระบบสายอักขระเฉพาะมาก เพราะให้ผลการทดลองประมาณ 90% แต่ระบบสายอักขระเฉพาะให้ผลแค่ประมาณ 50% และเมื่อนำระบบนี้ไปเทียบกับระบบเอ็นแกรม พบว่า ผลจากระบบนี้ให้ผลแยกว่าระบบเอ็นแกรมในทุกๆ แกรม ยกเว้น เมื่อระบุคำทับศัพท์ภาษาอังกฤษใน 2-แกรม ให้ผลดีกว่า

7.3 ข้อเสนอแนะ

หลังจากการวิเคราะห์ข้อผิดพลาดของระบบแล้ว ข้อผิดพลาดเกิดจากทั้งการใช้สายอักขระเฉพาะและแบบจำลองภาษา แม้ว่าข้อผิดพลาดจะเกิดขึ้นน้อยเมื่อเทียบกับคำที่ระบุได้ถูกต้องทั้งหมดในระบบ แต่ผู้วิจัยต้องการให้ระบบมีประสิทธิภาพมากที่สุด ผู้วิจัยจึงเสนอแนวทางในการพัฒนาและแก้ไขปัญหานี้ที่ระบบ ดังนี้

7.3.1 ระบบสายอักขระเฉพาะ

จากผลการทดลอง เป็นที่น่าสังเกตว่าระบบสายอักขระเฉพาะให้ผลถูกต้องแค่ประมาณ 50% และระบบสายอักขระเฉพาะให้ผลผิดพลาดบ้าง แม้จะเป็นส่วนน้อย ผู้วิจัยจึงเสนอแนวทางแก้ไข โดยการเพิ่มขนาดของคลังข้อมูลการฝึกให้มากขึ้น เพื่อใช้สร้างสายอักขระเฉพาะให้ครอบคลุมทั่วทั้งภาษา อย่างไรก็ตาม ในการเพิ่มขนาดคลังข้อมูลเพื่อหาสายอักขระเฉพาะ ผู้วิจัยก็ยังไม่ได้ศึกษาว่า ขนาดของคลังข้อมูลควรอยู่ในระดับเท่าใด ถึงจะเหมาะสมที่สุดสำหรับใช้ในการหาสายอักขระเฉพาะ จึงเป็นแนวทางที่นำไปพัฒนาระบบต่อไป

และเนื่องจากรูปแบบของภาษา (pattern) ที่กล่าวไว้ในระบบ ผู้วิจัยเป็นผู้วิเคราะห์หารูปแบบเอง ซึ่งอาจจะไม่ครอบคลุมทั่วทั้งภาษา จึงเสนอแนวทางที่จะเพิ่มรูปแบบภาษาให้มากขึ้น หรือหารูปแบบภาษาที่พบบ่อยๆ ในสายอักขระได้ โดยเลือกแทนที่ตัวอักษรบางตัวด้วย

วิธี dummy แล้วตรวจสอบดูว่าตรงกับข้อมูลสายอักขระเฉพาะมากน้อยแค่ไหน เช่น ทำให้เราพบว่าในภาษาอังกฤษมีรูปแบบภาษาที่พบบ่อย คือ _s ดังที่เห็นในตัวอย่าง

7.3.2 ระบบเอ็นแกรม

จากผลการทดลอง เป็นที่น่าสังเกตว่าหากต้องการให้ระบบสามารถระบุภาษาของคำไทยและคำทับศัพท์ภาษาต่างประเทศให้ได้ผลถูกต้องมากขึ้น ก็เป็นสิ่งที่เป็นไปได้เพียงแค่เพิ่มเติมในส่วนของคุณสมบัติการฝึกให้มากขึ้นเพราะจากผลการทดลอง การใช้ข้อมูลการฝึกมากระบบก็สามารถระบุภาษาได้ถูกต้องมากกว่าข้อมูลการฝึกจำนวนน้อย เห็นได้จากเมื่อเทียบกับกรณีการระบุภาษาของคำทับศัพท์ภาษาฝรั่งเศส แต่การเพิ่มขนาดของข้อมูลการฝึกเมื่อเพิ่มถึงจุดหนึ่งแล้ว อาจจะไม่ทำให้ผลลัพธ์ดีขึ้น ดังนั้นผู้วิจัยจึงเสนอให้ทดลองเพิ่มเติม โดยทดลองใช้คลังข้อมูลหลายๆ ขนาด และเพิ่มขนาดคลังข้อมูลการฝึกมากขึ้น เพื่อหาขนาดของคลังข้อมูลที่เหมาะสมที่สุด

และถึงแม้เป็นที่แน่ชัดแล้วว่าขนาดของเอ็นแกรม 3-แกรม ให้ผลการทดลองดีที่สุด แต่ผู้วิจัยหวังว่าจะสามารถพัฒนาระบบได้อีก โดยใช้ขนาดของเอ็นแกรมมากขึ้นเพื่อดูว่าการเปลี่ยนแปลงผลการทดลอง หากใช้ขนาดของเอ็นแกรมมากขึ้น อาจเป็นไปได้ว่า หากใช้คลังข้อมูลการฝึกมากขึ้นและขนาดเอ็นแกรมสูงๆ ผลการทดลองก็อาจจะดีขึ้นได้

7.3.3 ระบบประยุกต์ที่ใช้ระบบสายอักขระเฉพาะร่วมกับแบบจำลองเอ็นแกรม

เนื่องจากผลจากระบบประยุกต์นี้ให้ผลแยกจากระบบที่ใช้แบบจำลองภาษาเอ็นแกรม จึงเสนอว่า หลังจากที่แก้ปัญหาหาระบบที่ใช้สายอักขระเฉพาะและระบบเอ็นแกรม ให้ผลการทดลองถูกต้องมากขึ้นแล้ว จึงน่าจะนำระบบประยุกต์นี้มาประมวลผลอีกครั้งเพื่อดูว่า ผลของระบบประยุกต์นี้ดีกว่าระบบเอ็นแกรมหรือไม่ และให้ทดลองประสิทธิภาพความเร็วในการทำงานเพื่อหาว่า ระบบแบบเอ็นแกรมอย่างเดียวกับระบบแบบผสมทำงานด้วยความเร็วต่างกันอย่างไรมีนัยยะสำคัญหรือไม่

นอกจากนี้อาจนำระบบที่พัฒนาไปใช้ระบุภาษาของคำทับศัพท์ภาษาต่างประเทศอื่นๆ ได้อีก เพียงแค่ใช้คลังข้อมูลการฝึกของภาษานั้นเท่านั้น

รายการอ้างอิง

ภาษาไทย

- คณะอนุกรรมการปรับปรุงศัพท์เทคนิคทางวิศวกรรมไฟฟ้า. 2541. ศัพท์เทคนิควิศวกรรมไฟฟ้า
สื่อสาร. พิมพ์ครั้งที่ 4. กรุงเทพฯ: สมาคมวิศวกรรมสถานแห่งประเทศไทยในพระบรม
ราชูปถัมภ์,
- ราชบัณฑิตยสถาน. 2538. พจนานุกรม ฉบับราชบัณฑิตยสถาน พ.ศ. 2525. พิมพ์ครั้งที่ 5.
กรุงเทพฯ: ราชบัณฑิตยสถาน.
- ราชบัณฑิตยสถาน. 2535. หลักเกณฑ์การทับศัพท์ ภาษาฝรั่งเศส ภาษาเยอรมัน ภาษาอิตาลี
ภาษาสเปน ภาษารัสเซีย ภาษาญี่ปุ่น ภาษาอาหรับ ภาษามลายู. กรุงเทพฯ:
ราชบัณฑิตยสถาน.
- ราชบัณฑิตยสถาน. 2538. หลักเกณฑ์การทับศัพท์ภาษาอังกฤษ. พิมพ์ครั้งที่ 6. กรุงเทพฯ:
ราชบัณฑิตยสถาน.
- ศรีไพร ศักดิ์รุ่งพงศากุล. 2544. เทคโนโลยีคอมพิวเตอร์และสารสนเทศ. กรุงเทพฯ: ซีเอ็ดดูเคชั่น.
- สมชาย ชัยชนะตระกูล. 2544. รวมศัพท์ญี่ปุ่น. พิมพ์ครั้งที่ 2. กรุงเทพฯ: สำนักพิมพ์เพจพับ -
บลิชซิง.

ภาษาอังกฤษ

- Aizawa, A. 2003. Linguistic techniques to Improve the performance of automatic text categorization. In Proceedings of the Sixth Natural Language Processing Pacific Rim Symposium (NLPR2001), pp.307-314. Cited in Peng, Schuurmans, Wang and Huang. Text classification in Asian languages without word segmentation. In Proceedings of The Sixth International Workshop on Information Retrieval with Asian Languages (IRAL 2003). Association for Computational Linguistics, July 7. Sapporo, Japan.
- Cavnar, W. and Trenkle, J. 1994. N-gram-based text categorization. In Proceedings of Third Annual Symposium on Document Analysis and Information Retrieval (SDAIR). (11-13 April 1994): pp.161-175. Las Vegas, USA.

- Chen, S.F. and Goodman, J.T. 1998. An Empirical study of smoothing techniques for language modeling. In Proceedings of Thirty-Fourth Annual Meeting of the Association for Computational Linguistics, pp.310-318.
- Churcher Gavin. 1994. Distinctive character sequences. Personal communication.
- Churcher Gavin, Hayes Judith, Johnson Stephen, and Souter Clive. 1994. Bigraph and trigraph models for language identification and character recognition. In Proceedings of 1994 AISB Work-shop on Computational Linguistics for Speech and Handwriting Recognition, April. University of Leeds.
- Combrinck, H.P. and Botha, E.C. 1995. Text-based automatic language identification. In Proceedings of the Sixth Annual South African Workshop on Pattern Recognition, November. Rand Afrikaans University.
- Dunning, T. 1994. Statistical identification of language. In Technical report CRL MCCS-94-273, Computing Research Lab, March. New Mexcio State University.
- Harbeck, S Ohler, U, Noth, E and Niemann, H. 1999. Information theoretic based segments for language identification. In Proceeding of the Workshop on Text, Speech and dialog (TSD'99), pp.189-202. Berlin.
- Hayes Judith. 1993. Language recognition using two-and three-letter clusters. Technical Report. school of computer studies, University of Leeds.
- Johnson Stephen. 1993. Solving the problem of language recognition. Technical report, school of computer studies, University of Leeds.
- Peng, F. and Schuurmans D. 2003. Combing naïve bayes and n-gram language models for text classification. In F. Sebastiani (Eds.): Advances in Information Retrieval: Proceedings of The 25th European Conference on Information Retrieval Research (ECIR03), LNCS 2633: pp.335-350. Pisa, Italy.
- Peng F., Schuurmans D. and Wang S. 2003. Language and task independent text categorization with simple language models. In Proceedings of Human Language Technology Conference and North American Chapter of the Association for Computational Linguistics (HLT-NAACL'03), (May 27 - June 1 2003): pp.110-117. Edmonton, Canada.

- Peng F., Schuurmans D. and Wang S. 2004. Augmenting Naive Bayes text classifier with statistical language models. Information Retrieval 7 (3-4), pp.317–345.
- Peng, F., Schuurmans D., Wang S. and Huang X. 2003. Text classification in Asian languages without word segmentation. In Proceedings of The Sixth International Workshop on Information Retrieval with Asian Languages (IRAL 2003). Association for Computational Linguistics, July 7. Sapporo, Japan.
- Piotrowki Michael. 2000. Statistical language identification for NLP-supported full text retrieval. Master's Thesis, April 28.
- Ponte, M and Bruce Croft. 1998. A Language modeling approach to information retrieval. In Proceeding of the 21st annual international ACM SIGIR conference on Research and development in information retrieval, pp.275-281. Australia.
- Sibun Penelope and Reynar Jeffrey C. 1996. Language identification: examining the issues. In 5th symposium on document analysis and information retrieval, pp.125-135. Las Vegas, USA.
- Teahan, W. and Harper, D. 2001. Using compression-based language models for text categorization. In Proceedings of Workshop on LMIR.
- Yang, Y. 1999. An evaluation of statistical approaches to text categorization. Journal of Information Retrieval 1(1/2), pp.69-90.