

สำนักหอสมุดกลาง พระจอมเกล้าลาดกระบัง

การวิเคราะห์การตัดคำภาษาไทย

Thai Word Segmentation Analysis



น.ส. วราภรณ์ สภาวรรตนากุล

น.ส. สุภัตรา ศรีเกษมศาสตร์



ปริญญานิพนธ์นี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญาวิศวกรรมศาสตรบัณฑิต

ภาควิชาวิศวกรรมคอมพิวเตอร์

คณะวิศวกรรมศาสตร์

สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

ปีการศึกษา 2546

เลขหมู่.....

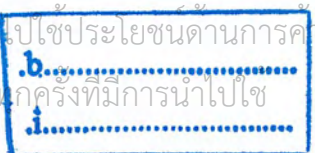
55120

เลขทะเบียน.....

8 เมษายน 2548

วัน,เดือน,ปี.....

เอกสารนี้เป็นเอกสารที่สงวนลิขสิทธิ์ไว้เพื่อการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดก็ตาม ผู้ที่นำเอกสารนี้ไปเผยแพร่โดยไม่ได้รับอนุญาตถือว่าผิดกฎหมาย และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



การวิเคราะห์การตัดคำภาษาไทย

Thai Word Segmentation Analysis



ปริญญานิพนธ์นี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญาวิทยาศาสตรบัณฑิต

ภาควิชาวิศวกรรมคอมพิวเตอร์

คณะวิศวกรรมศาสตร์

สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

ปีการศึกษา 2546

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ปริญญานิพนธ์ ปีการศึกษา 2546

ภาควิชาวิศวกรรมคอมพิวเตอร์

คณะวิศวกรรมศาสตร์ สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

เรื่อง

การวิเคราะห์การตัดคำภาษาไทย

Thai Word Segmentation Analysis

คณะผู้จัดทำ

น.ส. วราภรณ์ สกาวรัตนากุล รหัส 43010375

น.ส. สุพัตรา ตรีเกษตรศาสตร์ รหัส 43010488



..... อาจารย์ที่ปรึกษา
(ดร. วัชระ นัครวิริยะ)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

การวิเคราะห์การตัดคำภาษาไทย

น.ส. วราภรณ์ สกาวรัตน์กุล 43010375

น.ส. สุพิศรา ศรีเกษตรศาสตร์ 43010488

คร. วัชรระ นัครวิริยะ อาจารย์ที่ปรึกษา

ปีการศึกษา 2546

บทคัดย่อ

งานวิจัยนี้มีวัตถุประสงค์เพื่อวิเคราะห์และเปรียบเทียบประสิทธิภาพของหลักการตัดคำภาษาไทย ที่มีอยู่ โดยพิจารณาประสิทธิภาพของโปรแกรมในหลาย ๆ ด้านร่วมกัน เช่น ความเร็ว ความถูกต้อง ทรัพยากรที่ใช้เพื่อให้ผู้ใช้สามารถนำรูปแบบการตัดคำที่เหมาะสมมาประยุกต์ใช้กับงานที่ต้องการ ในการทำงานวิจัยนี้ ผู้วิจัยได้ศึกษาการตัดคำและคลังข้อความด้วยวิธีการต่าง ๆ แล้วนำมาสร้างเป็นโปรแกรม และทำการทดสอบประสิทธิภาพในสภาวะที่มีความใกล้เคียงกันมากที่สุด หลังจากนั้นจึงทำการสรุปวิเคราะห์ประสิทธิภาพในด้านต่าง ๆ

ทั้งนี้ผู้วิจัยได้ทำการเขียนโปรแกรมที่ใช้ในการทดลองประสิทธิภาพโดย MATLAB version 6.5 และข้อมูลและข้อมูลที่ได้ทำการทดลองได้แก่ คลังข้อความ และพจนานุกรม ได้นำมาจากศูนย์เทคโนโลยีอิเล็กทรอนิกส์และคอมพิวเตอร์แห่งชาติ (NECTEC)

จากผลการวัดประสิทธิภาพพบว่า การตัดคำโดยใช้พจนานุกรมมีความถูกต้องมากที่สุด การตัดคำโดยกฎมีความถูกต้อง, การใช้ทรัพยากรและเวลาในการตัดคำน้อยที่สุด การตัดคำโดยใช้คลังข้อความมีการใช้ทรัพยากรมากที่สุดและใช้เวลาในการตัดคำนานที่สุด

Thai Word Segmentation Analysis

Waraporn Sakaorattanakul 43010375

Supattra Trikasetsart 43010488

Dr. Watchara Chatwiriya Adviser

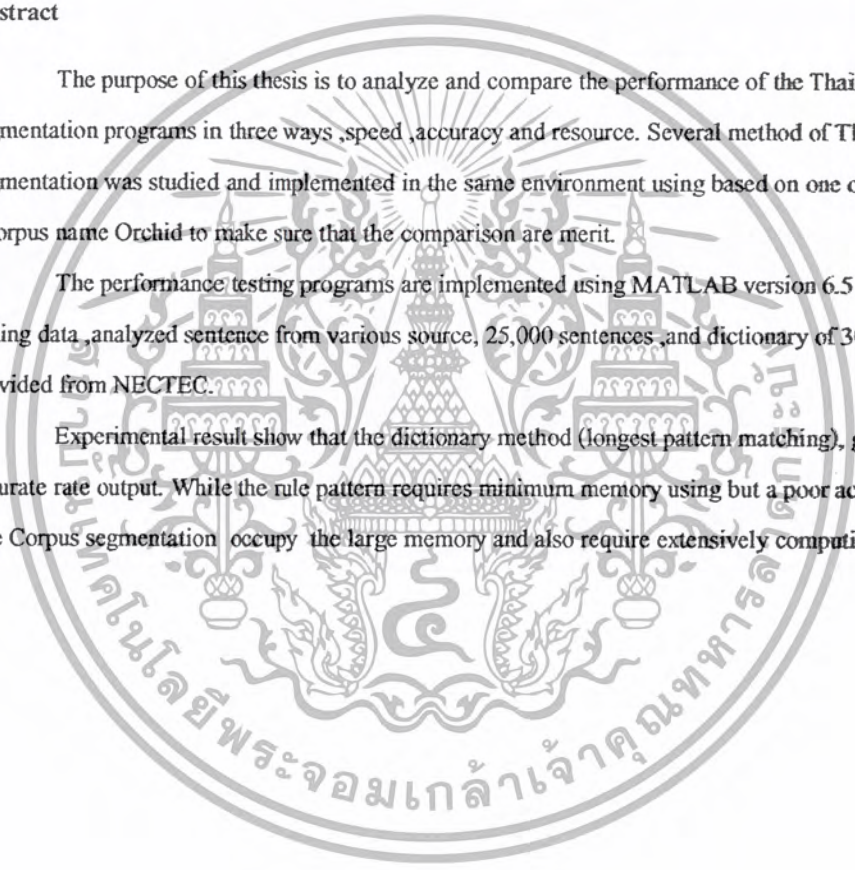
Academic Year 2004.

Abstract

The purpose of this thesis is to analyze and compare the performance of the Thai word segmentation programs in three ways ,speed ,accuracy and resource. Several method of Thai word segmentation was studied and implemented in the same environment using based on one of the analyze a corpus name Orchid to make sure that the comparison are merit.

The performance testing programs are implemented using MATLAB version 6.5 .Data and testing data ,analyzed sentence from various source, 25,000 sentences ,and dictionary of 30,000 are provided from NECTEC.

Experimental result show that the dictionary method (longest pattern matching), gives the most accurate rate output. While the rule pattern requires minimum memory using but a poor accurate rate. The Corpus segmentation occupy the large memory and also require extensively computing power.



กิตติกรรมประกาศ

ปริญญานิพนธ์ฉบับนี้สำเร็จลุล่วงได้อย่างดี ด้วยคำแนะนำ คำปรึกษาจากอาจารย์ที่ปรึกษา คร.วัชระ ภัทรวิริยะ ซึ่งต้องขอขอบพระคุณเป็นอย่างสูง ขอขอบคุณ คุณ ไพศาล เจริญพรสวัสดิ์ ที่ห้องปฏิบัติการวิจัยภาษาและวิทยาการความรู้(linguistics and knowledge science laboratory: LINKS) ของศูนย์วิจัยเทคโนโลยีอิเล็กทรอนิกส์และคอมพิวเตอร์แห่งชาติ(National Electronics and computer technology center:Nectec) ที่ให้คำแนะนำที่เป็นประโยชน์รวมทั้งคลังข้อความและพจนานุกรมที่ใช้ในงานวิจัยนี้ ขอขอบคุณภาควิชาวิศวกรรมคอมพิวเตอร์ที่มีอินเทอร์เน็ตความเร็วสูงให้บริการ สำหรับการค้นคว้าหาความรู้ต่างๆ ซึ่งท้ายที่สุดแล้วก็ประกอบกันเป็นส่วนหนึ่งของโครงการนี้ขอขอบคุณพี่ๆ เพื่อนๆ ในห้องปฏิบัติการเน็ตเวิร์กที่คอยให้ความช่วยเหลือและกำลังใจเสมอมา และสุดท้ายต้องขอขอบคุณบุคคลที่สำคัญที่สุดในชีวิตที่ทำให้ข้าพเจ้ามีวันนี้ นั่นคือ บิดา มารดาและบุคคลในครอบครัว อันเป็นที่เคารพรัก ซึ่งพร้อมให้โอกาสในการศึกษาอย่างเต็มที่ และยังให้กำลังใจ ความรักเสมอมา ข้าพเจ้าขอกราบพระคุณมา ณ ที่นี้ด้วย

วราภรณ์ สกาวรัตน์กุล
สุภัตรา ศรีเกษตราศาสตร์



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญ

บทคัดย่อภาษาไทย	I
บทคัดย่อภาษาอังกฤษ	II
กิตติกรรมประกาศ	III
สารบัญ	IV
สารบัญตาราง	VI
สารบัญภาพ	VII
บทที่ 1 บทนำ	1
1.1 ความเป็นมาของปัญหา	1
1.2 วัตถุประสงค์ของวิทยานิพนธ์	2
1.3 ขอบเขตงานวิจัย	2
1.4 ขั้นตอนงานวิจัย	3
1.5 ประโยชน์ที่คาดว่าจะได้รับจากงานวิจัย	3
บทที่ 2 แนวคิดและทฤษฎีที่เกี่ยวข้อง	4
2.1 การตัดคำ	4
2.2 เครื่องหมายและการแบ่งวรรคตอนภาษาไทยที่ถูกต้อง	4
2.3 การพัฒนาขั้นตอนวิธีการตัดคำแบบต่าง ๆ	5
2.3.1 การตัดคำด้วยกฎ	5
2.3.2 การตัดคำโดยใช้พจนานุกรม	6
2.3.3 การตัดคำโดยใช้คลังข้อความ	8
บทที่ 3 มาตรฐานประสิทธิภาพการตัดคำแบบต่าง ๆ	10
3.1 ความถูกต้องหลังจากตัดคำแล้ว	10
3.2 มาตรฐานประสิทธิภาพเชิงความเร็ว	11
3.3 มาตรฐานประสิทธิภาพการใช้ทรัพยากร	11
บทที่ 4 ข้อมูลที่ใช้	12
4.1 คลังข้อความ	12
4.2 พจนานุกรม	14

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 5 โปรแกรมและอัลกอริทึมในการตัดคำภาษาไทยที่นำมาทดสอบ	15
5.1 การตัดคำด้วยกฎ	15
5.2 การตัดคำโดยใช้พจนานุกรม	18
5.2.1 โครงสร้างพจนานุกรม	18
5.2.2 การแยกหน่วยคำ	19
5.3 การตัดคำโดยใช้คลังข้อความ	22
บทที่ 6 ขั้นตอนในการวัดประสิทธิภาพ	26
6.1 ปรับเปลี่ยนรูปแบบของอินพุต	26
6.2 วัดประสิทธิภาพด้วย Framework ที่สร้างขึ้น	26
6.2.1 เงื่อนไขของการใช้ framework	27
6.2.2 ขั้นตอนการใช้ framework	27
บทที่ 7 ผลการทดลอง	29
7.1 ความถูกต้องของคำ	29
7.2 ประสิทธิภาพเชิงความเร็ว	30
7.3 ประสิทธิภาพการใช้ทรัพยากร	31
บทที่ 8 การวิเคราะห์ผลจากการตัดคำแล้ว	32
8.1 โปรแกรมการตัดคำด้วยพจนานุกรม	32
8.2 โปรแกรมการตัดคำด้วยคลังข้อความ	33
8.3 โปรแกรมการตัดคำด้วยกฎ	34
บทที่ 9 บทวิจารณ์และสรุป	35
9.1 สรุปประสิทธิภาพของโปรแกรมการตัดคำ	35
9.2 ความเหมาะสมของการประยุกต์ใช้โปรแกรมตัดคำ	35
9.3 ประเภทเอกสารที่ผลของการเปรียบเทียบประสิทธิภาพการตัดคำ	36
9.4 ปัญหาต่างๆที่พบในงานวิจัย	36
9.5 ข้อเสนอแนะและงานวิจัยที่สามารถทำเพิ่มเติมต่อจากงานวิจัยนี้	36
ภาคผนวก	37
ภาคผนวก ก	37
ภาคผนวก ข	40
ภาคผนวก ค	76
บรรณานุกรม	86

สารบัญตาราง

ตารางที่ 2-1 เครื่องหมายวรรคตอนภาษาไทย	4
ตารางที่ 7-1 ค่าความถูกต้อง,จำนวนคำที่ตัดได้ถูกต้อง,จำนวนคำที่ตัดได้ทั้งหมด	29
ตารางที่ 7-2 ค่าความเร็วเฉลี่ยในการตัดคำ	30
ตารางที่ 7-3 จำนวนการใช้ทรัพยากรของโปรแกรมการตัดคำแบบต่าง ๆ	31



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญรูป

รูปที่ 4-1 แสดงลักษณะของคลังข้อความ	12
รูปที่ 4-2 ส่วนของรายละเอียดของโปรแกรมที่นำมาตัดคำ	13
รูปที่ 4-3 ส่วนของประโยคและคำที่ตัดไว้เรียบร้อยแล้ว	13
รูปที่ 4-4 ลักษณะของพจนานุกรม	14
รูปที่ 5-1 แสดงการตัดคำแบบทั่ว ๆ ไป ในส่วนการหาขอบเขตหลัง	16
รูปที่ 5-2 ส่วนของการหาขอบเขตหน้าของสระตัวหลัง	17
รูปที่ 5-3 แสดงโครงสร้างของพจนานุกรมไทย(ในรูปแบบ chart)	18
รูปที่ 5-4 แสดงโครงสร้างของพจนานุกรมไทย(ในรูปแบบ tree)	19
รูปที่ 5-5 แสดงการทำงานของโปรแกรมการตัดคำด้วยพจนานุกรม	21
รูปที่ 5-6 ขั้นตอนการนำวิธีวิเทอร์บีมาใช้ในการตัดคำภาษาไทย	23
รูปที่ 5-7 แสดงตัวอย่างการทำวิเทอร์บีร่วมกับวิธีการตัดคำภาษาไทย	24
รูปที่ 6-1 แสดงตัวอย่างของไฟล์ detail_dict.txt	27
รูปที่ 6-2 แสดงการทำงานในส่วนต่าง ๆ ของ framework	28
รูปที่ 7-1 แสดงการเปรียบเทียบโดยค่า f-measure	29
รูปที่ 7-2 แสดงการเปรียบเทียบความเร็วเฉลี่ย	30
รูปที่ 7-3 แสดงการเปรียบเทียบการใช้ทรัพยากร	31

บทนำ

หลักการตัดคำภาษาไทยได้มีการพัฒนาติดต่อกันมานานแล้ว และได้มีการพัฒนาวิธีการต่างๆ เพื่อให้เหมาะสมกับงานแต่ละงาน โดยในระบบต่างๆ ที่นำการตัดคำไปใช้นั้นก็มีความต้องการประสิทธิภาพในแต่ละด้านของการตัดคำไม่เท่ากัน เช่น ในงานบางอย่างอาจจะต้องการความถูกต้องในการตัดคำอย่างมาก มิฉะนั้นจะส่งผลให้ระบบไม่สามารถทำงานได้ถูกต้อง หรือบางงานอาจจะไม่ต้องการความถูกต้องมากนัก แต่ต้องการความเร็วในการตัดคำมากกว่า ส่งผลทำให้มีการพัฒนาการตัดคำในรูปแบบต่างๆ ขึ้น

1.1 ความเป็นมาของปัญหา

ปัจจุบันมีการนำคอมพิวเตอร์เข้ามาใช้ในงานด้านต่าง ๆ อย่างแพร่หลายรวมถึงการนำไปใช้ งานด้านการประมวลผลภาษาธรรมชาติ(Natural Language Processing) การประมวลผลภาษา ธรรมชาติคือกระบวนการที่ทำให้คอมพิวเตอร์สามารถที่จะเข้าใจภาษามนุษย์ได้

แต่เนื่องจากการเขียนคำในภาษาไทยมีรูปแบบการเขียนที่เรียงติดต่อกัน จะมีการแบ่งวรรค คอน ระยะเวลา เพื่อให้ผู้อ่านได้หยุดพัก ซึ่งต่างจากภาษาอื่น ๆ เช่น ภาษาอังกฤษ ภาษาฝรั่งเศส ที่ มีวรรคคอนเป็นตัวแบ่งคำอย่างชัดเจน ดังนั้นจึงทำให้การประมวลผลภาษาไทยด้วยคอมพิวเตอร์ จึงมีความแตกต่างออกไป โดยจำเป็นที่จะต้องมียุทธวิธีการในการหาขอบเขตของคำ เพื่อให้ได้คำที่ ถูกต้องตามหลักภาษาศาสตร์

ดังนั้นการตัดคำหรือการแยกคำออกจากประโยคจึงเป็นสิ่งที่มีความจำเป็นขั้นต้นต่อการ ประมวลผลทางภาษาธรรมชาติ และตัวอย่างงานประมวลผลทางภาษาธรรมชาติที่ต้องการการตัด คำนั้นเช่น

- การวิเคราะห์ไวยากรณ์ (syntax analysis)
- การแปลเอกสาร (machine translation)
- การสังเคราะห์เสียงพูด (speech synthesis)
- การทำดัชนีสำหรับเอกสาร (document indexing)
- การจัดรูปแบบเอกสาร ในงานประมวลผลคำ (word processing)

ในงานวิจัยของ ไพศาล เจริญพรสวัสดิ์ [4] ได้กล่าวถึงงานในด้านการประมวลผลภาษาธรรม ชาติไว้ว่า ในงานด้านนี้นอกจากที่จะต้องรู้ขอบเขตของคำแล้วบางงานยังมีความจำเป็นที่ต้องการ ทราบหน้าที่คำ (Part of Speed) หรือความหมาย (Semantic) ของคำด้วยเพื่อที่จะสามารถนำไปใช้ใน การประมวลผลให้มีประสิทธิภาพมากยิ่งขึ้น เช่นในการแปลภาษา การที่จะแปลให้ถูกต้องนั้น นอก

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับบริการเชิงงานเพื่อการศึกษาเท่านั้น ไม่นับผูกขาดเห็นไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ถูกต้องนั้น นอกจากจะต้องทราบขอบเขตของคำแล้ว การทราบหน้าที่ของคำจะเพิ่มความถูกต้องในการแปลด้วย เช่นคำว่า “เกาะ” อาจจะแปลเป็นภาษาอังกฤษได้เป็น “To attach” หรือ “Island” ซึ่งทั้ง 2 คำมีหน้าที่ต่างกัน ดังนั้นหากเราทราบถึงหน้าที่คำ เช่นถ้าต้องการแปลคำว่า “เกาะ” ที่มีหน้าที่คำเป็นคำนามเราก็จะแปลเป็น “Island” ดังนั้นการทราบหน้าที่คำจึงส่งผลทำให้การแปลภาษามีความถูกต้องมากยิ่งขึ้น

งานวิจัยของ ดร.รัตติกร วรากุลศิริพันธุ์และทีมงาน [13] ได้กล่าวถึงปัญหาของการประมวลผลภาษาไทยด้วยคอมพิวเตอร์ไว้ว่า ลักษณะพื้นฐานของภาษาไทยประกอบด้วยหน่วยคำ (morpheme) ที่เขียนติดกัน โดยไม่มีเครื่องหมายหรือช่องว่างบอกคำจบ ซึ่งทำให้เกิดความกำกวมเมื่อทำการประมวลผลประโยคภาษาไทยด้วยคอมพิวเตอร์จึงต้องอาศัยเทคนิคหรืออัลกอริทึมที่สามารถแยกหน่วยคำออกจากประโยคให้ได้ทั้งความถูกต้องและปราศจากความกำกวมไม่ว่าจะเป็นทางด้านไวยากรณ์หรือความหมาย

ในงานวิจัยหลายงานวิจัยนั้นได้ทำศึกษาวិธีการตัดคำ และได้กล่าวถึงประสิทธิภาพของโปรแกรมการตัดคำไว้แต่สังเกตได้ว่าข้อมูลที่ผู้วิจัยนั้น ๆ นำมาทดสอบประสิทธิภาพมีความแตกต่างกัน จึงทำให้การเปรียบเทียบประสิทธิภาพไม่มีมาตรฐานเพียงพอที่จะสรุปได้ว่าโปรแกรมการตัดคำในแบบต่าง ๆ ดีกว่ากันอย่างไร ความต้องการที่จะเปรียบเทียบประสิทธิภาพในแง่มุมต่าง ๆ ของขั้นตอนและวิธีการตัดคำแบบต่าง ๆ จึงจำเป็นในการนำไปสู่การพัฒนาต่อไป และยังช่วยให้เลือกวิธีการตัดคำให้เหมาะสมกับงานที่ต้องการได้

1.2 วัตถุประสงค์ของวิทยานิพนธ์

- 1.2.1 เพื่อศึกษาแนวทางที่ใช้ในการตัดคำแบบต่าง ๆ ที่มีอยู่ เพื่อหาแนวทางในการเปรียบเทียบประสิทธิภาพของการตัดคำภาษาไทย
- 1.2.2 เปรียบเทียบข้อดี ข้อเด่น ของประสิทธิภาพในการตัดคำ โดยใช้ Framework ที่สร้างขึ้น เช่นประสิทธิภาพด้านความถูกต้อง ด้านความเร็ว และด้านการใช้ทรัพยากร โดยใช้ข้อมูลเดียวกันเพื่อความเป็นมาตรฐาน
- 1.2.3 ทำการวิเคราะห์คำที่ได้จากโปรแกรมการตัดคำ เช่น สาเหตุที่ทำให้เกิดการตัดคำที่ผิดพลาดขึ้นเพื่อใช้ในการพัฒนาโปรแกรมการตัดคำต่อไป

1.3 ขอบเขตของวิทยานิพนธ์

- 1.3.1 ใช้ข้อมูลที่ได้จากเนคเทค (NECTEC) ซึ่งมีประโยคประมาณ 25000 ประโยคและได้ทำการตัดคำและกำกับหน้าที่ของคำไว้เรียบร้อยแล้ว โดยนักภาษาศาสตร์
- 1.3.2 ทำการทดลองโดยใช้วิธีการตัดคำด้วยกฎ การตัดคำแบบlongest และการตัดคำโดยใช้คลังข้อความ ในการเปรียบเทียบประสิทธิภาพในด้านต่าง ๆ
- 1.3.3 สร้าง framework เพื่อใช้ในการทดลอง

1.4 ขั้นตอนการวิจัย

- 1.4.1 กำหนดขอบเขตงานวิจัย
- 1.4.2 รวบรวมงานวิจัยทางด้านการตัดคำและงานที่เกี่ยวข้องกับการตัดคำที่ผ่านมา
- 1.4.3 ศึกษาการตัดคำวิธีการตัดคำแบบต่างๆที่ผ่านมา
- 1.4.4 สร้าง framework เพื่อใช้ในการทดลองวัดประสิทธิภาพและเปรียบเทียบผลของการตัดคำในรูปแบบต่างๆ
- 1.4.5 ทำการทดลองเพื่อวัดประสิทธิภาพและเปรียบเทียบผลของการตัดคำในรูปแบบต่าง ๆ
- 1.4.6 สรุปและวิเคราะห์ปัญหา
- 1.4.7 จัดทำเอกสาร

1.5 ประโยชน์ที่คาดว่าจะได้รับ

- 1.5.1 ได้ข้อมูลประสิทธิภาพของโปรแกรมการตัดคำภาษาไทย ในแง่มุมต่าง ๆ เพื่อให้เหมาะสมกับงาน
- 1.5.2 เป็นเครื่องมือในการวัดประสิทธิภาพของอัลกอริทึมที่ต้องการ
- 1.5.3 เป็นแนวทางนำไปสู่การพัฒนาการตัดคำภาษาไทยต่อไป



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 2

แนวคิดและทฤษฎีที่เกี่ยวข้อง

2.1 การตัดคำ

การตัดคำคือการแบ่งสายอักขระ (String) เพื่อหาขอบเขตของแต่ละหน่วยคำ (Morpheme) เนื่องจากโดยปกติแล้ว ภาษาไทยมีการเขียนในลักษณะที่ติดต่อกัน โดยไม่มีการใช้เครื่องหมายวรรคตอนใด ๆ ยกเว้นแต่มีการเว้นวรรคเป็นระยะ ๆ ให้ผู้อ่านได้หยุดพัก และทำความเข้าใจความหมายเป็นตอน ๆ ไปเท่านั้น แม้ว่าการมีเว้นวรรคในการเขียนบทความไม่ได้มีกฎเกณฑ์ที่ชัดเจน แต่ก็ช่วยลดความคลุมเครือของคำหรือประโยคได้ ในการตัดคำหรือการหาขอบเขตของหน่วยคำที่ต่อเนื่อง ถ้าหากเรารวบรวมคำทุกคำที่มีอยู่ในภาษานั้น ๆ ลงในพจนานุกรมทั้งหมดแล้วจากนั้นก็ค้นหาและเปรียบเทียบคำศัพท์นั้น ๆ ว่ามีอยู่ในพจนานุกรมหรือไม่ เท่านั้น ก็จะสามารถหาขอบเขตของคำได้ทั้งหมด แต่ในความเป็นจริง คำบางคำเป็นก็ไม่ได้บรรจุในพจนานุกรม เช่น ชื่อ เฉพาะ หรือคำที่เกิดจากการใช้คำใหม่

2.2 เครื่องหมายและการแบ่งวรรคตอนที่ต้อง

เครื่องหมายวรรคตอนตามเอกสารเรื่อง “หลักเกณฑ์การใช้เครื่องหมายวรรคตอน และเครื่องหมายอื่น ๆ หลักเกณฑ์การเว้นวรรคหลักเกณฑ์การเขียนคำย่อ” ของราชบัณฑิตยสถาน พ.ศ.2530 ดังตารางที่ 2.1

ชื่อภาษาไทย	ชื่อภาษาอังกฤษ	เครื่องหมายวรรคตอน
มหัพภาค	full stop ,period	.
จุดภาค	Comma	,
อัฒภาค	Semicolon	;
ทวิภาค	Colon	:
วิภัชภาค	-	:-
ขีดกึ่งคี่	Hyphen	-
นกลิขิต (วงเล็บ)	Parentheses	()
วงเล็บเหลี่ยม	Square brackets	[]
วงเล็บปีกกา	Braces	{}
ประศนี	Question mark	?
อัศเจรีย์	Exclamation mark	!
อัญประกาศ	Double quotation marks	" "

เอกสารนี้เป็นเอกสารสงวนไว้สำหรับการใช้ของบุคลากรเท่านั้น ไม่อนุญาตให้นำไปเผยแพร่ภายนอก การนำเอกสารนี้ไปใช้โดยไม่ได้รับอนุญาตถือว่าผิดกฎหมาย

ัญประกาศเดี่ยว	Single quotation marks	‘ ’
ไม้ยมก หรือ ยมก	-	๗
ไปยาลน้อย หรือ เปยยาลน้อย	-	๗
ไปยาลใหญ่ หรือ เปยยาลใหญ่	-	๗๗
ไข่ปลา หรือ จุดไข่ปลา	Elipsis, dotted line	...
เส้นประ	Dashed line	---
เสมอภาค หรือ เท่ากับ	Equals	=
ัญประกาศ	Underline	ขีดเส้นใต้
บุพัญญา	Ditto mark	
มหัตถัญญา	-	(ย่อหน้าขึ้นบรรทัดใหม่)
ทับ	Virgule, slant, slash	/

ตารางที่ 2-1 เครื่องหมายวรรคตอนภาษาไทย

2.3 การพัฒนาและขั้นตอนวิธีการตัดคำภาษาไทยแบบต่างๆ

เนื่องจากการตัดคำได้มีการพัฒนาคิดค้นการเป็นเวลานาน ทำให้มีงานวิจัยการตัดคำเกิดขึ้นมากมาย ซึ่งในช่วงแรกนั้น ได้มีการพัฒนาการตัดพยางค์ขึ้นมาก่อน หลังจากนั้นค่อยมีการพัฒนาด้านการคิดคำขึ้นตามมา ในบทนี้จะกล่าวถึงงานวิจัยที่เกี่ยวข้องกับการตัดคำที่ผ่านมาโดยแบ่งเป็น 3 ชุด คือ ชุดการใช้กฎ, ชุดการใช้พจนานุกรม และชุดการใช้คลังข้อความ

2.3.1 กฎการหาขอบเขตหน้า-ขอบเขตหลัง (สุรินทร์ จรรยาพรพงษ์ [14])

สมัยก่อนคอมพิวเตอร์ไม่มีความสามารถในการประมวลผลสูงมากนัก ประกอบกับหน่วยความจำในเครื่องคอมพิวเตอร์มีขนาดเล็ก จึงต้องพัฒนาการตัดพยางค์ขึ้นมาก่อน เนื่องจากพยางค์มีรูปแบบกฎเกณฑ์ที่ชัดเจนแน่นอนมากกว่าคำ

สุรินทร์ จรรยาพรพงษ์ [14] ได้ทำการวิจัยเกี่ยวกับการตัดคำภาษาไทย โดยใช้พยางค์ โดยกฎที่นำมาใช้นั้นได้นำมาจากหลักไวยากรณ์ภาษาไทย และได้ทำการวิเคราะห์ลักษณะต่างๆของพยางค์ภาษาไทยออกมาได้ 200 กว่ากฎ (ภาคผนวก ข) โดยลักษณะของกฎที่ได้นี้สามารถแบ่งได้เป็น 2 ชนิดคือ กฎการหาขอบเขตหน้า(Front boundary recognition) และกฎการหาขอบเขตหลัง (Tail boundary recognition) และในแต่ละกฎยังแบ่งออกเป็น 2 กลุ่มย่อย ๆ คือ แบ่งตามคุณสมบัติของตัวอักษร โดยกฎที่ได้เอามาจะจัดให้อยู่ในกลุ่มเอ (Group A) และแบ่งตามคุณสมบัติของรูปแบบการใช้สระ ซึ่งจัดกลุ่มนี้เป็นกลุ่มบี (Group B)

นอกจากนี้ยังมีการสร้างกฎสำหรับจัดการกับพยางค์ที่ไม่ใช่ลักษณะของพยางค์ภาษาไทย ซึ่ง

อาจประกอบไปด้วย อักษรพิเศษ ตัวเลข คำย่อ หรือ อักษรจากภาษาต่างชาติ

เอกสารนี้เป็นเอกสารที่สงวนลิขสิทธิ์ไว้เพื่อใช้ในการศึกษาเท่านั้น อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ในงานวิจัยของ ควงแก้ว สวามิภักดิ์ [5] ได้กล่าวถึงข้อดีของวิธีนี้ว่า การวิเคราะห์ห้าตัวสระเป็นเกณฑ์ เท่านั้น ไม่ได้นำตัวอักษรที่สามารถเป็นสระมาวิเคราะห์ด้วย เช่นตัวอักษร “อ” และ “ว” ดังนั้น หากนำไปใช้กับข้อความที่ไม่มีรูปสระเลยก็จะคัด ไม่ได้ เช่น “ขอคคคคคคคค” จากตัวอย่างนี้จะแสดงให้เห็นว่าการวิเคราะห์ด้วยรูปสระเพียงอย่างเดียวจะส่งผลที่สมบูรณ์แบบไม่ได้ สมควรจะต้องวิเคราะห์ถึงกฎเกณฑ์ที่ภาษาสามารถใช้ตัวอักษรสร้างสระด้วย

ในงานวิจัยนี้ได้ทดลองกับเอกสารต่าง ๆ จำนวน 100 เล่ม โดยเอกสารนั้น ได้นำมาจากเอกสารชนิดต่าง ๆ จำนวน 10 ชนิด เช่น หนังสือพิมพ์ ปรัชญา วิทยาศาสตร์ ศาสนา ภาษาศาสตร์ ฯลฯ และจากการทดสอบ ปรากฏว่าสามารถตัดพยางค์ได้ถูกต้องถึง 96%

2.3.2 การใช้พจนานุกรม

ในขณะนี้ถือว่าเป็นจุดเริ่มต้นแรกในการตัดคำเนื่องจากคอมพิวเตอร์มีการพัฒนาเพิ่มขึ้น และหน่วยความจำมากขึ้น ทำให้มีการคิดค้นวิธีการแบ่งคำโดยเอาพจนานุกรมเข้ามาใช้ในการตัดคำ

2.3.2.1 Longest Matching (ซิน ภู่วรรณและ วิวรรณ อิมฮารมน์ [8])

ในงานวิจัยนี้ได้ทำการจัดเก็บพยางค์ต่างๆ ไว้ในพจนานุกรม และมีกรนำไวยากรณ์ต่างๆจำนวน 18 กฎเข้ามาช่วยในกรณีที่ไม่มีพยางค์ในพจนานุกรม ได้ออกแบบการแยกพยางค์ โดยใช้หลักการดังนี้ โดยยึดหลักว่าคำที่สั้นกว่ามีขนาดเล็กพอที่จะบรรจุ ไว้ในหน่วยความจำหลัก

หลักการการทำงานคือ จะทำการตรวจสอบสายอักขระ (String) ที่เข้ามาจากซ้ายไปขวากับพยางค์ที่เก็บไว้ในพจนานุกรมเพื่อทำเครื่องหมายกับทุกสายอักขระที่สามารถเป็นหนึ่งคำให้เป็นจุดเลือก สายอักขระที่ยาวที่สุดเป็นตัวเลือกสำหรับคำแรก ถ้าตัวเลือกนี้สามารถทำให้อัลกอริทึมค้นหาคำที่เหลือได้สมบูรณ์ ตัวเลือกนี้ก็จะเป็นคำแรกจริง แต่กรณีที่เลือกพยางค์ที่ยาวที่สุดไปแล้ว ทำให้เกิดพยางค์ที่ไม่ปรากฏในพจนานุกรม ก็ยอมให้มีการย้อนรอย (Back Tracking) แล้ว ไปเลือกพยางค์ที่มีความยาวรองลงมาแทน และทำการค้นหาคำที่เหลือต่อไปเป็นเช่นนั้น ไปเรื่อยๆ จนจบสายอักขระ วิธีการนี้เป็นที่รู้จักกันในชื่อการตัดคำ(พยางค์)แบบเลือกคำ(พยางค์) ยาวที่สุด

ในงานวิจัยนี้ได้เปรียบเทียบความเร็วในการแบ่งพยางค์ ซึ่งสรุปผลได้ว่าเมื่อนำพจนานุกรมมาแบ่งพยางค์จะสามารถตัดพยางค์ได้รวดเร็วกว่าการใช้กฎ โดยที่มีความถูกต้องมากกว่า 99 % แต่วิธีการนี้ยังมีข้อเสียคือ ต้องเสียเนื้อที่ในการเก็บพจนานุกรมในหน่วยความจำเป็นหลักเป็นจำนวน 50 กิโลไบต์ แต่ก็ยังสามารถเก็บข้อมูลพจนานุกรมไว้ในเครื่องคอมพิวเตอร์สมัยนั้นได้

ต่อมา สัมพันธ์ ระรินรัมย์ [10] ได้นำวิธีการนี้มาพัฒนา โดยมีเป้าหมายเน้นที่ประสิทธิภาพในด้านความเร็วของขั้นตอนวิธีในการตัดคำและการลดขนาดพจนานุกรม เนื่องจากเมื่อนำพจนานุกรมเข้ามาใช้ในการตัดคำแล้วจะทำให้ความถูกต้องเพิ่มมากขึ้นกว่าการตัดคำโดยใช้กฎอย่างเดียว ดังนั้นในงานวิจัยจึงไม่ได้เน้นการเพิ่มประสิทธิภาพในด้านความถูกต้องมากนัก เพราะถือว่าการตัดคำโดยใช้พจนานุกรมจะให้ค่าความถูกต้องสูงอยู่แล้ว

โดยรายละเอียดของขั้นตอนวิธีการตัดคำนั้นจะมีวิธีการคล้ายกับงานวิจัยการแบ่งพยางค์โดยใช้ ดิกชันนารี (เช่น กุสุวรรณ และ วิวรรณ อัมอรณ [8]) แต่จะทำการจัดเก็บคำลงในพจนานุกรมแทน พยางค์

2.3.2.2 Maximal Matching (วิรัช ศรีเลิศวานิช [9])

ในงานวิจัยนี้ ได้มีการพัฒนาการตัดคำ โดยเรียกว่า “การตัดคำโดยเลือกแบบเหมือนกันมากที่สุด (Maximal Matching)” สามารถแก้ไขความบกพร่องของการตัดคำแบบเลือกคำยาวที่สุดได้ เพราะการเลือกคำที่ยาวที่สุดเมื่อเจอข้อความโดยไม่พิจารณาถึงข้อความถัดไป พิจารณาเฉพาะบริเวณใกล้เคียงเท่านั้น ทำให้เลือกคำที่ยาวเกินไปตั้งแต่แรก ข้อความที่ตามมาจึงเกิดข้อผิดพลาดได้ แต่วิธีการเลือกคำแบบเหมือนกันมากที่สุดจะพิจารณาข้อความทั้งหมดแทน หลักการของการตัดคำโดยเลือกแบบเหมือนกันมากที่สุดคือ ขั้นตอนแรกคือจะทำการตัดคำที่เป็นไปได้ทุกๆ แบบก่อนแล้วหลังจากนั้นจึงเลือกประโยคที่มีจำนวนคำน้อยที่สุด สำหรับในกรณีที่ตัดคำแล้วเกิดได้จำนวนคำที่เท่ากันก็ให้ทำการตัดคำแบบเลือกคำยาวที่สุดเข้ามาช่วยพิจารณา

แม้ว่าการค้นหาการแบ่งคำทุกทางที่เป็นไปได้ อาจทำให้ต้องเสียค่าใช้จ่าย (cost) มาก โดยทั่วไปวิธีนี้จะให้ความถูกต้องที่สูงกว่าวิธีตัดคำให้ยาวที่สุด แต่ก็มีประโยคจำนวนมากที่ไม่สามารถตัดได้ถูกต้อง โดยวิธีนี้ นอกจากนี้ ได้มีผู้พัฒนาการตัดคำ ของวิรัช เพิ่มเติม โดยแทนที่จะเปรียบเทียบคำที่ยาวสุดเสมอ จะมีการเปิดตารางค่าความยาวสูงสุดที่พบในพจนานุกรมของข้อความที่ค้นหา โดยใช้ตัวอักษร 2 ตัวแรกเป็นตัวเปรียบเทียบ ถ้ามีค่าน้อยกว่าความยาวของข้อความก็ให้เริ่ม แบ่งคำตั้งแต่คำที่น้อยที่สุด เช่น “ฉันมารอกราบ” มีความยาวข้อความเท่ากับ 11 และค่าที่ได้จากการเปิดตารางของ “ฉ” พบว่ามีค่าสูงสุดคือ 5 ฉะนั้นการแบ่งคำจึงเริ่มจากตัวอักษรที่ 5 ลงมา

ในงานวิจัยของคุณ ไพศาล เจริญพรสวัสดิ์ [4] ได้กล่าวถึงข้อดีของวิธีการนี้ว่า วิธีการนี้สามารถช่วยแก้ไขข้อบกพร่องของการตัดคำแบบเลือกคำยาวที่สุดเมื่อเจอข้อความที่กำลังจะมาก่อน โดยไม่มีการพิจารณาถึงข้อความถัดไปซึ่งมีลักษณะเหมือนการใช้ขั้นตอนวิธีแบบโลภ (Greedy Algorithm) ที่พิจารณาเฉพาะบริเวณที่ใกล้เคียงเท่านั้น แต่วิธีการตัดคำแบบเหมือนกันมากที่สุดนั้นจะเป็นการใช้ขั้นตอนวิธีแบบโลภโดยพิจารณาข้อความทั้งหมดแทน แต่วิธีนี้ก็ยังไม่สามารถตัดคำได้ถูกต้องทั้งหมด จำเป็นต้องมีการนำโครงสร้างทางไวยากรณ์หรือความสัมพันธ์ทางความหมายเข้ามาใช้ประกอบการพิจารณาด้วย

2.3.3 การตัดคำโดยใช้การใช้คลังข้อความ (บุญเสริม กิจศิริกุล [6])

ในชุดนี้ ได้มีการพัฒนาคลังข้อความ (Corpus) ขึ้นจำนวนมาก ทำให้มีการนำความรู้ต่างๆ จากคลังข้อความเข้ามาประยุกต์ใช้ด้วย ตัวอย่างความรู้ที่ได้จากคลังข้อความเช่น คำสถิติการใช้คำภายในคลังข้อความและลักษณะไวยากรณ์ที่ใช้ในคลังข้อความ เป็นต้น

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

วิธีการนี้เป็นวิธีการตัดคำโดยใช้ข้อมูลทางสถิติที่ได้จากคลังข้อความที่เตรียมไว้ โดยใช้ค่าสถิติที่เกิดขึ้นจากลำดับของหน้าที่คำหรือประเภทย่อยของคำ หรืออาจกล่าวได้ว่าเป็นการนำเอาส่วนหนึ่งของไวยากรณ์มาใช้ ซึ่งคลังข้อความนี้เราได้ให้นักภาษาศาสตร์ตัดคำและกำกับหมวดคำที่เหมาะสมไว้แล้ว ส่วนปัญหาการตัดคำของภาษาไทยเราสามารถเขียนแทนด้วยสมการดังนี้

$$\arg \max_{W_{1,n}} P(W_{1,n} / C_{1,m}) = \arg \max_{W_{1,n}} P(W_{1,n}) P(C_{1,m} / W_{1,n}) / P(C_{1,m}) \quad (2-1)$$

โดยที่

$C_{1,m} = c_1, c_2, c_3, \dots, c_m$ คือ คำสายอักขระที่เข้ามาเป็นอินพุตของเรา

$W_{1,n} = w_1, w_2, \dots, w_n$ คือ คำที่สามารถตัดได้

สมการที่ 2-1 จึงมีหมายความว่า เรามีสายอักขระ $C_{1,m}$ ที่เป็นอินพุตเข้ามา เราต้องการแบ่งสายอักขระนี้เป็นคำ w_1, w_2, \dots, w_n ($W_{1,n}$) และเนื่องจาก $P(C_{1,m} / W_{1,n})$ มีค่าเท่ากับ 1 และค่า $P(C_{1,m})$ ก็เป็นค่าคงที่ในทุก ๆ ครั้ง ดังนั้นทำให้เราสามารถเขียนสมการที่ 2-1 ในรูปสมการใหม่ได้เป็น

$$\arg \max_{W_{1,n}} P(W_{1,n} / C_{1,m}) = \arg \max_{W_{1,n}} P(W_{1,n}) \quad (2-2)$$

เมื่อพิจารณาถึงภาษาไทย การที่เราจะกำหนดขอบเขตของคำก็ทำได้นั้นเราต้องอาศัยสถิติของหน้าที่คำ N-gram และ สถิติการเกิดคำ ๆ นั้น จากนั้นก็หาค่าความน่าจะเป็นสูงสุดของคำหน้าที่คำ ในแต่ละทั้งประโยค เนื่องจากมีหน้าที่คำเข้ามาเกี่ยวข้องซึ่งช่วยให้สมการที่ 2-2 เปลี่ยนรูปแบบไปเป็นสมการที่ 2-3

$$\arg \max_{W_{1,n}} P(W_{1,n} / C_{1,m}) = \arg \max_{W_{1,n}} \sum_{T_i} P(W_{1,n} T_i) \quad (2-3)$$

โดยที่ $T_{1,n} = t_1, t_2, \dots, t_n$ คือหน้าที่คำของ w_1, w_2, \dots, w_n ตามลำดับ

และเพื่อให้ได้ค่าของ $P(W_i / C_{1,m})$ สูงที่สุด ในหลักการของโปรแกรมโมเดล เราจะตั้งสมมติฐานดังนี้

- 1) ความน่าจะเป็นของคำหนึ่ง ๆ จะปรากฏ ณ ตำแหน่งใด ๆ ในประโยคไม่ขึ้นกับสิ่งอื่น ๆ
- 2) หมวดคำหนึ่ง ๆ จะขึ้นอยู่กับหมวดคำก่อนหน้า 2 หมวดคำเท่านั้น

จากสมมติฐานทั้ง 2 ข้อนั้นจะทำให้สมการที่ 2-3 เปลี่ยนรูปแบบเป็นสมการที่ 2-4

$$\arg \max_{W_{1,n}} P(W_{1,n} / C_{1,m}) = \arg \max_{W_{1,n}} \sum_{T_i} \pi P(W_i T_i) * P(T_i / T_{i-1} T_{i-2}) \quad (2-4)$$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

โดยค่าของ $P(W_T)$ และ $P(T_i/T_{i-1}, T_{i-2})$ สามารถคำนวณจากข้อมูลในคลังข้อมูลที่มีอยู่ เราสามารถหาค่าของ $P(W_T)$ ได้โดยนับจำนวนของหมวดคำ w_i ที่เป็นหมวดคำของคำ w_i หารด้วยจำนวนของ w_i ที่เป็นหมวดคำของคำใดๆ (จำนวนของหมวดคำ w_i ทั้งหมด) ส่วน $P(T_i/T_{i-1}, T_{i-2})$ หาโดยนับจำนวนหมวดคำ w_i ที่มีหมวดคำ w_{i-1} และ w_{i-2} ตามหลังที่ปรากฏในคลังข้อมูลหารด้วยจำนวนของหมวดคำ w_{i-1} และ w_{i-2} ที่ปรากฏทั้งหมด จากค่าเหล่านี้ที่คำนวณได้เราจะนำมาใช้ในการตัดคำภาษาไทยและได้เขียนซอฟต์แวร์เพื่อทดลองผล โดยวิธีการทำงานจะกล่าวถึงในบทที่ 5 หัวข้อที่ 5.3



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 3

มาตรวัดประสิทธิภาพการตัดคำแบบต่าง ๆ

งานวิจัยฉบับนี้มีจุดประสงค์เพื่อวิเคราะห์การทำงานของโปรแกรมตัดคำภาษาไทย การที่เราจะวิเคราะห์ได้ว่าโปรแกรมทำงานมีประสิทธิภาพมากน้อยเพียงใดจำเป็นต้องกำหนดหลักเกณฑ์ที่ใช้ในการวัดประสิทธิภาพ การกำหนดมาตราวัดประสิทธิภาพที่เหมาะสมจะทำให้การเปรียบเทียบความสามารถของโปรแกรมมีความถูกต้องมากยิ่งขึ้น มาตราวัดประสิทธิภาพที่ใช้ในงานวิจัยนี้มี 3 หัวข้อ คือ ความถูกต้อง ความเร็ว และการใช้ทรัพยากร

3.1 ความถูกต้องในตัดคำ

การประเมินความถูกต้องทำได้โดยการตรวจสอบผลลัพธ์ที่ได้จากโปรแกรมตัดคำกับคำตอบซึ่งทำการตัดไว้แล้วโดยนักภาษาศาสตร์ในคลังข้อความ โดยถือว่าเป็นคำตอบที่มีความถูกต้อง 100 %

วิธีการประเมินความถูกต้องมีหลายรูปแบบ ในงานวิจัยฉบับนี้จะใช้มาตรวัดความถูกต้องโดยวิธีการวัดค่า F-Measure (Van Rijbergen, 1979 cited in Manning and Zuchuzel [15]) ซึ่งเป็นวิธีพื้นฐานที่ใช้กันโดยทั่วไปสำหรับวัดความแม่นยำที่ยังอิงความถูกต้องจากเอกสารที่ผ่านการกำกับจากผู้เชี่ยวชาญ โดยทำการพิจารณาจากค่าความแม่นยำและความครบถ้วน คำนวณ ได้ดังสมการที่ 3-1

$$F = \frac{2}{\frac{1}{P} + \frac{1}{R}} \quad (3-1)$$

ค่าความครบถ้วน (Recall) มีการคำนวณดังนี้

$$R = \frac{\text{จำนวนคำที่โปรแกรมตัดคำได้ถูกต้อง}}{\text{จำนวนคำทั้งหมดที่ตัดไว้ในคลังข้อความ}}$$

ค่าความแม่นยำ (Precision) มีการคำนวณดังนี้

$$P = \frac{\text{จำนวนคำที่โปรแกรมตัดคำได้ถูกต้อง}}{\text{จำนวนคำทั้งหมดที่โปรแกรมตัดได้}}$$

ซึ่งในการวัดประสิทธิภาพนี้พิจารณาไหลค่าทั้งสองมีน้ำหนักเท่าๆกัน ($\alpha = 0.5$) ทำให้สามารถแปลงสมการในการวัดค่าให้อยู่ในรูปที่ง่ายขึ้นได้ ดังสมการที่ 3-2 คือ

$$F = \frac{2 \times P \times R}{P + R} \quad (3-2)$$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

3.2 มาตรการวัดประสิทธิภาพเชิงความเร็ว

ในการนำโปรแกรมการตัดคำภาษาไทยไปประยุกต์ใช้กับการใช้งานจริง ปัจจัยที่ต้องนำมา
ร่วมพิจารณาคือความเร็วยังคือ ความเร็วในการทำงาน เนื่องจากงานในแต่ละรูปแบบจะมีความต้องการใน
เรื่องความเร็วที่แตกต่างกัน โดยมาตรการวัดประสิทธิภาพในเรื่องความเร็วนี้จะคำนวณออกมาในรูปแบบ
ความเร็วเฉลี่ยต่อคำ (sec/word)

3.3 มาตรการวัดประสิทธิภาพการใช้ทรัพยากร

การใช้งานทรัพยากรก็เป็นส่วนหนึ่งที่เราต้องคำนึงถึง ในงานบางอย่างที่มีหน่วยความจำมี
ความจำกัดต้องการอัลกอริทึมการตัดคำที่มีการใช้พื้นที่ในหน่วยความจำน้อย โดยหน่วยวัดของมาตร
วัดประสิทธิภาพนี้จะเป็น ขนาดของโค้ด (code) และ ข้อมูลที่เกี่ยวข้องในการตัดคำ ซึ่งมีหน่วยเป็น ไบท์
(byte)



บทที่ 4

ข้อมูลที่ใช้ในการเปรียบเทียบประสิทธิภาพและการตัดคำ

ข้อมูลที่น่ามาใช้ทั้งหมดนี้ ได้รับมาจากห้องวิจัยลิงค์ (LINKS) ที่เนคเทค (NECTEC) ซึ่งสามารถดาวน์โหลดได้ที่ <http://www.links.nectec.or.th/www-new/download.php> โดยคลังข้อความที่ใช้นั้นมีชื่อว่า Orchid Corpus ซึ่งมีประโยคประมาณ 25,000 ประโยค และ พจนานุกรม มีชื่อว่า RIWord ซึ่งเป็นพจนานุกรมลัพท์ราชบัณฑิตยสถานมีคำประมาณ 32,575 คำ

4.1 คลังข้อความ

คลังข้อความที่น่าใช้จะใช้ในการทดสอบและการหาค่าสถิติที่ใช้อ้างอิงในการตัดคำ โดยแบ่งลักษณะของคลังข้อความเป็น 2 ส่วน โดยส่วนแรกจะเป็นส่วนที่ใช้ในการทดสอบประมาณ 20 เปอร์เซ็นต์ เพื่อวัดประสิทธิภาพการตัดคำของแต่ละโปรแกรม ส่วนอีก 80 เปอร์เซ็นต์ ใช้หาค่าสถิติในการตัดคำ ลักษณะของคลังข้อความเป็นดังรูปที่ 4-1

```
%TTitle: การพัฒนาต้นแบบคอมพิวเตอร์ถ่ายภาพตัดขวางอวัยวะ (ระยะที่ 2)
%ETitle: COMPUTER X-RAY TOMOGRAPHY (PHASE 2)
%TAuthor: พิชัย ชัยพงษ์, อัครินทร์ คุณภักดี, นายพงษ์ รังสรรค์เสวี, จงชาย เกียรติกรกุล
%EAuthor:
%TInbook: การประชุมทางวิชาการ ครั้งที่ 5, โครงการวิจัยและพัฒนาอิเล็กทรอนิกส์และคอมพิวเตอร์ ปีงบประมาณ 2535
%Einbook: The 5th Annual Conference, Electronics and Computer Research and Development Project, Fiscal Year 1992
%TPublisher: ศูนย์เทคโนโลยีอิเล็กทรอนิกส์และคอมพิวเตอร์แห่งชาติ, กระทรวงวิทยาศาสตร์ เทคโนโลยีและสิ่งแวดล้อม
%EPublisher: National Electronics and Computer Technology Center, Ministry of Science, Technology and Environment
%Page:
%Year: 1993
%File:
#P1
#1
การพัฒนาต้นแบบคอมพิวเตอร์ถ่ายภาพตัดขวางอวัยวะ (ระยะที่ 2)//
การ/FIXN
พัฒนา/VACT
ต้นแบบ/NCMN
คอมพิวเตอร์/NCMN
ถ่าย/VACT
ภาพตัดขวาง/NCMN
อวัยวะ/NCMN
<space>/PUNC
<left_parenthesis>/PUNC
ระยะ/NCMN
ที่ 2/DONM
<right_parenthesis>/PUNC
//
```

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาค้นคว้าเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

รูปที่ 4-1 แสดงลักษณะของคลังข้อความ

สังเกตได้ว่าคั่งข้อความนั้นจะมีการบอกรายละเอียดที่มาของข้อมูล เช่น มาจากที่ใด ใครเป็นผู้แต่ง แต่ส่วนที่เราใช้นั้นจะเป็นส่วนของ ประโยค คำในประโยคที่ทำการตัดไว้ และหน้าที่ของคำเหล่านั้น ซึ่งทั้งหมดทำโดยนักภาษาศาสตร์ สำหรับสัญลักษณ์และอักษรย่อที่ใช้บอกหน้าที่ของคำสามารถดูรายละเอียดได้ในภาคผนวก ก

```

%TTitle: การพัฒนาต้นแบบคอมพิวเตอร์ถ่ายภาพตัดขวางอวัยวะ (ระยะที่ 2)
%ETitle: COMPUTER X-RAY TOMOGRAPHY (PHASE 2)
%TAuthor: ไพรัช ธัชยพงษ์, อัครินทร์ คุณภักดี, ยุทธพงษ์ รังสรรค์เสวี, สมชาย เกรียงอารีกุล
%EAuthor:
%TInbook: การประชุมทางวิชาการ ครั้งที่ 5, โครงการวิจัยและพัฒนาอิเล็กทรอนิกส์และคอมพิวเตอร์, ปีงบประมาณ 2535
%Einbook: The 5th Annual Conference, Electronics and Computer Research and Development Project, Fiscal Year 1992
%TPublisher: ศูนย์เทคโนโลยีอิเล็กทรอนิกส์และคอมพิวเตอร์แห่งชาติ, กระทรวงวิทยาศาสตร์ เทคโนโลยีและสิ่งแวดล้อม
%EPublisher: National Electronics and Computer Technology Center, Ministry of Science, Technology and Environment
%Page:
%Year: 1993
%File:
#P1
#1

```

รูปที่ 4-2 ส่วนของรายละเอียดของประโยคที่นำมาตัดคำ

```

การพัฒนาต้นแบบคอมพิวเตอร์ถ่ายภาพตัดขวางอวัยวะ (ระยะที่ 2)//
การ/FIXN
พัฒนา/VACT
ต้นแบบ/NCMN
:
อวัยวะ/NCMN
<space>/PUNC
<right_parenthesis>/PUNC
//

```

ประโยคที่นำมาตัดคำ

คำที่ตัดออกมาได้ เช่น "อวัยวะ" มีหน้าที่คำเป็น คำนามทั่วไป

รูป 4-3 ส่วนของประโยคและคำที่ตัดไว้เรียบร้อยแล้ว

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษานี้เท่านั้น เมื่อผู้ใช้งานเห็นใบใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ทั้งนี้ในการทดลองได้ทำการแก้ไขข้อความบางส่วนที่ผิดพลาด เช่น คำที่ตัดได้กับประโยคที่นำมาตัดมีตัวอักษรบางตัวขาดหายไป แต่มีการกำกับคำได้ถูกต้องก็จะทำการเพิ่มตัวอักษรให้เหมือนกัน เป็นต้น

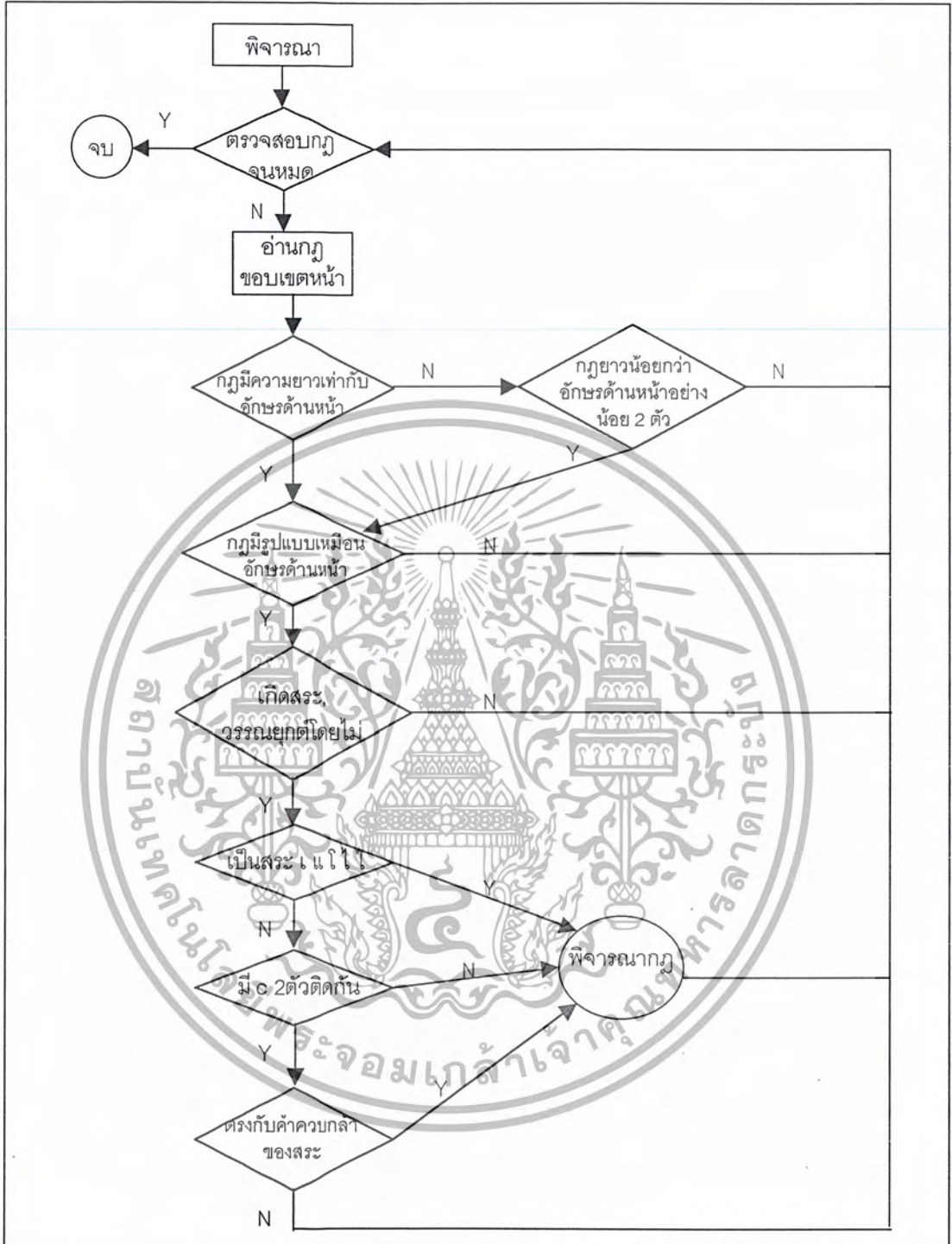
4.2 พจนานุกรม

พจนานุกรมที่ใช้เป็นพจนานุกรมฉบับราชบัณฑิตยสถานมีจำนวนคำทั้งหมด 32575 คำ มีขนาด 248 KB เนื่องจากในพจนานุกรมนี้มีคำซ้ำอยู่บ้างจึงได้คัดออกไปบางส่วน และได้คัดอักษรตัวเดียวออก ยกเว้น “ณ” “ธ” เนื่องจากอักษรอื่นที่นอกเหนือจากนี้ไม่มีการใช้งานเมื่ออยู่ตัวเดียว ลักษณะของพจนานุกรมเป็นดังรูปที่ 4-4



รูป 4-4 ลักษณะของพจนานุกรม

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 5-2 ส่วนของการหาขอบเขตหน้าของสระตัวหลัง

ในการพิจารณาขอบเขตหน้าก็มีหลักการคล้ายคลึงกัน โดยกฎที่ตัดได้ต้องไม่ทำให้เหลืออักษรตัวเดียว ไม่ทำให้เกิดวรรณยุกต์โดยไม่มีตัวอักษรต่อท้าย หากกฎมีพยัญชนะ 2 ตัวติดกัน เฉพาะสระ ใ, ใ, ใ เท่านั้นที่ไม่ต้องพิจารณาตัวควบ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลง 55120 ต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ในการเขียน โปรแกรมได้ใช้หลักการตัดพยางค์ของคุณ สุรินทร์ จรรยาพรพัฒน์ เพียงบางส่วนเท่านั้น เนื่องจาก มีเอกสารบางส่วนขาดหายไป เช่น รายละเอียดเกี่ยวกับตัวอักษรแต่ละประเภท เอกสารเกี่ยวกับคำที่ขกเว้นซึ่งเก็บอยู่ในพจนานุกรม ซึ่งอาจส่งผลกระทบต่อประสิทธิภาพความถูกต้องของการตัดพยางค์ได้

5.2 การตัดคำโดยใช้พจนานุกรม

การตัดคำโดยพจนานุกรมเป็นการใช้พจนานุกรมช่วยในการแบ่งอักษรที่เรียงติดต่อกันออกมาเป็นคำๆ โดยนำอักษรที่เรียงติดกันไปเปรียบเทียบกับพจนานุกรม ถ้าหากพบในพจนานุกรมก็แสดงว่าอักษรที่เรียงติดกันนั้นเป็นคำ ๆ หนึ่ง แต่เนื่องจากพจนานุกรมมีคำจำนวนมากและมีขนาดใหญ่ ดังนั้นการหาคำในพจนานุกรมจะทำการหาแบบ ไล่คู่ทีละตัว (Sequential) จึงไม่สามารถทำได้ เพราะจะใช้เวลาในการหาอย่างมาก ดังนั้นจึงต้องมีการสร้าง โครงสร้างพจนานุกรมขึ้นเพื่อใช้ในการค้นหาในพจนานุกรมมาเปรียบเทียบกับกลุ่มอักษรที่พิจารณาอยู่

5.2.1 โครงสร้างพจนานุกรม

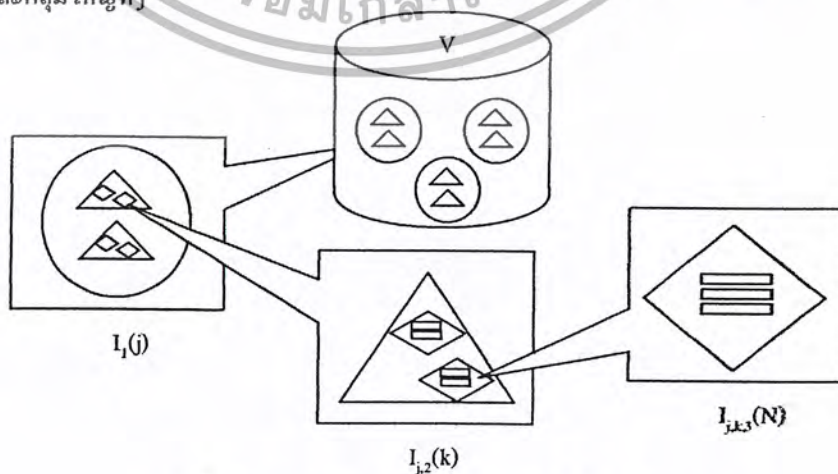
เราจะใช้ตรรกะนี้ในการจัดของคําคำศัพท์ในพจนานุกรมภาษาไทยทั้งหมด โดยจะทำการเก็บบันทึกคำศัพท์ตามค่าครรรชนี 3 ตัว ที่ได้จากค่าแอสกี้นัมเบอร์ (ASCII Number) ของอักษรตัวที่ 1 ค่าแอสกี้นัมเบอร์ของอักษรตัวที่ 2 และจำนวนอักษรทั้งหมดของคำศัพท์นั้น ๆ ซึ่งค่าแอสกี้นัมเบอร์ของตัวอักษรตัวแรกของคำศัพท์เป็นครรรชนีที่แบ่งคำศัพท์แบ่งคำศัพท์ออกเป็นกลุ่มใหญ่ และให้ค่าแอสกี้นัมเบอร์ของอักษรตัวที่ 2 ของคำศัพท์นั้นเป็นครรรชนีที่สองที่แบ่งคำศัพท์ภายในกลุ่มใหญ่ ให้เป็นกลุ่มย่อย และจำนวนของอักษรของคำศัพท์จะถูกใช้เป็นการแบ่งคำศัพท์ในกลุ่มย่อยอีกที

กำหนดให้ V เป็นคําคำศัพท์ทั้งหมดที่บรรจุอยู่ในพจนานุกรม

$I_1(j)$ เป็นครรรชนีตัวแรกที่แสดงกลุ่มคำศัพท์ในกลุ่มที่ j ($j = 1, 2, 3, \dots$)

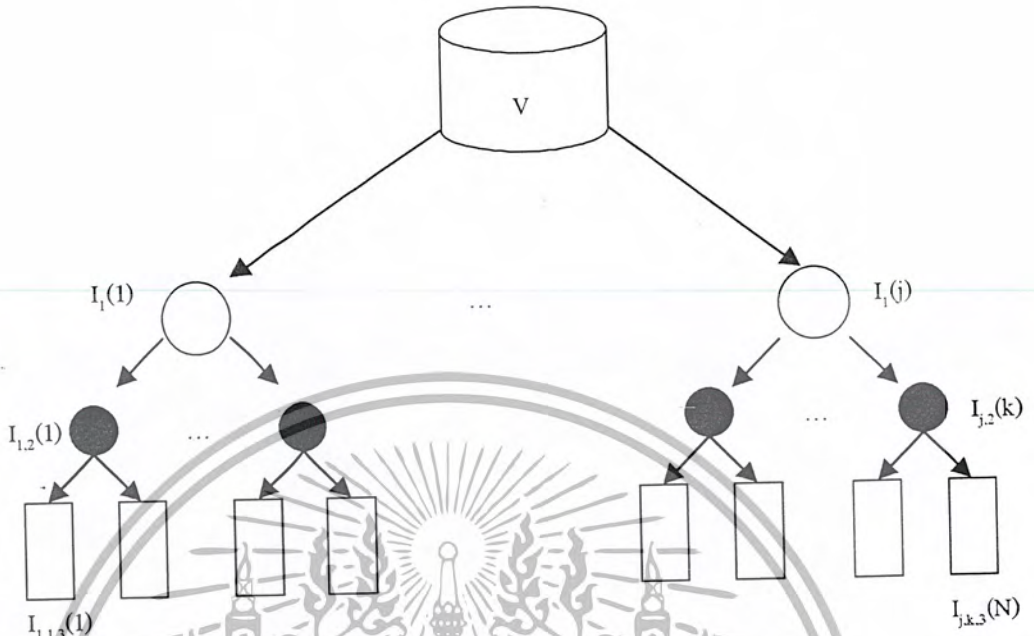
$I_{j,2}(k)$ เป็นครรรชนีตัวที่ 2 ที่แสดงกลุ่มคำศัพท์ย่อยกลุ่มที่ k ภายในกลุ่มใหญ่ที่ j ($k = 1, 2, 3, \dots$)

$I_{j,3}(N)$ เป็นครรรชนีตัวที่ 3 ที่แสดงกลุ่มคำศัพท์ย่อยกลุ่มที่ N โดยที่คำศัพท์นั้นอยู่ในกลุ่มย่อยที่ k และกลุ่มใหญ่ที่ j



รูปที่ 5-3 แสดงโครงสร้างของพจนานุกรมไทย (ในรูปแบบ chart)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 5-4 แสดงโครงสร้างของพจนานุกรมไทย (ในรูปแบบ tree)

5.2.2 การแยกหน่วยคำออกจากประโยค

อัลกอริทึมของการแยกหน่วยคำมีหลักการการทำงานดังนี้

ตัวแปรที่ใช้ในการตัดคำ

len = จำนวนอักขระที่ใช้เปรียบเทียบ

sen = ประโยคที่นำมาตัดคำ

buck = array ที่เก็บตำแหน่งของคำที่ตัดได้

ตัวแปรที่ใช้ชี้ข้อมูล

x = เป็นตัวชี้ตำแหน่งเริ่มต้นของคำที่จะตัด

y = เป็นตัวชี้ตำแหน่งเดิมของ x (หากตัดคำแล้ว x จะเลื่อนตำแหน่งไป หากไม่เลื่อน (y = x) แสดงว่าภายในคำนั้นไม่มีอยู่ในพจนานุกรม)

ตัวแปรที่ใช้เก็บข้อมูล

len = ตัวแปรที่เก็บความยาวของคำที่จะตัด

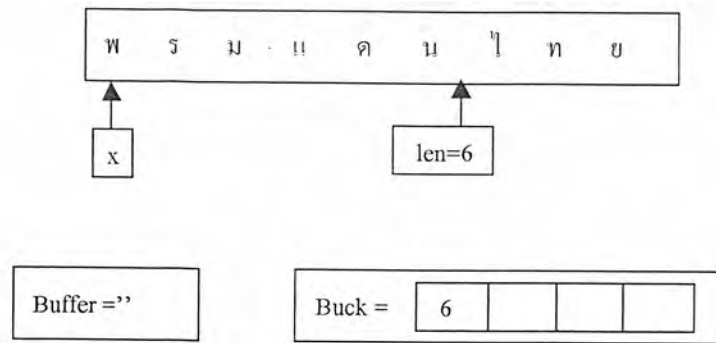
sen = array ที่เก็บประโยคที่จะนำมาตัดคำ

buck = array ที่เก็บตำแหน่งของคำที่ตัดได้

back = ตัวแปรที่บอกว่ามีการย้อนไปตัดคำก่อนหน้าไหม

buffer = ตัวแปรสำหรับเก็บจำนวนอักขระที่ตัดคำไม่ได้

เอกสารนี้เป็นเอกสารที่สงวนลิขสิทธิ์สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



- 1 ตรวจสอบว่ามีประโยคเข้ามาหรือไม่ หากมี ให้ตัวแปร `sen` รับค่าประโยคที่เข้ามาและกำหนดค่าเริ่มต้นให้กับ `x` ให้เท่ากับ 0 เพื่อชี้ที่ตำแหน่งแรกจากนั้นก็ทำตามขั้นตอนต่อไป หากไม่มีก็ออกจากโปรแกรม
- 2 ให้ `len` เท่ากับความยาวที่มากที่สุดของคำศัพท์ที่ขึ้นต้นด้วย `sen[x]` (ในตอนแรก `x = 0` ดังนั้น `sen[0]` ก็คืออักขระตัวแรกของประโยค) แต่หาก `len` มากกว่าความยาวทั้งหมดของประโยค ให้ `len` เท่ากับความยาวทั้งหมดของประโยค
- 3 ค้นหาคำศัพท์ในพจนานุกรมตั้งแต่อักขระตัวที่ `x` จนถึง `len` หากไม่พบในพจนานุกรม จะลด `len` ลงไปที่ละ 1 จนกระทั่งพบคำศัพท์นั้นหรือ `len = 0` จึงหยุดการเปรียบเทียบ
- 4 หากในขั้นตอนที่ 3 สามารถตัดคำได้ ให้นำค่า `len+x` เก็บลงใน `buck` แล้วไปขั้นตอนที่ 7 แต่หากตัดคำไม่ได้ (`x=y`) ให้ไปขั้นตอนต่อไป
- 5 ตรวจสอบว่าค่า `buffer` ไม่เท่ากับ 0 (นั่นคืออักขระตัวก่อนหน้าก็ไม่สามารถตัดเป็นคำได้) หรือ มี `pause` เท่ากับ 1 (ไม่สามารถย้อนกลับไปได้คำก่อนหน้าได้) หรือไม่ หากตรงกับเงื่อนไขใดเงื่อนไขหนึ่งให้เลื่อน `x` ไป 1 ตำแหน่งและให้ `buffer` เพิ่มขึ้นอีก 1 เป็นการข้ามอักขระตัวนั้น แล้วไปขั้นตอนที่ 7 หากไม่ตรงกับเงื่อนไขใดไปขั้นตอนต่อไป
- 6 ตรวจสอบว่าคำก่อนหน้าสามารถตัดได้สั้นกว่านี้หรือไม่ โดยให้ `cur` เก็บค่าตำแหน่งเดิมก่อนตัด หากคำก่อนหน้าไม่สามารถตัดได้สั้นกว่านี้ ให้เลื่อน `x` ไป 1 ตำแหน่ง และให้ `buffer` เพิ่มค่าไปอีกหนึ่ง และให้ค่า `pause = 1` เพื่อเป็นการบอกว่าไม่สามารถย้อนกลับไปได้ แต่หากว่าคำก่อนหน้าสามารถตัดได้เป็นคำสั้น ๆ 2 คำ แต่สุดท้ายก็ยังคงอยู่ที่ตำแหน่งเดิม (`x = cur`) เช่น คำว่า "ระหว่าง" สามารถแยกได้เป็น "ระ" และ "วาง" ก็ให้เลื่อนค่า `x` ไป 1 ตำแหน่งเพื่อข้ามอักขระนั้นไปและให้ `buck` เก็บค่าตำแหน่งเดิม
- 7 ให้ `y = x` เพื่อใช้ในการตรวจสอบต่อไป
- 8 ตรวจสอบว่าตำแหน่ง `x` เป็นตำแหน่งสุดท้ายของประโยคหรือไม่ ถ้าใช่ไปขั้นตอนต่อไป ถ้าไม่ใช่กลับไปขั้นตอนที่ 1 ใหม่เพื่อหาคำที่จะตัดต่อไป

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

- 9 ตรวจสอบว่า buffer มีค่าอยู่หรือไม่หากมีให้นำตำแหน่งสุดท้ายมาเก็บไว้ใน buck จากนั้น วนลูปนำค่าตำแหน่งต่าง ๆ ที่เก็บไว้ใน buck ขึ้นมาเพื่อแสดงค่าที่ตัดได้ จนหมดพร้อมทั้ง clear ค่าใน array จากนั้นไปขั้นตอนที่ 1 ใหม่



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

5.3 การตัดคำโดยใช้คลังข้อความ

จากที่ได้กล่าวมาแล้วในหัวข้อที่ 2.3.3.1 จะได้สมการในการตัดคำเป็นดังสมการที่ 2-3 คือ

$$\arg \max_{W_i} P(W_i | C_{1,m}) = \arg \max_{W_i} \sum_{T_i} \pi P(W_i | T_i) * P(T_i | T_{i-1}, T_{i+1})$$

หากเราทำการคำนวณตามสมการที่ 2-3 โดยตรง จะทำให้การกำกับหน้าที่คำทำได้ช้ามาก เนื่องจากต้องเสียเวลาในการคำนวณทุกรูปแบบที่เป็นไปได้ เช่น ถ้านำสายอักขระที่มีจำนวนคำทั้งหมด M คำ และมีจำนวนหมวดคำทั้งหมด N หมวด หากคิดในกรณี worst case คือ ทุกคำในสายอักขระสามารถเป็นได้ทั้ง N หมวด ดังนั้นจะต้องเสียเวลาในการคำนวณประมาณ $k * N^M$ ครั้งต่อสายคำเดียว โดย k คือค่าคงที่ (ไพศาล เจริญพรสวัสดิ์ [4])



จะเห็นว่าวิธีการนี้ใช้เวลาค่อนข้างนานมาก ซึ่งเวลาที่ใช้จะขึ้นอยู่กับจำนวนคำและจำนวนหมวดคำที่เป็นไปได้ของแต่ละคำในสายนั้น ซึ่งทำให้เวลาที่ใช้เป็นสัดส่วนแบบเอกซ์โปเนนเชียล (Exponential) ดังนั้นจึงได้นำเทคนิคเรื่องไดนามิกโปรแกรมมิ่งเข้ามาช่วย ซึ่งเทคนิคที่เลือกใช้นี้มีชื่อว่า “ขั้นตอนวิธีวิเทอร์บี (Viterbi Algorithm)” ซึ่งเป็นเทคนิคที่นิยมนำมาใช้กับแบบจำลองไครแกรม เมื่อนำมาประยุกต์กับการกำกับคำด้วยไครแกรมสามารถแสดงได้ดังรูปที่ 5-6

กำหนดให้

- w_1, w_2, \dots, w_m เป็นลำดับในประโยค
- M คือ จำนวนคำในประโยคที่นำมากำกับหน้าที่คำ
- t_1, t_2, \dots, t_n เป็นหน้าที่คำทั้งหมดที่เป็นไปได้
- N คือ จำนวนหน้าที่คำที่เป็นไปได้ทั้งหมด

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

- prob(w_i/t_i) คือค่าความน่าจะเป็นของการเกิดคำศัพท์ w_i ที่มีหน้าที่คำเป็น t_i หาได้โดยการนับจำนวนคำ w_i ที่มีหน้าที่คำเป็น t_i หารด้วย จำนวนหน้าที่คำ t_i ทั้งหมดในคลังข้อความ
- prob($t_i/t_{i-1}, t_{i-2}$) คือค่าความน่าจะเป็นของการเกิดลำดับของหน้าที่คำเป็น t_{i-2}, t_{i-1} และตามด้วย t_i หาได้โดยการนับจำนวนของหน้าที่คำที่มีลำดับเป็น t_{i-2}, t_{i-1}, t_i ทั้งหมดในคลังข้อความ หารด้วย จำนวนของหน้าที่คำที่มีลำดับเป็น t_{i-2}, t_{i-1} ทั้งหมดในคลังข้อความ

จะทำการสร้าง Array ขนาด $N \times N \times M$ จำนวน 2 Array คือ

- seqscore [i][j][r] ทำการเก็บค่าความน่าจะเป็นที่ดีที่สุดของการกำกับหน้าที่คำของ w_1, w_2, \dots, w_r ซึ่งคำที่ $r-1$ และ r มีการกำกับหน้าที่เป็น t_i และ t_j ตามลำดับ
- backptr [i][j][r] จะเก็บหน้าที่คำของคำ $r-2$ เมื่อคำที่ $r-1$ และ r มีหน้าที่คำเป็น t_i, t_j ตามลำดับ

หลังจากนั้นจะทำการหาลำดับของหน้าที่คำ C_1, C_2, \dots, C_M ที่เป็นลำดับของคำในประโยคที่มีความน่าจะเป็นมากที่สุด

Initialization step

For $i=1$ to N do

For $j=1$ to N do

seqscore[i][j][2] = prob(w_1/t_i)*prob(t_1/\emptyset)*prob(w_2/t_j)*prob(t_2/t_1)

backptr[1][j][2] = 0

Iteration step

For $r=3$ to M do

For $j=1$ to N do

For $k=1$ to N do

seqscore[j][k][r] = max _{$i=1, \dots, N$} (seqscore[i][j][r-1] * prob($t_k/t_i, t_i$)) * prob(w_r/t_k)

backptr[j][k][r] = ค่า i ที่ทำให้ค่าสมการที่ผ่านมาเป็นค่าที่มากที่สุด

Sequence identification step

$C[M] = k$ and $c[M-1] = j$ โดยที่ j และ k นั้นที่ทำให้ Seqscore[j][k][M] มีค่ามากที่สุด

For $i= M-2$ to 1 do

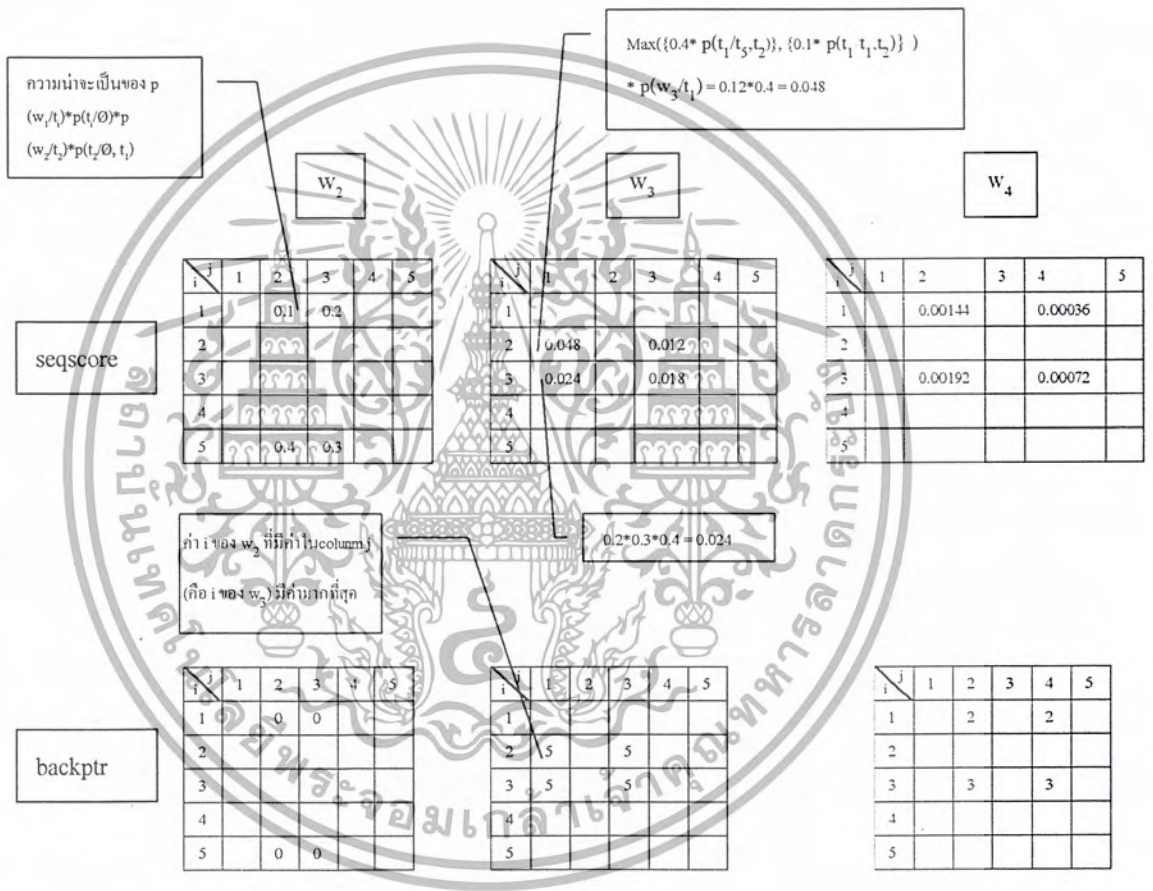
$C[i] = \text{backptr}[C[i+1]][C[i+2]][i+2]$

รูปที่ 5-6 ขั้นตอนการนำวิธีวิเทอร์บิมาใช้ในการตัดคำภาษาไทย

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตัวอย่าง สมมติว่ามีหน้าที่คำทั้งหมด 5 หมวดคำ และมีคำในประโยคทั้งหมด 4 คำ ซึ่งแต่ละคำมีการกำกับหมวดคำตามภาพที่แสดงด้านล่าง

w_1	w_2	w_3	w_4	$P(t_1/t_1, t_2) = 0.2$ $P(t_1/t_1, t_3) = 0.3$ $P(t_1/t_3, t_2) = 0.3$ $P(t_1/t_3, t_3) = 0.2$	$P(t_3/t_1, t_2) = 0.6$ $P(t_3/t_1, t_3) = 0.1$ $P(t_3/t_3, t_2) = 0.1$ $P(t_3/t_3, t_3) = 0.2$	$P(t_2/t_2, t_1) = 0.3$ $P(t_2/t_2, t_3) = 0.4$ $P(t_2/t_3, t_1) = 0.1$ $P(t_2/t_3, t_3) = 0.2$	$P(t_4/t_2, t_1) = 0.4$ $P(t_4/t_2, t_3) = 0.3$ $P(t_4/t_3, t_1) = 0.1$ $P(t_4/t_3, t_3) = 0.2$	$P(w_3/t_1) = 0.4$ $P(w_3/t_3) = 0.3$ $P(w_4/t_2) = 0.1$ $P(w_4/t_3) = 0.2$
t_1	t_2	t_1	t_2					
t_3	t_3	t_3	t_4					



รูปที่ 5-7 แสดงตัวอย่างการทำวิเทอริบีร่วมกับคำตัดภาษาไทย

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Initialization step

1. ทำการกำหนดค่าเริ่มต้นใน array ของ w_2 โดยให้มีค่าเท่ากับ

$$\text{seqscore}[i][j][2] = \text{prob}(w_1/t_1) * \text{prob}(t_1/o) * \text{prob}(w_2/t_1) * \text{prob}(t_1/t_1, o)$$

โดย $\text{prob}(w_1/t_1) * \text{prob}(t_1/o)$ หมายถึง ค่าความน่าจะเป็นที่คำ w_1 มีหมวดคำเป็น t_1 และความน่าจะเป็นที่หมวดคำ t_1 อยู่ต้นประโยค และ $\text{prob}(w_2/t_1) * \text{prob}(t_1/o, t_1)$ หมายถึงความน่าจะเป็นที่ w_2 มีหมวดคำเป็น t_1 และความน่าจะเป็นที่ t_1 ตามหลัง t_1 ที่อยู่ต้นประโยค

2. ใส่ค่า $\text{backptr}[i][j][2] = 0$ เพื่อให้รู้ว่าเป็นจุดเริ่มต้น

Iteration step

จะทำการคำนวณไครแกรม จากตัวอย่าง ทำการพิจารณาต่อที่ w_3 ซึ่งมีหมวดคำที่เป็นไปได้ 2 หมวดคือ t_1 และ t_2 โดยทำการพิจารณาที่หมวดคำ t_1 ก่อนซึ่งมีรูปแบบของหน้าที่คำมายัง t_1 ได้ 4 รูปแบบคือ



เราไม่รู้ว่าจะเส้นทางใดดีที่สุดระหว่าง $t_1, t_2, t_1, t_1, t_2, t_1$ และ t_2, t_2, t_1 ดังนั้นเราจึงดูความน่าจะเป็นของหมวดคำก่อนหน้าที่จะมาถึง w_3 ว่าค่าใดมีค่ามากที่สุด โดยที่เรารู้ค่า j ของ array w_2 ซึ่งก็คือหมวดคำก่อนหน้า w_3 (ค่า i ใน array w_3) จึงทำการหา i ที่ทำให้ $\text{column } j$ ของ array w_2 ที่ทำให้มีค่า $\text{seqscore}[i][j][2] * \text{prob}(t_k/t_i)$ มากที่สุด ถ้าเราพิจารณาหมวดคำ t_2 ของ w_2 (ค่า $i=2$ ใน array ของ w_3 แต่เป็นค่า $j=2$ ใน array ของ w_2) จะเห็นได้ว่า ในคอลัมน์ที่ 2 (ค่า j) ของ array w_2 มีค่า 0.1 ที่ $i=1$ และมีค่า 0.4 ที่ $i=5$ ดังนั้นเราจะต้องหาว่า $\text{seqscore}[i][2][2] * \text{prob}(t_k/t_i)$ ที่มากที่สุดระหว่างค่า i ทั้งสองของ w_3 ดังนั้นสุดท้ายแล้วจะได้ค่าของ $\text{Backptr}[2][1][3] = 5$ เมื่อทำงานจบประโยคแล้วจะทำการ Backtracking โดยใช้ Backptr เช่น $\text{Backptr}[2][1][3] = 5$ (ซึ่งหากค่านี้เป็นเส้นทางที่มีค่าความน่าจะเป็นมากที่สุด) จะหมายถึง คำที่ 3 มีหน้าที่คำเป็น 1 คำที่ 2 มีหน้าที่คำเป็น 2 และ คำที่ 1 มีหน้าที่คำเป็น 5 ซึ่งค่า C จะใช้ในการเก็บเส้นทางที่ได้จากการ Backtrack

บทที่ 6

ขั้นตอนในการวัดประสิทธิภาพ

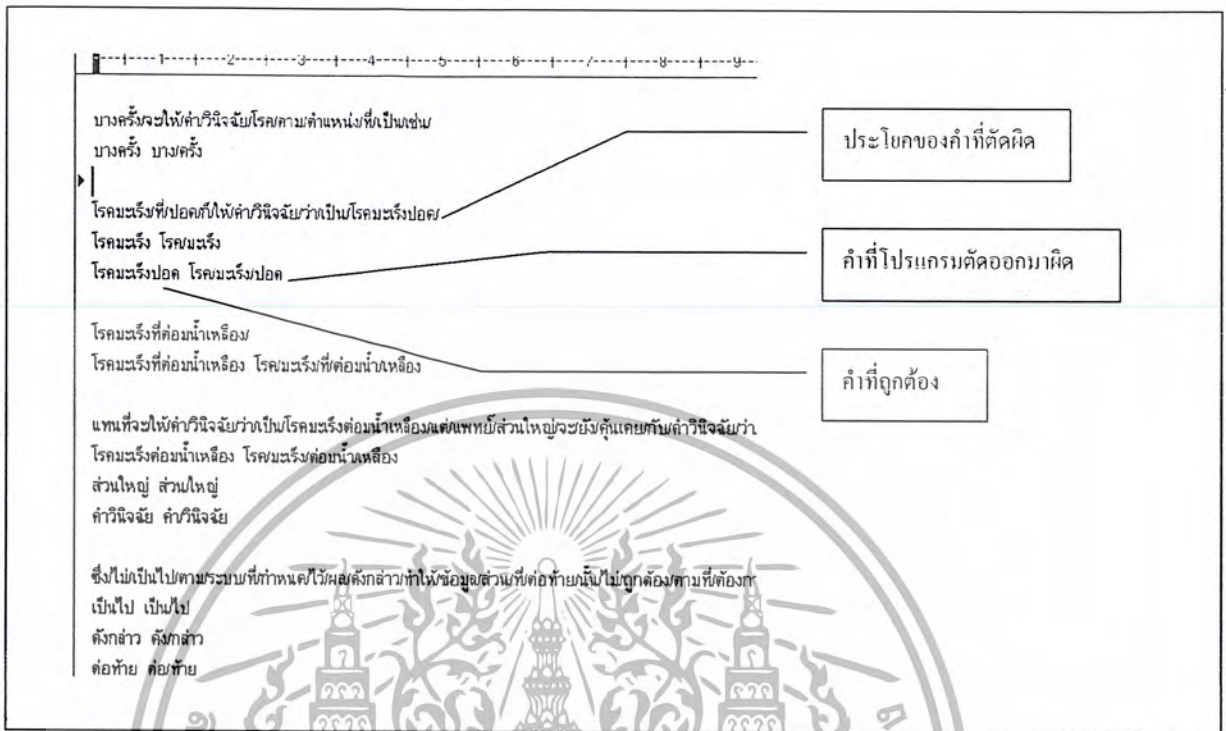
6.1 สร้างค่าของอินพุตและค่าตอบในการเปรียบเทียบ

อินพุตของการวัดประสิทธิภาพนั้นเราได้นำมาจาก 20% ของคลังข้อความซึ่งมีประมาณ 4881 ประโยค โดยที่ในแต่ละประโยคจะมีทั้งตัวอักษรภาษาไทย ภาษาอังกฤษ เครื่องหมายพิเศษ เครื่องหมายวันวรรค ซึ่งในส่วนนี้จะนำประโยคมาสร้างเป็นไฟล์ที่ใช้เป็นอินพุตในการทดสอบโปรแกรมตัดคำต่าง ๆ ส่วนคำที่ตัดไว้ในคลังข้อความของประโยคนั้นๆ จะนำมาเป็นไฟล์ที่ใช้เป็นค่าตอบในการเปรียบเทียบความถูกต้อง โดยจะถือค่าตอบในไฟล์นี้เป็นค่าตอบที่ถูกต้อง 100 %

6.2 วัดประสิทธิภาพด้วย Framework ที่สร้างขึ้น

การทำงานของ framework สามารถแบ่งได้เป็นขั้นตอนดังนี้

1. รับ path และ ชื่อของโปรแกรมที่ต้องการวัดประสิทธิภาพ นำมาเก็บไว้ใน name.txt
2. อ่านชื่อของ โปรแกรม ใน file name.txt ขึ้นมาเพื่อ run และเริ่มการจับเวลาเพื่อวัดประสิทธิภาพด้านเวลา
3. ทำการหาจำนวน byte ทั้งหมดที่อยู่ในmemory เพื่อวัดประสิทธิภาพด้านการใช้ทรัพยากร
4. นำเอาผลลัพธ์ที่ได้จากโปรแกรมการตัดคำไปทำการเปรียบเทียบกับไฟล์ที่ใช้เป็นค่าตอบที่กล่าวไว้ด้านบน จากนั้นคำนวณหาค่าความถูกต้องออกมา ตามสมการที่ 3-2 โดยมีกราฟแสดงจำนวนคำที่ถูกจำนวนคำทั้งหมดที่ตัดได้จากโปรแกรม จำนวนคำทั้งหมดที่ตัดไว้ในคลังข้อความ
5. รายละเอียดของประโยคที่ตัดคตินั้นจะอยู่ในไฟล์ detail <ชื่อโปรแกรมที่นำมาตัดคำ>.txt



รูปที่ 6-1 แสดงตัวอย่างของไฟล์ detail_dict.txt

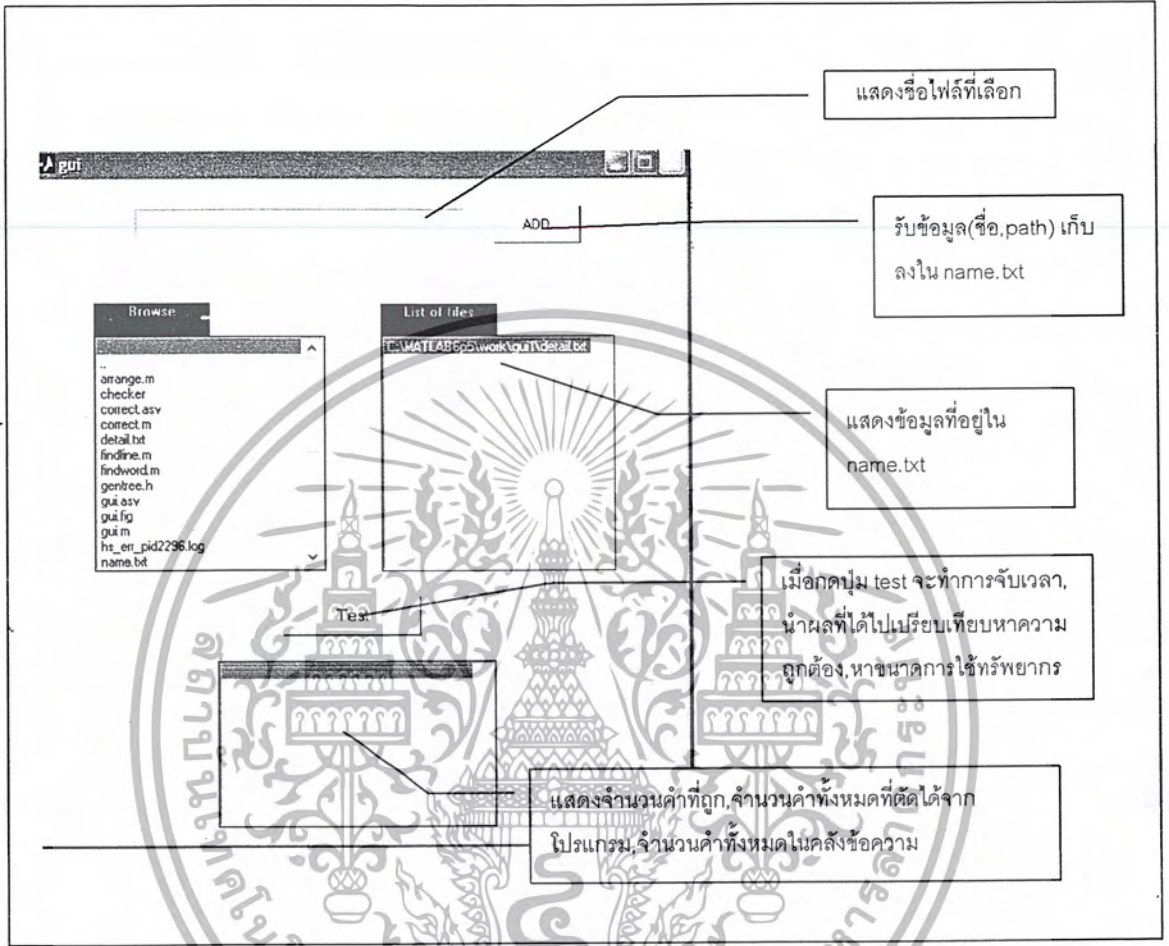
6.2.1 เงื่อนไข

1. ไฟล์ที่เป็นเอาต์พุต (output) ของ โปรแกรมที่คัดค้านั้นของคือเป็นชื่อ "out_ชื่อไฟล์ที่นำวัดประสิทธิภาพ" เช่น มีไฟล์ dict.dll ซึ่งเป็น โปรแกรมการคัดคำโดยใช้พจนานุกรม จะต้องมีการเขียนผลการคัดคำลงในไฟล์ชื่อ out_dict.txt เท่านั้น
2. รายละเอียดของคำที่ตัด ได้ผัดนั้นจะอยู่ในไฟล์ชื่อ "detail_ชื่อไฟล์ที่นำวัดประสิทธิภาพ" เช่น ในตัวอย่างของข้อที่1 เราจะ ได้ไฟล์ที่เก็บ รายละเอียดของคำที่ตัดผิดชื่อ detail_dict.txt
3. ไฟล์ที่จะนำมาเป็นอินพุตในการคัดค้านั้นให้มีชื่อว่า input.txt

6.2.2 ขั้นตอนการใช้ framework

1. เลือกชื่อไฟล์ที่ต้องการ วัดประสิทธิภาพใน listbox ด้านซ้าย ชื่อ ไฟล์ที่ต้องการจะปรากฏใน edit box ด้านบน
2. กดปุ่ม ADD เพื่อเพิ่มไฟล์ที่ต้องการ วัดประสิทธิภาพลงใน name.txt
3. กดปุ่ม Test เพื่อวัดประสิทธิภาพ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 6-2 แสดงการทำงานในส่วนต่างๆ ของ framework

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 7

ผลการทดลอง

ในบทนี้ ได้ทำการรวบรวมและเปรียบเทียบข้อมูลสถิติ ที่เกิดจากการทดลองวัดประสิทธิภาพของ โปรแกรมตัดคำภาษาไทยจาก framework ที่ได้สร้างขึ้น โดยที่ framework จะทำการคำนวณและสร้างกราฟ ต่าง ๆ ให้โดยอัตโนมัติ ซึ่งมีผลการทดลองดังนี้

7.1 เปรียบเทียบด้านความถูกต้องของคำ

จากการทดลองทำให้ได้ผลจำนวนคำที่โปรแกรมตัดคำแบบต่าง ๆ ตัดคำออกมาได้ดังตารางที่ 7-1



รูปที่ 7-1 แสดงการเปรียบเทียบโดยค่า F - Measure

มาตราวัดประสิทธิภาพ โปรแกรมการตัดคำ	ความถูกต้อง (f-measure)	จำนวนคำที่ตัดได้ถูก ต้อง	จำนวนคำที่ตัดได้ ทั้งหมด
การตัดคำด้วยคลังข้อความ	0.63064	38859	71084
การตัดคำด้วยพจนานุกรม	0.66231	39823	68103
การตัดคำด้วยกฎ	0.33803	22297	79772

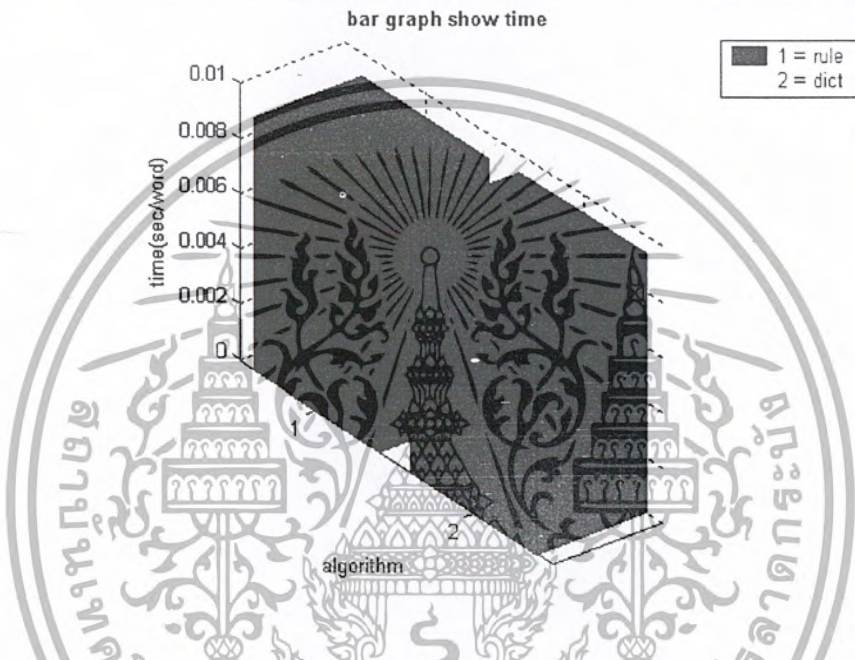
เอกสารนี้เป็นเอกสารที่สงวนลิขสิทธิ์ไว้สำหรับใช้ในเพื่อการวิจัยเท่านั้น เมื่อผู้ยูได้เห็น ใบนี้ขอประ โยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

7.2 เปรียบเทียบประสิทธิภาพเชิงความเร็ว

จากการทดลองทำให้ได้ผลจำนวนคำที่โปรแกรมตัดคำแบบต่าง ๆ ตัดคำออกมาได้ดังตารางที่ 7-2 เนื่องจากการตัดคำลั้งข้อความใช้เวลานานมาก เช่น ในประโยค

“โรคมะเร็งที่ปอดก็ให้คำวินิจฉัยว่าเป็น โรคมะเร็งปอด”

ใช้เวลาตัดประมาณ 3 นาที จึงไม่นำมาเปรียบเทียบ



รูปที่ 7-2 แสดงการเปรียบเทียบความเร็วเฉลี่ย

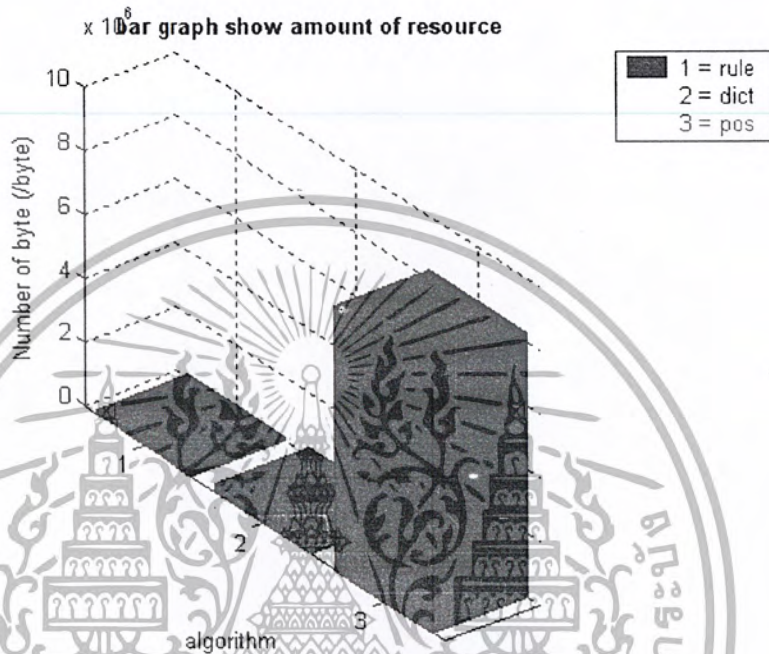
มาตรฐานประสิทธิภาพ	ความเร็ว
โปรแกรมการตัดคำ	
การตัดคำด้วยพจนานุกรม	0.0097224 sec/word
การตัดคำด้วยกฎ	0.0077636 sec/word

ตารางที่ 7-2 ค่าความเร็วเฉลี่ยในการตัดคำ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

7.3 เปรียบเทียบประสิทธิภาพการใช้ทรัพยากร

การใช้ทรัพยากรของโปรแกรมตัดคำแบบต่างๆมีผลมาจากขนาดของพจนานุกรม และคลังข้อความที่ใช้ งานด้วย ซึ่งขนาดของการใช้ทรัพยากรได้แสดงไว้ดังตารางที่ 7-3



รูปที่ 7-3 แสดงการเปรียบเทียบการใช้ทรัพยากร

มาตราวัดประสิทธิภาพ	ทรัพยากร
โปรแกรมการตัดคำ	
การตัดคำด้วยกฎ	123210 byte
การตัดคำด้วยพจนานุกรม	295027 byte
การตัดคำด้วยคลังข้อความ	8317641 byte

ตารางที่ 7-3 จำนวนการใช้ทรัพยากรของโปรแกรมการตัดคำแบบต่างๆ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 8

การวิเคราะห์ผลจากการตัดคำแล้ว

ผลจากการวัดประสิทธิภาพในด้านต่างๆสามารถวิเคราะห์สาเหตุของการตัดคำผิด โดยแยกตามอัลกอริทึมที่ใช้ ได้ดังต่อไปนี้

8.1 โปรแกรมการตัดคำด้วยพจนานุกรม

สาเหตุที่เกิดการตัดโดยพจนานุกรมมีความผิดพลาด โดยวิเคราะห์ผลลัพธ์ที่ผิดในการตัดคำ มีดังต่อไปนี้

ผ่านวัตถุแล้วไปตกกระทบบนอนุกรมรังสีซึ่งประกอบด้วย

▶ แล้วไป แล้วไป

ตกกระทบบ ตกกระทบบ

ประกอบด้วย ประกอบด้วย

เส้นเซอร์ เส้นเซอร์

จำนวนมาก จำนวนมาก

เส้นตรง เส้นตรง

คำที่ถูก => ตก/กระทบบ

แต่ตัดได้เป็น=>ตกกระทบบ

จากตัวอย่างด้านบนการตัดคำว่า “ตกกระทบบ” ไม่ถูกต้องนั้นเกิดจากการที่การตัดคำโดยใช้พจนานุกรมพบว่าคำว่า “ตกกระทบบ” เป็นคำที่ยาวที่สุดซึ่งทำให้คำข้างหลังสามารถตัดได้ ดังนั้นวิธีการนี้จึงตัดคำได้เป็น “ตกกระทบบ”-“บ” โดยที่ทั้ง 2 คำนี้ไม่มีความหมายเกี่ยวข้องกันทางด้าน ไวยากรณ์เลย ซึ่งคำกำกวมในลักษณะนี้จะทำให้ประสิทธิภาพการตัดคำประเภทนี้ค่อนข้างต่ำ

จะไม่มีส่วนที่ใช้สำหรับอินเทอร์เฟสเข้ากับคอมพิ

อินเทอร์เฟส อินเทอร์เฟส

เข้ากับ เข้ากับ

จำเป็นต้อง จำเป็นต้อง

ออกแบบ ออกแบบ

อินเทอร์เฟส อินเทอร์เฟส

เพื่อให้ เพื่อให้

คำที่ถูก => อินเทอร์เฟส

แต่ตัดได้เป็น=>อิน/เท/อร์/เฟส

การตัดคำโดยใช้พจนานุกรมนั้นหากไม่พบคำในพจนานุกรมจะทำให้ไม่สามารถตัดคำนั้นได้ถูกต้องซึ่งคำที่มักจะ ไม่พบในพจนานุกรมได้แก่ ชื่อคน คำทับศัพท์ภาษาต่างประเทศ ดังตัวอย่างด้านบน

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

กัมที่ละเอียดยตลอดแนวกลมเพื่อเก็บข้อมูลโปรเจกชันของวัตถุให้ครบ

ที่ละเอียย ที่ละเอียย

โปรเจกชันของ โปรเจ/คช/นข/อง

องศาในโครงการวิจัยระยะที่

โครงการวิจัย โครงการวิจัย

คำที่ถูก => โปรเจกชัน/ของ

แต่ตัดได้เป็น=> โปร/ร/เจ/คช/นข/อง

เราได้ทำการสร้างเครื่องต้นแบบเครื่องคอมพิวเตอร์/ถ่าย/ภาพตัดขวางได้สั

เครื่องต้นแบบ เครื่องต้นแบบ

จากตัวอย่างด้านบน สาเหตุอีกอย่างที่ทำให้การตัดคำผิดพลาดก็คือ ผลกระทบจากคำที่ไม่พบใน พจนานุกรม อาจส่งผลถึงคำใกล้เคียง เช่น โปรเจกชัน/ของ ตัดได้เป็น โปร/ร/เจ/คช/นข/อง จะเห็นว่าคำว่าของ จะถูกตัดผิดไปด้วย ซึ่งเป็นผลกระทบจากคำข้างหน้า คือคำว่า “โปรเจกชัน” เนื่องจากคำว่า “โปรเจกชัน” นั้น เป็นคำที่ไม่มีในพจนานุกรม

8.2 โปรแกรมตัดคำด้วยคลังข้อความ

จากผลการทดลอง สังเกตได้ว่าลักษณะการตัดคำของการตัดคำโดยใช้คลังข้อความ ส่วนใหญ่จะมีการผิดพลาดกันเกือบทั้งประโยค ดังภาพที่แสดงไว้ด้านล่าง

ซึ่งไม่เป็นไปตามระบบที่กำหนดไว้ผลจึงกล่าวทำให้อายุของส่วนที่ต่อท้ายนั้นไม่ถูกต้องตามที่ต้องการปัญหาจึงกล่าวสามารถแก้ไขโดยก่อน

เงินไป เงินไป

คิงสาว คิงสาว

ทำให ทำให

ต่อท้าย ต่อท้าย

ถูกต้อง ถูกต้อง

คามที่ คามที่

ต้องการ ต้องการ

คิงสาว คิงสาว

สามารถ สามารถ

แท้ไซ แท้ไซ

จัดทำ จัดทำ

คู่มือ คู่มือ

เอกสารอ้างอิง เอกสารอ้างอิง

อบรม อบรม

ผู้เกี่ยวข้อง ผู้เกี่ยวข้อง

ระดับ ระดับ

ฐานข้อมูล ฐานข้อมูล

หน่วยงาน หน่วยงาน

สาเหตุที่ทำให้เกิดการตัดคำผิดพลาดอาจเกิดจากการตัดคำที่ผิดตั้งแต่คำแรกๆ ซึ่งจะส่งผล ให้คำด้านหลังเกิดการผิดพลาดไปด้วยเนื่องจากการตัดคำประเภทนี้คำแต่ละคำมีการเกี่ยวข้องกัน ทั้งนี้การตัดคำที่ผิดพลาดอาจมีสาเหตุมาจาก ไม่เคยเจอความสัมพันธ์หมวดคำในรูปแบบนั้น ๆ มาก่อน และการกำกับหน้าที่คำที่ผิดพลาด ซึ่งความผิดพลาดนี้อาจส่งผลให้การตัดคำของคำต่อ ๆ ไปจำนวนมากผิดพลาดได้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

8.3 โปรแกรมการตัดคำด้วยกฎ

การตัดคำโดยใช้กฎนั้น เหมาะกับการตัดพยางค์มากกว่ามาตัดคำ เนื่องจากคำแต่ละคำไม่ได้มีรูปแบบตายตัวเหมือนเช่นพยางค์ ซึ่งธรรมชาติของภาษาไทยหากคำ 2 คำสามารถรวมเป็นกลุ่มคำได้ ก็มีความเป็นไปได้สูงที่คำที่ถูกต้องจะเป็นกลุ่มคำที่เกิดจากการรวมคำ 2 คำนั้น (นัฐวดี ไชยเจริญ [12]) ทำให้การใช้วิธีนี้ไม่ค่อยมีประสิทธิภาพมากนัก แต่กฎก็สามารถตัดคำที่เกิดจากการทับศัพท์ภาษาต่างประเทศได้ เช่น คำว่า “ริงส์เอ็กซ์” ตัดได้เป็น “ริง”-“ส์”-“เอ็กซ์”



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 9

บทสรุปและข้อเสนอแนะ

9.1 สรุปประสิทธิภาพของโปรแกรมการตัดคำ

โปรแกรมการตัดคำแบบพจนานุกรม มีความเร็วในการตัดคำมากที่สุด ในขณะที่การตัดคำโดยใช้คลังข้อความมีความเร็วในการตัดคำน้อยที่สุด

ในเรื่องของความถูกต้อง การตัดคำโดยใช้กฎมีความถูกต้องน้อยที่สุด การตัดคำโดยพจนานุกรมสามารถตัดคำได้มีความถูกต้องมากที่สุด

การตัดคำโดยใช้คลังข้อความมีการใช้ทรัพยากรมากกว่าการตัดคำอีก 2 แบบที่เหลือนมาก โดยการตัดคำแบบใช้กฎใช้ทรัพยากรน้อยที่สุด

9.2 ความเหมาะสมของการประยุกต์ใช้โปรแกรมตัดคำ

จากผลวิจัยนี้สามารถทำให้เห็นแนวทางในการพิจารณาในการนำโปรแกรมไปใช้กับงานประเภทต่าง ๆ ให้เหมาะสม สำหรับการพิจารณาในเบื้องต้นมีดังนี้

1. โปรแกรมตัดคำโดยคลังข้อความเหมาะสมกับงานที่ต้องการความถูกต้องด้านไวยากรณ์สูง เช่น งานแปลเอกสาร เนื่องจากมีการใช้น้ำหนักของคำมาร่วมพิจารณาในการตัดคำ ทำให้ผลการตัดคำที่ได้มีความถูกต้องทางด้านไวยากรณ์มากกว่ารูปแบบอื่นๆ

2. โปรแกรมการตัดคำโดยใช้พจนานุกรม สามารถนำมาประยุกต์ใช้ได้ในงาน เนื่องจากทั้งความเร็วและความถูกต้องในการตัดคำมากกว่าแบบอื่นๆ จึงอาจนำไปใช้ในงานประเภทจัดรูปแบบเอกสาร และการตรวจสอบคำสะกด ได้

3. โปรแกรมการตัดคำด้วยกฎ แม้ว่าจะมีความถูกต้องต่ำ แต่มีความสามารถในการตัดคำที่ไม่เคยพบมาก่อน เช่น คำทับศัพท์ภาษาต่างประเทศ ได้ถูกต้องมากกว่าการตัดคำในรูปแบบอื่นๆ ดังนั้นหากนำไปใช้ร่วมกับโปรแกรมการตัดคำแบบอื่นจะสามารถเพิ่มประสิทธิภาพในการตัดคำของโปรแกรมเหล่านั้นเพิ่มมากยิ่งขึ้น

9.3 ประเภทเอกสารกับผลของการเปรียบเทียบประสิทธิภาพการตัดคำ

เนื่องจากเอกสารที่ใช้ทดลองนั้นส่วนมากจะเกี่ยวข้องกับเรื่องของเทคโนโลยี ซึ่งมักจะมีคำที่เกิดจากการทับศัพท์ภาษาต่างประเทศเป็นจำนวนมาก ดังนั้นวิธีการตัดคำแบบต่าง ๆ อาจจะตัดคำออกมาได้ผิดพลาดมาก เนื่องจากเป็นคำที่ไม่รู้จัก ดังนั้นสามารถสรุปได้ว่าประเภทของเอกสารมีผลต่อความถูกต้องของการตัดคำ นั่นคือหากเอกสารมีคำพื้นฐานเป็นส่วนใหญ่ ความถูกต้องของการตัดคำก็จะมากขึ้น ในทางตรงกันข้าม หากเอกสารมีคำที่เกิดจากการทับศัพท์ภาษาต่างประเทศ หรือคำที่เฉพาะเจาะจงของแต่ละสาขาเป็นส่วนใหญ่ ความถูกต้องของการตัดคำก็จะลดลง

9.4 ปัญหาต่าง ๆ ที่พบในงานวิจัย

9.4.1 งานวิจัยได้พบปัญหาเกี่ยวกับการตัดคำโดยใช้คลังข้อความ เนื่องจากเวลาที่ใช้ในการตัดคำนั้น ใช้เวลาค่อนข้างนานมาก ดังนั้นผู้วิจัยจึงไม่ได้เปรียบเทียบ ในด้านเวลา

9.4.2 เนื่องจากเอกสารที่นำมา (คลังข้อความ) มีจุดผิดพลาดในการตัดคำค่อนข้างมาก และเนื้อหาที่ไม่หลากหลาย มักจะเกี่ยวข้องกับเทคโนโลยี ซึ่งในการใช้งานจริงเราจะเจอคำศัพท์ประเภทนี้น้อยมาก ทำให้ประสิทธิภาพการตัดคำในการใช้งานจริงอาจไม่ได้ผลตามที่ได้ทำการทดลอง

9.5 ข้อเสนอแนะและงานวิจัยที่สามารถทำเพิ่มเติมต่อจากงานวิจัยนี้

9.5.1 งานในด้านการพัฒนาคลังข้อความก็เป็นสิ่งที่จะต้องพัฒนาเพิ่มเนื่องจากเนื้อหาสาระที่มีนั้นค่อนข้างที่จะเกี่ยวกับเรื่องของเทคโนโลยีเป็นส่วนใหญ่ ดังนั้นอาจมีการพัฒนาเพิ่มข้อมูลลงในคลังข้อความในเนื้อหาสาระเรื่องอื่น ๆ ก็เป็นสิ่งที่จะต้องพัฒนาต่อ

9.5.2 ศึกษาการใช้งานการตัดคำแต่ละแบบ เพื่อความเหมาะสมในการเลือกใช้งาน

9.5.3 การพัฒนารูปแบบในการตัดคำใหม่ ๆ เพื่อความถูกต้องมากยิ่งขึ้น

9.5.4 ควรเพิ่มเติมคำย่อ คำเฉพาะที่พบบ่อย ลงในพจนานุกรม เพื่อความถูกต้องมากยิ่งขึ้น

ภาคผนวก ก

ความหมายของคำย่อและสัญลักษณ์ต่าง ๆ ในคลังข้อความ Orchid

Mark-up	Description
%TTitle:	Title of the document written in Thai.
%ETitle:	Title of the document written in English.
%TAuthor:	Author's name written in Thai.
%EAuthor:	Author's name written in English.
%TInbook:	Title of the book where the document exists, written in Thai.
%EInbook:	Title of the book where the document exists, written in English.
%TPublisher:	Publisher of the book, written in Thai.
%EPublisher:	Publisher of the book, written in English.
%Page:	Page number or the range of pages of the document.
%Year:	Published year (A.D.).
%File:	File number of the document. A long document may be separated into a number of files.

Mark-up	Description
#P[number]	Paragraph number of a text. The number in the bracket is shown in a sequence within a text.
#[number]	Sentence number of a paragraph. The number in the bracket is shown in a sequence within a paragraph.

Mark-up	Description
\	Line break symbol.
//	Sentence break symbol.
/[POS]	Tag marker for the appropriate POS of a word.

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

อักขระพิเศษ

Special characters	Defined strings	Special characters	Defined strings
	<space>	/	<slash>
!	<exclamation>	:	<colon>
"	<quotation>	;	<semi_colon>
#	<number>	<	<less_than>
\$	<dollar>	=	<equal>
%	<percent>	>	<greater_than>
&	<ampersand>	?	<question_mark>
'	<apostrophe>	@	<at_mark>
(<left_parenthesis>	[<left_square_bracket>
)	<right_parenthesis>]	<right_square_bracket>
*	<asterisk>	^	<circumflex_accent>
+	<plus>	_	<low_line>
,	<comma>	{	<left_curly_bracket>
-	<minus>	}	<right_curly_bracket>
.	<full_stop>	~	<tilda>

ตารางความหมายของคำย่อของหมวดคำต่าง ๆ

No.	POS	Description	Example
1	NPRP	Proper noun	วันโตวส์ 95, โดโรน่า, โลก, พระอาทิตย์
2	NCNM	Cardinal number	หนึ่ง, สอง, สาม, 1, 2, 3
3	NONM	Ordinal number	ที่หนึ่ง, ที่สอง, ที่สาม, ที่ 1, ที่ 2, ที่ 3
4	NLBL	Label noun	1, 2, 3, 4, ก, ข, a, b
5	NCMN	Common noun	หนังสือ, อาหาร, อาคาร, คน
6	NTIL	Title noun	ดร., พลเอก
7	PPRS	Personal pronoun	คุณ, เขา, ฉัน
8	PDMN	Demonstrative pronoun	นี้, นั้น, ที่นั่น, ที่นี้
9	PNTR	Interrogative pronoun	ใคร, อะไร, อย่างไร
10	PREL	Relative pronoun	ที่, ซึ่ง, อัน, ผู้
11	VACT	Active verb	ทำงาน, ร้องเพลง, กิน
12	VSTA	Stative verb	เห็น, รู้, คือ
13	VATT	Attributive verb	ฮ้วน, ดี, สวย
14	XVBM	Pre-verb auxiliary, before negator “ไม่”	เกิด, เกือบ, กำลัง
15	XVAM	Pre-verb auxiliary, after negator “ไม่”	ค่อย, นำ, ได้
16	XVMM	Pre-verb, before or after negator “ไม่”	ควร, เคย, ต้อง
17	XVBB	Pre-verb auxiliary, in imperative mood	กรุณา, จง, เชิญ, อย่า, ห้าม

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

18	XVAE	Post-verb auxiliary	ไป, มา, ขึ้น
19	DDAN	Definite determiner, after noun without classifier in between	นี้, นั้น, โน่น, ทั้งหมด
20	DDAC	Definite determiner, allowing classifier in between	นี้, นั้น, โน่น, บูน
21	DDBQ	Definite determiner, between noun and classifier or preceding quantitative expression	ทั้ง, อีก, เพียง
22	DDAQ	Definite determiner, following quantitative expression	พอดี, ถ้วน
23	DIAC	Indefinite determiner, following noun; allowing classifier in between	ไหน, อื่น, ต่างๆ
24	DIBQ	Indefinite determiner, between noun and classifier or preceding quantitative expression	บาง, ประมาณ, เกือบ
25	DIAQ	Indefinite determiner, following quantitative expression	กว่า, เสน
26	DCNM	Determiner, cardinal number expression	หนึ่งคน, เลือ 2 คำ
27	DONM	Determiner, ordinal number expression	ที่หนึ่ง, ที่สอง, ที่สุดท้าย
28	ADVN	Adverb with normal form	คง, ชั่ว, ช้า, สม่ำเสมอ
29	ADVI	Adverb with iterative form	เร็วๆ, หลายๆ, ซ้ำๆ
30	ADVP	Adverb with prefixed form	โดยเร็ว
31	ADVS	Sentential adverb	โดยปกติ, ธรรมดา
32	CNIT	Unit classifier	ตัว, อัน, เล่ม
33	CLTV	Collective classifier	คู่, กลุ่ม, ฝูง, เซ็น, ทาง, ด้าน, แบบ, รุ่น
34	CMTR	Measurement classifier	ดีใจครบ, แก้ว, ชั่วโมง
35	CFQC	Frequency classifier	ครั้ง, เพียง
36	CVBL	Verbal classifier	ม้วน, มัด
37	JCRG	Coordinating conjunction	และ, หรือ, แต่
38	JCMP	Comparative conjunction	กว่า, เหมือนกับ, เท่ากับ
39	JSBR	Subordinating conjunction	เพราะว่า, เนื่องจาก, ที่, แม้ว่า, ถ้า
40	RPRE	Preposition	จาก, ละ, ของ, ได้, บน
41	INT	Interjection	โอ๊ย, โอ้, เออ, เอ้, อ้อ
42	FIXN	Nominal prefix	การทำงาน, ความสนุกสนาน
43	FIXV	Adverbial prefix	อย่างรวดเร็ว
44	EAFF	Ending for affirmative sentence	จ๊ะ, จ๊ะ, ค่ะ, ครับ, นะ, ná, เอะ
45	EITT	Ending for interrogative sentence	หรือ, หรือ, ไหม, มั้ย
46	NEG	Negator	ไม่, ไม่ได้, ไม่ได้, มิ
47	PUNC	Punctuation	(,), “, ”, ;

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ภาคผนวก ข

กฎของคุณสุรินทร์ จรรยาพรพงษ์

+-----+

ความหมายของคำย่อที่ใช้

+-----+

- c คือ อักษรนำ
s คือ ศัพท์สะกด
t คือ วรรณยุกต์
g คือ อักษรที่มีการันต์ต่อท้ายได้
* คือ ขอบเขตของพยางค์
.. คือ ศัพท์อักษรใดใด

+-----+

รูปแบบของคุณลักษณะ

+-----+

- R คือ รูปแบบกฎที่ตายตัว
N คือ รูปแบบกฎที่อาจควได้ถูกหรือไม่ถูกต้องก็ได้
X คือ กรณียกเว้น

+-----+



PATTERN C-1

รูปแบบสระ	ขอบเขตของพยางค์	ตัวอย่าง	รูปแบบ คุณลักษณะ
- 1			
ขอบเขตหลัง	l s g g [◌] *	ศาสตร์ พราหมณ์	R
	l s g [◌] *	การ์ณ ปราชญ์ กาศย์	R
	l g [◌] *	การ์ บาร์	R
	l g [◌] s *	การ์ศ ฟาร์ม	X
	l s [◌] *	ชาติ ญาติ พยาธิ	X
	l s [◌] *	ชาติ	X
	l s s [◌] *	มาตร มารธ มารศ	X,N
	l s [◌] *คคค	การ มากคคปาก	N
	l s [◌] *คคค	ข้าม บ้าน	N
	l *	กา มา ขา	N
	t l *	ข้า ท้า	N
ขอบเขตหน้า	* c l	กา บาร์	N
	* c c l	มหา สถาน	N
	* c t l	น้ำ ย้าย	N
	* c c l l	หน้าย คล้าย	N

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางแสดงอักษรควบของรูปสระ : - ๑

พยัญชนะตัวที่ 2 (c2)	พยัญชนะตัวแรก (c1)
น	อท ห ส ข ฐ
ร	ก อ ท ค ส ป ค พ ข ฉ ม
อ	ส
ง	ห ส ผ
ม	ห ส
ย	อ ห ส พ ข
ว	ค ท ห ส ค ข ผ ถ
ล	ก อ ท ค ห ส ป ค พ ข ผ ถ ฉ
ท	ป
ค	ส
ห	ม ท ส
บ	ส ข
ป	ส
น	ส
ญ	ห
ฎ	ม ล
ฏ	ข



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

PATTERN C-2

รูปแบบสระ	ขอบเขตของพยางค์	ตัวอย่าง	รูปแบบ คุณลักษณะ
ขอบเขตหลัง	s g g *	รินทร์ มินทร์	R
	s g *	สิทธิ์	R
	s g *	สิงห์ คิชูร์	R
	s s *	จิตร มิตร	X,N
	s *	คิด จิต ขจิต	N
	*	จิ คี ปฎิ	N
	l s *	สิง หิง ทิม	R
ขอบเขตหน้า	* c	คิด จิบ	N
	* c c	หยัง ปฎิ	N
	* c c c	อมรินทร์ สปริง อคคี	X,N

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางแสดงอักษรควบของรูปสระ : -

พยัญชนะตัวที่ 2 (c2)	พยัญชนะตัวแรก (c1)
น	อ ว ห ส ช ฅ
ร	ก อ ห ป บ จ ด พ ฅ
ง	ห
ม	อ ท ห ส
ย	ห
ว	อ ท ห ส ค ฅ
ล	ก อ ค ห ส ป พ ฅ
ค	ส
ท	อ
ช	ว
ฅ	ส
ณ	ค พ
ธ	อ
ญ	ห
ภ	อ
ม	ป

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

PATTERN C-3

รูปแบบสระ	ขอบเขตของพยางค์	ตัวอย่าง	รูปแบบ คุณลักษณะ
ขอบเขตหลัง	g ^๙ *	คีย์ รัย ทรีย	R
	s *	อิก รับ มีค	N
	t *	ถึ นึ หมึ ชึ	R
	s *	จึค กรึค	R
	*	มิ คี ตี ปี	N
ขอบเขตหน้า	* c	คิ คี คี	N
	* c c	วติ หนี ตรี	N
	* c c c	สคริ	X,N
ตารางแสดงอักษรควบของรูปสระ : -			
พยัญชนะตัวที่ 2 (c2)		พยัญชนะตัวแรก (c1)	
น		อ ห	
ร		ก ฆ ท ฬ ป ล ศ ฟ	
ก		ค	
ม		ห	
ย		ห	
ว		ก ท ห ส ค	
ล		ก ว ฬ ป ล พ ศ ถ	
ค		บ ค	
บ		ก	
พ		ร	
ณ		ม	

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาติให้นำไปใช้ประโยชน์ทางการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

PATTERN C-4

รูปแบบสระ	ขอบเขตของพยางค์	ตัวอย่าง	รูปแบบ คุณลักษณะ
ขอบเขตหลัง	s ɔ̌ ɔ̌ *	สุทธี	R
	s ɔ̌ *	มนุษย์ สุกร์	R
	s ɔ̌ *	วุฒิ พุฒิ	X
	ss *	บุตร ชูท อุนห สุมาร	X
	s *	ตุก สุมค ตุก	N
	*	ตุ อุ ตุ ทุ	N
	ts *	ยง ปุย มุง	N
	t *	ตุ	N
ขอบเขตหน้า	*c	ชค จค หัน	N
	*cc	สมค ปกค หนุม	N

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางแสดงอักษรควบของรูปสระ : -

พยัญชนะตัวที่ 2 (c2)	พยัญชนะตัวแรก (c1)
น	อ ม ห ส
ร	ก อ ว ศ ป ค ผ
ก	อ ม ส ศ
อ	ช
ง	อ
ม	ห ส
ย	อ ห ส พ
ว	ห
ล	ก ห ส ค พ ผ ฟ
ท	ป
ค	ส ผ
ส	อ พ
ถ	ด
ธ	ม

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

PATTERN C-5

รูปแบบสระ	ขอบเขตของพยางค์	ตัวอย่าง	รูปแบบ คุณลักษณะ
ขอบเขตหลัง	s g * s s * s * * t s * t *	ศูนย์ บุรณ สุคร มุคร สูบ ลูก รูป คู ปู หู ปู่ ตู้ อยู่	R X,N N N N N
ขอบเขตหน้า	* c * c c * c c c	ตู้ คู หนุ อยู่ สบู ศกรู	N N X,N
ตารางแสดงอักษรควบของรูปสระ :			
พยัญชนะตัวที่ 2 (c2)		พยัญชนะตัวแรก (c1)	
น	ม ห ฐ		
ร	ห ป ค		
ม	ห		
ย	อ ห		
ว	ห		
ล	อ ห ป ฟ		
ค	ศ		
ท	พ		
บ	ส		
ณ	อ		

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

PATTERN C-6

รูปแบบสระ	ขอบเขตของพยางค์	ตัวอย่าง	รูปแบบ คุณลักษณะ
ขอบเขตหลัง	s g g g ^๑ *	ลักษณะ	R
	s g ^๒ g ^๑ *	กษัตริย์	R
	s g g ^๑ *	ยันตร์	R
	s g ^๒ *	ศักดิ์	R
	s g ^๑ *	พันธุ	R
	s g ^๑ *	ศัพท ยันต์	R
	s *	บติ วิติ ณีติ	R
	s s *	บัตริ มณฑา สัมคร	X,N
	s *	กัค มณฑา หัก	N
	ts *	ตง ยง ขน	R
	v *	บัว วิว หิว	R
	tv *	ริ้ว ยัว มิว	R
	y *	วิย นัย ศัย	R
	ty *	ไม่พบคำที่สอดคล้อง	R
	v ^๑ *	หิวะ ทลิวะ	R
ขอบเขตหน้า	* c	ัคค จั๊บ รั๊ก	N
	* c c	หยัง สังค หวัด	N
	* c c c	มฆวัน	X

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางแสดงอักษรควมของรูปสระ : - ๘

พยัญชนะตัวที่ 2 (c2)	พยัญชนะตัวแรก (c1)
----------------------	--------------------

น	อ ห ส ป ท ผ ล
ร	ก อ ท ค ห ส ป บ จ ค พ ผ
ก	ส
ง	ห ส ผ
ม	ห ส
ย	ห ข
ว	ร ค ห ส ค ข ฐ
ด	ก อ ค ห ส ป ล พ ผ
ท	ห
ค	ด
ห	ร อ ม
ส	อ ว
บ	ค
ภ	อ



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

PATTERN C-7

รูปแบบสระ	ขอบเขตของพยางค์	ตัวอย่าง	รูปแบบ คุณลักษณะ
ขอบเขตหลัง	Δ s g Δ * s * * Δ t s * Δ t *	ถึงค์ คิงส์ คริ่งส์ ตึก ตึก คิง รี หี ซี ขบ คิง ำง ไม่พบคำที่สอดคล้อง	R N X,N R N
ขอบเขตหน้า	* c * c c	คิก จิง ำง หน่ง ตกลง ำงค	N N
ตารางแสดงอักษรควบของรูปสระ :	พยัญชนะตัวที่ 2 (c2)	พยัญชนะตัวแรก (c1)	
	น	อ ู ค ฃ	
	ร	ค ำ ค ฬ	
	ม	ห	
	ด	ก อ ฃ ค ฃ ฃ	

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

PATTERN C-8

รูปแบบสระ	ขอบเขตของพยางค์	ตัวอย่าง	รูปแบบ คุณลักษณะ
ขอบเขตหลัง	s *	คิ่น ยืม ปิ่น	R
	t *	ไม่พบคำที่สอดคล้อง	R
	t s *	พิ่น คิ่น คลิ่น	R
	อ *	คือ มือ หรือ	R
	ล อ *	เรือ มือ ลือ	R
ขอบเขตหน้า	* c	มีด คิ่น	N
	* c ๑	หมิ่น คลิ่น หรือ	N
ตารางแสดงอักษรควบของรูปสระ :-			
	พยัญชนะตัวที่ 2 (c2)	พยัญชนะตัวแรก (c1)	
	น	ห	
	ร	ห ป ค	
	ม	ห	
	ว	ห	
	ล	ก ห ป ค	

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

PATTERN C-9

รูปแบบสระ - ^๕ , - ^๕ อ-	ขอบเขตของพยางค์	ตัวอย่าง	รูปแบบ คุณลักษณะ
ขอบเขตหลัง	^๕ อ s *	ช้อค น้อค ล้อค	R
	^๕ *	กี้ (เป็นคำเดียวที่สอดคล้อง)	R
ขอบเขตหน้า	* c	มี ^๕ ค คี ^๕ บ	N
	* c c	หมี ^๕ น กลี ^๕ น หรือ	N
ตารางแสดงอักษรควบของรูปสระ :-			
	พยัญชนะตัวที่ 2 (c2)	พยัญชนะตัวแรก (c1)	
	ม	ณ	
	ถ	พ	

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

PATTERN C-10

รูปแบบสระ - ะ	ขอบเขตของพยางค์	ตัวอย่าง	รูปแบบ คุณลักษณะ
ขอบเขตหลัง	ะ *	กะ อะ พระ ชนะ	R
ขอบเขตหน้า	* ะ * ะ ะ * ะ ะ * ะ ะ ะ * ะ ะ ะ ะ	มะ ะ สระ ะ ะ ะ คะ ะ ะ ะ ไม่พบคำที่สอดคล้อง สระ ะ ะ ะ ะ	N N N N X,N
ตารางแสดงอักษรควบของรูปสระ : - ะ			
	พยัญชนะตัวที่ 2 (c2)	พยัญชนะตัวแรก (c1)	
	น	ช	
	ร	ก ฑ ฑ ฑ ฑ ฑ ฑ	
	ว	ท ฑ	
	ล	ฬ ฑ ฑ ฑ	
	ณ	ณ ฑ ฑ	

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

PATTERN C-11

รูปเบบสระ - ำ	ขอบเขตของพยางค์	ตัวอย่าง	รูปแบบ คุณลักษณะ
ขอบเขตหลัง	*	รำ คำ กำ ล้ำ	R
ขอบเขตหน้า	* c ำ * c c ำ * c t ำ * c c t ำ	คำ จำ นำ หย่า ทราย คำ น้ำ ย้ำ ปล้ำ สม่่า หน้า	N N N N
ตารางแสดงอักษรควบของรูปสระ : - ำ			
	พยัญชนะตัวที่ 2 (c2)	พยัญชนะตัวแรก (c1)	
	น	ห	
	ร	ก ค ฝ ค พ	
	ม	ต ข	
	ย	ท	
	ว	ห ค	
	ด	ก ห ฝ ค พ ถ	

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

PATTERN C-12

รูปแบบสระ - ฤ	ขอบเขตของพยางค์	ตัวอย่าง	รูปแบบ คุณลักษณะ
ขอบเขตหลัง	ฤ s g ิ ุ *	ฤทธิ	X,R
	ฤ s g ุ *	ฤกษ์	R
	ฤ g ุ *	คฤห์	R
	ฤ s *	กฤษ พฤษ มฤต สฤง	N
	ฤ *	พฤ คฤ นฤ	N
	ฤ	ฤสุ ฤลี ฤทัย ฤทธิ ฤทธิ	X,R
ขอบเขตหน้า	* c ฤ	พฤ สฤงศ์ ฤ	R
	* c c ฤ	อมฤต	X,R
	* ฤ	ฤษี ฤายี ฤทธิ (พบเพียง 9 คำที่สอดคล้อง)	X,R

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

PATTERN C-13

รูปแบบสระ	ขอบเขตของพยางค์	ตัวอย่าง	รูปแบบ คุณลักษณะ
เ -			
ขอบเขตหลัง	l c s g * l c t g * l c g * l c s * l c s s * l c s * l c * l c t * l c t s * l c c * l c c s * l c c s s * l c c t * l c c t s *	เกณฑ์ เจศก์ เศษณ์ เถ่ห้ เท่ห้ เรย์ เนย์ เหตุ เกตุ เมรุ เมคร เนคร เอก เมฆ เวร เท เท เก เร เข็ เร้ง เข้มเวก่ง เปลล เกรล อเนก เปลล เกรล อเนก เศรยฐ เกษคร เศร์ เปลลจ เขม่น	R R R X N N N N N N N X,N N N
	เ - (t) ใ . . .	ดู pattern C-13.a	
	เ - (t) อ . . .	ดู pattern C-13.b	
	เ - (t) ฎ . . .	ดู pattern C-13.c	
	เ - (t) ฮ . . .	ดู pattern C-13.c	
	เ - ^า . . .	ดู pattern C-13.d	
	เ - (t) ะ . . .	ดู pattern C-13.e	
	เ - ^า . . .	ดู pattern C-13.f	
	เ - ^า . . .	ดู pattern C-13.g	
	เ - ^า . . .	ดู pattern C-13.h	

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ขอบเขตหน้า	* เ	พบเป็นส่วนมาก	N
	* c เ	พเนจร อเวจี อเมริกา (พบเพียงไม่กี่คำ)	X,N

ตารางแสดงอักษรควบของรูปสระ : เ -



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

PATTERN C-13.a

รูปแบบสระ	ขอบเขตของพยางค์	ตัวอย่าง	รูปแบบ คุณลักษณะ
เ-า , เ-าะ			
ขอบเขตหลัง	เ-า g * เ-า c *	เข่า เข่า เมาท์ เคาะห์ เคะว้	R X,N
	เ-า *	เกาะ เกาะ เปราะ	R
	เ-าะ g * เ-าะ *	เกณฑ์ เทศก์ เสน่ห์ เวีย เเนย	R R
ขอบเขตหน้า		ดู pattern C-13	
ตารางแสดงอักษรควบของรูปสระ : เ-า			
พยัญชนะตัวที่ 2 (c2)		พยัญชนะตัวแรก (c1)	
ร		ป ค ท ส	
ง			
ม		ห ข	
ย		ห ข	
ว		ห	
ด		ก ห ป ค พ ข	

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางแสดงอักษรควบของรูปสระ : เ - ำ ะ

พยัญชนะตัวที่ 2 (c2)	พยัญชนะตัวแรก (c1)
ร	ก ป ค พ
ม	ห
ย	ห
ล	ห ป พ ผ



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

PATTERN C-13.b

รูปแบบสระ เ-อ ,เ-อะ	ขอบเขตของพยางค์	ตัวอย่าง	รูปแบบ คุณลักษณะ
------------------------	-----------------	----------	---------------------

ขอบเขตหลัง	เ - อ g s *	เทอร์ม เอร์ค	X,N
	เ - อ g *	เบอร์ เฟอร์ เทอร์	R
	เ - อ s *	เท็ด เทอม เทอญ เซอร์	X,N
	เ - อ *	เออ เจอ เท่อ เสมอ เกล่อ	N

	เ - อ ะ *	เคอะ เกรอะ เปราะ เทอะ	R
--	-----------	-----------------------	---

ขอบเขตหน้า	* เ	ดู pattern C-13	
------------	-----	-----------------	--

ตารางแสดงอักษรควบของรูปสระ : เ-อ

พยัญชนะตัวที่ 2 (c2)	พยัญชนะตัวแรก (c1)
น	ด
ร	ก ห ป
ม	ท ต
ย	ห
ล	ก ห ผ

ตารางแสดงอักษรควบของรูปสระ : เ-อะ

พยัญชนะตัวที่ 2 (c2)	พยัญชนะตัวแรก (c1)
ร	ก ป ค

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

PATTERN C-13.c

รูปแบบสระ เ-ว ,เ-ย	ขอบเขตของพยางค์	ตัวอย่าง	รูปแบบ คุณลักษณะ
ขอบเขตหลัง	เ - ว ย *	เทอร์ม เตอร์ก	X,N
	เ - ว *	เบอร์ เฟอร์ เทอร์	R
	เ - ย *	เท็ด เทอม เทอญ เฮอร์	X,N
ขอบเขตหน้า	* เ	ดู pattern C-13	
ตารางแสดงอักษรควบของรูปสระ : เ-ว			
	พยัญชนะตัวที่ 2 (c2)	พยัญชนะตัวแรก (c1)	
	ล	ฬ	
ตารางแสดงอักษรควบของรูปสระ : เ-ย			
	พยัญชนะตัวที่ 2 (c2)	พยัญชนะตัวแรก (c1)	
	น	จ	
	ว	ห	
	ล	ป ช ฉ	
	บ	ส	

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

PATTERN C-13.d

รูปแบบสระ	ขอบเขตของพยางค์	ตัวอย่าง	รูปแบบ คุณลักษณะ
เ- ย ,เ- ยะ			
ขอบเขตหลัง	เ - ฎ ย ฎ * เ - ฎ ย ส * เ - ฎ ย * เ - ฎ เ ย * เ - ฎ เ ย ส *	เกียรติ์ เขียว์ เกลียว์ เสียง เรียง เหนียว เสียด เียด เมียด เสียด เียด เมียด เสียง เขียว เหนียว	R N N N N
	เ - ฎ ย ฎ * เ - ฎ เ ย ฎ *	เปรี้ยว เจี้ยว เเพี้ยว	R R
ขอบเขตหน้า	* เ	คู pattern C-13	
ตารางแสดงอักษรควบของรูปสระ : เ- ย			
พยัญชนะตัวที่ 2 (c2)		พยัญชนะตัวแรก (c1)	
น		ห ส พ	
ร		ก ด ห ล พ	
ม		ห ส	
ย		ห	
ว		ก ห ฉ	
ล		ก ห ป พ ฉ	
บ		ส	
ษ		ก	

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางแสดงอักษรควบของรูปสระ : - ยะ

พยัญชนะตัวที่ 2 (c2)	พยัญชนะตัวแรก (c1)
ร	ป
ล	พ ผ



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

PATTERN C-13.e

รูปแบบสระ เ - ะ	ขอบเขตของพยางค์	ตัวอย่าง	รูปแบบ คุณลักษณะ
ขอบเขตหลัง	! - ะ *	กะ เอะ เตะ เค้ะ เประ เฆละ เพละ	R N
ขอบเขตหน้า	* !	ดู pattern C-13	
ตารางแสดงอักษรควบของรูปสระ : เ - ะ			
พยัญชนะตัวที่ 2 (c2)			
ร ก			
ป ค			
พยัญชนะตัวแรก (c1)			

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

PATTERN C-13.f

รูปแบบสระ เ - -	ขอบเขตของพยางค์	ตัวอย่าง	รูปแบบ คุณลักษณะ
ขอบเขตหลัง	เ - s *	เกิน เซ็ญ เป็ด เล็ก เพ็ลิด เสริฐ เผล็ญ	R
	เ - t s *	เพ็ง เน็น เพ้ม	R
	เ - . . .	เซ็ค เอิ๊ตซ์	X
ขอบเขตหน้า	* เ	ดู pattern C-13	
ตารางแสดงอักษรควบของรูปสระ : เ - -			
พยัญชนะตัวที่ 2 (c2)		พยัญชนะตัวแรก (c1)	
ร		ภ ส จ พ	
อ		ฬ	
ย		ห ผ	
ถ		ต ท ล พ ฉ	
ด		ศ	
ช		ศ	

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

PATTERN C-13.g

รูปแบบสระ เ - ๙ อ	ขอบเขตของพยางค์	ตัวอย่าง	รูปแบบ คุณลักษณะ
ขอบเขตหลัง	เ - อ s *	เลือก เหมือน เปลี่ยน	N
	เ - อ *	ถือ เกถือ เหลือ	N
	เ - t อ *	ถือ เมื่อ เชื้อ	N
	เ - t อ s *	เชื่อง เครื่อง เหนื่อย	N
	เ - (t) อ ๙ *	ไม่มีใช้	R
ขอบเขตหน้า	* เ	ดู pattern C-13	
ตารางแสดงอักษรควบของรูปสระ : เ - อ			
	พยัญชนะตัวที่ 2 (c2)	พยัญชนะตัวแรก (c1)	
	น	ห	
	ร	ป ค พ	
	ง	ห	
	ม	ห ล	
	ย	ห	
	ล	ก อ ม ห ป ค ถ	

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

PATTERN C-13.h

รูปแบบสระ เ - -	ขอบเขตของพยางค์	ตัวอย่าง	รูปแบบ คุณลักษณะ
ขอบเขตหลัง	เ - s g * เ - s *	เซ็นต์ เต็นท์ เล็ก เย็น เจ็บ เอ็ด เสร็จ เมล็ด เค็ลค เกล็จ	R N N
ขอบเขตหน้า	* เ	ดู pattern C-13	
ตารางแสดงอักษรควบของรูปสระ : เ - -			
	พยัญชนะตัวที่ 2 (c2)	พยัญชนะตัวแรก (c1)	
	น ร ม ว ถ ด	ห ก ค ฑ ค ฑ ห ส ข ค จ ก ม ค ค ด ฝ	

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

PATTERN C-14

รูปแบบสระ แ -	ขอบเขตของพยางค์	ตัวอย่าง	รูปแบบ คุณลักษณะ
ขอบเขตหลัง	แ c c s g ^๑ *	แสดคมบี	R
	แ c s g ^๑ *	แพทช	R
	แ c g ^๑ *	แคร่ แพรี่	R
	แ c s *	แสดง แขน แดค	N
	แ c *	แก แห แพ	N
	แ c t *	แก็ แม่ แท้	N
	แ c t s *	แก็ง แร้ง แต้้ว	N
	แ c c s *	แถมดง แดคัง แขนง	N
	แ c c c *	แคว แพร แคร	N
	แ c c t s *	แคว้ว แกร้ง แสรั้ง	N
	แ c c t *	แพร แคร	N
	แ c ะ *	แทะ แลละ แคะ	R
	แ c c ะ *	แกลละ แครละ	R
	แ c t ะ *	แนะ	R
	แ c c t ะ *	ไม่พบคำที่สอดคล้อง	R
	แ c s *	แแข็ง แล็บ แต้ป	R
	แ c c s *	แผล็บ แเข็ก แหมีบ	R
ขอบเขตหน้า	* แ	พบเป็นส่วนมาก	
	* c แ	แหง สเน็ค สแก็คส์ (พบเพียงไม่กี่คำ)	

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางแสดงอักษรควบของรูปสระ : แ -

พยัญชนะตัวที่ 2 (c2)	พยัญชนะตัวแรก (c1)
น	ห พ ศ
ร	ก ท ส ป ค พ
ง	ห
ม	ห ข
ย	ท ห ข
ว	ก ส ค ข
ล	ก ม ห ป ค พ ช ถ ญ
ค	ส

ตารางแสดงอักษรควบของรูปสระ : แ - ะ

พยัญชนะตัวที่ 2 (c2)	พยัญชนะตัวแรก (c1)
ร	ค
ม	ห
ย	ห ส
ว	ห ข
ล	ห พ ศ

ตารางแสดงอักษรควบของรูปสระ : แ - ะ

พยัญชนะตัวที่ 2 (c2)	พยัญชนะตัวแรก (c1)
ร	ค
ม	ห
ย	ห ส
ว	ห ข
ล	ห พ ศ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ทางการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

PATTERN C-15

รูปแบบสระ โ -	ขอบเขตของพยางค์	ตัวอย่าง	รูปแบบ คุณลักษณะ
ขอบเขตหลัง	โ c s g ^๑ *	โยชน์ โภชน์ โรจน์	R
	โ c t g ^๑ *	โถ่	R
	โ c g ^๑ *	โพธิ์	R
	โ c g ^๑ *	โขว้ โป้รี่ โบ้ว	R
	โ c s ^๑ *	โชติ	X
	โ c s s *	โศคร โอยฐ	X
	โ c s *	โศค โจอก โจน	N
	โ c *	โศ โห โท	N
	โ c t *	โช้ โก่อ้ โง้ง	N
	โ c t s *	โศ้ง โมง โกอง	N
	โ c c *	โปร โทล	N
	โ c c s *	โปรล โกรธ	N
	โ c c s s *	ไม่พบคำที่สอดคล้อง	R
	โ c c t *	โศล โห้ว โทล	N
	โ c c t s *	โปร้ง โทม้ง	N
	โ c ๕ *	โละ โปะ	R
	โ c t ๕ *	โละ โปะ	R
ขอบเขตหน้า	* โ	พบเป็นส่วนมาก	R
	* c โ	นโยบาย สโมสร อโสค มโหพาร (พบเพียงไม่กี่คำ)	X,R

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางแสดงอักษรควบของรูปสระ : โ -

พยัญชนะตัวที่ 2 (c2)	พยัญชนะตัวแรก (c1)
น	ห ฉ
ร	ท ห ฬ ค พ ศ
ม	ห ข
ย	ห พ ษ
ว	ห
ล	ก ห ค พ ศ



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

PATTERN C-16

รูปแบบสระ ๆ -	ขอบเขตของพยางค์	ตัวอย่าง	รูปแบบ คุณลักษณะ
ขอบเขตหลัง	ไ c c g ʳ *	ไฟลท์	R
	ไ c g ʳ *	ไมล์ ไกท์ สไลด์ สไลด์	R
	ไ c *	ไป ไท โย ไฟ	N
	ไ c t *	ไร่ ไร่ ไร่ ไร่	R
	ไ c c *	ไหน ไหม ไกว	N
	ไ c c t *	ไร่ ไร่ ไร่ ไร่	N
	ไ c ย *	ไทย ไคย	N
ขอบเขตหน้า	* โ	พบเป็นส่วนมาก	R
	* c โ	ต๋าย ชไมพร ต๋ายด์ มโหศวรรษ (พบเพียงไม่กี่คำ)	X,R
ตารางแสดงอักษรควบของรูปสระ : ไ -			
	พยัญชนะตัวที่ 2 (c2)	พยัญชนะตัวแรก (c1)	
	น	ห น	
	ร	ก ต ห ป พ ช	
	ม	ห	
	ว	ก ห ข	
	ถ	ก ห ค พ ถ	
	ท	ผ	

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

PATTERN C-17

รูปแบบสระ ๑ -	ขอบเขตของพยางค์	ตัวอย่าง	รูปแบบ คุณลักษณะ
ขอบเขตหลัง	๑ c *	จ ๑ ๒ ๓ ๔ ๕ ๖	R
	๑ c t *	ห้ ไร่ ไร่ ไร่ ไร่	R
	๑ c c *	หล ๑ ๒	N
	๑ c c t *	ใหญ่ ใหญ่ ๑ ๒ ๓ ๔ (พบคำที่ขึ้นต้นด้วย ๑ เพียง 20 คำ)	N
ขอบเขตหน้า	* ๑	พบในทุกคำ	R



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ภาคผนวก ก

ตัวอย่างรายละเอียดของคำที่ตัดผิดจากโปรแกรมการตัดคำโดยใช้คลังข้อความ

บางครั้ง/จะ/ให้/คำ/วินิจฉัย/โรค/ตาม/ตำแหน่ง/ที่/เป็น/เช่น/
วินิจฉัย/โรค วินิจฉัยโรค

โรคมะเร็ง/ที่/ปอด/ก็/ให้/คำ/วินิจฉัย/ว่า/เป็น/โรคมะเร็งปอด/
โรคมะเร็ง โรคมะเร็ง
โรคมะเร็งปอด โรคมะเร็ง/ปอด

โรคมะเร็ง/ที่/ต่อมน้ำเหลือง/
โรคมะเร็ง/ที่/ต่อมน้ำเหลือง โรคมะเร็ง/ที่/ต่อมน้ำเหลือง/ง

แทนที่จะ/ให้/คำ/วินิจฉัย/ว่า/เป็น/โรคมะเร็งต่อมน้ำเหลือง/แต่/แพทย์/ส่วนใหญ่/จะ/ยัง/คุ้นเคย/กับ/คำ/วินิจฉัย/
ว่า/เป็น/

แทนที่จะ แทน/ที่/จะ
โรคมะเร็งต่อมน้ำเหลือง โรคมะเร็ง/ต่อมน้ำเหลือง/ง
ส่วนใหญ่ ส่วน/ใหญ่
คำวินิจฉัย คำ/วินิจฉัย

ซึ่ง/ไม่/เป็น/ไป/ตาม/ระบบ/ที่/กำหนด/ไว้/ผล/ดังกล่าว/ทำให้/ข้อมูล/ส่วน/ที่/ต่อท้าย/นั้น/ไม่/ถูกต้อง/ตามที่/
ต้องการ/ปัญหา/ดังกล่าว/สามารถ/แก้ไข/โดย/ก่อน/การ/นำ/ระบบ/ไป/ใช้/นั้น/จะ/ต้อง/มี/การ/จัดทำ/คู่มือ/
อธิบาย/เพื่อ/ไว้/เป็น/เอกสารอ้างอิง/และ/มี/การ/จัด/อบรม/ให้/ผู้/เกี่ยวข้อง/ทุก/ระดับ/ทำ/ความ/เข้าใจ/ใน/ระบบ/
ที่/ได้/พัฒนา/ขึ้น/การ/พัฒนา/ระบบ/ฐานข้อมูล/ของ/หน่วยงาน/ต่างๆ/

เป็นไป เป็น/ไป
ดังกล่าว ดัง/กล่าว
ทำให้ ทำ/ให้
ต่อท้าย ต่อ/ท้าย
ถูกต้อง ถูก/ต้อง
ตามที่ ตาม/ที่

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ต้องการ ต้อง/การ

ตั้งกล่าว ตั้ง/กล่าว

สามารถ สามารถ

แก้ไข แก้/ไข

จัดทำ จัด/ทำ

คู่มือ คู่มือ

เอกสารอ้างอิง เอก/สาร/อ้างอิง

อบรม อบรม

ผู้เกี่ยวข้อง ผู้/เกี่ยวข้อง

ระดับ ระดับ

ฐานข้อมูล ฐาน/ข้อมูล

หน่วยงาน หน่วยงาน

แม้/จะ/มี/ความ/ยุ่งยาก/สลับซับซ้อน/และ/มี/ความ/หลากหลาย/แตกต่าง/กัน/แต่/ก็/ไม่ใช่/เป็น/สิ่ง/ที่/เหนือ/วิสัย/ความ/สามารถ/ของ/มนุษย์/ปัญหา/สำคัญ/อยู่ที่/ผู้ใช้/จะ/ยอม/ใช้/ระบบ/ที่/พัฒนา/มา/อย่าง/ดี/แล้ว/หรือ/ไม่/ต่างหาก/ซึ่ง/โดย/ธรรมชาติ/ของ/มนุษย์/สลับซับซ้อน สลับ/ซับซ้อน

ไม่ใช่ ไม่ใช่

สามารถ สามารถ

ผู้ใช้ ผู้/ใช้

อย่างดี อย่าง/ดี

หรือไม่ หรือ/ไม่

ต่างหาก ต่าง/หาก

ธรรมชาติ ธรรม/ชาติ

จะ/มี/ความ/ไม่/เป็น/ระบบ/ใน/ตัวเอง/

ตัวเอง ตัว/เอง

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ขอขอบพระคุณ ขอ/ขอบ/พระ/คุณ

สำนักงานพัฒนาวิทยาศาสตร์และเทคโนโลยีแห่งชาติ สำนักงาน/พัฒนา/วิทยาศาสตร์/และ/เทคโนโลยี/แห่ง/ชาติ

ศูนย์เทคโนโลยีอิเล็กทรอนิกส์และคอมพิวเตอร์แห่งชาติ ศูนย์/เทคโนโลยี/อิเล็กทรอนิกส์/และ/คอมพิวเตอร์/แห่ง/ชาติ

ที่/ได้/เห็น/ความ/สำคัญ/ของ/โครงการ/นี้/และ/ให้/การ/สนับสนุน/ใน/การ/วิจัย/พัฒนา/ครั้งนี้/เอกสารอ้างอิง/กรมการปกครอง/

เอกสารอ้างอิง เอก/สาร/อ้างอิง

กรมการปกครอง กรม/การ/ปก/ครอง

เอกสาร/สิ่งพิมพ์/รายชื่อ/จังหวัด/

เอกสาร เอก/สาร

สิ่งพิมพ์ สิ่ง/พิมพ์

รายชื่อ ราย/ชื่อ

ไทย/พร้อมกับ/รหัส/ที่/เกี่ยวข้อง/

พร้อมกับ พร้อม/กับ

เอกสาร/พิมพ์/จาก/ฐานข้อมูล/ใน/เครื่องคอมพิวเตอร์/

เอกสาร เอก/สาร

ฐานข้อมูล ฐาน/ข้อมูล

เครื่องคอมพิวเตอร์ เครื่อง/คอมพิวเตอร์

กรุงเทพมหานคร/

กรุงเทพมหานคร กรุงเทพมหานคร

กระทรวงมหาดไทย/

กระทรวงมหาดไทย กระทรวง/มหาด/ไทย

การสื่อสารแห่งประเทศไทย/

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตัวอย่างรายละเอียดของคำที่ตัดผิดจากโปรแกรมการตัดคำโดยใช้พจนานุกรม

บางครั้ง/จะ/ให้/คำ/วินิจฉัย/โรค/ตาม/ตำแหน่ง/ที่/เป็น/เช่น/
 วินิจฉัย/โรค วินิจฉัยโรค

โรคมะเร็ง/ที่/ปอด/ก็/ให้/คำ/วินิจฉัย/ว่า/เป็น/โรคมะเร็งปอด/

โรคมะเร็ง โรค/มะเร็ง

โรคมะเร็งปอด โรค/มะเร็ง/ปอด

โรคมะเร็งที่ค่อมน้ำเหลือง/

โรคมะเร็งที่ค่อมน้ำเหลือง โรค/มะเร็ง/ที่/ค่อม/น้ำ/เหลือง/

แทนที่จะ/ให้/คำ/วินิจฉัย/ว่า/เป็น/โรคมะเร็งค่อมน้ำเหลือง/แต่/แพทย์/ส่วนใหญ่/จะ/ยัง/คุ้น/เคย/กับ/คำ/วินิจฉัย/
 ว่า/เป็น/

แทนที่จะ แทน/ที่/จะ

โรคมะเร็งค่อมน้ำเหลือง โรค/มะเร็ง/ค่อม/น้ำ/เหลือง/

ส่วนใหญ่ ส่วน/ใหญ่

คำวินิจฉัย คำ/วินิจฉัย

ซึ่ง/ไม่/เป็น/ไป/ตาม/ระบบ/ที่/กำหนด/ไว้/ผล/ดังกล่าว/ทำให้/ข้อมูล/ส่วน/ที่/ต่อท้าย/นั้น/ไม่/ถูกต้อง/ตามที่/
 ต้องการ/ปัญหา/ดังกล่าว/สามารถ/แก้ไข/โดย/ก่อน/การ/นำ/ระบบ/ไป/ใช้/นั้น/จะ/ต้อง/มี/การ/จัดทำ/คู่มือ/
 อธิบาย/เพื่อ/ไว้/เป็น/เอกสารอ้างอิง/และ/มี/การ/จัด/อบรม/ให้/ผู้/เกี่ยวข้อง/ทุก/ระดับ/ทำ/ความ/เข้าใจ/ใน/ระบบ/
 ที่/ได้/พัฒนา/ขึ้น/การ/พัฒนา/ระบบ/ฐานข้อมูล/ของ/หน่วยงาน/ต่าง/ๆ/

เป็นไป เป็น/ไป

ดังกล่าว ดัง/กล่าว

ทำให้ ทำ/ให้

ต่อท้าย ต่อ/ท้าย

ถูกต้อง ถูก/ต้อง

ตามที่ ตาม/ที่

ต้องการ คือ/งการ

ดังกล่าว ดัง/กล่าว

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สามารถ สามารถ

แก้ไข แก้ไข

จัดทำ จัดทำ

คู่มือ คู่มือ

เอกสารอ้างอิง เอกสารอ้างอิง

อบรม อบรม

ผู้เกี่ยวข้อง ผู้เกี่ยวข้อง

ระดับ ระดับ

ฐานข้อมูล ฐานข้อมูล

หน่วยงาน หน่วยงาน

แม้จะมี/ความ/ยุ่งยาก/สลับซับซ้อน/และมี/ความ/หลากหลาย/แตกต่างกัน/แต่ก็/ไม่ใช่/เป็น/สิ่ง/ที่/เหนือ/
วิสัย/ความ/สามารถ/ของ/มนุษย์/ปัญหา/สำคัญ/อยู่ที่/ผู้ใช้/จะ/ยอม/ให้/ระบบ/ที่/พัฒนา/มา/อย่าง/ดี/แล้ว/หรือ/ไม่/
ต่าง/หาก/ซึ่ง/โดย/ธรรมชาติ/ของ/มนุษย์/
สลับซับซ้อน สลับซับซ้อน

ไม่ใช่ ไม่ใช่

สามารถ สามารถ

ผู้ใช้ ผู้ใช้

อย่างดี อย่างดี

หรือไม่ หรือไม่

ต่างหาก ต่างหาก

ธรรมชาติ ธรรมชาติ

จะ/มี/ความ/ไม่/เป็น/ระบบ/ใน/ตัวเอง/
ตัวเอง

ตัวเอง

ทำตามใจ/ตัวเอง/และ/ไม่/ชอบ/ถูก/บังคับ/และ/สิ่ง/นี้/คือ/จุด/สำคัญ/ของ/ผู้ที่/จะ/พัฒนา/ระบบ/ฐานข้อมูล/ต้อง/
ทำ/ความ/เข้าใจ/กับ/ผู้ใช้/อย่าง/ดั่ง/แท้/เพื่อ/ให้/ระบบ/ฐานข้อมูล/ที่/พัฒนา/แล้ว/มี/ผู้ใช้/ตอบสนอง/อย่าง/คุ้มค่า/
สูงสุด/กิตติกรรมประกาศ/คณะ/ผู้วิจัย/พัฒนา/ใคร่/ขอ/ขอบพระคุณ/
ความใจ

ความใจ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตัวเอง คัว/เอ/ง

จุดสำคัญ จุด/สำคัญ

ผู้ที่ ผู้/ที่

ฐานข้อมูล ฐาน/ข้อมูล

ผู้ใช้ ผู้/ใช้

อย่างต้องแท้ อย่าง/ทอ/งแท้

เพื่อให้ เพื่อให้

ฐานข้อมูล ฐาน/ข้อมูล

ผู้ใช้ ผู้/ใช้

อย่างคุ้มค่า อย่าง/คุ้ม/ค่า

สูงสุด สูง/สุด

กิตติกรรมประกาศ กิตติ/กร/รม/ป/ระ/กาศ

ผู้วิจัยพัฒนา ผู้/วิจัย/พัฒนา

ขอขอบพระคุณ ขอ/ขอบ/พระ/คุณ

วิวัฒนาการ

สุทธิตถวิวัฒนาการ สุทธิ/พล/วิวัฒนาการ

ที่/ได้/ให้/ข้อเสนอแนะ/ตลอดจน/ข้อมูล/ต่าง/อัน/มี/ประโยชน์/อย่างมาก/ต่อ/การ/วิจัย/พัฒนา/ใน/ครั้ง/นี้/นอกจาก/นี้/คณะ/ผู้/วิจัย/พัฒนา/ใคร่/ขอ/ขอบ/พระ/คุณ/สำนักงาน/พัฒนา/วิทยาศาสตร์/และเทคโนโลยี/แห่งชาติ/และ/ศูนย์/เทคโนโลยี/อิเล็กทรอนิกส์/และ/คอมพิวเตอร์/แห่งชาติ/ข้อเสนอแนะ/ข้อ/เสนอ/แนะ

ตลอดจน ตลอด/จน

มีประโยชน์ มี/ประโยชน์

อย่างมาก อย่าง/มาก

นอกจากนี้ นอก/จาก/นี้

ผู้วิจัยพัฒนา ผู้/วิจัย/พัฒนา

ขอขอบพระคุณ ขอ/ขอบ/พระ/คุณ

สำนักงานพัฒนาวิทยาศาสตร์และเทคโนโลยีแห่งชาติ สำนักงาน/พัฒนา/วิทยาศาสตร์/และ/เทคโนโลยี/แห่งชาติ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ศูนย์เทคโนโลยีอิเล็กทรอนิกส์และคอมพิวเตอร์แห่งชาติ ศูนย์เทคโนโลยีอิเล็กทรอนิกส์/และ/คอมพิวเตอร์/
แห่ง/ชาติ

ตัวอย่างรายละเอียดของคำที่ตัดผิดจากโปรแกรมการตัดคำโดยใช้กฎ

บางครั้ง/จะ/ให้/คำ/วินิจฉัย/โรค/ตาม/ตำแหน่ง/ที่/เป็น/เช่น/
บางครั้ง บางครั้ง

วินิจฉัย วินิจฉัย

ตำแหน่ง ตำแหน่ง

โรคมะเร็ง/ที่/ปอด/ก็/ให้/คำ/วินิจฉัย/ว่า/เป็น/โรคมะเร็ง/ปอด/
โรคมะเร็ง โรคมะเร็ง

โรคมะเร็ง โรคมะเร็ง

วินิจฉัย วินิจฉัย

โรคมะเร็ง/ปอด โรคมะเร็ง/ปอด

โรคมะเร็ง/ที่/ต่อมน้ำเหลือง/
โรคมะเร็ง/ที่/ต่อมน้ำเหลือง โรคมะเร็ง/ที่/ต่อมน้ำเหลือง

โรคมะเร็ง/ที่/ต่อมน้ำเหลือง โรคมะเร็ง/ที่/ต่อมน้ำเหลือง

ซึ่ง/มี/เซลล์/มะเร็ง/ชื่อ/
มะเร็ง มะเร็ง

มะเร็ง มะเร็ง

แทน/ที่จะ/ให้/คำ/วินิจฉัย/ว่า/เป็น/โรคมะเร็ง/ต่อมน้ำเหลือง/แต่/แพทย์/ส่วนใหญ่/จะ/ยัง/คุ้น/เคย/กับ/คำ/วินิจฉัย/
ว่า/เป็น/

แทน/ที่จะ แทน/ที่จะ

วินิจฉัย วินิจฉัย

โรคมะเร็ง/ต่อมน้ำเหลือง โรคมะเร็ง/ต่อมน้ำเหลือง

ส่วนใหญ่ ส่วนใหญ่

คุ้นเคย คุ้นเคย

คำวินิจฉัย คำวินิจฉัย

ซึ่ง/ไม่/เป็น/ไป/ตาม/ระบบ/ที่/กำหนด/ไว้/ผล/ดังกล่าว/ทำให้/ข้อมูล/ส่วน/ที่/ต่อท้าย/นั้น/ไม่/ถูกต้อง/ตามที่/
ต้องการ/ปัญหา/ดังกล่าว/สามารถ/แก้ไข/โดย/ก่อน/การ/นำ/ระบบ/ไป/ใช้/นั้น/จะ/ต้อง/มี/การ/จัดทำ/คู่มือ/

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

อธิบาย/เพื่อ/ไว้/เป็น/เอกสารอ้างอิง/และ/มี/การ/จัด/อบรม/ให้/ผู้/เกี่ยวข้อง/ทุก/ระดับ/ทำ/ความ/เข้าใจ/ใน/ระบบ/
ที่/ได้/พัฒนา/ขึ้น/การ/พัฒนา/ระบบ/ฐานข้อมูล/ของ/หน่วยงาน/ต่าง/ๆ/

เป็นไป เป็นไป

ระบบ/ที่ ระบบ/ที่

กำหนด กำหนด

ผล/ดังกล่าว ผลดังกล่าว

ทำให้/ข้อมูล ทำให้/ข้อมูล

ค่อยๆ ค่อยๆ

ถูกต้อง ถูกต้อง

ตามที่ ตามที่

ต้องการ ต้องการ

ปัญหา ปัญหา

ดังกล่าว ดังกล่าว

สามารถ สามารถ

แก้ไข แก้ไข

ระบบ ระบบ

จัดทำ จัดทำ

คู่มือ คู่มือ

อธิบาย อธิบาย

เอกสารอ้างอิง เอกสาร/อ้างอิง

จัด/อบรม จัด/อบรม

ผู้/เกี่ยวข้อง ผู้/เกี่ยวข้อง

ระดับ ระดับ

เข้าใจ เข้าใจ

ระบบ/ที่ ระบบ/ที่

พัฒนา พัฒนา

พัฒนา พัฒนา

ระบบ ระบบ

ฐานข้อมูล/ของ/หน่วยงาน ฐาน/ข้อมูล/ของ/หน่วย/งาน

ต่างๆ ต่างๆ



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

แม้/จะ/มี/ความ/ยุ่งยาก/สลับซับซ้อน/และ/มี/ความ/หลากหลาย/แตกต่างกัน/แต่/ก็/ไม่ใช่/เป็น/สิ่ง/ที่/เหนือ/
 วิสัย/ความ/สามารถ/ของ/มนุษย์/ปัญหา/สำคัญ/อยู่/ที่/ผู้ใช้/จะ/ยอม/ใช้/ระบบ/ที่/พัฒนา/มา/อย่าง/ดี/แล้ว/หรือ/ไม่/
 ต่างหาก/ซึ่ง/โดย/ธรรมชาติ/ของ/มนุษย์/

มี/ความ มีค/วาม

ยุ่งยาก/สลับซับซ้อน ยุ่ง/ยากส/ลับ/ซับซ้อน

มี/ความ/หลากหลาย มีค/วามห/ลากห/ลาย

แตกต่างกัน แดก/ต่าง

ไม่ใช่ ไม่/ใช่

วิสัย วิ/สัย

สามารถ สา/มารถ

ของ/มนุษย์ ของม/นุษย์

ปัญหา ปัญ/หา

สำคัญ/อยู่ สำ/คัญอยู่

ผู้ใช้ ผู้/ใช้

ระบบ/ที่ ระ/บบ/ที่

พัฒนา พัฒ/นา

มา/อย่าง/ดี มา/อย่าง/ดี

หรือ/ไม่ หรือ/ไม่

ต่างหาก ต่าง/หาก

โดย/ธรรมชาติ/ของ โดย/รรรม/ชาติ/ของ

จะ/มี/ความ/ไม่/เป็น/ระบบ/ใน/ตัวเอง/

มี/ความ มีค/วาม

ระบบ ระ/บบ

ตัวเอง คั/วเอง

ทำ/ตาม/ใจ/ตัวเอง/และ/ไม่/ชอบ/ถูก/บังคับ/และ/สิ่ง/นี้/คือ/จุด/สำคัญ/ของ/ผู้ที่/จะ/พัฒนา/ระบบ/ฐาน/ข้อมูล/ต้อง/
 ทำ/ความ/เข้าใจ/กับ/ผู้ใช้/อย่าง/ถ่อง/แท้/เพื่อ/ให้/ระบบ/ฐาน/ข้อมูล/ที่/พัฒนา/แล้ว/มี/ผู้ใช้/ตอบสนอง/อย่าง/คุ้มค่า/
 สูง/สุด/กิตติ/กรรม/ประกาศ/คณะ/ผู้/วิจัย/พัฒนา/ใคร/ขอ/ขอบ/พระ/คุณ/

ตาม/ใจ ตาม/ใจ

ตัวเอง คั/วเอง



บรรณานุกรม

ภาษาไทย

- [1] สมปรารถนา รัชยานนท์ “โครงสร้างข้อมูลสำหรับพจนานุกรมอิเล็กทรอนิกส์ภาษาไทย” วิทยานิพนธ์
ปริญญาวิทยาศาสตรมหาบัณฑิต ภาควิชาวิศวกรรมคอมพิวเตอร์ คณะวิศวกรรมศาสตร์ จุฬาลงกรณ์
มหาวิทยาลัย พ.ศ. 2535
- [2] ชื่น ภู่วรรณ และ นายวิวรรณ อิมอรณม์ “การตรวจสอบตัวสะกดด้วยคอมพิวเตอร์”
บทความทางวิชาการ การประชุมทางวิชาการ วิศวกรรม ไฟฟ้า 2530
- [3] พิสิทธิ์ พรหมจันทร์ “การวิเคราะห์แนวทางการเปรียบเทียบสมรรถนะของโปรแกรมแยกคำภาษาไทย”
วิทยานิพนธ์ ปริญญาวิทยาศาสตรมหาบัณฑิต ภาควิชาวิศวกรรมคอมพิวเตอร์ คณะวิศวกรรมศาสตร์
จุฬาลงกรณ์มหาวิทยาลัย พ.ศ. 2540
- [4] ไพศาล เจริญพรสวัสดิ์ “การตัดคำภาษาไทยโดยใช้คุณลักษณะ” วิทยานิพนธ์ ปริญญาวิทยาศาสตรมหาบัณฑิต
ภาควิชาวิศวกรรมคอมพิวเตอร์ คณะวิศวกรรมศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย พ.ศ. 2541
- [5] ดวงแก้ว สวามิภักดิ์ “การสร้างซอฟต์แวร์วิเคราะห์ไวยากรณ์ไทยภายใต้ระบบยูนิกซ์” กรุงเทพฯ : สำนัก
พิมพ์มหาวิทยาลัยธรรมศาสตร์ พ.ศ.2533
- [6] บุญเสริม กิ่งศิริกุล “การกำกับหมวดคำสำหรับข้อความภาษาไทย” กรุงเทพฯ: สถาบันวิจัยและ
พัฒนา คณะวิศวกรรมศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย พ.ศ.2540
- [7] ชมพูนุช คุปต์วิมล “การตัดคำกำกับในข้อความภาษาไทยด้วยการโปรแกรมตรรกะเชิงอุปนัย”
ปริญญาวิทยาศาสตรมหาบัณฑิต สาขาวิทยาศาสตร์คอมพิวเตอร์ ภาควิชาวิศวกรรมคอมพิวเตอร์ บัณฑิต
วิทยาลัย จุฬาลงกรณ์มหาวิทยาลัย พ.ศ.2542.
- [8] ชื่น ภู่วรรณ และ วิวรรณ อิมอรณม์. “การแบ่งแยกพยางค์ไทยด้วยคิกซ์มาร์ค.” รายงานการ
ประชุมวิชาการวิศวกรรมไฟฟ้า ครั้งที่ 9 พ.ศ.2529.
- [9] วิรัช ศรีเลิศล้ำนิช “การตัดคำภาษาไทยในระบบแปลภาษา” การแปลภาษาด้วยคอมพิวเตอร์.
หน้า 50-55. กรุงเทพฯ : ศูนย์เทคโนโลยีอิเล็กทรอนิกส์และคอมพิวเตอร์แห่งชาติ พ.ศ.2536
- [10] สัมพันธ์ ธีรรัตน์ “การแบ่งคำไทยด้วยพจนานุกรม” โครงงานวิศวกรรม ภาควิชาวิศวกรรม
คอมพิวเตอร์ คณะวิศวกรรมศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย พ.ศ. 2534.
- [11] สมศักดิ์ จันทน์ “ระบบวิเคราะห์โครงสร้างภาษาไทยด้วยคอมพิวเตอร์” วิทยานิพนธ์มหาบัณฑิต
ภาควิชาวิศวกรรมคอมพิวเตอร์ คณะวิศวกรรมศาสตร์ บัณฑิตวิทยาลัย สถาบันเทคโนโลยีพระ
จอมเกล้าเจ้าคุณทหารลาดกระบัง พ.ศ. 2534.
- [12] นัฐวุฒิ ไชยเจริญ. “การตัดคำและกำกับหมวดคำภาษาไทยแบบเบ็ดเสร็จด้วยคอมพิวเตอร์”
วิทยานิพนธ์ปริญญาอักษรศาสตรมหาบัณฑิต ภาควิชาภาษาศาสตร์ คณะอักษรศาสตร์
จุฬาลงกรณ์มหาวิทยาลัย พ.ศ. 2544.

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

- [13] รัตติกร วรากุลศิริพันธุ์,สง่า คงสุพานิช “การวิเคราะห์เลือกประโยคที่ถูกต้องจากความถี่ของการใช้คำ” Papers on Natural Language Processing , Compiled by Virach Sornlertlamvanich 1995

ภาษาอังกฤษ

- [14] Charnyapornpong, S. “ A Thai Syllable Separation Algorithm” Master Thesis.Asian Institute of Technology. 1983.
- [15] VAN Rijsebergen , C.J. “Information Retrieval”.1979 London: Butterworths.cited in Manning ,Chistoper D. And Hinrich Schutze. Foundations of statistical natural language processing .Cambridge:MIT Press. 1999.



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้