

การเลือกข้อความออนไลน์โดยอัตโนมัติเพื่อสร้างคลังข้อความตามการกระจายตัวหน่วยเสียง
ที่กำหนดได้

นายสุรพล วรรณัทธาทร

วิทยานิพนธ์นี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญาวิทยาศาสตรมหาบัณฑิต

สาขาวิชาวิศวกรรมคอมพิวเตอร์ ภาควิชาวิศวกรรมคอมพิวเตอร์

คณะวิศวกรรมศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย

ปีการศึกษา 2554

ลิขสิทธิ์ของจุฬาลงกรณ์มหาวิทยาลัย

บทคัดย่อและแฟ้มข้อมูลฉบับเต็มของวิทยานิพนธ์ตั้งแต่ปีการศึกษา 2554 ที่ให้บริการในคลังปัญญาจุฬาฯ (CUIR)

เป็นแฟ้มข้อมูลของนิสิตเจ้าของวิทยานิพนธ์ที่ส่งผ่านทางบัณฑิตวิทยาลัย

The abstract and full text of theses from the academic year 2011 in Chulalongkorn University Intellectual Repository(CUIR)
are the thesis authors' files submitted through the Graduate School.

AUTOMATIC ONLINE TEXT SELECTION FOR CONSTRUCTING
TEXT CORPUS WITH CUSTOM PHONEME DISTRIBUTION

Mr. Surapol Vorapatratorn

A Thesis Submitted in Partial Fulfillment of the Requirements
for the Degree of Master of Engineering Program in Computer Engineering

Department of Computer Engineering

Faculty of Engineering

Chulalongkorn University

Academic Year 2011

Copyright of Chulalongkorn University

หัวข้อวิทยานิพนธ์	การเลือกข้อความออนไลน์โดยอัตโนมัติเพื่อสร้างคลังข้อความตามการกระจายตัวหน่วยเสียงที่กำหนดได้
โดย	นายสุรพล วรรณัทธาทร
สาขาวิชา	วิศวกรรมคอมพิวเตอร์
อาจารย์ที่ปรึกษาวิทยานิพนธ์หลัก	ผู้ช่วยศาสตราจารย์ ดร.โปรดปราน บุญยพุกกณะ
อาจารย์ที่ปรึกษาวิทยานิพนธ์ร่วม	ผู้ช่วยศาสตราจารย์ ดร.อดิวงค์ สุชาติ

คณะวิศวกรรมศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย อนุมัติให้หัวข้อวิทยานิพนธ์ฉบับนี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญาวิทยาศาสตรบัณฑิต

..... คณบดีคณะวิศวกรรมศาสตร์
(รองศาสตราจารย์ ดร.บุญสม เลิศหิรัญวงศ์)

คณะกรรมการสอบวิทยานิพนธ์

..... ประธานกรรมการ
(ผู้ช่วยศาสตราจารย์ ดร.เศรษฐา ปานงาม)

..... อาจารย์ที่ปรึกษาวิทยานิพนธ์หลัก
(ผู้ช่วยศาสตราจารย์ ดร.โปรดปราน บุญยพุกกณะ)

..... อาจารย์ที่ปรึกษาวิทยานิพนธ์ร่วม
(ผู้ช่วยศาสตราจารย์ ดร.อดิวงค์ สุชาติ)

..... กรรมการภายนอกมหาวิทยาลัย
(ดร.ชัย วุฒิวิวัฒน์ชัย)

สุรพล วัชรพร : การเลือกข้อความออนไลน์โดยอัตโนมัติเพื่อสร้างคลังข้อความตามการกระจายตัวหน่วยเสียงที่กำหนดได้. (AUTOMATIC ONLINE TEXT SELECTION FOR CONSTRUCTING TEXT CORPUS WITH CUSTOM PHONEME DISTRIBUTION)

อ. ที่ปรึกษาวิทยานิพนธ์หลัก : ผศ.ดร.โปรดปราน บุญยพุกกณะ, อ. ที่ปรึกษาวิทยานิพนธ์ร่วม : ผศ. ดร.อดิวงค์ สุชาโต, 106 หน้า.

ประสิทธิภาพของ ระบบรู้จำเสียงพูดอัตโนมัติและ ระบบสังเคราะห์เสียงพูด ขึ้นอยู่กับความครอบคลุมของหน่วยเสียงจากคลังข้อความที่เหมาะสม วิทยานิพนธ์นี้เสนอการสร้างคลังข้อความอัตโนมัติ จากการกระจายตัวของหน่วยเสียงตามที่กำหนดการกระจายตัวของหน่วยเสียงตามที่กำหนดนั้น สามารถกำหนดได้จากชนิดของหน่วยเสียง ขนาดของคลังข้อความ เกณฑ์ขั้นต่ำของจำนวนหน่วยเสียง และรูปแบบของการกระจายตัวเป้าหมาย ได้คัดเลือกข้อความมาจากข้อมูลจากอินเทอร์เน็ต โดยข้อความนั้นจะถูก จัดเก็บมาอย่างต่อเนื่องโดย กระบวนการดึงบทความจากหน้าเว็บบนอินเทอร์เน็ต จนกระทั่งได้คลังข้อความที่เหมาะสม ในวิทยานิพนธ์นี้ยังได้ประยุกต์ใช้วิธีการเชิงละเอียด เพื่อเลือกประโยคที่เหมาะสมที่จะทำให้เกิดการกระจายตัวของหน่วยเสียงตามเป้าหมาย ในการทดลองได้ใช้ข้อความจากฐานข้อมูล Large Vocabulary Continuous Speech Recognition (LVCSR) corpus for Thai language ในการสร้างเป้าหมายของการกระจายตัวหน่วยเสียง ผลการทดลองที่ได้คือ จำนวนของข้อมูลข้อความที่ดึงมาจากอินเทอร์เน็ตที่เพิ่มขึ้นสามารถทำให้การกระจายตัวของหน่วยเสียงเป็นไปตามเป้าหมายได้ และเกิดความครอบคลุมทางหน่วยเสียงคู่ ถึง 99.13% คลังข้อความที่ถูกสร้างขึ้นนี้ จึงสามารถนำไปใช้ในการสร้างคลังเสียงพูดได้อย่างมีประสิทธิภาพ

ภาควิชา.....วิศวกรรมคอมพิวเตอร์.....ลายมือชื่อนิสิต.....
 สาขาวิชา.....วิศวกรรมคอมพิวเตอร์.....ลายมือชื่อ อ.ที่ปรึกษาวิทยานิพนธ์หลัก.....
 ปีการศึกษา.....2554.....ลายมือชื่อ อ.ที่ปรึกษาวิทยานิพนธ์ร่วม.....

5370514421 : MAJOR COMPUTER ENGINEERING

KEYWORDS : TEXT SELECTION / ONLINE CORPUS / PHONETICALLY BALANCED / GREEDY ALGORITHM / PHONETIC / SENTENCE SEGMENTATION

SURAPOL VORAPATRATORN : AUTOMATIC ONLINE TEXT SELECTION FOR CONSTRUCTING TEXT CORPUS WITH CUSTOM PHONEME DISTRIBUTION.

ADVISOR : ASST.PROF. PROADPRAN PUNYABUKKANA, PH.D.,

CO-ADVISOR : ASST.PROF. DR.ATIWONG SUCHATO, PH.D., 106 pp.

Performance of Automatic Speech Recognition (ASR) and Text-to-Speech (TTS) systems depend on appropriate text corpus. This article explains about the automated text corpus generating method using custom phonetic distribution. This distribution is defined by phonemes type, corpus size, minimum criterion number of phonemes, and target phonetic distribution. Generally, the system selects text data from the internet by continuously downloading them using web crawler. The greedy algorithm is applied to extract the proper sentences, in order to fit with the target phonetic distribution until the appropriate text corpus is established. The experiment is done by using the text from Large Vocabulary Continuous Speech Recognition (LVCSR) corpus for Thai language to generate target phonetic distribution. The result shown that, the increased number of data drawn from the internet is able to accomplish target phonetic distribution and generate diphone coverage for 99.13%. This text corpus then generate speech corpus efficiently.

Department : Computer Engineering Student's Signature.....

Field of Study : Computer Engineering Advisor's Signature.....

Academic Year : 2011 Co-advisor's Signature.....

กิตติกรรมประกาศ

การศึกษาและการทำวิทยานิพนธ์ในครั้งนี้ เกิดปัญหาและอุปสรรคมากมาย ทั้งในเรื่องของระยะเวลาที่จำกัดด้วย แต่ด้วยความช่วยเหลือจากทุกคนจึงทำให้วิทยานิพนธ์นี้เสร็จสิ้นอย่างสมบูรณ์ด้วยคุณภาพที่ดีในระยะเวลาที่กำหนดได้ ข้าพเจ้าขอขอบพระคุณ อาจารย์ที่ปรึกษาวิทยานิพนธ์หลัก ผู้ช่วยศาสตราจารย์ ดร.โปรดปราน บุญยพุกกณะ และ อาจารย์ที่ปรึกษาวิทยานิพนธ์ร่วม ผู้ช่วยศาสตราจารย์ ดร.อดิวงค์ สุชาติ ที่ให้ทั้งความรู้และโอกาสที่ดีต่าง ๆ ในระหว่างระยะเวลาที่ศึกษานี้ และได้สละเวลาอันมีค่าเพื่อให้คำปรึกษาเมื่อข้าพเจ้าต้องการได้เสมอ ข้าพเจ้าขอขอบคุณ รองศาสตราจารย์ ดร.วิโรจน์ อรุณมานะกุล, ผู้ช่วยศาสตราจารย์ ดร.เศรษฐา ปานงาม อาจารย์ ดร.ชัย วุฒิวิวัฒน์ชัย และ อาจารย์ ศิวบุธ อัมพูช ที่ให้คำแนะนำแนวคิดอันมีคุณค่าและถูกนำมาใช้ในการทำวิทยานิพนธ์นี้

ข้าพเจ้าขอขอบคุณ หน่วยปฏิบัติการวิจัยวิทยาการมนุษยภาษา ศูนย์เทคโนโลยีอิเล็กทรอนิกส์และคอมพิวเตอร์แห่งชาติ: NECTEC และภาควิชาภาษาศาสตร์ คณะอักษรศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย ที่ให้เครื่องมือที่จำเป็นในการทดลองของการทำวิทยานิพนธ์ครั้งนี้ ขอขอบคุณเพื่อนร่วมงานในห้องปฏิบัติการระบบภาษาพูด ภาควิชาวิศวกรรมคอมพิวเตอร์ คณะวิศวกรรมศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย ที่ให้คำแนะนำที่ดีกับข้าพเจ้า ให้กำลังใจและสร้างรอยยิ้มให้ข้าพเจ้าในขณะที่ข้าพเจ้าท้อแท้ได้เสมอ ขอขอบคุณ คุณพ่อ คุณแม่ และครอบครัวที่เข้าใจและสนับสนุนการทำวิจัยในครั้งนี้มาโดยตลอด ทำให้การจัดทำวิทยานิพนธ์ในครั้งนี้ประสบความสำเร็จได้

สารบัญ

	หน้า
บทคัดย่อภาษาไทย.....	ง
บทคัดย่อภาษาอังกฤษ.....	จ
กิตติกรรมประกาศ.....	ฉ
สารบัญ.....	ช
สารบัญตาราง.....	ฅ
สารบัญภาพ.....	ญ
บทที่ 1 บทนำ.....	1
ความเป็นมาและความสำคัญของปัญหา.....	1
วัตถุประสงค์ของการวิจัย.....	2
ขอบเขตของการวิจัย.....	2
ลำดับขั้นตอนในการเสนอผลการวิจัย.....	3
ประโยชน์ที่คาดว่าจะได้รับ	3
ผลงานตีพิมพ์จากวิทยานิพนธ์.....	5
ลำดับการจัดเรียงเนื้อหาในวิทยานิพนธ์.....	5
บทที่ 2 ทฤษฎีและงานวิจัยที่เกี่ยวข้อง.....	6
แนวคิดและทฤษฎี.....	6
1. ระบบการเขียนภาษาไทย.....	6
2. ระบบการออกเสียงภาษาไทย.....	6
3. รูปแบบหน่วยเสียง.....	10
4. การแปลงข้อความภาษาไทยเป็นสัทอักษร.....	11
5. การค้นหาสารสนเทศบนอินเทอร์เน็ต.....	12
6. โครงสร้างคำสั่งในภาษา HTML.....	13
7. ขั้นตอนวิธีเชิงละโมภ: Greedy Algorithm.....	18
8. การวัดระยะทางแบบยูคลิด: Euclidean distance.....	18
เอกสารและงานวิจัยที่เกี่ยวข้อง.....	20

บทที่ 3 ขั้นตอนการดำเนินงานวิจัย.....	23
เครื่องมือที่ใช้ในการวิจัย.....	23
ขั้นตอนการดำเนินงานวิจัย.....	24
1. ขั้นตอนการสร้างการกระจายตัวทางหน่วยเสียงเป้าหมาย.....	25
2. ขั้นตอนการสร้างประโยคออนไลน์.....	30
3. ขั้นตอนการเลือกประโยคและจัดเก็บคลังข้อความ.....	39
บทที่ 4 การทดลอง และอภิปรายผล.....	50
การทดลอง.....	50
1. การสร้างคลังข้อความจากรูปแบบการกระจายตัวทางหน่วยเสียงจากการกระจายตัวเป้าหมาย.....	50
2. การสร้างคลังข้อความโดยใช้การกำหนดรูปแบบการกระจายตัวทางหน่วยเสียงเอง.....	58
3. ทดสอบผลของการสร้างคลังข้อความจากการปรับพารามิเตอร์ค่าคงที่ความสำคัญรูปแบบหน่วยเสียงหายาก (α).....	60
4. เปรียบเทียบผลจากวิธีการให้คะแนนในวิทยานิพนธ์กับการให้คะแนนของงานวิจัย Automatic Construction for a TTS Corpus with Limited Text (W. Zhang, 2010).....	63
บทที่ 5 บทสรุปผลการวิจัย และข้อเสนอแนะ.....	67
สรุปผลการวิจัย.....	67
ข้อเสนอแนะ.....	69
รายการอ้างอิง.....	71
ภาคผนวก.....	75
ภาคผนวก ก ตารางเทียบสัทอักษรสากลกับหน่วยเสียงไทยที่ใช้ในวิทยานิพนธ์.....	76
ภาคผนวก ข ตัวอย่างประโยคที่ได้จากคลังข้อความที่สร้างขึ้น 300 ประโยค.....	79
ประวัติผู้เขียนวิทยานิพนธ์.....	95

สารบัญตาราง

ตารางที่		หน้า
2.1	การออกเสียงของพยัญชนะภาษาไทย.....	7
2.2	การออกเสียงของพยัญชนะควบกล้ำไทย.....	8
2.3	การออกเสียงของพยัญชนะหน่วยเสียงยืมภาษาต่างประเทศ.....	8
2.4	การออกเสียงสระไทย.....	9
2.5	ตัวอย่างรูปแบบหน่วยเสียง.....	11
2.6	คำสั่งหัวเรื่องโครงสร้าง HTML.....	15
2.7	คำสั่งส่วนเนื้อความโครงสร้าง HTML.....	16
3.1	ตัวอย่างการแปลงรูปเขียนเป็นรูปอ่านจากฐานข้อมูล LOTUS.....	26
3.2	ตัวอย่างจำนวนของแต่ละรูปแบบหน่วยเสียงคู่ที่เกิดขึ้นจากฐานข้อมูล LOTUS.	28
3.3	ตัวอย่างผลลัพธ์จากการแปลงข้อความรูปแบบบทความเป็นรูปแบบลำดับ พยางค์.....	33
3.4	ตัวอย่างผลการตัดแยกประโยคตามอักขระนิพจน์ปรกติที่ต้องการ.....	34
3.5	ตัวอย่างผลการตัดแยกประโยคตามอักขระนิพจน์ปรกติที่ไม่ต้องการออก.....	35
3.6	ตัวอย่างผลการแปลงรูปเขียนเป็นรูปอ่านของประโยคออนไลน์.....	37
3.7	ข้อมูลทางสถิติการกระจายตัวทางหน่วยเสียงที่ใช้ในการคำนวณจากฐานข้อมูล	43
3.8	ผลคะแนนประโยคจากขั้นตอนการให้คะแนน.....	47
4.1	คุณลักษณะของคลังข้อความที่สร้างขึ้น.....	53
4.2	รูปแบบหน่วยเสียงที่ไม่ปรากฏจากการแปลงรูปเขียนเป็นรูปอ่านที่ผิด.....	53
4.3	รูปแบบหน่วยเสียงที่ไม่ปรากฏจากชื่อเฉพาะ.....	54
5.1	ผลของการเปลี่ยนแปลงค่าพารามิเตอร์.....	67

สารบัญภาพ

ภาพที่		หน้า
1.1	แผนการดำเนินงานวิจัย.....	4
2.1	ลักษณะของเสียงวรรณยุกต์ในภาษาไทย.....	10
2.2	ระบบการแปลงข้อความเป็นสัทอักษรภาษาไทย.....	11
3.1	แผนภาพความสัมพันธ์ของเครื่องมือวิจัย.....	23
3.2	ภาพรวมการทำงานหลักของระบบ.....	24
3.3	การสร้างการกระจายตัวทางหน่วยเสียงเป้าหมาย.....	25
3.4	แผนภาพการกระจายตัวทางสถิติทางหน่วยเสียงของฐานข้อมูล LOTUS.....	29
3.5	การดึงข้อความออนไลน์ด้วยตัวดึงหน้าเว็บแยกวิเคราะห์ HTML.....	30
3.6	ตัวอย่างหน้าเว็บจากอินเทอร์เน็ตที่มีเนื้อความที่เหมาะกับการดึงข้อความ.....	31
3.7	ตัวอย่างข้อความออนไลน์จากอินเทอร์เน็ตที่เก็บไว้ในฐานข้อมูล.....	31
3.8	การแปลงข้อความในรูปแบบบทความเป็นรูปแบบประโยค.....	32
3.9	การแปลงประโยคจากรูปเขียนเป็นรูปอ่าน.....	36
3.10	ตัวอย่างประโยคออนไลน์ในรูปเขียนและรูปอ่านที่เก็บไว้ในฐานข้อมูล.....	37
3.11	หน้าต่างโปรแกรมขั้นตอนการสร้างประโยคออนไลน์.....	39
3.12	แผนภาพขั้นตอนการเลือกประโยคและจัดเก็บคลังข้อความ.....	40
3.13	ใช้ค่าความต้องการรูปแบบหน่วยเสียงในการให้คะแนนประโยค.....	45
4.1	แผนภาพการกระจายตัวทางสถิติทางหน่วยเสียงของฐานข้อมูล LOTUS.....	52
4.2	ผลของค่าเฉลี่ยผลต่างของการกระจายตัวทางหน่วยเสียงเป้าหมายเทียบกับจำนวน ประโยคออนไลน์ที่เข้ามาใหม่.....	55
4.3	ความครอบคลุมทางหน่วยเสียงเทียบกับจำนวนประโยคออนไลน์ที่เข้ามาใหม่..	56
4.4	แผนภาพการกระจายตัวของหน่วยหลังจากการคัดเลือกประโยค.....	57
4.5	แผนภาพการกระจายตัวของหน่วยหลังจากการคัดเลือกประโยคใช้การกระจาย ตัวหน่วยเสียงเป้าหมายกำหนดเอง.....	59
4.6	ความครอบคลุมทางหน่วยเสียงเมื่อกำหนดให้ค่า α เท่ากับ 0.0.....	60
4.7	ผลของค่าเฉลี่ยผลต่างของการกระจายตัวทางหน่วยเสียงเป้าหมายเมื่อกำหนดให้ ค่า α เท่ากับ 0.0.....	61

4.8	ความครอบคลุมทางหน่วยเสียงเมื่อกำหนดให้ค่า α เท่ากับ 1.5.....	61
4.9	ผลของค่าเฉลี่ยผลต่างของกระจายตัวทางหน่วยเสียงเป้าหมายเมื่อกำหนดให้ ค่า α เท่ากับ 1.5.....	62
4.10	ความครอบคลุมทางหน่วยเสียงเมื่อกำหนดให้ค่า α เท่ากับ 3.0.....	62
4.11	ผลของค่าเฉลี่ยผลต่างของกระจายตัวทางหน่วยเสียงเป้าหมายเมื่อกำหนดให้ ค่า α เท่ากับ 3.0.....	63
4.12	ความครอบคลุมทางหน่วยเสียงเมื่อเปรียบเทียบวิธีการเลือกประโยคของ วิทยานิพนธ์กับวิธีการเลือกประโยคแบบเก่า.....	64
4.13	ผลของค่าเฉลี่ยผลต่างของกระจายตัวทางหน่วยเสียงเป้าหมายเมื่อเปรียบเทียบ วิธีการเลือกประโยคของวิทยานิพนธ์กับวิธีการเลือกประโยคแบบเก่า.....	65
4.14	แผนภาพการกระจายตัวของหน่วยหลังจากการคัดเลือกประโยคของ วิทยานิพนธ์เปรียบเทียบกับวิธีของ Zhang.....	66

บทที่ 1

บทนำ

ความเป็นมาและความสำคัญของปัญหา

เทคโนโลยีการโต้ตอบกันระหว่างคอมพิวเตอร์กับมนุษย์กำลังเป็นที่สนใจของบรรดานักวิจัยและหน่วยงานต่าง ๆ ทำให้เกิดนวัตกรรมด้านเสียงพูดออกมามากขึ้น เช่น ระบบรู้จำเสียงพูดที่ใช้งานเพื่อใช้ในการโต้ตอบกันระหว่างมนุษย์กับหุ่นยนต์ [1] โปรแกรมช่วยอ่านหน้าจอคอมพิวเตอร์ [2] ระบบรายงานสภาพการจราจรด้วยเสียงอัตโนมัติ: Traffic Voice Information Services, TVIS [3] หรือการสังเคราะห์เสียงพูดบนอุปกรณ์พกพา [4] สิ่งที่เป็นสำหรับเทคโนโลยีการรู้จำเสียงพูดอัตโนมัติและการสังเคราะห์เสียงพูดนี้คือ “คลังข้อความ” ซึ่งประโยคจากคลังข้อความนี้จะถูกนำไปอ่านบันทึกเสียงโดยอาสาสมัครผู้ให้เสียง เพื่อสร้างระบบการรู้จำเสียงพูดอัตโนมัติหรือระบบสังเคราะห์เสียงพูดดังกล่าว อย่างไรก็ตามประสิทธิภาพระบบรู้จำเสียงพูดและคุณภาพของเสียงพูดสังเคราะห์ขึ้นอยู่กับจำนวนหน่วยเสียงที่ครอบคลุมในฐานข้อมูล [5] จึงจำเป็นต้องใช้เวลาในการบันทึกเสียงค่อนข้างมากเนื่องจากต้องการให้เกิดความครอบคลุมของหน่วยเสียงที่มากที่สุด ทั้งนี้คำศัพท์และสำนวนใหม่ ๆ เกิดขึ้นทุกวัน เราจึงต้องการคลังข้อความที่ทันสมัย เพื่อที่จะนำมาสังเคราะห์เสียงพูดหรือรู้จำเสียงพูดในปัจจุบันได้ การเลือกประโยคที่เหมาะสม สามารถลดเวลาในขั้นตอนการบันทึกเสียงเพื่อการสร้างคลังเสียงพูด และได้รูปแบบหน่วยเสียงที่ครอบคลุมและเกิดการกระจายตัวของหน่วยเสียงที่ดี สามารถคาดเดาคุณภาพของคลังข้อความได้ ทำให้เกิดงานวิจัยด้านการเลือกประโยค เพื่อสร้างคลังข้อความที่เหมาะสม

มีงานวิจัยที่ใช้การสร้างประโยค ด้วยวิธีการสุ่มประโยคเพื่อสร้างคลังข้อความ [6] คือ การเลือกประโยคจากหลาย ๆ บทความ มาบรรจุในคลังข้อความแบบสุ่ม เป็นวิธีที่นิยมให้กันอย่างแพร่หลายในสมัยก่อน ปัญหาที่เกิดขึ้นจากวิธีนี้คือ การขาดความหลากหลายของรูปแบบหน่วยเสียง และใช้เวลาในการอัดเสียงที่มากและไม่สามารถคาดการณ์คุณภาพของคลังเสียงได้ ส่งผลให้เกิดงานวิจัยด้านการเลือกประโยคเพื่อสร้างคลังข้อความ คำนี้ถึงมีการกระจายตัวของหน่วยเสียงสมดุลโดยใช้ขั้นตอนเชิงละโมภ: Greedy algorithm [7] ในการเลือกประโยคที่เหมาะสม [8, 9,10] โดยการเลือกประโยคที่ประกอบด้วยหน่วยเสียงที่ทำให้เกิดความครอบคลุมทางหน่วยเสียงมากที่สุดจากคลังข้อความแม้ คลังข้อความแม่มายังถึงแหล่งข้อมูลขนาดใหญ่ที่ประกอบไปด้วยประโยคจำนวนมาก แต่ปัญหาของวิธีนี้คือ ไม่สามารถควบคุมการกระจายตัวของหน่วยได้ และด้วยขนาดของข้อมูลที่ใช้ มีขนาดคงที่ ทำให้ประสิทธิภาพของการเลือกประโยคถูกจำกัด จึงมีบาง

งานวิจัยที่ใช้ข้อมูลอินเทอร์เน็ตมาสร้างคลังข้อความ [11] แต่งานวิจัยนี้ได้เก็บข้อมูลจากอินเทอร์เน็ตมาจำนวนหนึ่งโดยไม่มีเก็บเพิ่มอีก ทำให้การเลือกประโยคถูกจำกัดด้วยจำนวนข้อมูล

ในงานวิจัยนี้เราจึงเสนอ การสร้างคลังข้อความจากการกระจายตัวของหน่วยเสียงที่กำหนด และใช้วิธีคัดเลือกประโยคจากอินเทอร์เน็ตอย่างต่อเนื่อง การกระจายตัวทางหน่วยเสียงที่กำหนด ในที่นี้หมายถึง คลังข้อความที่ได้ มีจำนวนรูปแบบของแต่ละหน่วยเสียง เป็นไปตามอัตราส่วนเป้าหมาย และยังสามารถกำหนดเกณฑ์ของจำนวนขั้นต่ำของหน่วยเสียง การคัดเลือกประโยคจะใช้วิธีการประยุกต์ขั้นตอนเชิงละโมบ (Greedy Algorithm) คือ ต้องการประโยคที่ทำให้การกระจายตัวทางหน่วยเสียงเป็นไปตามที่ต้องการและเกิดความครอบคลุมทางหน่วยเสียงด้วย ใช้ค่าความต้องการของรูปแบบหน่วยเสียงที่เกิดขึ้นในประโยคในการให้คะแนน เพื่อทำให้เกิดการกระจายตัวตามที่กำหนดได้จากรูปแบบหน่วยเสียงที่อยู่ในประโยคที่ถูกเลือก ทั้งนี้ยังสามารถกำหนดลักษณะของประโยคในคลังข้อความได้ คือ สามารถกำหนดจำนวนประโยคในคลังข้อความ กำหนดลักษณะการกระจายตัวทางหน่วยเสียง กำหนดรูปแบบหน่วยเสียง กำหนดขอบเขตความยาวของประโยค เพื่อให้ได้คลังข้อความตามที่ต้องการใช้ต้องการได้

วัตถุประสงค์ของการวิจัย

เพื่อเสนอวิธีในการสร้างคลังข้อความ จากการกระจายตัวรูปแบบหน่วยเสียงที่กำหนด เพื่อได้คลังข้อความที่มีค่าการกระจายตัวทางสถิติทางหน่วยเสียงตามที่ต้องการ สร้างประโยคออนไลน์จากบทความอินเทอร์เน็ตเพื่อความทันสมัยของประโยค ใช้วิธีการให้คะแนนและเลือกประโยคที่มีค่าความต้องการรูปแบบหน่วยเสียงมาก เพื่อให้ลักษณะการกระจายตัวทางหน่วยเสียงและความครอบคลุมหน่วยเสียงเป็นไปตามที่กำหนดได้

ขอบเขตของการวิจัย

1. ในการทดลองประโยคที่คัดเลือกเป็นมาต้องเป็นประโยคที่ประกอบด้วยตัวอักษรภาษาไทยเท่านั้น
2. ใช้รูปแบบหน่วยเสียงชนิด หน่วยเสียงคู่ ในการทดลองสร้างฐานข้อมูลข้อความ
3. ในการทดลองใช้ตัวแปลงรูปเขียนเป็นรูปอ่านไทยแบบ Probabilistic Generalized LR parser [12]

4. ในการทดลองใช้ชุดประโยคจากคลังข้อมูลข้อความ ฐานข้อมูล Large Vocabulary Continuous Speech Recognition: LVCSR [13] เลือกใช้ชุดประโยคที่ครอบคลุมภาษาไทย 5,000 คำ ในการสร้างการกระจายตัวของหน่วยเสียงเป้าหมาย

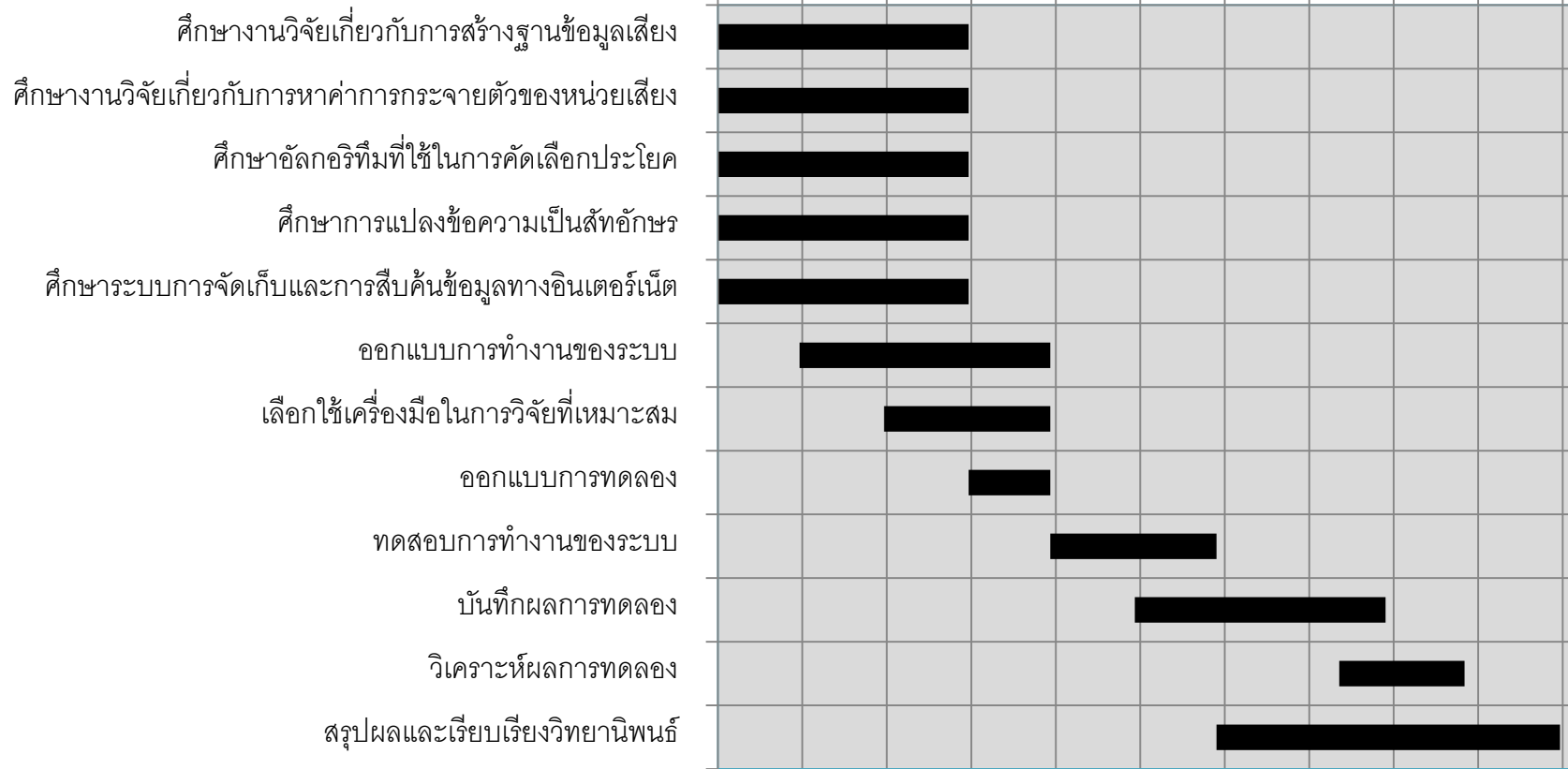
ลำดับขั้นตอนในการเสนอผลการวิจัย

1. ขั้นตอนการศึกษาเบื้องต้น
 - 1.1. ศึกษางานวิจัยเกี่ยวกับการสร้างฐานข้อมูลเสียง
 - 1.2. ศึกษางานวิจัยเกี่ยวกับการหาค่าการกระจายตัวของหน่วยเสียง
 - 1.3. ศึกษาอัลกอริทึมที่ใช้ในการคัดเลือกประโยค
 - 1.4. ศึกษาการแปลงข้อความเป็นรูปอ่าน
 - 1.5. ศึกษาระบบการจัดเก็บและการสืบค้นข้อมูลทางอินเทอร์เน็ต
2. ขั้นตอนการออกแบบระบบ
 - 2.1. ออกแบบการทำงานของระบบ
 - 2.2. เลือกใช้เครื่องมือในการวิจัยที่เหมาะสม
 - 2.3. ออกแบบการทดลอง
3. ขั้นตอนการทดลอง
 - 3.1. ทดสอบการทำงานของระบบ
 - 3.2. บันทึกผลการทดลอง
 - 3.3. วิเคราะห์ผลการทดลอง
4. สรุปผลและเรียบเรียงวิทยานิพนธ์

ประโยชน์ที่คาดว่าจะได้รับ

1. ได้ระบบการสร้างคลังข้อความอัตโนมัติตามรูปแบบการกระจายตัวของหน่วยเสียงที่กำหนดได้
2. ประโยคที่ได้มีใจความต่อเนื่องและทันสมัยเหมาะกับการนำไปใช้
3. คาดเดาคุณภาพคลังข้อความได้
4. ช่วยลดเวลาการบันทึกเสียง เมื่อใช้ประโยคจากคลังข้อความ
5. สามารถนำระบบนี้ไปประยุกต์ใช้กับสร้างฐานข้อมูลข้อความในภาษาอื่นได้

มี.ย. 54 ก.ค. 54 ส.ค. 54 ก.ย. 54 ต.ค. 54 พ.ย. 54 ธ.ค. 54 ม.ค. 55 ก.พ. 55 มี.ค. 55 เม.ย. 55



ภาพที่ 1.1 แผนการดำเนินงานวิจัย

ผลงานตีพิมพ์จากวิทยานิพนธ์

ส่วนหนึ่งของวิทยานิพนธ์นี้ได้ตีพิมพ์เป็นบทความทางวิชาการในหัวข้อเรื่อง “Automatic Online Text Selection for Constructing Text Corpus with Custom Phonetic Distribution” จัดทำโดย “Surapol Vorapatratom, Atiwong Suchato, Proadpran Punyabukkana” ถูกนำเสนอในงานประชุมวิชาการ “The Ninth International Joint Conference on Computer Science and Software Engineering: JCSSE'2012” ณ มหาวิทยาลัยหอการค้าไทย ประเทศไทย ในวันที่ 31 พฤษภาคม 2555

ลำดับการจัดเรียงเนื้อหาในวิทยานิพนธ์

ในวิทยานิพนธ์นี้ได้แบ่งเนื้อหาออกเป็น 5 บท คือ บทที่ 1 บทนำ กล่าวถึง ความเป็นมาและความสำคัญของปัญหา วัตถุประสงค์ของการวิจัย ขอบเขตของการวิจัย ลำดับขั้นตอนในการเสนอผลการวิจัย และผลงานตีพิมพ์จากวิทยานิพนธ์ บทที่ 2 ทฤษฎีและงานวิจัยที่เกี่ยวข้อง กล่าวถึง แนวคิดและทฤษฎี ประกอบด้วย ระบบการเขียนภาษาไทย ระบบการออกเสียงภาษาไทย รูปแบบหน่วยเสียง การแปลงข้อความภาษาไทยเป็นสัทอักษร การค้นหาสารสนเทศบนอินเทอร์เน็ต โครงสร้างคำสั่งในภาษา HTML ขั้นตอนวิธีเชิงละโมบ และการวัดระยะทางแบบยุคลิด อีกส่วนคืองานวิจัยที่เกี่ยวข้อง บทที่ 3 ขั้นตอนการดำเนินงานวิจัย กล่าวถึง เครื่องมือที่ใช้ในงานวิจัย ขั้นตอนการดำเนินงานวิจัย ประกอบด้วย ขั้นตอนการสร้างการกระจายตัวทางหน่วยเสียงเป้าหมาย ขั้นตอนการสร้างประโยคออนไลน์ ขั้นตอนการเลือกประโยคและจัดเก็บคลังข้อความ บทที่ 4 การทดลอง และอภิปรายผล กล่าวถึง การทดลองของวิทยานิพนธ์นี้ และผลการทดลอง บทที่ 5 บทสรุปผลการวิจัยและข้อเสนอแนะ กล่าวถึง การสรุปผลการวิจัย ข้อเสนอในจุดเด่นจุดด้อยของงานวิจัย และงานวิจัยในอนาคต

บทที่ 2

ทฤษฎีและงานวิจัยที่เกี่ยวข้อง

แนวคิดและทฤษฎี

ในงานวิจัยนี้ทฤษฎีที่เกี่ยวข้องแบ่งออกเป็น ระบบการเขียนภาษาไทย, ระบบการออกเสียงภาษาไทย, การแปลงข้อความภาษาไทยเป็นรหัสอักษร, การค้นหาสารสนเทศบนอินเทอร์เน็ต, ขั้นตอนวิธีเชิงละโมบ: Greedy algorithm และการวัดระยะทางแบบยูคลิด: Euclidean distance

1. ระบบการเขียนภาษาไทย

ระบบการเขียนของภาษาไทยเป็นภาษาที่มีระดับของเสียง [14] พยัญชนะไทยมี 44 ตัว ได้แก่ ก ข ฃ ค ฅ ง จ ฉ ช ฌ ญ ฎ ฏ ฐ ฑ ฒ ด ต ถ ท ธ น บ ป ผ ฝ พ ฟ ภ ม ย ร ล ว ศ ษ ส ห พ อ ฮ มีสระพื้นฐาน 15 ตัว ได้แก่ -ะ -า ิ ี ื ึ ุ ู -แ -เ -โ -ไ -อ -ฮ มีเครื่องหมายบอกระดับเสียงหรือวรรณยุกต์ 4 ตัว (่, ้, ๊, ๋) ข้อความถูกเขียนจากทางด้านซ้ายไปขวาตามแนวนอน ไม่มีการเว้นวรรคแต่ละคำ สระสามารถอยู่ด้านบน, ด้านล่าง, ด้านซ้ายหรือด้านขวาของพยัญชนะ หลักไวยากรณ์ของภาษาไทยเมื่อเทียบกับภาษาของทางโลกตะวันตกมีความง่ายกว่ามาก ในแต่ละประโยคประกอบด้วย “ประธาน+กริยา+กรรม” ไม่มีคำนำหน้านาม ไม่ต้องเปลี่ยนรูปกริยาจากคำบอกเวลาหรือจำนวน ไม่มีเครื่องหมายจบประโยค ไม่มีการเว้นวรรคระหว่างคำ เช่น “น้ำหอมอบอวล” ประโยคนี้สามารถตัดคำได้อย่างหลากหลาย “น้ำ-หอม-อบ-อวล” หรือ “น้ำ-หอม-อบ-อวน” ซึ่งการแก้ปัญหาในส่วนใหญ่นี้จะใช้การดูความหมายจากประโยครอบข้าง “กลิ่นน้ำหอมอบอวนเต็มห้อง” เป็นต้น [15] งานวิจัยนี้ใช้คำในภาษาไทยในการทดลองจึงเลือกใช้รูปแบบการเขียนภาษาไทยดังกล่าวเป็นมาตรฐานการเขียนประโยคภาษาไทย

2. ระบบการออกเสียงในภาษาไทย

การออกเสียงแต่ละพยางค์ในภาษาไทย [14] โดยทั่วไปแล้วมี 4 ส่วนคือ เสียงพยัญชนะต้น เสียงสระ เสียงตัวสะกด และเสียงวรรณยุกต์ พยัญชนะนำหน้าทีแทนด้วย “C” คือพยัญชนะเดี่ยว พยัญชนะนำหน้าทีแทนด้วย “CC” แทนด้วยพยัญชนะควบกล้ำ “V” แทนด้วยสระ, พยัญชนะท้ายทีแทนด้วย “Cf” คือตัวสะกดและตัวเลข (0-4) แทนด้วยเสียงวรรณยุกต์

โครงสร้างของพยางค์มี 4 รูปแบบคือ CV เช่น ปา [p-aa-0], วี [r-ii-0], CCV เช่น ปลา [pl-a-0], ครู [khr-uu-0], CVCf เช่น ปาด [p-aa-t⁻¹], กาก [k-aa-k⁻¹], CCVCf เช่น ปราบ [pr-aa-p⁻¹], กวาด [kw-aa-t⁻¹] พยัญชนะต้นในภาษาไทยมี 44 ตัว 21 เสียง มีสามระดับเสียงคือ อักษรสูง, อักษรกลางและอักษรต่ำ พยัญชนะต้นมีทั้งพยัญชนะเดี่ยว เช่น (ก, บ, ผ) พยัญชนะควบกล้ำ (ปล-, คว-, ตร-) และบางตัวยืมมาจากต่างประเทศ (ดร-, ฟร-, บร-) สามารถดูการออกเสียงของพยัญชนะไทยได้ดังตารางที่ 2.1 พยัญชนะควบกล้ำ ตารางที่ 2.2 และ หน่วยเสียงยืมภาษาต่างประเทศ ดังตารางที่ 2.3

ตารางที่ 2.1 การออกเสียงของพยัญชนะภาษาไทย

ตัวอักษรไทย	หน่วยเสียงพยัญชนะต้น	หน่วยเสียงตัวสะกด
ก	k	k [^]
ข ฃ ค ฅ ฆ	kh	-
ง	ng	ng [^]
จ	c	-
ฉ ฌ	ch	-
ซ ศ ษ ส	s	-
ญ ย	j	j [^]
ฎ ด	d	t [^]
ฏ ต	t	-
ฐ ฑ ฒ ถ ท ฒ	th	-
ณ น	n	n [^]
บ	b	p [^]
ป	p	-
พ ภ ผ	ph	-
ฟ ฝ	f	-
ม	m	m [^]
ร	r	-
ล ฬ	l	-
ว	w	w [^]
ห ฮ	h	-
อ	z	-

ตารางที่ 2.2 การออกเสียงของพยัญชนะควบกล้ำไทย

ตัวอักษรไทย	หน่วยเสียงพยัญชนะต้น	หน่วยเสียงตัวสะกด
ปร-	pr	-
ปล-	pl	-
ตร-	tr	-
กร-	kr	-
กล-	kl	-
กว-	kw	-
พร-, ผล-	phr	-
ทร	phl	-
คร-, ขร-	khr	-
คล-, ขล-	khl	-
คว	khw	-

ตารางที่ 2.3 การออกเสียงของพยัญชนะหน่วยเสียงยืมภาษาต่างประเทศ

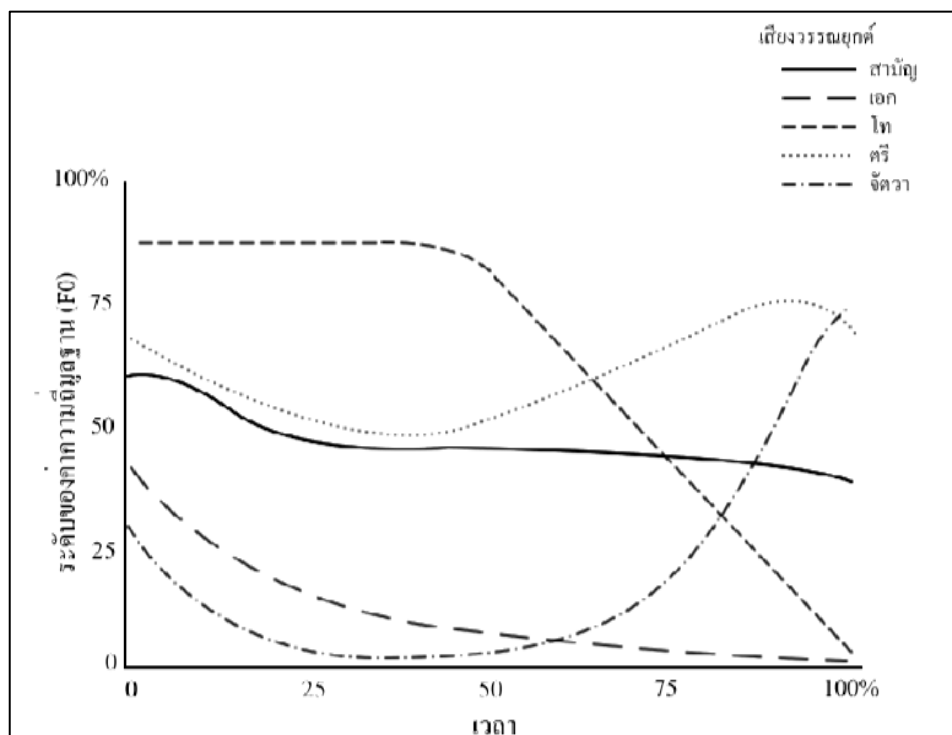
ตัวอักษรไทย	หน่วยเสียงพยัญชนะต้น	หน่วยเสียงตัวสะกด
бр-	br	-
บล-	bl	-
ฟร-	fr	-
ฟล-	fl	-
ดร-	dr	-
-ฟ	-	f [^]
-ล	-	l [^]
-ส	-	s [^]
-ช	-	ch [^]

สระในภาษาไทยมีสองชนิดคือ สระเดี่ยวและสระผสม สระเดี่ยวมี 9 ตัวที่เป็นเสียงยาว เช่น -า/a/, - ู /-ยู, /,แ-x / และอีก 9 ตัวเป็นสระเสียงสั้นเช่น -ะ/a/, - ุ /u/,แ-x;/ สระผสมมี ทั้งหมด 3 ตัวที่ออกเสียงยาว คือ - ี ย / i; a/ , - ี อ / v; a/ , - ัว / u; และอีก 3 ตัวที่ออกเสียงสั้นคือ - ี ยะ / ia/ , - ี อะ / va/ , - ัวะ / ua/ สามารถดูการออกเสียงของสระจาก ตารางที่ 2.4 และสามารถเทียบวิธีการอ่านหน่วยเสียงทั้งหมดกับสัทอักษรสากล จากภาคผนวก ก

ตารางที่ 2.4 การออกเสียงสระไทย

สระเดี่ยว				สระผสม			
เสียงสั้น		เสียงยาว		เสียงสั้น		เสียงยาว	
รูปสระ	หน่วยเสียง	รูปสระ	หน่วยเสียง	รูปสระ	หน่วยเสียง	รูปสระ	หน่วยเสียง
ะ	a	า	aa	ียะ	ia	ี ย	ii a
ิ	i	ี	ii	ีอะ	va	ี อ	vva
ึ	v	ือ	vv	ัวะ	ua	ัว	uuv
ุ	u	ู	uu				
เะ	e	เ	ee				
แะ	x	แ	xx				
โะ	o	โ	oo				
เาะ	@	-อ	@@				
เอะ	q	เ-อ	qq				
สระผสม							
เสียงสั้น		เสียงยาว					
รูปสระ	หน่วยเสียง	รูปสระ	หน่วยเสียง				
ำ	am	-	-				
ไ, ไ	ai	-	-				
เา	aw	-	-				

เสียงวรรณยุกต์ในภาษาไทยแบ่งออกเป็น 5 เสียงแสดงในลำดับตัวเลข 0-4 ได้แก่ เสียงสามัญ (Mid Tone: 0), เสียงเอก (Low Tone: 1), เสียงโท (Falling tone: 2), เสียงตรี (High Tone, 3) และเสียงจัตวา (Rising Tone: 4) ซึ่งแต่ละเสียงมีลักษณะเด่นดังแสดงในภาพที่ 2.1



ภาพที่ 2.1 ลักษณะของเสียงวรรณยุกต์ในภาษาไทย [7]

ในงานวิจัยนี้เลือกใช้ระบบการออกเสียงภาษาไทยดังกล่าวเป็นมาตรฐานในการแปลงรูปเขียนเป็นรูปอ่าน ผู้ที่ต้องการนำงานวิจัยนี้ไปประยุกต์ใช้กับหน่วยเสียงภาษาอื่นก็ทำได้เช่นกัน

3. รูปแบบหน่วยเสียง

หน่วยเสียงคือหน่วยที่เล็กที่สุดของการออกเสียงเสียงในแต่ละภาษา รวมถึงหน่วยเสียงเงียบ "silence" ที่แสดงถึงการที่ไม่ได้เปล่งเสียงใด ๆ ในช่วงเวลานั้น รูปแบบหน่วยเสียงโดยทั่วไปที่มีลักษณะต่างกันไปตามการใช้งาน ตัวอย่างรูปแบบหน่วยเสียงชนิดต่าง ๆ แสดงดังตารางที่ 2.5 หน่วยเสียงเดี่ยว ประกอบไปด้วยหน่วยที่เล็กที่สุดของหน่วยเสียง 1 หน่วย หน่วยเสียงคู่ ประกอบไปด้วยลำดับของหน่วยเสียงเดี่ยว 2 หน่วยเสียง หน่วยเสียงสาม ประกอบไปด้วยลำดับของหน่วยเสียงเดี่ยว 3 หน่วยเสียง และหน่วยเสียงที่มากกว่า 3 หน่วยเสียงขึ้นไป ซึ่งทุกรูปแบบอาจมีรูปแบบ

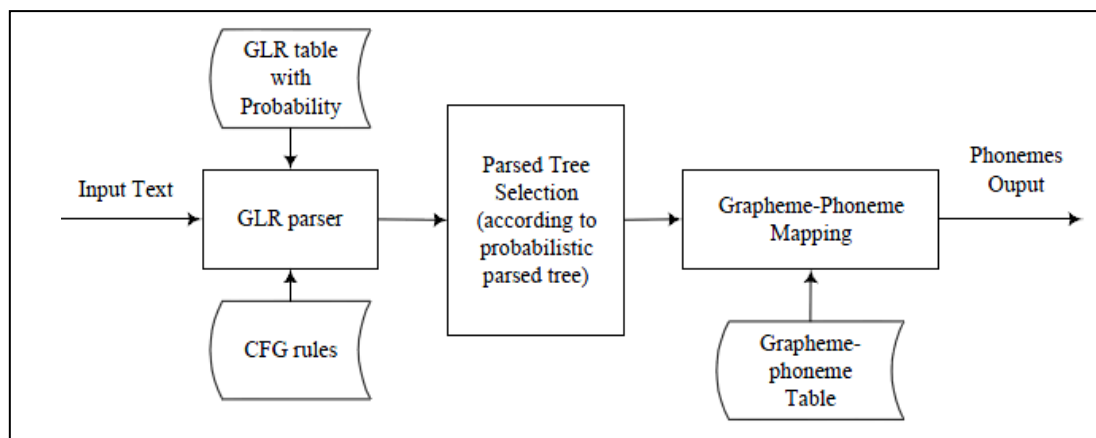
เสียงวรรณยุกต์กำกับด้วย ในการทดลองของงานวิจัยนี้ ได้ใช้ “หน่วยเสียงคู่” แบบไม่มีวรรณยุกต์ เป็นรูปแบบหน่วยเสียง

ตารางที่ 2.5 ตัวอย่างรูปแบบหน่วยเสียง

รูปแบบหน่วยเสียง	ตัวอย่างคำ	ลำดับหน่วยเสียง
หน่วยเสียงเดี่ยว	ฮัลไหล	sil, h, e, l [^] , l, o:, sil
หน่วยเสียงคู่		sil+h, h+e, e+l [^] , l [^] +l, l+o:, o:+sil
หน่วยเสียงสาม		sil-h+e, h-e+l [^] , e-l [^] +l, l [^] -l+o:, l-o:+sil...
หน่วยเสียงห้า		sil_sil-h+e=l [^] , sil_h-e+l [^] =l, h_e-l [^] =l...
หน่วยเสียงที่มี วรรณยุกต์กำกับ	สวัสดี	sil, s1, a1, w1, a1, ^t1, d0, i:0, sil

4. การแปลงข้อความภาษาไทยเป็นสัทอักษร

การแปลงรูปเขียนเป็นสัทอักษร คือ การแปลง ลำดับรูปเขียนของคำใด ๆ เป็นรูปแบบคำอ่าน เป็นหนึ่งในสิ่งสำคัญในการใช้กับการพัฒนา การสังเคราะห์ มีหลายวิธีในการแปลงแปลงรูปเขียนเป็นสัทอักษร เช่น การใช้พจนานุกรม (Dictionary-based) เป็นการแปลงรูปเขียนเป็นสัทอักษรโดยการดูคำอ่านจากพจนานุกรม วิธีนี้ต้องการ คลังคำศัพท์ที่ใหญ่และอาจเกิดปัญหาในการหาคำศัพท์ไม่เจอได้อีกด้วย การใช้กฎ (Rule-based) เป็นการแปลงรูปเขียนเป็นสัทอักษรโดยการจัดเตรียม โครงสร้างของกลุ่มพยางค์ นำไปใช้กับเครื่องสถานะจำกัด (Finite state machine) แต่วิธีนี้อาจเกิดข้อผิดพลาดได้เมื่อเจอการออกเสียงที่คลุมเครือของภาษาไทยได้ การใช้ต้นไม้ตัดสินใจ (Decision-tree) เป็นการแปลงรูปเขียนเป็นสัทอักษรโดยการใช้ต้นไม้ตัดสินใจ วิธีนี้ยังคงเกิดปัญหาในการแปลงการอ่านตามขอบเขตของพยางค์อยู่ การใช้ PGLR (Probabilistic Generalized LR) [12] เป็นเทคนิคการนำค่าความน่าจะเป็นมารวมกับต้นไม้ตัดสินใจของประโยค เป็นวิธีการที่ได้ผลการแปลงแม่นยำที่สุด เมื่อเทียบกับวิธีการใช้พจนานุกรมการใช้กฎหรือการใช้ต้นไม้ตัดสินใจ หลักการทำงานแสดงดังภาพที่ 2.2



ภาพที่ 2.2 ระบบการแปลงข้อความเป็นสัทอักษรภาษาไทย [10]

ในงานวิจัยนี้จึงนำกระบวนการแปลงข้อความภาษาไทยเป็นสัทอักษรดังกล่าว มาใช้ในการแปลงรูปเขียนเป็นรูปอ่าน ทั้งในขั้นตอนการสร้างการกระจายตัวเป้าหมายและขั้นตอนการให้คะแนนประโยคออนไลน์

5. การค้นหาสารสนเทศบนอินเทอร์เน็ต

ข้อมูลสารสนเทศบนอินเทอร์เน็ตเป็นแหล่งข้อมูลทางอิเล็กทรอนิกส์ที่สำคัญและใหญ่ที่สุดที่มีการเปลี่ยนแปลงอยู่ตลอดเวลา [11] ดังนั้น ในการสืบค้นข้อมูลสารสนเทศบนอินเทอร์เน็ต ควรกำหนด วัตถุประสงค์ในการสืบค้น ให้ชัดเจนกำหนด Search engine ที่เหมาะสมกำหนด ช่วงเวลาที่ต้องการค้นหา จะทำให้การสืบค้นใช้เวลาสั้นๆ นอกจากนั้นเราควรระบุประเภทข้อมูลสารสนเทศที่ต้องการสืบค้น เช่น ข้อมูลที่เป็นข้อความ ภาพวาด ภาพเขียนหรือภาพลายเส้น ภาพไดอะแกรม ภาพถ่าย เสียง เสียงสังเคราะห์ เช่น เสียงดนตรี ภาพยนตร์ ภาพเคลื่อนไหวอะนิเมชัน เครื่องมือหรือโปรแกรม สำหรับการสืบค้น (Search Engine) มีอยู่มากมายและมีให้บริการ อยู่ตามเว็บไซต์ต่างๆ ที่ให้บริการการสืบค้นข้อมูลโดยเฉพาะ การเลือกใช้นั้นขึ้นกับประเภทของข้อมูลสารสนเทศที่ต้องการสืบค้น Search Engine ต่างๆ จะให้ข้อมูลที่มีความลึกในแง่มุมหรือศาสตร์ต่างๆ ไม่เท่ากัน

ข้อดีของการค้นหาข้อมูลจากอินเทอร์เน็ตคือ ขอบเขตของข้อมูลสารสนเทศกว้างขวาง มีความหลากหลาย ไร้พรมแดน ข้อมูลสารสนเทศที่สืบค้นได้มีความทันสมัย เนื่องจากผู้สร้างข้อมูลสามารถแก้ไข ปรับปรุงได้ง่ายและตลอดเวลา สะดวก ไม่มีข้อจำกัดในแง่ของเวลาและสถานที่

สามารถสืบค้นเวลาได้ทีใดก็ได้ สามารถสืบค้นได้ง่ายและรวดเร็วโดยอาศัย Search Engine การได้มาซึ่งข้อมูลผ่านอินเทอร์เน็ตใช้เวลาสั้นมาก เมื่อเทียบกับการส่งเอกสารวิธีอื่น ๆ การได้มาซึ่งข้อมูลนั้น ประหยัดทั้งเวลาและทรัพยากร จัดเป็นห้องสมุดที่ใหญ่ที่สุดในโลก ข้อมูลสารสนเทศที่สืบค้นมา มีประโยชน์มาก สามารถนำไปจัดหมวดหมู่ ทำฐานข้อมูล บรรณานุกรม และจัดการต่อได้โดยง่าย และมีความทันสมัย งานวิจัยนี้จึงใช้ การค้นหาข้อมูลทางอินเทอร์เน็ตในการค้นหา ประโยคที่จะนำมาพิจารณาในการเลือกคู่คลังข้อความ

6. โครงสร้างคำสั่งในภาษา HTML

โครงสร้างคำสั่งในภาษา HTML [17] จะมีรูปแบบโครงสร้างการเขียนแบ่งออกเป็น 3 ส่วน คือ ส่วนประกาศ ส่วนเนื้อเรื่อง และส่วนเนื้อหา

ส่วนประกาศ เป็นส่วนที่กำหนดให้เบราว์เซอร์ทราบว่า นี่คือภาษาเอชทีเอ็มแอล และจะต้องทำการแปลผลอย่างไรมีคำสั่งคู่เดียวคือ `<html>` และ `</html>` ประกาศที่หัวและท้ายไฟล์

ส่วนหัวเรื่อง (head) เป็นส่วนที่แสดงผลข้อความบนไตเติ้ลบาร์ของเบราว์เซอร์ และอาจมีคำสั่งสำหรับกำหนดรายละเอียดด้านเทคนิคอื่น ๆ อีก แทรกอยู่ระหว่างคำสั่ง `<head>` และ `</head>`

ส่วนเนื้อหา (body) เป็นส่วนที่มีความซับซ้อนมากที่สุด และสามารถใส่เทคนิคลูกเล่นเพื่อดึงดูดความสนใจจากผู้ชมได้มาก ความแตกต่างระหว่างเว็บไซต์ต่างๆ แสดงความมีฝีมือของผู้จัดทำ ศิลปะในการออกแบบจะอยู่ในส่วนนี้ทั้งหมด ซึ่งจะแทรกอยู่ระหว่างคำสั่ง `<body>` และ `</body>`

โครงสร้าง พื้นฐานของภาษา Computer เป็นส่วนที่สำคัญที่สุด ของการเขียนภาษา Computer โดยทั่วไปแล้ว มันจะต้องถูกเขียนขึ้นทุกครั้ง ภาษา HTML ก็เหมือนกับภาษา Computer ทั่วไป ที่มี โครงสร้าง พื้นฐานเฉพาะ ของมันคำสั่งของ HTML ส่วนมากจะถูกกำหนด อยู่ในเครื่องหมาย `<` และ `>` ซึ่งถูกเรียกว่า Tag สำหรับในส่วนของคำสั่ง Tag ภายในคำสั่ง ซึ่ง Tag แบ่งออกเป็น 2 ประเภท คือ แท็กเดี่ยว และแท็กคู่

แท็กเดี่ยว คือ คำสั่งที่มีคำสั่งเพียงอย่างเดียว ซึ่งสามารถใช้และสิ้นสุดคำสั่งได้ด้วยตัว
ของมันเอง เช่น

```
ข้อความ... <br>
```

```
<hr>
```

```
<! - ข้อความ - >
```

แท็กคู่ คือ คำสั่งที่ต้องมีส่วนเริ่มต้นและส่วนจุดจบของคำสั่งนั้นๆ โดยแท็กที่เป็นส่วนจบ
นั้นจะมีเครื่องหมาย slash / ติดเอาไว้ เช่น

```
<html> ส่วนของเนื้อหา ..... </html>
```

```
<center> ข้อความ..... </center>
```

```
<p> ข้อความ.... </p>
```

*** ถ้าหากมีการใช้แท็กคู่หลายๆ คำสั่ง เช่น คำสั่งตัวขีดเส้นใต้ <U> </U> และตาม
ด้วย คำสั่ง ตัวเอียง <I>....</I> จะต้องปิดคำสั่งตัวเอียงก่อน แล้วจึงจะมาปิดคำสั่งตัวหนา***

```
</I> <U> ข้อความ.... </U> </I>
```

สำหรับในส่วนของคำสั่ง Tag ภายในคำสั่ง โครงสร้างพื้นฐานพอที่จะอธิบายคร่าว ๆ ได้
ดังนี้

Title(ชื่อหัวเรื่อง) จะถูกกำหนด อยู่ภายใน Tag คำสั่ง

```
<HEAD> <TITLE> ชื่อหัวเรื่อง </TITLE> </HEAD>
```

ข้อมูลที่ถูก เขียนอยู่ใน Tag จะแสดงผล ออกมาให้เห็น ที่บนบาร์ของเว็บเบราว์เซอร์

ข้อมูลที่ต้องการแสดงผล จะเป็นส่วนที่แสดงให้เราเห็นไม่ว่าจะเป็น ตัวอักษร, รูปภาพ,
ตาราง ฯลฯ (คำสั่งที่ต้องการแสดงผลจะอยู่ระหว่าง tag BODY ทั้งหมด) ซึ่งถูกกำหนดอยู่ ระหว่าง
คำสั่ง

```
<BODY> จนถึงคำสั่ง </BODY>
```

คำสั่ง Comment Tag เป็นคำสั่งที่ใช้ในการอธิบายอยู่ใน HTML จะไม่มีการแสดงผล
ออกมาที่ Browser จะมีประโยชน์สำหรับผู้ที่ทำการแก้ไขโปรแกรมในภายหลัง

```
<!--ใส่ข้อความใดๆก็ได้ เพื่อใช้ในการ อธิบาย-->
```

คำสั่งขึ้นบรรทัดใหม่ เป็นคำสั่งที่ใช้กำหนดให้ข้อความที่เราพิมพ์ลงไปเอกสารขึ้นบรรทัดใหม่ได้ตามที่เราต้องการ เพราะถ้าเราไม่ใช้คำสั่งสั่งให้เอกสารแสดงผลขึ้นบรรทัดใหม่การแสดงผลของข้อความจะแสดงต่อกัน แม้ว่าเราจะพิมพ์ข้อความขึ้นบรรทัดใหม่ก็ตาม

คำสั่งการย่อหน้าใหม่ รูปแบบคำสั่ง

<P>..... </P> หรือ <P>

มีลักษณะคล้ายกับคำสั่ง < BR > แต่คำสั่งนี้จะมีการเว้นบรรทัดว่างให้หนึ่งบรรทัด เพราะบางครั้งเราต้องการเว้นบรรทัดว่างหนึ่งบรรทัดแต่โปรแกรม Web Browser จะไม่เข้าใจการพิมพ์ บรรทัดเปล่า

เส้นคั่นบรรทัด เป็นคำสั่งที่ใช้แบ่งข้อความของจอภาพให้เป็นส่วน ๆ รูปแบบคำสั่ง

<HR>

สำหรับโครงสร้างคำสั่งภาษา HTML สามารถแบ่งออกได้เป็น 2 ส่วนคือ ส่วนคำสั่งหัวเรื่อง และส่วนคำสั่งเนื้อความ

คำสั่งหัวเรื่อง เป็น คำสั่งชื่อโฮมเพจและข้อความอธิบายข้อมูลที่เกี่ยวข้องกับเว็บไซต์เวอร์ชันที่เป็นเจ้าของ โดยชื่อโปรแกรมดังกล่าวจะไปปรากฏบนเมนูของโปรแกรมเบราว์เซอร์ในขณะที่โปรแกรมได้รับการเชื่อมโยงแบบไฮเปอร์เท็กซ์ ดังนั้น หัวเรื่องจึงหมายถึงเอกลักษณ์ประจำโฮมเพจ เพราะเนื่องจากโปรแกรม HTML เป็นโปรแกรมของโฮมเพจ สามารถดูคำสั่งหัวเรื่องจากตารางที่ 2.6

ตารางที่ 2.6 คำสั่งหัวเรื่องโครงสร้าง HTML

คำสั่งหัวเรื่อง	คำอธิบาย
<TITLE>...</TITLE>	เพื่อแสดงไฟล์เอกสาร HTML หรือ ชื่อ HOMEPAGE
<ISINDEX>	เพื่อแสดงไฟล์เอกสาร HTML ชนิดที่สืบค้นได้
<NEXTID>	เพื่อแสดงเลขประจำตัวของไฟล์เอกสาร HTML
<LINK>	เพื่อกำหนดความสัมพันธ์ระหว่างไฟล์เอกสารนี้กับไฟล์เอกสารอื่นที่เกี่ยวข้อง
<BASE>	เพื่ออ้างอิงรหัสสืบค้น URL
<META>	เพื่อแสดงข่าวสารของไฟล์เอกสาร HTML

ส่วนคำสั่งเนื้อความ เป็นคำสั่งแสดงข้อความบนโฮมเพจ ซึ่งประกอบด้วยคำสั่งแสดงแบบของตัวอักษรของคำที่ใช้ในการอธิบาย คำสั่งที่ใช้ในการจัดวางหน้าของข้อความ คำสั่งเพื่อการเชื่อมโยงแบบไฮเปอร์ลิงค์ และคำสั่งเชื่อมโยงรูปภาพ เป็นต้น แสดงดังตารางที่ 2.7

ตารางที่ 2.7 คำสั่งส่วนเนื้อความโครงสร้าง HTML

คำสั่งเนื้อความ	คำอธิบาย
<H1>...</H1>	เพื่อกำหนดแบบหัวข้อให้เป็นตัวอักษรขนาดใหญ่ที่สุด
<H2>...</H2>	เพื่อกำหนดแบบหัวข้อให้เป็นตัวอักษรขนาดใหญ่อันดับสอง
<H3>...</H3>	เพื่อกำหนดแบบหัวข้อให้เป็นตัวอักษรขนาดใหญ่อันดับสาม
<H4>...</H4>	เพื่อกำหนดแบบหัวข้อให้เป็นตัวอักษรขนาดขนาดกลาง
<H5>...</H5>	เพื่อกำหนดแบบหัวข้อให้เป็นตัวอักษรขนาดเล็ก
<H6>...</H6>	เพื่อกำหนดแบบหัวข้อให้เป็นตัวอักษรขนาดเล็กที่สุด
<A>...	เพื่อสร้างไฮเปอร์สำหรับการเชื่อมโยง
<P>	เพื่อกำหนดย่อหน้าของข้อความ
 	เพื่อเว้นบรรทัดเมื่อจบข้อความ
<HR>	เพื่อขีดเส้นใต้ในแนวนอน
<PRE>...</PRE>	เพื่อกำหนดแบบอักษรของข้อความ
...	เพื่อแสดงรายการโดยไม่ต้องเรียงลำดับ
...	เพื่อแสดงรายการโดยเรียงลำดับ
	เพื่อแสดงข้อความแต่ละบรรทัดตามคำสั่งและ
<DL>...</DL>	เพื่อแสดงการอธิบายรายการ
<DT>	เพื่อแสดงคำที่ต้องการอธิบายภายใต้คำสั่ง<DL>
<DD>	เพื่อแสดงข้อความอธิบายคำที่กำหนดโดยคำสั่ง<DT>
...	เพื่อกำหนดตัวอักษรเป็นหนา
<I>...</I>	เพื่อกำหนดตัวอักษรเป็นเอน
<A>...	เพื่อกำหนดการเชื่อมโยงไฟล์HTMLสำหรับการโอนย้าย
	เพื่อแสดงภาพจากการเชื่อมโยง
<TABLE>..</TABLE>	เพื่อกำหนดการสร้างตาราง
<caption>...</caption>	เพื่อกำหนดข้อความอธิบายตาราง

คำสั่งเนื้อความ	คำอธิบาย
<TH>...</TH>	เพื่อกำหนดข้อความหัวเรื่องของตาราง
<TR>...</TR>	เพื่อกำกับการแสดงข้อความแต่ละแถวของตาราง
<TD>...</TD>	เพื่อกำหนดข้อความในแถวของตาราง
<FORM>...</FORM>	เพื่อกำหนดช่องว่างสำหรับกรอกข้อมูล
<INPUT>	เพื่อกำหนดข้อมูลสำหรับใช้ในคำสั่ง<FORM>...</FORM>
<textarea>..</textarea>	เพื่อกำหนดช่องว่างสำหรับกรอกข้อมูล

สำหรับการเชื่อมข้อมูล (LINK) ของโครงสร้าง HTML โดยทั่วไป มี 2 ประเภทคือ แบบ INTERNAL LINKS คือการเชื่อมโยงข้อมูลภายในไฟล์ของเราเอง และ แบบ EXTERNAL LINKS คือ การเชื่อมโยงข้อมูลออกไปสู่ภายนอกไฟล์ หรือ นอกโฮมเพจของเรา

คำสั่งการเชื่อมโยงข้อมูลภายใน (Internal Links) จะแบ่งออกเป็น 2 ส่วนด้วยกันคือ ชื่อเป้าหมายที่ต้องการแสดงผล รูปแบบคำสั่ง < A NAME="ชื่อเป้าหมาย" >.ข้อความ หรือ รูปภาพที่จะใช้เป็นตัวเป้าหมาย... < /A> และคำสั่งในการกำหนดการแสดงผลการเชื่อมโยงของเป้าหมาย รูปแบบคำสั่ง < A HREF="#ชื่อเป้าหมาย">.ข้อความ หรือ รูปภาพที่จะใช้เป็นตัว ลิงค์..

คำสั่งการเชื่อมโยงข้อมูลออกไปสู่ภายนอกไฟล์(External Links) มี 2 แบบ คือ การเชื่อมโยงภายในโฮมเพจ การกำหนดเป้าหมายก็เพียงแต่กำหนด ชื่อของไฟล์ก็เพียงพอแล้ว และการเชื่อมโยงออกไปสู่ภายนอกโฮมเพจ ซึ่งจะต้องกำหนดชื่อของของโฮมเพจที่เราต้องการเชื่อมโยงโดยละเอียดรูปแบบคำสั่ง < A HREF="URL"> ..ข้อความ หรือ รูปภาพที่จะใช้เป็นตัว ลิงค์...

สำหรับงานวิจัยนี้ได้ใช้โครงสร้าง HTML ในการเข้าถึงเนื้อข้อมูลในรูปแบบบทความของเว็บไซต์ต่าง ๆ และยังใช้ในการเชื่อมโยงไปยังหน้าถัดไปภายในเว็บไซต์นั้น ในขั้นตอนการสร้างประโยชน์ออนไลน์

7. ขั้นตอนวิธีเชิงละโมบ: Greedy Algorithm

เป็นขั้นตอนวิธีการแก้ปัญหาที่คิดแบบง่าย ๆ และตรงไปตรงมา โดยพิจารณาว่าข้อมูลที่มีอยู่ในขณะนั้นมีทางเลือกใดที่ ให้ผลตอบแทนคุ้มค่าที่สุด ขั้นตอนวิธีจะหาทางเลือกที่ดีที่สุด ในขณะที่ถ้าข้อมูลนั้นพอเพียงที่จะทำให้สรุปคำตอบที่ดีที่สุด เราจะได้ขั้นตอนวิธีที่มีประสิทธิภาพ โดยทั่วไปเราจะใช้ Greedy algorithm กับปัญหาเหมาะสมที่สุด Optimization problem เพราะว่าเราต้องการการตัดสินใจว่าทางเลือกในปัจจุบันมีค่าตอบแทนมากที่สุดหรือน้อยที่สุดหรือไม่ [7]

ในงานวิจัยนี้ ขั้นตอนวิธีเชิงละโมบ (Greedy Algorithm) ถูกนำมาใช้ในขั้นตอนการให้คะแนนของแต่ละประโยคเพื่อให้ได้ประโยคที่มีคะแนนสูงและมีการกระจายตัวในแต่ละประโยคที่เหมาะสม

8. การวัดระยะทางแบบยูคลิด: Euclidean distance

ระยะทางแบบยูคลิด (Euclidean distance, Euclidean metric) [17] คือการวัดระยะทางปกติระหว่างจุดสองจุดในแนวเส้นตรง นิยามของระยะทางแบบยูคลิดคือ

ระยะทางแบบยูคลิดระหว่างจุดสองจุด p และ q คือความยาวของส่วนของเส้นตรง pq ถ้า $p = (p_1, p_2, \dots, p_n)$ และ $q = (q_1, q_2, \dots, q_n)$ ในระบบพิกัดคาร์ทีเซียน เป็นจุดสองจุดบนปริภูมิยูคลิด n มิติ ระยะทางระหว่างจุด p กับ q คำนวณได้จากสมการที่ 1

$$d(\mathbf{p}, \mathbf{q}) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + \dots + (p_n - q_n)^2} = \sqrt{\sum_{i=1}^n (p_i - q_i)^2} \quad (1)$$

ค่าประจำแบบยูคลิด คือระยะทางจากจุดหนึ่งจุด p ไปยังจุดกำเนิด $(0, 0, \dots, 0)$ บนปริภูมิยูคลิด ได้ ดังสมการที่ 2

$$\|\mathbf{p}\| = \sqrt{p_1^2 + p_2^2 + \dots + p_n^2} = \sqrt{\mathbf{p} \cdot \mathbf{p}} \quad (2)$$

ซึ่งสมการตัวหลังเกี่ยวข้องกับผลคูณจุด เป็นขนาดของเวกเตอร์ p จากจุดกำเนิด ระยะทางแบบยูคลิดจึงอาจนิยามได้อีกแบบหนึ่ง ดังสมการที่ 3

$$\|\mathbf{p} - \mathbf{q}\| = \sqrt{(\mathbf{p} - \mathbf{q}) \cdot (\mathbf{p} - \mathbf{q})} = \sqrt{\|\mathbf{p}\|^2 + \|\mathbf{q}\|^2 - 2\mathbf{p} \cdot \mathbf{q}} \quad (3)$$

กรณีพิเศษในหนึ่งมิติ ระยะทางระหว่างจุดสองจุดบนเส้นจำนวนจริงคือค่าสัมบูรณ์ของผลต่างของสองค่า นั้น ดังนั้นถ้าให้ p และ q เป็นจุดสองจุด (หรือจำนวนสองจำนวน) บนเส้นจำนวนจริงแล้ว ระยะทางระหว่าง p และ q คำนวณได้ ดังสมการที่ 4

$$d(\mathbf{p}, \mathbf{q}) = \sqrt{(p - q)^2} = |p - q| \quad (4)$$

ในสองมิติแบบยูคลิด ถ้า $p = (p_1, p_2)$ และ $q = (q_1, q_2)$ แล้ว ระยะทางระหว่าง p และ q สามารถคำนวณได้ดังนี้ ซึ่งมีสูตรเหมือนกับทฤษฎีบทพีทาโกรัส ดังสมการที่ 5

$$d(\mathbf{p}, \mathbf{q}) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2} \quad (5)$$

จากนิยามแบบที่สองของระยะทางแบบยูคลิด ถ้าหาก $p = (r_1, \theta_1)$ และ $q = (r_2, \theta_2)$ ในระบบพิกัดเชิงขั้ว จะสามารถคำนวณระยะทางได้ ดังสมการที่ 6

$$\|\mathbf{p} - \mathbf{q}\| = \sqrt{r_1^2 + r_2^2 - 2r_1r_2 \cos(\theta_1 - \theta_2)} \quad (6)$$

ในสามมิติแบบยูคลิด ระยะทางระหว่าง p และ q ได้ ดังสมการที่ 7

$$d(\mathbf{p}, \mathbf{q}) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + (p_3 - q_3)^2} \quad (7)$$

ในงานวิจัยนี้เลือกใช้การวัดระยะทางแบบยูคลิดแบบการเปรียบเทียบค่าสัมบูรณ์ ของผลต่างของแต่ละจำนวนหน่วยเสียง จากสมการที่ 4 เพื่อใช้ในการหาค่าเบี่ยงเบนเฉลี่ยร้อยละจากเป้าหมาย (Percent average deviation of target) ในการเปรียบเทียบความแตกต่างกันระหว่างค่าการกระจายตัวทางหน่วยเสียงกำหนดกับค่าการกระจายตัวทางหน่วยเสียงปัจจุบัน

เอกสารและงานวิจัยที่เกี่ยวข้อง

Methods of Sentences Selection for Read-Speech Corpus Design (V. Radova, 1999) [18] งานวิจัยนี้เสนอวิธีการเลือกข้อความเพื่อนำไปสร้างระบบรู้จำเสียงพูดในภาษาเช็ก ใช้ประโยคจากหนังสือพิมพ์รายวันภาษาเช็ก จำนวนทั้งสิ้น 24,442 ประโยค มีการคัดเลือกประโยคเบื้องต้นที่ต้องประกอบด้วยคำมากกว่า 3 คำแต่ห้ามเกิน 15 คำ เพื่อความไม่เหลื่อมกันของคะแนน ไม่มีตัวเลข สัญลักษณ์ หรือภาษาอื่น ๆ ปนอยู่ ให้คะแนนแต่ละประโยคจากการกระจายตัวทางสถิติของหน่วย Tri-Phone คัดเหลือ 40 เหลือประโยคผลลัพธ์ ได้การครอบคลุมทางหน่วยเสียง Tri-Phone 73.66% จากงานวิจัยนี้จึงเกิดแนวคิดว่าการคัดเลือกประโยคที่มีการกระจายตัวทางหน่วยเสียงครอบคลุมจะทำให้ลดเวลาในการอัดเสียงได้

LOTUS: Phonetically Distributed Continuous Speech Corpus for Thai Language (C. Wutiwivatjai, 2002) [13] มีการเลือกข้อความเพื่อสร้างฐานข้อมูลเสียงพูดต่อเนื่องภาษาไทยครอบคลุมคำศัพท์จำนวน 5,000 คำ ใช้หลักการออกแบบฐานข้อมูลเสียงเดียวกับฐานข้อมูล JNAS [22] และ WSJAMO [24] การสร้างฐานข้อมูลเริ่มจากเลือกข้อความจากประโยคในฐานข้อมูลบทความ ORCHID (V. Somlertlamvanich, 1998) [20] ซึ่งประกอบไปด้วย 27,634 ประโยคโดยการเลือกประโยคมีขั้นตอนการสกัดประโยคที่ประกอบด้วยตัวอักษรภาษาต่างประเทศ ออก และแปลงประโยคในรูปเขียนให้เป็นรูปอ่าน ใช้รูปแบบหน่วยเสียงคู่ จากนั้นให้คะแนนประโยคจากความครอบคลุมทางหน่วยเสียงที่เกิดขึ้น จนกระทั่งความครอบคลุมทางหน่วยเสียงผ่านเกณฑ์ที่กำหนดไว้ ผลคือได้ความครอบคลุมทางหน่วยเสียงรูปแบบหน่วยเสียงคู่เป็น 90.9% ทั้งหมด 1,628 หน่วยเสียงคู่ งานวิจัยนี้ทำให้เกิดแนวคิดเรื่อง การใช้ความครอบคลุมทางหน่วยเสียงเป็นหลักในการให้คะแนนประโยคจนกระทั่ง ผลความครอบคลุมทางหน่วยเสียงผ่านเกณฑ์ที่กำหนดไว้

TSynC-1:Thai tagged speech corpus for speech synthesis (C. Hansakunbuntheung, 2003) [19] เลือกข้อความนำไปสร้างระบบสังเคราะห์เสียงพูดภาษาไทย ใช้ประโยคจากฐานข้อมูล ORCHID (V. Somlertlamvanich, 1998) [20] สารานุกรมสำหรับเด็ก บทความวิชาการจากงานสัมมนาด้านเทคโนโลยี 43,340 ประโยค มีการกำหนด

ขอบเขตจำนวนพยางค์ในแต่ละประโยค แต่ละประโยคประกอบด้วยอักษรไทยเท่านั้น ให้คะแนนแต่ละประโยคจากการกระจายตัวทางสถิติของหน่วยเสียงสาม ใช้ Greedy Algorithm คัดเลือกประโยค คัดเหลือ 5,200 ประโยค ผลลัพธ์ ได้การครอบคลุมทางหน่วยเสียงแบบ หน่วยเสียงสามเท่ากับ 39.6% งานวิจัยนี้ทำให้เกิดแนวคิดเรื่อง การนำ Greedy Algorithm มาใช้ในการคัดเลือกประโยคเป็นวิธีที่ดีในการคัดเลือก คลังข้อความที่มีขนาดจำกัดอาจเกิดปัญหาในการเพิ่มความครอบคลุมของหน่วยเสียง

Constructing a Phonetic-Rich Speech Corpus while Controlling Time-Dependent Voice Quality Variability for English Speech Synthesis (N. Jinfu, 2006) [25] การเลือกข้อความนำไปสร้างระบบสังเคราะห์เสียงพูดภาษาอังกฤษใช้ประโยคจากฐานข้อมูล BTEC 749,500 ประโยค และ หนังสือพิมพ์อังกฤษ 4,985,200 ประโยค ในหนึ่งประโยคมีจำนวนคำ 10-25 คำ ให้คะแนนแต่ละประโยคจากการกระจายตัวทางสถิติของหน่วยเสียงสาม คัดประโยคเหลือ 3,100 ประโยค ผลลัพธ์ได้การครอบคลุมทางหน่วยเสียงสาม 99.33% งานวิจัยนี้ทำให้เกิดแนวคิดเรื่องการเลือกข้อความจากแหล่งข้อมูลขนาดใหญ่ เพิ่มค่าความครอบคลุมทางหน่วยเสียงได้มากกว่า

TSynC-2: An intensive design of a Thai speech corpus (C. Wutiw WATCHAI, 2007) [14] การเลือกข้อความนำไปสร้างระบบสังเคราะห์เสียงพูดภาษาไทย ใช้ประโยคจากฐานข้อมูล TSynC-1 (C. HANSAKUNBUNTHEUNG, 2003) เป็นประโยคตั้งต้นในการเลือกประโยค 5,386 ประโยค เลือกประโยคโดยเลือกคะแนนที่สูงที่สุด ให้คะแนนโดยการเพิ่มคะแนนให้หน่วยเสียงคู่ที่มีน้อยและคะแนนยังถูกรบกวนด้วยจำนวนความยาวของหน่วยเสียงด้วย ผลลัพธ์คือได้ประโยคมาทั้งสิ้น 3,581 ประโยค 5,823 ชนิดของรูปแบบหน่วยเสียงคู่ จากนั้นนำประโยคที่ได้ไปสร้างโมเดลเสียงและสังเคราะห์ออกมา จากการประเมินได้คะแนนความพึงพอใจ อยู่ที่ระดับ 3.09 ซึ่งมากกว่าเสียงสังเคราะห์จาก TSynC-1 งานวิจัยนี้ทำให้เกิดแนวคิดเรื่อง การนำให้คะแนนกับหน่วยเสียงที่มีน้อยหรือหายาก จะทำให้การเลือกประโยคได้ผลลัพธ์ตามเป้าได้ดี

Automatic Construction for a TTS Corpus with Limited Text (W. Zhang, 2010) [11] เลือกข้อความนำไปสร้างระบบสังเคราะห์เสียงพูดภาษาอังกฤษ ใช้ประโยคจากฐานข้อมูลออนไลน์ VOA special English คัดมาแบบ Random 4,621 ประโยค (ใช้เป็น Mother Text)

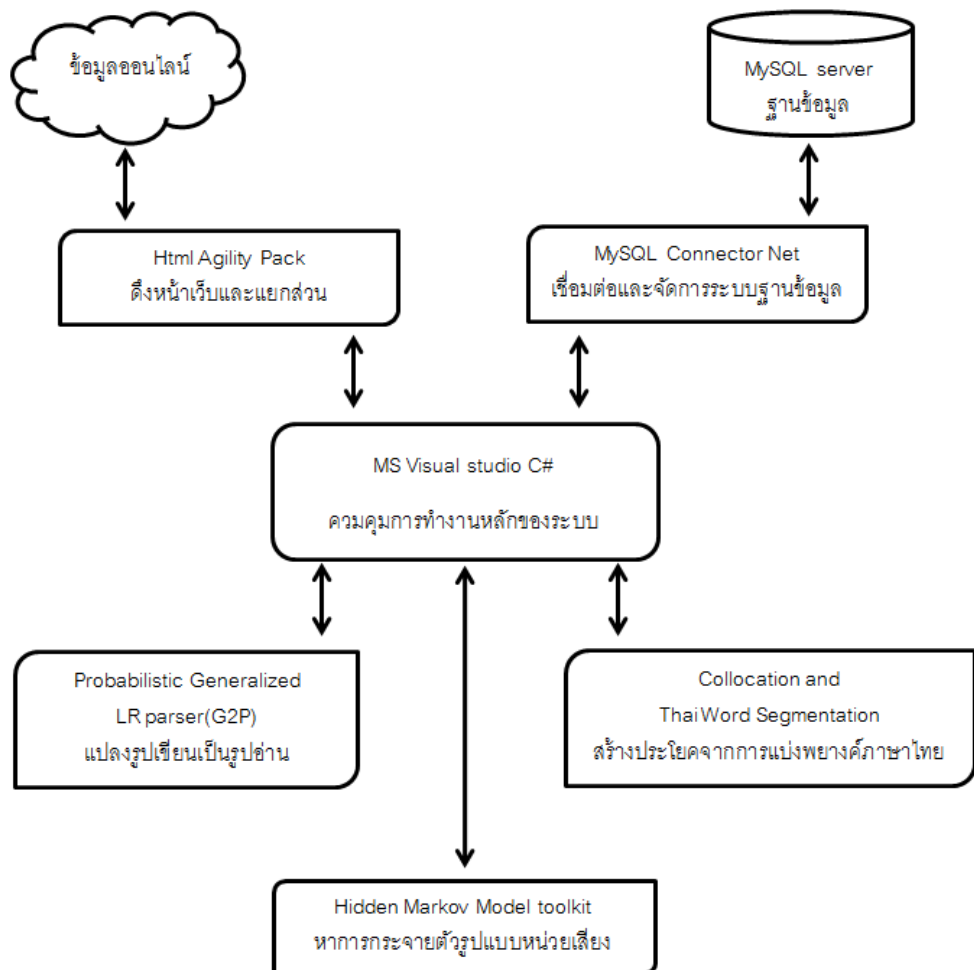
แปลง text เป็น Di-phone ด้วย HTK+CMU dictionary ใช้ Greedy algorithm Weight คะแนนให้หน่วยเสียงที่หายาก (Di-phone) ให้คะแนนแต่ละ sentence ด้วย Okapi formula คัดเหลือ 1,000 ประโยค ผลลัพธ์ ได้การครอบคลุมทางหน่วยเสียงคู่ 93.52% งานวิจัยนี้ทำให้เกิดแนวคิดเรื่องการศึกษาในเรื่องความสั้นยาวของแต่ละประโยค เนื่องจากคะแนนขึ้นอยู่กับจำนวนความยาวของประโยคด้วย การ Weight คะแนนให้ หน่วยเสียงหายากทำให้เข้าถึงเป้าหมายได้เร็ว เกิดแนวคิดเรื่องการดึงข้อความจากฐานข้อมูลออนไลน์

บทที่ 3

ขั้นตอนการดำเนินงานวิจัย

เครื่องมือที่ใช้ในการวิจัย

1. ตัวดึงหน้าเว็บและแยกส่วน HTML, Html Agility Pack (HAP) [30]
2. ตัวแปลงรูปเขียนเป็นรูปอ่าน, Probabilistic Generalized LR parser [12]
3. ตัวนับความถี่รูปแบบหน่วยเสียง, Hidden Markov Model toolkit [31]
4. ตัวแบ่งพยางค์ภาษาไทย, Collocation and Thai Word Segmentation [32]
5. ระบบเชื่อมต่อระบบฐานข้อมูล, MySQL Connector Net 6.2.4 [33]
6. ระบบฐานข้อมูล, MySQL Server 5.1 [34]
7. เครื่องมือพัฒนาโปรแกรมภาษา C#, Microsoft Visual Studio 2010 [35]

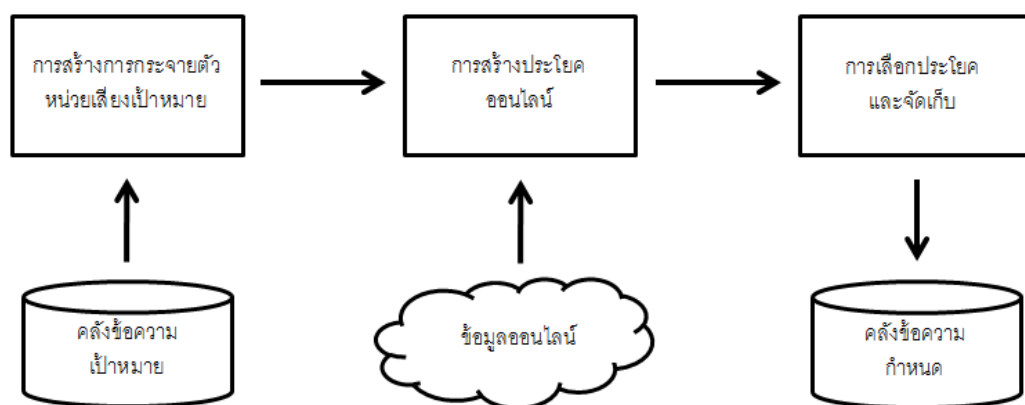


ภาพที่ 3.1 แผนภาพความสัมพันธ์ของเครื่องมือวิจัย

สำหรับเครื่องมือที่ใช้ในงานวิจัย ซึ่งการทำงานมีความสัมพันธ์กัน ดังภาพที่ 3.1 ตัวดึงหน้าเว็บและแยกส่วน มีหน้าที่ดึงหน้าเว็บในรูปแบบ HTML จากข้อมูลออนไลน์บนอินเทอร์เน็ตมาวิเคราะห์แยกส่วนที่เป็นข้อความภาษาไทยและส่วนที่เป็น URL ไปยังหน้าเว็บเพจอื่น ๆ เพื่อส่งให้ระบบหลัก ตัวแปลงรูปเขียนเป็นรูปอ่านมีหน้าที่รับประโยคในรูปเขียนมาแปลงเป็นรูปอ่านและส่งให้ระบบหลัก ตัวนับความถี่รูปแบบหน่วยเสียง มีหน้าที่นำชุดลำดับรูปอ่านทั้งหมดมาแปลงให้อยู่ในรูปแบบหน่วยเสียงที่ต้องการและนับจำนวนที่เกิดขึ้นของแต่ละรูปแบบเพื่อส่งไปให้ระบบหลัก สร้างการกระจายตัวรูปแบบหน่วยเสียง ตัวแบ่งพยางค์ภาษาไทย มีหน้าที่นำชุดบทความมาแปลงเป็นลำดับของประโยคซึ่งในแต่ละประโยคมีจำนวนพยางค์ตามที่ต้องการ และส่งให้ระบบหลัก ระบบเชื่อมต่อฐานข้อมูลมีหน้าที่ ส่งชุดคำสั่ง SQL จากระบบหลักไปจัดการฐานข้อมูล สามารถส่งข้อมูลหรืออ่านข้อมูลจากฐานข้อมูลได้ ส่วนเครื่องมือพัฒนาระบบหลักใช้ชุดคำสั่งภาษา C# ในการควบคุมระบบทั้งหมด

ขั้นตอนการดำเนินงานวิจัย

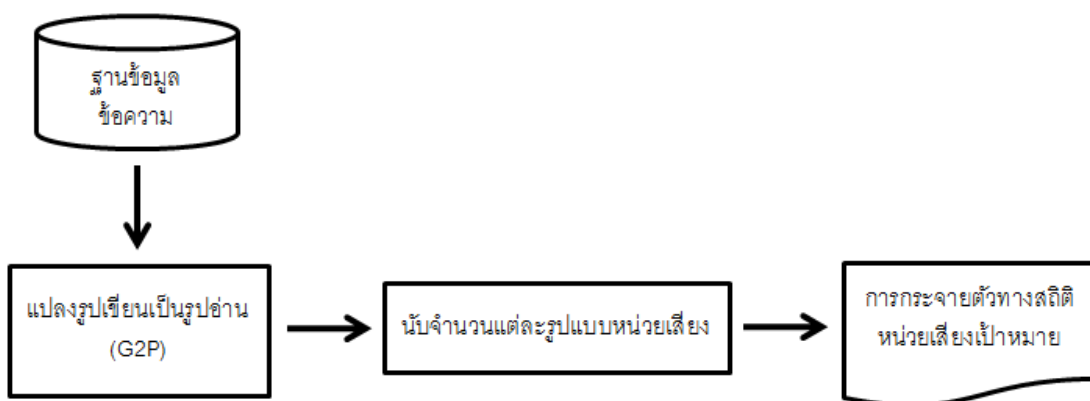
ในบทนี้จะกล่าวถึงขั้นตอนการดำเนินงานวิจัยซึ่งประกอบด้วยสามส่วนหลัก ส่วนแรกคือการหาการกระจายตัวทางหน่วยเสียงเป้าหมาย โดยสามารถหาการกระจายตัวของหน่วยเสียงเป้าหมายได้จากการกำหนดจำนวนความถี่ของแต่ละหน่วยเสียงเอง หรือกำหนดจากรูปแบบการกระจายตัวทางหน่วยเสียงของคลังข้อความอื่น ๆ ส่วนที่สองคือส่วนของการสร้างประโยค ในส่วนนี้จะรวมถึงขั้นตอนการดึงข้อความจากข้อมูลทางอินเทอร์เน็ต ขั้นตอนการแบ่งประโยคจากกลุ่มของข้อความและขั้นตอนการคัดประโยคที่มีส่วนประกอบของอักขระที่ไม่ต้องการออก โดยขั้นตอนนี้จะทำงานไปเรื่อย ๆ จนกระทั่งคลังข้อความที่ต้องการ ส่วนที่สามคือส่วนของการเลือกประโยคและการจัดเก็บสู่คลังข้อความ การทำงานหลักของระบบดังกล่าว แสดงดังภาพที่ 3.2



ภาพที่ 3.2 ภาพรวมการทำงานหลักของระบบ

1. ขั้นตอนการสร้างการกระจายตัวทางหน่วยเสียงเป้าหมาย

การกระจายตัวทางหน่วยเสียงเป้าหมายคือ รูปแบบของการกระจายตัวทางหน่วยเสียงที่ต้องการให้เกิดในคลังข้อความที่จะถูกสร้างขึ้น ในวิทยานิพนธ์นี้สามารถใช้รูปแบบการกระจายตัวได้สองรูปแบบ รูปแบบแรกคือ การกระจายตัวทางหน่วยเสียงที่สร้างจากการกำหนดความถี่จำนวนที่เกิดขึ้นของแต่ละรูปแบบหน่วยเสียงเอง อีกรูปแบบคือ ใช้การกระจายตัวทางหน่วยเสียงที่สร้างจากความถี่จำนวนที่เกิดขึ้นของแต่ละรูปแบบหน่วยเสียงจากคลังข้อความในฐานข้อมูลอื่น ๆ สำหรับวิทยานิพนธ์นี้ได้เลือกใช้ข้อความจากฐานข้อมูล Large Vocabulary Continuous Speech Recognition: LVCSR [13] หรือที่ถูกเรียกง่าย ๆ ว่า LOTUS พัฒนาโดยหน่วยปฏิบัติการวิจัยวิทยาการมนุษยภาษา ศูนย์เทคโนโลยีอิเล็กทรอนิกส์และคอมพิวเตอร์แห่งชาติ (NECTEC) เลือกใช้ชุดประโยคที่ครอบคลุมภาษาไทย 5,000 คำ ซึ่งมาจากบทความทั่วไปและบทความข่าวซึ่งเป็นภาษาไทย จำนวน 3,007 ประโยค ขั้นตอนการสร้างการกระจายตัวทางหน่วยเสียงเป้าหมายสามารถแสดงภาพการทำงาน ดังภาพที่ 3.3



ภาพที่ 3.3 การสร้างการกระจายตัวทางหน่วยเสียงเป้าหมาย

ขั้นตอนแรกนำข้อความจากฐานข้อมูล LOTUS มาแปลงรูปเขียนให้เป็นรูปอ่าน โดยในวิทยานิพนธ์นี้ได้ใช้ ตัวแปลงรูปเขียนเป็นรูปอ่านของ Probabilistic Generalized LR parser [12] ผลลัพธ์จากการแปลงรูปเขียนของข้อความจากฐานข้อมูล LOTUS เป็นรูปอ่านสามารถดูได้จากตารางที่ 3.1 การแปลงรูปเขียนเป็นรูปอ่านจากประโยคทั้งหมด 3,007 ประโยคได้หน่วยเสียงทั้งหมดรวมเป็น 241,199 หน่วยเสียง หลังจากได้รูปอ่านของแต่ละประโยคมาแล้ว ขั้นตอนต่อไปคือการนับจำนวนแต่ละรูปแบบหน่วยเสียงที่เกิดขึ้นทั้งหมดจากรูปอ่านที่แปลงมา โดยในวิทยานิพนธ์นี้ ได้เลือกใช้ รูปแบบหน่วยเสียงแบบหน่วยเสียงคู่ (Diphone) ในการทดลองขั้นตอนการนับจำนวนแต่ละรูปแบบหน่วยเสียงที่เกิดขึ้น จะใช้เครื่องมือ Hidden Markov Model toolkit,

HTK [31] มาช่วยจัดการในขั้นตอนนี้ โดยใช้ฟังก์ชัน HLEd ในการหารูปแบบของหน่วยเสียงคู่ ที่เกิดขึ้นจากรูปอ่านทั้งหมดในที่นี่มีรูปแบบหน่วยเสียงคู่ที่เกิดขึ้น 1,383 รูปแบบ จากนั้นนับการเกิดขึ้นของแต่ละรูปแบบหน่วยเสียงคู่ด้วยฟังก์ชัน HLStats หลังจากขั้นตอนนี้เราจะได้การกระจายตัวทางสถิติทางหน่วยเสียงของเป้าหมาย เพื่อนำไปใช้ในขั้นตอนการคัดเลือกประโยคต่อไป สามารถดูการกระจายตัวทางหน่วยเสียงเป้าหมายในรูปแบบหน่วยเสียงคู่ได้จากตารางที่ 3.2 และภาพที่ 3.4

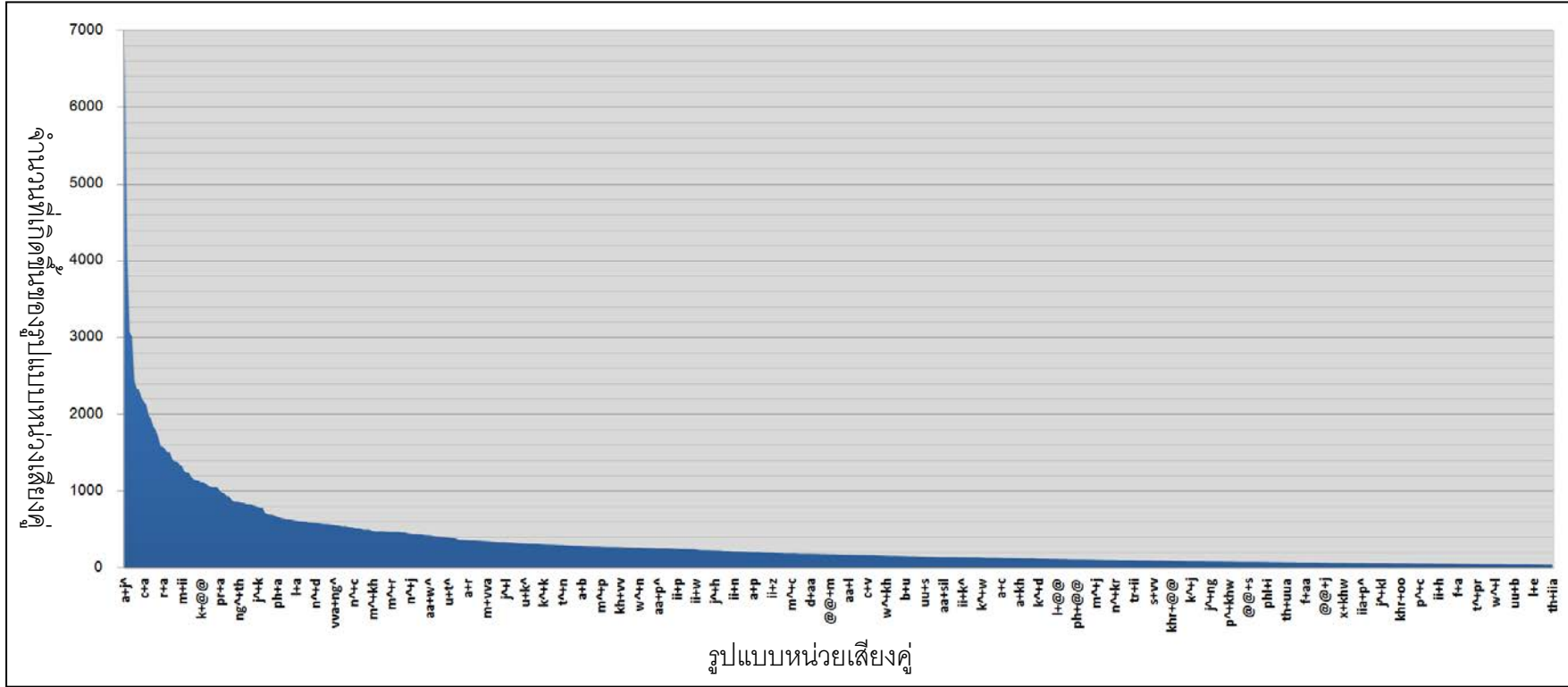
ตารางที่ 3.1 ตัวอย่างการแปลงรูปเขียนเป็นรูปอ่านจากฐานข้อมูล LOTUS

ประโยคในรูปเขียน	ประโยคในรูปอ่าน
อย่างไรก็ตามเหล่านักธุรกิจนักลงทุนที่เข้าร่วมประชุมสัมมนาครั้งนี้มีความเห็นว่าความเติบโตทางเศรษฐกิจในภูมิภาคเอเชียแปซิฟิกจะขึ้นอยู่กับทรัพยากรบุคคลที่ได้รับการฝึกฝนจนมีความชำนาญชำนาญในหน้าที่การงานอย่างสูง	sil j aa ng^ r a j^ k @ @ t aa m^ l a w^ n a k^ t h u r a k i t^ n a k^ l o n g^ t h u n^ t h i i k h a w^ r u u a m^ p r a c h u m^ s a m^ m a n a a k h r a n g^ n i i m i i k h w a a m^ h e n^ w a a k h w a a m^ t q q p^ t o o t h a a n g^ s e e t^ t h a k i t^ n a j^ p h u u m i p h a a k^ z e e c h i i a p x x s i f i k^ c a k h v n^ j u u k a p^ s a p^ p h a j a a k @ @ n^ b u k^ k h o n^ t h i i d a j^ r a p^ k a a n^ f v k^ f o n^ c o n^ m i i k h w a a m^ c h a m^ n i c h a m^ n a a n^ n a j^ n a a t h i i k a a n^ n g a a n^ j a a n g^ s u u n g^ s i l
ก็ต้องแก้ด้วยการศึกษาอบรมด้วยวิธีแปลความหมายใหม่ในความเชื่อที่เป็นพื้นเพเดิมของเขาดังจะกล่าวต่อไปในตอนหลัง	sil k @ @ t @ @ n g^ k x x d u u a j^ k a a n^ s v k^ s a a z o p^ r o m^ d u u a j^ w i t h i i p l x x k h w a a m^ m a a j^ m a j^ n a j^ k h w a a m^ c h v v a t h i i p e n^ p h v v n^ p h e e d q q m^ k h @ @ n g^ k h a w^ d a n g^ c a k l a a w^ t @ @ p a j^ n a j^ t @ @ n^ l a n g^ s i l
ซึ่งเป็นการสร้างภาพตนเองให้เป็นบุคคลสำคัญในลักษณะที่สงสารตัวเอง	sil s v n g^ p e n^ k a a n^ s a a n g^ p h a a p^ t o n^ z e e n g^ h a j^ p e n^ b u k^ k h o n^ s a m^ k h a n^ n a j^ l a k^ s a n a t h i i s o n g^ s a a n^ t u u a z e e n g^ s i l
แต่มันก็บอกไปอะไรไม่ได้มากกว่านั้น	sil t x x m a n^ k @ @ b @ @ k^ b a j^ z a r a j^ m a j^ d a j^ m a a k^ k w a a n a n^ s i l
แต่ก็ตกอยู่ระหว่างเวลาของคัมภีร์ทั้งสองที่กล่าวมาแล้ว	sil t x x k @ @ t o k^ j u u r a w a a n g^ w e e l a a k h @ @ n g^ k h a m^ p h i i t h a n g^ s @ @ n g^ t h i i k l a a w^ m a a l x x w^ s i l

ประโยคในรูปเขียน	ประโยคในรูปอ่าน
หรือชาวชนบทอพยพเข้าสู่ตัวเมืองที่ ทวีจำนวนมากขึ้นในแต่ละประเทศได้ ก่อให้เกิดปัญหายุ่งยากมากทำให้บาง ประเทศเช่นสหรัฐอเมริกาต้องออก กฎหมายห้ามคนเข้าเมืองในปี คริสต์ศักราชหนึ่งพันเก้าร้อยยี่สิบ	sil r v ch aa w^ ch o n^ b o t^ z o p^ p h a j o p^ kh a w^ s u u t u a m v v a n g^ t h i i t h a w i i c a m^ n u u a n^ m a a k^ kh v n^ n a j^ t x x l a p r a t h e e t^ d a j^ k @ @ h a j^ k q q t^ p a n^ h a a j u n g a j a a k^ m a a k^ t h a m^ h a j^ b a a n g^ p r a t h e e t^ c h e e n^ s a r a t^ z a m e e r i k a a t @ @ n g^ z @ @ k^ k o t^ m a a j^ h a a m^ k h o n^ k h a w^ m v v a n g^ n a j^ p i i k h r i t^ s a k^ r a a t^ n v n g^ p h a n^ k a w^ r @ @ j^ j i i s i p^ s i l
กับคำถามแรกที่ให้วิจารณ์บทบาท ของสภารักษาความสงบเรียบร้อย แห่งชาติและรัฐบาลในรอบสองเดือน ที่ผ่านมา	sil k a p^ kh a m^ t h a a m^ r x x k^ t h i i h a j^ w i c a a n^ b o t^ b a a t^ kh @ @ n g^ s a p h a a r a k^ s a a k h w a a m^ s a n g o p^ r i i a p^ r @ @ j^ h x x n g^ c h a a t^ l x r a t^ t h a b a a n^ n a j^ r @ @ p^ s @ @ n g^ d v v a n^ t h i i p h a a n^ m a a s i l
แต่ในการประชุมเรื่องโลกร้อนขึ้นที่ ประเทศเนเธอร์แลนด์เมื่อเดือน พฤศจิกายนปีที่แล้ว	sil t x x n a j^ k a a n^ p r a c h u m^ r v v a n g^ l o o k^ r @ @ n^ kh v n^ t h i i p r a t h e e t^ n e e t h q q l x x n^ m v v a d v v a n^ p h r v t^ s a c i k a a j o n^ p i i t h i i l x x w^ s i l
และมีชีวิตอยู่ในโลกด้วยการร้องเพลง ให้ความรื่นรมย์เท่านั้น	sil l x m i i c h i i w i t^ j u u n a j^ l o o k^ d u u a j^ k a a n^ r @ @ n g^ p h l e e n g^ h a j^ k h w a a m^ r v v n^ r o m^ t h a w^ n a n^ s i l
ไม่ใช่ความจริงที่จับต้องได้เดี๋ยวนี	sil m a j^ c h a j^ k h w a a m^ c i n g^ t h i i c a p^ t @ @ n g^ d a j^ d i i a w^ n i i s i l
คงเป็นที่ทราบกันดีว่าการจัด การศึกษาในปัจจุบันคุณภาพ การศึกษายังได้ผลเป็นที่น่าพอใจ	sil k h o n g^ p e n^ t h i i s a a p^ k a n^ d i i w a a k a a n^ c a t^ k a a n^ s v k^ s a a n a j^ p a t^ c u b a n^ k h u n^ n a p h a a p^ k a a n^ s v k^ s a a j a n g^ d a j^ p h o n^ p e n^ t h i i n a a p h @ @ c a j^ s i l
ซึ่งได้แก่บุคคลที่ไม่อาจจะเรียนรู้ ระเบียบของสังคมได้ถูกต้องหรือผู้ที่ ต่อต้านคุณค่าทางสังคมและหรือ ความเชื่อทางสังคม	sil s v n g^ d a j^ k x x b u k^ k h o n^ t h i i m a j^ z a a t^ c a r i i a n^ r u u r a b i i a p^ kh @ @ n g^ s a n g^ k h o m^ d a j^ t h u u k^ t @ @ n g^ r v v p h u u t h i i t @ @ t a a n^ k h u n^ k h a a t h a a n g^ s a n g^ k h o m^ l x r v v k h w a a m^ c h v v a t h a a n g^ s a n g^ k h o m^ s i l

ตารางที่ 3.2 ตัวอย่างจำนวนของแต่ละรูปแบบหน่วยเสียงคู่ที่เกิดขึ้นจากฐานข้อมูล LOTUS

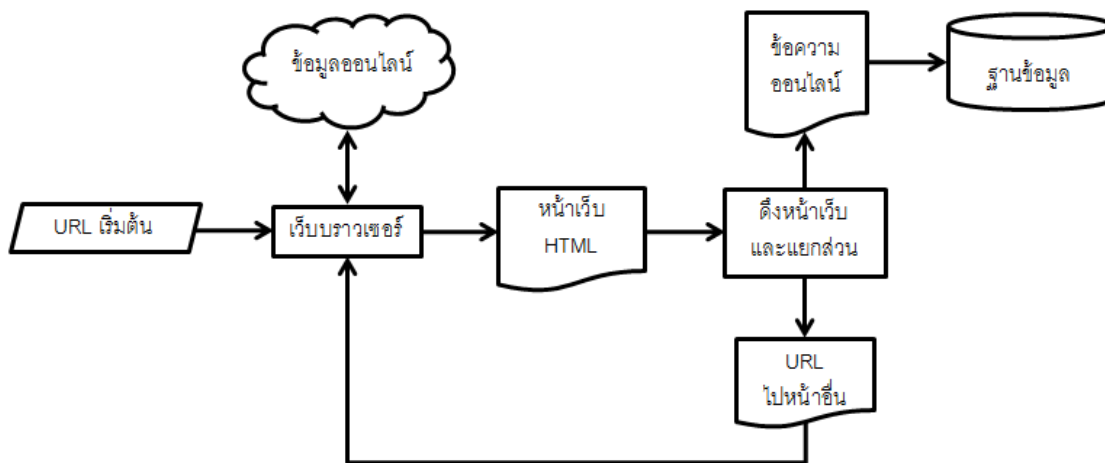
รูปแบบ	จำนวน	รูปแบบ	จำนวน	รูปแบบ	จำนวน	รูปแบบ	จำนวน
a+j^	6792	a+n	1061	l+a	617	j^+d	476
aa+n^	4136	v+ng^	1056	j^+sil	613	ng^+r	476
n+a	3076	s+aa	1055	k^+s	609	j^+kh	473
k+aa	3023	p+a	1050	ng^+m	607	ii+s	472
th+ii	2454	pr+a	1007	a+m	606	m^+r	471
a+n^	2335	j+aa	982	j^+r	596	m^+s	470
a+m^	2326	th+aa	969	n^+ph	593	t^+s	470
@@+ng^	2226	aa+k^	937	j^+th	592	z+aa	469
c+a	2165	n^+n	933	n^+d	590	ii+m	466
th+a	2131	kh+o	886	n^+r	586	ng^+ph	463
aa+ng^	1998	n^+k	870	z+a	583	a+s	462
s+a	1947	ng^+k	865	s+v	576	j+a	452
aa+m^	1844	ng^+th	863	s+i	575	n^+j	449
e+n^	1803	n+aa	851	c+aa	574	ii+k	443
o+n^	1714	n^+kh	848	o+ng^	569	aa+k	442
d+a	1595	h+a	836	j^+s	568	w+i	442
r+a	1583	n^+s	829	vva+ng^	563	a+ph	440
p+e	1548	l+x	828	n^+z	556	t+xx	440
k+a	1506	ng^+kh	817	o+m^	556	k^+kh	431
kh+@@	1504	a+k^	814	t+a	548	n^+h	429
a+ng^	1416	j^+k	791	a+k	546	aa+w^	428
m+a	1392	a+th	785	h+aa	546	j+uu	424
khw+aa	1389	t+@@	784	ng^+t	538	kh+aa	417
n^+th	1347	n+ii	715	j^+n	533	aa+th	412
m+ii	1334	ch+a	702	n^+c	529	a+d	409
kh+a	1264	t+aa	698	ph+uu	517	b+aa	405
a+t^	1242	ng^+s	695	i+ng^	516	t^+c	405
m+aa	1240	n^+m	688	xx+ng^	514	uua+n^	401
a+w^	1189	ph+a	666	n^+l	504	u+t^	397
a+p^	1152	ph+aa	661	n^+t	500	s+o	395
t^+th	1147	ee+t^	647	a+j	499	aa+s	390
i+t^	1140	n^+sil	645	n^+p	499	t^+kh	388
k+@@	1123	w+a	636	m^+kh	485	ii+c	373
w+aa	1122	ng^+n	632	ch+aa	478	uua+j^	371
aa+t^	1098	j^+m	631	ph+o	478	ng^+c	366



ภาพที่ 3.4 แผนภาพการกระจายตัวทางสถิติทางหน่วยเสียงของฐานข้อมูล LOTUS

2. ขั้นตอนการสร้างประโยคออนไลน์

ขั้นตอนนี้มีหน้าที่สร้างประโยคเพื่อนำไปใช้ในขั้นตอนการเลือกประโยคเพื่อจัดเก็บสู่คลังข้อความกำหนด แผนภาพการทำงานโดยรวมของขั้นตอนดูได้จากภาพที่ 3.5 ในวิทยานิพนธ์นี้ได้เลือกใช้บทความจากอินเทอร์เน็ตมาสร้างประโยค เครื่องมือการดึงข้อความจากหน้าเว็บจึงถูกนำมาใช้ ในที่นี้เราใช้ตัวดึงหน้าเว็บแยกวิเคราะห์ HTML ของ Html Agility Pack (HAP) [30] เนื่องจากเป็นฟังก์ชันที่ใช้ง่ายและสามารถดึงข้อมูลจากหน้าเว็บได้เร็ว หลักการทำงานคือ URL เริ่มต้น จะถูกส่งไปให้เว็บเบราว์เซอร์สิ่งที่ได้คือ โครงสร้าง HTML ของหน้าเว็บเพจนั้นซึ่งอาจประกอบไปด้วยข้อมูลข้อความ ข้อมูลภาพ ข้อมูล URL หรือข้อมูลอื่น ๆ จากนั้นตัวดึงหน้าเว็บแยกวิเคราะห์ HTML จะทำการคัดกรองส่วนที่เราต้องการออกมา ในที่นี้ได้ดึงส่วนที่เป็น URL เก็บไว้เพื่อใช้ในการเข้าสู่หน้า HTML อื่น ๆ และดึงส่วนที่เป็นข้อความไปจัดเก็บไว้ที่ฐานข้อมูลเพื่อเตรียมไว้ในส่วนของการจัดการประโยคต่อไป ในงานวิทยานิพนธ์นี้ได้เลือกใช้บทความจากเว็บไซต์ สารานุกรมไทยสำหรับเยาวชน [36] ตัวอย่างหน้าเว็บแสดงดังภาพที่ 3.6 และจากเว็บคลังข้อมูลข้อความภาษาไทยแห่งชาติ [37] เนื่องจากมีโครงสร้าง HTML ที่ไม่ซับซ้อนและบทความมีหลากหลายลักษณะ เช่น บทความวิชาการ หนังสือแบบเรียน หนังสือพิมพ์ นิตยสารและอื่น ๆ สามารถดูตัวอย่างของข้อความที่เก็บไว้ในฐานข้อมูลได้ดังภาพที่ 3.7



ภาพที่ 3.5 การดึงข้อความออนไลน์ด้วยตัวดึงหน้าเว็บแยกวิเคราะห์ HTML

สารานุกรมไทยสำหรับเยาวชน
โดยพระราชประสงค์ในพระบาทสมเด็จพระเจ้าอยู่หัว

หน้าหลัก พระราชดำริ โครงการ เรื่องในสารานุกรม ลับสมอง ค้นข้อมูล ติดต่อ ระบบงาน

สารานุกรมไทย สำหรับเยาวชน เล่มที่ ๑

- เรื่องที่ ๑ ดวงอาทิตย์
- เรื่องที่ ๒ อุปราคา
- เรื่องที่ ๓ ห้องฟิวชัน
- เรื่องที่ ๔ นก
- เรื่องที่ ๖ เครื่องจักรกล
- เรื่องที่ ๗ พลังงาน
- เรื่องที่ ๘ อากาศยาน
- เรื่องที่ ๙ ดนตรีไทย

กลับหน้าแรก

Facebook Share สารานุกรมไทยสำหรับเยาวชน / เล่มที่ ๑ / เรื่องที่ ๑ ดวงอาทิตย์ / สภาพภายในดวงอาทิตย์

สภาพภายในดวงอาทิตย์

ดังที่ได้กล่าวมาแล้วว่า ดวงอาทิตย์เป็นก้อนสสารใหญ่ร้อนจัด และรวมตัวเป็นสสารพื้นฐาน ทรงกลมอยู่ได้ โดยแรงดึงดูดระหว่างอะตอมและโมเลกุล แรงดึงดูดหรือแรงโน้มถ่วง (gravitational force) นี้ มีทิศทางเข้าหาจุดศูนย์กลาง เนื้อสารของดวงอาทิตย์ ซึ่งอยู่ที่ระดับใดระดับหนึ่ง ภายในดวงอาทิตย์จะถูกบีบอัดโดยเนื้อสารที่อยู่ชั้นนอกมา จึงเป็นธรรมชาติที่จะต้องมีความดันและความหนาแน่นมากกว่าเนื้อสารในชั้นที่สูงกว่าตามลักษณะที่อื่นนี้ กล่าวได้ว่า ความดันและความหนาแน่นของเนื้อสารเพิ่มขึ้นในระดับลึกลงไปภายในดวงอาทิตย์ อนึ่งภายใต้ความกดดันสูงนั้นก๊าซหรือไอจะถูกบีบให้ปริมาตรลดลงเรียก ก๊าซอะตอมหรือโมเลกุลของก๊าซหรือไอนั้น ไม่มีความเร็วในตัว พอดีจะเคลื่อนที่ต่อสู่ไว้ ความเร็วที่กล่าวถึงนี้ได้จากการมีอุณหภูมิสูง ทั้งนี้เพราะอุณหภูมิของวัตถุก็คือพลังงานของการเคลื่อนที่ และการสั่นสะเทือนของโมเลกุลอะตอมในสสารนั้นๆ โดยเหตุนี้เอง เราจึงได้ว่า เนื้อสารที่ระดับใดระดับหนึ่ง ภายในดวงอาทิตย์อยู่ในสภาพสมดุล เมื่ออุณหภูมิความกดดัน และความหนาแน่นพอเหมาะแก่กัน ซึ่งจะมีค่าสูงขึ้นเรื่อยๆ สำหรับระดับที่ลึกกลงไปภายในดวงอาทิตย์

ที่แก่นกลางของดวงอาทิตย์ในอาณาบริเวณรูปทรงกลมมีรัศมีประมาณ ๒ แสนกิโลเมตร ซึ่งมีอุณหภูมิสูงเพียงพอนั้น มีปฏิกิริยาเทอร์โมนิวเคลียร์เกิดขึ้น และให้พลังงานในลักษณะของรังสีแกมมา ซึ่งมีขนาดคลื่นสั้น รังสีนี้แผ่กระจายโดยการถ่ายเทผ่านเนื้อสารของดวงอาทิตย์ออกมา จนถึงระดับลึกประมาณ ๑๐๐,๐๐๐ กิโลเมตร จากที่นั่นความ การถ่ายเทพลังงานก็จะแปรวิธจากจากการแผ่รังสี (radiation) มาเป็น การนำความร้อน (convection) ไปด้วยก๊าซที่ร้อนจะลอยตัวขึ้นมาสู่ระดับสูง จนถึงระดับผิวดวงอาทิตย์ก็จะแผ่รังสีแสงสว่าง และความร้อน ออกสู่อวกาศ ครั้นแล้วเมื่ออุณหภูมิของมันลดลงก็จะจับตัวเป็นเม็ดกรมนิวไคลนิกคล้ายคลึงกับการเดือดของของเหลว เช่น น้ำหรือน้ำมันที่ใส่ภาชนะต้มบนเตาไฟให้ร้อนนั่นเอง

หัวข้อก่อนหน้า หัวข้อถัดไป

โครงการสารานุกรมไทยสำหรับเยาวชน โดยพระราชประสงค์ในพระบาทสมเด็จพระเจ้าอยู่หัว
โครงการสารานุกรมไทยฯ สนามเสือป่า ถนนศรีอยุธยา เขตดุสิต กรุงเทพฯ 10300

ภาพที่ 3.6 ตัวอย่างหน้าเว็บจากอินเทอร์เน็ตที่มีเนื้อความที่เท่ากับการดึงข้อความ

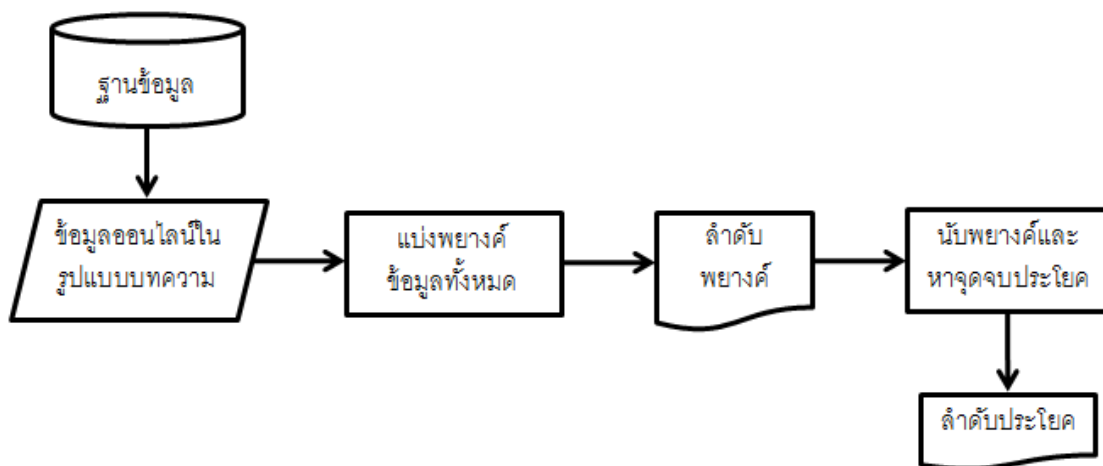
Data_ID	Data_Text	Data_Source
1	กิลเมรชิวา แหมลมพรหมเทพ จุลชนวิวัฒน์อินฉนวนัน มงะทีเฉลิมและประเสริฐ เหมระสำหรับเขียนรู้มหาสมุทร พังสามความหลากหลาย	http://ling.arts.chula.ac.th/tnc
2	คล้ายไม่จริงเท่านั้นโตแห่งนิคมลียงใหญ่ แทนทะเลและเปียก ที่กำลังละลาย โลกก็ขยับจากผืนทราย ตามด้วยดอกไม้อัลไพน์อันงดงาม	http://ling.arts.chula.ac.th/tnc
3	ดิเอสอีได้พูดถึงการใช้สาร AF2 ไซในน้ำดื่มเพื่อป้องกันมลพิษเป็นอันตรายต่อมนุษย์ นอกจากนี้ตัวไอพ่นของยาฆ่าเชื้อที่เข้าสู่อากาศและน้ำ	http://ling.arts.chula.ac.th/tnc
4	เขียนให้ดีขึ้นเป็นเสียงแรกบนเวที และดีขึ้นอีกเป็นระยะๆ และดังตอบโต้เรื่องที่น่าสนใจ เช่น ตอนที่ถูกขายกับกิมมาล่าเรื่อง ดอกบัวสีขาว	http://ling.arts.chula.ac.th/tnc
5	แม่เลี้ยงพ่ออย่างน่าประหลาด แม่พูดให้ใครต่อใครฟังหลายปีแล้วว่าทุกอย่างในหัวฉันเกิดขึ้ไม่อยู่ในสตรีน้อยอย่างเธอเลย โดยเฉพาะฉัน	http://ling.arts.chula.ac.th/tnc
6	คนไข้เด็ก เมื่อพ่อแม่พาเด็กมาปรึกษาการตรวจที่โรงพยาบาล เด็กโดยมากมักมีตั้งต้นตกใจมากกว่าผู้ใหญ่ และกลัวว่าพ่อแม่จะทิ้งเขาไป จึงไม่	http://ling.arts.chula.ac.th/tnc
7	อาชีพทุกอาชีพเป้าหมายคือการหาเงินเลี้ยงชีพ ต่างคนต่างมีวิถีชีวิตที่เลือกเองไม่ได้ ถ้าพ่อแม่ไม่รวย อาชีพของหญิงบริการในเมืองไทยเป็น	http://ling.arts.chula.ac.th/tnc
8	เจ้าหญิงผู้พิศมัย น. เข็มใหม่ เจ้าพระรัตนคำ น. เข็มใหม่ ศิษย์งาม ก้านดอกขาวลุ่มมากิน ก็คงว่าไม่เอาไหน ด้วยอึ้งแฉ่งแฉ่ง แต่คนภาคกลาง	http://ling.arts.chula.ac.th/tnc
9	ย้อนกลับมารับประทานแล้ว หลังจากกราบขอพรใหม่จากคุณยายเสร็จ แม่กลับเข้าไปนอนเคลตต่อน้ำ หกโมงเย็นตรงปะ คุณยายมาประจวบ	http://ling.arts.chula.ac.th/tnc
10	การดื่มเหล้ามีจุดหมายเพิ่มโปรตีนและแคลเซียมเข้าในหัว แต่มีสัดส่วนของแอลกอฮอล์สูงเกินไป แทนที่จะได้ผลในการป้องกันกระดูกพรุน	http://ling.arts.chula.ac.th/tnc
11	ก่อนเดินทางจากไป เจ้าหญิงกมลจุกหามาว่าโตพวกเขาจะมีชีวิตที่เรียบง่ายหรือชีวิตที่หรูหรา ก็คงไม่ได้คำตอบและเรียกเจ้าหญิงว่า	http://ling.arts.chula.ac.th/tnc
12	งานเขียนของวินทร์ เลียววาริณ ที่มีลักษณะเป็นวรรณกรรม อย่างเช่น เรื่องชยะ ซึ่งเป็นภาพเรขาคณิต สิ่งที่ถูกดูเข้าไปก็คือคำพูดของนักก	http://ling.arts.chula.ac.th/tnc
13	จนพ่ายแพ้คุณศัพท์ เป็นศัพท์ที่วิจิตร ภาษพงศ์ 2526 นักภาษาศาสตร์ในสำนักโบราณคดีโครงสร้าง structural school เป็นผู้เริ่มใช้โดย	http://ling.arts.chula.ac.th/tnc
14	ความรู้จักชั้นหนึ่งนั้น อยู่ที่ความสามารถของผู้แต่งที่ได้พยายามสอดเยื้องหลับประวัติศาสตร์เข้าไปในเนื้อเรื่อง ได้พอเหมาะ ได้สาระ	http://ling.arts.chula.ac.th/tnc
15	เหมือน กับผู้หญิงคนหนึ่ง ๆ จากหมู่บ้านภาคอีสานที่เลือกแต่งงานกับฝรั่ง ดันหนีจากครอบครัวของตน ท่าจะรักไม่สำเร็จ ไปหาจับข้างลุง	http://ling.arts.chula.ac.th/tnc
16	การทหารหรือฐานของรัฐบาล เช่นเดียวกับเรื่องผู้สังหาร ทฤษฎีวิพากษ์ ได้แปลงสภาพของรัฐบาลที่มีลักษณะนามธรรม หรือถูกแบ่งด้วย	http://ling.arts.chula.ac.th/tnc
17	งานเขียนของวินทร์ เลียววาริณ ที่มีลักษณะเป็นวรรณกรรม อย่างเช่น เรื่องชยะ ซึ่งเป็นภาพเรขาคณิต สิ่งที่ถูกดูเข้าไปก็คือคำพูดของนักก	http://ling.arts.chula.ac.th/tnc
18	ความแหวกแนวอาจเกิดขึ้นในบันเทิงคดีได้จากการใช้ผู้เล่าเรื่องที่ไม่ธรรมดา เล่าเรื่องกันพบเห็น อาจใช้ผู้เล่าเรื่องที่เป็นเด็ก สัตว์ หรือ	http://ling.arts.chula.ac.th/tnc
19	เขารับใหม่จะไม่ขอลืมชื่อของเรา เราไม่อยากจะเดินทางต่อไป เมื่อเปิดเปลือกความขึ้นมา จับรู้ว่ามีคนรอบข้างกำลังเคลื่อนไหว ชีวิตดำเนิน	http://ling.arts.chula.ac.th/tnc
20	การจิตใจของดิสนีย์และประพันธ์ จิตสำนึกกำลังเปิดฉากแห่งความทรงจำภายใน โลกเล่นเข้าสู่ความฝัน ดันห้องเที่ยวไปในโลกแห่งความฝัน	http://ling.arts.chula.ac.th/tnc

SELECT * FROM 't_data'

Record 1 of 853

ภาพที่ 3.7 ตัวอย่างข้อความออนไลน์จากอินเทอร์เน็ตที่เก็บไว้ในฐานข้อมูล

ต่อมาระบบจะนำข้อความออนไลน์ที่ถูกจัดเก็บไว้ในฐานข้อมูล มาทำการแปลงข้อความในรูปแบบบทความให้เป็นรูปแบบของประโยค โดยแต่ละประโยคต้องมีความยาวที่เหมาะสมที่จะนำไปใช้ในขั้นตอนการเลือกประโยค โดยความยาวประโยคในที่นี้สามารถกำหนดได้ โดยระบุขอบเขตค่าต่ำสุดและสูงสุดของจำนวนพยางค์ที่อยู่ในประโยค การทำงานของขั้นตอนนี้ได้จากภาพที่ 3.8 เริ่มจากระบบจะดึงบทความออนไลน์ในฐานข้อมูลมาทีละหนึ่งบทความ จากนั้นจะทำการแบ่งพยางค์ด้วยเครื่องมือ ตัวแบ่งพยางค์จากคำไทย Collocation and Thai Word Segmentation [32] เมื่อแบ่งบทความในรูปแบบของพยางค์แล้วจะได้ลำดับของพยางค์ ตัวอย่างผลลัพธ์ของการแบ่งพยางค์ได้จากตารางที่ 3.3 จากนั้นระบบจะทำการนำลำดับพยางค์มาจัดเรียงใหม่ตามลำดับ โดยที่จะเพิ่มพยางค์เพื่อสร้างประโยคใหม่ที่ละลำดับทำการนับพยางค์โดยที่ ถ้าผลรวมของลำดับพยางค์ มากกว่าค่าต่ำสุดที่กำหนดไว้ของจำนวนพยางค์ให้เพิ่มพยางค์ไปเรื่อยๆ จนพบกับลำดับที่เป็นเครื่องหมายวรรคตอน จึงตัดเป็นหนึ่งประโยค ถ้าผลรวมของจำนวนพยางค์ในประโยคน้อยกว่าหรือมากกว่าขอบเขตค่าต่ำสุดและสูงสุดของพยางค์ที่กำหนดไว้ ให้ตัดประโยคนั้นทิ้งไป ถือว่าใช้ไม่ได้ แต่ถ้าผลรวมของจำนวนพยางค์ในประโยคอยู่ในขอบเขตค่าต่ำสุดและสูงสุดของพยางค์ที่กำหนดไว้ ประโยคนั้นจะถือว่าใช้ได้และเข้าสู่ขั้นตอนต่อไป โดยในการทดลองของงานวิจัยนี้ได้กำหนดขอบเขตของจำนวนพยางค์ในหนึ่งประโยค คือ ต้องมากกว่า 20 พยางค์แต่ไม่เกิน 40 พยางค์ในหนึ่งประโยค เพื่อความเหมาะสมและง่ายต่อการนำไปบันทึกเสียง



ภาพที่ 3.8 การแปลงข้อความในรูปแบบบทความเป็นรูปแบบประโยค

ตารางที่ 3.3 ตัวอย่างผลลัพธ์จากการแปลงข้อความรูปแบบบทความเป็นรูปแบบลำดับพยางค์

ข้อความในรูปแบบลำดับพยางค์
<p>จะว่าไปแล้ว ก็สุภานั่นแหละที่ทำให้ชีวิตจันทร์เกิดอยากจะได้สิ่งดวงแบบนี้ขึ้นมาบ้าง เพราะเมื่อวันก่อนตอนเธอก้าวเข้าประตูออฟฟิศมาก็เห็นสาวๆ ทั้งหลายกำลังมุงอยู่รอบโต๊ะของสุภา ชีวิตจันทร์ผู้มาทีหลังเลยอยากรู้ขึ้นมาบ้างว่าอะไรเป็นอะไร พอแทรกตัวเข้าไปได้ก็เห็นหน้าหมวยๆ จืดๆ ของสุภากำลังยิ้มซะจนตาหยี หล่อมนุ่มเองที่เป็นสาเหตุของความฮือฮาในตอนเช้านี้ สุภาถูกล็อตเตอรี่ ชันซื้อเพราะเป็นเรื่องบังเอิญไง ที่อยู่ๆ ก็มีรถของจันทร์มารับให้ไม่ต้องเปียกฝน มันก็ต้องเป็นรถนำโชคของฉันทันแหละ จริงไหมเล่า ชีวิตจันทร์ไม่ค่อยขับรถมาทำงานหรอก เพราะรถนั้นก็ไม่ใช่ของเธอ แต่เป็นรถของพ่อ เมื่อก่อนพ่อต้องขับรถไปทำงานทุกวัน พ่อเกษียณแล้วก็เลยจอดเอาไว้ขับไปซื้อกับข้าวบ้าง ไปสังสรรค์กับเพื่อนเก่าๆ บ้าง ชีวิตจันทร์ก็ได้อาศัยขับมาทำงานในวันที่ฝนตกหรือต้องออกไปพบลูกค้าข้างนอก มันไม่ใช่รถใหม่ป้ายแดงอย่างที่คนเล่นห่วยชอบซื้อเลขตามทะเบียนเสียหน่อย แต่แล้วว่ามันก็ทำให้สุภาถูกรางวัลเข้าจนได้ แถมไม่ใช่เงินน้อยๆ เพราะเลขทะเบียนรถของพ่อชีวิตจันทร์นะ ผิดไปจากรางวัลที่ 1 แค่ตัวเดียว สุภาจึงได้เงินร่วมครึ่งแสนมาแบบสบายๆ ด้วยโชคที่ควรจะเป็นของชีวิตจันทร์ หลังจากเข้านั้น ชีวิตจันทร์ผู้แค่นี้ก็เลยตั้งเป้าหมายไว้ว่าจะซื้อล็อตเตอรี่กับเค้าบ้าง ทั้งที่เธอไม่เคยสนใจเรื่องการเสี่ยงดวงแบบนี้เลยแม้แต่หน่อย แต่ไอ้การที่เธอขับรถที่มีเงินครึ่งแสนติดท้ายฝาฝนมาให้คนอื่นเชิดไปง่ายๆ แบบนี้ก็ไม่ใช่เรื่องชีวิตจันทร์จะยอมได้เหมือนกัน ชีวิตจันทร์เลยคิดว่า ถ้าทะเบียนรถพ่อออกไปแล้ว คราวนี้ก็มาถึงทะเบียนรถของแม่บ้างละ เธอเก็บเงินไม่ยอมบอกเรื่องนี้กับใคร แล้วค่อยๆ หาซื้อล็อตเตอรี่ที่มีเลขตรงกับทะเบียนรถของแม่เท่านั้นมาเก็บไว้ เลข 6482 นั้นหาไม่ยากนัก แต่ทุกครั้งที่เธอยืนยันจะเอาแต่เลขจำนวนนี้ ก็คงทำให้คนขายล็อตเตอรี่มองว่าเธอคงไปได้เลขเด็ดจากไหนเป็นแน่ แต่ถึงจะหายากยังไง เวลา 15 วันก่อนจะประกาศผลทำให้ชีวิตจันทร์มีเวลาซื้อเก็บไว้ได้พอสมควร</p>

เมื่อได้สร้างประโยคแล้วประโยคเหล่านั้นจะถูกนำไปแปลงรูปเขียนเป็นรูปอ่าน สามารถดูแผนภาพขั้นตอนนี้ได้จาก ภาพที่ 3.9 โดยจะนำไปผ่านขั้นตอนการกรอกรูปแบบตัวอักษรที่ไม่ต้องการออกก่อน ซึ่งในการทดลองของวิทยานิพนธ์นี้ได้กำหนดให้ ประโยคที่ใช้ประกอบด้วยตัว

อักษรไทยเท่านั้น โดยในที่นี้จะใช้หลักการระบุช่วงอักขระของนิพจน์ปกติ (Regular Expression) ในที่นี้ระบุไว้ที่ช่วง [ก-ฮ], [ะ-ุ], [เ-็] และ [็-็] เพื่อต้องการตัดอักขระ ภาษาต่างประเทศ ตัวเลข และสัญลักษณ์พิเศษออก ในการคัดกรองนี้ยังนำประโยคที่ไม่สามารถแปลงรูปเขียนเป็นรูปอ่านสำเร็จออกไปด้วย เช่น ประโยคที่มีเครื่องหมาย [?] จะทำให้การแปลงรูปเขียนเป็นรูปอ่านไม่สำเร็จ สามารถดูผลการคัดกรองประโยคได้จากตารางที่ 3.4 และตารางที่ 3.5

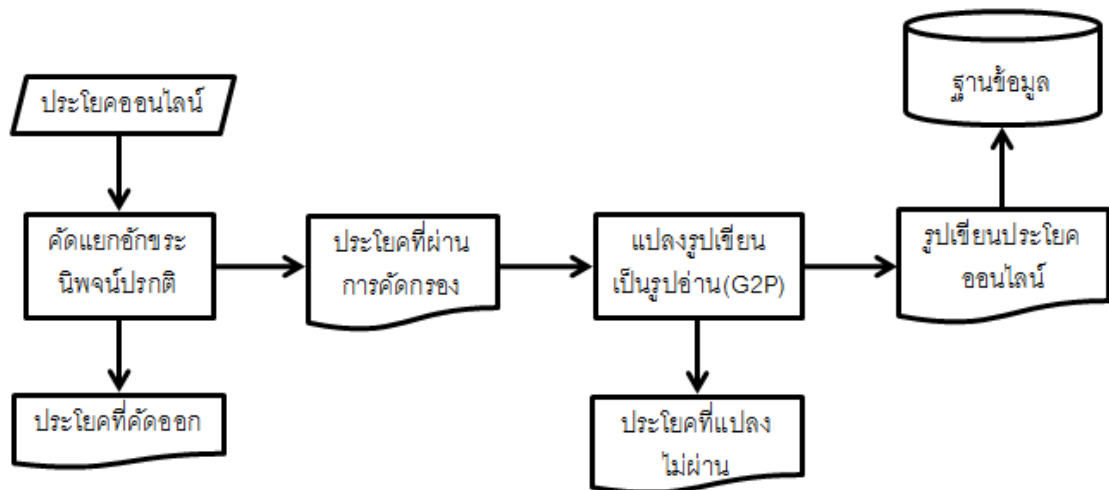
ตารางที่ 3.4 ตัวอย่างผลการคัดแยกประโยคตามอักขระนิพจน์ปกติที่ต้องการ

ไปขึ้นรถญี่ปุ่นของเขาที่จอดอยู่หน้าโรงสี ระหว่างทางชายหนุ่มชวน
มีของกินสารพัดอย่างกับผลไม้ตลอดปี ถ้ามีเงินซื้อกิน เขาภูมิใจที่เกิดมาเป็นคนไทย
เดินมา ฉันทยับไปไหนไม่ได้ขาซ้ายเป็นตะคริวเพราะย่อตัวนานเกินไป
ประกอบ น้ำเขียวออกอย่างนั้น ไม่ต้องติดป้ายว่าห้ามกินห้ามอาบก็ไม่มีใครกล้าเข้าใกล้อยู่แล้ว
ยกเว้นครั้งที่ผมเห็นมันทำท่าเหมือนจะเปลี่ยนสีขน หมอดูแลรักษาผมด้วยยาสายลมอยู่หนึ่งเดือนเต็ม
สหายเพียงบอกว่าประเดี๋ยวจะพาไปตกปลา เราไม่ขัดข้อง หลังกินข้าวเข้าเสร็จเราก็เดินออกมาพร้อมกับคันเบ็ด
พ่อเฒ่าอาจเอาเยื่อปอกถัสด้านแปดพันอยู่ตามร่มไม้ชายคา แม่เฒ่าอาจต้มข้าวหมู
เมื่อเห็นเพื่อนสาว ไม่พอใจ ฮิลลารีจึงต้องเงียบ จัดการกับ แฮมเบอร์เกอร์ตรงหน้าจนหมดอัน
ดวงตาบวมซ้ำเพราะเพิ่งผ่านการร้องไห้อย่างหนัก เธอคุยกับคิงด้วยสีหน้าซึ่งยากแก่การคาดเดาว่าคิดอะไรอยู่
สำหรับเจ้าชายตะโดโมเียงไปยืนอยู่ชายคากัญญา พวกชั้นที่กับนางข้าหลวงชาวมอญไปนั่งอยู่ที่เรือแซ
จึงตัดสินใจเดินเท้าจากบ้านเกิด อำลาญาติมิตรเดินเท้าไปชั่วยาว
ภคตราตะโกนบอกหัวหน้าโรงเพาะที่เป็นชายวัยเกือบห้าสิบท่าทางใจดี
เมื่อถึงพระตำหนักอันเคยเป็นที่ประทับของพระมารดา ก็เสด็จพระราชดำเนินลิ่วไปยังห้องบรรทมทันที
เขามขก็ไม่ชอบ บางคนก็อยากจะให้คนชมทำเดี๋ยว ผู้กองบ่น
ทั้งเมืองโกลาหล ลูกไฟขนาดใหญ่ราวจุดโคมลอยกระจายไปทุกทิศ ตกไหนติดนั้นคล้ายมีคนจุดไฟรับ
แต่หล่อนกลับเกลียดฉันมากขึ้นพอรู้ว่าฉันเป็นลูกของคนรักของคุณวี เธอไม่ได้เกลียดคุณหรือ

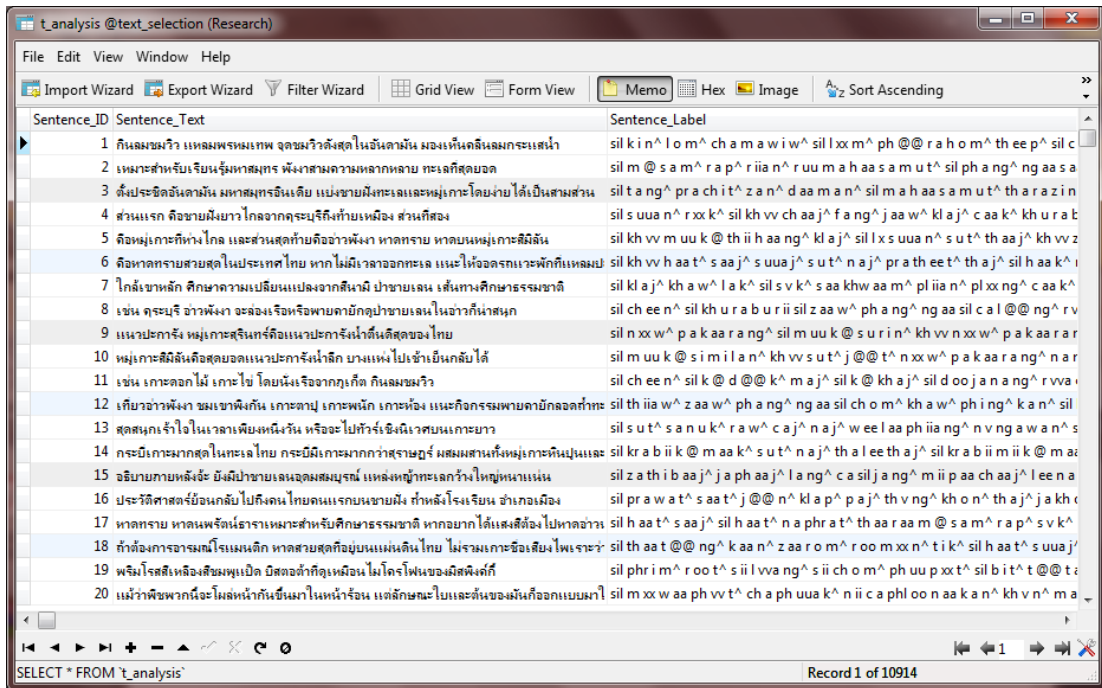
ตารางที่ 3.5 ตัวอย่างผลการคัดแยกประโยคตามอักขระนิพจน์ปรกติที่ไม่ต้องการออก

อายุ ๑๗ ปี นั่นมันชื่อของผม แต่อายุมากกว่าผมถึงห้าปี
ออกวิ่งไปเรื่อยๆ ตามถนนที่ทอดยาวออกจากหมู่บ้าน โลกผมกว้างขึ้นอีกแน่นอนครับ
เรื่องนิวยอร์ค 9-11 เป็นโศกนาฏกรรมครั้งยิ่งใหญ่ของโลกครั้งหนึ่ง และเป็นที่มาของการสร้างสรรค์วรรณูปชิ้นนี้ของวินทร์
หันไปดู 'เจ้าชายเรื่องมาก' เขากำลังจัดการกับม้า ซึ่งผ้าใบกันน้ำค้างผืนใหญ่ให้
Collage ภาพปะติด งานศิลปะแบบหนึ่งที่ใช้วิธีปะแผ่นวัสดุบนพื้นในลักษณะแบนราบหรืออ่อนสูงจากพื้นภาพเพียงเล็กน้อย
ขังแล้วค่อยเรียบบร้อยหน่อยง.มันก็มาทำให้หนูขังไม่ได้ ผมเจ็บกริบ
และมีศิลปะจะต้องผิดหวัง เวิร์ล ทำใจสบายๆ เกอะ อย่กั่วงวลไปเลย
จะต้องมีคุณสมบัติดังนี้ ๑ เป็นหน่วยงานภาครัฐ หรือองค์กรภาคเอกชน
เลี้ยงเพื่อนเมื่อทำบุญอายุ 6 รอบ 72 ปี เมื่อเดือนตุลาคม
4 จะต้องมึระเบียบ ความคิดรวบยอดและหลักการจะต้องเป็นจริงกระจ่างชัด
ชายหนุ่มจ้องหลอนตาเขม็ง ไม่เกรงไม่กลัวใดๆ หากเจ้าอยากแก๊งค์ข้าเจ้า
ที่ไหนสักแห่งในโลกคู่ขนานเหมือนในนิยายวิทยาศาสตร์ โลกแห่ง <Fail>ถ้าหากว่า.</Fail>
นึกแล้วว่าเธอต้องมา เรื่องแยๆ ที่แปลงเกษตรทำให้เธอแข็งใช้มียะ
ที่จริงก็น่าจะทำได้เพียงแต่ต้องเหนื่อยหน่อยเพราะจะต้องขุดหลุมกันคนละหลุมกว้างๆจนพอที่จะลงไปนอนได้
nonlover ซึ่งจะมีโอกาสและสถิติสัมปัญญะที่จะเลือก เพื่อนคู่ใจ ที่มีคุณภาพสูงกว่าผู้ที่ตกอยู่ในห้วงรัก
ข้อสงสัยต่างๆที่มีอยู่เดิมจึงหายไปหมดสิ้น พวกเราผลัดกันอาบน้ำ
บุตรบุญธรรมของข้าพเจ้าได้เรียนทำโดยตรงจากหม่อมเจ้าหญิงแย้มเยี่ยม เพื่อนฝูงใครๆที่เคยได้กิน
เค้ารู้ทันหรือกั งั้นพี่ตายก็แล้วกัน ไก่ย่างแน่เลย สวัสดิ์มึม ะ
เท่านั้น งานวิจัยเรื่อง การเมืองของชนชั้นกลางๆ ของพัสนัย ได้พิสูจน์ให้เห็นอย่างชัดเจนว่า
โดยเฉพาะฉินตะโรซึ่งเป็นลูกชายคนเดียวได้ฟังมานับครั้งไม่ถ้วนว่าแม่เกลียดพ่อที่สวมชุดสีฟ้าในวันเข้าพิธีแต่งงานเพียงใด
ที่สุดความรู้สึกหลังก็เป็นฝ่ายชนะ เต้าปุยค่อยๆ คีบเข้าใกล้ชายหนุ่มผู้ส่ายหน้าน้อยๆ

จากนั้นประโยคที่ผ่านการคัดกรองแล้วจะถูกนำไปแปลงรูปเขียนเป็นรูปอ่าน ในการทดลองของงานวิจัยนี้ได้ใช้ ตัวแปลงรูปเขียนเป็นรูปอ่านของ Probabilistic Generalized LR parser [12] ในงานวิจัยนี้ได้ทำการประยุกต์วิธีการแปลงรูปเขียนเป็นรูปอ่านในเรื่องของการเพิ่มหน่วยเสียงเงียบ (Silence) เข้าไปในประโยคด้วยการเพิ่มหน่วยเสียงเงียบไปที่ส่วนเริ่มของประโยค ส่วนที่เป็นวรรคของประโยค และส่วนท้ายของประโยค เทคนิคนี้สามารถช่วยเพิ่มรูปแบบหน่วยเสียงที่มีหน่วยเสียงเงียบเป็นส่วนประกอบได้ จากนั้นรูปเขียนของประโยคนั้น ๆ จะถูกเก็บเข้าสู่ฐานข้อมูล และในขั้นตอนนี้ประโยคออนไลน์ที่ไม่สามารถแปลงรูปเขียนเป็นรูปอ่านสำเร็จ จะถูกคัดออกอีกด้วยดังตารางที่ 3.5 ตัวอย่างผลลัพธ์ของประโยคที่แปลงรูปเขียนเป็นรูปอ่านสำเร็จในฐานข้อมูล ในภาพที่ 3.10 และจากตารางที่ 3.6



ภาพที่ 3.9 การแปลงประโยคจากรูปเขียนเป็นรูปอ่าน

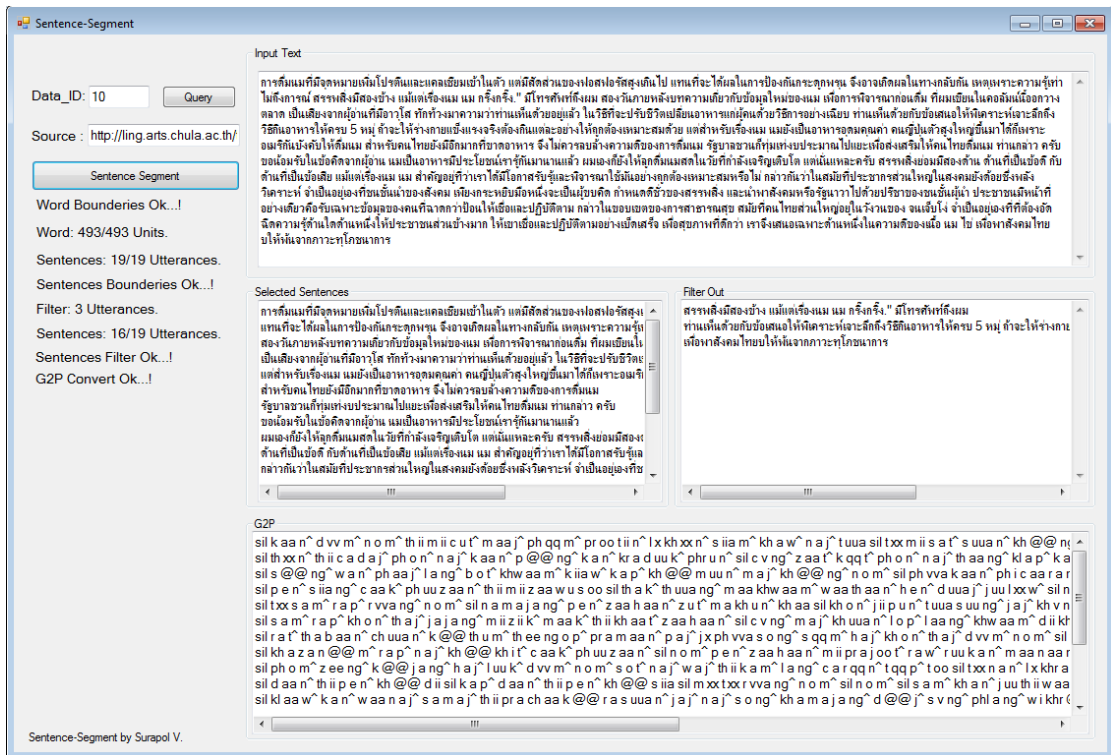


ภาพที่ 3.10 ตัวอย่างประโยคออนไลน์ในรูปแบบเขียนและรูปอ่านที่เก็บไว้ในฐานข้อมูล

ตารางที่ 3.6 ตัวอย่างผลการแปลงรูปเขียนเป็นรูปอ่านของประโยคออนไลน์

ประโยคในรูปเขียน	ประโยคในรูปอ่าน
เธอขอให้เพลงจบ จะเล่าเรื่องเด็กหนุ่มนักเปียโนให้สามีฟัง แต่ประพจน์เหมือนไม่ได้ยินอะไรทั้งสิ้น	sil th qq r @ @ h a j^ phl ee ng a c o p^ sil ca l a w^ r vva ng^ de k^ n u m^ n a k^ p iia n oo h a j^ s aa m i i f a ng^ sil txx pra ph o t^ m vva n^ m a j^ d a j^ j i n^ z a r a j^ th a ng^ s i n^ sil
เจอแล้ว ผมร้องออกมาแล้วก็เปิดประตู ผลัวะเข้าไป ลืมแม้กระทั่งว่าจะต้องเคาะประตูก่อน	sil c qq lxx w^ sil ph o m^ r @ @ ng^ z @ @ k^ m aa lxx w^ k @ @ p qq t^ pra t uu phl ua kh a w^ pa j^ sil l vv m^ m xx kra th a ng^ w aa ca t @ @ ng^ kh @ pra t uu k @ @ n^ sil
ทั้งที่จัดของเกือบเสร็จเรียบร้อย ของฝากญาติโกโหติกา ผองเพื่อนก็เตรียมไว้เสร็จสรรพ	sil th a ng^ th ii ca t^ kh @ @ ng^ k vva p^ s e t^ riia p^ r @ @ j^ sil kh @ @ ng^ faa k^ j aa t^ k oo h oo ti k aa sil ph @ @ ng^ ph vva n^ k @ @ triia m^ w a j^ s e t^ s a p^ sil
อ้าว คุณปู่ตอบแถมหัวเราะ นี่แหละบทวีของกวีเอกชาว อิตาลี	sil z aa w^ sil kh u n^ p uu t @ @ p^ k xx m^ h uua r @ sil n ii lxx bo t^ kw ii kh @ @ ng^ kw ii z ee k^ ch aa w^ sil z i t aa l i i sil
กานบัวสะบัดปาก ปล่อยเหยื่อปลา นกกวก กู้ฟ้าลั่นบึงน้ำ กระสาขาส่ายคล้ายยืนต้น	sil k aa p^ b uua s a ba t^ p aa k^ sil pl @ @ j^ j vva pl aa sil n o k^ k a w a k^ k uu faa l a n^ b v ng^ n a m^ sil kra s aa kh aa s aa j^ kh l aa j^ j a j vv n^ s a n^ sil
ถ้าฮิตเลอร์ไม่หันปลายปืนไปทางมอสโก	sil th aa hi t^ l qq m a j^ h a n^ pl aa j^ p vv n^ pa j^ th aa ng^ m @ @ t^ k oo sil j qq r a m a n^ k @ @ kh o

ประโยคในรูปเขียน	ประโยคในรูปอ่าน
เยอรมันก็คงไม่แพ้สงคราม	ng [^] m a j [^] ph xx s o ng [^] khr aa m [^] sil
ล้างจาน ปิดไฟนอน คืนนั้นทั้งคู่ร่วมรักกัน อย่างดุเดือด เปญคว่ำหน้า	sil l aa ng [^] c aa n [^] sil p i t [^] f a j [^] n @@ n [^] sil kh vv n [^] n a n [^] th a ng [^] kh uu r uua m [^] r a k [^] k a n [^] j aa ng [^] d u d vva t [^] sil p a j ee khw a m [^] n aa sil
พี่จิวทะเลาะกับพี่พาย พี่สะใภ้มีเวรดูแล คุณพ่อ แต่แผลอหับจนน้ำเกลือหมดขวด	sil ph ii c i w [^] th a l @ k a p [^] ph ii ph aa j [^] sil ph ii s a ph a j [^] m ii w ee n [^] d uu l xx kh u n [^] ph @@ sil t xx phl qq l a p [^] c o n [^] n a m [^] kl vva m o t [^] kh uua t [^] sil
ปรอดปราดเหล็กเสี้ยปีกบัด สาลิกากูร้อง ก่อนล่องลัด เขาชวาคล้ายหวัดหวาดขัน ร้อง	sil pr @@ t [^] pr aa t [^] l ii k [^] s iia p ii k [^] b a t [^] sil s aa l i k aa k uu r @@ ng [^] k @@ n [^] l @@ ng [^] l a t [^] sil kh a w [^] ch a w aa khl aa j [^] w a t [^] w aa t [^] kh a n [^] r @@ ng [^] sil
กลับกลายเป็นตรงกันข้าม ต่อมาทั้งสอง ได้รับความช่วยเหลือจากเบดูอินกลุ่มหนึ่ง	sil kl a p [^] kl aa j [^] p e n [^] tr o ng [^] k a n [^] kh aa m [^] sil t @@ m aa th a ng [^] s @@ ng [^] d a j [^] r a p [^] khw aa m [^] ch uua j [^] l vva c aa k [^] b ee d uu z i n [^] kl u m [^] n v ng [^] sil
แบบว่าซัดซีไคลในเวลาเดียวกัน แต่ก็หาย เร็วกว่าความนุ่มนวลโดยการทาขมิ้นแบบ ของยาย	sil b xx b a w aa kh a t [^] kh ii khl a j [^] n a j [^] w ee l aa d iia w [^] k a n [^] sil t xx k @@ h aa j [^] r e w [^] kw aa khw aa m [^] n u m [^] n uua n [^] d oo j [^] k aa n [^] th aa kh a m i n [^] b xx p [^] kh @@ ng [^] j aa j [^] sil
เหตุการณ์ได้เข้มข้นยิ่งขึ้นเมื่อทั้งสองถูก พวกโจรเบดูอินปล้นและทำร้ายบาดเจ็บ	sil h ee t [^] t u k aa n [^] d a j [^] kh ee m [^] kh o n [^] n a j i ng [^] kh v n [^] m vva th a ng [^] s @@ ng [^] th uu k [^] ph uua k [^] c oo n [^] b ee d uu z i n [^] pl o n [^] l x th a m [^] r aa j [^] b aa t [^] c e p [^] sil
ผมจะทำได้ดีป่า บอกหน่อย จะบ้าตายอยู่ แล้ว จะทำอะไรนองป้อปก็ตามประกบจน ผมทำอะไรไม่ถูกแล้ว	sil ph o m [^] c a th a m [^] ng a j [^] d ii p aa sil b @@ k [^] n @@ j [^] sil c a b aa t aa j [^] j uu l xx w [^] sil c a th a m [^] z a r a j [^] n @@ ng [^] p @@ p a k @@ t aa m [^] pr a k o p [^] c o n [^] ph o m [^] th a m [^] z a r a j [^] m a j [^] th uu k [^] l xx w [^] sil
แถมย่อตัวไปมาด้วยท่าทางเหมือนพิธีกร รายการสาระแนโชว์เสียหลายรอบ	sil th a m xx j @@ t uua p a j [^] m aa d uua j [^] th aa th aa ng [^] m vva n [^] ph i th ii k @@ n [^] r aa j [^] k aa n [^] s aa r a n xx ch oo s iia l aa j [^] r @@ p [^] sil
สมน้ำหน้าอย่างแรง คือดีนี่กว่าเป็นหนุ่ม เนื้อหอม ไข่ละแก คือกไม่ออก	sil s o m [^] n a m [^] n aa j aa ng [^] r xx ng [^] sil kh u j [^] d ii n a k [^] w aa p e n [^] n u m [^] n vva h @@ m [^] sil ng a j [^] l a k xx sil kh v k [^] m a j [^] z @@ k [^] sil
สงครามของหมาวัดค่อนข้างโหดร้ายและ โหด หมาในสภาพต่างเอาตัวรอดสั่งสม ความเลวมากกว่าความดี	sil s o ng [^] khr aa m [^] kh @@ ng [^] m aa w a t [^] kh @@ n [^] kh aa ng [^] h oo t [^] r aa j [^] l x ch oo t [^] sil m aa n a j [^] s a ph aa p [^] t aa ng [^] z a w [^] t uua r @@ t [^] s a ng [^] s o m [^] khw aa m [^] l ee w [^] m aa k [^] kw aa khw aa m [^] d ii sil
ชนเหลือบเขี้ยว อีนวล ไข่ลาย ชนบั้ง อี เหลือง ไข่มัด ดำมีด	sil kh o n [^] l vva p [^] kh iia w [^] sil z ii n uua n [^] sil z a j [^] l aa j [^] sil kh o n [^] b a ng [^] sil z ii l vva ng [^] sil z a j [^] m i t [^] sil d a m [^] m vv t [^] sil
อีขาว ไข่ต่าง ไข่แต้ม อีจุด ไข่ด้วน หางกุด ไข่ไม่ง	sil z ii kh aa w [^] sil z a j [^] d aa ng [^] sil z a j [^] t xx m [^] sil z ii c u t [^] sil z a j [^] d uua n [^] sil h aa ng [^] k u t [^] sil z a j [^] m oo ng [^] sil

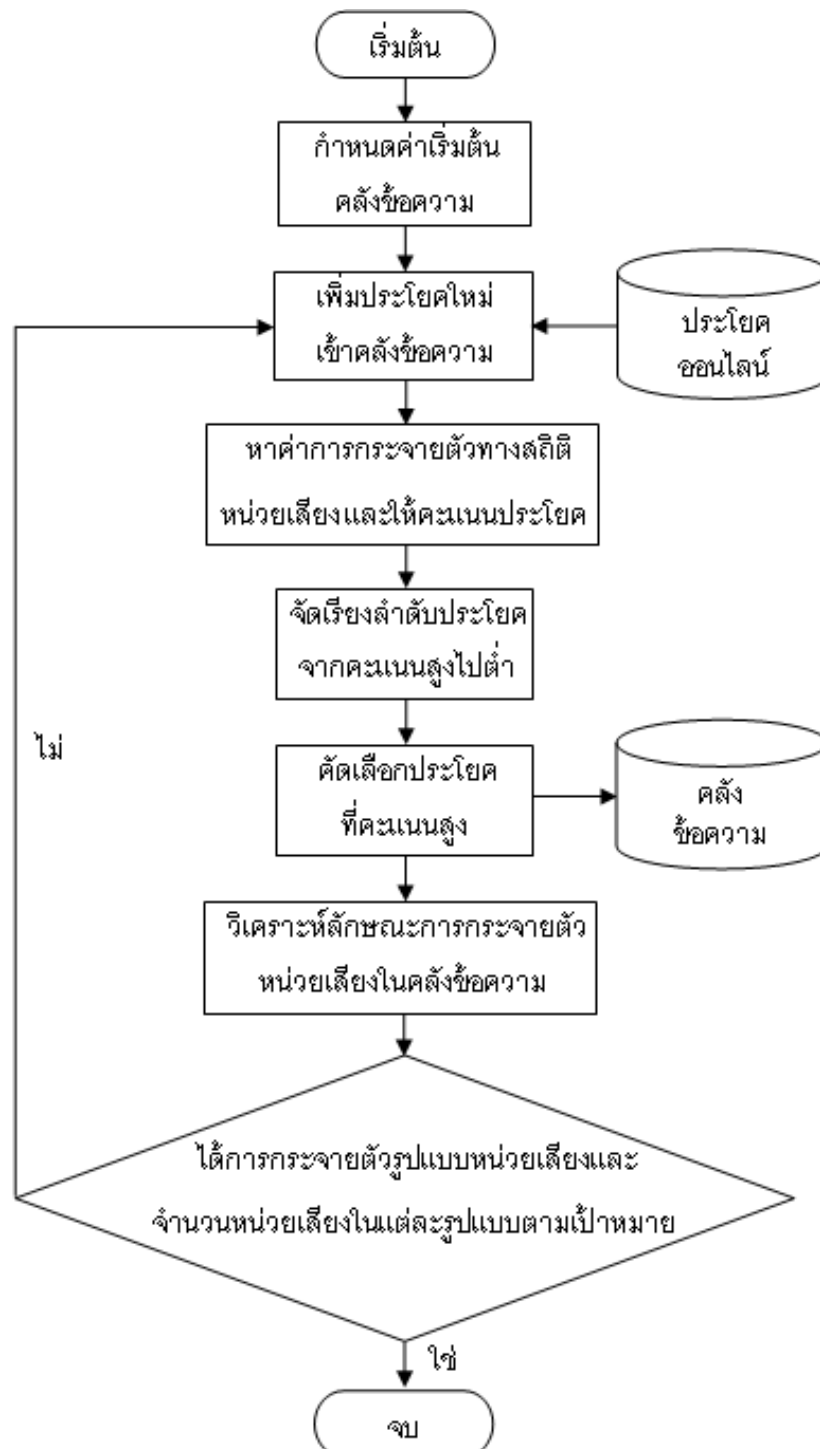


ภาพที่ 3.11 หน้าต่างโปรแกรมขั้นตอนการสร้างประโยคออนไลน์

หลังจากที่บทความจากฐานข้อมูลข้อความออนไลน์ได้ถูกจัดเก็บลงสู่ฐานข้อมูลรูปแบบประโยคและแปลงรูปเขียนของแต่ละประโยคเป็นรูปอ่านแล้ว รูปอ่านเหล่านี้จะถูกนำไปวิเคราะห์และให้คะแนนเพื่อคัดเลือกประโยคดังกล่าวเข้าสู่คลังข้อความที่ต้องการสร้าง ในขั้นตอนการเลือกประโยคและจัดเก็บข้อมูลต่อไป สามารถดูหน้าต่างโปรแกรมนี้จากภาพที่ 3.11

3. ขั้นตอนการเลือกประโยคและจัดเก็บคลังข้อความ

ขั้นตอนการเลือกประโยคและจัดเก็บคลังข้อความนี้ เป็นขั้นตอนที่สำคัญที่สุดในการทำวิทยานิพนธ์นี้ ขั้นตอนนี้มีหน้าที่สร้างฐานข้อมูลข้อความตามที่ต้องการ โดยการนำเอารูปเขียนของแต่ละประโยคจากขั้นตอนการสร้างประโยคออนไลน์ มาให้คะแนนแต่ละประโยค ขั้นตอนวิธีเชิงละโมภ (Greedy Algorithm) ถูกนำมาใช้ในการให้คะแนนดังกล่าว เพื่อให้คลังข้อความที่ถูกสร้าง มีจำนวนรูปแบบหน่วยเสียงที่ต้องการและให้การกระจายตัวทางหน่วยเสียงของคลังข้อความที่ถูกสร้างขึ้นคล้ายกับการกระจายตัวทางหน่วยเสียงของต้นแบบมากที่สุด ประโยคที่ได้คะแนนสูงจะถูกดึงเข้าสู่ฐานข้อมูล และประโยคที่ได้คะแนนต่ำจะถูกดึงออกไป ขั้นตอนนี้จะถูกทำแบบวนซ้ำไปเรื่อย ๆ จนกระทั่งคลังข้อความที่ต้องการสร้าง มีคุณสมบัติตรงกับสิ่งที่กำหนด แผนภาพการทำงานของขั้นตอนการเลือกประโยคและการจัดเก็บลงสู่ฐานข้อมูลนี้ สามารถดูได้จากภาพที่ 3.12



ภาพที่ 3.12 แผนภาพขั้นตอนการเลือกประโยคและจัดเก็บคลังข้อความ

ขั้นตอนการเลือกประโยคและจัดเก็บข้อมูล เริ่มจากการกำหนดค่าเริ่มต้นต่าง ๆ เพื่อใช้ในการสร้างคลังข้อความที่ต้องการ สิ่งที่ต้องกำหนดในขั้นตอนนี้ใช้ในการ กำหนดรูปแบบของคลังข้อความที่ต้องการสร้าง ซึ่งได้แก่

1. *รูปแบบของหน่วยเสียง* คือ ชนิดของรูปแบบหน่วยเสียงที่จะใช้ในการคำนวณทางสถิติในการเลือกประโยค เช่น หน่วยเสียงเดี่ยว (Monophone) หน่วยเสียงคู่ (Diphone) หน่วยเสียงสาม (Triphone) หรือหน่วยเสียงรูปแบบอื่น ๆ สำหรับวิทยานิพนธ์นี้ได้ใช้ รูปแบบหน่วยเสียงชนิดหน่วยเสียงคู่ ในการทดลอง
2. *จำนวนขั้นต่ำของรูปแบบหน่วยเสียง* คือ จำนวนของแต่ละรูปแบบหน่วยเสียงขั้นต่ำที่ต้องการให้เกิดขึ้นในคลังข้อความที่ต้องการสร้าง สำหรับในการทดลองของวิทยานิพนธ์นี้ ได้กำหนดให้รูปแบบของหน่วยเสียงเกิดอย่างน้อยหนึ่งครั้ง
3. *ขอบเขตของจำนวนพยางค์ในแต่ละประโยค* คือ การกำหนดความยาวของประโยคออนไลน์ด้วยจำนวนพยางค์ ซึ่งในการทดลองของงานวิจัยนี้ ในแต่ละประโยคกำหนดให้ประกอบด้วยจำนวนพยางค์ที่มากกว่า 20 พยางค์แต่ไม่เกิน 40 พยางค์
4. *ขนาดของจำนวนประโยคในคลังข้อความ* คือ จำนวนของประโยคที่จะเกิดขึ้นในคลังข้อความที่ต้องการ จำนวนของประโยคในคลังข้อความนี้มีผลต่อจำนวนความครอบคลุมของรูปแบบหน่วยเสียง โดยถ้าจำนวนประโยคมีน้อยเกินไปอาจทำให้ความครอบคลุมของแต่ละรูปแบบหน่วยเสียงและจำนวนขั้นต่ำของรูปแบบหน่วยเสียง ไม่ได้ตามเป้าหมายที่ต้องการได้ สำหรับการทดลองในวิทยานิพนธ์นี้ได้กำหนดจำนวนประโยคในคลังข้อความไว้ที่ 1,000 ประโยค
5. *ค่าเบี่ยงเบนเฉลี่ยร้อยละจากเป้าหมาย* คือ ค่าผลต่างที่ยอมรับได้ของการกระจายตัวรูปแบบหน่วยเสียงคลังข้อความที่กำลังสร้างเทียบกับการกระจายตัวรูปแบบหน่วยเสียงของเป้าหมายในรูปแบบของร้อยละ ในที่นี้ถ้าค่าเบี่ยงเบนเฉลี่ยจากเป้าหมายมีค่าน้อยจะทำให้คลังเสียงที่ต้องการสร้างมีการกระจายตัวทางหน่วยเสียงที่คล้ายกับเป้าหมายมากขึ้น แต่อย่างไรก็ตาม ถ้ากำหนดค่าเบี่ยงเบนเฉลี่ยจากเป้าหมายน้อยเกินไปอาจทำให้การสร้างคลังข้อความเป็นไปได้ยากหรือใช้ระยะเวลาามากไปได้ ในการทดลองของวิทยานิพนธ์นี้ค่าเบี่ยงเบนเฉลี่ยจากเป้าหมายที่ ร้อยละ 3.0

เมื่อกำหนดรูปแบบของคลังข้อความที่ต้องการแล้ว ขั้นตอนต่อไประบบจะทำการดึงประโยคจากฐานข้อมูลประโยคออนไลน์เพิ่มเข้าไปในฐานข้อมูลวิเคราะห์ประโยค ฐานข้อมูลนี้เป็นฐานข้อมูลชั่วคราว ใช้ในการคำนวณหาค่าทางสถิติทางการกระจายตัวหน่วยเสียง เพื่อหาความต้องการในแต่ละรูปแบบของหน่วยเสียงเพื่อให้คะแนนประโยค และใช้ในการหาค่าเบี่ยงเบนเฉลี่ย

จากเป้าหมาย วิธีเพิ่มประโยคออนไลน์เข้าไปในฐานข้อมูลวิเคราะห์ประโยคนี้ ในรอบแรกของขั้นตอนนี้จะเราจะเพิ่มไปเท่ากับ 110% ของขนาดคลังข้อความ เช่นถ้าคลังข้อความที่ต้องการสร้างกำหนดขนาดไว้ที่ 1,000 ประโยคในการเพิ่มประโยคครั้งแรก 1,100 ประโยค และในรอบต่อไปของขั้นตอนนี้จะเพิ่มประโยคเข้าไปอีกทีละ 10% ของจำนวนคลังข้อความที่ต้องการสร้างเช่นถ้าคลังข้อความมีขนาด 1,000 จะเพิ่มประโยคเข้าไปอีก 100 เป็นต้น ทั้งนี้อีก 10% ของประโยคจะถูกคัดเลือกทิ้งไปในขั้นตอนการคัดเลือกประโยค เมื่อเพิ่มประโยคเข้าไปแล้วระบบจำทำการคำนวณหาค่าทางสถิติของการกระจายตัวทางหน่วยเสียง จากฐานข้อมูลวิเคราะห์ประโยค ตัวอย่างข้อมูลทางสถิติของการกระจายตัวทางหน่วยเสียงดังกล่าว ดูได้จากตารางที่ 3.7

ข้อมูลทางสถิติของการกระจายตัวทางหน่วยเสียงนี้ ใช้เพื่อการศึกษาในการหาค่าความต้องการในแต่ละรูปแบบของหน่วยเสียงและใช้ในการหาค่าเบี่ยงเบนเฉลี่ยจากเป้าหมาย ซึ่งได้มาจากการการนำรูปอ่านของประโยคออนไลน์ในคลังเสียงที่กำลังสร้างมาแจกแจงและนับจำนวนเทียบกับรูปแบบหน่วยเสียงของเป้าหมาย ข้อมูลทางสถิติทั้งหมดนี้ ประกอบด้วย

1. **Patterns** คือ รูปแบบของหน่วยเสียงที่ปรากฏ สร้างมาจากขั้นตอนการสร้างการกระจายตัวทางหน่วยเสียงเป้าหมายได้มาจากการใช้ฟังก์ชัน HLEd จากเครื่องมือ Hidden Markov Model toolkit, HTK [31] ในการหารูปแบบของหน่วยเสียงคู่ ที่เกิดขึ้นจากรูปอ่านทั้งหมดจากคลังข้อความเป้าหมาย ในที่นี้ปรากฏรูปแบบหน่วยเสียงคู่ที่เกิดขึ้น 1,383 หน่วยเสียง
2. **Target_Freq** คือ จำนวนความถี่ที่เกิดขึ้นของหน่วยเสียงในเป้าหมาย สร้างมาจากขั้นตอนการสร้างการกระจายตัวทางหน่วยเสียงเป้าหมายเช่นกัน โดยใช้ฟังก์ชัน HLStats จากเครื่องมือ Hidden Markov Model toolkit, HTK [31] ในการนับจำนวนที่เกิดขึ้นในแต่ละรูปแบบของหน่วยเสียงในคลังข้อความเป้าหมาย
3. **Target_Prob** คือ โอกาสที่รูปแบบของแต่ละหน่วยเสียงจะเกิดขึ้นในคลังข้อความเป้าหมาย เพื่อให้อัตราการเกิดรูปแบบหน่วยเสียงของเป้าหมายและอัตราการเกิดรูปแบบหน่วยเสียงในคลังข้อความที่ต้องการสร้างอยู่ในบรรทัดฐานเดียวกัน ใช้ในการเปรียบเทียบการกระจายตัวรูปแบบหน่วยเสียงของทั้งสอง ค่าวนได้จากสมการที่ (1)

ตารางที่ 3.7 ข้อมูลทางสถิติการกระจายตัวทางหน่วยเสียงที่ใช้ในการคำนวณจากฐานข้อมูล

Patterns	Target_Freq	Target_Prob	Custom_Freq	Custom_Prob	Weight	Dellta_Prob	Deviation
ng ⁺ z	269	-2.89743	73	-2.98	0.0825672	0.0825672	0.0284967
ph+uu	517	-2.61407	162	-2.63543	0.0213612	0.0213612	0.00817161
b+u	157	-3.1307	17	-3.60325	0.472546	0.472546	0.15094
n+v	282	-2.87697	92	-2.88014	0.00317279	0.00317279	0.00110282
vv+s	64	-3.51843	26	-3.42304	-0.0953913	0.0953913	0.0271119
t ⁺ +ng	11	-4.26729	13	-3.71595	-0.551339	0.551339	0.129201
ii+c	373	-2.75569	61	-3.05741	0.301719	0.301719	0.10949
m ⁺ +pr	70	-3.4798	16	-3.6288	0.149001	0.149001	0.0428187
@+s	16	-4.1105	8	-3.91687	-0.193634	0.193634	0.0471073
k ⁺ +k	311	-2.83453	91	-2.88486	0.0503334	0.0503334	0.0177572
i+d	35	-3.77776	15	-3.65595	-0.121807	0.121807	0.0322432
ee+l	84	-3.40113	29	-3.37646	-0.0246675	0.0246675	0.00725273
i+kh	16	-4.1105	6	-4.03337	-0.0771288	0.0771288	0.0187639
m ⁺ +m	290	-2.86484	133	-2.7208	-0.144037	0.144037	0.0502774
uu+l	61	-3.53911	22	-3.4941	-0.045008	0.045008	0.0127173

$$P(\text{Target}_i) = \log_{10} \frac{F_i + 0.5}{n + 0.5} \quad (1)$$

โดยที่ F_i คือ จำนวนความถี่ที่รูปแบบหน่วยเสียงที่ i เกิดขึ้นในคลังข้อความเป้าหมาย และ n คือผลรวมทั้งหมดของหน่วยเสียงในคลังข้อความเป้าหมาย และด้วยลักษณะการกระจายตัวของหน่วยเสียงในภาษารวมชาติมีความสัมพันธ์กันแบบลอการิทึม จึงใช้ฟังก์ชันลอการิทึมฐานสิบ ส่วนค่าคงที่ 0.5 ใช้ในการบังคับค่าจากการคำนวณไม่ให้เกิดความผิดพลาดในกรณีที่ $F_i = 0$

4. **Custom_Freq** คือ จำนวนความถี่ที่เกิดขึ้นของหน่วยเสียงในคลังเสียงที่กำลังสร้าง สร้างจากการนับจำนวนที่เกิดขึ้นในแต่ละรูปแบบของหน่วยเสียงในคลังข้อความที่กำลังสร้าง โดยใช้ ฟังก์ชัน HLStats จากเครื่องมือ Hidden Markov Model toolkit, HTK [31]
5. **Custom_Prob** คือ โอกาสที่หน่วยเสียงนี้เกิดขึ้นในคลังข้อความที่กำลังสร้าง คำนวณได้จากสมการที่ (2)

$$P(\text{Custom}_i) = \log_{10} \frac{F_i + 0.5}{n + 0.5} \quad (2)$$

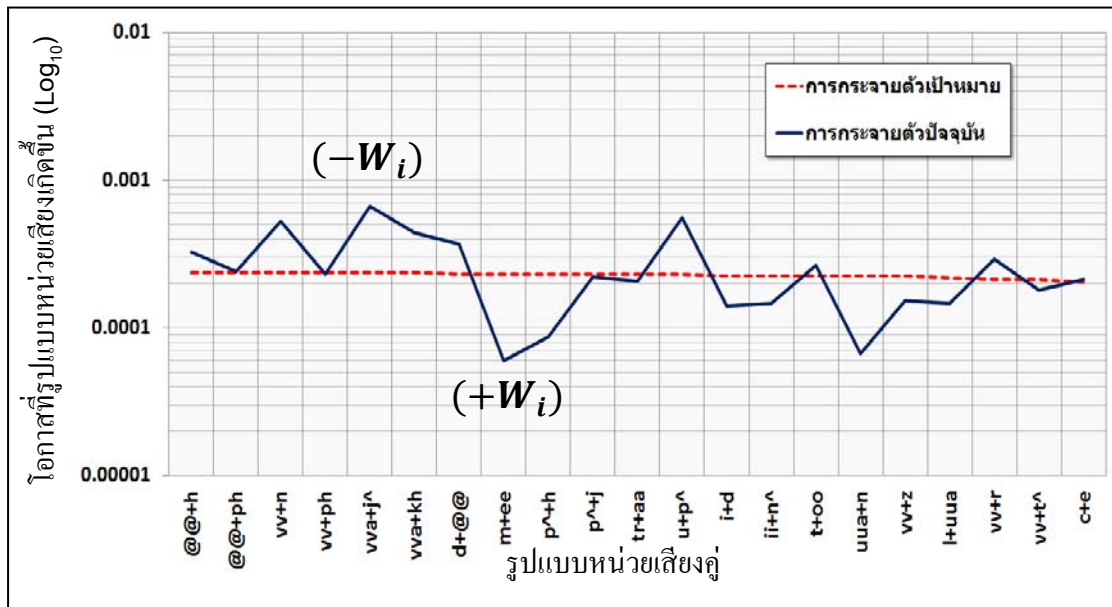
โดยที่ F_i คือ จำนวนความถี่ที่รูปแบบหน่วยเสียงที่ i เกิดขึ้นในคลังข้อความที่กำลังสร้าง และ n คือผลรวมทั้งหมดของหน่วยเสียงในคลังข้อความที่กำลังสร้าง

6. **Weight** คือ ค่าของความต้องการในของรูปแบบหน่วยเสียง เกิดจากผลต่างของโอกาสการเกิดของรูปแบบหน่วยเสียงนั้น การคำนวณค่าของความต้องการรูปแบบหน่วยเสียงสำหรับงานวิทยานิพนธ์นี้มีสองกรณี กรณีแรกคือ จำนวนของรูปแบบหน่วยเสียงผ่านเกณฑ์ขั้นต่ำที่กำหนดแล้ว จะคำนวณค่าของความต้องการรูปแบบหน่วยเสียงได้จากสมการที่ (3)

$$W_i = P(\text{Target}_i) - P(\text{Custom}_i) \quad (3)$$

โดยที่ $P(\text{Target}_i)$ คือ โอกาสที่รูปแบบของหน่วยเสียงที่ i เกิดขึ้นในคลังข้อความเป้าหมาย และ $P(\text{Custom}_i)$ คือ โอกาสที่รูปแบบของหน่วยเสียงที่ i เกิดขึ้นในคลังข้อความที่กำลังสร้าง ถ้าผลของ W_i เป็นบวก (+) แสดงว่าหน่วยเสียงนั้นในคลังข้อความที่ต้องการสร้างยังขาดอยู่ ดังนั้นระบบจึงต้องการรูปแบบหน่วยเสียงนี้เพิ่มอีก แต่ถ้าค่าของ W_i ออกมาเป็นค่าลบ (-) แสดงว่าหน่วยเสียงนั้นในคลังข้อความที่ต้องการ

สร้าง มีมากเกินไป ระบบจึงต้องการลดรูปแบบหน่วยเสียงนี้ออก เพื่อความเข้าใจมากขึ้นสามารถดูภาพที่ 3.13 ประกอบได้



ภาพที่ 3.13 ใช้ค่าความต้องการรูปแบบหน่วยเสียงในการให้คะแนนประโยค

ส่วนการคำนวณค่าความต้องการรูปแบบหน่วยเสียง ในกรณีที่สอง คือ จำนวนรูปแบบหน่วยเสียงนั้นยังต่ำกว่าเกณฑ์ขั้นต่ำที่กำหนด เนื่องจากเราต้องการรักษาหน่วยเสียงที่หายากไว้เพื่อให้ได้ความครอบคลุมหน่วยเสียงตามเป้าหมาย การคำนวณค่าของความ ต้องการรูปแบบหน่วยเสียง จากสมการที่ (4) ในการคำนวณ

$$W_i = |P(\text{Custom}_i)| \times \alpha \quad (4)$$

หลังจากทำการหาค่าสมบูรณ์ของโอกาสที่หน่วยเสียงนี้เกิดขึ้นในคลังข้อความที่กำลังสร้างแล้ว เราจะทำการการเพิ่มระดับความสำคัญให้รูปแบบหน่วยเสียงหายากด้วยการคูณด้วยค่าคงที่ α ถ้ากำหนดให้ค่า α มีค่ามากจะทำให้คลังข้อความที่ถูกสร้างเข้าถึงการครอบคลุมทางหน่วยเสียงเป้าหมายได้เร็วขึ้น แต่อย่างไรก็ตามถ้ากำหนดค่า α มากเกินไปจะทำให้ การกระจายตัวของคลังข้อความที่ต้องการสร้างเป็นไปตามเป้าหมายยากมากขึ้น ในวิทยานิพนธ์นี้ได้ใช้ ค่า α ที่ 1.5

7. **Delta_Prob** คือ ผลต่างระหว่างโอกาสที่เกิดขึ้นของรูปแบบหน่วยเสียงเป้าหมายและปัจจุบัน ใช้เพื่อการคำนวณค่าเบี่ยงเบนจากเป้าหมายของแต่ละหน่วยเสียง คำนวณจากสมการที่ (5)

$$\Delta P_i = |P(Target_i) - P(Custom_i)| \quad (5)$$

โดยที่ $P(Target_i)$ คือ โอกาสที่รูปแบบของหน่วยเสียงที่ i เกิดขึ้นในคลังข้อความเป้าหมาย และ $P(Custom_i)$ คือ โอกาสที่รูปแบบของหน่วยเสียงที่ i เกิดขึ้นในคลังข้อความที่กำลังสร้าง

8. **Deviation** คือ ผลต่างระหว่างโอกาสที่เกิดขึ้นของรูปแบบหน่วยเสียงเป้าหมายและปัจจุบันเทียบกับโอกาสที่เกิดขึ้นของเป้าหมาย ใช้คำนวณ ค่าเบี่ยงเบนเฉลี่ยร้อยละจากเป้าหมาย (Percent average deviation of target) ดูได้จากสมการที่ (6)

$$\sigma_i = \frac{\Delta P_i}{P(Target_i)} \quad (6)$$

โดยที่ σ_i คือ ผลต่างระหว่างโอกาสที่เกิดขึ้นของรูปแบบหน่วยเสียงเป้าหมายและปัจจุบัน และ $P(Target_i)$ คือ โอกาสที่รูปแบบของหน่วยเสียงที่ i เกิดขึ้นในคลังข้อความเป้าหมาย ถ้าค่าของ σ_i มีค่ามากรูปแบบหน่วยเสียงนั้นห่างจากเป้าหมายอยู่มากและถ้าค่าของ σ_i มีค่าน้อย รูปแบบหน่วยเสียงนั้นเข้าใกล้เป้าหมายมากขึ้น

ขั้นตอนต่อไปคือการให้คะแนนประโยคจากข้อมูลการกระจายตัวทางสถิติ เพื่อใช้พิจารณาในการเลือกประโยคที่จะเก็บไว้ในคลังข้อความที่ต้องการสร้าง การให้คะแนนประโยคในวิทยานิพนธ์นี้ ได้ใช้สมการที่ (7) ในการให้คะแนน

$$Score_c = \sum_{i=1}^n \frac{F_i \times W_i}{L_c} \quad (7)$$

ในที่นี้ F_i คือ จำนวนรูปแบบของหน่วยเสียง i ที่เกิดในคลังข้อความที่กำลังสร้างปัจจุบัน ทั้งนี้ในประโยคที่กำลังพิจารณาอาจไม่ปรากฏรูปแบบของหน่วยเสียงที่ i ก็ได้ ส่วน W_i คือ ค่าของความต้องการในของรูปแบบหน่วยเสียงนั้นจากสมการที่ (3) และ L_c คือ ความยาวของประโยคนับจากผลรวมของจำนวนหน่วยเสียงในประโยคที่ C ค่า L_c ใช้ในเพื่อทำให้การให้คะแนนไม่ขึ้นอยู่กับความยาวของประโยค n คือจำนวนรูปแบบของหน่วยเสียงทั้งหมดที่ปรากฏในคลังข้อความเป้าหมาย สำหรับคะแนนประโยคนี้ ถ้าค่า $Score_c$ ในประโยคมีค่ามาก ประโยคนั้นจะประกอบด้วยรูปแบบหน่วยเสียงที่ต้องการหรือหน่วยเสียงที่ขาดอยู่มาก แต่ถ้าประโยคมีคะแนนน้อย

ประโยคนั้นอาจประกอบด้วยรูปแบบหน่วยเสียงไม่จำเป็นหรือมีมากเกินไป สำหรับการสร้าง คลังข้อความที่ต้องการ สามารถดูตัวอย่างผลการให้คะแนนประโยคได้จากตารางที่ 3.8

ตารางที่ 3.8 ผลคะแนนประโยคจากขั้นตอนการให้คะแนน

ประโยคในฐานข้อมูล	คะแนน
ถึงจะเหนื่อย ฉันก็ไม่มีวันถอย ฉันจะบุกไปจนกว่าจะถึงบึงหญ้าใหญ่	0.1113
ในแง่สุขภาพทางจิต ปรากฏว่า คนชราที่รู้สึกว่าคุณถูกทอดทิ้ง	-0.0107
หากเมื่อเจ้าตัวไม่ยอมเอ่ยปากพูดอะไรสักคำ ชายหนุ่มก็สนใจ ได้แต่หวังว่าสาวไทยจะ สัมผัสได้ถึงความห่วงใยเหลือล้นของเขา	0.0257
มือหนึ่งถือมีดเหลาด้ามยาว อีกมือถือชอล์กไม้แห้งผิวเหลืองมัน หลวงลุงใช้ให้ถากเหลา เกลารูปชอล์กไม้เป็นไม้เหน็บประดับธรรม	-0.0151
น่าเสียดายนะที่แกได้เรียนน้อยไปหน่อย แต่ฉันดูแกฉลาดเฉลียวกว่าคนอื่น	0.1585
ให้คำปรึกษาแนะนำเกี่ยวกับการเสริมสร้างวินัยและพัฒนาให้ข้าราชการมีวินัยแก่	-0.0138
ผู้ส่งสาร กลายเป็นสิ่งที่มีตัวตน มีมิติด้านเศรษฐกิจการเมือง	0.0495
เดินออกมาจากห้อง ผมใช้ทิศทางเดิมที่ผมสำรวจไว้และทดลองมาแล้วเมื่อคืนก่อน	-0.0338
การผดุงครรภ์ หรือการพยาบาลและการผดุงครรภ์ ตามระเบียบที่รัฐมนตรีกำหนดโดย ประกาศในราชกิจจานุเบกษา	0.0864
ความมีชีวิตชีวา สดใส แม้จะทำให้เขาปวดหัวในช่วงแรกที่ได้พบ	-0.0184
ชอบนึกว่าตัวเองดีจนคนอื่นต้องคอยอิจฉา ฉันพยักหน้า แอมมิเลียจะไม่เชื่อกี่ช่าง หล่อน	-0.0067
ทำไงดีมาเธียส คุณเป็นซูบิให้ฉันไม่ได้หรอก นี่ชบ่นเสียงแหลม เมื่อเห็นประกาศจาก ทางคณะตั้งแต่ต้นเทอมสอง	-0.0329
เข้ามาแต่เต็มหัวใจที่เว้าแหว่งเพราะพิษความทุกข์ได้อย่างคาดไม่ถึง จากที่เคยคิดว่า เป็นเพียงเด็กสาวสอดรู้สอดเห็น	0.0533
ในใจได้บ้าง หากภาพของผู้คนแม้จะบางตา แต่ก็มักเป็นครอบครัวส่วนใหญ่	-0.0032
ข้ากำเนิดดินน้ำลมและไฟหนึ่งเกิดอาจดับไหม้ในพริบตา ข้าท่องฟ้าเมฆหมอกนอกหน หาวหนึ่งดวงดาวพรราวแสงแห่งเวหา	0.2096
ในการอำนวยความสะดวก คุ่มครองสิทธิเสรีภาพประชาชน ป้องกันและควบคุม อาชญากรรม	0.0151

ประโยคในฐานข้อมูล	คะแนน
เพื่อนคนหนึ่งเลยชักชวนไปขายผลไม้สดซึ่งที่ภูเก็ต ฉันไม่กล้าไปทำงานแบบนั้นหรอก	0.0773
เป้าหมายอย่างมีประสิทธิภาพและประสิทธิผล ปฏิบัติงานร่วมกับหรือสนับสนุนการปฏิบัติงานของหน่วยงานอื่นที่เกี่ยวข้องหรือที่ได้รับมอบหมาย	0.0161
หรือหน่วยกรรมรอกก็ได้ วิจินตน์เรียกคำเชื่อมที่เชื่อมอนุพากย์คุณศัพท์กับคำนามที่มันขยายว่า	-0.0175
ผลดีอีกอย่างคือได้มีเวลาดูแลนักเรียนใกล้ชิดยิ่งขึ้น ครูวิชัยพยายามรวบรวมนักเรียนให้มากินข้าวด้วยกันที่ระเบียบ	-0.0179
ไม่กี่นาทีก็รู้กันทั้งโรงเรียนแล้ว ว่าแต่เธอเห็นนั่นมั๊ย รินซีไปที่พอลกับแซมซึ่งยังคงมองไปที่เวที	-0.0111
อยากเป็นหอเป็นเรือนเยาะไร้ไถนา มีลูกสักสองสามคน หน้าแดงขึ้นวูบหนึ่ง	0.1036
เกี่ยวกับขั้นตอนในการดำเนินงานของจังหวัด นายอารีย์ วงศ์อารยะ	-0.0160
หญิงสาวไม่เข้าใจว่ากำลังเกิดอะไรขึ้นกับตัวเอง พอรู้ตัวตั้งสติได้	-0.0270
หาดทราย หาดนพรัตนธาราเหมาะสำหรับศึกษาธรรมชาติ หากอยากได้แสงสีต้องไปหาดอ่าวนาง	-0.0095
ท่าทางจะอายุจริง อาย คือตอนนี่ก็กำลังติดหนังสืออย่างรุนแรงนะคะพี่แหวน	0.1767
ขัดกับภาพสาวมาดแครงที่มาพร้อมกับแว่นตากรอบดำทรงกลม ก็ทำเอาคนโดนถามอดมยืมกับจานสปาเกิดตีตรงหน้าไม่ได้	0.0238
ริมถนน รถเข็นขายขนมสี่ชั้นดูจาดจอนิ่งอยู่ สีสนของขนมบนรถดูน่ากลัวมากกว่าน้ำกิน	0.1445
เพื่อการพิจารณาก่อนดื่ม ที่ผมเขียนในคอลัมน์นี้ออกวางตลาด เป็นเสียงจากผู้อ่านที่มีอาวุโส	0.0650
ลามใหม่ทั้งด้านเหนือและใต้มุ่งไปทางทิศตะวันตกเข้าหาลำน้ำแม่ยมซึ่งมีวัดราชธานีที่เต็มไปด้วยโบราณวัตถุวางอยู่	0.1996
ยังผู้หญิงทั้งคู่นี้อีกเล่า อาจเสแสร้งแกล้งเข้ามาตีสนิท เพื่อล้วงเอาความสงสัยออกมาให้สิ้น	0.1254
โครงการศึกษาภาษาและวัฒนธรรมท้องถิ่นของข้าราชการจังหวัดชายแดนภาคใต้	0.1184
ทิพกฤตาเห็นฝายนั่นเงยหน้าขึ้นมองท้องฟ้าสีหม่น มือใหญ่ยกขึ้นปาดหยดน้ำซึ่งเริ่มจับตัวบนศีรษะได้รูปนั้น	0.1077

เมื่อระบบคำนวณคะแนนของแต่ละประโยคจนครบทุกประโยคแล้ว ประโยคเหล่านั้นจะถูกทำการจัดอันดับประโยคตามคะแนนของประโยคจากมากไปหาน้อย เพื่อความรวดเร็วในการจัดอันดับประโยค ได้ใช้ฟังก์ชันการจัดเรียงตามลำดับประโยคอัตโนมัติด้วยคำสั่งจากฐานข้อมูล ,MySQL Server 5.1 [30]

เมื่อจัดเรียงประโยคจากมากไปน้อยแล้วระบบจะทำการคัดประโยคที่คะแนนตกอันดับออก ในที่นี้ได้ตัดประโยคออก 10% จากคลังข้อความปัจจุบันออก เช่น ถ้าปัจจุบันในคลังข้อความที่กำลังสร้างมีขนาด 1,100 ประโยค ประโยคที่ได้คะแนน 100 อันดับสุดท้ายจะถูกตัดออก เมื่อคัดประโยคออกแล้ว ระบบจะคัดลอกประโยคที่เหลือในฐานข้อมูลคลังข้อความวิเคราะห์นำไปใส่ในฐานข้อมูลคลังข้อความที่ต้องการสร้าง

จากนั้นระบบจะทำการวิเคราะห์ลักษณะการกระจายตัวหน่วยเสียงในคลังข้อความที่สร้าง เพื่อตรวจสอบการกระจายตัวรูปแบบหน่วยเสียงของเป้าหมายและการกระจายตัวที่เกิดขึ้นในคลังข้อความที่สร้าง ดังนั้นขั้นตอนนี้จึงทำการหาค่าเบี่ยงเบนเฉลี่ยร้อยละจากเป้าหมาย (Percent average deviation of target) จากสมการที่ (8)

$$Avg\ Deviation(\%) = \frac{\sum_{i=1}^n \sigma_i \times 100}{n} \quad (8)$$

โดยที่ σ_i คือ ผลต่างระหว่างโอกาสที่เกิดขึ้นของรูปแบบหน่วยเสียงเป้าหมายและปัจจุบันเทียบกับโอกาสที่เกิดขึ้นของเป้าหมาย จากสมการที่ (6) ส่วน n คือจำนวนรูปแบบหน่วยเสียงที่เกิดขึ้นทั้งหมดของเป้าหมาย

ขั้นตอนต่อไป ระบบจะทำการตรวจสอบ ว่าคลังข้อความที่สร้างขึ้นมีคุณสมบัติตามที่ต้องการสร้างหรือไม่ โดยการเปรียบเทียบการกระจายตัวรูปแบบหน่วยเสียง จากค่าเบี่ยงเบนเฉลี่ยร้อยละจากเป้าหมาย และจำนวนขั้นต่ำของรูปแบบของแต่ละหน่วยเสียง ว่าผ่านเกณฑ์ที่กำหนดไว้หรือไม่ ถ้าค่าดังกล่าวยังไม่ผ่านเกณฑ์ ระบบจะกลับเข้าสู่ขั้นตอนเพิ่มประโยคออนไลน์เข้าสู่คลังข้อความใหม่ แต่ถ้าค่าของทั้งสองผ่านเกณฑ์ระบบจะหยุดสร้างประโยคออนไลน์เพิ่มและหยุดกระบวนการเลือกประโยค เป็นการเสร็จสิ้นการทำงาน คลังข้อความที่ถูกสร้างขึ้นนี้สามารถนำไปใช้ได้ทันที สามารถดูประโยคที่ถูกเลือก 300 อันดับแรกได้จากภาคผนวก

บทที่ 4

การทดลอง และอภิปรายผล

การทดลอง

ในการทดลองระบบการเลือกข้อความอัตโนมัติเพื่อนสร้างคลังข้อความตามการกระจายตัวหน่วยเสียงที่กำหนดได้นั้น ได้แบ่งการทดลองเป็นสี่การทดลองคือ การทดลองแรกเป็นการทดสอบระบบกับรูปแบบการกระจายตัวทางหน่วยเสียงของเป้าหมายจากคลังข้อความต้นแบบ การทดลองที่สองเป็นการทดสอบระบบกับรูปแบบการกระจายตัวทางหน่วยเสียงแบบกำหนดเอง โดยไม่ได้นำรูปแบบการกระจายตัวมาจากรูปแบบการกระจายตัวของคลังข้อความอื่น ๆ ซึ่งในที่นี้ กำหนดให้ทุกรูปแบบหน่วยเสียงมีโอกาสเกิดเท่ากันทุกหน่วย การทดลองที่สามเป็นการทดลองเพื่อพิจารณาผลของการปรับพารามิเตอร์ค่าคงที่ความสำคัญรูปแบบหน่วยเสียงหายาก (α) ซึ่งใช้ในการคำนวณค่าความต้องการรูปแบบหน่วยเสียง การทดลองที่สี่เพื่อทดสอบประสิทธิภาพของงานวิจัยนี้วัดผลเทียบกับการเลือกประโยคแบบเก่าจากงานวิจัย W. Zhang, 2010: Automatic Construction for a TTS Corpus with Limited Text [11] ซึ่งเป็นงานวิจัยที่ใช้หลักการเลือกประโยคคล้ายกันแต่ต่างกันที่วิธีการคิดคะแนน

การสร้างคลังข้อความจากรูปแบบการกระจายตัวทางหน่วยเสียงจากการกระจายตัวเป้าหมาย

ในการทดลองนี้ ต้องการสร้างฐานข้อมูลข้อความภาษาไทยขนาด 1,000 ประโยค แต่ละประโยคประกอบด้วยจำนวนพยางค์มากกว่า 20 พยางค์แต่ไม่เกิน 40 พยางค์ ชนิดของหน่วยเสียงใช้รูปแบบ หน่วยเสียงคู่ ความครอบคลุมทางหน่วยเสียง 100% เทียบจากจำนวนรูปแบบหน่วยเสียงที่เกิดขึ้นจากรูปแบบหน่วยเสียงของเป้าหมาย แต่ละหน่วยเสียง ต้องเกิดมากกว่า 1 ครั้ง ค่าคงที่ความสำคัญรูปแบบหน่วยเสียงหายาก (α) ไว้ที่ 1.5 ค่าเบี่ยงเบนเฉลี่ยร้อยละจากเป้าหมายต้องน้อยกว่า 3.0 % ของรูปแบบของการกระจายตัวเป้าหมาย การทดลองมีขั้นตอนดังนี้

1. การสร้างการกระจายตัวทางหน่วยเสียงเป้าหมาย

ในการทดลองได้เลือกใช้การกระจายตัวของหน่วยเสียงเป้าหมาย จากคลังเสียงพูด Large Vocabulary Continuous Speech Recognition (LVCSR) corpus for Thai language[18] เลือกใช้ชุดครอบคลุมคำศัพท์ 5,000 คำ ฐานข้อมูลข้อความนี้คัดเลือกประโยคมา

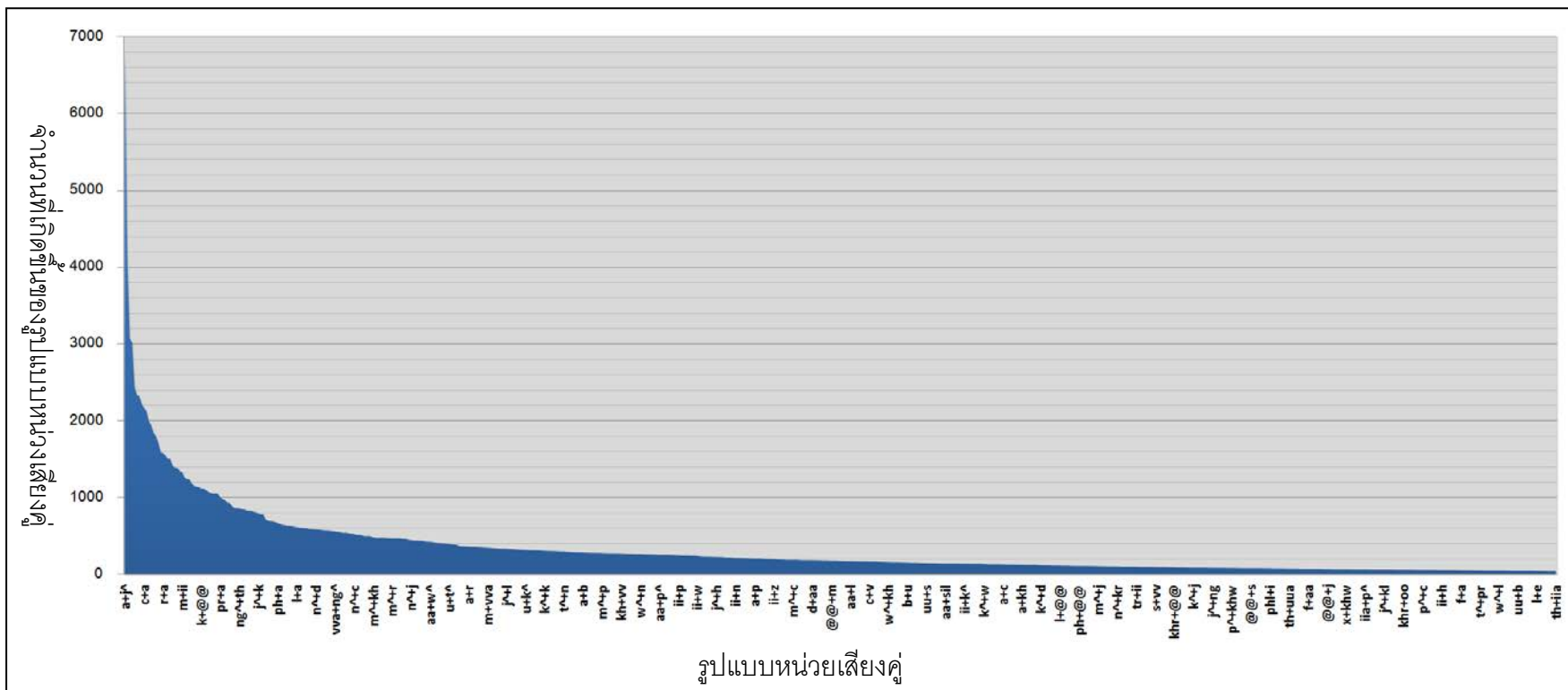
จากฐานข้อมูลบทความภาษาไทย ORCHID (Open Linguistic Resource Channelled toward InterDisciplinary research) [20] และบทความจากฐานข้อมูลอื่น ๆ ประกอบด้วย 3,007 ประโยค ระบบจะคำนวณหาค่าการกระจายตัวของหน่วยเสียง โดยการเริ่มจากนำข้อความทั้งหมดเหล่านั้นไปแปลงรูปเขียนให้เป็นคำอ่านไทยโดยใช้ PGLR (Probabilistic generalized LR parser) [112] สิ่งที่ได้คือหน่วยเสียงทั้งหมดจำนวน 246,199 หน่วยเสียง จากนั้นจะทำการหาจำนวนการเกิดของหน่วยเสียงคู่ทั้งหมด โดยใช้เครื่องมือ Hidden Markov toolkit (HTK) [31] ได้หน่วยเสียงคู่ที่เกิดขึ้นทั้งหมดคือ 1,383 รูปแบบ ใช้เวลาคำนวณเพียง 1.4 วินาทีความถี่ของจำนวนการเกิดขึ้นบางส่วนของแต่ละหน่วยเสียง ดูได้จากภาพที่ 4.1 การกระจายตัวนี้สามารถนำไปใช้เป็นรูปแบบการกระจายตัวของหน่วยเสียงเป้าหมายได้ทันที

2. การสร้างประโยคออนไลน์

ขั้นตอนนี้จะเริ่มทำการดึงข้อมูลข้อความจากหน้าเว็บบนอินเทอร์เน็ต โดยกลุ่มข้อมูลข้อความที่ได้สืบค้นมาจาก เว็บไซต์สารานุกรมไทยออนไลน์ [36] และจากคลังข้อมูลภาษาไทยออนไลน์ [37] ซึ่งประกอบไปด้วยข้อมูลจากบทความวิชาการ หนังสือแบบเรียน หนังสือพิมพ์ และนิตยสาร หาดขอบเขตของประโยคภาษาไทยโดยนับจากจำนวนพยางค์ด้วยเครื่องมือ Collocation and Thai Word Segmentation [32] โดยที่ในแต่ละประโยคต้องประกอบด้วยจำนวนพยางค์ที่มากกว่า 20 พยางค์ แต่น้อยกว่า 40 พยางค์ จากนั้นประโยคจะถูกนำไปคัดกรอง โดย ในที่นี้ตัวอักษรที่ยอมให้อยู่ในประโยคคือ ตัวอักษรภาษาไทยเท่านั้น จากนั้นประโยคผ่านการคัดกรองจะถูกนำไปแปลงแปลงรูปเขียนให้เป็นคำอ่านไทยโดยใช้ PGLR (Probabilistic generalized LR parser) [12] จากนั้นระบบจะทำการดึงข้อมูลจากอินเทอร์เน็ตไปเรื่อย ๆ อย่างต่อเนื่องรอการสั่งหยุดในกรณีที่ระบบทำการจัดทำคลังข้อมูลที่เสร็จสมบูรณ์แล้ว ประโยคดังกล่าวจะถูกนำไปเก็บไว้ใน ฐานข้อมูลเพื่อเตรียมนำไปคัดเลือกประโยคต่อไป

3. การเลือกประโยคและจัดเก็บคลังข้อความ

ขนาดของฐานข้อมูลคลังข้อความที่กำหนดไว้คือคือ 1,000 ประโยค ข้อความในแต่ละประโยคต้องไม่ซ้ำกัน จากนั้นนำประโยคจากฐานข้อมูลที่เก็บมาจากอินเทอร์เน็ตถูกนำเข้ามายังคลังข้อความปัจจุบัน อีกครั้งละ 100 ประโยค จากนั้นจะทำการให้คะแนนของประโยค จากการทดลองนี้การให้คะแนนแต่ละประโยคใช้เวลาเฉลี่ย 0.38 วินาทีต่อประโยค (CPU: Intel Core i3 processor 2.13 GHz, Memor: RAM DDR3 4 GB) ขั้นตอนนี้ใช้เวลาไปทั้งสิ้น 11 ชั่วโมง 37 นาที



ภาพที่ 4.1 แผนภาพการกระจายตัวทางสถิติทางหน่วยเสียงของฐานข้อมูล LOTUS [18]

ผลของการเปรียบเทียบคลังข้อความที่ถูกสร้างขึ้นกับคลังข้อความเป้าหมาย จากการเลือกดึงข้อมูลจากอินเทอร์เน็ตไปทั้งสิ้น 100 รอบได้มา 10,000 ประโยค แสดงในตารางที่ 4.1 ความยาวของประโยคในคลังข้อความที่ถูกสร้างขึ้นมาได้ความยาวที่เหมาะสม ผลของการครอบคลุมทางหน่วยเสียงคู่อยู่ที่ 99.13 % จากการดึงประโยคจากอินเทอร์เน็ต พบว่ารูปแบบหน่วยเสียงที่หาไม่ได้ คิดเป็น 0.87 % ของรูปแบบหน่วยเสียงทั้งหมด ส่วนมากเกิดจากการแปลงรูปเขียนเป็นคำอ่านที่ผิดพลาด ดังตารางที่ 4.2 และปัญหาจากชื่อเฉพาะที่ไม่พบ ตารางที่ 4.3 ปัญหาด้านการแปลงรูปเขียนเป็นคำอ่านที่ผิดพลาดสามารถแก้ไขได้โดยหาใช้เครื่องมือแปลงรูปเขียนเป็นคำอ่านที่มีประสิทธิภาพมากขึ้น ส่วนชื่อเฉพาะที่หาไม่พบ คาดว่าถ้าเลือกประโยคจากอินเทอร์เน็ตเข้ามาเพิ่มอีก อาจทำให้เจอคำศัพท์นั้นจากแหล่งข้อมูลทางอินเทอร์เน็ตเพิ่มขึ้นได้

ตารางที่ 4.1 คุณลักษณะของคลังข้อความที่สร้างขึ้น

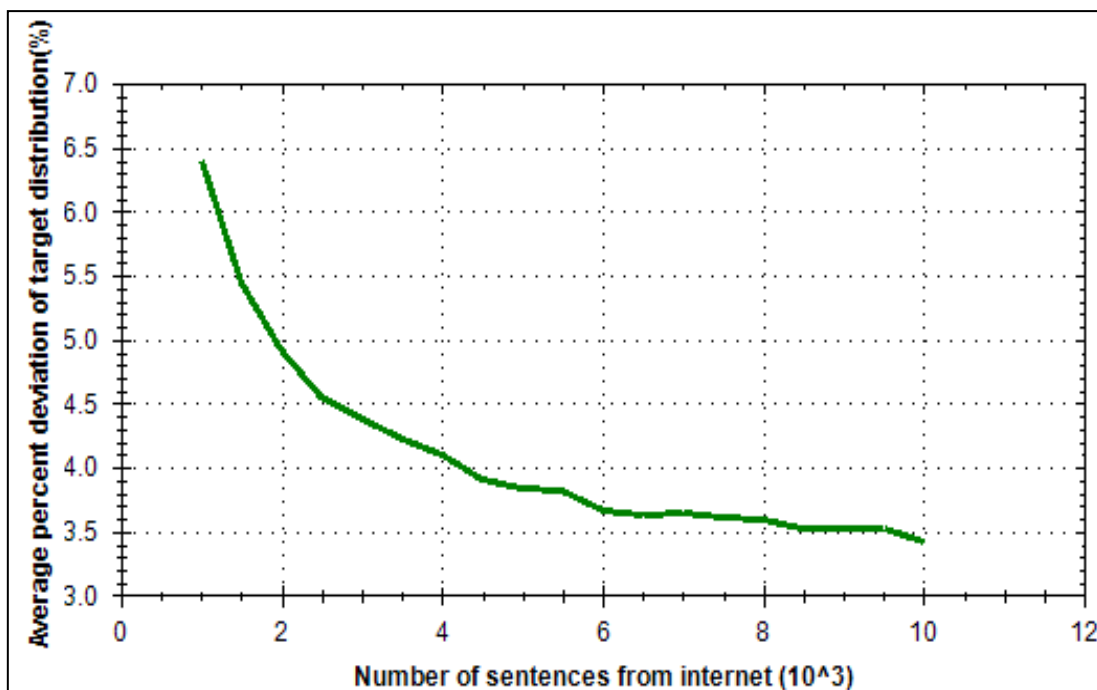
คุณลักษณะ	คลังข้อความที่สร้างขึ้น
จำนวนประโยค	1,000
จำนวนพยางค์ที่มากที่สุดในประโยค	40
จำนวนพยางค์ที่น้อยที่สุดในประโยค	20
ค่าเฉลี่ยจำนวนพยางค์ต่อประโยค	25.1
จำนวนหน่วยเสียงคู่ทั้งหมด	72,458
จำนวนรูปแบบหน่วยเสียงคู่	1,371
ความครอบคลุมทางหน่วยเสียง(%)	99.13

ตารางที่ 4.2 รูปแบบหน่วยเสียงที่ไม่ปรากฏจากการแปลงรูปเขียนเป็นรูปอ่านที่ผิด

รูปแบบที่ไม่พบ	ตัวอย่าง	แปลงรูปเสียงเป็นรูปอ่าน	
		คลังข้อความเป้าหมาย	คลังข้อความกำหนด
i+khw	ตีความ	t i khw a m ^	t ii khw a m ^
uu+br	ผู้บริโภค	ph uu br i ph oo k ^	ph uu b @ @ r i ph oo k ^
@+br	เพราะบริษัท	phr @ br i s a t ^	phr @ b @ @ r i s a t ^
ii+br	ที่บริเวณ	th ii br i w ee n ^	th ii b @ @ r i w ee n ^

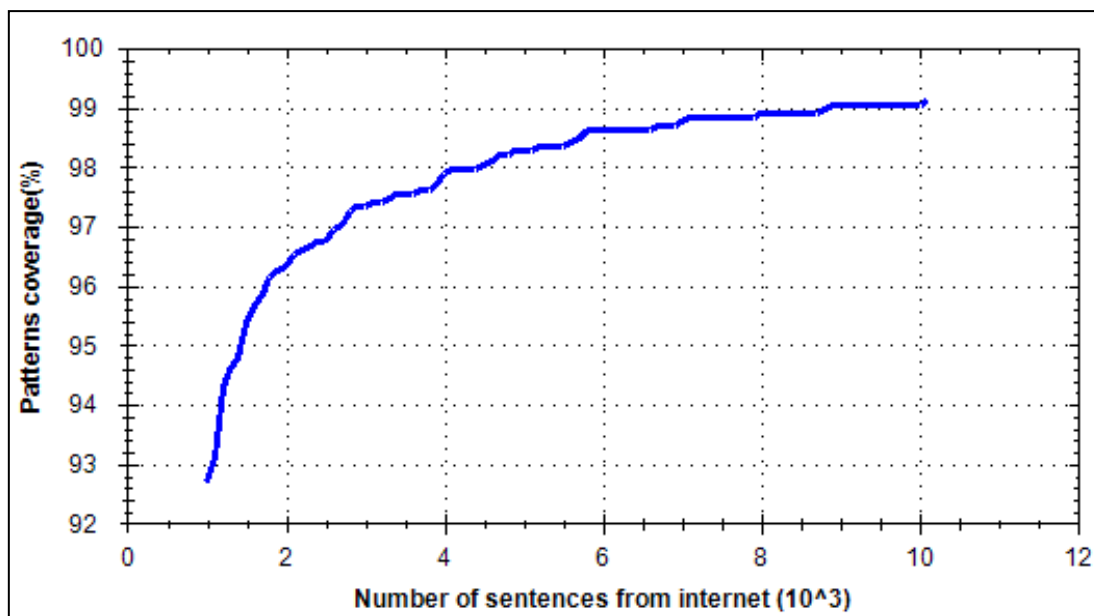
ตารางที่ 4.3 รูปแบบหน่วยเสียงที่ไม่ปรากฏจากชื่อเฉพาะ

รูปแบบที่ไม่พบ	ตัวอย่าง	แปลงรูปเสียงเป็นรูปอ่าน	
		คลังข้อความเป้าหมาย	คลังข้อความกำหนด
t+@	เตอะแตะ	t@ t x	ไม่พบ
q+b	เดอะบอสส์	d q b a z o t ^	ไม่พบ
dr+ia	อเล็กซานเดรีย	z a l e k ^ s a a n ^ d r i i a	ไม่พบ
q+kl	เดอะแกลเลอร์	d q k l x x l q q r i i	ไม่พบ
vva+kw	เมื่อกว่านซื่อ	m v v a k w a a n ^ s v v	ไม่พบ
ia+pr	เสียเปรียบ	s i i a p r i i a p ^	ไม่พบ
ng+@	เจ้าเงาะ	c a w ^ n g @	ไม่พบ
kw+ee	วัตตันเกว่น	w a t ^ t o n ^ k w e e n ^	ไม่พบ
kl+oo	กระดุกรูปโกลน	k r a d u u k ^ r u u p ^ k l o o n ^	ไม่พบ
bl+@	หนึ่งบล็อค	n v n g ^ b l @ k ^	ไม่พบ
uua+kr	ตัวกระตุ้น	t u u a k r a t u n ^	ไม่พบ
uua+kw	นำกลั้วกว่า	n a a k l u u a k w a a	ไม่พบ
pr+uua	แปรปรวน	p r x x p r u u a n ^	ไม่พบ
f+ia	มาเฟีย	m a a f i i a	ไม่พบ
kw+xx	แกว่ง	k w x x n g ^	ไม่พบ
h+@@	เหาะ	h @ @	ไม่พบ
ia+khw	เสียความ	s i i a k h w a a m ^	ไม่พบ



ภาพที่ 4.2 ผลของค่าเฉลี่ยผลต่างของกระจายตัวทางหน่วยเสียงเป้าหมายเทียบกับจำนวนประโยคออนไลน์ที่เข้ามาใหม่

ผลของค่าเฉลี่ยผลต่างของกระจายตัวทางหน่วยเสียงเป้าหมายจากประโยคที่ดึงมาจาก internet ดังภาพที่ 4.2 จากการทดลองผลของค่าเฉลี่ยความต่างจากเป้าหมาย จะลดลงมากในช่วงแรกและค่อยๆลดลงอย่างช้าๆในช่วงหลัง ส่วนการครอบคลุมทางหน่วยเสียง แสดงออกมดั่งภาพที่ 4.3 จะเห็นได้ว่าการครอบคลุมทางหน่วยเสียงเพิ่มขึ้นอย่างมากในช่วงแรกและค่อย ๆ เพิ่มขึ้นช้า ๆ ในช่วงหลัง เนื่องจากข้อมูลข้อความจากอินเทอร์เน็ตที่ได้มาในแต่ละช่วงการรับเข้าสู่ระบบมีผลต่อค่าเฉลี่ยผลต่างของกระจายตัวทางหน่วยเสียงเป้าหมายและค่าการครอบคลุมทางหน่วยเสียง ถ้าข้อมูลช่วงนั้นประกอบด้วยประโยคที่มีหน่วยเสียงที่ระบบต้องการมากจะทำให้ค่าของทั้งสองดีขึ้นอย่างมาก แต่ถ้าข้อมูลที่รับเข้ามาประกอบด้วยประโยคที่มีหน่วยเสียงที่ต้องการน้อยระบบจะรับประโยคใหม่เข้าแทนที่ได้น้อย จึงทำให้ค่าเฉลี่ยผลต่างของกระจายตัวทางหน่วยเสียงเป้าหมายค่อย ๆ ลดลง และค่าของการครอบคลุมทางหน่วยเสียงค่อย ๆ เพิ่มขึ้น



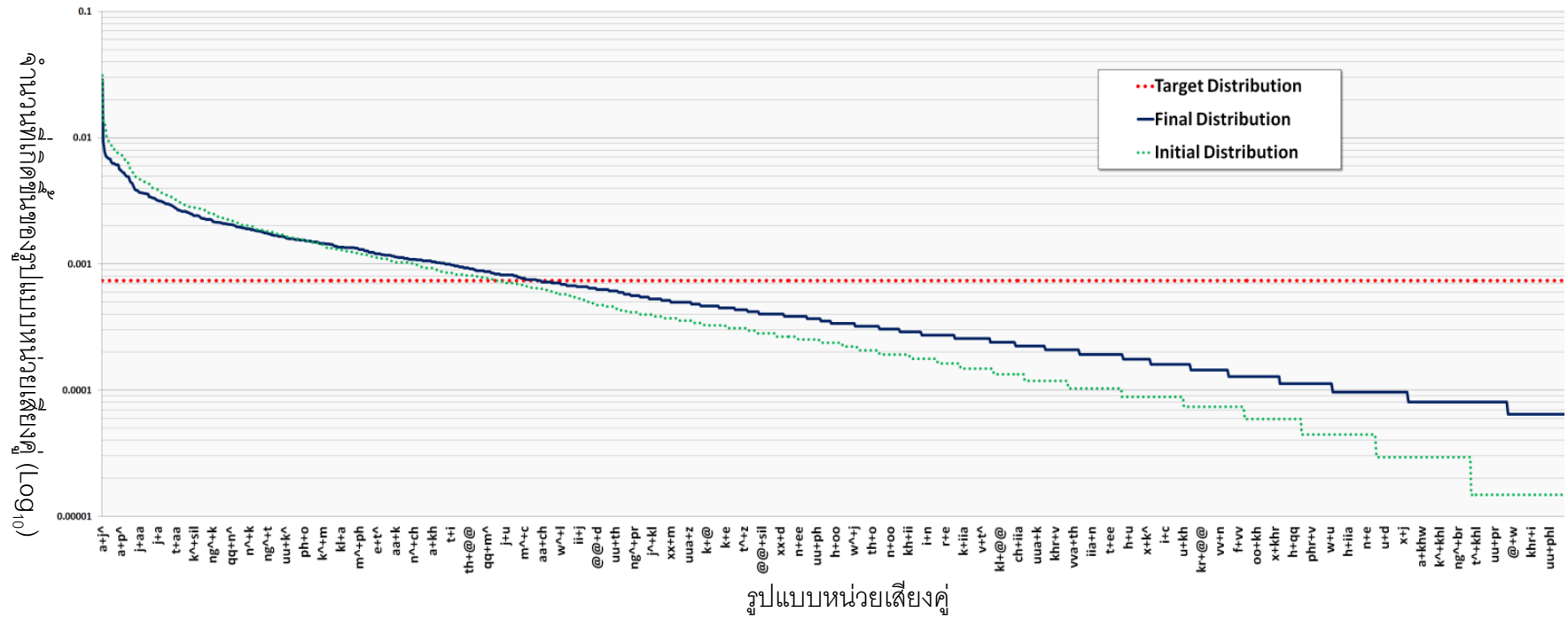
ภาพที่ 4.3 ความครอบคลุมทางหน่วยเสียงเทียบกับจำนวนประโยคออนไลน์ที่เข้ามาใหม่

เมื่อพิจารณาการกระจายตัวของหน่วยหลังจากการคัดเลือกประโยคแล้ว ดังภาพที่ 4.4 จะเห็นได้ว่าการกระจายตัวของหน่วยเสียงได้ถูกปรับ ให้ได้เป็นไปตามรูปแบบของการกระจายตัวของหน่วยเป้าหมายตามอัตราส่วนของจำนวนขนาดหน่วยเสียงรวมระบบ ได้ดีขึ้น เมื่อเทียบค่าจากการกระจายตัวตั้งต้น และระบบยังสามารถรักษาจำนวนของหน่วยเสียงขั้นต่ำให้อยู่จำนวนที่ต้องการได้อีกด้วย จึงเป็นข้อพิสูจน์ได้ว่าการเลือกประโยคจากฐานข้อมูลออนไลน์ ซึ่งเป็นคลังข้อความที่ไม่จำกัด สามารถสร้างคลังข้อความที่มีการกระจายตัวของหน่วยเสียงตามที่ต้องการได้อย่างมีประสิทธิภาพ

การสร้างคลังข้อความโดยใช้การกำหนดรูปแบบการกระจายตัวทางหน่วยเสียงเอง

ในการทดลองนี้ต้องการที่จะเปลี่ยนมาใช้ในการกระจายตัวทางหน่วยเสียงแบบกำหนดเอง โดยที่ใช้รูปแบบหน่วยเสียงคู่ 1,383 รูปแบบที่ได้มาจากการสร้างคลังข้อความจากรูปแบบการกระจายตัวทางหน่วยเสียงจากการกระจายตัวเป้าหมายในการทดลองแรก มากำหนดการกระจายตัวเป้าหมายโดยให้โอกาสที่ทุกรูปแบบของหน่วยเสียงจะเกิดขึ้นได้เท่ากันเพื่อที่จะได้การกระจายตัวเป้าหมายที่ต้องการ สามารถดูรูปแบบการกระจายตัวเป้าหมายได้จากภาพที่ 4.5 การกำหนดขนาดคลังข้อความไว้ที่ 1,000 ประโยค แต่ละประโยคประกอบด้วยจำนวนพยางค์มากกว่า 20 พยางค์แต่ไม่เกิน 40 พยางค์ ชนิดของหน่วยเสียงใช้รูปแบบ หน่วยเสียงคู่ ความครอบคลุมทางหน่วยเสียง 100% เทียบจากจำนวนรูปแบบหน่วยเสียงที่เกิดขึ้นจากรูปแบบหน่วยเสียงของเป้าหมาย แต่ละหน่วยเสียง ต้องเกิดมากกว่า 1 ครั้ง ค่าคงที่ความสำคัญรูปแบบหน่วยเสียงหายาก (α) ไว้ที่ 1.5 ค่าเบี่ยงเบนเฉลี่ยร้อยละจากเป้าหมายต้องน้อยกว่า 3.0 % ของรูปแบบของการกระจายตัวเป้าหมาย การทดลอง ใช้ลำดับขั้นตอนการทดลองเหมือนการทดลองที่แรกทั้งหมด

จากการทดลองกำหนดการกระจายตัวเป้าหมายโดยให้โอกาสที่ทุกรูปแบบของหน่วยเสียงจะเกิดขึ้นได้เท่ากันเพื่อที่จะได้การกระจายตัวเป้าหมายที่ต้องการนั้น พบว่าเมื่อพิจารณาการกระจายตัวทางหน่วยเสียงจากประโยคออนไลน์ในคลังข้อความที่ถูกสร้างเริ่มต้น การกระจายตัวนั้นเป็นตามลักษณะรูปแบบคล้ายกราฟของฟังก์ชันลอการิทึม ดังภาพที่ 4.5 และเมื่อทำการเลือกประโยคจากอินเทอร์เน็ตไปเรื่อย ๆ พบว่าการกระจายตัวของรูปแบบหน่วยเสียงนั้น ใกล้เคียงการกระจายตัวทางหน่วยเสียงเป้าหมายได้ไม่มากเท่าที่ควรทั้งนี้เพราะว่าลักษณะธรรมชาติของการกระจายตัวทางหน่วยเสียงภาษาไทยนั้น ไม่ได้เป็นแบบการกระจายตัวที่เรากำหนดไว้ เพราะฉะนั้นในการทดลองนี้อาจต้องใช้เวลา มากกว่าจะได้รูปแบบการกระจายตัวในคลังข้อความตามที่ต้องการ เพราะฉะนั้นการกำหนดลักษณะการกระจายตัวจึงมีผลต่อระยะเวลาในการสร้างคลังข้อความ ถ้าลักษณะการกระจายตัวขัดแย้งกับธรรมชาติของการกระจายตัวของภาษามากเกินไป จะทำให้การสร้างคลังข้อความเป็นไปได้ยากมากขึ้น

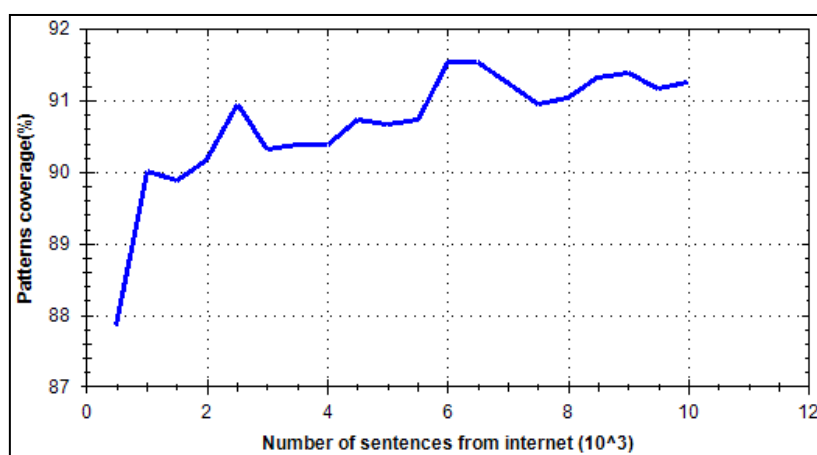


ภาพที่ 4.5 แผนภาพการกระจายตัวของหน่วยหลังจากการคัดเลือกประโยคใช้การกระจายตัวหน่วยเสียงเป้าหมายกำหนดเอง

ทดสอบผลของการสร้างคลังข้อความจากการปรับพารามิเตอร์ค่าคงที่ความสำคัญ รูปแบบหน่วยเสียงหายาก (α)

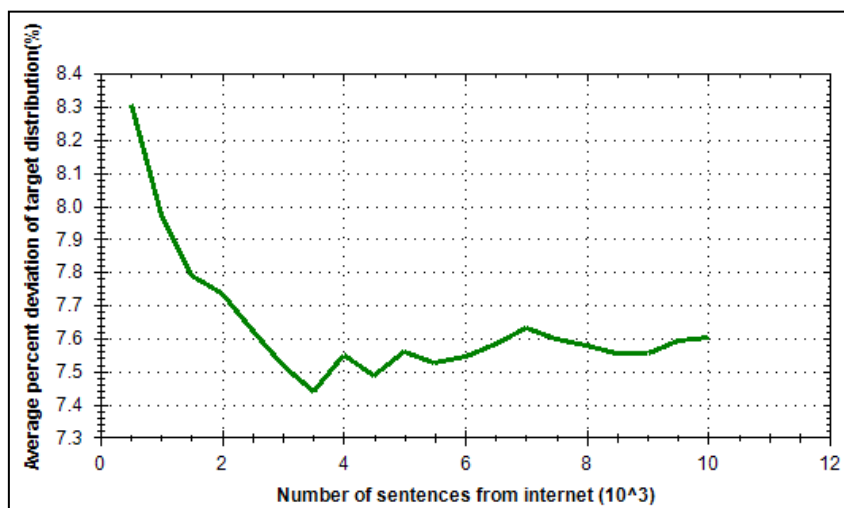
ในการทดลองนี้ต้องการวัดผลกระทบของการเปลี่ยนแปลงในการปรับพารามิเตอร์ค่าคงที่ความสำคัญรูปแบบหน่วยเสียงหายาก ทดลองโดยการสร้างคลังข้อความ 1,000 ประโยค แต่ละประโยคประกอบด้วยจำนวนพยางค์มากกว่า 20 พยางค์แต่ไม่เกิน 40 พยางค์ ชนิดของหน่วยเสียงใช้รูปแบบ หน่วยเสียงคู่ ความครอบคลุมทางหน่วยเสียง 100% เทียบจากจำนวนรูปแบบหน่วยเสียงที่เกิดขึ้นจากรูปแบบหน่วยเสียงของเป้าหมาย แต่ละหน่วยเสียง ต้องเกิดมากกว่า 1 ครั้ง ค่าเบี่ยงเบนเฉลี่ยร้อยละจากเป้าหมายต้องน้อยกว่า 3.0 % ของรูปแบบของการกระจายตัวเป้าหมาย การทดลอง ใช้ลำดับขั้นตอนการทดลองเหมือนการทดลองที่แรกทั้งหมด แต่เปลี่ยนแปลงค่าคงที่ความสำคัญรูปแบบหน่วยเสียงหายาก (α) โดยการทดลองครั้งแรกกำหนดไว้ที่ 0.0 การทดลองครั้งที่สองกำหนดไว้ที่ 1.5 และการทดลองครั้งที่สามกำหนดไว้ที่ 6.0

ผลการทดลองครั้งที่หนึ่งใช้ค่าคงที่ความสำคัญรูปแบบหน่วยเสียงหายาก (α) ที่ 0.0 เป็นการลดความสำคัญของหน่วยเสียงหายากว่า ความครอบคลุมทางหน่วยเสียงเข้าสู่เป้าหมายได้ช้าและเป็นไปได้ยากที่ผลลัพธ์ของความครอบคลุมทางหน่วยเสียงจะเข้าสู่เป้าหมายได้ในระยะเวลาสั้น ๆ สามารถดูการเปลี่ยนแปลงความครอบคลุมทางหน่วยเสียงได้จากภาพที่ 4.6



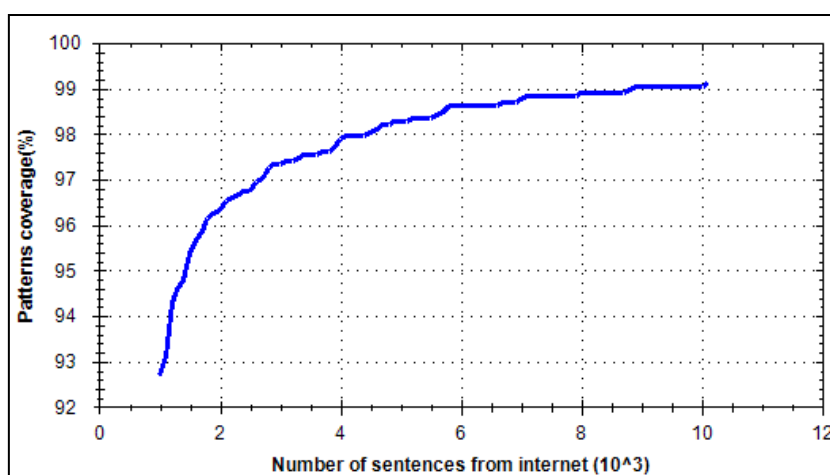
ภาพที่ 4.6 ความครอบคลุมทางหน่วยเสียงเมื่อกำหนดให้ค่า α เท่ากับ 0.0

เมื่อพิจารณาค่าเบี่ยงเบนเฉลี่ยร้อยละจากการกระจายตัวทางหน่วยเสียงเป้าหมาย พบว่าสามารถเข้าสู่การกระจายตัวทางหน่วยเสียงที่ต้องการได้ดีในระยะเวลาสั้น ๆ สามารถดูการเปลี่ยนแปลงได้ ดังภาพที่ 4.7



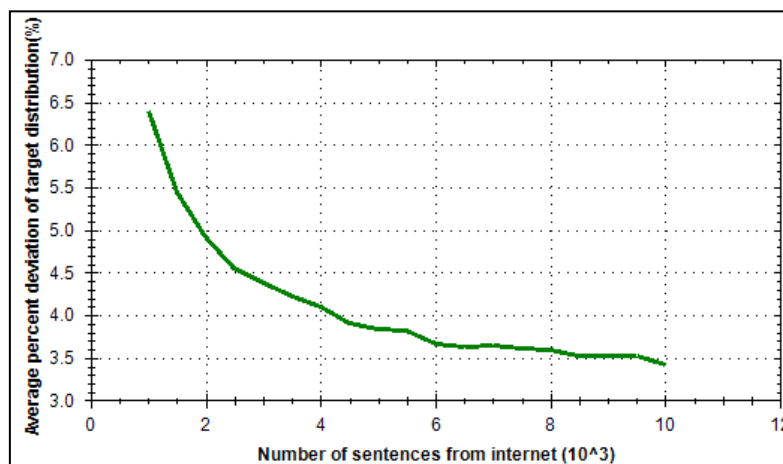
ภาพที่ 4.7 ผลของค่าเฉลี่ยผลต่างของการกระจายตัวทางหน่วยเสียงเป้าหมายเมื่อกำหนดให้ค่า α เท่ากับ 0.0

ผลการทดลองครั้งที่สองใช้ค่าคงที่ความสำคัญรูปแบบหน่วยเสียงหายาก (α) ที่ 1.5 เป็นค่าที่ใช้ในการทดลองหลักของวิทยานิพนธ์นี้ พบว่าความครอบคลุมทางหน่วยเสียงเข้าสู่เป้าหมายได้เร็วกว่าการกำหนดค่าคงที่ความสำคัญรูปแบบหน่วยเสียงหายาก (α) ที่ 0.0 และสามารถเข้าถึงความครอบคลุมทางหน่วยเสียงเป้าหมายได้ แสดงดังภาพที่ 4.8



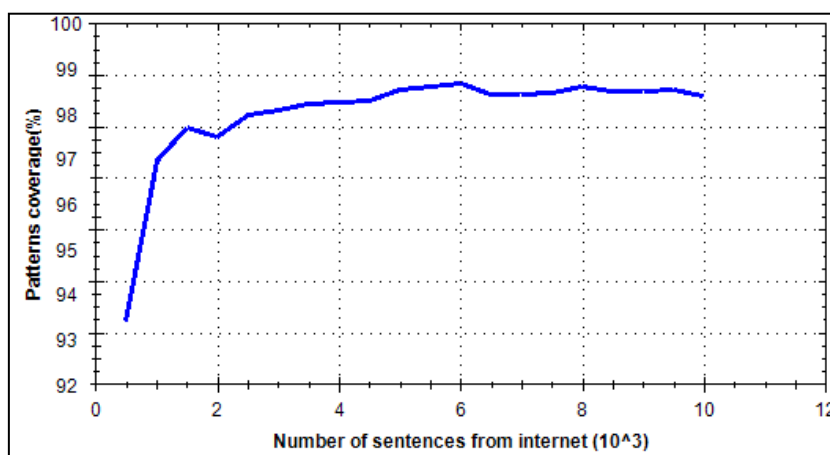
ภาพที่ 4.8 ความครอบคลุมทางหน่วยเสียงเมื่อกำหนดให้ค่า α เท่ากับ 1.5

ผลของความเปลี่ยนแปลงค่าเบี่ยงเบนเฉลี่ยร้อยละจากการกระจายตัวทางหน่วยเสียงเป้าหมาย แสดงดังภาพที่ 4.9 พบว่าสามารถเข้าสู่การกระจายตัวทางหน่วยเสียงที่ต้องการได้ และได้รูปแบบการกระจายตัวทางหน่วยเสียงที่โตกว่าการใช้ค่า ค่า α เท่ากับ 0.0



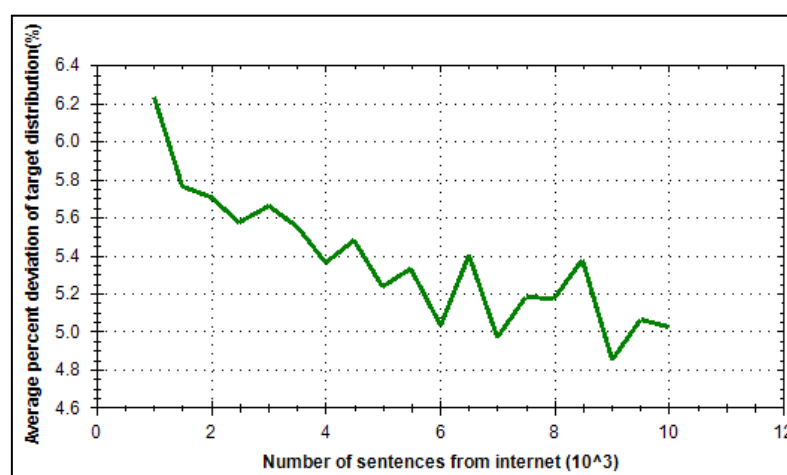
ภาพที่ 4.9 ผลของค่าเฉลี่ยผลต่างของการกระจายตัวทางหน่วยเสียงเป้าหมายเมื่อกำหนดให้ค่า α เท่ากับ 1.5

ผลการทดลองครั้งที่สามกำหนดค่าคงที่ความสำคัญรูปแบบหน่วยเสียงหายาก (α) ที่ 3.0 พบว่าความครอบคลุมทางหน่วยเสียงเข้าสู่เป้าหมายได้เร็วมากที่สุด เทียบกับกว่าการกำหนดค่าคงที่ความสำคัญรูปแบบหน่วยเสียงหายาก (α) ที่ 0.0 และ 1.5 และสามารถเข้าถึงความครอบคลุมทางหน่วยเสียงเป้าหมายได้ แสดงดังภาพที่ 4.10



ภาพที่ 4.10 ความครอบคลุมทางหน่วยเสียงเมื่อกำหนดให้ค่า α เท่ากับ 3.0

แต่เมื่อพิจารณาผลของความเปลี่ยนแปลงค่าเบี่ยงเบนเฉลี่ยร้อยละจากการกระจายตัวทางหน่วยเสียงเป้าหมาย แสดงดังภาพที่ 4.11 พบว่าความสามารถในการเข้าสู่การกระจายตัวทางหน่วยเสียงที่ต้องการขาดความต่อเนื่องและขึ้น ๆ ลง ๆ อย่างคาดการณ์ไม่ได้ ทั้งนี้เนื่องจากการกำหนดค่าคงที่ความสำคัญรูปแบบหน่วยเสียงหายาก (α) ที่มากเกินไปทำให้การคัดเลือกประโยคมุ่งเน้นการให้คะแนนไปที่ประโยคที่มีหน่วยเสียงหายากมากเกินไป



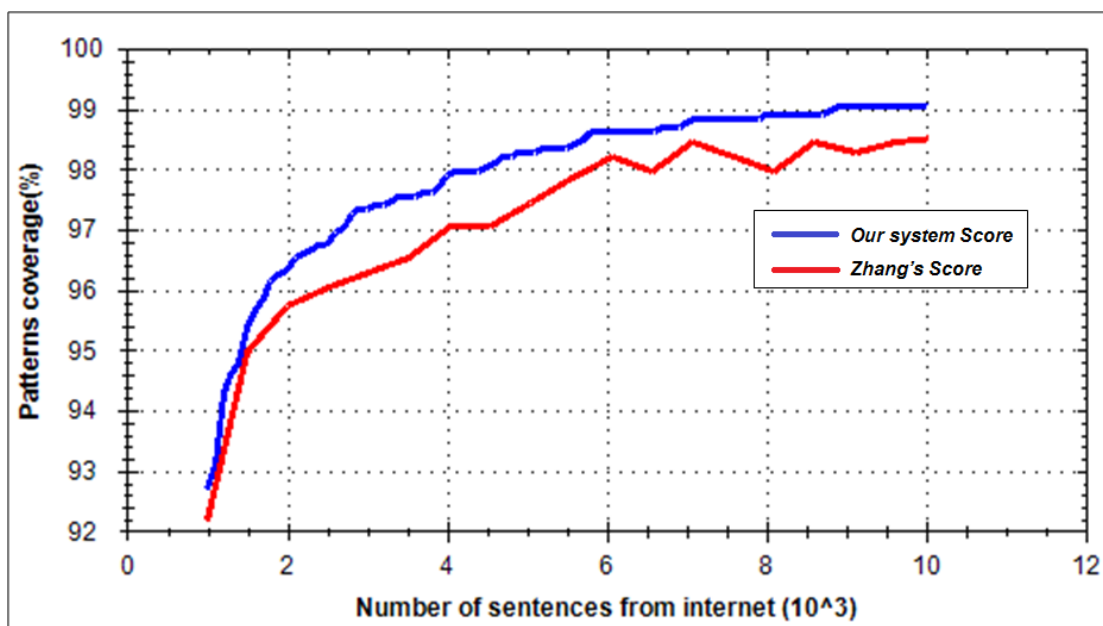
ภาพที่ 4.11 ผลของค่าเฉลี่ยผลต่างของการกระจายตัวทางหน่วยเสียงเป้าหมายเมื่อกำหนดให้ค่า α เท่ากับ 3.0

ดังนั้นในการสร้างคลังข้อความแต่ละครั้งควรกำหนดค่าคงที่ความสำคัญรูปแบบหน่วยเสียงหายาก (α) ให้พอดี ถ้าเรากำหนดน้อยเกินไปจะทำให้ความครอบคลุมทางหน่วยเสียงเข้าสู่เป้าหมายได้ช้า แต่ถ้าเรากำหนดค่าคงที่ความสำคัญรูปแบบหน่วยเสียงหายาก (α) มากเกินไปจะทำให้กระจายตัวทางหน่วยเสียงในคลังข้อความที่สร้างเข้าสู่การกระจายตัวทางหน่วยเสียงเป้าหมายได้ช้าเกินไปอีกด้วย

เปรียบเทียบผลจากวิธีการให้คะแนนในวิทยานิพนธ์กับการให้คะแนนของงานวิจัย Automatic Construction for a TTS Corpus with Limited Text (W. Zhang, 2010)

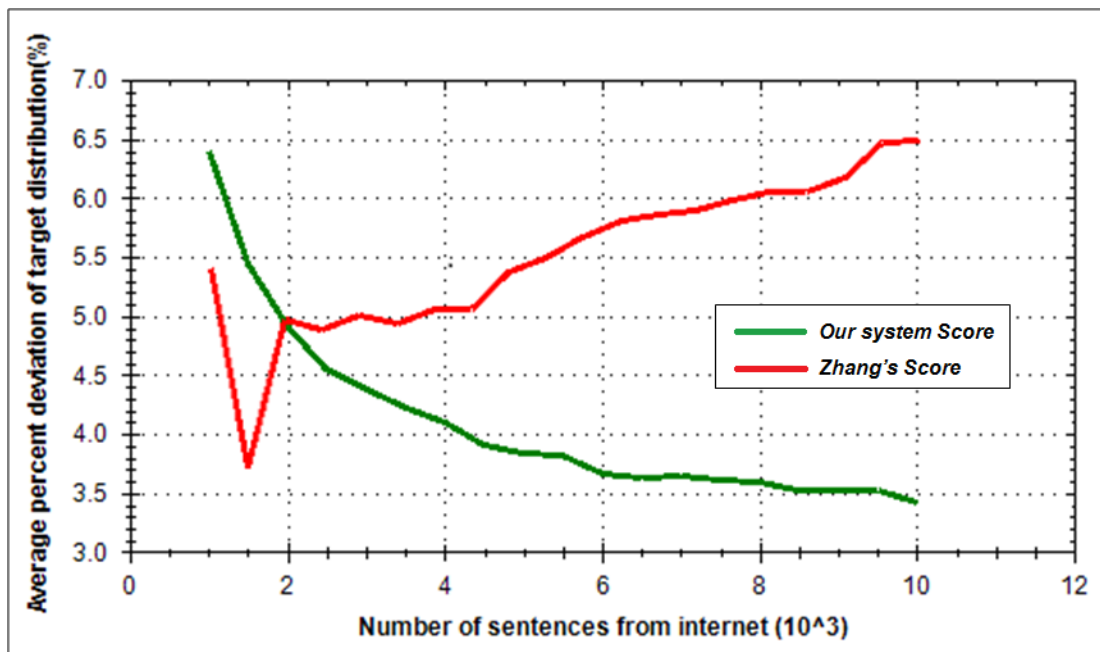
ในการทดลองนี้เป็นการเปรียบเทียบผลลัพธ์ของการเลือกประโยคโดยเทียบกับงานวิจัย Automatic Construction for a TTS Corpus with Limited Text (W. Zhang, 2010) [11] ซึ่งเป็นงานวิจัยที่ใช้วิธีการเลือกประโยคคล้ายกลับวิธีการเลือกประโยคของวิทยานิพนธ์นี้ ต่างกันที่การให้คะแนนประโยคโดยวิธีของ Zhang [11] ให้คะแนนประโยคโดยการมุ่งเน้นไปที่

หน่วยเสียงหายากเพียงอย่างเดียว การทดลองนี้ทำโดยการกำหนดการสร้าง คลังข้อความขนาด 1,000 ประโยค แต่ละประโยคประกอบด้วยจำนวนพยางค์มากกว่า 20 พยางค์แต่ไม่เกิน 40 พยางค์ ชนิดของหน่วยเสียงใช้รูปแบบ หน่วยเสียงคู่ ครอบคลุมทางหน่วยเสียง 100% เทียบจากจำนวนรูปแบบหน่วยเสียงที่เกิดขึ้นจากรูปแบบหน่วยเสียงของเป้าหมาย แต่ละหน่วยเสียงต้องเกิดมากกว่า 1 ครั้ง ค่าคงที่ความสำคัญรูปแบบหน่วยเสียงหายาก (α) ไว้ที่ 1.5 ค่าเบี่ยงเบนเฉลี่ยร้อยละจากเป้าหมายต้องน้อยกว่า 3.0 % ของรูปแบบของการกระจายตัวเป้าหมาย ผลการทดลองออกมาดังภาพที่ 4.12 และ 4.13



ภาพที่ 4.12 ครอบคลุมทางหน่วยเสียงเมื่อเปรียบเทียบวิธีการเลือกประโยคของวิทยานิพนธ์กับวิธีการเลือกประโยคแบบเก่า

ผลครอบคลุมทางหน่วยเสียงเมื่อเปรียบเทียบวิธีการเลือกประโยคของวิทยานิพนธ์กับวิธีการเลือกประโยคแบบเก่าแสดงดังภาพที่ 4.12 จะเห็นว่าวิธีการให้คะแนนประโยคของ Zhang [11] สามารถทำให้ครอบคลุมทางหน่วยเสียงเป็นไปตามที่ต้องการได้เนื่องจากใช้ระบบการให้คะแนนไปที่หน่วยเสียงหายากจึงทำให้ครอบคลุมทางหน่วยเสียงได้ออกมาตามเป้าหมายได้



ภาพที่ 4.13 ผลของค่าเฉลี่ยผลต่างของกระจายตัวทางหน่วยเสียงเป้าหมายเมื่อเปรียบเทียบวิธีการเลือกประโยคของวิทยานิพนธ์กับวิธีการเลือกประโยคแบบเก่า

แต่เมื่อพิจารณาผลของความเปลี่ยนแปลงค่าเบี่ยงเบนเฉลี่ยร้อยละจากการกระจายตัวทางหน่วยเสียงเป้าหมาย แสดงดังภาพที่ 4.13 พบว่าวิธีของ Zhang [11] ไม่สามารถเข้าถึงการกระจายตัวทางหน่วยเสียงตามที่ต้องการได้ ดูจากการเปรียบเทียบผลของความเปลี่ยนแปลงค่าเบี่ยงเบนเฉลี่ยร้อยละจากการกระจายตัวทางหน่วยเสียงเป้าหมายที่มีค่ามากขึ้นเรื่อย ๆ เมื่อมีประโยคใหม่เข้ามาทั้งนี้เนื่องจากว่าการให้คะแนนประโยคของ Zhang [11] มุ่งสู่ความครอบคลุมทางหน่วยเสียงเท่านั้น จึงไม่สามารถเข้าถึงการกระจายตัวตามหน่วยเสียงเป้าหมายได้ ส่วนวิธีการให้คะแนนของวิทยานิพนธ์นี้สามารถเข้าถึงการกระจายตัวทางหน่วยเสียงตามเป้าหมายได้ เนื่องจากวิธีการให้คะแนนประโยค ที่นอกจากจะให้คะแนนประโยคที่ประกอบด้วยหน่วยเสียงหายากแล้วยังให้คะแนนประโยคโดยการพิจารณาว่าประโยคนั้นจะทำให้การกระจายตัวทางสถิติเข้าสู่เป้าหมายได้หรือไม่อีกด้วย สามารถดูลักษณะการกระจายตัวทางหน่วยเสียงของวิธีการให้คะแนนประโยคของวิทยานิพนธ์นี้เปรียบเทียบกับวิธีการให้คะแนนประโยค Zhang [11] ได้ แสดงดังภาพที่ 4.14

บทที่ 5

บทสรุปผลการวิจัย และข้อเสนอแนะ

สรุปผลการวิจัย

บทความนี้ได้เสนอวิธีการการสร้างคลังข้อความอัตโนมัติ จากการกระจายตัวของหน่วยที่กำหนด โดยใช้การเลือกประโยคจากอินเทอร์เน็ตเพื่อให้ได้คลังข้อความตามที่ต้องการ จากการทดลองจึงสรุปได้ว่าการเลือกประโยคจากอินเทอร์เน็ตมาสร้างคลังข้อความสามารถสร้างคลังข้อความตามการกระจายตัวของหน่วยที่กำหนดได้อย่างมีประสิทธิภาพ และได้ประโยคที่ทันสมัยในการทดสอบระบบในเรื่องการกำหนดค่าพารามิเตอร์ต่าง ๆ ซึ่งผู้ใช้สามารถกำหนดได้เอง สามารถสรุปผลการปรับค่าพารามิเตอร์ได้ดังตารางที่ 5.1

ตารางที่ 5.1 ผลของการเปลี่ยนแปลงค่าพารามิเตอร์

พารามิเตอร์	หน่วย	ผลของการปรับค่าพารามิเตอร์	
		เพิ่มค่าพารามิเตอร์	ลดค่าพารามิเตอร์
รูปแบบของหน่วยเสียง	หน่วยเสียง/ รูปแบบ	- พบรูปแบบหน่วยเสียงได้ยาก - ใช้เวลาในการสร้างคลังข้อความมากขึ้น - ใช้กับงานสังเคราะห์เสียงได้ดีขึ้น เสียงเป็นธรรมชาติมากขึ้น	- พบรูปแบบหน่วยเสียงได้ง่ายกว่า - ใช้เวลาในการสร้างคลังข้อความน้อยลง - น้อยเกินไปจนทำให้ไม่เหมาะกับงานสังเคราะห์เสียง เสียงขาดความธรรมชาติ
จำนวนขั้นต่ำในแต่ละรูปแบบหน่วยเสียง	หน่วยเสียง	- กำหนดมากเกินไป - สะสมรูปแบบหน่วยเสียงที่หายากได้ไม่พอกับความครอบคลุมทางหน่วยเสียงที่ต้องการ - ใช้เวลาในการสร้างคลังข้อความมากขึ้น	- ใช้เวลาหารูปแบบหน่วยเสียงที่หาบ่อยกว่า เข้าถึงความครอบคลุมทางหน่วยเสียงได้ดีกว่า - ใช้เวลาในการสร้างคลังข้อความน้อยลง

พารามิเตอร์	หน่วย	ผลของการปรับค่าพารามิเตอร์	
		เพิ่มค่าพารามิเตอร์	ลดค่าพารามิเตอร์
ขนาดของจำนวน ประโยชน์ในคลังข้อความ	ประโยชน์	<ul style="list-style-type: none"> - ความครอบคลุมทางหน่วยเสียงเพิ่มขึ้นได้เร็วขึ้น - ใช้เวลาในการสร้างคลังข้อความมากขึ้น - การกระจายตัวทางหน่วยเสียงเป็นไปตามเป้าหมายได้ดีกว่า 	<ul style="list-style-type: none"> - ความครอบคลุมทางหน่วยเสียงเข้าถึงเป้าหมายได้ช้าลง - ใช้เวลาในการสร้างคลังข้อความน้อยลง - การกระจายตัวทางหน่วยเสียงเป็นไปตามเป้าหมายได้น้อยลง
ขอบเขตขั้นต่ำจำนวนพยางค์ในแต่ละประโยค	พยางค์/ ประโยค	<ul style="list-style-type: none"> - จำนวนพยางค์ที่มากเกินไปอาจทำให้ การให้คะแนนประโยคผิดพลาด 	<ul style="list-style-type: none"> - จำนวนพยางค์ที่น้อยเกินไปอาจทำให้ประโยคไม่ได้ใจความ
ค่าคงที่ความสำคัญ รูปแบบหน่วยเสียงหายาก(α)	-	<ul style="list-style-type: none"> - เข้าสู่เป้าหมายของความครอบคลุมทางหน่วยเสียงได้เร็ว - รักษาประโยคที่มีหน่วยเสียงที่หายากไว้ได้ - การกระจายตัวทางหน่วยเสียงเข้าสู่เป้าหมายได้ยากขึ้น 	<ul style="list-style-type: none"> - เข้าสู่เป้าหมายของความครอบคลุมทางหน่วยเสียงได้ช้ากว่า - รักษาประโยคที่มีหน่วยเสียงที่หายากไว้ไม่ได้ - ต้องไปหาใหม่ - การกระจายตัวทางหน่วยเสียงเข้าสู่เป้าหมายได้ดีขึ้น
ความครอบคลุมทางหน่วยเสียง	เปอร์เซ็นต์	<ul style="list-style-type: none"> - คลังข้อความมีความครอบคลุมทางหน่วยเสียงที่ดี เหมาะสำหรับการใช้งานในทุก ๆ ด้าน - ใช้เวลาในการสร้างคลังข้อความมากขึ้น 	<ul style="list-style-type: none"> - คลังข้อความมีความครอบคลุมทางหน่วยเสียงที่ดี เหมาะสำหรับการใช้งานในทุก ๆ ด้าน - ใช้เวลาในการสร้างคลังข้อความน้อยลง

พารามิเตอร์	หน่วย	ผลของการปรับค่าพารามิเตอร์	
		เพิ่มค่าพารามิเตอร์	ลดค่าพารามิเตอร์
ค่าเบี่ยงเบนเฉลี่ยร้อยละ จากเป้าหมาย	เปอร์เซ็นต์	- เข้าสู่การกระจายตัวทาง หน่วยเสียงเป้าหมายได้ น้อยลง - ใช้เวลาในการสร้างคลัง ข้อความน้อยลง	- เข้าสู่การกระจายตัวทาง หน่วยเสียงเป้าหมาย ได้มากขึ้น - ใช้เวลาในการสร้างคลัง ข้อความมากขึ้น

การกำหนดพารามิเตอร์ต่าง ๆ ดังที่กล่าวมาเป็นสิ่งสำคัญที่จะทำให้คลังข้อความออกมาตามเป้าหมายที่เราต้องการ ในบางงานอาจต้องการสร้างคลังข้อความที่มีการกระจายตัวทางสถิติและความครอบคลุมทางหน่วยเสียงในระดับดีมาก แต่ในบางงานอาจต้องการสร้างคลังข้อความในระยะเวลาที่กำหนดเนื่องจากเวลาน้อย ในที่นี้ผู้ใช้จึงต้องกำหนดพารามิเตอร์ให้ดี

คลังข้อความที่ได้นี้จึงเป็นคลังข้อความที่ดีสามารถคาดการณ์ประสิทธิภาพได้ หลักการเลือกข้อความของงานวิจัยนี้ ยังสามารถนำไปประยุกต์ใช้การสร้างคลังข้อความกับภาษาอื่น ๆ ได้ อีกด้วย โดยการกำหนดแหล่งข้อมูลออนไลน์ให้ตรงกับภาษาที่ต้องการ กำหนดตัวกรองประโยคใหม่ และเลือกใช้เครื่องมือตัวแปลงรูปเขียนเป็นรูปอ่านของภาษาให้ตรงกับภาษาที่กำหนด

ข้อเสนอแนะ

ในขั้นตอนการคัดเลือกประโยคของวิทยานิพนธ์นี้ เมื่อพิจารณาประโยคที่ถูกคัดออกในการคัดเลือกส่วนใหญ่แล้วเป็นประโยคที่น่าจะนำไปใช้ได้ แต่ประโยคเหล่านั้นประกอบไปด้วยตัวอักษรหรือสัญลักษณ์พิเศษต่าง ๆ ที่ไม่สามารถแปลงรูปเขียนเป็นรูปอ่านได้ ทำให้ระบบต้องทิ้งประโยคที่ดีเหล่านั้นออกไปเป็นจำนวนมาก จึงอยากเสนอแนะให้ปรับปรุงตัวแปลงรูปเขียนเป็นรูปอ่านที่ครอบคลุมการแปลงของตัวอักษรที่มีสัญลักษณ์พิเศษที่ป้อนอยู่ในแต่ละภาษาด้วย โดยอาจใช้เทคนิคการแยกสัญลักษณ์พิเศษประเภทที่จะทำให้การแปลงรูปเขียนเป็นรูปอ่านไม่สำเร็จออกก่อนเมื่อทำการแปลงรูปเขียนเป็นรูปอ่านส่วนประโยคหลักผ่านแล้ว ค่อยนำสัญลักษณ์เหล่านั้นมาใส่ในประโยคหลักทีหลัง

จากปัญหาการค้นหารูปแบบหน่วยเสียงบางหน่วยที่หายาก อาจเกิดมาจากชื่อเฉพาะที่เกิดขึ้นในคลังข้อความเป้าหมาย ทำให้การเข้าถึงความครอบคลุมทางหน่วยเสียงเป็นไปได้ช้า จึงอยากเสนอแนะให้ใช้วิธีการนำคำศัพท์ที่เกิดจากรูปแบบหน่วยเสียงที่มาจากชื่อเฉพาะเหล่านั้น ใช้

เป็นคำสำคัญในการค้นหาประโยคจากอินเทอร์เน็ต จะทำให้การได้มาของหน่วยเสียงที่หายากหรือหน่วยเสียงที่มาจากชื่อเฉพาะเป็นไปอย่างรวดเร็ว ทำให้สามารถเข้าถึงความครอบคลุมทางหน่วยเสียงเร็วขึ้น

ปัญหาที่คาดว่าจะเกิดขึ้นจากความต้องการสร้างคลังข้อความในการใช้งานจริงคือ การสร้างคลังข้อความที่ต้องการนั้น อาจจะไม่สามารถหาการกระจายตัวทางหน่วยเสียงของเป้าหมายที่เหมาะสมได้ งานวิจัยในอนาคตของเรามุ่งไปที่การสร้างการกระจายตัวเป้าหมายที่เหมาะสมอัตโนมัติ โดยพิจารณารูปแบบและจำนวนหน่วยเสียงตามโอกาสที่เกิดขึ้นในแต่ละภาษา

มีแนวคิดเรื่องการทดสอบคุณภาพของคลังข้อความที่ถูกสร้างขึ้นเพิ่มเติม โดยนำข้อความจากคลังข้อความที่สร้างขึ้นจากกระบวนการของวิทยานิพนธ์นี้ นำไปอ่านบันทึกเสียงโดยอาสาสมัครเพื่อสร้างคลังเสียงจริง และนำไปสร้างระบบรู้จำเสียงพูดอัตโนมัติหรือระบบสังเคราะห์เสียงพูด เพื่อทดสอบคุณภาพ และประเมินผลของคลังข้อความ

รายการอ้างอิง

- [1] Wijaya, T., and Wijaya, D. Limited Speech Recognition for Controlling Movement of Mobile Robot Implemented on ATmega162 Microcontroller. International Conference on Computer and Automation Engineering (2009) : 47–350.
- [2] NV Access Inc. NonVisual Desktop Access (NVDA) [Online]. 2011. Available from: <http://www.nvda-project.org/> [2011,July 6]
- [3] NECTEC. TVIS Traffic Talk [Online]. 2011. Available from: <http://www.tvis.nectec.or.th/> [2011, August 8]
- [4] Stan, A., Yamagishi, J., King, S., and Aylett, M. The Romanian Speech Synthesis (RSS) corpus: building a high quality HMM-based speech synthesis system using a high sampling rate. Speech Communication 53,3 (2011): 442-450.
- [5] Wei, Z., Ranran, D., Minhui, P., and Qihong, W. Automatic Speech Corpus Construction from Broadcasting Speech Databases. Computational Intelligence and Security (CIS), 2010 International Conference (2010) : 639–643.
- [6] Hansakunbuntheung, C., Rugchatjaroen, A., and Wutiwiwatchai, C. An intensive design of a Thai speech corpus. International Symposium on Natural Language Processing (SNLP) (2007) : 27-132.
- [7] Cormen, Leiserson, and Rivest. Introduction to Algorithms. Chapter 16 Greedy Algorithms (1990) : 329-360.
- [8] Mandal, S., Das, B., Mitra, P., and Basu, A. Developing Bengali Speech Corpus for Phone Recognizer Using Optimum Text Selection Technique. Asian Language Processing (IALP), 2011 International Conference (2011) : 268 –271.
- [9] Chulalongkorn University. ChulaDaisy: Daisy full-text audio book creator [Online]. 2010. Available from:<http://www.chuladaisy.eng.chula.ac.th> [2011,July 1]

- [10] Kominek, J., and Black, A.W. CMU ARCTIC database for speech synthesis, Technical Report CMU LTI-03-177 (2003).
- [11] Zhang W., Liu Y., Deng Y., Pang M., Automatic Construction for a TTS Corpus with Limited Text. International Conference on Measuring Technology and Mechatronics Automation (2010) : 707-710.
- [12] Tarsaku, P., Sornlertlamvanich, V., and Thongpresirt, R., Thai Grapheme-to-Phoneme Using Probabilistic GLR Parser. Eurospeech 2001 2 (2001) : 1057-1060.
- [13] Kasuriya, S., Sornlertlamvanich, V., Cotsomrong, P., Kanokphara, S., Thatphithakkul, N. Thai Speech Corpus for Speech Recognition. Proc. Of Oriental-COCOSDA 2003 (2003) : 54-61.
- [14] Wutiwiwatchai, C., Saychum, S., Rugchatijaroen, A. An Intensive Design of a Thai Speech Synthesis Corpus. Proceeding of the SNLP 2007 (2007).
- [15] Mittapiyanurak, P., Hansakunbuntheung, C. and Teprasit, V. and Sornlertlamvanich, V. Issues in Thai Text-to-Speech Synthesis: The NECTEC Approach. Proceedings of NECTEC Annual Conference (2000) : 483-495.
- [16] Wiboon, T. The Engineering Training Report at NECTEC. Engineering Faculty, Khasetsart University (1990).
- [17] Deza, E., Deza, M. Dictionary of Distances, Elsevier, ISBN 0444520872 (2006)
- [18] Vlasta V., Petr V. Methods of Sentences Selection for Read-Speech Corpus Design TSD'99, LNAI 1692 (1999) : 165-170.
- [19] Hansakunbuntheung, C., Tesprasit V., and Sornlertlamvanich, V. Thai Tagged Speech Corpus for Speech Synthesis. the Oriental COCOSDA 2003 (2003) : 97-104.
- [20] Sornlertlamvanich, V., Charoenporn, T., and Isahara, H., ORCHID: Thai Part-of Speech Tagged Corpus. National Electronics and Computer Technology Center Technical Report (1997) : 5-19.

- [21] Sornlertlamvanich, V. Probabilistic Language Modeling for Generalized LR Parsing. Technical Report, Department of Computer Science, Tokyo Institute of Technology (September 1998).
- [22] Itou, K., Yamamoto, M., Takeda, K., Takezawa, T., Matsuka, T., Kobayashi, T., Shikano, K., and Itahashi, S. JNAS: Japanese Speech Corpus for Large Vocabulary Continuous Speech Recognition Research. In Journal of Acoustic Society of Japan 20,3 (1999).
- [23] Fransen, J., Pye, D., Robinson, T., Woodland, P., and Young, S. WSJCAM0 Corpus and Recording Description, Cambridge University (1994).
- [24] Luksaneeyanawin, S. A Thai Text to Speech System. Proceeding of the Conference of the Regional Workshops on Computer Processing of Asian Languages (1990) : 305-315.
- [25] Ni, J.F., Hirai, T., and Kawai, H. Constructing a Phonetic-Rich Speech Corpus while Controlling Time-Dependent Voice Quality Variability for English Speech Synthesis. Proc.ICASSPU (2006) : 881-884.
- [26] Wutiw WATCHAI, C., and Furui, S. Thai speech processing technology: a review. Speech Communication 49,1 (2007) : 8-27.
- [27] Mittrapiy-nuruk, P., Hansakunbuntheung, C., Tesprasit, V., and Sornlertlamvanich, V. Issues in Thai-to-Speech Synthesis The NECTECT Approach. NECTEC Technical Journal 2,7 (2000) : 36-47.
- [28] Chitturi, R., Sebsibe H.M, Kumar, R., Elissa, K. Rapid methods for optimal text selection, [Online]. 2011. Available from: <http://www.rohitkumar.net/publications/RANLP2005.pdf> [2011,July 6]
- [29] Ling, S., Yu, H., Ren-hua, W. Corpus design for the Chinese speech synthesis system. In ICSLP-2000 2 (2000) : 391-394.
- [30] Microsoft Public License (Ms-PL). HTML parser Html Agility Pack, [Online]. 2011. Available from: <http://htmlagilitypack.codeplex.com/> [2011,July 6]

- [31] Young, S., Evermann, G., Galse, M., Kershaw, D., Moore, G. Hidden Markov model toolkit – speech recognition toolkit. [Online]. 2011. Available from: <http://htk.eng.cam.ac.uk/> [2011,July 12]
- [32] Aroonmanakun, W. Collocation and Thai Word Segmentation, Proc. SNLP and Oriental COCOSDA Workshop (2002) : 68-75.
- [33] Oracle Corporation, MySQL Connector Net 6.2.4, [Online]. 2011. Available from: <http://dev.mysql.com/downloads/connector/net/> [2011,July 12]
- [34] Oracle Corporation, MySQL Community Server 5.1, [Online]. 2011. Available from: <http://dev.mysql.com/downloads/mysql/5.1.html/> [2011,Jan 12]
- [35] Microsoft Corporation, Microsoft Visual Studio 2010, [Online]. 2011. Available from: <http://www.microsoft.com/visualstudio/th-th/download> [2011,Jan 12]
- [36] Royal Command of H.M. the King, Thai junior encyclopedia project, [Online]. 2011. Available from: <http://kanchanapisek.or.th/kp6/New/index.php> [2011,Feb 13]
- [37] Department of Linguistics Chulalongkorn University. TNC: Thai national corpus. [Online]. 2011. Available from: <http://ling.arts.chula.ac.th/tnc2> [2011,Jan 17]
- [38] Mohammad, A., Abushariahl, M., Raja, N.A., Zainuddinl, R., Moustafa, E. Phonetically rich and balance speech corpus for Arabic speaker-independent continuous automatic speech recognition system. 10th International Conference on Information Science, Signal Processing and their Applications (2010) : 65-68.

ภาคผนวก

ภาคผนวก ก

ตารางเทียบสัทอักษรสากลกับหน่วยเสียงไทยที่ใช้ในวิทยานิพนธ์

Thai Phoneme

พยัญชนะต้น			
เดี่ยว	ตัวอย่าง	ผสม	ตัวอย่าง
p	ปาก	pr	ประสาน
t	เต็น, กฏี	ph r	พหราน
c	ละ	tr	เตรียม
k	ก่อน	kr	กราบ
z	อาน	kh r	คร้า
ph	พบ, ภัย, ย่าน	pl	ปลา
th	ทิ้ง, ชุง, เต้า, ลาน, มณโฑ	phl	พลาด
ch	ชอบ, เจอ	thr	จันทรา
kh	คน, เขิน, ฆ่า	kl	เกลอ
b	บอก	khl	เคลื่อน
d	ด้าน, ชญา	kw	กวาง
m	ไม่	kh w	ขวา
n	นาน, เณร	เสียงทับศัพท์	
ng	เงิน	br	เบรน
l	เลน, กีฬา	bl	บล
r	รอ, ภัย	fr	ฟราย
f	ฝน, ฟัน	fl	เฟลม
s	สาย, ศิลา, รักษา, ซ่อน	dr	ดราคอน
h	โหน, เฮฮา	17 หน่วย	
w	วา		
j	ย่อน, หญิง		
21 หน่วย			

สระ			
เดี่ยว	ตัวอย่าง	ผสม	ตัวอย่าง
a	อะ	ia	เอียะ
aa	อา	iaa	เอีย
i	อี	va	เอือะ
ii	เอี	vva	เอือ
v	อื	ua	อัวะ
vv	อือ	uua	อิว
u	อุ	6 หน่วย	
uu	อู		
e	เอะ		
ee	เอ		
x	แอะ		
xx	แอ		
o	โอะ		
oo	โอ		
@	เออะ		
@@	ออ		
q	เออะ		
qq	เออ		
18 หน่วย			

ตัวสะกด	
เดี่ยว	ตัวอย่าง
p^	พบ
t^	เทริด
k^	ปาก
g^	หาร
m^	ลม
ng^	ฟาง
j^	ยาย
w^	กาวิ
เสียงทับศัพท์	
f^	กราฟ
l^	แอล
s^	เอส
ch^	คลัช
12 หน่วย	

THE INTERNATIONAL PHONETIC ALPHABET (2005)

CONSONANTS (PULMONIC)

	LABIAL		CORONAL				DORSAL			RADICAL		LARYNGEAL
	Bilabial	Labio-dental	Dental	Alveolar	Palato-alveolar	Retroflex	Palatal	Velar	Uvular	Pharyngeal	Epi-glottal	Glottal
Nasal	m	ɱ	n				ɳ	ɲ	ŋ	ɴ		
Plosive	p b	ɸ β	t d				ʈ ɖ	c ɟ	k ɡ	q ɢ	ʔ	ʕ
Fricative	ɸ β	f v	θ ð	s z	ʃ ʒ	ʂ ʐ	ç ʝ	x ɣ	χ ʁ	ħ ʕ	ħ ʕ	h ɦ
Approximant		ʋ	ɹ				ɻ	j	ɰ			
Trill	ʙ		r							ʀ		ʀ
Tap, Flap		ⱱ	ɾ				ɽ					
Lateral fricative			ɬ ɮ			ɮ	ɬ	ɮ				
Lateral approximant			l			ɭ	ʎ	ʎ				
Lateral flap			ɺ			ɻ						

Where symbols appear in pairs, the one to the right represents a modally voiced consonant, except for murmured *ɦ*. Shaded areas denote articulations judged to be impossible. Light grey letters are unofficial extensions of the IPA.

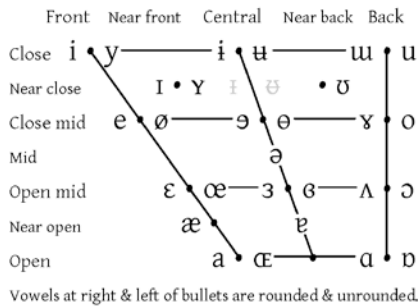
CONSONANTS (NON-PULMONIC)

Anterior click releases (require posterior stops)	Voiced implosives	Ejectives
ʘ Bilabial fricated	ɓ Bilabial	ʼ <i>Examples:</i>
ǀ Laminar alveolar fricated ("dental")	ɗ Dental or alveolar	ɓ' Bilabial
ǃ Apical (post)alveolar abrupt ("retroflex")	ɟ Palatal	ɗ' Dental or alveolar
ǂ Laminar postalveolar abrupt ("palatal")	ɠ Velar	ɠ' Velar
ǁ Lateral alveolar fricated ("lateral")	ʄ Uvular	ʄ' Alveolar fricative

CONSONANTS (CO-ARTICULATED)

- ɱ Voiceless labialized velar approximant
- ʋ Voiced labialized velar approximant
- ɥ Voiced labialized palatal approximant
- ç Voiceless palatalized postalveolar (alveolo-palatal) fricative
- ʝ Voiced palatalized postalveolar (alveolo-palatal) fricative
- ɧ Simultaneous x and ʃ (disputed)
- kp ts Affricates and double articulations may be joined by a tie bar

VOWELS



SUPRASEGMENTALS

- ˈ Primary stress
- ˌ Secondary stress [ˌfoʊnəˈtʃən]
- eː Long
- eˑ Half-long
- e Short
- ˑ Extra-short
- ˑ Syllable break
- ˑ Linking (no break)
- INTONATION
- ˊ Minor (foot) break
- ˋ Major (intonation) break
- ↗ Global rise
- ↘ Global fall

TO NE

- ˊ Level tones
- ˋ Contour-tone examples:
- ˊ Top
- ˋ Rising
- ˊ High
- ˋ Falling
- ˊ Mid
- ˋ High rising
- ˊ Low
- ˋ Low rising
- ˊ Bottom
- ˋ High falling
- ˊ Tone terracing
- ˋ Low falling
- ˊ Upstep
- ˋ Peaking
- ˋ Downstep
- ˋ Dipping

DIACRITICS

Diacritics may be placed above a symbol with a descender, as ɲ̥. Other IPA symbols may appear as diacritics to represent phonetic detail: ɾ̥ (fricative release), ɓ̥ (breathy voice), ʔ̥ (glottal onset), ə̥ (epenthetic schwa), o̥ (diphthongization).

SYLLABICITY & RELEASES	PHONATION	PRIMARY ARTICULATION	SECONDARY ARTICULATION
ɲ̥ ɲ̥	Syllabic	ɲ̥ ɲ̥	Voiceless or Slack voice
ɓ̥ ɓ̥	Non-syllabic	ɓ̥ ɓ̥	Modal voice or Stiff voice
ɰ̥ ɰ̥	(Pre)aspirated	ɰ̥ ɰ̥	Breathy voice
ɲ̥̚	Nasal release	ɲ̥̚ ɲ̥̚	Creaky voice
ɲ̥̚	Lateral release	ɲ̥̚ ɲ̥̚	Strident
ɲ̥̚	No audible release	ɲ̥̚ ɲ̥̚	Linguolabial
ɸ̥ β̥	Lowered (β̥ is a bilabial approximant)	ɸ̥ ɸ̥	Raised (ɻ̥ is a voiced alveolar non-sibilant fricative, ʀ̥ a fricative trill)

ภาคผนวก ข
ตัวอย่างประโยคจากคลังข้อความที่สร้างขึ้น 300 อันดับแรก

ลำดับ	ประโยคในฐานข้อมูล
1	เธอรอให้เพลงจบ จะเล่าเรื่องเด็กหนุ่มนักเปียโนให้สามีฟัง แต่ประพจน์เหมือนไม่ได้ยินอะไรทั้งสิ้น
2	เจอแล้ว ผมร้องออกมาแล้วก็เปิดประตูลัดเข้าไป ลืมแม่กระทั่งว่าจะต้องเคาะประตูก่อน
3	ทั้งที่จัดของเกือบเสร็จเรียบร้อย ของฝากญาติโกโหติกา ผองเพื่อนก็เตรียมไว้เสร็จสรรพ
4	อ้าว คุณปู่ตอบเกมหัวเราะ นี่แหละบทกวีของกวีเอกชาว อิตาลี
5	กานบัวสะบัดปาก ปล่อยเหยื่อปลา นกกวักกู่ฟาลันบึงน้ำ กระจาสาสายคล้ายเย็นสัน
6	ถ้าอีตเลอริไม่หันปลายปีกไปทางมอสโก เยอรมันก็คงไม่แพ้สงคราม
7	ล้างจาน ปิดไฟนอน คินนันทังคู่วรกกัณ้อย่างดุเดือด เปญคว่าหน้า
8	พื้จิวทะเลาะกับพื้พาย พื้สะใ้มีเวรคุณพ่อ แต่ผลอหลับจนน้ำเกลือหมดขวด
9	ปรอดปราดหลีกเสียดปีกบัต สาลิกาถูรังก่อนล่องลัด เขาชวาคคล้ายหวัดหวาดขันร้อง
10	กลับกลายเป็นตรงกันข้าม ต่อมาทั้งสองได้รับความช่วยเหลือจากเบดูอินกลุ่มหนึ่ง
11	แบบว่าซัดซีโคลในเวลาเดียวกัน แต่ก็หายเร็วกว่าความนุ่มนวลโดยการทาทามีนแบบของยาย
12	เหตุการณ์ได้เข้มข้นยิ่งขึ้นเมื่อทั้งสองถูกพวกโจรเบดูอินปล้นและทำร้ายบาดเจ็บ
13	ทั้งนี้ ในความควบคุมของเจ้าหน้าที่ผู้ฝึกหัดหรือผู้ให้การฝึกอบรมซึ่งเป็นผู้ประกอบวิชาชีพการพยาบาล
14	รวมทั้งน้ำเสียงอบอุ่นที่ปลุกปลอบเด็กหญิงพิการนำสงสารคนนั้นอยู่ในหัวใจเธอแล้ว
15	นางถือกระเช้าหวายบรรจุอาหารเข้าประเภทแซนด์วิช และขนมปังแครกเกอร์
16	เคยเห็นเต่าทะเลไซ้ใหม่ครับ ถ้าเทียบกับเต่าบก เต่าแม่น้ำเต่าทะเลกระดองจะกลมกว่า
17	สงสารมิ่งนัก กูตัดใจ กูยกมือนี้ให้มิ่ง นี่คือนพ้อมมือที่มีนิ้วครบ
18	ความล้มเหลวทางการเมืองที่ประเทศไม่สามารถผลักดันข้อประวัติศาสตร์ทางการเมืองไปสู่ยุคปฏิรูปได้
19	หากผู้ผลิตรายการเป็นหน่วยงานรัฐ ก็น่าจะผลิตรายการที่ให้บริการแก่คนทุกชั้น
20	เดือนเธอให้ระลึกถึงวันที่เขาเอ่ยปากขอเธอแต่งงาน แต่เพราะวันทวยวิษัหมั้วแต่คิดแผนหลอกเธออยู่นั้น
21	เพราะยังไม่กลับตัวกัน แค่ความซื่อสัตย์สุจริตก็เป็นคุณภาพที่ขาดตกบกพร่องกันเสียแล้ว

ลำดับ	ประโยคในฐานข้อมูล
22	หนึ่งชั่วโมงครึ่งฉันคือไฟ ข้าเปลวไฟในตะวันกัณฑ์จักรหนึ่งแกนหลักจักรวาลซับซ้อนไซ
23	แกแหละจำเขาไม่ได้ ผมโกรธฟูบขึ้นมาทันทีที่ได้ยินประโยคนี
24	เขาบอกว่าดีใจที่ได้เจอฉันซะที แม่เล่าเรื่องของฉันให้ฟังเยอะมาก
25	ถึงดวงหทัยที่สดใสและอยู่ใกล้ฉัน สังคมมีทั้งความกลมกลืนและความแตกต่าง
26	ทีเซอร์วังฉันไว้แน่นไม่ยอมปล่อย เธอมากับฉันก็ต้องกลับไปกับฉันสิ
27	ซึ่งเพลโตได้นำชื่อพีธีร์มาเป็นชื่อเรื่องหนังสือของเขา กล่าวถึงผู้ไม่ตกอยู่ในห้วงรัก
28	ว่าเกิดอะไรขึ้นกับหญิงสาวที่เป็นความหวังของมิสซิสคลีฟ คุณเห็นวิกตอเรียใช่ไหม
29	ผมเข้าห้องน้ำซึ่งตอนนั้นว่างคน คราวนี้ผมต้องเงี่ยหูฟังเป็นพิเศษเพราะกลัวจะไม่ได้ยินเสียงวิทยุจากห้องข้าง
30	แต่เมืองไทยเป็นอย่างนี้ ข่างอุดมสมบูรณ์ด้วยความเป็นธรรมชาติที่นำกลางแควงและชวนให้คิด
31	เชิญฟังเพลงบรรเลงหวานประสานเสียง เสนาะสำเนียงได้ต้นไทรพุ่มใบหนา เมื่อแสงเรือแจอชมพู่สู่ท้องฟ้า
32	ดวงตาคมกริบผลอทอดนิ่งมายังภาพสาวตรงหน้า ก่อให้ความเงียบโรยตัวโอบล้อมทั้งห้อง
33	พ่อ ฉันมีลาภ เขาชูขึ้นอวด อ้วนปีเลย นี่ไง มันเซ่อ
34	เสียง ฉันล่องท่องไพรกว้างฟังเสียงค่างบางซะนี่ เสียงนกพลอดจู้จู้จุดนตรีในป่าดอย
35	โลกนับสิบจะออกมาคล้ายคลึงกันเป็นโลกแฝดซึ่งสุดท้ายจะถูกกลืนรวมเป็นหนึ่งเดียว
36	พญาเมงฮายสร้างเสียงใหม่ขึ้นมาเพื่อเป็นศูนย์กลางล้านนาเนื้อลูกเนื้อ ก็ดังที่พ่อได้บอกเจ้าแล้วว่าพญาเมงฮายสร้างเสียงใหม่เมื่อ
37	น้องหมวยได้ยินเข้าก็หันไปทางอื่นทันที เพราะเพราะเรื่องนี้เป็นแผลในใจมาโดยตลอด
38	กระบี่เกาะมากสุดในทะเลไทย กระบี่มีเกาะมากกว่าสุราษฎร์ ผสมผสานทั้งหมู่เกาะหินปูนและหมู่เกาะหินแกรนิต
39	สัตยาปรากฏแม้เกลี้ยกลบ สัตว์อย่ามา หมายถึงวองคัพพ สัตว์อย่าหมาย
40	บนฝั่งเกาะเกาะถุน เลี้ยงลูกสาวคนเดียวจนเติบโตเป็นสาว เงินหยวนทุกเหรียญ
41	เพลโตใช้วิธีการนี้เพราะคิดว่าเป็นวิธีเดียวที่จะทำให้ค้นพบความจริง
42	หกโมงเย็นตรงเป๊ะ คุณสามาเคาะประตูเรียกให้ไปกินข้าว พอโผล่หน้าเข้าไปคุณยายนั่งที่โต๊ะกินข้าวเรียบร้อยแล้ว

ลำดับ	ประโยคในฐานข้อมูล
43	หากเมื่อบทสนทนาดำเนินไปได้ไม่นาน มันก็เวียนวนกลับไปอยู่ในรูปเดิม
44	นานและนานหลายนานผมเป็นผีนคั่นเนื่องจากเหงื่อและกระแฉีก ความหมักหมม แม่กลับจากสามชุก
45	รอยยิ้มกว้างปรากฏบนใบหน้าที่อ่อนระโหยจากการร่ำศึกสอบ เมื่อนึกถึง
46	พริมโรสสีเหลืองสีชมพูแปด บิสตอต้าที่ดูเหมือนไมโครโฟนของมิสพิงคี้
47	ข้ามีวิถีข้าเยี่ยงภูผาเกินพลัดเซ กว้างเทือกเกินถ้ำเทเกินเหลิงเล่ห์หลงลมหลง
48	และไม่รู้ตัวเลยว่าอีกไม่กี่วินาทีคนส่วนใหญ่จะสิ้นชีวิตเพราะตีกระเบิด
49	พีติดณณ์ไม่ชอบให้มีคนนอกมาวุ่นวายที่นี่ คราวก่อนองศาทำเรื่องไว้ก็โกรธแทบแย่
50	เลยไปถึงบริเวณบาร์ข้างสระที่มีเก้าอี้หมุนี่นั่น ฝรั่งก็แหงนมองความงามตระหง่านของ
51	อีกฝ่ายก็สูญเสียเองและชวนเพนกวินคุยอย่างสนิทสนมเสียด้วย มันทำให้เขาเสียตายช่วงเวลานี้มาก
52	ยุคนั้นไม่มีตราเรือใบและตรามะลิ หมาสาวหมาหนุ่มหรือหมาแก่ไม่เคยถูกจูงไปอาบน้ำ
53	ในทันทีหลังจากนั้นถ้าได้กลิ่นของกินมันไปคลุกและคั่นกองขยะเน่า เป็อนขนตัวเหม็น
54	ฝรั่งหัวเราะแล้วไม่พูดพาลำทำเพลงอีกต่อไป วางสบู่เอาผ้าเช็ดตัวพาดเอาไว้
55	ฉันละก็ไม่ได้ ถึงไหนถึงกันเชียว เหม็นกลิ่นกระถางไฟคะ ยุพราตอบเรียบ
56	จิต ทรามต่ำแค่คอยรั้งความ น้ำ ฝนหล่นฟ้ามาฝาก มิตร ภาพพลันหลากเลี้ยงย่าน
57	เสียงเครื่องจักรที่ด้านสมัยและฝุ่นข้างเปลือกที่กระจายคั่งอยู่ทั่วไปในโรงสีมิได้ทำให้เข้าเกิดความย่ำแย่แต่อย่างใด
58	อู๊ดังมาเดินข้างฉัน เธอบอกฉันว่า ช่างมันเถอะ แม่นี้ก็เรื่องมากอย่างนี้แหละ
59	หากว่าสองวันก่อนเธอจะไม่ล้มเลิกแผนการกลับเมืองไทยอย่างกะทันหันเสียก่อน
60	แต่งตัว กินอาหารเข้าเมื่อได้เวลาทำงาน เขาจะเดินออกไปตรวจดูกระบวนการสีข้าวทุกขั้นตอนอย่างละเอียด
61	มินะมะตะ ปลาเล็กกินแพลงตัน ปลาขนาดกลางกินปลาเล็ก ปลาใหญ่กินปลาขนาดกลาง
62	คะ คือแบบว่าไอ้ย อายจ้ง กิรดาลอบยืมกับท่าทางของหนูน้อยตรงหน้า
63	มุขยัยจะก้าวพรวดไปข้างหน้าเพื่อดูให้แน่ใจอีกครั้ง หากเซเสียหลักเสียก่อน
64	ส่งโน่นส่งนี่มาให้ปอกให้หันยกครัวชนิดนอนสตีออป กะจ้งหวะไม่ให้ลุกอย่างเทียมโหด
65	มองเฟิร์นเพลินพิศใบดุดจลู่ไวกามน่าดู เก้งปามาเป็นคู่ช่างแสนรู้ระวังภัย

ลำดับ	ประโยคในฐานข้อมูล
66	แต่วิธีของยายนั้นละลายกำแพงกันน้ำหน่อไม่ต้องเบียดขโหลม เผลอไม่นานหันมองก็พบว่าโรคนั้นล่อนออกเกลี้ยงเกลา
67	สงครามของหมาวัดค่อนข้างโหดร้ายและโหด หมาในสภาพต่างเอาตัวรอดสั่งสมความเลวมากกว่าความดี
68	พื้นที่ที่มีแสงสว่างก็เห็นนาฬิกาทรงโบราณพร้อมทั้งเสียงที่ดังขึ้น
69	และได้สูญเสียมิตรภาพไป เขากลายเป็นคนเงียบขี้อาย ซึมเศร้านานเท่าไรหนอ
70	ตั้งชื่อลูกสาวว่า จาง หรือ ลินลี่ หลังจากคลอดลูกได้ไม่นานเสียงก็เสียชีวิตไป
71	การเลี้ยงสัตว์ก็ตกอยู่ในชะตากรรมเดียวกันกับการเพาะปลูก หลังสงคราม
72	และแนวทางของคณะกรรมการพัฒนาระบบราชการ ซึ่งครอบคลุมถึงการปรับบทบาท
73	อาจอึดปวยเป็นโรคมะเร็งหมอวินิจฉัยว่าอาการขั้นสุดท้ายจะมีชีวิตอยู่ไม่เกินหกเดือน
74	มีเมียน้อยอีกคนละซี มีเมียน้อยอีกสองคนกระมัง เขาแกลังกระซิบถามอีก
75	ด้านหลังมหาวิทยาลัย ผ่านพิพิธภัณฑศาล ดูโอโมโดยไม่แะชม เพราะเธอบอกว่าได้เข้าไปดูภายในวิหารแล้วเมื่อวาน
76	การประสานงานกับหน่วยรับตรวจและหน่วยตรวจสอบภายนอก และงานประชาสัมพันธ์งานคณะกรรมการตรวจสอบ
77	และซารีฟได้กลับไปช่วยกษัตริย์ของประเทศเขาปราบกบฏสำเร็จ ทั้งสองได้เข้าพิธีแต่งงาน
78	เพราะต่อมวัดระดับความหล่อของฉันทันร่องเตือนว่าความหล่อของเขาถึงขีดมาตรฐาน
79	แทบไม่รู้ตัวว่ามาถึงบ้านได้อย่างไร ก่อนจะก้าวขึ้นเรือน พ่อชะงักเหมือนถูกตรึงกับแผ่นดินเมื่อได้ยินบุญเพ็งเรียกเสียงดัง
80	เป็นการใส่ร้ายทั้งนั้น สาบานได้ ให้ฟ้าผ่าลิเข้า ถ้าผมเลวจริง
81	จะเห็นปรากฏการณ์นี้ได้ชัดเจน ทางด้านเหนือฝั่งฝรั่งเศสจะมีป่ามากกว่า
82	เรานอนตากลมเล่นที่ได้ต้นไม้ชายทุ่งอีกสักพักก็กลับ เพื่อให้ทันอาหารกลางวัน
83	หลายคนหมิ่นแคลนเสียแค้นคับ แต่มียอมอ่อนรับสิ่งอัปรีภัย ฟ้า เมฆายน
84	ค่านิยมของนักการเมือง ซึ่งเป็นตัวแปรสำคัญในเรื่องนี้ บางคนยอมรับว่า
85	ก็จบการศึกษาระดับปริญญาตรีทางด้านธุรกิจและการบัญชีจากมหาวิทยาลัยมีชื่อที่สุดแห่งหนึ่งของเกาะฮ่องกงด้วยวัย
86	มันตะแล้วกระที่บฉันทรายเกือบตาย สู้มันบนหาดทรายไม่ได้ พอมันผลอดอน

ลำดับ	ประโยคในฐานข้อมูล
	อาบน้ํา
87	อะริโยะฉิได้พูดถึงสารเคมีที่ใช้ในการ เพาะปลูกว่ามีมากมายหลายชนิด
88	ซึ่งมีใบอนุญาตเป็นผู้ประกอบวิชาชีพการพยาบาล การผดุงครรภ์
89	เพื่อให้การปฏิบัติงานเกี่ยวกับการพัฒนาระบบราชการบรรลุตามวัตถุประสงค์
90	ผมพลิกตัวกลับมาทันที ตามองไปที่กรงนกที่แขวนลอยอยู่ เจ้าจิวโผล่หน้าทำตาบ้องแป้ว อยู่
91	คือตัวกลางที่จะต้องเอานโยบายมาแปลงเป็นภาคปฏิบัติให้ได้ตามที่รัฐบาลเขาต้องการ
92	ผมชันสายลูกชิ้นนั่งมองผ่านหน้าต่างไปทางก้อนน้ำ ฉับพลันที่เห็นภาพ
93	ลูกหามาไม่ได้ให้คำตอบและเรียกเจ้าหญิงว่า เจ้าหญิง ชื่อเจ้าหญิงซึ่งเป็นชื่อหนังสือรวม เรื่องสั้นเล่มนี้จึงมาจากเรื่องสั้นเรื่องนี้นั่นเอง
94	รัฐบาลญี่ปุ่นเร่งที่จะผลิตเนื้อสัตว์เพื่อเป็นอาหารให้ได้จำนวนมาก ทำให้เกิดปัญหาต่าง
95	และจะพยายามกีดกันคนรักของตนออกจากปรัชญาที่เทพเจ้าประทาน วาตะชินที่
96	รอยยิ้มแกมยิ้มที่ไม่เชิงว่ามีไมตรีหรือรังเกียจก้ำกึ่งกันอยู่บนริมฝีปากแดงเข้ม
97	เจ้าหกไปเป็นตราให้กระทรวงเกษตร ส่วนเจ้าเจ็ดเป็นนายแบบแฟชั่นถ่ายรูปลงในนิตยสาร ซีดีไฟ
98	โดยเฉพาะอย่างยิ่งสามารถให้บริการกับนักเรียน นักศึกษา เด็ก เยาวชน
99	หลักสูตรการศึกษาในโรงเรียนทุกระดับยังไม่สอดคล้องกับความเป็นอยู่ของคนในชนบท
100	กับข้าวเยอะเยาะ ถ้าไม่ตีว่าเป็นปิ่นโต ก็มากินด้วยกัน เย็นวันนั้นคุณยายเจริญอาหาร เป็นพิเศษ
101	ข้ากำเนิดดินน้ำลมและไฟหนึ่งเกิดอาจดับใหม่ในพริบตา ข้าท้องฟ้าเมฆหมอกนอกหน หาวหนึ่งดวงดาวพราวแสงแห่งเวหา
102	ข้าอวดโอ้อวดหึ่งพ่ายพังเพหนึ่งเวลาก็บกัลป์อาสัณโลก ข้าเป็นหนึ่งในมวลล้วนซากศพ หนึ่งจักรภพประสพสุขทุกซีก
103	เกี่ยวข้องกับแข็งแรงละเอียดหรือแสดงหลักฐานเพิ่มเติม หรือจัดส่งเจ้าหน้าที่ไปตรวจสอบ ข้อเท็จจริงยังสำนักงาน
104	ประกอบกับลักษณะทางภูมิศาสตร์ที่เต็มไปด้วยเขาสูงและหุบลึก มีสังคมพืชเขตร้อน เกือบทุกระดับ
105	ป่า ยังยิ่งใหญ่ใกล้ตึก คือ ความเคร่งเครียดโครมครึก เมือง หลับหลังตึกยาวนาน

ลำดับ	ประโยคในฐานข้อมูล
106	ต้องไล่ตะครุบ ส่วนสหายเพียงวันนั้นจะดำเนินาจะว่าอย่างไรผ้าก็ไม่หลุด
107	เมื่อวิเคราะห์กันต่อไปถึงทางออก กลับเป็นเรื่องค่อนข้างแปลก ที่ผู้ซึ่งมองเห็นว่า
108	และกลัวว่าพ่อแม่จะทิ้งตนไป จึงไม่ยอมห่างพ่อแม่ ในเรื่องของคนไข้เด็กนี้
109	เพราะเป็นกระบวนการสื่อสารแบบทางเดียว นอกจากนั้นยังมีข้อเท็จจริงเพิ่มเติมที่ว่า
110	เนื่องด้วยมีคอกกั้นอยู่ปลายเท้า เป็นคอกอาบน้ำดังกล่าวนั้น แลแคว่ด้านซ้ายก็กั้นด้วย แผงไม้ไผ่
111	จนถึงหุบเขาเล็กแคบรูปตัววีที่ถูกร่องน้ำเซาะ เพมาโคจึงเป็นแหล่งรวมพันธุ์ไม้ที่ หลากหลายอย่างน่าอัศจรรย์
112	เกี่ยวกับการโอนนักโทษ ให้เป็นไปตามที่กำหนดในกฎกระทรวง มาตรา
113	กับข้าวของคนจีนจะต้องมีแกงจืดขามใหญ่ขามหนึ่งไว้ชดน้ำด้วยเสมอและไม่มีช้อน กลาง
114	เหตุเพราะความรู้เท่าไม่ถึงการณ์ สรรพสิ่งมีสองข้าง แม้แต่เรื่องนม นม
115	เพื่อกระชับการบรรยายให้กะทัดรัด เข้าใจง่าย และเห็นภาพว่าการบรรยายธรรมดา
116	กล่าวกันว่าในสมัยที่ประชากรส่วนใหญ่ในสังคมยังด้อยซึ่งพลังวิเคราะห์
117	อาการยังทรงอยู่ แต่หมอบอกว่าให้เตรียมใจไว้อาจไม่ถึงสัปดาห์หน้า
118	ลามใหม่ทั้งด้านเหนือและได้มุ่งไปทางทิศตะวันตกเข้าหาลำน้ำแม่ยมซึ่งมีวัดราชธานีที่ เต็มไปด้วยโบราณวัตถุวางอยู่
119	ไซ่ นั้นเสียงซึ่งอ้ายอินตา ข้าพเจ้าเจียหุฟ้ง เสียงซึ่งดังฉ่ง ฉ่งไกลเข้ามา
120	จำชื่อนี้ได้ แก่นแก้วนัก ก็เพราะความรู้ภาษาอังกฤษที่เคยร่ำเรียนมากับอาจารย์
121	หรือองค์กรของผู้สูงอายุ ที่ดำเนินการตามวัตถุประสงค์ และกิจกรรมเกี่ยวกับการ คุ้มครอง
122	แม้ท้องฟ้าเบื้องนอกจะยังคงคลุ้มด้วยม่านกำมะหยี่สีดำอยู่ หากเขาก็ไม่อาจชมตาหลับ ได้อีก
123	บำบัดแก้ไขฟื้นฟูและสงเคราะห์ผู้กระทำผิด สนับสนุนการแก้ไขปัญหาเศรษฐกิจและ สังคมของประเทศ
124	จึงได้ทราบว่าจะขอร้องยังต้องออกไปหางานทำในเมืองใหญ่เพื่อนำเงินมาซื้อเครื่องจักร และปุ๋ยสำหรับใช้ในการเพาะปลูก
125	แล้วควายมันก็ลงน่านอน โผล่แต่หัวที่มีเขาเกะกะขึ้นมาสะบัด

ลำดับ	ประโยคในฐานข้อมูล
126	โดนดอกเบี๋ยไล่จีทุกวัน บางคนต้องเปลี่ยนอาชีพมาทำงานกลางคืนเพื่อความอยู่รอด
127	โดยมิเชลล์ยังไหวหวั่น กลัวซารีฟจะมีภรรยาอีกตามที่ศาสนาอนุญาต
128	ความไม่สบายใจจากสิ่งที่เพิ่งได้ยินเข้าเกาะกุมจิตใจของเขาจนหนักอึ้ง
129	คุณปู่จึงต้องแก้แค้นโดยการจับหัวประสิทธิ์ศักดิ์อีกแล้ว น้องหมวยเห็นภาพนี้เป็นครั้งที่สอง
130	เคย์ ต่อ บรรยายภาคที่นี้คล้ายตลาดนัดในเมืองไทย แต่สินค้าหลากหลายกว่า
131	ผมแทบช็อค เพราะฝรั่งพเนจรคนนั้นกำลังเปลือยกายล่อนจ้อนอาบน้ำอยู่คนเดียว
132	นิทรรศการ ศูนย์รับบริจาคหนังสือ เป็นหน่วยงานที่ทำหน้าที่บริการรับบริจาคหนังสือ
133	ธนวรรค์ ไม่พอใจแค่การผสมผสานผู้นำหรือลดความรู้สึกต่อต้านในหมู่ผู้นำต่ออำนาจรัฐ
134	ชายหนุ่มยังไม่สามารถทำอะไรลงได้ คงต้องอาศัยเวลาอีกนานทีเดียว
135	ผมมากินก๋วยเตี๋ยว ไม่ได้มาดูโชว์ กลางคืนของสีลม ตามฟุตบาททางเดินตั้งแต่หน้ากรุงเทพคริสเตียนถึงหน้าห้างโรบินสัน
136	และกระทรวงมหาดไทย โดยมีเนื้อหารายละเอียดของการดำเนินการที่ดูจะแตกต่างกันออกไป
137	ทรงเชี่ยวชาญเรื่องดนตรีไทยและเพลงไทย จอห์นกระซิบ อัลธามองตามสายตาของเขา
138	โดยบรรยายประกอบรูปไว้ว่า พญายักษ์ทั้งเจ็ดต่างสำนึกตัวกลับเป็นงูดี
139	เรนอดรู้สึกไม่ได้ว่ามันทั้งสกปรก ทั้งขุ่นข้น ดูน่ากลัวพิลึก
140	แล้วก็ได้เห็นไฟประลัยกัลป์ผลาญเผา เข้มหมัดเท่าที่คนสุขุขทัยเคยประสบกลางตลาดใหญ่เป็นห้องแถวเรือนไม้เก่าสร้างต่อเนื่องกัน
141	คิดแต่ว่าชีวิตของตัวเองจบสิ้นตามการจากไปของคนที่รักยิ่ง หากมาวันนี้เขาสามารถปลดปล่อยความทุกข์โศกนั้นได้
142	นกร้องถี่ แล้วเราก็ไหลถึงที่โล่ง หล้าดำ มันเขียวอะไรอย่างนั้น
143	ศึกษา ค้นคว้า วิเคราะห์ วิจัย รวบรวมสถิติข้อมูลเพื่อปรับปรุงกฎหมายและระเบียบ
144	ยังชีพด้วยการเป็นพ่อค้าขายสมุนไพรจีนจากจีนแผ่นดินใหญ่แถวย่านชิมซาซุย
145	ยังไม่กำหนดไว้เด็ดขาดตายตัว คงจะต้องเตรียมอะไรบางอย่างให้พร้อมเสียก่อน
146	และงอกขนใหม่สวย การอาบและฟอกสบู่ยากกับไฉ่มีดอีหรือยวันนั้นจำเป็นต้องเรียกมันขึ้นบ้าน

ลำดับ	ประโยคในฐานข้อมูล
147	ท่าทางจะอายุจริง อายุ คือตอนนี่ก็กำลังติดหนังสืออย่างรุนแรงนะคะพี่แหวน
148	ออกจากซาเคร์ท สปริง ก็เคลื่อนเดินไปดูบริเวณที่เรียกว่าเท็มเพิล คอร์ทยาร์ด
149	พวกเครื่องบวงสรวงทั้งหลายก็เอาไปตั้งตรงบริเวณที่เป็นแท่นบูชา ที่นี้มีจอโทรทัศน์ให้ดูภาพและเรื่องราวประกอบด้วย
150	เสียงระเบิดกึ่งก้อง แผ่นสังกะสีคมกริบนับไม่ถ้วนปลิวว่อนลอยควัดเฉวียนขึ้นไปกับกลุ่มควันนำหวาดเสียว
151	ดูคล้ายหมอนปักเข็ม การกอดติดกันช่วยให้นั่นเก็บความอบอุ่นไว้ข้างในพุ่ม
152	ขุมขันเงินเปล่า เหยหน้าตะเบ็งเสียงตอบ หญิงวัยกลางคนท่าทางเรียบร้อย
153	กระดาษ ผ้า เศษวัสดุ พื้นระนาบที่ใช้แผ่นวัสดุปะติด อาจเป็นแผ่นไม้
154	ตัวอย่างเช่น ระบบมัธยมศึกษาเป็นหลักสูตรที่ถูกบีบลงมาจากระบบอุดมศึกษา
155	ใจอยากปฏิเสธ หากความมืดและเปลี่ยวเล็กน้อยบริเวณนั้นทำให้รู้สึกว่าต้องพยักหน้ารับความช่วยเหลือที่
156	ผมนึกถึงด็อกเตอร์ปัทมา เธอเกาะแขนเหมือนต้องการเครื่องยึดเหนี่ยวชีวิตของคนเรามีวิธีดำเนินที่พิสดารอย่างไม่น่าจะเป็นไปได้
157	พลอยเดินเข้าไปใกล้ ได้ยินเสียงอีกฝ่ายบรรยาย แต่ก่อนชาวเลเช่นหลอโบบังด้วยเต้าทะเล
158	สาวหน้าหวานที่ไม่เคยอับอายกับการอ่านนิยายโรมานซ์ซึ่งเธอถือว่ามันคือนิยายแปลประเภทหนึ่งทำหน้าที่แปลกใจ
159	ซึ่งล้วนเป็นอันตรายต่อมนุษย์ทั้งสิ้น และได้พูดถึงสภาพความเป็นอยู่ของเกษตรกรว่ายังลำบากอยู่
160	ซึ่งเป็นย่านการค้าของหมู่บ้านประมาณครึ่งกิโล ตรงนั้นมีร้านขายของและมีผู้ขายอาหารประเภทส้มตำ
161	เพราะการเดินทางรอบโลกแบบอนาถาอย่างเขาไม่เตรียมของพวกนี้มาด้วยก็แย่เต็มที่ โดยเฉพาะมีโครงการผ่านเมืองยูงอย่างเมืองไทยด้วยแล้ว
162	หน้าเธอก็ไว้ที่เดิมแหละ ไม่รู้จะย้ายทำไม ทีเซอร์กว่นประสาทฉันอีก
163	ชายวัยหกสิบร่างผอมผิวคล้ำ ผมเป็นสีดอกเลาเกือบทั้งศีรษะ ความเก่าแก่ของโรงสีที่
164	แต่ความเห็นส่วนใหญ่ยังคงเน้นอยู่ที่ การทุจริตในกาหาเสียงเลือกตั้งของนักการเมืองกับประสิทธิภาพของ
165	เครือข่ายการสื่อสาร ชนิดหนึ่งเท่านั้น เนื่องจากการก่อกำเนิดความเจริญ เดิบโตของอิน

ลำดับ	ประโยคในฐานข้อมูล
	เทอร์เน็ต
166	และจะค้นพบได้ก็โดยการทำให้เหตุผลเท่านั้นเอง นอกจากวิธีภาษาวีธี
167	บางทีก็มีล้ามมาด้วยเป็นคนไทยที่เป็นแฟนของลูกชายเพื่อน รู้จักกันที่ภูเก็ตได้อาทิตย์กว่า
168	แต่นาฬิกาบนเวทีเดินเร็วกว่าเวลาจริงเพราะต้องให้เนื้อเรื่องจบภายในเวลาที่กำหนด
169	แต่ผมคัดค้านบอกว่า อย่าเลยเพราะอาบน้ําเสร็จกว่าจะเดินถึงบ้านได้ก็ต้องร้อนอีก
170	หนึ่งศุูนย์แปด แปดศุูนย์หนึ่ง แล้วทำไมจึงกลายเป็นเจ้าจิวไปได้ แต่ผมก็ดีใจเหลือหลายแล้ว
171	เสียงไซข้างตัวบอก เมื่อทั้งสองเดินมาถึงอ่างน้ำใหญ่มหึมานอกอาคารที่เรียกว่าเกรท บาบ
172	ทันใดนั้น โรงหนังทั้งหลายก็ลอยขึ้น ปีกหลังคาสายคล้ายนกยักษ์โคลงเคลง
173	แต่ก็กล่าวต่อ รู้สึกว่าคุณไม่ค่อยเต็มใจมาด้วยซ้ำ สีหน้านางสวาทเองก็พลอยเจื่อน
174	และ همینอย่างว่าถึงโดนเจ้าของตวาดขับไล่ไปไกลจุมกรอว่าหาย همینแล้วจึงให้เข้าบ้าน
175	ได้อารมณ์ได้สาระ และได้กลิ่นอายประวัติศาสตร์ เช่น ให้ตัวละครชาย
176	ประเดี๋ยวกี่ขึ้นนี่แล้วทำกิจกรรมสืบเชื้อสืบชาติ ประเดี๋ยวกี่ลง แล้วก็แหเข้าหา
177	ทำไมเครื่องบินโคลงเคลง ทำไมฉันดูหนังไม่รู้เรื่อง ทำไมพวกเขาบอกว่าให้รัดเข็มขัดอีก ทั้งที่ฉันรู้สึกอึดอัด
178	น่าเสียดายนะที่แกได้เรียนน้อยไปหน่อย แต่ฉันดูแกฉลาดเฉลียวกว่าคนอื่น
179	เล็บตีนเล่าหลุดเรียจมนุ่นอยู่ที่ใด ยังดีแขนข้างที่เหลือยังมีมือ
180	หรือผู้รับการฝึกอบรม ในความควบคุมของสถาบันการศึกษาวิชาการพยาบาล
181	จะได้ฟังคำสอน หรือไม่ก็บทกวีไพเราะของกรมหมื่นพิทยาลงกรณ์
182	ผลที่เกิดขึ้นทันที คือ คนที่โดนสารปรอทจะมีอาการเปลือกตาบวม
183	เมื่อไม่ติดยึดอยู่กับชื่อสมมุติเหล่านี้ หมา หมู หมู หนู ก็สามารถใส่สติปัญญา
184	ด้วยการให้ประชาชนมีส่วนร่วมในการปกครองและตรวจสอบการใช้อำนาจรัฐ
185	ข้อนี้ใครข้อนี้มันจ้วงตักน้ำแกงจืดใส่ปากได้เลย คนจีนกินข้าวไม่ค่อยได้ใช้ข้อนี้
186	สำนักทฤษฎีวิพากษ์ได้วิจารณ์ทฤษฎีสื่อมวลชนกระแสหลัก ที่กล่าวถึงการทำหน้าที่ของสื่อมวลชนว่า
187	ไม่มีบ้านพัก ต้องอาศัยวัดอยู่ ที่วัดคนก็มาก ไหนจะต้องเสียค่าใช้จ่ายแพงอีก

ลำดับ	ประโยคในฐานข้อมูล
188	ประโยคสุดท้ายกิริตาแน่นเสียงหนัก มือบางบิบกระชับมันกับมือของคนที่นั่งข้างตัว
189	ความนึกคิดคำนึงล่องลอยมือเลยข้างลง เอ็งอย่าใจลอย ทำอะไรอย่าใจลอย
190	ศิระหนนผ้าโพกผมต่างหมอน ผ้าห่มหนาคลุมจากเจ้าขึ้นมาถึงคอ แต่ก็แสนขัดข้องอยู่ในใจ
191	แม้ว่าพืชพวกนี้จะเผลอหน้ากันขึ้นมาในหน้าร้อน แต่ลักษณะใบและต้นของมันก็ออกแบบมาให้กันน้ำค้างแข็งแม่คะนึ่งและลมหนาว
192	มาสยบมากลบทับ ข้าพร้อมยอมสิ้นถูกหมิ่นค่า ด้อยต่ำธรรมดาไร้ระดับ
193	คล้ายติดคราบตมขุนต่างวุ่นวาย ห่วงชีวิตติดวนไปจนวาย กว่าชีพตรมล้มตายมิได้พบ
194	เที่ยงค่อยยังชั่วขึ้นใหม่ลูก นางสาวทถามเมื่อตามกันมาจนทะเลถึงด้านหน้าโรงแรม
195	เอกสาร และสิ่งพิมพ์ที่มีคุณค่าจากประชาชนที่มีจิตศรัทธา เพื่อจัด
196	อย่างไรก็ตาม นับว่ารัฐบาลสมัยจอมพลสฤษดิ์ ธนะรัชต์ ได้แสดงให้เห็นประจักษ์ว่ามีความสนใจต่อปัญหาจังหวัดชายแดนภาคใต้
197	การปฏิรูปการเมือง โดยรวมจะปรากฏเกินขีดว่าการมองโลกในแง่ดีว่าความสำเร็จของการปฏิรูปการเมืองอยู่ใกล้แค่เอื้อม
198	สังคมและวัฒนธรรม ไม่ว่าจะเป็นผู้ส่งสารในระดับปัจเจกบุคคล
199	ไปในคืนนี้ ยิ่งเมื่อดวงหน้าที่คุ้งเคยนิ่วหน้าครุ่นคิดหนักอย่างไม่รู้ตัวขณะที่เก็บกวาดจานบนโต๊ะ
200	หญิงสาวรีบบอก เมื่อน้ำเสียงของเจ้านายราวกับจะดำนิว่าเป็นความผิดของเธออย่างงั้นล่ะ
201	และองค์การเอกชนที่เป็นไปในรูปแบบของโรงงานหรือห้างร้านบริษัทใหญ่
202	เรื่องอะไรจะหนักกลับไปคนเดียวได้ไง แล้วใครจะช่วยฉันขึ้นไต่แบทแมนล่ะ
203	หรือการให้ยาอันตราย ยาควบคุมพิเศษ วัตถุออกฤทธิ์ต่อจิตและประสาท
204	ชูกระดาดาร้อน อาไต้งหีบกระดาดของประสิทธิ์ศักดิ์มาดู มองเส้นยุกยิกกลางกระดาดอยู่พักหนึ่ง
205	ขอเครตีสได้พูดถึงผู้ที่ตกอยู่ในห้วงรัก จากญาติที่พี่ศรีสตั้งให้
206	ด้านใต้ทางเพมาโคหันรับมรสุมเต็มที เพมาโคจึงเป็นพื้นที่ที่มีฝนตกตลอดปี
207	สักสามสี่วันเสียก่อนจึงค่อยเริ่มลุยทำวิทยานิพนธ์ คิดว่ามาเถี่ยสอาจารย์ที่ปรึกษาคงเข้าใจ

ลำดับ	ประโยคในฐานข้อมูล
208	โดยไม่ได้พูดปิดเลยสักคำเดียว และผู้กองก็ไม่ได้เฉลียวใจสักนิดว่ามีอะไรเคลือบคลุม วาจาเหล่านั้นอยู่บ้าง
209	ข้าเกาะแก่งแอ่งผาพนาโตรกหนึ่งหุบเหวต่ำโขดซ้ำชะตากรรม ข้าภูผามหาสมุทรสุดคะเน หนึ่งว่าเหวระอุปะทุระล่ำ
210	ริมถนน รถเข็นขายขนมสี่ชั้นจอดจอดนิ่งอยู่ สีสันของขนมบนรถดูน่ากลัวมากกว่าน่ากิน
211	มองผ่านหน้าต่างก็เห็นบีเอ็มดับเบิลยูสีน้ำเงินซึ่งคุ้นตาเหลือเกิน แม่ฟุ้งเหมือนติดจรวด ไปที่ประตูพร้อมรอยยิ้มกว้าง
212	เมื่อเปิดเปลือกตาขึ้นมา รับรู้ว่าคุณคนรอบข้างกำลังเคลื่อนไหว ชีวิตดำเนินไปตามวิถีของตนเอง
213	ใบหน้าหวานเชิดหน้าราวกับหมดเรื่องพูดแล้ว และแน่นอนเธอต้องเพิกเฉยต่อเสียงนกเสียงกาที่ร้องวิ๊ดวิ๊ดกับท่าทางของเธอและตาเกือบอูญ
214	พูดพลางเอื้อมมือไขประตู กลับไปอยู่ป่า พบคู่ ช่วยเขากกไข่
215	ผู้พระธรรม ถ้วยชามสังข์โลก ลวดลายปูนปั้น งานแกะสลักไม้ พระพุทธรูปใหญ่หน่อย
216	จอกแหนลอยฟองกำลังกระเพื่อมคลื่น เราตะลึงยั้งนิ่ง แต่ว่าตรงหน้านั้นยังมีป่าอ้อกว้างหลายไร่บังเว้งน้ำไว้
217	คืนนั้นผมยังโดนคนเชียร์แซกตามตื้อให้คุณโชว์หลายครั้งต้องบอกว่าไม่ ไม่
218	ก็จะให้ทำยังไงละ เรื่องมันเกิดไปแล้ว เธอไม่ได้ใช้งานคอมพิวเตอร์อยู่นะ
219	ตอบข้อหาหรือและให้คำปรึกษาแก่ส่วนราชการที่เกี่ยวข้อง จัดทำนิติกรรมสัญญาแก่กระทรวงและหน่วยงานในสังกัด
220	เพมาโค ตำนานแห่งพันธุ์พืช หุบเขาหลังดาซิงลาเป็นเทือกหิมาลัยที่หันหน้าเข้าทิศใต้
221	ความสำคัญในการดำรงชีวิตของมนุษย์ คือการต่อสู้เพื่อความอยู่รอด
222	แต่ไม่มีเงินซื้อตัวเครื่องบินกลับบ้าน นี่ก็กำลังส่งข่าวคราวไปบอกลูก
223	แล้วชำระล้างฝุ่นโคล พวงนี้กองทัพจะหยุดพักหนึ่งวัน แรมคืนอีกหนึ่งคืน
224	ดูแล้วไม่น่าเชื่อว่าเขาจะเดินทางมาได้ถึงค่อนโลกเช่นนี้ รุ่งเช้าผมตื่นนอนเช้ากว่าปกติ
225	นอกจากจะเพิ่มความน่าสนใจให้ผู้อ่านคิดไตร่ตรองตามแล้ว ยังเป็นหนทางไปสู่ความเข้าใจแนวคิดของเรื่องได้อย่างลึกซึ้งขึ้นด้วย
226	ฝนเริ่มโปรยเม็ดหนาขึ้นจนร่มไม่ใหญ่ก็ดูจะปกป้องเจ้าของเสื้อเชิ้ตสีฟ้าอมเทาไว้ไม่ได้

ลำดับ	ประโยคในฐานข้อมูล
227	พาไปวัดไทยให้รู้จักเพื่อนคนไทยที่มาจากภาคอีสาน หลังแต่งงานแล้วเขาส่งฉันไปเรียนภาษาเยอรมันเกือบ
228	การบริหารภายในสำนักงานปลัดกระทรวงและกระทรวง ให้สอดคล้องกับนโยบายของรัฐบาล
229	นี่ก็อยากจะปรับปรุงกิจการโรงสีให้ดีขึ้นและทันสมัยกว่าที่เป็นอยู่ในปัจจุบัน
230	ถ้าต้องการอารมณ์โรแมนติก หาดสวยสุดที่อยู่บนแผ่นดินไทย ไม่รวมเกาะชื่อเสียงไพเราะว่า
231	เป็นบ้า ฉันเอาไม้หวดเอา เด็กพูดพลางเดินรีเขามาวางตัวโซดลงในมือผู้เป็นพ่อ
232	มีกองเพลิงอยู่ใกล้กายพร้อมคบไฟโชนแสง บักเรียงกันตลอดแนว แล้วจะเลยดำหัว
233	เขาก็เพียงปรายตามองเธออย่างรำคาญใจก่อนจะส่งสายตาท่องไปดูว่าป้ายนั้นเสียหายหรือเปล่า
234	น้ำมะเฟืองทำเสียงมีลับลมคมนัย เสียงกรรณตะโกนเสริมขึ้น น้ำมะเฟืองเรื่องมาก
235	จริงอย่างที่รินว่า วันนี้แซมใส่เสื้อยืดแขนยาวสีดำกับกางเกงผ้าสีเดียวกับเสื้อ
236	เขาเปิดกระโปรงรถไฟร์วีล ตรวจเช็คความเรียบร้อยอยู่ขณะที่ ในชุดเสื้อกางเกงยืนทั้งคู่ดูหิวกระป๋องเดินทางใบยอมตามไปที่รถ
237	โดยคำนึงถึงประโยชน์สูงสุดที่ประชาชนจะได้รับเป็นสำคัญและ มาตรา
238	แล้วอีกอย่าง เราก็ไม่ใช่แฟนกันสักหน่อย เพราะฉะนั้นมันไม่ใช่เรื่องธรรมดาแน่นอน
239	สายน้ำหันเข้าหาฟ้าสีหมากสุกด้านตะวันตก กอไผ่บนฝั่งที่สูงเด่นดูโปร่งเบาเหนือทิวไม้อื่น
240	โดยปรับปรุงให้สมบูรณ์ชัดเจน และสอดคล้องกับสภาวะการณ์ในปัจจุบันยิ่งขึ้น
241	ผมว่าคนผู้นั้นอันตรายมาก อันตรายต่อประเทศชาติ ต่อความมั่นคง
242	ความเข้าใจในสิ่งที่เปลี่ยนแปลงทางด้านวิทยาศาสตร์และเทคโนโลยี จึงได้จัดแหล่งความรู้ที่สำคัญยิ่งในการจัดการศึกษาทางด้าน
243	หรือว่าฝูงวัวควายลงมากินน้ำเพ็งจะขึ้น ข้างบนมีบ้านคน เมื่อโรปุ้จะชื่อชนมมากินบ้างล่ะ
244	กลุ่มคว้นเต็มท้องฟ้า พระดอนเก็บกลดมุ่งหน้าข้ามทุ่งไปทางทิศนั้นด้วยความอยากรู้
245	ได้แปรสภาพฐานะของ ผู้ส่งสาร ที่เป็นรูปธรรมที่มีตัวตนจับต้องได้
246	พอยืนซึ่มมองร่างที่กำลังเดินโผล่ใกล้เข้ามา แม่ก้าวลงจากบันไดและเริ่มร้องไห้

ลำดับ	ประโยคในฐานข้อมูล
247	แต่ถ้าสาวคนนั้นนุกมาอีกรอบ ก็ตัวใครตัวมัน ธรรมชาติยกหน้าอย่างโล่งอกที่หญิงสาว เข้าใจแต่โดยดี
248	ก่อนจะปล่อยตัวให้เข้าสู่การพักผ่อนอย่างแท้จริง เขาสะดุ้งตื่นตอนที่ได้ยินเสียงแม่เรียก หา
249	ยังผู้หญิงทั้งคู่นี้อีกเล่า อาจเสแสร้งแกล้งเข้ามาตีสนิท เพื่อล้วงเอาความสงสัยออกมาให้ สิ้น
250	เข้าสู่วันรุ่งขึ้น ยังแต่งชุดนอนเสื้อแขนสั้นกางเกงขายาว ผมดำยุ่งเพราะเจ้าตัวเสยมือ อย่างลวก
251	วิทยาศาสตร์และเทคโนโลยีให้กับประชาชนทั้งในระบบโรงเรียนและนอกระบบโรงเรียน
252	ไฟที่เจ็ดยังทางเดินในตึกเปิดสว่างตามปกติ แต่ไม่มีผู้คน ผมถอนใจอย่างโล่งอก
253	ผมเห็นท่าทางเขาอ่อนเพลียมาก จึงบอกให้เขานอนก่อน เขากลางคำสวัสดิ์กับผมแล้ว คลานเข้ามางู้ง
254	ไม่อาจจำนำของได้ หรือได้ก็คงต้องมีวิธีการยุ่งยากพิลึก
255	ประชาชนมีหน้าที่อย่างเดียวคือรับเฉพาะข้อมูลของคนที่ฉลาดกว่าป้อนให้เชื่อและปฏิบัติ ตาม
256	นี่เธอคนนั้นนะ คุรุภารดีที่มาอยู่กับเราเพียงสามเดือนก็ได้ย้ายเข้าอำเภอก็เพราะเป็น หลานป่าไม้จังหวัดเรียกค่าผาง
257	ชนิดที่ไม่มีใครคาดฝันรวมทั้งตัวเธอเองด้วย ต้นเหตุของโศกนาฏกรรมเกิดขึ้นในวันก่อน ออกเดินทาง
258	อีเบี้ยว ปากเบี้ยว อีพลอย ตาแดง อีนิล ดำ ไข่ตูบ ไข่หมอก
259	วิทยาศาสตร์เพื่อการศึกษา เป็นการจัดการศึกษาตามอัธยาศัยในรูปแบบของ วิทยาศาสตร์เพื่อเป็นการพัฒนาคนให้มีความรู้
260	ไปอยู่เฉพาะที่การทำให้การเลือกตั้งสะอาดบริสุทธิ์เพียงอย่างเดียว เป็นผลให้ภารกิจของ
261	เช่นในตอนหนึ่งที่ทวนเข้ามาในความทรงจำของฉินตะโรก็คือตอนที่สงคราม
262	เหลียงซึ่งมีอายุน้อยกว่าจางอิมันนับสิบปีจึงให้กำเนิดบุตรสาวคนหนึ่ง
263	หัวโตและหูตก ไข่เทา อีสำลี ขาวปุย อีเตี้ย ขาหน้าหรือขาหลังกะเผลก
264	ทักท้วงมาความที่ท่านเห็นด้วยอยู่แล้ว ในวิธีที่จะปรับชีวิตเปลี่ยนอาหารแก่ผู้คนด้วย วิธีการอย่างเฉียบ

ลำดับ	ประโยคในฐานข้อมูล
265	แต่บัดนี้มันก็ได้มาทอดท่อนนอนสนิทเป็นระเบียบอยู่เฉพาะหน้าแล้ว เปลี่ยนสภาพจากต้นสักบนดอยสูงมาเป็นแพชุงที่มีค่ามหาศาล
266	ฉินตะโรคคิดว่าตัวเองไม่ชอบพ่อเพราะได้รับอิทธิพลจากแม่ การที่ตัวเองไม่ชอบพ่อจะต้องได้รับอิทธิพลจากแม่อย่างแน่นอน
267	รวมตลอดถึงเป็นปัญหาสืบเนื่องมาจากนโยบายของรัฐบาลที่ขาดประสิทธิภาพ
268	และไม่มีองค์การควบคุมการเลือกตั้งอย่างในปัจจุบันเสียอีก จึงเป็นเรื่องปกติที่ความสิ้นหวังต่อ
269	การปฏิรูปการเมืองจะต้องไปเริ่มที่เด็กเกิดใหม่ ด้วยการให้การศึกษาสอนเรื่องประชาธิปไตยยบวิสุทธิ
270	โครงการศึกษาภาษาและวัฒนธรรมท้องถิ่นของข้าราชการจังหวัดชายแดนภาคใต้
271	กำหนดดีชั่วของสรรพสิ่ง และนำพาสังคมหรือรัฐนาวาไปด้วยปรีชาของชนชั้นผู้นำ
272	นอกจากเต้าหู้ อาหารอย่างอื่นที่ทำขึ้นเพื่อเป็นการถนอมอาหารให้อยู่ได้นาน
273	คือครูใหม่ซึ่งผมขอออกตรงนี้เลยว่าชื่อครูวิชัย ครูวิชัยพักที่บายครูใหญ่ได้ชวนครูใหญ่ห่อข้าวไปกินที่โรงเรียน
274	เจ้าหญิงที่ไม่ติดยึดกับความเป็นเจ้าหญิงก็สามารถเป็นครูสอนหนังสือเด็กยากจนได้
275	กลุ่มผู้ส่งสาร นั้น เป็นเพียงกลุ่มคนกลุ่มน้อยที่ได้รับการฝึกฝนเรื่องการส่งสาร
276	เกือบไม่ทัน ทำอย่างนั้นยิ่งอยากมอง ไม่รู้หรือหรือ เขาแก้ม้วยั่วเล่นและจ้องมองหน้าขาว
277	อันใดก็ถูกใจ กลัวไปหมด แต่ก็อยากมีชีวิตที่ดีกว่าเดิม
278	ของจังหวัดว่าเราควรมีอุดมการณ์หรือแผนการทำงานอย่างไร เพื่อให้แผนสำเร็จผลหรือประสบความสำเร็จ
279	สุขาภิบาล องค์การบริหารส่วนท้องถิ่นตามที่รัฐมนตรีประกาศในราชกิจจานุเบกษา
280	มองเห็นประชาชนเป็นเพียงนั่งร้านที่จะใช้ปีนขึ้นสู่จุดหมาย เมื่อประสบผลตามความต้องการแล้วก็สลัดนั่งร้านทิ้ง
281	ลมเย็นยามบ่ายทำให้หัวใจที่อัดแน่นหนักอึ้งของเขาผ่อนคลายลง มีเรื่องอะไรเกี่ยวกับพ่อหรือฮะ
282	วินาทีนี้เธอผู้นั้นรวมทั้งตัวเขาด้วยได้เป็นอิสระแล้ว อิสระจากความเศร้าหมองที่พันธนาการใจเขาไว้นับจากวันที่เธอจากไป

ลำดับ	ประโยคในฐานข้อมูล
283	หากเป็นเอกชนก็น่าจะผลิตรายการที่เป็นสินค้า เพื่อขายให้แก่ผู้ที่มีอำนาจซื้อมากกว่า
284	อยู่เมืองจีนในปัจจุบันนี้มีของกินอุดมสมบูรณ์ ขอให้เงินซื้อ
285	ใกล้เขาหลัก ศึกษาความเปลี่ยนแปลงจากสีนามิ ป่าชายเลน เส้นทางศึกษารวมชาติ
286	ถึงจะเหนื่อย ฉันก็ไม่มีวันถอย ฉันจะบุกไปจนกว่าจะถึงบึงหญ้าใหญ่
287	งานวิจัยของอารยาได้พบว่ากระบวนการควบคุมให้การสื่อความหมายเป็นไปตามความตั้งใจของผู้ส่งสารนั้น
288	คู่สนทนารู้สึกตัวเบาเกิดความเข้าใจแจ้ง ว่าเราทุกคนล้วนเกิดมาร่วมทุกข์
289	ทิพกฤตาเห็นฝายนั่นเงยหน้าขึ้นมองท้องฟ้าสีหม่น มือใหญ่ยกขึ้นปาดหยดน้ำซึ่งเริ่มจับตัวบนศีรษะได้รูปนั้น
290	เพื่อมิให้อาหารเสียเร็ว นั่นคือ เน้นความสะดวกที่จะเก็บอาหารไว้นาน
291	แม่เป็นบุตรสาวของพนักงานธนาคาร เกิดในโตเกียว และมาเติบโตที่โอซากา
292	แต่อาชีพของเธอยังคงดูดีกว่าพวกมีการศึกษาสูงใส่สูทผูกไท เข้าไปนั่งในสภาและรวมหัวทุจริตโกงกินบ้านเมืองจนฉิบหาย
293	น้อย มาจากการที่เกษตรกรขายผลผลิตในราคาต่ำ ไม่คุ้มต่อการซื้อหาปัจจัยการผลิตมาใช้
294	แม่มีความลังเลที่จะพูดความจริงและไม่ยอมให้ฉินตะโรพูดอะไรด้วย
295	ใจของเขาได้รับการปลดปล่อยและหญิงสาวในกรอบรูปก็เช่นกัน รอยยิ้มที่สวยงามจึงปรากฏให้เห็นเป็นครั้งสุดท้าย
296	และช่วยกันทำนุบำรุงรักษารวมชาติเหล่านี้ให้ดูสวยงามตามธรรมชาติ
297	อะริโยะฉิได้พูดถึงสารปรอทในยาฆ่าแมลงของญี่ปุ่นว่ามีมากกว่าของอเมริกาถึง
298	จนคนที่นั่งอ่านหนังสือรออยู่ที่โต๊ะทำงานเพียงคนเดียวในห้องสี่เหลี่ยมกว้างใหญ่นั้นอดหัวเราะออกมาไม่ได้
299	เช่น ตลาดที่วุ่นวายจอแจ การจราจรที่ติดขัด สภาพบ้านเมืองที่เก่าแก่
300	หรือที่ได้รับอนุญาตจากทางราชการให้จัดตั้ง หรือสถาบันการศึกษาที่คณะกรรมการรับรอง

ประวัติผู้เขียนวิทยานิพนธ์

นายสุรพล วรรณาทราทร เกิดวันจันทร์ที่ 6 กรกฎาคม พ.ศ.2530 ที่จังหวัดขอนแก่น เข้าศึกษาระดับประถมศึกษาโรงเรียนเคหะบางพลี จังหวัดสมุทรปราการ สำเร็จการศึกษาระดับมัธยมศึกษาตอนต้นที่ โรงเรียนสันติวิทยาคม เชียงราย สำเร็จการศึกษาระดับมัธยมศึกษาตอนปลายจาก โรงเรียนสามัคคีวิทยาคม เชียงราย จากนั้นในปี 2549 ได้เข้าศึกษาระดับปริญญาบัณฑิต จนได้รับเกียรตินิยมอันดับสอง จากหลักสูตรปริญญาวิศวกรรมศาสตรบัณฑิต สาขาวิชาวิศวกรรมคอมพิวเตอร์ สำนักวิชาเทคโนโลยีสารสนเทศ มหาวิทยาลัยแม่ฟ้าหลวง จังหวัดเชียงราย และจากนั้นในปี 2553 ได้เข้าศึกษาต่อระดับปริญญาโท สาขาวิชาวิศวกรรมศาสตรมหาบัณฑิต สาขาวิชาวิศวกรรมคอมพิวเตอร์ ภาควิชาวิศวกรรมคอมพิวเตอร์ คณะวิศวกรรมศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย จังหวัดกรุงเทพมหานคร