

การเชื่อมต่อลายเส้นที่ขาดหายไปของอักขรตัวพิมพ์ภาษาไทย

CONNECTION OF BROKEN LINE OF THAI PRINTED CHARACTERS



วิทยานิพนธ์นี้เป็นส่วนหนึ่งของการศึกษาคำหลักศิลาจารึกปริศนาวิภาษศาสตร์มหาบัณฑิต

สาขาวิชาเทคโนโลยีสารสนเทศ

บัณฑิตวิทยาลัย

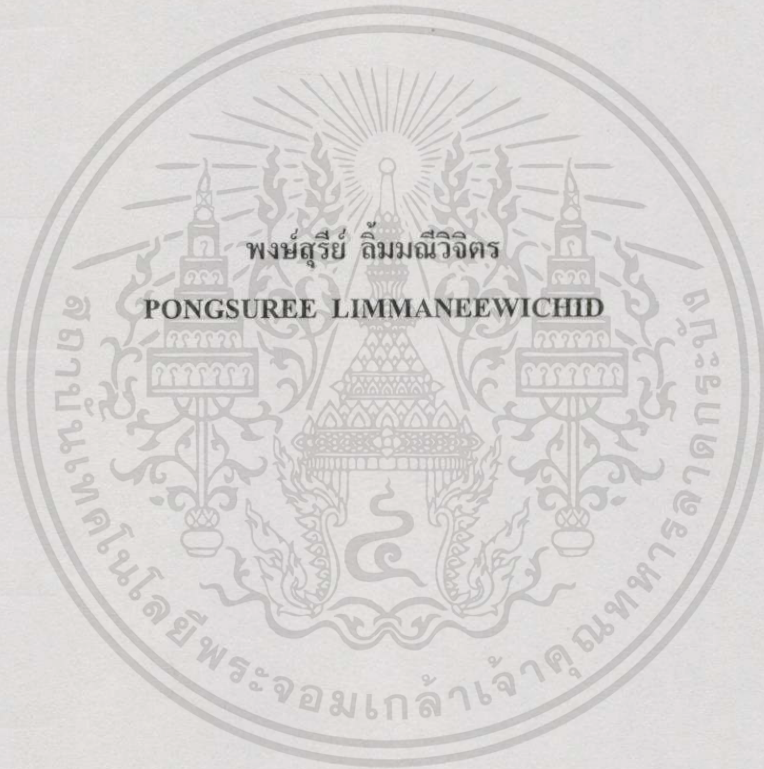
สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

พ.ศ. 2548

ISBN 974-622-931-1

การเชื่อมต่อสายเส้นที่ขาดหายไปของอักษรตัวพิมพ์ภาษาไทย

CONNECTION OF BROKEN LINE OF THAI PRINTED CHARACTERS



วิทยานิพนธ์นี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญาวิทยาศาสตรมหาบัณฑิต

สาขาวิชาเทคโนโลยีสารสนเทศ

บัณฑิตวิทยาลัย

สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

พ.ศ. 2543

38038

ISBN 974-622-931-1

ไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

20 พ.ย. 2543

CONNECTION OF BROKEN LINE OF THAI PRINTED CHARACTERS



PONGSUREE LIMMANEEWICHID

A THESIS SUBMITTED IN PARTIAL FULFULLMENT

OF THE REQUIREMENT FOR THE DEGREE OF

MASTER OF SCIENCE IN INFORMATION TECHNOLOGY

SCHOOL OF GRADUATE STUDIES

KING MONKUT'S INSTITUTE OF TECHNOLOGY LADKRABANG

2000

ISBN 974-622-931-1



COPYRIGHT 2000

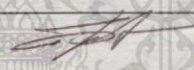



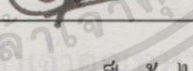
SCHOOL OF GRADUATE STUDIES

KING MONKUT'S INSTITUTE OF TECHNOLOGY LADKRABANG

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บัณฑิตวิทยาลัย
สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง
ใบรับรองวิทยานิพนธ์

หัวข้อวิทยานิพนธ์ การเชื่อมต่อลายเส้นที่ขาดหายไปของอักษรตัวพิมพ์ภาษาไทย
CONNECTION OF BROKEN LINE OF THAI PRINTED
CHARACTERS
ชื่อนักศึกษา นายพงษ์สุรีย์ ลิ้มมณีวิจิตร
รหัสประจำตัว 41067013
ปริญญา วิทยาศาสตรมหาบัณฑิต
สาขาวิชา เทคโนโลยีสารสนเทศ
อาจารย์ผู้ควบคุมวิทยานิพนธ์ รศ.ดร.วิเชียร เปรมชัยสวัสดิ์

คณะกรรมการสอบวิทยานิพนธ์		ลายมือชื่อ
รศ.ดร.วิเชียร	เปรมชัยสวัสดิ์	
รศ.นุชรี	เปรมชัยสวัสดิ์	
ผศ.ดร.บุญฉวีร์	เกรือตราขู	
ดร.รัฐการ	อภิวัฒน์วาท	
ดร.ประจวบ	วานิชชัชวาล	

วัน/เดือน/ปี ที่สอบ 16 สิงหาคม 2543 เวลา 10.30 น. เป็นต้นไป

สถานที่สอบ ณ ห้อง 231-232 ชั้น 2 อาคารสำนักวิจัยและบริการคอมพิวเตอร์

บัณฑิตวิทยาลัยรับรองแล้ว

(รศ.ดร.บุญฉวีร์ อัทชู)

รักษาราชการแทนคณบดีบัณฑิตวิทยาลัย

วันที่ 25 เดือน กันยายน พ.ศ. 2543

หัวข้อวิทยานิพนธ์

การเชื่อมต่อสายเส้นที่ขาดหายไปของอักษรตัวพิมพ์ภาษาไทย

ชื่อนักศึกษา

นาย พงษ์สุรีย์ ลิ้มฉวีจิตร

รหัสประจำตัว

41067013

ปริญญา

วิทยาศาสตรมหาบัณฑิต

สาขาวิชา

เทคโนโลยีสารสนเทศ

พ.ศ.

2543

อาจารย์ผู้ควบคุมวิทยานิพนธ์

รศ. ดร. วิเชียร เปรมชัยสวัสดิ์

บทคัดย่อ

วิทยานิพนธ์นี้นำเสนอวิธีการใหม่ทำการซ่อมแซมอักษรตัวพิมพ์ไทยที่ขาด ก่อนที่จะนำไปสู่กระบวนการรู้จำตัวอักษร วิธีการที่นำเสนอนี้จะใช้การพิจารณาคำแหน่งการเหลื่อมล้ำกันของกรอบภาพที่ถูกตัดขาดและการพิจารณาลักษณะเด่นของตัวอักษรภาษาไทยเป็นหลัก จากการศึกษาพบว่าลักษณะการขาดของตัวอักษรแบ่งเป็น 2 ประเภทคือ การขาด โดยมีการเหลื่อมล้ำกันของกรอบภาพที่ขาดและการขาด โดยที่ไม่มีมีการเหลื่อมล้ำกันของกรอบภาพ วิธีการซ่อมแซมจะขึ้นกับลักษณะการขาดของตัวอักษรทั้งสองแบบ ตำแหน่งของการเชื่อมต่อหาได้จากผลการพิจารณาค่าฮิสโตแกรมของภาพตัวอักษร กระบวนการซ่อมแซมตัวอักษรขาดแบ่งออกเป็น 2 ขั้นตอนคือการตรวจสอบลักษณะตัวอักษรขาดและการซ่อมแซมตัวอักษรขาด ข้อมูลภาพตัวอักษรขาดได้ทำการทดสอบด้วยโปรแกรมรู้จำภาษาไทยที่มีขายในปัจจุบัน พบว่าโปรแกรมเหล่านี้ไม่สามารถรู้จำตัวอักษรขาดได้เลย แต่หลังจากที่ได้รับการซ่อมแซมตัวอักษรขาดด้วยวิธีการที่นำเสนอนี้ทำให้สามารถรู้จำเพิ่มขึ้นได้

Thesis Title	Connection of Broken Line of Thai Printed Characters
Student	Mr. Pongsuree Limmaneewichid
Student ID.	41067013
Degree	Master of Science
Programme	Information Technology
Year	2000
Thesis Advisor	Assoc. Prof. Dr. Wichian Premchaiswadi

ABSTRACT

This thesis presents a new scheme for repairing broken characters before passing to the character recognition process. The overlapping area of broken image and specific features of character are employed in this purpose. From the study of broken images, it can be divided into 2 categories: having an overlapping area and non-overlapping area. The method of repairing of these broken images is depended on the type of broken images. Histogram is used to detect connecting position of broken images. There are two steps in the proposed scheme: determination of broken characters and repairing of broken characters. The broken characters images have been tested with commercially available Thai OCR softwares. These softwares can not recognized the broken characters at all.

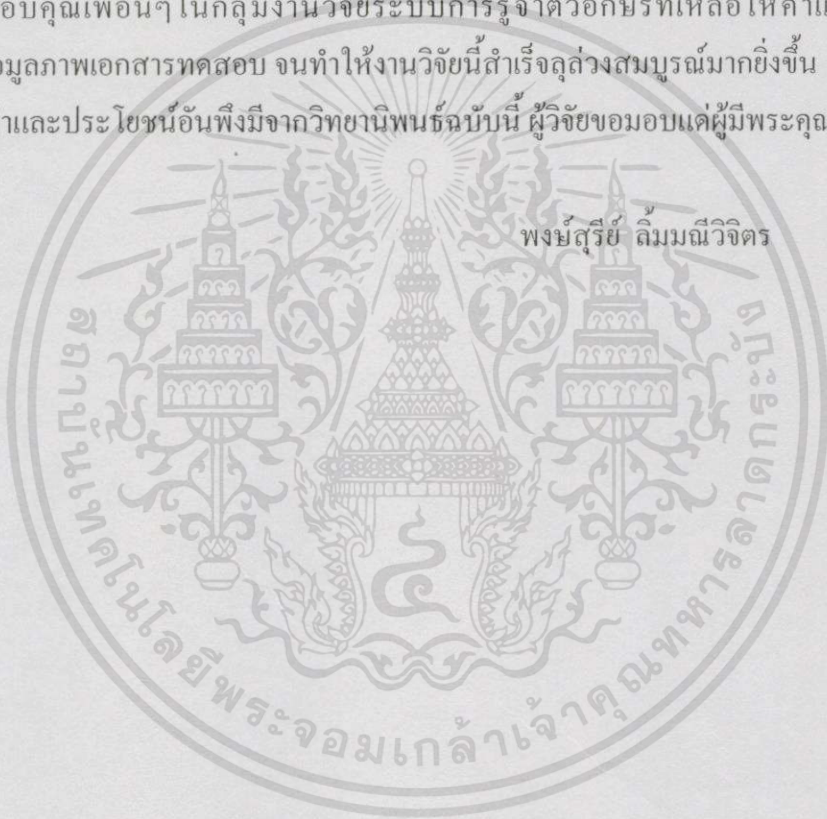
กิตติกรรมประกาศ

วิทยานิพนธ์ฉบับนี้สำเร็จลุล่วงได้อย่างดี ด้วยคำแนะนำและคำปรึกษาจาก รศ. ดร. วิเชียร เปรมชัยสวัสดิ์ ซึ่งเป็นอาจารย์ผู้ควบคุมวิทยานิพนธ์ ผู้วิจัยรู้สึกซาบซึ้งในความอนุเคราะห์จากท่าน และขอกราบขอบพระคุณเป็นอย่างสูง

ขอขอบคุณเจ้าหน้าที่คณะเทคโนโลยีสารสนเทศทุกท่าน ที่ได้ช่วยดำเนินงานในการจัดการสอบและให้ข้อมูลที่เป็นประโยชน์ในการเตรียมตัวในการเสนอวิทยานิพนธ์

ขอขอบคุณเพื่อนๆในกลุ่มงานวิจัยระบบการรู้จำตัวอักษรที่เหลือให้คำแนะนำและอนุเคราะห์ข้อมูลภาพเอกสารทดสอบ จนทำให้งานวิจัยนี้สำเร็จลุล่วงสมบูรณ์มากยิ่งขึ้น

คุณค่าและประโยชน์อันพึงมีจากวิทยานิพนธ์ฉบับนี้ ผู้วิจัยขอบแต่ผู้มีพระคุณทุกท่าน



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญ

	หน้า
บทคัดย่อภาษาไทย	I
บทคัดย่อภาษาอังกฤษ	II
กิตติกรรมประกาศ	III
สารบัญ	IV
สารบัญตาราง	VI
สารบัญรูป	VII
บทที่ 1 บทนำ	1
1.1 ความเป็นมาและความสำคัญของปัญหา	1
1.2 วัตถุประสงค์ของงานวิจัย	2
1.3 ทฤษฎีและหลักการที่เกี่ยวข้อง	2
1.3.1 นิยามตัวอักษรขาด	3
1.3.2 กระบวนการซ่อมแซมตัวอักษรขาด	3
1.4 แผนการดำเนินงาน	7
บทที่ 2 การประมวลผลภาพเบื้องต้น	8
2.1 โครงสร้างของตัวอักษรภาษาไทย	8
2.2 การประมวลผลภาพเบื้องต้น	9
2.2.1 การแยกบรรทัด	10
2.2.2 การหาระดับของตัวอักษร	11
2.2.3 การหาตำแหน่งและขนาดของภาพตัวอักษร	12
บทที่ 3 การซ่อมแซมตัวอักษรตัวพิมพ์ภาษาไทยที่ขาด	15
3.1 นิยามตัวอักษรขาด	15
3.2 การวิเคราะห์ตัวอักษรขาด	16
3.3 ลักษณะตัวอักษรขาด	18
3.3.1 ตัวอักษรขาด 2 ส่วน โดยมี Character Frame เหลือมั่วกัน	19
3.3.2 การวิเคราะห์ตัวอักษรครึ่งหลัง	20
3.3.3 การวิเคราะห์ตัวอักษรครึ่งหน้า	20

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษเท่านั้น เมื่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดทั้งสิ้น ถือว่าห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญ (ต่อ)

	หน้า
3.3.4 ตัวอักษรขนาด 2 ส่วน โดยไม่มีส่วนที่เหลื่อมล้ำกัน	23
3.3.5 ตัวอักษรขนาด 3 ส่วน	24
3.3.5.1 ส่วนที่กว้างอยู่ครึ่งบน	24
3.3.5.2 ส่วนที่กว้างอยู่ครึ่งล่าง	26
3.3.6 ตัวอักษรขนาด 4 ส่วน	24
3.3.6.1 การเชื่อมต่อ Seg[0],Seg[2]	30
3.3.6.2 การเชื่อมต่อ Seg[0],Seg[3]	30
3.3.6.3 การเชื่อมต่อ Seg[1],Seg[3]	30
บทที่ 4 ผลการทดลอง	32
บทที่ 5 สรุปผลการวิจัยและข้อเสนอแนะ	36
เอกสารอ้างอิง	38
ภาคผนวก ก ตัวอย่างเอกสารที่ใช้ในการทดลอง	39
ภาคผนวก ข การใช้งานโปรแกรมเชื่อมต่อตัวอักษรขนาด	43
ภาคผนวก ค ภาพตัวอย่างการเชื่อมต่อตัวอักษรขนาด	50
ประวัติผู้เขียน	53

สารบัญตาราง

ตารางที่	หน้า
1.1 แสดงผลการทดสอบโปรแกรม OCR ในการรู้จำตัวอักษรขาด	2
2.1 ตารางแสดงการแบ่งตัวอักษรภาษาไทยตามระดับโครงสร้างของคำ	9
4.1 แสดงผลการทดสอบโปรแกรม OCR ในการรู้จำตัวอักษรขาด	35
4.2 แสดงผลการทดสอบโปรแกรม OCR ในการรู้จำตัวอักษรขาดที่ซ่อมแซมแล้ว	35
4.3 แสดงสถิติของการพบตัวอักษรขาด	35



สารบัญรูป

รูปที่	หน้า
1.1 แสดงข้อมูลภาพตัวอักษรขาด	1
1.2 แสดงตัวอย่างข้อมูลภาพที่ใช้ในการทดสอบ	2
1.3 แสดงกระบวนการซ่อมแซมตัวอักษรขาด	3
1.4 แสดงการแยกบรรทัดและตัวอักษร โดยใช้ฮิสโตแกรมและการหาขอบภาพ	4
1.5 แสดงกรอบภาพจากการหากรอบภาพ	5
1.6 แสดงตัวอักษรขาดทั้ง 2 แบบ	5
1.7 แสดงขอบเขตของตัวอักษรและส่วนที่เหลือมัลลา	6
1.8 แสดงการเชื่อมจุดขาด	6
1.9 แสดงการใช้ฮิสโตแกรมหาค่าแห่งการเชื่อม	7
2.1 พยัญชนะและวรรณยุกต์ไทย	8
2.2 ระดับพยัญชนะไทย	8
2.3 แสดงกระบวนการประมวลผลภาพเบื้องต้น	9
2.4 แสดงฮิสโตแกรมและการแบ่งบรรทัด โดยการวิเคราะห์ค่าฮิสโตแกรมในแนวนอน	10
2.5 แสดงระดับของตัวอักษรทั้ง 3 ส่วน	11
2.6 แสดงค่า Horizontal Histogram เปรียบเทียบกับ 90% ของฮิสโตแกรมเฉลี่ย	12
2.7 การ Histogram ในแกน x	13
2.8 แสดงการทำ Character Block	13
2.9 แสดง Character Frame	14
3.1 แสดงภาพตัวอักษรขาด	15
3.2 แสดง Character Frame ที่ได้จากการประมวลผลภาพเบื้องต้น	15
3.3 แสดงการขาดที่เกิดขึ้นกับตัวอักษรภาษาไทย	16
3.4 ผังแสดงการทำงานโดยรวมการเชื่อมต่อตัวอักษรขาด	17
3.5 แสดงลักษณะการขาดของตัวอักษรในภาษาไทย	18
3.6 แสดงลักษณะการขาด 2 ส่วน โดยมีการเหลือมัลลาของ Character Frame	18
3.7 แสดงขอบเขตของตัวอักษรและ Character Frame ที่เหลือมัลลา	19
3.8 ผังแสดงการวิเคราะห์ตัวอักษรขาด 2 ส่วน โดยมีส่วนของ Character Frame เหลือมัลลากัน	20
3.9 แสดงการเชื่อมต่อจุดขาดของส่วนเหลือมัลลาครึ่งบน	20
3.10 แสดงการเชื่อมต่อจุดขาดของส่วนเหลือมัลลาครึ่งล่าง	21

สารบัญรูป (ต่อ)

รูปที่	หน้า
3.11 แสดงการใช้ฮิสโตแกรมในการหาตำแหน่งเชื่อม	21
3.12 แสดงผลการใช้ฮิสโตแกรมในการเชื่อม	22
3.13 แสดงการวิเคราะห์ตัวอักษรครึ่งหน้า	22
3.14 แสดงการวิเคราะห์จุดเชื่อมต่อโดยใช้ฮิสโตแกรม	23
3.15 แสดงการเชื่อมต่อที่เสร็จสมบูรณ์	23
3.16 แสดงการเชื่อมต่อตัวอักษรขาด 2 ส่วนโดยไม่มีส่วนที่เหลื่อมล้ำกัน	24
3.17 แสดงการขาดที่ส่วนกว้างอยู่ครึ่งบน	25
3.18 แสดงการเปลี่ยนพิกัดของ Seg[0].ymax ให้เป็น y_2 เพื่อให้เกิดการเหลื่อมล้ำ	25
3.19 แสดงการจัดการ Seg[0] กับ Seg[1] เพื่อทำการเชื่อม	26
3.20 แสดงการจัดการ Seg[0] กับ Seg[2] เพื่อทำการเชื่อม	26
3.21 แสดงการเชื่อมต่อที่เสร็จสมบูรณ์	26
3.22 แสดงส่วนที่กว้างอยู่ครึ่งล่าง	27
3.23 แสดงการเปลี่ยนพิกัดของ Seg[2].ymax ให้เป็น y_1 เพื่อให้เกิดการเหลื่อมล้ำ	27
3.24 แสดงการจัดการ Seg[2] กับ Seg[0] เพื่อทำการเชื่อมส่วนที่ขาด	28
3.25 แสดงการจัดการ Seg[2] กับ Seg[1] เพื่อทำการเชื่อมส่วนที่ขาด	28
3.26 แสดงการเชื่อมต่อที่เสร็จสมบูรณ์	28
3.27 แสดงส่วนที่ขาดทั้ง 4 ส่วน	29
3.28 แสดงการเปลี่ยนตำแหน่งของ Seg[0].ymax ให้เป็น y_2	30
3.29 แสดงการเปลี่ยนพิกัดของ Seg[0].xmin และ Seg[3].ymin	30
3.30 แสดงการเหลื่อมล้ำของ Seg[1] กับ Seg[3]	31
3.31 แสดงการเชื่อมต่อที่เสร็จสมบูรณ์	31
4.1 แสดงตัวอย่างข้อมูลภาพที่ใช้ในการทดสอบ	32
4.2 แสดงตัวอย่างข้อมูลภาพที่ใช้ในการทดสอบหลังการเชื่อมต่อ	32
5.1 แสดงการเชื่อมต่อที่ผิดพลาด.....	36
5.2 แสดงพิกัดของ Character Frame ส่วนบนและส่วนล่างของภาพตัวอักษร ‘ร’.....	37
5.3 แสดงแนวการเชื่อมต่อที่ผิดพลาด.....	37

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 1

บทนำ

1.1 ความเป็นมา และความสำคัญของปัญหา

ในระบบการรู้จำตัวอักษร ความถูกต้องของการรู้จำจะขึ้นกับข้อมูลภาพตัวอักษรซึ่งจะต้องครบถ้วน ซึ่งในงานวิจัยที่ผ่านมาจะมีการกล่าวถึงเฉพาะ การแก้ไขตัวอักษรที่ติดกัน[4][5][6] การลดสัญญาณรบกวน[6] การทำขอบตัวอักษรให้เรียบ[7] แต่ไม่ได้มีการกล่าวถึงการแก้ไขข้อมูลที่ไม่สมบูรณ์ก่อนที่จะนำไปสู่กระบวนการรู้จำตัวอักษร โดยเฉพาะตัวอักษรขาดแสดงไว้ดังรูปที่ 1.1 ดังนั้นในงานวิจัยนี้จึงมีจุดประสงค์หลักเพื่อทำการปรับปรุงซ่อมแซมข้อมูลภาพตัวอักษรที่ขาด ซึ่งอาจเกิดขึ้นเนื่องจากการสแกน หรือเกิดจากเอกสารคุณภาพต่ำ เป็นผลให้กระบวนการรู้จำเกิดความผิดพลาดหรือไม่สามารถรู้จำได้เลย

ที่มาของเอกสาร

รูปที่ 1.1 แสดงข้อมูลภาพตัวอักษรขาด

จากการทดสอบซอฟต์แวร์ OCR สำหรับภาษาไทยที่มีขายในขณะนี้ คือ AmThai และ ThaiOCR พบว่าซอฟต์แวร์ทั้งสองไม่สามารถรู้จำข้อมูลตัวอักษรขาดได้ ตัวอย่างข้อมูลภาพในการทดสอบแสดงไว้ดังรูปที่ 1.2 ผลการรู้จำของซอฟต์แวร์แสดงไว้ในตารางที่ 1.1

เป้าหมายของโครงการ

ที่มาของเอกสาร

รูปที่ 1.2 แสดงตัวอย่างข้อมูลภาพที่ใช้ในการทดสอบ

ตารางที่ 1.1 แสดงผลการทดสอบโปรแกรม OCR ในการรู้จำตัวอักษรขาด

ข้อความทดสอบ	ซอฟต์แวร์ทดสอบ	
	ThaiOCR	ArnThai
เป้าหมายของโครงการ	๘ 6บ๑'Aมาปี่ ซึ่ N ค.ฐึNNนี้ศึ	บจจมาयरข๑จก๑๑ขจขจจ
ที่มาของเอกสาร	ที่	ฐึมา.คขเจจ.ขกส๑

งานวิจัยนี้ได้นำเสนอแนวทางเพื่อแก้ปัญหาภาพตัวอักษรขาดอันเป็นสาเหตุหนึ่งที่ทำให้กระบวนการรู้จำตัวอักษรผิดพลาดหรือไม่สามารถรู้จำได้ ซึ่งจะทำให้ระบบการรู้จำตัวอักษรภาษาไทยมีความสมบูรณ์และมีประสิทธิภาพเพิ่มขึ้น

1.2 วัตถุประสงค์ของงานวิจัย

1. เพื่อศึกษาวิธีการแยกข้อมูลภาพตัวอักษรออกจากเอกสาร
2. เพื่อศึกษาวิธีการในการระบุภาพตัวอักษรว่าเป็นตัวอักษรขาดหรือไม่
3. เพื่อศึกษาโครงสร้างของตัวอักษรภาษาไทย ในการแบ่งประเภทการขาดของตัวอักษร
4. เพื่อศึกษาแนวทางการซ่อมแซมตัวอักษรขาดในแต่ละประเภท
5. เพื่อพัฒนาระบบการรู้จำตัวอักษรภาษาไทยมีความสมบูรณ์และมีประสิทธิภาพเพิ่มขึ้น

1.3 ทฤษฎีและหลักการที่เกี่ยวข้อง

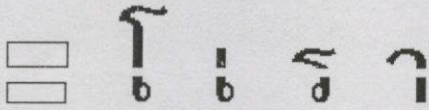
1.3.1 นิยามตัวอักษรขาด

ตัวอักษรขาด หมายถึงภาพตัวอักษรที่ถูกแยกออกเป็นชิ้นส่วนภาพย่อยๆ โดยขอบเขตของงานวิจัยนี้มีดังนี้

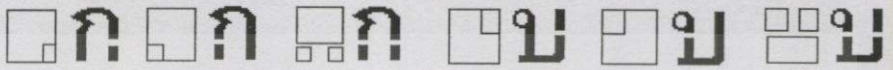
- ไม่ขึ้นกับฟอนต์ของตัวอักษร (Multi-font)
- ใช้กับตัวอักษรตัวปรกติ (ตัวตรง)
- เป็นการขาดในแนวนอน
- วิเคราะห์เฉพาะตัวอักษรใน Central Zone ซึ่งประกอบด้วยพยัญชนะไทย 44 ตัว และสระ 6 ตัว คือ ๑, ๒, ๓, ๔, ๕ และ ๖
- รูปแบบและตำแหน่งของการขาดคือขาดในส่วนที่เป็นขาของตัวอักษร

จากนิยามดังกล่าวทำให้พบลักษณะการขาดของตัวอักษรพิมพ์ภาษาไทยมีลักษณะดังนี้

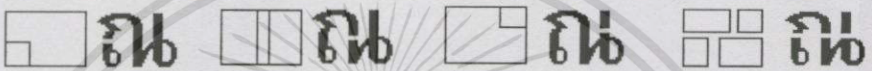
- ตัวอักษรที่มีขาเดียว จะมีลักษณะของกรอบเป็นดังนี้



- ตัวอักษรที่มีสองขา จะมีลักษณะของกรอบเป็นดังนี้

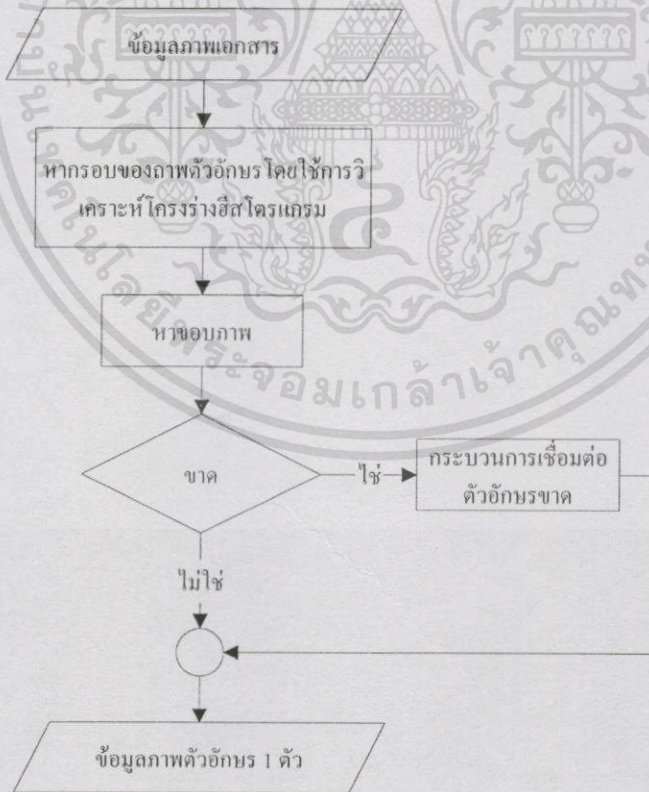


- ตัวอักษรที่มีสามขา จะมีลักษณะกรอบภาพเป็นดังนี้



1.3.2 กระบวนการการซ่อมแซมอักษรตัวอักษรขาด

กระบวนการการซ่อมแซมอักษรตัวอักษรขาดแสดงได้ดัง รูปที่ 1.3



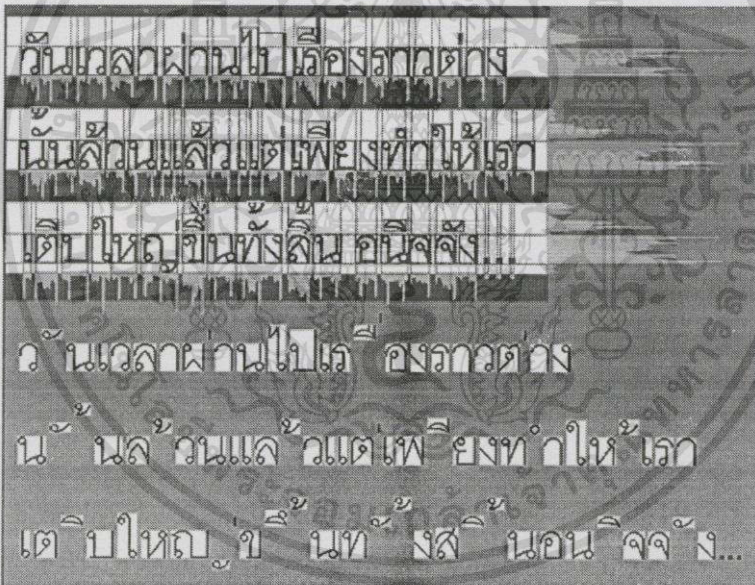
รูปที่ 1.3 แสดงกระบวนการซ่อมแซมตัวอักษรขาด

ทฤษฎีและหลักการที่ใช้ในงานวิจัยจะแบ่งเป็น 3 ส่วนหลักๆ คือ

1. การประมวลผลภาพเบื้องต้น
2. การวิเคราะห์การขาดของตัวอักษร
3. การเชื่อมต่อตัวอักษรที่ขาด

1) การประมวลผลภาพเบื้องต้น

ภาพเอกสารที่เป็นอินพุตของระบบจะประกอบด้วยหลายบรรทัด ดังนั้นในขั้นตอนนี้ จะทำการแยกบรรทัดและแยกตัวอักษรแต่ละตัวออกจากบรรทัดโดยใช้ฮิสโตแกรม (Histogram) [8] และการหาขอบภาพ (Contour Algorithm) [9] ซึ่งในขั้นตอนนี้จะได้ ตำแหน่งและขนาดของภาพตัวอักษรแต่ละตัว เรียกว่า กรอบตัวอักษร แล้วจัดเรียงลำดับ ภาพตัวอักษร แสดงดังรูปที่ 1.4

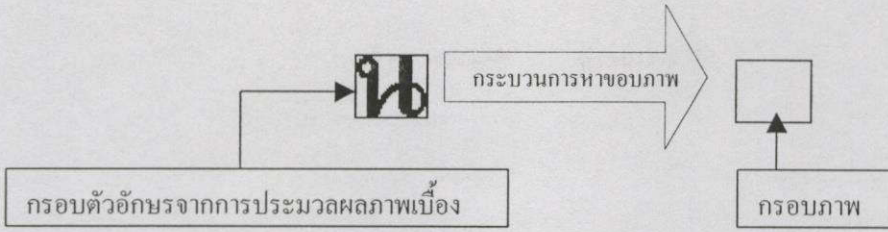


รูปที่ 1.4 แสดงการแยกบรรทัดและตัวอักษร โดยใช้ฮิสโตแกรม และการหาขอบภาพ

2) การวิเคราะห์การขาดของตัวอักษร

ขั้นตอนนี้จะนำกรอบตัวอักษรจากการประมวลผลภาพเบื้องต้นมาผ่านกระบวนการหาขอบภาพ หลักเกณฑ์ในการพิจารณาว่าภาพตัวอักษรขาดหรือไม่นั้นจะใช้วิธีการนับจำนวน กรอบภาพที่พบภายในกรอบตัวอักษร ถ้าในกรอบตัวอักษรนั้นพบกรอบภาพเพียงหนึ่ง กรอบแสดงว่าในกรอบตัวอักษรไม่พบตัวอักษรขาด ในทางตรงข้ามหากพบกรอบภาพมากกว่า 1 กรอบ แสดงว่าพบตัวอักษรขาดในกรอบตัวอักษร แสดงดังรูปที่ 1.5

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับกรใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



ก) แสดงกรอบภาพที่ได้จากตัวอักษรที่ไม่ขาดซึ่งจะพบเพียง 1 กรอบ

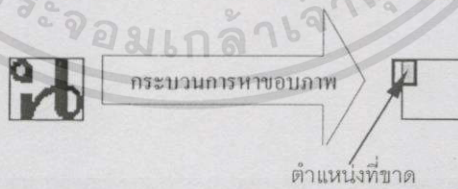


ข) แสดงกรอบภาพ 2 กรอบเนื่องจากตัวอักษรขาด

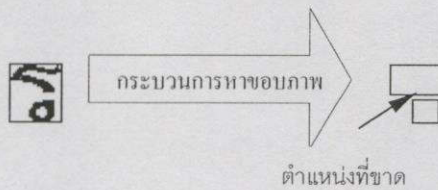
รูปที่ 1.5 แสดงกรอบภาพจากการหาขอบภาพ

ผลลัพธ์ที่ได้จากการหาขอบภาพคือจำนวนกรอบภาพและขอบเขตของกรอบภาพภายในกรอบตัวอักษรทำให้สามารถระบุตำแหน่งที่ขาดของตัวอักษรได้ ลักษณะการขาดของตัวอักษรในภาษาไทยสามารถแบ่งเป็นกลุ่มหลักๆ ได้เป็น

- 1) การขาดโดยมีการเหลื่อมล้ำกันของกรอบภาพ ดังรูปที่ 1.6 (ก)
- 2) การขาดโดยที่ไม่มีการเหลื่อมล้ำกันของกรอบภาพ ดังรูปที่ 1.6 (ข)



ก) แสดงลักษณะการขาดที่กรอบภาพเหลื่อมล้ำกัน



ข) แสดงลักษณะการขาดที่ไม่มีการเหลื่อมล้ำกันของกรอบภาพ

รูปที่ 1.6 แสดงตัวอักษรขาดทั้ง 2 แบบ

3) กระบวนการเชื่อมตัวอักษรขาด

ในการเชื่อมต่อจะนำข้อมูลจากการวิเคราะห์การขาดของตัวอักษรมาพิจารณาการเชื่อมต่อ เนื่องจากการขาดของตัวอักษรจะพบที่ตำแหน่งที่เป็นขาของตัวอักษรซึ่งมีลักษณะเป็นแท่ง ดังนั้นการเชื่อมต่อตัวอักษรจะเป็นการลากเส้นเชื่อมส่วนที่ขาดในแนวตั้ง ซึ่งจากลักษณะการขาดของตัวอักษรทั้งสองแบบดังที่กล่าวมาในการวิเคราะห์การขาดของตัวอักษรจะทำให้เกิดการเชื่อมมี 2 แบบเช่นกันคือ

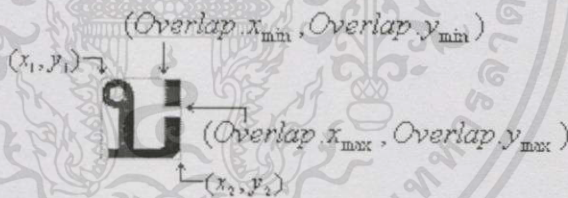
- 1) การเชื่อมตัวอักษรขาดโดยมีการเหลื่อมล้ำกันของกรอบภาพ
- 2) การเชื่อมตัวอักษรขาดโดยไม่มีการเหลื่อมล้ำกันของกรอบภาพ

3.1) การเชื่อมตัวอักษรขาดโดยมีการเหลื่อมล้ำกันของกรอบภาพ

กำหนดพิกัดกรอบตัวอักษรที่ได้จากการประมวลผลภาพเบื้องต้นเป็น (x_1, y_1) และ (x_2, y_2) ส่วนที่เหลื่อมล้ำที่ได้จากการวิเคราะห์การขาดของตัวอักษรกำหนดเป็น

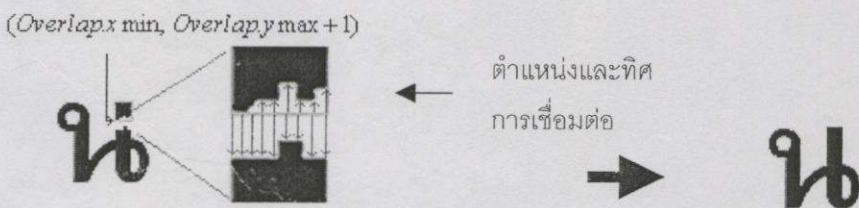
$$(Overlap.x_{min}, Overlap.y_{min}) \text{ และ } (Overlap.x_{max}, Overlap.y_{max})$$

ดังรูปที่ 1.7



รูปที่ 1.7 แสดงขอบเขตของตัวอักษรและส่วนที่เหลื่อมล้ำ

การเชื่อมต่อจะยึดถือส่วนที่เหลื่อมล้ำเป็นหลักกล่าวคือ ถ้าส่วนที่เหลื่อมล้ำอยู่ครึ่งบนก็จะทำการเชื่อมเส้นจากจุด *Overlap* จนกระทั่งพบจุดสีดำซึ่งเป็นเนื้อของตัวอักษรแสดงได้ดังรูปที่ 1.8



รูปที่ 1.8 แสดงการเชื่อมต่อจุดขาด

3.2) การเชื่อมตัวอักษรขาดโดยไม่มีกรล้อมล้ากันของกรอบภาพ

สำหรับการขาดของตัวอักษรที่ไม่มีส่วนที่เชื่อมกันนั้นจะใช้การพิจารณาจากค่าฮิสโตแกรมของข้อมูลภาพตัวอักษร โดยตำแหน่งที่พบค่าฮิสโตแกรมที่สูงที่สุดจะเป็นตำแหน่งขาของตัวอักษร หรือเป็นตำแหน่งที่จะเชื่อม ดังรูปที่ 1.9



รูปที่ 1.9 แสดงการใช้ฮิสโตแกรมหาตำแหน่งการเชื่อม

1.4 แผนการดำเนินงาน

1. ศึกษาบทความและผลงานวิจัยต่างๆ ที่มีความเกี่ยวข้องกับงานวิจัยนี้
2. เก็บข้อมูลตัวอย่างของตัวอักษรขาด พร้อมจัดเก็บลงคอมพิวเตอร์
3. ศึกษาลักษณะ โครงสร้างของตัวอักษร ไทยเพื่อนำไปวิเคราะห์ลักษณะและตำแหน่งการขาด
4. ออกแบบอัลกอริทึมในการวิเคราะห์ภาพตัวอักษรขาด และทำการเชื่อมต่อตัวอักษรขาด
5. เขียนโปรแกรมเพื่อวิเคราะห์ภาพตัวอักษรขาด และทำการเชื่อมต่อตัวอักษรขาด
6. ทดลองเชื่อมต่อตัวอักษรขาดกับข้อมูลที่จัดเก็บ
7. ทดสอบผลจากการเชื่อมต่อว่าสามารถนำไปใช้งานได้จริงกับซอฟต์แวร์ ThaiOCR และ ArnThai โดยทดสอบข้อมูลก่อนและหลังการเชื่อมต่อ
8. สรุปผลการดำเนินการ และรวบรวมนำจัดทำเอกสารนำเสนอเป็นงานวิจัย

2.2.1 การแยกบรรทัด (Line Separation)

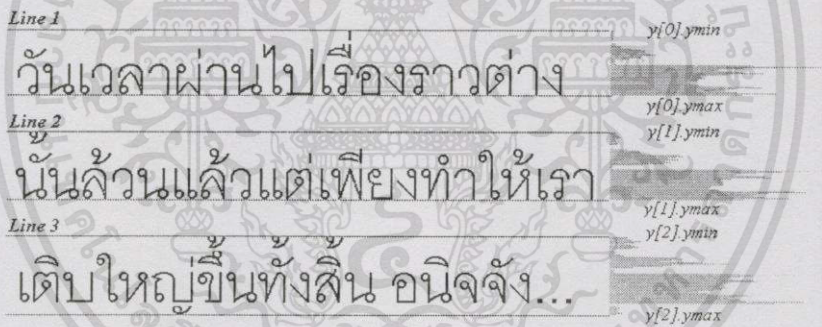
เป็นการนำทฤษฎี Histogram มาช่วยในการหาระดับของภาพตัวอักษรในแต่ละบรรทัด เพื่อจะสามารถนำไปวิเคราะห์หาบรรทัดที่ปรากฏตัวอักษรในหน้าเอกสาร โดยใช้ Histogram ทางด้านแกน y หรือ Horizontal Histogram เพื่อหาหาบรรทัดของหน้าเอกสาร แล้วทำการเก็บจุดเริ่มต้น และจุดสิ้นสุดในแนวแกน y ของบรรทัด แล้วทำการเก็บเป็นลิสต์ของการแบ่งที่เกิดขึ้นทั้งหมดเพื่อนำไปวิเคราะห์หาบรรทัด

การแบ่งบรรทัดจะพิจารณาจากการ หา Horizontal Histogram ดังสมการ 2.1

$$yHis(y) = \sum_{x=0}^{x=X \max} P(x, y)$$

เมื่อ $P(x,y)$ เป็นจุดของภาพ และ $Xmax$ เป็นความกว้างของภาพ

รูปที่ 2.4 แสดงฮิสโตแกรมและการแบ่งบรรทัดโดยการวิเคราะห์ค่าฮิสโตแกรมในแนวนอน



รูปที่ 2.4 การ Histogram ในแกน y

จากรูปที่ 2.4 จะพบว่าส่วนที่อยู่ระหว่างบรรทัดจะเป็นส่วนที่มีค่าฮิสโตแกรมเป็นศูนย์ และมีระยะที่กว้าง ซึ่งเงื่อนไขนี้จะทำให้สามารถแบ่งบรรทัดออกมาได้ ดังภาพตัวอย่างรูปที่ 2.4 จะได้ 3 บรรทัดแต่ละบรรทัดก็จะได้ของเขตของตำแหน่ง y เริ่มต้น $ymin$ และตำแหน่ง y สุดท้าย $ymax$ ของบรรทัด ซึ่งจะถูกจัดเก็บในรูปแบบของอาร์เรย์

เมื่อได้ของเขตของแต่ละบรรทัดแล้วจะนำขอบเขตของแต่ละบรรทัดไปทำการแบ่งระดับของตัวอักษรในแต่ละบรรทัดโดยยึดถือตามโมเสกโครงสร้างของตัวอักษรภาษาไทย ซึ่งจะกล่าวในหัวข้อถัดไป

2.2.2 การหาระดับของตัวอักษร (Level Separation)

ขั้นตอนนี้จะทำการวิเคราะห์ค่า Horizontal Histogram ที่ได้จากขั้นตอนก่อนหน้า เพื่อทำการแบ่งระดับของภาพตัวอักษรในแต่ละบรรทัด ซึ่งระดับของตัวอักษรจะแบ่งออกเป็น 3 ส่วนคือ

- Upper Zone ประกอบด้วย ตัวอักษรในระดับ Tonal Line Level และ Upper Line Level
- Central Zone ประกอบด้วย ตัวอักษรในระดับ Consonant Line Level
- Lower Zone ประกอบด้วย ตัวอักษรในระดับ Lower Vowel Line Level

ตัวอย่างระดับของตัวอักษรในหนึ่งบรรทัดแสดงดังรูปที่ 2.5



รูปที่ 2.5 แสดงระดับของตัวอักษรทั้ง 3 ส่วน

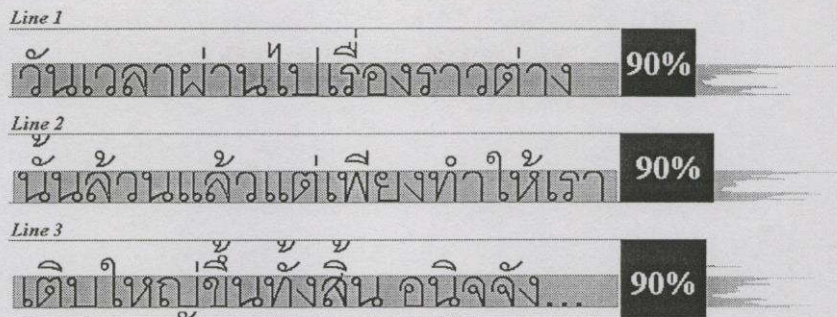
การแบ่งระดับจะเริ่มด้วยการหาค่าเฉลี่ยของค่า Horizontal Histogram ที่ได้ในแต่ละบรรทัดซึ่งแสดงได้ดังสมการที่ 2.2

$$Average\ Hist\ (line) = \frac{\sum_{i=y[line].y\ min}^{y[line].y\ max} yHis(i)}{y[line].y\ max - y[line].y\ min}$$

เมื่อ $yHist$ เป็นค่าฮิสโตแกรมตามแนวนอน

$y[line].ymin$ และ $y[line].ymax$ เป็นพิกัดขอบบนและขอบล่างของบรรทัด

หลักในการแบ่งระดับในแต่ละบรรทัดคือตำแหน่งที่มีค่า Horizontal Histogram มากกว่า 90 เปอร์เซ็นต์ของค่าเฉลี่ยฮิสโตแกรมจะถือว่าเป็นส่วน Central Zone รูปที่ 2.6 แสดงค่า Horizontal Histogram ค่า 90 เปอร์เซ็นต์ของค่าเฉลี่ยฮิสโตแกรมในแต่ละบรรทัด



รูปที่ 2.6 แสดงค่า Horizontal Histogram เปรียบเทียบกับ 90% ของฮิสโตแกรมเฉลี่ย

จากรูปที่ 2.6 จะพบว่ากลุ่มหรือช่วงที่มีค่าฮิสโตแกรมมากกว่า 90 เปอร์เซ็นต์ของค่าเฉลี่ยฮิสโตแกรม (จากรูปคือแถบสีดำ)จะเป็น Central Zone ในกรณีบรรทัดที่ 3 จะพบว่า มีฮิสโตแกรมบางส่วนที่มากกว่า 90 เปอร์เซ็นต์ของค่าเฉลี่ยฮิสโตแกรมแต่จะไม่นำมาเป็น Central Zone เนื่องจากในแต่ละบรรทัดจะต้องมี Central Zone แค่ 1 ส่วนเท่านั้น ดังนั้นจะเลือกช่วงของฮิสโตแกรมที่มีค่ากว้างที่สุดให้เป็นส่วน Central Zone ดังนั้นส่วนที่อยู่เหนือ Central Zone ก็จะกลายเป็น Upper Zone และส่วนที่อยู่ต่ำกว่า Central Zone ก็จะเป็น Lower Zone

ระดับของตัวอักษรจะมีประโยชน์ในการระบุระดับของภาพตัวอักษรที่ได้จากการหาขอบเขตในขั้นตอน Character Segmentation ในขั้นตอนถัดไปและมีประโยชน์ในการจัดเรียงตัวอักษรด้วย

สำหรับงานวิจัยนี้ จะวิเคราะห์การขาดของตัวอักษรเฉพาะตัวอักษรใน Consonant Level หรือส่วนที่เป็น Central Zone ซึ่งประกอบด้วยพยัญชนะไทย 44 ตัวและสระ 6 ตัว คือ ๑, ๒, ๓, ๔ และ ๕

2.2.3 การหาตำแหน่งและขนาดของภาพตัวอักษร (Character Segmentation)

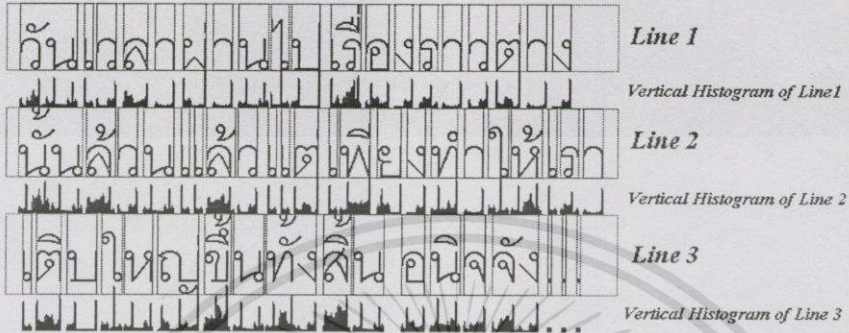
ขั้นตอนนี้จะทำการหาตำแหน่งและขอบเขตของตัวอักษรแต่ละตัวพร้อมระบุระดับของตัวอักษรว่าอยู่ในระดับใดใน 3 ส่วนคือ Upper Zone, Central Zone และ Lower Zone จากนั้นจะทำการจัดเก็บลงคลังคลังพร้อมจัดเรียงลำดับให้ถูกหลักตามลักษณะการพิมพ์ภาษาไทย

เพื่อให้การจัดเรียงลำดับทำได้ง่ายขึ้นในขั้นตอนแรกจะทำการแบ่งภาพตัวอักษรแบบหยายๆก่อนโดยใช้ Vertical Histogram แสดงได้ดังสมการ 2.3

$$xHis(x) = \sum_{y=0}^{y=Y \max} P(x, y)$$

เมื่อ $P(x,y)$ เป็นจุดภาพ และ $Ymax$ เป็นความสูงของบรรทัดที่แบ่งได้

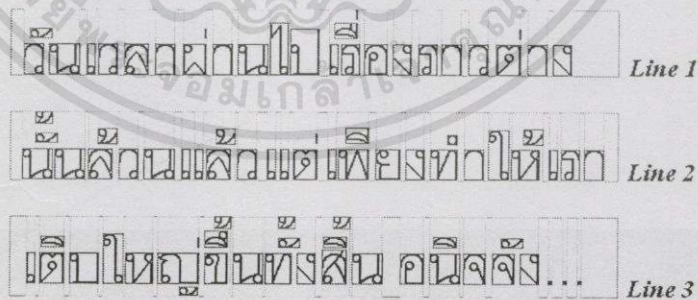
ตำแหน่งที่มีค่า Vertical Histogram มากกว่าศูนย์จะเป็นตำแหน่งที่จะตัดแบ่ง ซึ่งผลของการตัดแบ่งจะอาจจะได้ตัวอักษรเพียงตัวเดียวหรือเป็นกลุ่มของตัวอักษรก็ได้ ขอบเขตที่ตัดแบ่งได้ในขั้นตอนนี้จะเรียกว่า Character Block แสดงดังรูปที่ 2.7



รูปที่ 2.7 การ Histogram ในแกน x

ขั้นตอนต่อไปจะทำการไต่ขอบ (tracing contour) ของภาพตัวอักษรที่อยู่ภายใน Character Block เพื่อหาขอบเขตของตัวอักษรแต่ละตัว แล้วทำการคัดลอกภาพตัวอักษรที่ไต่ขอบแล้วไปใส่ในโครงสร้างข้อมูลแบบลิงค์ลิสต์พร้อมจัดเรียงตามลักษณะการพิมพ์ตัวอักษรในภาษาไทย ซึ่งภาพตัวอักษรแต่ละตัวที่คัดลอกออกมาจะเรียกว่า Character Frame

ผลที่ได้จากการไต่ขอบแสดงดังรูปที่ 2.8 และ Character Frame ที่เก็บในลิงค์ลิสต์แสดงดังรูปที่ 2.9



รูปที่ 2.8 แสดงการหา Character Block

ด้^๑นเวล^๑ก^๑ผ^๑ก^๑น^๑ไป^๑เร^๑็^๑อง^๑ร^๑า^๑ว^๑ต^๑า^๑ง
 น^๒้^๒น^๒ล^๒ว^๒น^๒แ^๒ล^๒ว^๒แต่^๒เพ^๒็^๒ย^๒ง^๒ท^๒่า^๒ใ^๒ห^๒้^๒เร^๒า
 เ^๓็^๓บ^๓ใ^๓ห^๓เ^๓ญ^๓ว^๓ข^๓ี^๓น^๓ท^๓ี่^๓ง^๓ส^๓ี^๓น^๓อ^๓น^๓ี^๓จ^๓จ^๓ัง...

รูปที่ 2.9 แสดง Character Frame

ผลลัพธ์สุดท้ายที่ได้คือ Character Frame จะนำไปใช้ในโครงการวิจัยเพื่อวิเคราะห์ว่า
 ภาพตัวอักษรที่อยู่ใน Character Frame ขาดหรือไม่ ซึ่งจะกล่าวรายละเอียดในบทถัดไป



บทที่ 3

การซ่อมแซมอักษรตัวพิมพ์ภาษาไทยที่ขาด

3.1 นิยามตัวอักษรขาด

ตัวอักษรขาด หมายถึงภาพตัวอักษรที่ถูกแยกออกเป็นชิ้นส่วนภาพย่อยๆซึ่งในงานวิจัยนี้จะเน้นเฉพาะการขาดที่ส่วนขาของตัวอักษรเท่านั้น ตัวอย่างของตัวอักษรขาดแสดงในรูปที่ 3.1

เป้าหมายของโครงการ ที่มาของเอกสาร

รูปที่ 3.1 แสดงภาพตัวอักษรขาด

เมื่อนำภาพจากรูปที่ 3.1 ที่มีตัวอักษรขาดผ่านกระบวนการประมวลผลภาพเบื้องต้น (กล่าวในบทที่ 2) จะได้ผลลัพธ์ที่เป็น Character Frame ดังรูปที่ 3.2



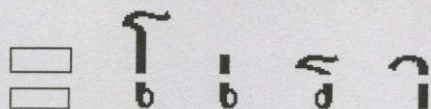
รูปที่ 3.2 แสดง Character Frame ที่ได้จากการประมวลผลภาพเบื้องต้น

Character Frame ทั้งหมดจะถูกนำมาตรวจสอบว่าเป็นตัวอักษรขาดหรือไม่และถ้าขาดจะมีลักษณะการขาดรูปแบบใด โดยใช้หลักเกณฑ์คือตรวจสอบการเหลื่อมล้ำกันของ Character Frame ซึ่งจะใช้พิกัดจริงที่ภาพตัวอักษรปรากฏอยู่บนเอกสาร โดยพิกัดจริงจะได้ในกระบวนการประมวลผลภาพเบื้องต้นแล้ว

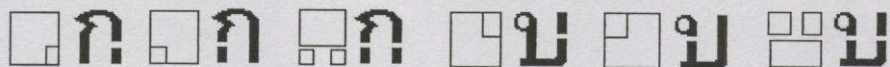
จากการศึกษาพบว่าลักษณะการขาดของตัวอักษรแบ่งเป็น 2 ประเภทคือ การขาดโดยมีการเหลื่อมล้ำกันของ Character Frame และการขาดโดยที่ไม่มีมีการเหลื่อมล้ำกันของ Character Frame และเมื่อนำลักษณะการขาดทั้ง 2 แบบมาพิจารณาร่วมกับลักษณะตัวอักษรในภาษาไทยคือ จำนวนขา จะพบว่าลักษณะการขาดของตัวอักษรพิมพ์ภาษาไทยมีลักษณะดังนี้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่นอนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างถึงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

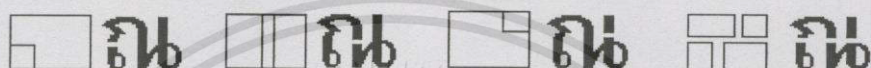
- ตัวอักษรที่มีขาเดียว จะมีลักษณะของ Character Frame เป็นดังนี้



- ตัวอักษรที่มีสองขา จะมีลักษณะของ Character Frame เป็นดังนี้



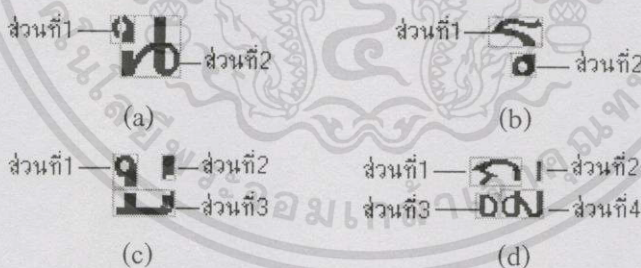
- ตัวอักษรที่มีสามขา จะมีลักษณะ Character Frame เป็นดังนี้



3.2 การวิเคราะห์ตัวอักษรขาด

ลักษณะการขาดของตัวอักษรตามที่กล่าวมาข้างต้นสามารถก่อให้เกิดการขาดของตัวอักษรภาษาไทยดังนี้

- ขาด 2 ส่วน ดังรูปที่ 3.3(a) และ 3.3(b)
- ขาด 3 ส่วน ดังรูปที่ 3.3(c)
- ขาด 4 ส่วน ดังรูปที่ 3.3(d)

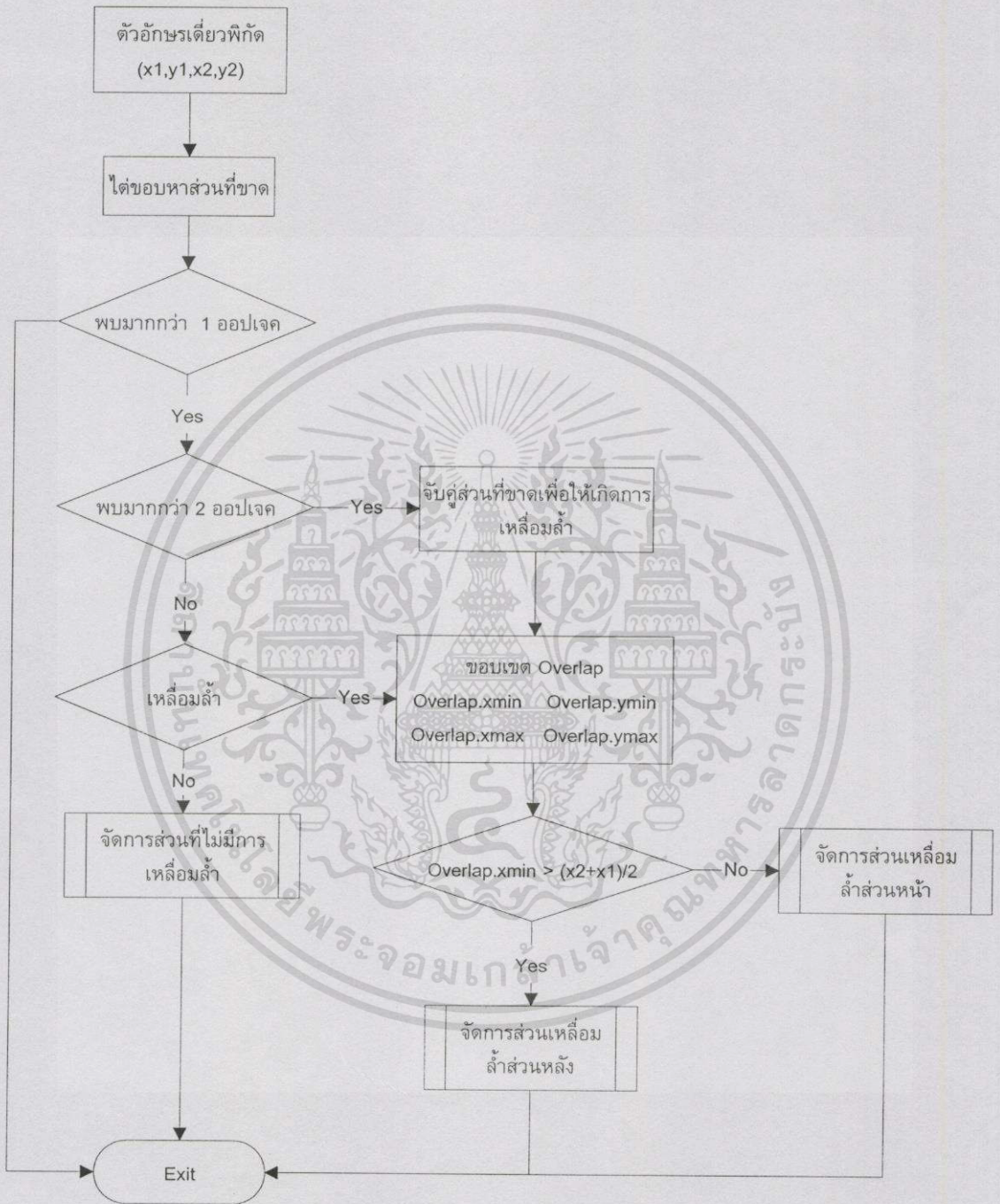


รูปที่ 3.3 แสดงการขาดที่จะเกิดขึ้นกับตัวอักษรภาษาไทย

ในส่วนการออกแบบงานวิจัยจะแบ่งการวิจัยออกเป็น

1. ตัวอักษรขาด 2 ส่วน โดยมี Character Frame เหลือมล้ากัน
2. ตัวอักษรขาด 2 ส่วน โดยไม่มี Character Frame เหลือมล้ากัน
3. ตัวอักษรขาด 3 ส่วน
4. ตัวอักษรขาด 4 ส่วน

การทำงาน โดยรวมสามารถแสดงได้ดังรูปที่ 3.4

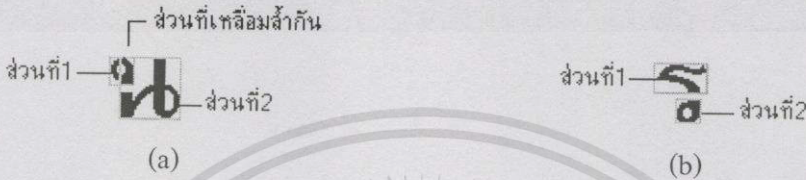


รูปที่ 3.4 ผังแสดงการทำงาน โดยรวมการเชื่อมต่อตัวอักษรขาด

3.3 ลักษณะตัวอักษรขาด

ลักษณะการขาดของตัวอักษรในภาษาไทยที่พบจะสามารถแบ่งเป็นกลุ่มหลักๆ ได้เป็น

- ขาด โดยมี Character Frame เหลื่อมล้ำกัน ดังรูปที่ 3.5(a)
- ขาด โดยที่ไม่มี Character Frame เหลื่อมล้ำกันเลย ซึ่งจะเกิดกับตัวอักษรที่มีขาเดียว ดังรูปที่ 3.5(b)



รูปที่ 3.5 แสดงลักษณะการขาดของตัวอักษรในภาษาไทย

จากลักษณะของการเหลื่อมล้ำดังกล่าวอาจกล่าวได้ว่า ส่วนที่จะต่อหรือส่วนที่ขาดนั้นคือส่วนที่เหลื่อมล้ำกัน ในส่วนของลักษณะตัวอักษรที่ไม่มีส่วนที่เหลื่อมล้ำกันจะพิจารณาจากส่วนล่าง (จากรูปที่ 2 คือส่วนที่2) เป็นหลักเนื่องจากลักษณะของตัวอักษรภาษาไทยจะเริ่มต้นจากส่วนล่างขึ้นไปด้านบน

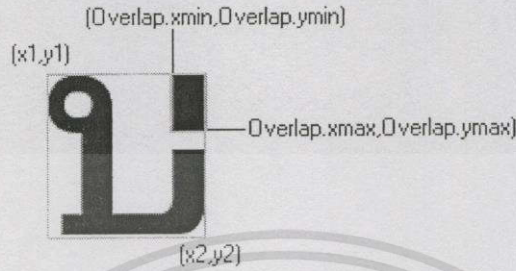
3.3.1 ตัวอักษรขาด 2 ส่วน โดยมี Character Frame เหลื่อมล้ำกัน

ลักษณะนี้จะนำส่วนที่เกิดการเหลื่อมล้ำมาพิจารณา โดยแบ่งการพิจารณาส่วนที่เหลื่อมล้ำเป็นเหลื่อมล้ำในส่วนครึ่งหน้าของตัวอักษรดังรูปที่ 3.6(a) และเหลื่อมล้ำในส่วนครึ่งหลังของตัวอักษรรูปที่ 3.6(b) และครึ่งบนกับครึ่งล่าง



รูปที่ 3.6 แสดงลักษณะการขาด 2 ส่วน โดยมีการเหลื่อมล้ำของ Character Frame

กำหนดพิกัดขอบเขตของตัวอักษรเป็นตัวอักษร (x_1, y_1) และ (x_2, y_2) ส่วนของ Character Frame ที่เหลื่อมล้ำกำหนดเป็น $(Overlap.x_{min}, Overlap.y_{min})$ และ $(Overlap.x_{max}, Overlap.y_{max})$ ดังรูปที่ 3.7



รูปที่ 3.7 แสดงขอบเขตของตัวอักษรและ Character Frame ที่เหลื่อมล้ำ

ส่วนของตัวอักษรที่เหลื่อมล้ำกันจะถือว่าอยู่ส่วนหลังถ้าเป็นไปตามเงื่อนไข

$$Overlap.x_{min} > (x_2 + x_1) / 2$$

และจะถือว่าอยู่ส่วนหลังถ้าเป็นไปตามเงื่อนไข

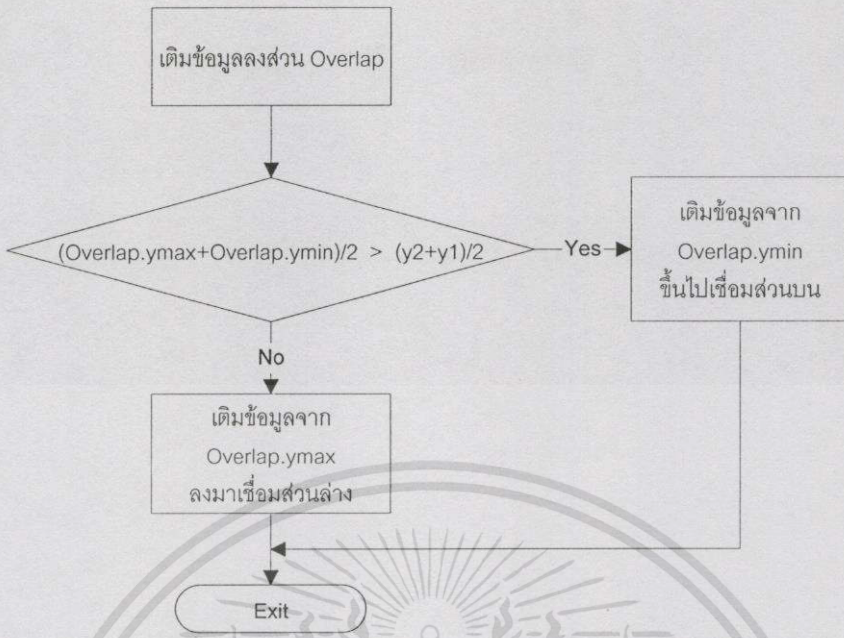
$$(Overlap.x_{max} + Overlap.x_{min}) / 2 > (x_2 + x_1) / 2$$

เงื่อนไขการแบ่งครึ่งบนคือ

$$(Overlap.y_{max} + Overlap.y_{min}) / 2 > (y_2 + y_1) / 2$$

ถ้านอกเหนือจากนี้จะถือว่าส่วนที่เหลื่อมล้ำอยู่ครึ่งล่าง

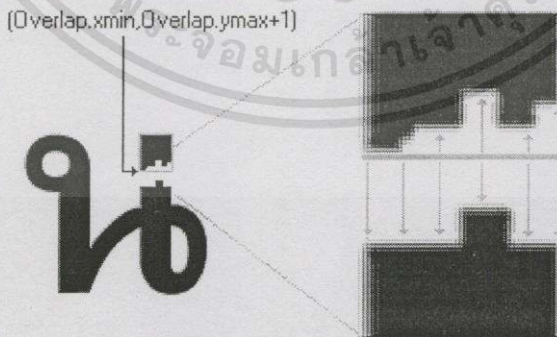
การทำงานโดยรวมของการวิเคราะห์ตัวอักษรขาด 2 ส่วน โดยมีส่วนของ Character Frame เหลื่อมล้ำกันแสดงได้ดังรูปที่ 3.8



รูปที่ 3.8 แสดงการวิเคราะห์ตัวอักษรขนาด 2 ส่วนโดยมีส่วนของ Character Frame เหลื่อมล้ำกัน

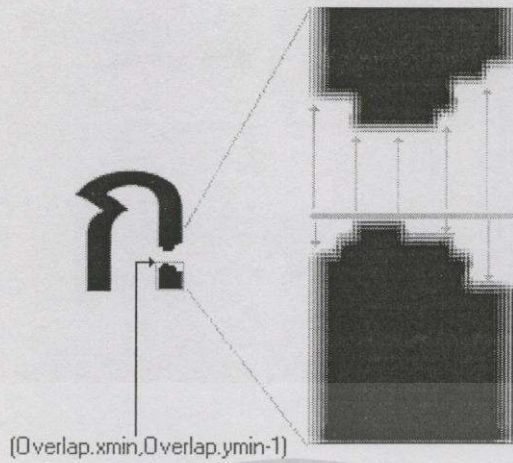
3.3.2 การวิเคราะห์ตัวอักษรครึ่งหลัง

ในการเชื่อมต่อจะยึดถือส่วนที่เหลื่อมล้ำเป็นหลักกล่าวคือ ถ้าส่วนที่เหลื่อมล้ำอยู่ครึ่งบนก็จะทำการเชื่อมเส้นจากจุด $(Overlap.x\ min, Overlap.y\ max+1)$ ลงมาจนกระทั่งพบจุดสีดำ ซึ่งเป็นเนื้อของตัวอักษรและในขณะเดียวกันก็ทำการเชื่อมจากจุดเดียวกันนี้ขึ้นไปด้านบนด้วย เนื่องจากการขาดของตัวอักษรปลายจะไม่เรียงดังแสดงได้ดังรูปที่ 3.9



รูปที่ 3.9 แสดงการเชื่อมต่อจุดขาดของส่วนเหลื่อมล้ำครึ่งบน

ในทางกลับกันถ้าส่วนที่เป็น Character Frame เหลื่อมล้ำอยู่ด้านล่างจะเริ่มต้นเชื่อมจากจุด $(Overlap.x\ min, Overlap.y\ min-1)$ เป็นหลักในการเริ่มต้นเชื่อมเส้นขึ้นไปครึ่งบน ดังรูปที่ 3.10



รูปที่ 3.10 แสดงการเชื่อมต่อจุดขาดของส่วนเหลื่อมล้ำครึ่งล่าง

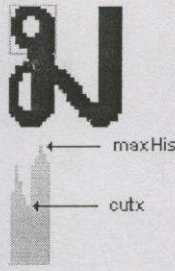
3.3.3 การวิเคราะห์ตัวอักษรครึ่งหน้า

ในส่วนของตัวอักษรที่อยู่ครึ่งหน้าจะพบว่าส่วนที่เหลื่อมล้ำมีโอกาสที่จะเกิดเป็นส่วนหัวของตัวอักษรได้ ดังนั้นในการเชื่อมจะต้องทำการวิเคราะห์โดยใช้กระบวนการฮีสโตรแกรมในการวิเคราะห์ว่าตำแหน่งใดควรจะเชื่อม (ตำแหน่งที่เป็นขาของตัวอักษร) แสดงได้ดังรูปที่ 3.11



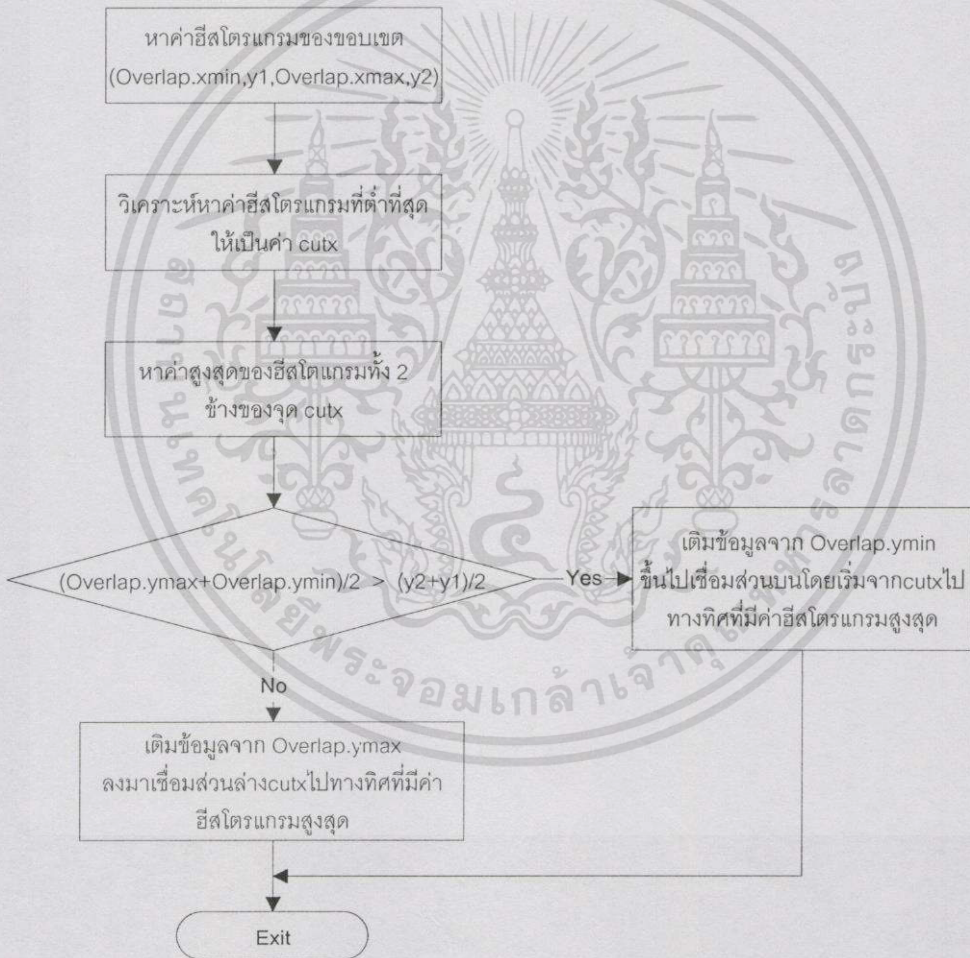
รูปที่ 3.11 แสดงการใช้ฮีสโตรแกรมในการหาตำแหน่งเชื่อม

จากรูปแสดงจุดสูงสุดและต่ำสุดของฮีสโตรแกรม จากตัวอย่างนี้จะสรุปว่าจุดที่สูงที่สุดจะเป็นตำแหน่งที่เป็นขาของตัวอักษร ดังนั้นจุดที่จะทำการเชื่อมคือ จุดต่ำสุดไปจนกระทั่งถึง $Overlap.x\ max$ ในทางกลับกันถ้าเป็นลักษณะของตัวอักษรหัวเข้าจุดที่จะทำการเชื่อมจะเป็นจุด $Overlap.x\ min$ ถึงจุดต่ำสุด จะได้ผลลัพธ์ดังรูปที่ 3.12



รูปที่ 3.12 แสดงผลการใช้ฮิสโตแกรมในการเชื่อมต่อ

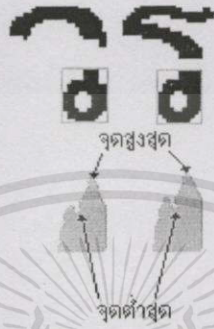
การทำงานโดยรวมของวิเคราะห์ตัวอักษรครึ่งหน้าสามารถแสดงได้ดังรูปที่ 3.13



รูปที่ 3.13 แสดงการวิเคราะห์ตัวอักษรครึ่งหน้า

3.3.4 ตัวอักษรขาด 2 ส่วนโดยไม่มีส่วนที่เหลื่อมล้ำกัน

ลักษณะการขาดของตัวอักษรแบบนี้จะเป็นการขาดของตัวอักษรที่มีขาเดียว ซึ่งจะพิจารณาการเชื่อมต่อจากครึ่งล่างเป็นหลัก โดยใช้การวิเคราะห์ทางฮิสโตแกรมช่วยจะทำให้สามารถวิเคราะห์ตำแหน่งที่น่าจะเชื่อมต่อได้ ดังรูปที่ 3.14



รูปที่ 3.14 แสดงการวิเคราะห์จุดเชื่อมต่อโดยใช้ฮิสโตแกรม

จากรูปแสดงให้เห็นจุดต่ำสุดที่อยู่ระหว่างจุดที่มีค่าฮิสโตแกรมสูง 2 ค่า โดยจะเลือกเชื่อมจากฮิสโตแกรมต่ำสุดไปที่ศที่มีฮิสโตแกรมสูงสุด จากตัวอย่างจะต้องเชื่อมจากจุดต่ำสุดไปยัง $Overlap.x_{max}$ โดยพิกัดของ y ที่จะเชื่อมจะเริ่มจาก $Overlap.x_{min}-1$ เชื่อมขึ้นไปจนพบจุดค่าขณะเดียวกันก็เชื่อมลงมาด้านล่างด้วยลักษณะเช่นเดียวกับการเชื่อมตัวอักษรที่เหลื่อมล้ำกัน ผลลัพธ์การเชื่อมต่อแสดงได้ดังรูปที่ 3.15



รูปที่ 3.15 แสดงการเชื่อมต่อที่เสร็จสมบูรณ์

การทำงานโดยรวมของการวิเคราะห์เชื่อมตัวอักษรขาด 2 ส่วนโดยไม่มีส่วนที่เหลื่อมล้ำกันสามารถแสดงได้ดังรูปที่ 3.16



รูปที่ 3.16 แสดงการเชื่อมตัวอักษรขนาด 2 ส่วน โดยไม่มีส่วนที่เหลื่อมล้ำกัน

3.3.5 ตัวอักษรขนาด 3 ส่วน

สำหรับตัวอักษรที่ขนาด 3 ส่วนจะเกิดกับตัวอักษร 2 ขา ซึ่งสามารถแบ่งการพิจารณาเป็น 2 ส่วนย่อยๆ โดยจะพิจารณาส่วนที่กว้างที่สุดเป็นหลักคือ

3.3.5.1 ส่วนที่กว้างอยู่ครึ่งบน

3.3.5.2 ส่วนที่กว้างอยู่ครึ่งล่าง

สมการของการหาส่วนที่กว้างที่สุดของส่วนที่ขนาดแต่ละส่วนจะเป็นดังสมการที่ 3.1

$$width = x_{max} - x_{min} + 1$$

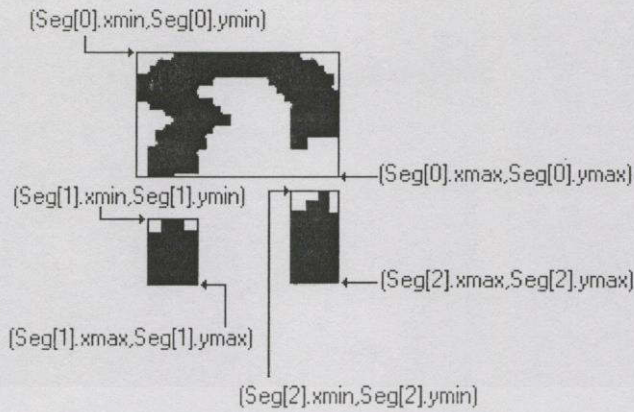
ส่วนการหาระบุตำแหน่งครึ่งบนจะเป็นดังเงื่อนไข

$$(y_{max} + y_{min}) / 2 > (y_2 + y_1) / 2$$

ถ้าเงื่อนไขเป็นจริงแสดงว่าส่วนที่กว้างที่สุดจะอยู่ครึ่งบนถ้าเป็นเท็จแสดงว่าอยู่ครึ่งล่าง

3.3.5.1 ส่วนที่กว้างอยู่ครึ่งบน

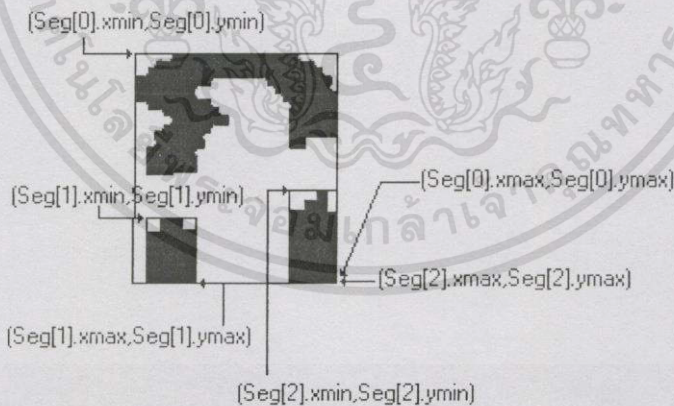
การขาดในลักษณะนี้จะพบกับตัวอักษรที่มีลักษณะปิดบน เช่น ก, ภ และ ถ เป็นต้น ลักษณะส่วนของตัวอักษรที่ขาดในแบบนี้แสดงได้ดังรูปที่ 3.17



รูปที่ 3.17 แสดงการขาดที่ส่วนกว้างอยู่ครึ่งบน

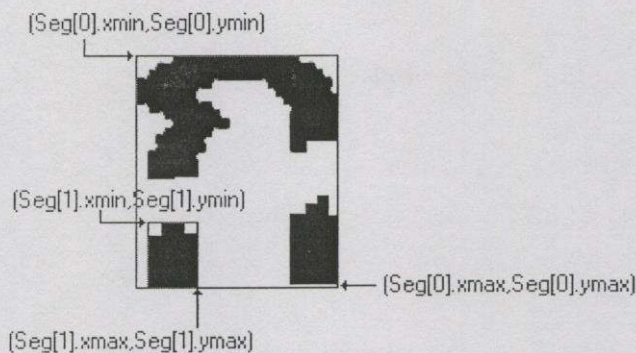
จากรูปส่วนที่กว้างที่สุดคือ $Seg[0].x_{max} = Seg[0].x_{min} + 1$ ดังนั้นส่วนที่จะยึดเป็นหลักคือ $Seg[0]$ และ $(Seg[0].y_{min} + Seg[0].y_{max}) / 2 > (y_2 + y_1) / 2$ ดังนั้นส่วนที่กว้างที่สุดจะอยู่ครึ่งบน

ในการพิจารณาจะใช้การพิจารณาเช่นเดียวกับตัวอักษรที่เหลื่อมล้ำกัน ดังนั้นจะต้องทำให้ส่วนที่ขาดเกิดการเหลื่อมล้ำกัน โดยขยายส่วน $Seg[0].y_{max}$ ให้เป็น y_2 (จุดต่ำสุดของภาพตัวอักษร) จะทำให้ภาพเกิดการเหลื่อมล้ำกันดังรูปที่ 3.18

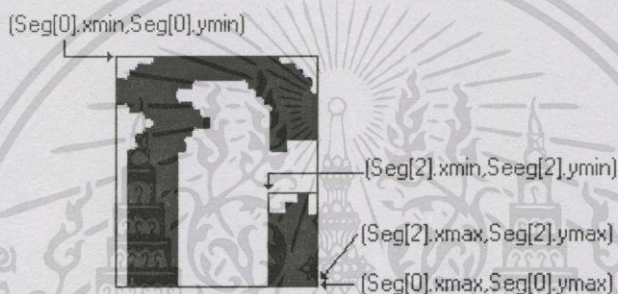


รูปที่ 3.18 แสดงการเปลี่ยนพิกัดของ $Seg[0].y_{max}$ ให้เป็น y_2 เพื่อให้เกิดการเหลื่อมล้ำ

จากนั้นทำการเชื่อมต่อ $Seg[0]$ กับ $Seg[1]$ เข้าด้วยกัน และเชื่อมต่อ $Seg[0]$ กับ $Seg[2]$ ตามกระบวนการการเชื่อมต่อตัวอักษรขาด 2 ส่วน โดยมีส่วนที่เหลื่อมล้ำกัน ในข้อ 3.3.1 ได้ ตามลำดับ ดังรูปที่ 3.19 และ รูปที่ 3.20 และผลลัพธ์จะเป็นดังรูปที่ 3.21



รูปที่ 3.19 แสดงการจัดการ Seg[0] กับ Seg[1] เพื่อทำการเชื่อม

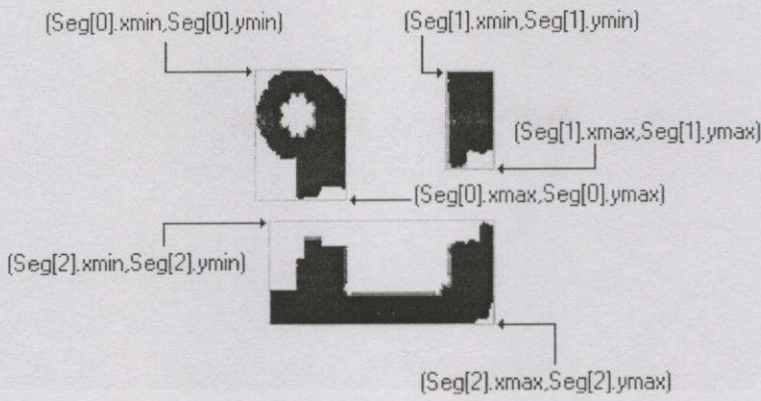


รูปที่ 3.20 แสดงการจัดการ Seg[0] กับ Seg[2] เพื่อทำการเชื่อม

รูปที่ 3.21 แสดงการเชื่อมต่อเสร็จสมบูรณ์

3.3.5.2 ส่วนที่กว้างอยู่ครึ่งล่าง

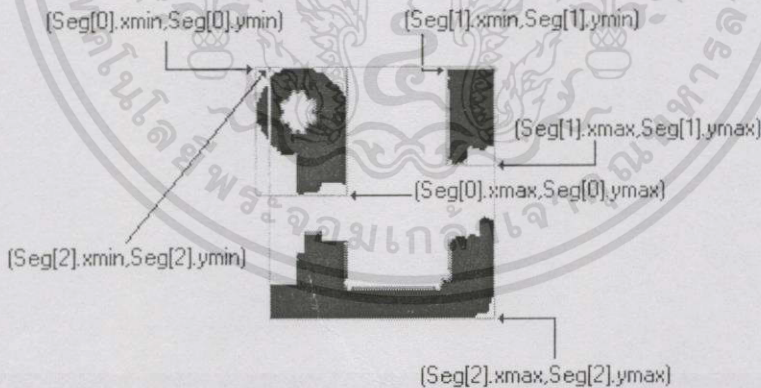
การขาดในลักษณะนี้จะพบกับตัวอักษรที่มีลักษณะปิดล่าง เช่น บ, ข และ ช เป็นต้น ลักษณะส่วนของตัวอักษรที่ขาดในแบบนี้แสดงได้ดังรูปที่ 3.22



รูปที่ 3.22 แสดงส่วนที่กว้างอยู่ครึ่งล่าง

จากรูปที่ 3.22 ส่วนที่กว้างที่สุดคือ $Seg[2].x_{max} - Seg[2].x_{min} + 1$ ดังนั้นส่วนที่จะยึดเป็นหลักคือ $Seg[2]$ และ $(Seg[2].y_{min} + Seg[2].y_{max}) / 2 < (y_2 + y_1) / 2$ ดังนั้นส่วนที่กว้างที่สุดจะอยู่ครึ่งล่าง

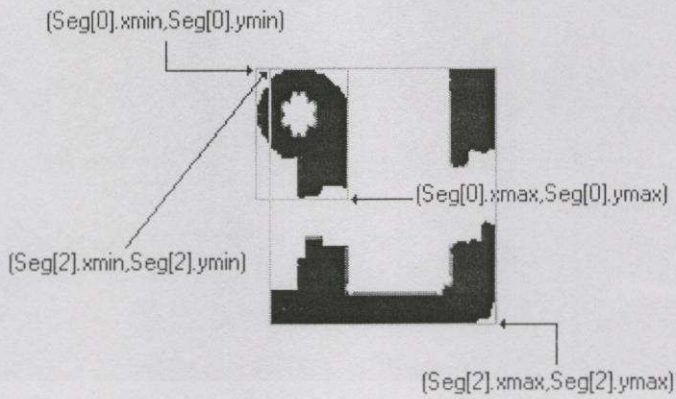
ในการพิจารณาจะใช้การพิจารณาเช่นเดียวกับตัวอักษรที่เหลื่อมล้ำกัน ดังนั้นจะต้องทำให้ส่วนที่ขาดเกิดการเหลื่อมล้ำกัน โดยขยายส่วน $Seg[2].y_{min}$ ให้เป็น y_1 (จุดสูงสุดของภาพตัวอักษร) จะทำให้ภาพเกิดการเหลื่อมล้ำกันดังรูปที่ 3.23



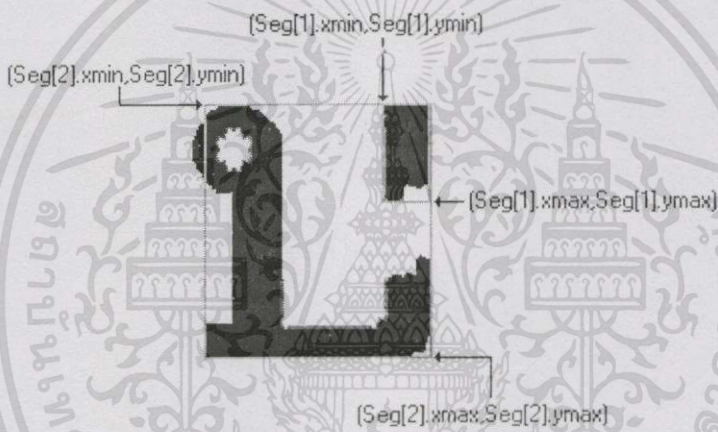
รูปที่ 3.23 แสดงการเปลี่ยนพิกัดของ $Seg[2].y_{max}$ ให้เป็น y_1 เพื่อให้เกิดการเหลื่อมล้ำ

จากนั้นทำการเชื่อมต่อ $Seg[2]$ กับ $Seg[0]$ เข้าด้วยกัน และเชื่อมต่อ $Seg[2]$ กับ $Seg[1]$ ตามกระบวนการการเชื่อมต่อตัวอักษรขาด 2 ส่วนโดยมีส่วน Character Frame เหลื่อมล้ำกัน ในข้อ 3.3.1 ได้ ตามลำดับ ดังรูปที่ 3.24 และรูปที่ 3.25 จะได้ผลลัพธ์ดังรูปที่

3.26



รูปที่ 3.24 แสดงการจัดการ Seg[2] กับ Seg[0] เพื่อทำการเชื่อมส่วนที่ขาด



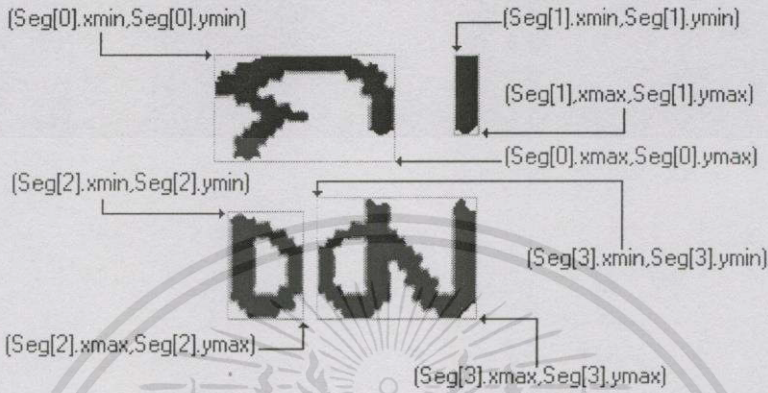
รูปที่ 3.25 แสดงการจัดการ Seg[2] กับ Seg[1] เพื่อทำการเชื่อมส่วนที่ขาด



รูปที่ 3.26 แสดงการเชื่อมต่อเสร็จสมบูรณ์

3.3.6 ตัวอักษรขาด 4 ส่วน

การขาดในลักษณะนี้จะพบกับตัวอักษรที่มีลักษณะคือมี 3 ขา เช่น ณ, ณ และ ณ เป็นต้น ลักษณะส่วนของตัวอักษรที่ขาดในแบบนี้แสดงได้ดังรูปที่ 3.27



รูปที่ 3.27 แสดงส่วนที่ขาดทั้ง 4 ส่วน

การพิจารณาจะต้องรู้ลำดับของส่วนที่ขาดทั้งหมดก่อน ซึ่งจะช่วยให้สามารถจัดการส่วนที่เชื่อมต่อได้ง่ายขึ้น ในการจัดลำดับจะใช้พิกัดกึ่งกลางทางแกน x ของทั้ง 4 ส่วนมาเปรียบเทียบ จัดลำดับจากน้อยไปมาก ดังสมการที่ 3.2

$$x_{center}[n] = (Seg[n].x_{min} + x_{Seg[n].max}) / 2$$

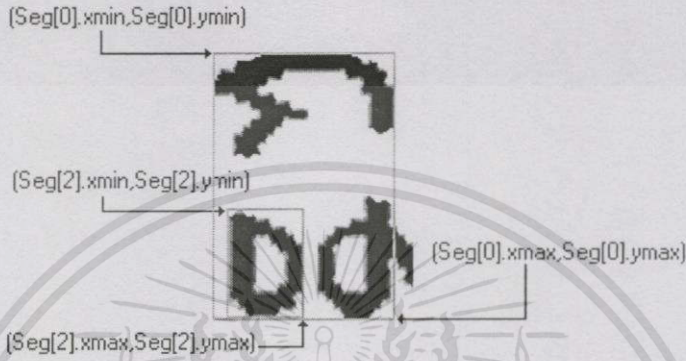
เมื่อได้ค่า x_{center} ทั้ง 4 ค่าแล้ว จะนำมาจัดเรียงค่าพิกัดจากน้อยไปมากซึ่งจะจัดเรียงได้ดังนี้

$$X_{center}[2] < X_{center}[0] < X_{center}[3] < X_{center}[1]$$

ทำให้สามารถจัดเรียงส่วนที่ขาดกันได้เป็น $Seg[2]$, $Seg[0]$, $Seg[3]$ และ $Seg[1]$ ตามลำดับ และจากลักษณะตัวอักษรที่มี 3 ขาจะทำให้ทราบว่าผลจากการจัดเรียงตำแหน่งที่ 2 และ 3 ในที่นี้คือ $Seg[0]$ และ $Seg[3]$ จะเป็นส่วนที่มีความกว้างมากกว่าส่วนอื่น ในการเชื่อมต่อจะนำส่วนที่ขาดแต่ละคู่มาทำการวิเคราะห์ทีละคู่ตามลำดับที่จัดเรียงได้ แต่จะต้องทำให้เกิดการเหลื่อมล้ำกันเพื่อให้สามารถใช้วิธีการเชื่อมต่อตัวอักษรขาด 2 ส่วน โดยมีส่วนที่เหลื่อมล้ำกันในข้อ 3.3.1 ได้

3.3.6.1 การเชื่อมต่อ Seg[0], Seg[2]

ทำการเปลี่ยนตำแหน่ง Seg[0].ymax ให้เป็น y_2 เพื่อให้เกิดการเหลื่อมล้ำกับ Seg[2] และนำ Seg[0] กับ Seg[2] ไปเชื่อมต่อโดยใช้วิธีการเชื่อมต่อตัวอักษรขาด 2 ส่วนโดยมีส่วนที่เหลื่อมล้ำกัน ในข้อ 3.3.1 แสดงได้ดังรูปที่ 3.28



รูปที่ 3.28 แสดงการเปลี่ยนตำแหน่งของ Seg[0].ymax ให้เป็น y_2

3.3.6.2 การเชื่อมต่อ Seg[0], Seg[3]

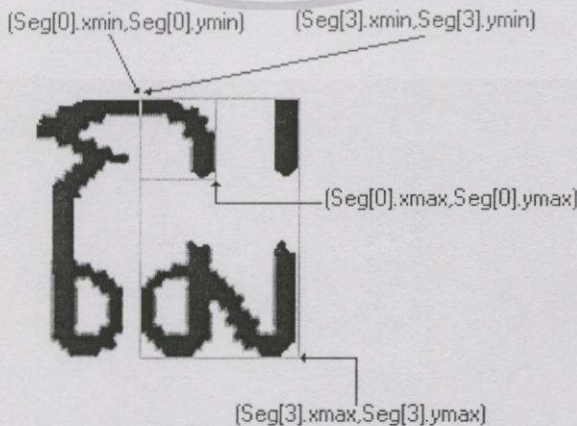
เพื่อให้เกิดการเหลื่อมล้ำของส่วนที่ขาดจะต้องทำการเปลี่ยนพิกัดดังนี้

$$Seg[0].x\ min = Seg[3].x\ min$$

$$Seg[3].y\ min = Seg[0].y\ min$$

จะได้ส่วนที่เหลื่อมล้ำคือ Seg[0] และนำ Seg[0] กับ Seg[3] ไปเชื่อมต่อโดยใช้วิธีการเชื่อมต่อตัวอักษรขาด 2 ส่วนโดยมีส่วนที่เหลื่อมล้ำกัน ในข้อ 3.3.1 แสดงได้ดังรูปที่

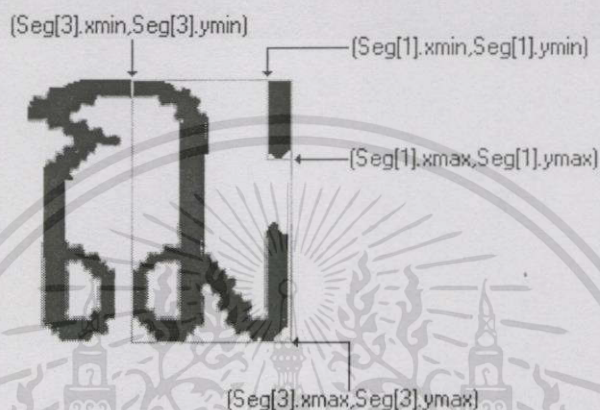
3.29



รูปที่ 3.29 แสดงการเปลี่ยนพิกัดของ Seg[0].x min และ Seg[3].y min

3.3.6.3 การเชื่อมต่อ $Seg[1], Seg[3]$

จากการปรับ $Seg[3].y_{min}$ ของขั้นตอนก่อนหน้าทำให้ $Seg[3]$ เกิดการเหลื่อมล้ำกับ $Seg[1]$ และส่วนที่เหลื่อมล้ำคือ $Seg[1]$ ดังนั้นสามารถนำ $Seg[1]$ กับ $Seg[3]$ ไปเชื่อมต่อโดยใช้วิธีการเชื่อมต่อตัวอักษรขาด 2 ส่วนโดยมีส่วนที่เหลื่อมล้ำกัน ในข้อ 3.3.1 ได้ ดังรูปที่ 3.30



รูปที่ 3.30 แสดงการเหลื่อมล้ำของ $Seg[1]$ กับ $Seg[3]$

ผลที่ได้จากการนำชิ้นส่วนทั้ง 4 ทำการจับคู่เพื่อเชื่อมต่อจะทำให้การเชื่อมต่อเสร็จสมบูรณ์ดังรูปที่ 3.31

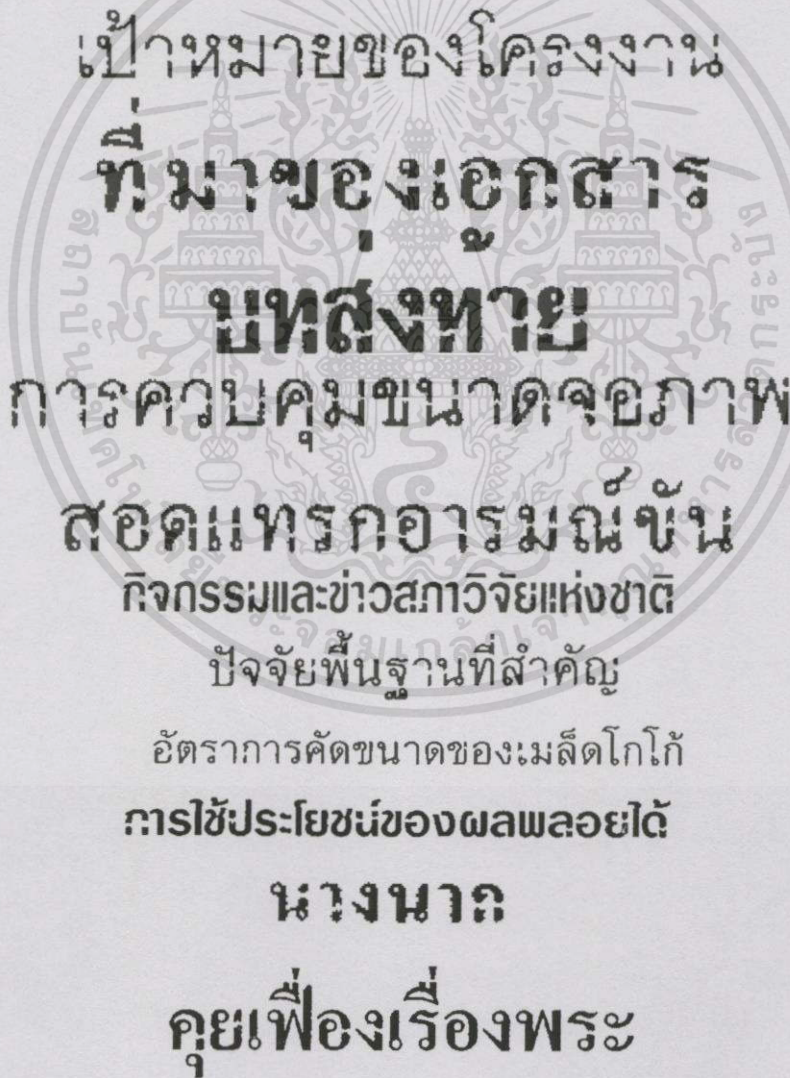


รูปที่ 3.31 แสดงการเชื่อมต่อที่เสร็จสมบูรณ์

บทที่ 4

ผลการทดลอง

งานวิจัยนี้ได้พัฒนาโปรแกรมโดยใช้ Microsoft Visual C++ เวอร์ชัน 6.0 และทำการทดลองกับข้อมูลที่ได้มาจากข้อมูลเอกสารที่มีการสแกนภาพที่ความละเอียด 600 จุด และนำข้อมูลภาพที่ได้รับการซ่อมแซมแล้วทำการทดสอบและเปรียบเทียบผลด้วยโปรแกรมการรู้จำภาษาไทยที่มีขายตามท้องตลาด เพื่อแสดงให้เห็นว่าวิธีการซ่อมแซมตัวอักษรนี้สามารถนำไปใช้งานได้จริง ข้อมูลที่ใช้ในการทดสอบมีหลายรูปแบบด้วยกันทั้งตัวหนาและตัวบาง และขนาดที่แตกต่างกันไป โดยมีตัวอย่างในการทดลองแสดงดังรูปที่ 4.1



รูปที่ 4.1 แสดงตัวอย่างข้อมูลภาพที่ใช้ในการทดสอบ

เป้าหมายของโครงการ
ที่มาของเอกสาร
บทส่งท้าย
การควบคุมขนาดจอภาพ
สอศแทรกอารมณ์ขัน
กิจกรรมและข่าวสภาวิจัยแห่งชาติ
ปัจจัยพื้นฐานที่สำคัญ
อัตราการคัดขนาดของเมล็ดโกโก้
การใช้ประโยชน์ของผลพลอยได้
นางนาก
คุยเฟื่องเรื่องพระ

รูปที่ 4.2 แสดงตัวอย่างข้อมูลภาพที่ใช้ในการทดสอบหลังการเชื่อมต่อ

ตารางที่ 4.2 แสดงผลการทดสอบโปรแกรม OCR ในการรู้จำตัวอักษรขนาดที่ซ่อนแซมแล้ว

ข้อความทดสอบ	ซอฟต์แวร์ทดสอบ	
	ThaiOCR	ArnThai
เป้าหมายของโครงการ	เป้าหมายของโครงการ	เป้าหมายของโครงการ
ที่มาของเอกสาร	ที่ มาของเอกสาร	ที่มาของเอกสาร
บทส่งท้าย	ส่งท้าย	บทสงหาย
การควบคุมขนาดจอภาพ	การควบคุมขนาดจอภาพ	การควบคุมขนาดจอภาพ
สอดแทรกอารมณ์ขัน	สอดแทรกอารมณ์ขัน	สอดแทรกอารมณ์ขัน
กิจกรรมและข่าวสภาวะวิจัยแห่งชาติ		ค่าฤๅษฤๅษ.วาวสกาว่ายห่งชาด
ปัจจัยพื้นฐานที่สำคัญ	ปัจจัยพื้นฐานที่สำคัญ	ปัจจัยพื้นฐานที่สำคัญ
อัตราการคัดขนาดของเมล็ดโกโก้	อัตราการคัดขนาดของเมล็ดโกโก้	อัตราการคัดขนาดของเมล็ดโกโก้
การใช้ประโยชน์ของผลพลอยได้		ถไรใช้ประโยชน์องผลพลอยได้
นางนาก	นางนา	นางนาก
คุยเฟื่องเรื่องพระ	คุยเฟื่องเรื่องพระ	คุยเฟื่องเรื่องพระ

จากการเก็บรวบรวมสถิติของการพบตัวอักษรขาด โดยแยกตามประเภทของเอกสารและสิ่งพิมพ์แสดงดังตารางที่ 4.3

ตารางที่ 4.3 แสดงสถิติของการพบตัวอักษรขาด

ประเภทของเอกสาร/สิ่งพิมพ์	จำนวนตัวอักษรทั้งหมด	จำนวนตัวอักษรขาด	เปอร์เซ็นต์การพบตัวอักษรขาด
วารสาร	10000	5	0.05
หนังสือพิมพ์	10354	7	0.068
ถ่ายเอกสาร	20533	99	0.48
เอกสารโรเนียว	10281	1600	15.56

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไปว่ากรณีใดทั้งสี่ มีลักษณะให้คัดลอกไปลงเน็ตหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สรุปผลการวิจัยและข้อเสนอแนะ

บทความนี้ได้นำเสนอแนวทางเพื่อแก้ปัญหาภาพตัวอักษรขาดอันเป็นสาเหตุหนึ่งที่ทำให้กระบวนการรู้จำตัวอักษรผิดพลาดหรือไม่สามารถรู้จำได้ สำหรับงานวิจัยนี้จะอยู่ในส่วนหนึ่งของกระบวนการแยกภาพตัวอักษร โดยขั้นตอนแรกจะต้องทำการหากรอบตัวอักษร โดยใช้วิธีการหาโครงร่างฮิสโตแกรมและการหาขอบภาพ ซึ่งจะทำได้ตำแหน่งและขอบเขตของภาพข้อมูลตัวอักษรแต่ละตัว ภายในกรอบตัวอักษรอาจพบทั้งภาพตัวอักษรที่ขาดและไม่ขาด ดังนั้นหลังจากหากรอบตัวอักษรได้แล้ว จะต้องนำกรอบตัวอักษรไปวิเคราะห์ว่าภาพตัวอักษรในกรอบขาดหรือไม่ โดยทำการส่งขอบเขตของภาพตัวอักษรแต่ละตัวไปวิเคราะห์การขาด โดยมีเงื่อนไขในการพิจารณาถ้าในกรอบตัวอักษรนั้นพบกรอบภาพเพียงหนึ่งกรอบแสดงว่าในกรอบตัวอักษรไม่พบตัวอักษรขาด ในทางตรงข้ามหากพบกรอบภาพมากกว่า 1 กรอบ แสดงว่าพบตัวอักษรขาดในกรอบตัวอักษร

จากการทดลองพบว่าลักษณะการขาดของตัวอักษรในภาษาไทยสามารถแบ่งเป็นกลุ่มหลักๆ ได้เป็น

1. การขาดโดยมีการเหลื่อมล้ำกันของกรอบภาพ
2. การขาดโดยที่ไม่มีการเหลื่อมล้ำกันของกรอบภาพ

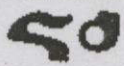
หลังจากการวิเคราะห์การขาดของตัวอักษรจะได้ตำแหน่ง ขอบเขตและลักษณะการขาด ซึ่งข้อมูลนี้จะนำไปใช้ในกระบวนการเชื่อมต่อตัวอักษรขาด

จากผลการทดลองสรุปได้ว่าวิธีการที่นำเสนอเป็นวิธีการที่มีประสิทธิภาพที่ดีและสามารถนำไปใช้กับการเชื่อมต่อตัวอักษรที่ขาดได้จริงอันจะทำให้ผลการรู้จำตัวอักษรมีความถูกต้องมากยิ่งขึ้น

สำหรับกรณีที่ทำให้การเชื่อมต่อผิดพลาดที่พบคือ กรณีที่ภาพเป็นตัวอักษรขาดเดียว โดยมี Character Frame ตรงส่วนขอบมีพิสัยไม่ตรงกันจะทำให้การเชื่อมต่อผิดพลาด แสดงดังรูปที่ 5.1

Broken Type : Non-Overlaped Segments

Character Frame



Broken Classification



Overlapped at a half-down

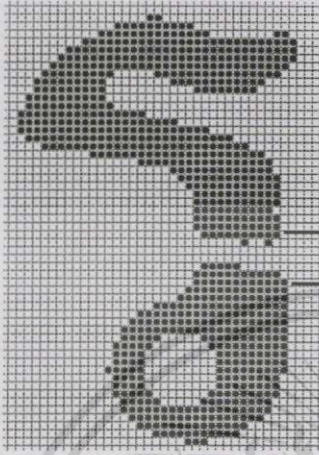
Overlapped at a half-right

A Connected Broken Character



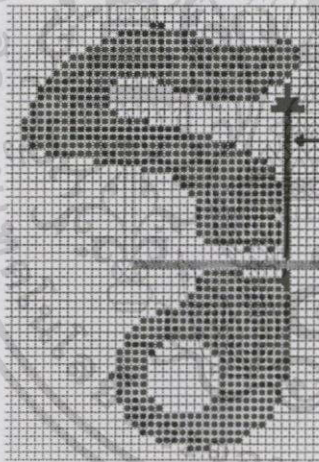
การเชื่อมต่อผิดพลาด

เมื่อทำการขยายภาพตัวอย่างตัวอักษร 'ร' จะพบว่าขอบของ Character Frame ส่วนล่าง มีขนาดกว้างกว่า ดังรูปที่ 5.2 เมื่อทำการเชื่อมต่อโดยทำการลากเส้นตรงจากแนวอ้างอิงในการเชื่อมต่อในทิศทางขึ้น-ลงจะทำให้เกิดส่วนของเส้นเชื่อมต่อที่ผิดพลาดขึ้นดังรูปที่ 5.3



แสดงขอบของ Character Frame ส่วนล่างซึ่งมีขนาดกว้างกว่าขอบของ Character Frame ส่วนบนอยู่ 1 จุด

รูปที่ 5.2 แสดงพิกัดของ Character Frame ส่วนบนและส่วนล่างของภาพตัวอักษร 'ร'



เส้นแสดงแนวการเชื่อมต่อที่ผิดพลาดเนื่องจากความกว้างของ Character Frame ไม่เท่ากัน

ระดับอ้างอิงในการลากเส้นเชื่อม

รูปที่ 5.3 แสดงแนวการเชื่อมต่อที่ผิดพลาด

เอกสารอ้างอิง

- [1] N. Premchaiswadi, W. Premchaiswadi, P. Limmaneewichid and S. Narita, "Reconstruction of Broken Character Images for Thai Character Recognition Systems," International Conference on Digital Image Computing, Techniques and Applications, DICTA'99, Perth, Australia on December, pp.222-226, 1999.
- [2] P. Limmaneewichid, W. Premchaisawadi and N. Premchaisawadi, "Repairing Broken Thai Printed Characters Using Feature Extraction," The 1999 National Computer Science and Engineering Conference (NCSEC'99), Bangkok, Thailand, December, pp. 152-157, 1999.
- [3] N. Premchaisawadi, W. Premchaisawadi, A. Thammano and S. Narita, "Merged and Broken Printed Thai Characters Segmentation," the 1999 International Conference on Artificial Neural Network In Engineering, St. Louis, USA, on November, pp.893-898, 1999.
- [4] N. Premchaiswadi, W. Premchaiswadi, Seinosuke Narita, "Segmentation Of Horizontal and Vertical Touching Thai Character.," ITC-CSCC'99 International Technical Conference on Circuit Systems, Computers and Communications, Niigata, Japan.
- [5] Wicha Panich, Somchai Jitapunkul, Prasert Choruengwiwat, "Segmentation of Connected Characters Using Distinctive Feature Of The Character in Thai Character Recognition System." Electrical Engineering Conference on Circuits and systems, pp.338-342, 1997.
- [6] Shunji, Ching Y. Suen and Kazuhiko Yamamoto, "Historical Review of OCR Research and Development", Proceedings of the IEEE, Vol. 80,7 July 1992.
- [7] D.G. Elliman and I.T. Lancaster, "A Review of Segmentation and Contextual Analysis Techniques for Text Recognition", Pattern Recognition, Vol. 23. No. 3/4 , pp. 337-346, 1990.
- [8] Rafael C. Gonzalez, Richard E. Woods, "Digital Image Processing.," Addison-Wesley, 1993
- [9] E. R. Davies, "Machine Vision", Academic Press, 1997.

ภาคผนวก ก
ตัวอย่างเอกสารที่ใช้ในการทดลอง



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไปว่ากรณีโดยทั้งสิ้น ลึกทั้งห้ามิให้ตัดแปลงเนื้อหา และต้องอ้างถึงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สัญญาเช่าคู่สายโทรศัพท์หรือวงจรรื่นๆ

สัญญานี้ทำขึ้น ณ องค์การโทรศัพท์แห่งประเทศไทย

วันที่ _____ เดือน _____ พ.ศ. _____

สัญญานี้ทำขึ้นระหว่างองค์การโทรศัพท์แห่งประเทศไทย โดย _____

ทำการแทน ผู้อำนวยการองค์การโทรศัพท์แห่งประเทศไทย ซึ่งต่อไปในสัญญานี้เรียกว่า "ผู้ให้เช่า" ฝ่ายหนึ่ง กับ _____ โดย _____

กรรมการผู้มีอำนาจลงนามผูกพันบริษัท / ผู้รับมอบอำนาจตามหนังสือมอบอำนาจลงวันที่ _____

ทะเบียนนิติบุคคลเลขที่ _____ ภูมิลำเนาเลขที่ _____

ซึ่งต่อไปในสัญญานี้เรียกว่า "ผู้เช่า" อีกฝ่ายหนึ่ง ทั้งสองฝ่ายได้ตกลงทำสัญญากันดังมีข้อความต่อไปนี้

ข้อ 1. ผู้เช่าตกลงเช่าวงจรรุ่นเช่าแบบจุดต่อจุด

จำนวน _____ วงจร ความเร็วที่ _____ Kbps/Mbps

ต้นทางจากที่ _____ เลขที่ _____ ซอย _____

ถนน _____ แขวง _____ เขต _____

จังหวัด _____

ปลายทางถึงที่ _____ เลขที่ _____ ซอย _____

ถนน _____ แขวง _____ เขต _____

จังหวัด _____

หมายเลขวงจร _____

โดยมีอัตราค่าเช่ารายเดือนเป็นเงินดังนี้:

ค่าเช่า คู่สายวงจรรุ่นเช่าแบบจุดต่อจุด รวมระยะเวลา _____ เดือน / ปี

ความเร็ว _____ Kbps / Mbps ผ่าน _____ ชุมสาย เป็นเงิน _____ บาท

รวมเป็นเงินทั้งสิ้นที่ต้องชำระ _____ บาท / เดือน

ข้อ 2. ในวันที่ทำสัญญานี้ ผู้เช่าได้ชำระเงินค่าขอใช้จำนวน _____ บาท

(_____) และค่าเช่าล่วงหน้าเป็นเวลา _____ เดือน

จำนวน _____ บาท (_____) และสัญญา

ว่าจะชำระค่าเช่าต่อไปตามอัตราที่กำหนดในสัญญาข้อ 1. ภายในกำหนด 7 วัน นับแต่วันที่ผู้ให้เช่าเรียกเก็บ

ข้อ 3. ผู้เช่าตกลงเช่าวงจรรุ่นเช่าตามสัญญาข้อ 1. มีกำหนดระยะเวลา _____ เดือน / ปี

รวมเป็นเงินทั้งสิ้น _____ บาท (_____)

และผู้ให้เช่าตกลงให้ค่าตามระยะเวลาดังกล่าว และสัญญานี้มีผลบังคับตั้งแต่วันที่ลงนามในสัญญาเป็นต้นไป ถ้าผู้เช่าประสงค์จะเลิกสัญญาก่อนกำหนดจะต้องแจ้งให้ผู้ให้เช่าทราบล่วงหน้าเป็นลายลักษณ์อักษรอย่างน้อย 7 วันและผู้เช่ายินยอมชำระค่าเช่าตามระยะเวลาที่เช่าจริงตามอัตราที่กำหนดไว้ในระเบียบหรือข้อบังคับ

ข้อ 4. ในกรณีที่ถ้าเป็นต้องนำคู่สายโทรศัพท์ หรือวงจรรุ่นเช่ามาให้เช่าตามสัญญานี้ไปใช้ในทางราชการ ผู้ให้เช่าอาจแจ้งผลการใช้คู่สายโทรศัพท์ หรือวงจรรุ่นเช่าตามสัญญานี้เป็นกรณีชั่วคราว หรือยกเลิกสัญญานี้เสียก็ได้ โดยแจ้งผู้ให้

เช่าทราบเป็นลายลักษณ์อักษรล่วงหน้าอย่างน้อย 7 วัน แต่ต้นเกิดความจำเป็นของทางราชการโดยด่วน ซึ่งผู้ให้เช่าไม่สามารถแจ้งให้ทราบล่วงหน้าได้ ผู้ให้เช่ามีสิทธิจัดการให้เช่าคู่สายโทรศัพท์หรือวงจรงดค้างได้ทันที

ข้อ 5. ผู้เช่าสัญญาว่าจะปฏิบัติตามพระราชบัญญัติ กฎบังคับ และระเบียบอื่นเกี่ยวกับการเช่าใช้คู่สายโทรศัพท์หรือวงจรงดค้างการ โทรศัพท์แห่งประเทศไทย ที่ใช้อยู่หรือที่จะเปลี่ยนแปลงแก้ไข หรือที่จะตราขึ้นใหม่ทุกประการ

ข้อ 6. ในกรณีที่ผู้ให้เช่านำคู่สายโทรศัพท์ หรือวงจรงดค้างที่ได้ให้เช่าตามสัญญาไปใช้ในกิจการอื่น หรือเกิดอุปสรรคในการใช้บริการ ผู้ให้เช่าจะคิดส่วนลดค่าเช่าให้ผู้เช่า โดยผู้เช่าจะเรียกร้องค่าเสียหายอย่างไรไม่ได้

ข้อ 7. ในกรณีที่เกิดอุปสรรค หรือเหตุขัดข้องอันเนื่องมาจากอุปกรณ์ทางเทคนิค หรือความผิดพลาดหรือความบกพร่อง อันเนื่องมาจากการกระทำของผู้เช่า จนเป็นเหตุให้ผู้เช่าใช้บริการไม่ได้ ผู้เช่าจะเรียกร้องค่าเสียหายใดๆ จากผู้ให้เช่าไม่ได้ และผู้ให้เช่าก็จะไม่คิดค่าตัวมูลค่าให้แก่ผู้เช่าด้วย

ข้อ 8. ผู้ให้เช่าสงวนไว้ซึ่งสิทธิในการให้เช่าคู่สายโทรศัพท์ หรือวงจรงดค้างเฉพาะในกิจการของผู้เช่าที่ได้แจ้งวัตถุประสงค์ไว้ ข้อ 1. เท่านั้น ผู้เช่าจะนำไปใช้วัตถุประสงค์อื่น หรืองานประเภทอื่น หรือนำไปให้ผู้อื่นเช่าพ่วงมิได้ เว้นแต่ผู้ให้เช่าจะได้ตกลงยินยอมเป็นลายลักษณ์อักษร

ข้อ 9. ถ้าผู้เช่าไม่ปฏิบัติตามสัญญา ผู้ให้เช่าจะจัดการให้เช่าคู่สายโทรศัพท์หรือวงจรงดค้าง หรือจะเลิกเสียก็ได้ ผู้เช่ายอมรับว่าการกระทำใดๆ อันเนื่องมาจากการเลิกสัญญาตามข้อนี้ ไม่เป็นการละเมิดต่อผู้เช่าแต่ประการใดและผู้เช่ายังคงต้องปฏิบัติตามหน้าที่ที่จะต้องชำระค่าเช่าที่ค้างอยู่ให้ครบถ้วนด้วย

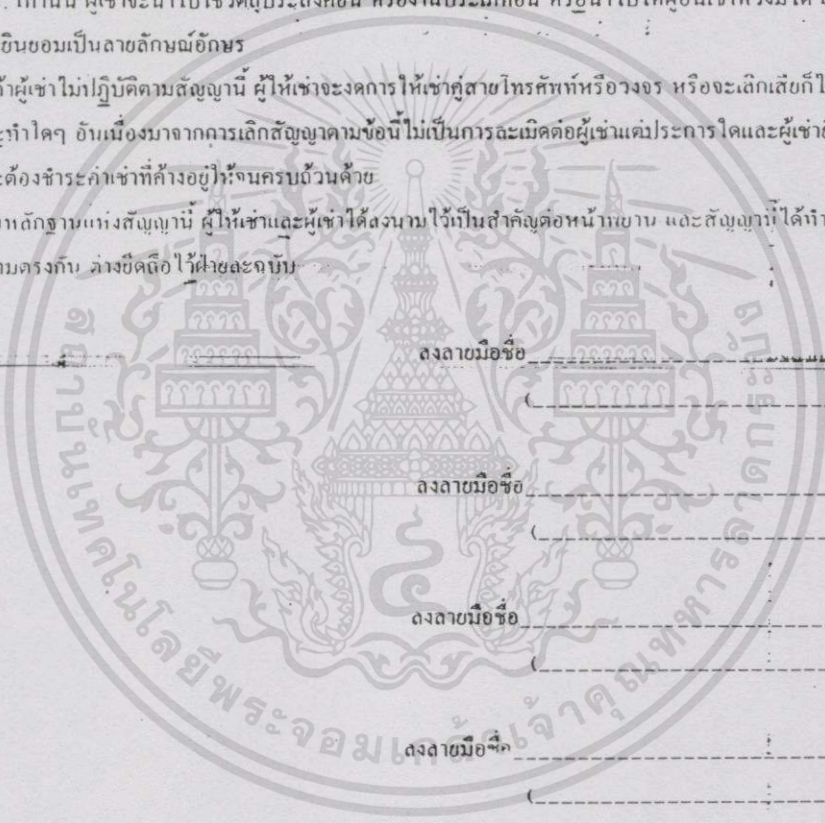
เงื่อนไขหลักฐานแห่งสัญญา ผู้ให้เช่าและผู้เช่าได้สงวนไว้ในสำคัญต่อหน้าพยาน และสัญญาที่ได้ทำไว้ทั้งสองฉบับมีข้อความตรงกับ ร่างข้อนี้ไว้โดยละเอียด

_____ ลงลายมือชื่อ _____ ผู้เช่า

_____ ลงลายมือชื่อ _____ ผู้ให้เช่า

_____ ลงลายมือชื่อ _____ พยาน

_____ ลงลายมือชื่อ _____ พยาน



ข้อปฏิบัติของการดำเนินงานโครงการวิศวกรรม ภาควิชาคอมพิวเตอร์

1. เสนอ “หัวข้อ”

สำหรับวิชาโครงการวิศวกรรม 1 เมื่อนักศึกษาลงทะเบียนเรียบร้อยแล้ว ขั้นตอนแรกคือพิจารณาเลือก “หัวข้อโครงการฯ” โดยการเสนอหัวข้อโครงการนี้ ตัวนักศึกษาจะเป็นคนตัดสินใจเองตามความถนัดหรือความสนใจ การเสนอหัวข้อโครงการฯ จะมีอยู่ 2 แบบ คือเสนอหัวข้อที่มีอยู่แล้วจากภาควิชาหรือเสนอหัวข้อที่คิดเอง ขั้นตอนนี้จะสมบูรณ์ก็ต่อเมื่อนักศึกษายื่น “แบบฟอร์มเสนอหัวข้อโครงการที่มีลายเซ็นอาจารย์ที่ปรึกษาเซ็นกำกับ” ส่งธุรการภาคฯ และโครงการวิศวกรรม 2 ให้ปฏิบัติเหมือนกัน แต่หัวข้อที่เสนอจะเป็นหัวข้อที่ต่อเนื่องจากโครงการ วิศวกรรม 1

2. เขียน “รายงานเสนอโครงการ” (Project Proposal)

เฉพาะโครงการวิศวกรรม 1 หลังจากที่นักศึกษาเสนอหัวข้อโครงการแล้ว นักศึกษาต้องเขียน “รายงานเสนอโครงการ” ส่งที่ธุรการภาคฯ การเขียนรายงานเสนอฯ นักศึกษาต้องปรึกษากับอาจารย์ที่ปรึกษารูปแบบดูจาก “ระเบียบ การดำเนินการวิชาวิศวกรรม 1” (เอกสารหมายเลข 2, หน้า 16-18)

3. “ดำเนินการ”

ทั้งวิชาโครงการวิศวกรรม 1 และ 2 ช่วงนี้จะเป็นช่วงที่นักศึกษาดำเนินงานต่างๆ เช่น ค้นคว้า, หาข้อมูล, ออกแบบ, เขียนโปรแกรม, ทดลองและบันทึกผล โดยนักศึกษาดำเนินงานต้องพบอาจารย์ที่ปรึกษาพร้อมกับ “รายงานความก้าวหน้าของโครงการ” อย่างน้อยสัปดาห์ละ 1 ครั้ง (ข้อกำหนดนี้อาจแก้ไขได้ขึ้นกับอาจารย์ที่ปรึกษาเป็นสำคัญ)

4. “ตรวจสอบรูปแบบและเนื้อหา”

การตรวจสอบรูปแบบและเนื้อหา จะกำหนดให้อยู่ในเวลา 2 สัปดาห์ โดยปกติช่วงเวลาดังกล่าวจะไม่พอ ดังนั้นในช่วง 2-3 สัปดาห์สุดท้ายของช่วงดำเนินการ นักศึกษาควรนำรายงานบางส่วน เช่น บทที่ 1-2 ให้อาจารย์ที่ปรึกษาอ่านเพื่อตรวจสอบรูปแบบและเนื้อหา ก่อน

5. “ขอขึ้นสอบ”

ทั้งวิชาโครงการวิศวกรรม 1 และ 2 เมื่อนักศึกษาประสงค์ขึ้นสอบ ต้องได้รับความเห็นชอบจากอาจารย์ที่ปรึกษา ก่อน โดยนักศึกษาต้องนำ “ใบอนุมัติขอขึ้นสอบ” ให้อาจารย์ที่ปรึกษาเซ็นอนุมัติแล้วจึงนำใบอนุมัติขอขึ้นสอบแนบกับ รายงาน 5 ชุด ส่งธุรการภาคฯ

6. “สารัตถ์โครงการ”

สำหรับโครงการวิศวกรรม 2 กระบวนการขอขึ้นสอบจะสมบูรณ์ได้ ต้อง “ผ่านการสารัตถ์โครงการ” จึงจะมีสิทธิสอบ

7. “การเตรียมตัวสอบ”

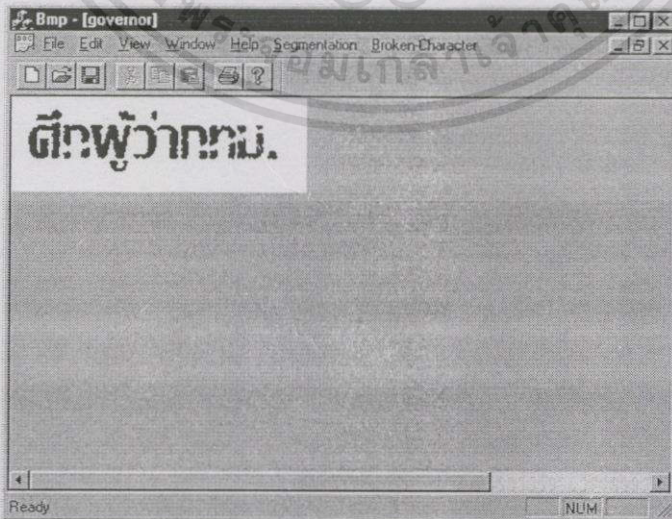
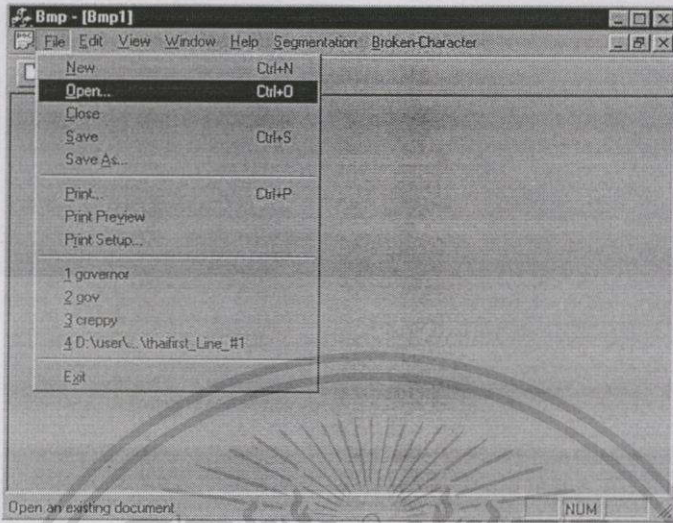
การเตรียมตัวสอบ ให้ดูตามรายละเอียด “เกี่ยวกับ สอบวิชาโครงการ”

ภาคผนวก ข

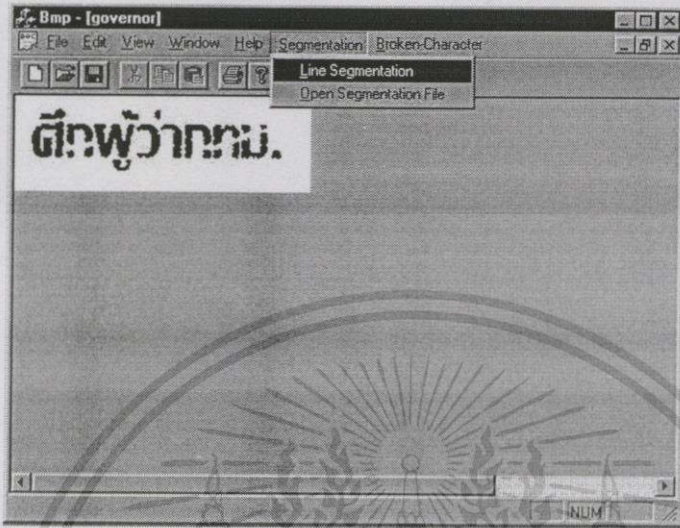
การใช้งานโปรแกรมเชื่อมต่อตัวอักษรขาด



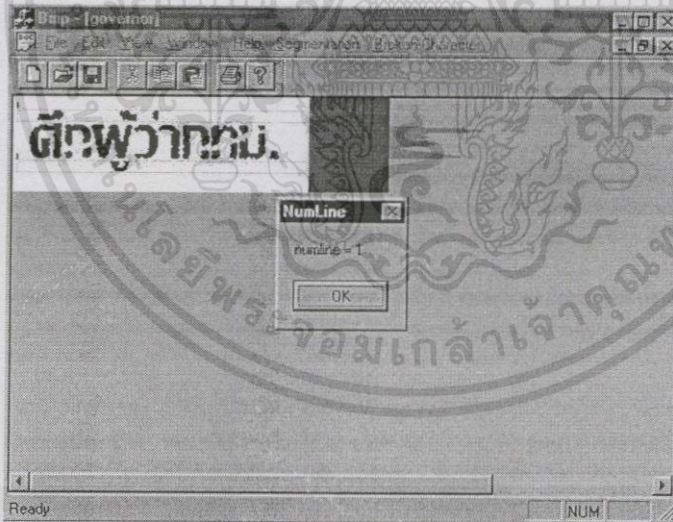
1. ไปที่เมนู File และเลือก Open เพื่อทำการเปิดไฟล์ที่มีตัวอักษรขาด



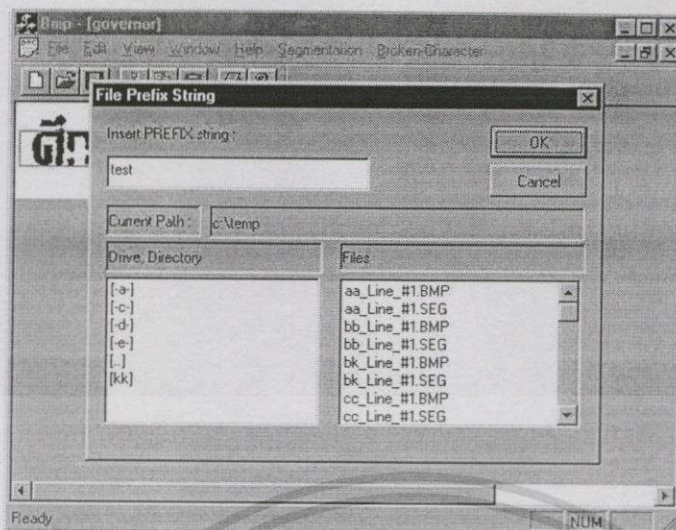
2. ขั้นตอนแรกจะต้องทำการ Segmentation เพื่อทำการหาขอบเขตและข้อมูลของภาพตัวอักษรแต่ละตัว โดยเลือกเมนู Segmentation เลือก Line Segmentation



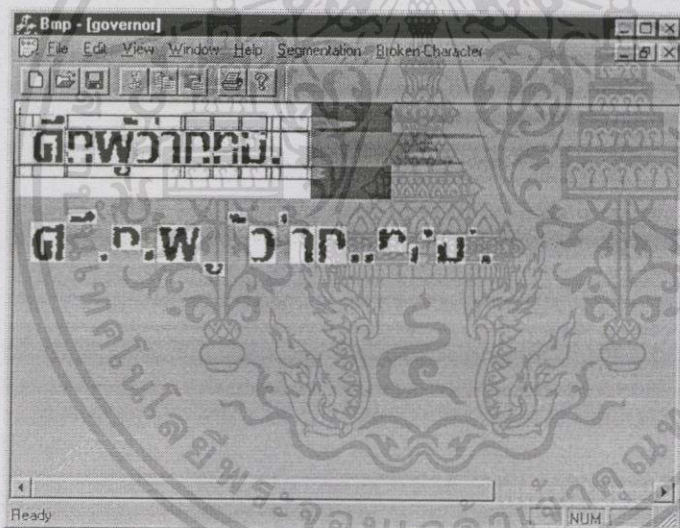
โปรแกรมจะแสดงจำนวนบรรทัดที่ได้จากกระบวนการ Segmentation



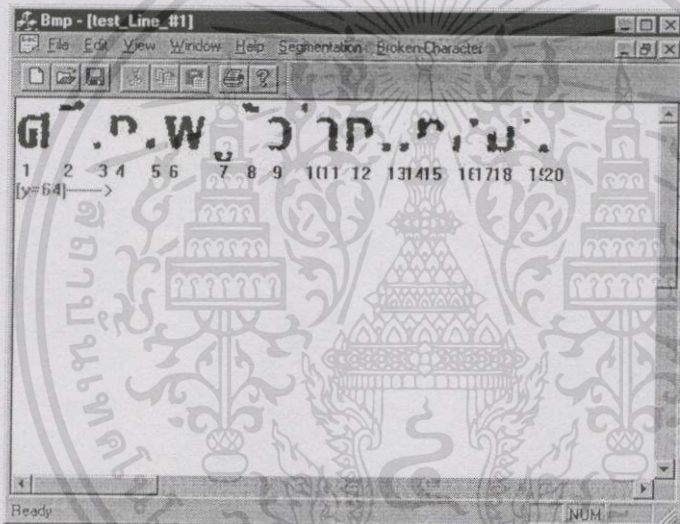
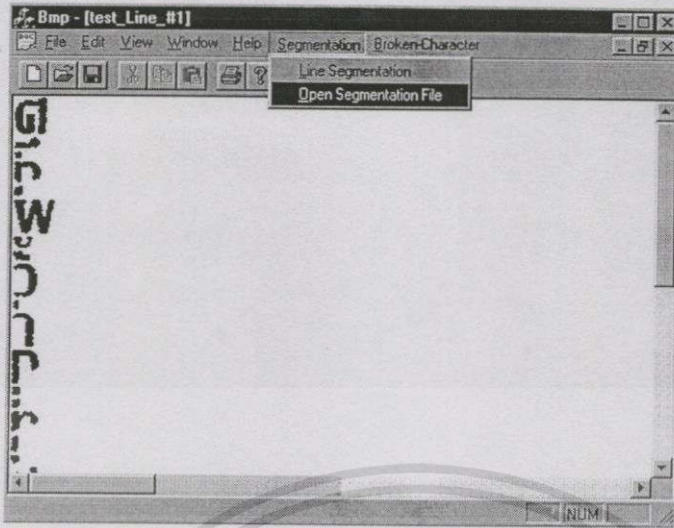
จากนั้นจะเป็นขั้นตอนในการจัดเก็บข้อมูลที่ผ่านมากระบวนการ Segmentation เรียบร้อยแล้ว โดยทำการระบุข้อความขึ้นต้นไฟล์ จากนั้นโปรแกรมจะทำการต่อท้ายชื่อไฟล์ด้วย "Line_# หมายเลขบรรทัด" เพื่อความสะดวกในการทำงาน



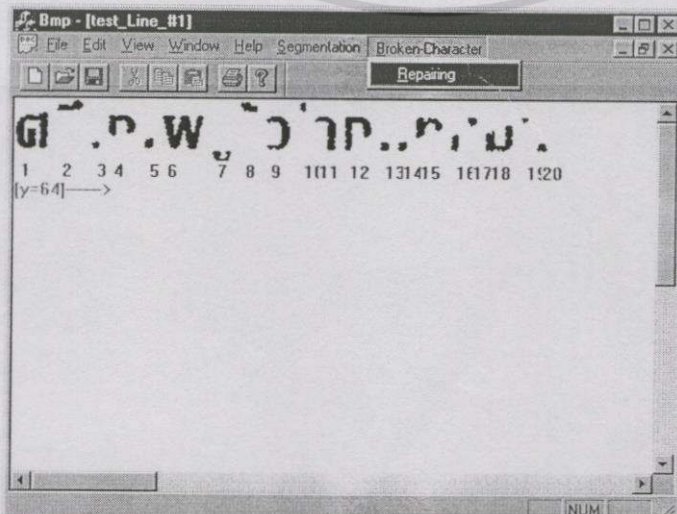
ผลลัพธ์สุดท้ายของกระบวนการ Segmentation แสดงได้ดังรูป



3. ในขั้นตอนต่อไปจะเป็นการโหลดภาพที่ทำ Segmentation แล้ว มาทำการซ่อมแซมตัวอักษรที่ขาด โดยทำการเปิดไฟล์เช่นเดียวกับในขั้นตอนที่ 1. จากนั้นให้ไปที่เมนู Segmentation เลือก Open Segmentation File เพื่อทำการอ่านข้อมูลของภาพตัวอักษรแต่ละตัว เช่น Level, ขอบเขตของตัวอักษร ความกว้าง และความสูงของตัวอักษร เป็นต้น

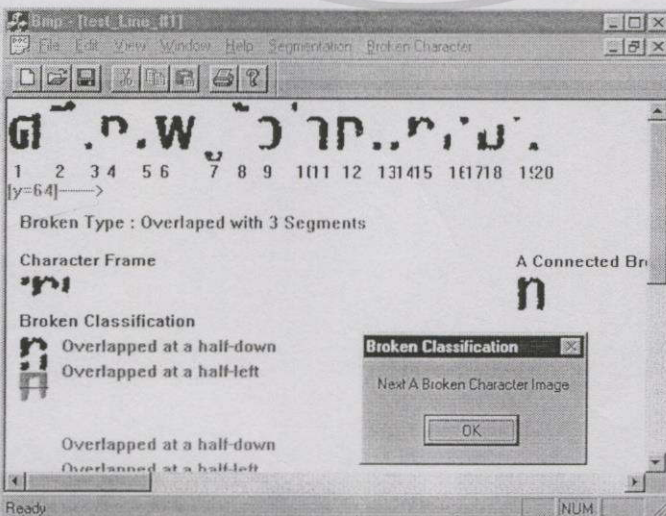
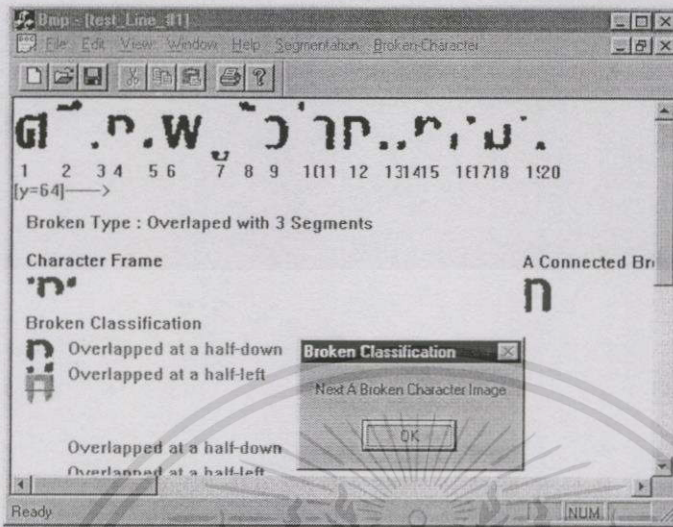


4. ในขั้นตอนนี้จะสามารถเชื่อมต่อตัวอักษรขาดได้แล้ว โดยเลือกที่เมนู Broken-Character เลือก Repairing



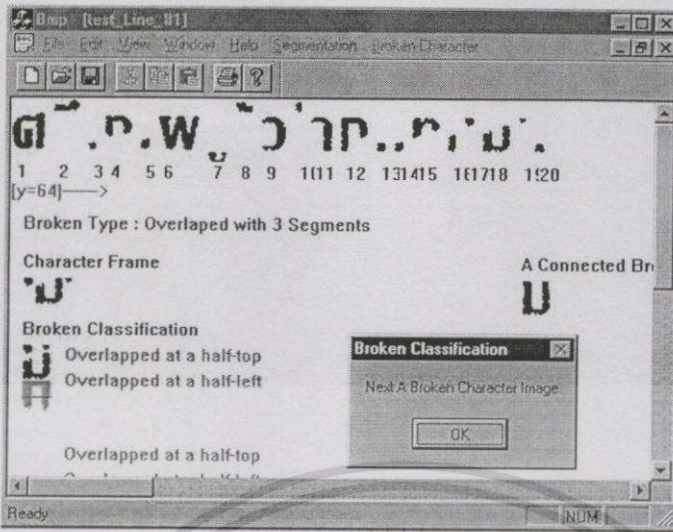
เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น เมื่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไปว่ากรณียุคทั้งสี่ปี ลึกทั้งห้าปีให้ดัดแปลงแก้ไข และต้องอ้างถึงที่มาของเอกสารทุกครั้งที่มีการนำไปใช้

หากโปรแกรมตรวจสอบพบว่ามีตัวอักษรขาดก็จะทำการเชื่อมต่อแล้วแสดงผลให้ทราบ

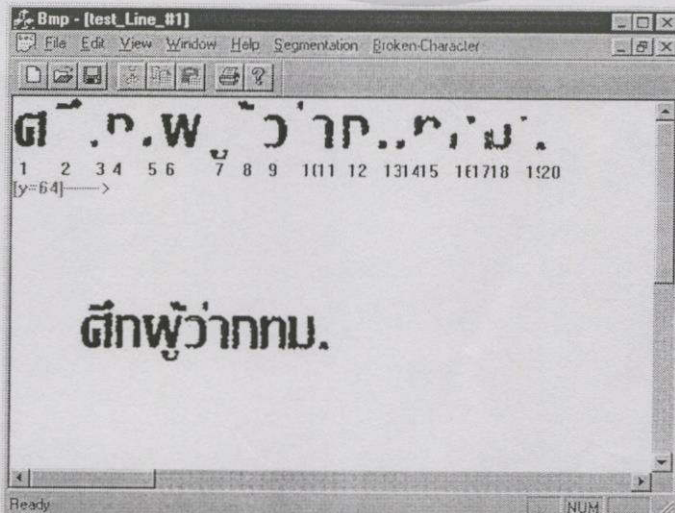
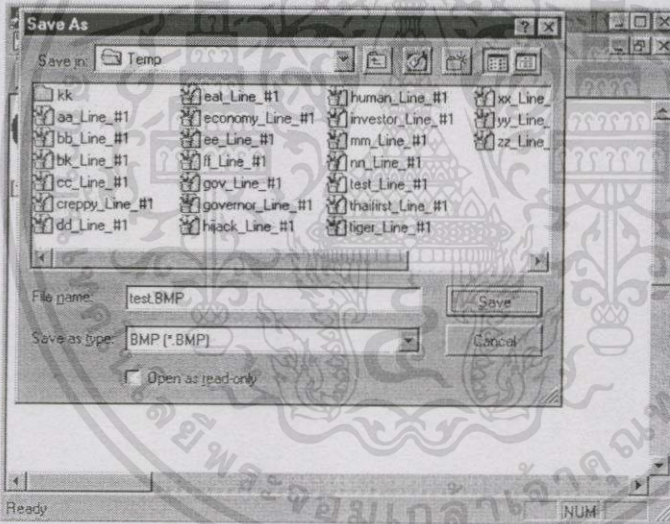


เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไปว่ากรณีใดบ้าง ซึ่งสิ่งเหล่านี้ช่วยลดแรงปะทะ และต้องอ้างอิงถึงว่าของเอกสารทว่าจริงที่มีการแก้ไขได้



เมื่อโปรแกรมตรวจสอบครบหมดทุกตัวอักษรแล้ว เราสามารถจัดเก็บผลการเชื่อมต่อเป็นไฟล์ใหม่ได้ทันที



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น เมื่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ภาคผนวก ค

ภาพตัวอย่างการเชื่อมต่ตัวอักษรขาด



ภาพตัวอักษรขาด	ผลลัพธ์จากการเชื่อมต่อ
เป้าหมายของโครงการ ที่มาของเอกสาร	เป้าหมายของโครงการ ที่มาของเอกสาร
บทส่งท้าย	บทส่งท้าย
การควบคุมขนาดจอภาพ	การควบคุมขนาดจอภาพ
สอดคล้องอารมณ์ขัน กิจกรรมและข่าวสภาวิจัยแห่งชาติ	สอดคล้องอารมณ์ขัน กิจกรรมและข่าวสภาวิจัยแห่งชาติ
ปัจจัยพื้นฐานที่สำคัญ	ปัจจัยพื้นฐานที่สำคัญ
อัตราการคัดขนาดของเมล็ดโกโก้	อัตราการคัดขนาดของเมล็ดโกโก้
การใช้ประโยชน์ของผลพลอยได้	การใช้ประโยชน์ของผลพลอยได้
นางนาก	นางนาก
คุยเฟื่องเรื่องพระ	คุยเฟื่องเรื่องพระ
สำนักงานเขตกรุงเทพฯ	สำนักงานเขตกรุงเทพฯ
ชนหัวลูกกันทั้งเมือง	ชนหัวลูกกันทั้งเมือง
กินอยู่แบบไทยๆ	กินอยู่แบบไทยๆ
ศก.ไม่พินคนเมิน ใช้บัตรพลาสติก	ศก.ไม่พินคนเมิน ใช้บัตรพลาสติก
ทวงจับมือ ศธ.แก้เกรดเพื่อ	ทวงจับมือ ศธ.แก้เกรดเพื่อ
ตึกพู่วากรม.	ตึกพู่วากรม.

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไปหาครีโดยทั้งสี่ ลึกทั้งห้ามิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ประวัติผู้เขียน

นายพงษ์สุรีย์ ลิ้มมณีวิจิตร เกิดเมื่อวันที่ 31 มกราคม 2518 สำเร็จการศึกษาระดับปริญญาตรี (คอมพิวเตอร์) จากมหาวิทยาลัยเทคโนโลยีมหานคร ปีการศึกษา 2539 ปัจจุบันเป็นอาจารย์ประจำภาควิชาคอมพิวเตอร์ คณะวิศวกรรมศาสตร์ และหัวหน้าศูนย์วิจัยและพัฒนาระบบเครือข่าย มหาวิทยาลัยเทคโนโลยีมหานคร

