

การใช้คอมพิวเตอร์ตรวจรู้อักษรภาษาไทย



นายชมพ พรพนชัย

วิทยานิพนธ์นี้ เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรบริษัทฯ สาขาสารสนเทศพัฒนาบัณฑิต

ภาควิชาบริหารธุรกิจคอมพิวเตอร์

บัณฑิตวิทยาลัย จุฬาลงกรณ์มหาวิทยาลัย

พ.ศ. 2529

ISBN 974-567-193-2

ลิขสิทธิ์ของบัณฑิตวิทยาลัย จุฬาลงกรณ์มหาวิทยาลัย

012099

๑๖/๖๔๕

THAI CHARACTER RECOGNITION BY COMPUTER

Mr. Chomtip Pornpanomchai, 1959-

A Thesis Submitted in Partial Fulfillment of the Requirements

for the Degree of Master of Science

Department of Computer Engineering

Graduate School

1986

ISBN 974-567-193-2

หัวข้อวิทยานิพนธ์	การใช้คอมพิวเตอร์ตรวจสอบอักษรภาษาไทย
โดย	นายชุมพิพ พรพนนชัย
ภาควิชา	วิศวกรรมคอมพิวเตอร์
อาจารย์ที่ปรึกษา	ดร.ศุภชัย ตั้งวงศ์ศานต์ ผู้ช่วยศาสตราจารย์ สุเมธ วัชระชัยสุรพล



บัณฑิตวิทยาลัย จุฬาลงกรณ์มหาวิทยาลัย อนุมัติให้นับวิทยานิพนธ์ฉบับนี้
เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญามหาบัณฑิต

..... *.....* คณบดีบัณฑิตวิทยาลัย
(ศาสตราจารย์ ดร.กาว วัชระกัย)

คณะกรรมการสอบวิทยานิพนธ์

..... ประธานกรรมการ
(ดร.ยรรยง เต็งอานวย)

..... *.....* กรรมการ
(ดร.ศุภชัย ตั้งวงศ์ศานต์)

..... *.....* กรรมการ
(ผู้ช่วยศาสตราจารย์ สุเมธ วัชระชัยสุรพล)

..... *.....* กรรมการ
(ดร.จารุมาตรา บันทอง)

หัวชื่อวิทยานิพนธ์	การใช้คอมพิวเตอร์ตรวจรู้อักษรภาษาไทย
ชื่อนิสิต	นายชนกพงษ์ พราหมณชัย
อาจารย์ที่ปรึกษา	ดร. ศุภชัย ตั้งวงศ์ศานต์
ภาควิชา	ผู้ช่วยศาสตราจารย์ สุเมธ วัชระชัยศุรพล วิศวกรรมคอมพิวเตอร์
ปีการศึกษา	2529

บทคัดย่อ



ในวิทยานิพนธ์นี้มีวัตถุประสงค์ที่จะศึกษาและพัฒนาวิธีการตรวจรู้อักษรภาษาไทย ซึ่งอยู่ในรูปแบบของตัวพิมพ์ดิจิต วิธีการและเทคนิคในการตรวจรู้อักษรภาษาไทยที่ใช้ในการวิจัยสามารถนำมาประยุกต์ใช้กับอุปกรณ์หรือเครื่องมือต่างๆ ได้ เช่น เครื่องอ่านข้อมูล (optical character reader) เป็นต้น ซึ่งจะทำให้เราสามารถใช้เครื่องมือชนิดนี้อ่านข้อความในเอกสารจากต้นฉบับได้โดยตรง อันเป็นการลดขั้นตอนและเวลาในการเตรียมข้อมูลเพื่อใช้ในการประมวลผลด้วย เครื่องคอมพิวเตอร์ และยังเป็นการลดความผิดพลาดในการเตรียมข้อมูลได้อีกด้วย

อักษรที่ใช้ในการตรวจรู้จะสมมุติให้มีการจัดเตรียมอยู่ในรูปของภาพบิตเมทริกซ์ตัวอักษร (bit image matrix) ที่มีขนาด 20×20 หน่วย โดยอักษรจะเหล่านี้จะปราศจากสัญญาณรบกวน (noise free) งานเบื้องต้นได้สร้างอักษรภาษาไทยในรูปตัวพิมพ์จำนวน 5 ชุด ไว้ใช้เป็นต้นแบบในการหาผลลัพธ์ของการตรวจรู้อักษร โดยที่วิธีการพัฒนาระบบการตรวจรู้อักษรนี้จะแบ่งเป็นชั้น ตอนดังนี้

1. ขั้นตอนการเปลี่ยนภาพบิทเมทริกซ์ตัวอักษรให้เป็นโครงร่างของอักษร (skeletal form) ขั้นตอนนี้อาจเรียกว่าเป็นขั้นตอนการลดความหนา (thinning process) ของภาพบิทเมทริกซ์ตัวอักษรที่เตรียมไว้ก่อนการนำเข้ามาใช้ในการตรวจรู้อักษร

2. ขั้นตอนการเปลี่ยนโครงร่างของอักษรให้อยู่ในรูปของรหัส ซึ่งรหัสเหล่านี้จะหมายถึงลักษณะของตัวอักษรตามแนวแกวและแนวสอดม

3. ขั้นตอนการสร้างความสัมพันธ์ระหว่างรหัสในขั้นที่ 2 กับอักษรต้นแบบที่สร้างไว้เพื่อการตรวจรู้ โดยความสัมพันธ์นี้จะมีโครงสร้างเป็นรูปต้นไม้ที่เรียกว่า "recognizer tree"

จากการพัฒนาวิธีการตรวจรู้อักษรนี้ได้เขียนเป็นชุดคำสั่ง (program) และทดสอบชุดคำสั่งบนเครื่องไมโครคอมพิวเตอร์ โดยใช้廉子微处理器 8088 ภาษาที่ชุดคำสั่งที่ใช้จัดระบบงาน MS-DOS 3.1 และยังใช้ไฟรเซสเซอร์ร่วม (co-processor) 8087 เข้าร่วมด้วย

ผลการวิจัยพบว่าเทคนิคการตรวจรู้อักษรที่พัฒนาขึ้นสามารถตรวจรู้อักษรต้นแบบได้ทุกรูป ด้วยอัตราความเร็วเฉลี่ย 72 ตัวอักษรต่อนาที จากนั้นได้ทำการทดสอบอักษรเพิ่มเติมโดยการสร้างอักษรชุดใหม่เพิ่มอีกจำนวน 2 ชุด ผลปรากฏว่าสามารถตรวจรู้ได้ประมาณ 70% ของอักษรทั้งหมด

Thesis Title THAI CHARACTER RECOGNITION BY COMPUTER
Name Mr. Chomtip Pornpanomchai
Thesis Advisors Dr. Supachai Tangwongsan
 Assistant Professor Sumet Vacharachaisurapol
Department Computer Engineering
Academic Year 1986

ABSTRACT



This thesis is concerned with the development of a technique for recognizing Thai characters, particularly in printed form. The technique can be applied to devices, such as optical character reader, so that they are capable of reading printed source documents. Thus, time and error can be greatly reduced in the stage of data preparation for computer input.

It is assumed that input characters are clearly separated, noise free and pre-processed in the form of bit image matrix with the size of 20 X 20. In the proposed technique, we preliminarily use five sets of Thai printed characters as basic patterns to develop the scheme of recognizer. The technique essentially consists of the following procedures:

First, transform each bit image matrix which represents a base character into a skeletal form. This step is so-called thinning process. Next, convert the skeletal form into a string of codes, which represents the pattern characteristics of rows and columns. Finally, establish links between base characters and their code strings, a recognizer tree is then formed in this step.

The technique is programmed and implemented on 16-bit 8088 based microcomputer. The machine is running under MS-DOS 3.1 operating system and also with 8087 co-processor. It is found that the perposed technique is able to recognize all base characters with the speed of 72 characters per minute. Further test is performed by preparing two new sets of printed Thai characters, the result shows that only 70% of correct recognition is obtained.

กิตติกรรมประกาศ

วิทยานิพนธ์ลับบันไดสาเร็จลุล่วงไปได้ด้วยความช่วยเหลืออย่างดีเยี่ยม
ของ ดร.ศุภชัย ตั้งวงศ์ศานต์ และผู้ช่วยศาสตราจารย์ สุเมธ วัชระชัยสุรพล
อาจารย์ที่ปรึกษาวิทยานิพนธ์ ซึ่งท่านทั้งสองได้กรุณาสละเวลาอันมีค่าของท่านให้
ดำเนินงานและข้อคิดเห็นต่างๆ ของการวิจัยมาด้วยดี

การวิจัยครั้งนี้ได้ใช้อุปกรณ์บริการของสำนักคอมพิวเตอร์ มหาวิทยาลัย
นิตล ได้มี คุณอรรเทพ สกิรบัญญา ผู้ช่วยในการจัดเตรียมอุปกรณ์บริการที่สำหรับ
งานวิจัย คุณบริชา งามเสาวรส ผู้ให้คำปรึกษาด้านชุดคำสั่งภาษาแอลเซมบล
ตลอดจน คุณกอบกุล อายุรพัน และคุณศรรารุติ ทรงเจริญ ผู้ให้คำแนะนำด้านการ
ใช้ภาษาสำหรับการเขียนวิทยานิพนธ์ ผู้เชี่ยวชาญด้านภาษาและคุณมา ณ ที่นี่

ชนพิพ พราหมณชัย



สารบัญ

หน้า

บทคัดย่อภาษาไทย	๔
บทคัดย่อภาษาอังกฤษ	๕
กิจกรรมประจำ	๖
รายการตารางประจำ	๘
รายการรูปประจำ	๙
บทที่	
1. บทนำ	1
1.1 ความเป็นมาของนักภาษา	1
1.2 วัตถุประสงค์ของการวิจัย	1
1.3 ขอบเขตของงานวิจัย	2
1.4 วิธีดำเนินการวิจัย	10
2. แนวความคิดเกี่ยวกับการตรวจรู้อักษร	11
2.1 การรับภาพของมุขย์	11
2.2 ระบบการตรวจรู้อักษร	12
2.3 ประวัติการค้นคว้าเกี่ยวกับการตรวจรู้อักษรภาษาไทย	16
2.4 เทคนิคการตรวจรู้อักษร	18
3. การตรวจรู้อักษรภาษาไทย	21
3.1 โครงสร้างของอักษรและคำภาษาไทย	21
3.2 สมाचิกของภาพนิทเมทริกซ์ตัวอักษร	22
3.3 การแปลงภาพนิทเมทริกซ์ตัวอักษรที่เป็นโครงร่าง ..	22
3.4 การกำหนดรหัสตัวอักษรแต่ละรูป	27
3.5 การค้นหารหัสต้นแบบเพื่อการตรวจรู้	30
3.6 ผลลัพธ์ที่ได้จากการตรวจรู้อักษร	36

บทที่		หน้า
4.	การพัฒนาระบบการตรวจรู้อักษร	38
4.1	ผังงานของระบบการตรวจรู้อักษร	38
4.2	ลักษณะของจอภาพที่แสดงรายการต่าง ๆ	38
4.3	การสับเปลี่ยนชื่อเมืองภายในหน่วยความจำหลัก	40
4.4	การตรวจรู้อักษรภาษาไทย	49
4.5	การแสดงผลของการตรวจรู้อักษร	49
5.	สรุปผลการวิจัยและ ข้อเสนอแนะ	51
5.1	ผลการวิจัย	51
5.2	สรุปผลการวิจัย	53
5.3	ข้อเสนอแนะ	54
	เอกสารอ้างอิง	56
ภาคผนวก ก.	ตารางแสดงรหัสต้นแบบของอักษรแต่ละรูป	58
ช.	ผังงานแสดงส่วนของชุดคำสั่งที่ใช้ในการตรวจรู้อักษร ..	67
ค.	แสดงอักษรทดสอบภาษาไทย	71
	ประวัติผู้เขียน	74

รายการตารางประกอบ

ตารางที่	หน้า
3.1 ตารางก้านดรหสสานรับการตรวจรู้อักษรภาษาไทย	28
3.2 ตารางแสดงขนาดของข้อมูลที่บรรจุในหน่วยความจำหลัก ..	31
4.1 ตารางแสดงตัวแหน่งของข้อมูลที่บรรจุในหน่วยความจำหลัก	47
5.1 ตารางแสดงผลของการตรวจรู้อักษรทดสอบภาษาไทย ...	51
ก ตารางแสดงรหัสต้นแบบของอักษรแต่ละรูป	58

รายการรูปประกอบ

รูปที่	หน้า
1.1 แสดงภาพนิทเมทริกซ์ตัวอักษร "ก"	2
1.2 แสดงอักษรภาษาไทย	3
1.3 แสดงอักษรต้นแบบภาษาไทยแบบพิมพ์ที่ 1	4
1.4 แสดงอักษรต้นแบบภาษาไทยแบบพิมพ์ที่ 2	5
1.5 แสดงอักษรต้นแบบภาษาไทยแบบพิมพ์ที่ 3	6
1.6 แสดงอักษรต้นแบบภาษาไทยแบบพิมพ์ที่ 4	7
1.7 แสดงอักษรต้นแบบภาษาไทยแบบพิมพ์ที่ 5	8
1.8 แสดงอักษรต้นแบบภาษาอังกฤษแบบพิมพ์ที่ 1	9
2.1 แสดงโครงสร้างของนัยน์ตามนุชชาร์ด	11
2.2 แสดงส่วนประกอบของเครื่อง OCR ทั้วไป	13
2.3 แสดงส่วนรับภาพของเครื่อง OCR	13
2.4 แสดงการแบ่งอักษร "ก" เพื่อการตรวจจับข้อมูลทางสถิติ ของอักษร	19
3.1 แสดงคำในภาษาไทย	21
3.2 แสดงตัวแన่ของ N_0 ถึง N_8	23
3.3 แสดงแบบการคานวณที่ใช้พิเคราะห์ค่า N_0 จากด้านขวาไป ด้านซ้าย	24
3.4 แสดงแบบการคานวณที่ใช้พิเคราะห์ค่า N_0 จากด้านซ้ายไป ด้านขวา	24
3.5 (ก)-(จ) แสดงภาพนิทเมทริกซ์ตัวอักษร "ก" ก่อนและหลัง แปลงให้เหลือเพียงโครงร่าง	25
3.6 แสดงโครงร่างของอักษร "ก"	29
3.7 แสดงโครงสร้างของการค้นหาแบบต้นไม้	32
3.8 แสดงการเก็บข้อมูลในแต่ละระเบียนของแฟ้มข้อมูลดังนี้	33

รูปที่

หน้า

3.9 แสดงการเก็บข้อมูลในแต่ละระเบียนของแฟ้มข้อมูลอักษรฯ ..	33
3.10 แสดงการค้นหารหัสรวม "BDDGDAC" ของอักษร "ก" ...	34
3.11 แสดงลักษณะการจัดเรียงอักษรภาษาไทยในแต่ละระดับให้อยู่ในระดับเดียวกัน	37
3.12 แสดงการจัดเรียงอักษรทั้ง 4 ระดับของคำว่า "คนมุ่งมั่น" ให้อยู่ในระดับเดียวกัน	37
4.1 แสดงผังงานของระบบตรวจรู้อักษรที่พัฒนาขึ้น	39
4.2 แสดงจuxtaposition รายการหลัก	41
4.3 แสดงจuxtaposition รายการจัดคำภาษาไทยทั้ง 4 ระดับ	42
4.4 แสดงจuxtaposition รายการปรับปรุงแก้ไขรูปแบบของอักษร ..	43
4.5 แสดงจuxtaposition รายการตารางแอกซ์ของอักษรภาษาไทย ..	44
4.6 แสดงจuxtaposition แสดงผลของอักษรที่สามารถตรวจรู้ได้ ..	45
4.7 แสดงชุดคำสั่งที่ใช้ในการสับเปลี่ยนข้อมูลในหน่วยความจำหลัก	48
4.8 แสดงจuxtaposition ที่แสดงผลของอักษรที่ไม่สามารถตรวจรู้ได้ ...	50
ช.1 แสดงผังงานของส่วนของชุดคำสั่งที่ใช้ในการจัดเตรียมคำในภาษาไทย	68
ช.2 แสดงผังงานของส่วนของชุดคำสั่งที่ใช้ในการคัดเลือกอักษรเพื่อบรรจุแก้ไขรูปแบบ	69
ช.3 แสดงผังงานของส่วนของชุดคำสั่งที่ใช้ในการเปลี่ยนรูปแบบของอักษร	70
ค.1 แสดงอักษรทดสอบภาษาไทยแบบพิมพ์ที่ 1	72
ค.2 แสดงอักษรทดสอบภาษาไทยแบบพิมพ์ที่ 2	73