

บทที่ 1

บทนำ



ความเบื้องต้น

การใช้คอมพิวเตอร์สำหรับงานเอกสาร และการประมวลผลข้อมูล เป็นหนึ่งในแนวโน้มหลักของสำนักงานอัตโนมัติ (OA) ในปัจจุบัน ซึ่งการรู้จำตัวอักษรด้วยคอมพิวเตอร์ก็เริ่มเข้ามามีบทบาทในด้านนี้เพิ่มขึ้นอย่างมาก เพื่อให้เป็นหนทางเลือกในการป้อนข้อมูลเข้าสู่คอมพิวเตอร์โดยการอ่านข้อความจากเอกสาร หรือ รับรู้จากการเขียนที่จอภาพได้โดยตรง โดยมนุษย์ไม่จำเป็นต้องป้อนข้อมูลผ่านทางแป้นพิมพ์เพียงอย่างเดียว

การรู้จำรูปแบบตัวอักษรได้มีการประยุกต์ใช้กับหลายภาษาเช่น ตัวอักษรภาษาอังกฤษซึ่งได้มีการพัฒนาก้าวหน้าอย่างมากจนกระทั่งสามารถนำมาประยุกต์ใช้ และ จำหน่ายในเชิงพาณิชย์ได้ทั้งแบบตัวพิมพ์และ ตัวเขียน ส่วนภาษาอื่นเช่น ตัวอักษรภาษาอาหรับ , ตัวอักษรภาษาจีน , ตัวอักษรคันจิในภาษาญี่ปุ่น เป็นต้น ก็ได้มีงานวิจัยศึกษาค้นคว้าเพื่อพัฒนาระบบการรู้จำของภาษานั้น ๆ ออกเผยแพร่ และ ตีพิมพ์ หลายงานวิจัย แต่สำหรับตัวอักษรภาษาไทยได้มีการศึกษาค้นคว้ากันมาบางพอสมควร แต่เนื่องจากตัวอักษรภาษาไทยมีโครงสร้างที่ค่อนข้างซับซ้อนประกอบด้วยลักษณะที่เป็นเส้นตรงผสมเส้นโค้ง และ วงกลม ซึ่งบางครั้งมีการตัดกันของเส้น รวมทั้งตัวอักษรก็มีหลายระดับ ทำให้การศึกษาค้นคว้าไม่แพร่หลาย และ ต่อเนื่องมากนัก ดังเช่นงานวิจัยการรู้จำตัวอักษรไทย และ ตัวเลขไทยเท่าที่พบคือ

สรุพันธ์ เอื้อไพบุลย์ (2531)[1] นำเสนอวิธีการรู้จำตัวอักษรลายมือเขียนภาษาไทย โดยการนำหัวของตัวอักษรนำมาพิจารณาเพื่อทำการจำแนกกลุ่มของตัวอักษรออกเป็นกลุ่มย่อย โดยให้ตัวอักษรที่มีหัวอยู่บริเวณเดียวกันอยู่กลุ่มเดียวกัน และใช้คุณสมบัติทางการจำแนกชนิด (topology) ของตัวอักษร ในการแยกตัวอักษรออกจากกลุ่มย่อย รวมทั้งยังนำเอาเทคนิคหลาย ๆ วิธี เพื่อให้แยกตัวอักษรที่คล้ายคลึงกันออกจากกันด้วย ซึ่งข้อดีคือการเปรียบเทียบแบบนี้ทำให้เปรียบเทียบได้อย่างรวดเร็ว

พิพัฒน์ และ มนลดา[2] นำเสนอการรู้จำตัวอักษรไทยหลายรูปแบบโดยใช้วิธีการไดนามิกโปรแกรมมิ่ง วิธีการนี้อาศัยการพิจารณาเส้นแสดงขอบของอักขระโดยนำรหัสทิศทางแบบลูกโซ่

ของฟรีแมน กับ ความแตกต่างของทิศทางของเส้นแสดงขอบของอักขระ มาใช้ในการตัดแบ่งเส้นแสดงขอบของอักขระออกเป็นส่วนโค้งย่อย คือ ส่วนโค้งเว้า และ ส่วนโค้งนูน โดยจะนำมาเปรียบเทียบหาส่วนโค้งที่คล้ายกันมากที่สุดกับอักขระต้นแบบ โดยใช้เทคนิคของการเปรียบเทียบแบบไดนามิกโปรแกรมมิ่ง ผลอัตราการรู้จำถูกต้อง 94.7 เปอร์เซ็นต์

อนันต์ เอกวงศวิริยะ (2537)[3] นำเสนอเรื่องการเรียนรู้จำตัวเลขไทยแบบตัวพิมพ์โดยวิธีซินแทกติก โดยวิธีการเริ่มต้นที่ทำการปรับแต่งรูปให้มีคุณภาพดีขึ้น และ เหมาะสมเพื่อที่จะนำไปใช้ในส่วนของการแทนรูปด้วยโครงสร้างของภาษาซึ่งกำหนดโดยไวยากรณ์ต้นไม้ (tree grammar) และ จำแนกตัวเลขภาษาไทยโดยใช้เทคนิคในการหาระยะห่างระหว่างต้นไม้ที่เป็นอินพุต กับต้นไม้ต้นแบบ โดยวิธีการซินแทกติกนี้มีจุดเด่นอยู่ที่นำเอาโครงสร้างของภาพตัวอักษรมาใช้ประโยชน์ในการรู้จำ ฉะนั้นจึงสามารถรู้จำภาพของตัวอักษรแบบต่าง font และต่างขนาดได้ดี

สนธยา เมรินทร์ (2537)[4] นำเสนอเรื่องการศึกษาการเรียนรู้จำตัวอักษรพิมพ์ภาษาไทยโดยวิธีซินแทกติกซึ่งเป็นการพิจารณาที่โครงสร้างของตัวอักษร โดยมีการอธิบายโครงสร้างของตัวอักษรในรูปของประโยคที่ประกอบด้วย primitive ทำให้สามารถจำแนกตัวอักษรที่มีโครงสร้างที่แตกต่างกันออกจากกันได้อย่างง่ายและรวดเร็ว รวมทั้งยังใช้วิธีการเปรียบเทียบทาง feature ควบคู่กันไปด้วยสำหรับตัวอักษรภาษาไทยบางกลุ่มที่มีลักษณะคล้ายคลึงกันมาก

วิทยานิพนธ์ฉบับนี้ได้จัดทำขึ้นเพื่อเสนอแนวทางในการแก้ไขปรับปรุง การรู้จำตัวอักษรไทยด้วยวิธีการซินแทกติกของ สนธยา เมรินทร์[4] โดยได้นำเอาเทคนิคแบบพีชซีโลจิกเข้ามาช่วยในการปรับปรุงให้มีความแม่นยำในการรู้จำสูงขึ้นโดยยังคงการวิเคราะห์โครงสร้างของตัวอักษรในวิธีทางซินแทกติกไว้ซึ่งจะเป็นการแก้ปัญหามีอยู่คือการไม่สามารถรู้จำตัวอักษร หรือ รู้จำผิดพลาดในกลุ่มพยัญชนะ 20 ตัว ได้แก่ ข 2 ตัว , ม 1 ตัว , ฉ 1 ตัว , ช 3 ตัว , ฎ 1 ตัว , ฐ 3 ตัว , ถ 1 ตัว , น 1 ตัว , ศ 5 ตัว , ส 1 ตัว , ห 1 ตัว และกลุ่มสระ 11 ตัว ได้แก่ อี 5 ตัว , อี 1 ตัว , อี 2 ตัว , ไม้เอก 1 ตัว , ไม้โท 2 ตัว ใน font แบบ Eucrosia และ แบบ Cordia ขนาด 20,22,24,28,32,36,48

วัตถุประสงค์ของงานวิจัย

- ก. เพื่อพัฒนาโปรแกรมสำหรับการรู้จำตัวอักษรพิมพ์ภาษาไทย
- ข. เพื่อศึกษาวิธีการปรับปรุงการรู้จำตัวอักษรพิมพ์ภาษาไทย โดยเทคนิคแบบ ฟัชซีโลจิก
- ค. เพื่อเป็นแนวทางสำหรับการพัฒนาวิธีการรู้จำในเรื่อง ๆ อื่นต่อไป

ประโยชน์ที่คาดว่าจะได้รับจากงานวิจัย

- ก. พัฒนาความรู้ในเรื่องการรู้จำตัวอักษรพิมพ์ภาษาไทย
- ข. ผลิตโปรแกรมที่ทำให้คอมพิวเตอร์รู้จำตัวอักษรพิมพ์ภาษาไทยได้
- ค. เป็นแนวทางในการพัฒนาเครื่องมืออ่านตัวอักษรภาษาไทยสำหรับผู้พิการทางสายตา และเพิ่มประสิทธิภาพในการทำงานด้านสำนักงาน

ขอบเขตของงานวิจัย

1. ปรับปรุงให้อัตราการรู้จำมากกว่า 98% ขึ้นไป[4]
2. เครื่องคอมพิวเตอร์ที่ใช้ : เครื่อง PC รุ่น 486DX2-66 ใช้ไมโครโปรเซสเซอร์ 80486 ความเร็ว 66 เมกกะเฮิร์ตซ์ มีหน่วยความจำ 8 เมกกะไบต์ และจานแม่เหล็กแบบแข็งขนาด 330 เมกกะไบต์
3. ภาษาคอมพิวเตอร์ที่ใช้ในการปรับปรุง : ภาษา C
4. อุปกรณ์อ่านเอกสาร : เครื่องสแกนเนอร์ Microtek รุ่น ScanMaker II ใช้ความละเอียดของการอ่านเอกสารที่ 300 dpi
5. รูปแบบตัวอักษรที่ใช้ : ตัวอักษรพิมพ์ภาษาไทยในงานวิทยานิพนธ์ของ สนธยา เมรินทร์ ซึ่งนำรูปแบบตัวอักษรชื่อ EUCROSIA และ CORDIA ที่เป็นตัวธรรมดา ไม่เป็นตัวเอียง หรือมีเส้นใต้ มาใช้เก็บลงไฟล์ 1 ตัวอักษรต่อ 1 ไฟล์ในรูปแบบ BMP หรือ PCX โดยมีตัวอักษรที่สามารถรู้จำได้ 79 ตัวต่อ 1 รูปแบบ ดังนี้
- 5.1 พยัญชนะ พยัญชนะไทยมีทั้งหมด 44 ตัว คือ

ก ข ช ค ต พ ง

จ ฉ ช ซ ฌ ญ

