



บทนำ

ความเป็นมาของปัญหา

ในปัจจุบันการจัดพิมพ์งานต่างๆ ที่มีลักษณะทางคณิตศาสตร์เข้ามามากขึ้น ด้วยนั้นคุณภาพของงานที่ออกมากจะไม่สวยงาม บางครั้งอาจทำให้ผู้พบเห็นมีความเชื่อใจคลาดเคลื่อนจากความต้องการของผู้จัดพิมพ์ถึงแม้ว่าจะมีโปรแกรมประมวลผลคำต่างๆ ให้เลือกใช้มากมาขั้นต่ำ ก็ยังไม่มีประสิทธิภาพดีเท่าไหร่นัก นอกจากนี้ โปรแกรมประมวลผลคำต่างๆ ดังกล่าวที่มักจะถูกจำกัดด้วยกฎหมายทางลิขสิทธิ์ ทำให้ต้องเสียค่าใช้จ่ายสูงขึ้น

ในปี ก.ศ. 1977 ศาสตราจารย์ 朵南德 อี. ค努ธ (Donald E. Knuth) ได้ตระหนักรถึงปัญหาการพิมพ์หนังสือที่เกี่ยวข้องกับผลงานทางวิทยาศาสตร์ด้วยโปรแกรมประมวลผลคำดังกล่าวจะมีคุณภาพต่ำไม่เป็นที่น่าพอใจ ดังนั้นเขาจึงออกแบบตัวอักษรใหม่โดยใช้โปรแกรมเมตาฟอนต์ (METAFONT) หลังจากนั้นก็ได้พัฒนาโปรแกรมเท็กซ์ (TeX) ขึ้นเพื่อใช้กับตัวอักษรที่สร้างขึ้นด้วยโปรแกรมเมตาฟอนต์ดังกล่าว ซึ่งโปรแกรมเท็กซ์ที่พัฒนาขึ้นนี้จะให้ตัวอักษรที่มีคุณภาพดีกว่ามาก หลังจากที่การพัฒนาโปรแกรมเท็กซ์ได้สิ้นสุดลง เขายังได้อุทิศโปรแกรมนี้ให้สามารถใช้ได้ทั่วไปโดยไม่คิดค่าลิขสิทธิ์แต่อย่างใด ด้วยเหตุผลดังกล่าวนี้ ทำให้เท็กซ์มีการพัฒนาต่อไปเพื่อให้สามารถใช้งานได้กับหลายภาษา เช่น เยอร์มัน ญี่ปุ่น เป็นต้น

สำหรับในประเทศไทยของเรา การใช้งานเท็กซ์ยังไม่แพร่หลายนัก คงใช้กันแต่ในกลุ่มของนักวิจัยเท่านั้นทั้งนี้เนื่องจากการใช้งานโปรแกรมเท็กซ์ค่อนข้างยุ่งยากและซับซ้อนกว่าการใช้โปรแกรมประมวลผลคำทั่วไปอีกทั้งยังไม่มีการพัฒนาให้สามารถแสดงผลเป็นภาษาไทยได้

ในปัจจุบันกลุ่มนักเรียนไทยที่ศึกษาอยู่ที่สถาบันเทคโนโลยีแห่งโตเกียว (Tokyo Institute of Technology) ประเทศญี่ปุ่น ได้ร่วมกันศึกษาและพัฒนาโปรแกรมเท็กซ์ภาษาไทยขึ้น โดยอาศัยชุดคำสั่งต่างๆ ที่มีมากับตัวโปรแกรมลาเท็กซ์ (LaTeX) นั้นทำการจัดบรรทัดเพื่อแสดงผลเป็นภาษาไทย และสร้างโปรแกรมที่

สามารถอ่านเอกสารที่เป็นภาษาไทยที่สร้างขึ้นโดยโปรแกรมบรรณาริการภาษาไทยต่างๆ แปลงเอกสารภาษาไทยเหล่านั้นให้อยู่ในรูปแบบของไฟล์ที่มีนามสกุลเป็น .tex และสามารถใช้กับโปรแกรมเทิกซ์ธรรมด้าได้

เนื่องจากการพัฒนาดังกล่าวยังไม่สามารถที่จะใช้งานได้อย่างสมบูรณ์แบบ จึงยังไม่สามารถนำมาใช้งานได้อย่างจริงจัง ขั้นนี้ข้อมูลของที่ต้องปรับปรุงอีกมาก เช่น ในเรื่องของการตัดคำภาษาไทยที่ยังไม่มีการใช้โปรแกรมตัดคำที่มีประสิทธิภาพดีเท่าที่ควร ในเรื่องของตัวอักษรภาษาไทยที่ยังไม่มีการพัฒนาขึ้นมาใช่องค์วิทยาโปรแกรมเมต้าฟอนต์แต่ต้องใช้ฟอนต์ที่พัฒนาขึ้นจากระบบอื่นแทน เหล่านี้เป็นต้น การศึกษาและแก้ไขปัญหาเหล่านี้จะทำให้โปรแกรมเทิกซ์สำหรับภาษาไทยสมบูรณ์ยิ่งขึ้น

ด้วยเหตุผลดังกล่าวข้างต้นนี้ การสร้างโปรแกรมแสดงผลตัวอักษรไทยบนเทิกซ์จะเป็นประโยชน์สำหรับงานพิมพ์ทางด้านวิชาการที่จำเป็นต้องใช้ภาษาไทยซึ่งจะพร้อมถ่ายต่อไปในอนาคต

อนึ่งเพื่อความสะดวกในการอ้างอิง จะใช้ชื่อเรียกโปรแกรมแสดงผลตัวอักษรภาษาไทยบนระบบเทิกซ์ว่า ไทยเทิกซ์ (ThaiTeX)

แนวคิด ทฤษฎี หรือสมมติฐาน ประกอบการวิจัย

การทำงานของโปรแกรมเทิกซ์นั้น โดยทั่วไปแล้วมีหลักการคือจะสร้างเอกสาร (document) ขึ้นมาจากโปรแกรมบรรณาริการใดก็ได้ โดยเอกสารที่สร้างขึ้นมาจะต้องมีรูปแบบ (syntax) ตามกฎเกณฑ์ของเทิกซ์ หลังจากนั้นก็จะใช้โปรแกรมเทิกซ์จัดการแปลงเอกสารนั้นให้เป็นไฟล์ข้อมูลที่สามารถแสดงผลตามที่ได้ออกแบบไว้

ความหมายของเอกสาร (document) ก็คือข้อมูลที่ต้องการพิมพ์ลงกับคำสั่งต่างๆ ที่มีในเทิกซ์ ซึ่งหลักเกณฑ์นี้จะต่างกับโปรแกรมประมวลผลคำประเภท WYSIWYG (What You See Is What You Get) ที่ใช้กันอยู่ทั่วไปซึ่งจะมีเฉพาะเนื้อข้อมูลที่เราต้องการพิมพ์ล้วนๆ เท่านั้น หน้าที่ของเทิกซ์คือการแยกระหว่างเนื้อข้อมูลที่ต้องการจะพิมพ์กับคำสั่งที่ควบคุมการทำงาน ดังนั้นถ้าต้องการตัวอักษร "A" ที่เป็นตัวหนา กับตัวอักษร "A" ที่เป็นตัวอ่อน สิ่งที่ต้องทำคือพิมพ์ข้อมูลเป็นตัว "A" เซ็นเดิน แต่จะต่างกันตรงคำสั่งควบคุมการทำงานของเทิกซ์ดังนี้

\bf A จะได้ตัวอักษรที่เป็นตัวหนา A

\sl A จะได้ตัวอักษรที่เป็นตัวอ่อน A

นอกจากนี้หากซึ่งสามารถรับข้อมูลที่กำหนดให้อยู่ในรูปแบบของรหัสແອສกีได้อีกด้วย เช่น แทนที่จะพิมพ์ A ก็อาจจะพิมพ์ 41 ผลลัพธ์ที่ได้ก็ยังคงเป็นเช่นเดิม

ด้วยหลักการคังกล่าว�ี การสร้างโปรแกรมไทยเท็กซ์จึงมีความเป็นได้โดยที่สร้างตัวอักษรไทยขึ้นและแทนตัวอักษรที่สร้างขึ้นนั้นด้วยรหัส เช่นเดียวกับการแทนตัวอักษรในภาษาอังกฤษด้วยรหัสແອສกี โดยรหัสของอักษรภาษาไทยที่ใช้นี้จะเป็นเช่นเดียวกับตัวรหัสอักษรภาษาไทยที่ใช้กันอยู่โดยทั่วไป เมื่อนำรหัสของอักษรเหล่านี้มาจับคู่เข้ากับรูปภาพของตัวอักษรที่สร้างขึ้นโดยโปรแกรมเมตาฟอนต์ โดยทำการแทรกคำสั่งของเท็กซ์เพื่อจัดการกับรหัสของตัวอักษรภาษาไทยเหล่านั้นก็จะสำเร็จเป็นโปรแกรมไทยเท็กซ์

การพัฒนาโปรแกรมไทยเท็กซ์จะสามารถแบ่งเป็นขั้นตอนการทำงานต่างๆ ได้ดังนี้

1. สร้างตัวอักษรภาษาไทยสำหรับใช้ในการแสดงผล โดยใช้โปรแกรมเมตาฟอนต์สร้างแฟ้มที่มีนามสกุล .mf ตามที่ทำการออกแบบไว้ให้เป็นแฟ้มข้อมูลในรูปแบบที่โปรแกรมไครเวอร์ (driver program) ของอุปกรณ์การแสดงผลแต่ละชนิดสามารถนำໄไปเปลี่ยนเป็นพิกเซลสำหรับแต่ละ อุปกรณ์ได้

เนื่องจากตัวอักษรแต่ละตัวก็คือรูปภาพที่ประกอบไปด้วยเส้นตรง เส้นโค้ง วงกลม วงรีต่างๆ ดังนั้น ถ้าสามารถกำหนดจุดต่างๆ ให้เหมาะสมเพื่อสร้างส่วนประกอบอย่างๆ เหล่านี้ เมื่อนำองค์ประกอบบ่ายโยมมาประกอบเข้าด้วยกันก็จะได้ตัวอักษรตามท้องการ

การสร้างองค์ประกอบของตัวอักษรดังกล่าววนีจะใช้วิธีการแทนที่ด้วยเวคเตอร์ การนำองค์ประกอบมารวมกันก็คือการนำเอาเวคเตอร์เหล่านี้มาร่วมไว้ที่ตำแหน่งต่างๆ อย่างเหมาะสม ดังนั้นตัวอักษรตัวเดียวกันไม่ว่าจะมีขนาดใดก็ตามจะประกอบด้วยเวคเตอร์กี่ลุ่มเดียวกันเสมอแต่ เวคเตอร์เหล่านั้นจะมีขนาดไม่เท่ากันการเปลี่ยนขนาดของตัวอักษรที่ทำได้โดยเปลี่ยนขนาดของเวคเตอร์ท่านั้น วิธีการทำเช่นนี้จะทำให้ตัวอักษรที่ได้มีคุณภาพเหมือนเดิมเสมอ ไม่ว่าจะขยายหรือ ย่อขนาดลงเท่าใดก็ตาม

ในหลักการทำงานของโปรแกรมเทกซันจะใช้หลักการของกล่องและการ (box and glue) ซึ่งก็คือ การนำเอาตัวอักษรแต่ละตัวมาเรียงต่อกัน โดยไม่จำกัดว่าจะต้องเรียงให้อยู่ในแนวเดียวกันหรือไม่ ตัวอักษรแต่ละตัวที่คือกล่อง 1 กล่อง หรืออาจจะนำกล่องเล็กๆเหล่านี้มาเรียงต่อๆกันแล้วทำการกำหนดเป็นกล่องที่ใหญ่ขึ้นก็ได้ การเรียงต่อ กันก็คือการนำมาติดกาวที่ด้านหลังและแปะลงไปบนกระดาษเปล่านั่นเอง แต่กาวนี้จะมีความยืดหยุ่นของตัวเอง โดยสามารถที่จะหดแคบลง หรือขยายขนาดให้กว้างขึ้นได้ตามความต้องการของผู้ออกแบบการที่จะนำสาระหรือวรรณยุกต์ไปเรียงกันหนึ่งหรือได้แนวบรรทัดปกติที่จำเป็นต้องใช้หลักการนี้กล่าวคือหากซึ่งไม่สนใจว่าตัวอักษรในกล่องนั้นจะมีหน้าตาเป็นอย่างไร เพียงแต่จะทำหน้าที่วางกล่องไว้ในตำแหน่งที่เหมาะสมเท่านั้น ดังนั้นการออกแบบตัวอักษรในกล่องจึงสามารถทำได้ตามความต้องการ

การกำหนดกล่องนั้นจะมี 3 มิติคือ ความกว้าง (width) ความสูง (height) และ ความลึก (depth) ความกว้างคือระยะทางตามแนวระดับ (horizontal dimension) ความสูง คือระยะทางตามแนวดิ่งเหนือเส้นฐาน (baseline) ความลึกคือระยะทางตามแนวดิ่งใต้เส้นฐาน ถึงแม้ว่าจะมีการกำหนดกล่องไว้ เช่นนี้แล้วตัวอักษรที่ออกแบบก็ไม่จำเป็นที่จะต้องอยู่ในกล่องเสมอไป อาจจะมีการออกแบบกล่องไปบ้างก็ได้ เช่นการออกแบบตัวอักษรที่เป็นตัวอักษรไทยที่อยู่ข้างหนึ่งของบรรทัดปกติของภาษาไทย จะมีบางส่วนของตัวอักษรนั้นอยู่นอกกล่อง เป็นต้น

การกำหนดค่ามิติ (dimension) ของกล่องนั้นจะกำหนดในรูปแบบของปริมาณชาร์ป (sharp quantities) กล่าวคือ ขนาดของกล่องจะไม่ขึ้นอยู่กับอุปกรณ์การแสดงผล โปรแกรมมาฟอนต์จะเป็นตัวจัดการให้อุปกรณ์นั้นรู้จักกับเอกสารที่สร้างขึ้นโดยการสร้างรูปแบบการทำงาน (mode) ใหม่ แต่ตัวอักษรนั้นจะยังคงมีมิติของกล่องเท่าเดิมเมื่อได้กล่องตามที่ต้องการแล้วก็จะกำหนดพิกเซล (pixel) ในกล่องนั้นการสร้างเวคเตอร์ต่างๆก็จะใช้ค่าของพิกเซลเหล่านี้เป็นจุดอ้างอิง เมื่อทำการสร้างตัวอักษรขึ้นมาตัวแล้วก็จะต้องมีการให้คำสั่งรับตัวอักษรตัวที่สร้างขึ้นนี้ด้วยเพื่อที่เทกซ์จะได้สามารถอ้างอิงได้ เช่นตัวอักษร "ก" จะแทนที่ด้วยค่า 161 เป็นต้น

หลังจากที่สามารถเขียนโปรแกรมสร้างตัวอักษรได้แล้วก็จะทำการแปลงชุดคำสั่งของโปรแกรมที่เขียนขึ้นให้อยู่ในรูปตัวอักษรแบบทั่วไป (generic font) ซึ่งสามารถเปลี่ยนเป็นตัวอักษรสำหรับอุปกรณ์แสดงผลใดๆก็ได้อย่างไรก็ตามตัวอักษรที่ได้นี้อาจจะยังไม่มีความละเอียดเหมาะสมกับอุปกรณ์แสดงผล ดังนั้นจำเป็นจะต้องอาศัยซอฟท์แวร์เปลี่ยนข้อมูลตัวอักษรแบบทั่วไปเหล่านี้ให้เป็นรูปแบบที่โปรแกรมไคร์เวอร์ช่วย

ในการแสดงผลสำหรับอุปกรณ์แสดงผลแต่ละชนิดจะรู้จักและนำออกไปสร้างผลลัพธ์ได้ในที่สุด ตัวอักษรที่ได้จากนั้นตอนนี้จะนำไปแสดงผลและใช้อ้างอิงในขั้นตอนต่อไป

2. สร้างโปรแกรมขึ้นมาเพื่ออ่านเอกสารที่เป็นภาษาไทย จัดการแยกระหว่างข้อมูลและคำสั่งของเทกซ์ แยกข้อมูลที่เป็นภาษาไทยและอังกฤษหรือตัวเลข จากนั้นทำการแทรกสัญลักษณ์พิเศษ เพื่อบอกการสินสุดคำภาษาไทย หลังจากนั้นเปลี่ยนตัวอักษรภาษาไทยให้อยู่ในรูปแบบ "\cXXX" (สำหรับเทกซ์ในเวอร์ชันแรกๆ จะไม่สามารถรู้จักตัวรหัสแอสกีโดยตรงต้องมีการใช้คำสั่ง "\charXXX" นำหน้า แต่ในเวอร์ชันต่อๆ มาเทกซ์จะสามารถรู้จักรหัสแอสกีและสามารถนำรหัสนั้นไปค้นหาตัวอักษรในฟอนต์ต่างๆ ได้เลย) เพื่อให้เข้ากันได้กับตัวอักษรภาษาไทยที่สร้างขึ้น โดยโปรแกรมเมตาฟ่อนต์ เมื่อสามารถเปลี่ยนให้อยู่ในรูปแบบดังกล่าวได้แล้ว จะทำการล้อมประโยคภาษาไทยนั้นด้วยเครื่องหมาย {\thai...} เพื่อบอกว่ารหัสเหล่านี้ต้องใช้กับฟอนต์ภาษาไทย

เนื่องจากการทำการตัดคำเพื่อแทรกสัญลักษณ์พิเศษนั้นอัลกอริธึมที่ใช้งานไม่มีประสิทธิภาพเทียบเท่ากับอัลกอริธึมตัดคำที่ใช้อยู่ในโปรแกรมประมวลผลคำ CU Writer ทำให้รูปประโยคที่ได้ไม่มีความสมบูรณ์เพียงพอ ดังนั้นถ้าสามารถนำอัลกอริธึมตัดคำของ CU Writer มาใช้ก็จะทำให้การตัดคำมีประสิทธิภาพดียิ่งขึ้น

การทำงานของโปรแกรมแบ่งพยางค์ของ CU Writer จะใช้วิธีการค้นหาจุดแบ่งพยางค์ให้ได้มากที่สุด และผิดพลาดน้อยที่สุดด้วยพจนานุกรมที่มีโครงสร้างแบบทรรย (trie) โดยใช้วิธีการสร้างต้นไม้แบ่งคำ (separation tree) สำหรับคำที่จะแบ่ง ทรรย (trie) คือ ต้นไม้ (tree) ชนิดหนึ่งซึ่งจะมีลักษณะของข้อมูลเป็นลำดับ จากเซ็ตของตัวอักษร โดยตรง แทนที่จะมองแต่ละคำเป็นข้อมูลหนึ่งชิ้น ทำให้มีความสามารถจัดการกับข้อมูลที่มีความยาวไม่คงที่ได้ดีและเข้ากันได้ดีกับโปรแกรมที่ทำงานกับอักษรที่ละตัวอักษร เช่น โปรแกรมแบ่งคำนี้ นอกจากนี้ทรรยยังมีความสามารถในการจัดการกับคำนำหน้า (prefix) ร่วมของคำในพจนานุกรมโดยการยุบรวมกันหมวด ทำให้ลดขนาดของพจนานุกรมลงได้เป็นอันมาก

การเขียนต้นไม้แบ่งคำจะเป็นคำที่ต้องการจะแบ่งทีล่ะ 1 ตัวอักษรจากบนลงล่าง เมื่อถึงจุดที่สามารถแบ่งคำได้ก็จะเกิดการตัดสินใจว่าจะแบ่งหรือไม่ ถ้าตัวอักษรที่เขียนลงไปทางกิ่งด้านขวาของต้นไม้ แบ่งคำก็จะยังไม่แบ่งคำและตัวอักษรตัดไปก็ยังเป็นคำเดียวกัน แต่ถ้าไปทางซ้ายมีหมายถึงการแบ่งคำ และตัวอักษรตัดไปจะถือเป็นตัวเริ่มต้นคำใหม่และหมายถึงการได้จุดแบ่งคำเพิ่มขึ้น อีกหนึ่งจุดด้วย

สำหรับอัลกอริธึมในการค้นหาข้อมูลในต้นไม้ของ CU Writer จะค้นหาแบบทางลึกก่อน (depth-first search) ร่วมกับการทำงานแบบย้อนกลับ (backtrack) เรียกว่า BWS ในงานวิจัยนี้จะนำเอาอัลกอริธึมนี้มาประยุกต์ให้สามารถลืมกรอบคำที่ต้องการเบ่งด้วยชุดคำสั่งของเทิกซ์

3. ทำการแทรกชุดคำสั่งของเทิกซ์เข้าไปให้เหมาะสมกับงานที่ต้องการ เช่น การจัดระยะห่างบรรทัดของภาษาไทย การจัดช่องไปต่างๆ การจัดระยะห่างของช่องว่างของแต่ละย่อหน้า การจัดบรรทัดสำหรับหมายเหตุที่เป็นภาษาไทย สิ่งต่างๆเหล่านี้จะมีความแตกต่างจากการใช้ตัวอักษรที่เป็นภาษาอังกฤษ นอกจากนี้ต้องทำการตรวจสอบสระหรือวรรณยุกต์ที่ต้องมีการปรับระดับให้อยู่เหนือหรือใต้พยัญชนะนั้น แทรกชุดคำสั่งของเทิกซ์เพื่อวางตำแหน่งของสระและวรรณยุกต์เหล่านั้นให้เหมาะสม

หลักการวางแผนของสร่านั้นจะสามารถทำได้ในขั้นตอนของการออกแบบตัวอักษรด้วยโปรแกรมเมตาฟอนต์โดยจะกำหนดให้ตัวอักษรอยู่นอกกรอบไปทางด้านซ้ายมือของจุดเริ่มต้นและกำหนดพิกัดสำหรับตัวอักษรเหล่านั้นด้วยค่าที่เป็นลบ สำหรับในกรณีของสระที่อยู่ข้างบนและมีวรรณยุกต์ ประกอบด้วยนั้นก็จะต้องทำการปรับวรรณยุกต์ให้ไม่ทับซ้อนกับตำแหน่งของสร่านั้นด้วย

วัตถุประสงค์ในการทำวิจัย

1. ศึกษาการทำงานของโปรแกรมเทิกซ์
2. ศึกษาวิธีการสร้างตัวอักษรโดยการใช้โปรแกรมเมตาฟอนต์
3. สร้างตัวอักษรไทยโดยใช้โปรแกรมเมตาฟอนต์
4. ออกแบบและพัฒนาโปรแกรมแสดงผลตัวอักษรไทยระบบเทิกซ์

ขอบเขตการทำวิจัย

1. ดัดแปลงอัลกอริธึมตัดคำภาษาไทยของโปรแกรม CU Writer ให้สามารถใช้กับโปรแกรมแสดงผลตัวอักษรไทยบนเทิกซ์ได้

2. โปรแกรมสามารถทำงานบนไมโครคอมพิวเตอร์ที่ใช้ระบบปฏิบัติการ MS-DOS

ขั้นตอนการทําวิจัย

1. ศึกษาการทำงานของ โปรแกรมเท็กซ์
2. ศึกษาการทำงานของ โปรแกรมเมตาฟอนต์
3. สร้างตัวอักษรภาษาไทยสำหรับใช้ในการทดสอบ โปรแกรมไทยเท็กซ์
4. ศึกษาอัลกอริธึมตัดคำภาษาไทยของ CU Writer
5. ศึกษาวิธีการสร้างชุดคำสั่งของเท็กซ์ เพื่อจัดการกับปัญหาของภาษาไทย 3 ระดับ
6. สร้าง โปรแกรมเปลี่ยนแฟ้มบรรณाचิกร ให้เป็นแฟ้มข้อมูลในรูปแบบของเท็กซ์
7. ทดสอบการใช้งานของ โปรแกรมที่สร้างขึ้น โดยใช้ โปรแกรมเท็กซ์
8. แก้ไขในรายละเอียดของ โปรแกรมที่สร้างขึ้นรวมทั้งตกแต่งตัวอักษรให้สวยงาม
9. ทดสอบอีกรอบเพื่อความถูกต้อง
10. สรุปผลการวิจัยและจัดทำเอกสารประกอบการวิจัย

ประโยชน์ที่คาดว่าจะได้จากการวิจัย

1. จะได้ฟอนต์ตัวอักษรภาษาไทยที่สร้างขึ้นด้วย โปรแกรมเมตาฟอนต์
2. จะได้โปรแกรมแสดงผลตัวอักษรไทยบนระบบเท็กซ์

**ศูนย์วิทยทรพยากร
จุฬาลงกรณ์มหาวิทยาลัย**