

บทที่ 2

แนวคิดและทฤษฎี

ในการที่จะออกแบบและพัฒนาโปรแกรมแสดงผลตัวอักษรไทยด้วยระบบเท็กซ์นั้นจำเป็นที่จะต้องมีความรู้พื้นฐานเกี่ยวกับองค์ประกอบที่จำเป็นในการสร้างโปรแกรมไทยเท็กซ์เสียก่อน

องค์ประกอบที่จำเป็นในการทำงานของโปรแกรมเท็กซ์โดยทั่วไปนั้นสามารถแบ่งออกได้เป็น 4 ส่วน ดังนี้

1. ส่วนที่ทำหน้าที่จัดการกับแฟ้มเอกสารให้อยู่ในรูปแบบที่ต้องการจะพิมพ์ ในส่วนนี้จะหมายถึงตัวโปรแกรมของเท็กซ์เองซึ่งจะทำหน้าที่เปลี่ยนแฟ้มเอกสารที่อยู่ในรูปแบบไวยากรณ์ที่เท็กซ์เข้าใจให้ไปอยู่ในรูปที่อุปกรณ์แสดงผลรู้จักได้

2. ส่วนของตัวอักษร (font) ที่ใช้ ซึ่งในส่วนนี้จะมีโปรแกรมเมตาฟอนต์เป็นตัวจัดการสร้างขึ้น การเรียกใช้ฟอนต์ของโปรแกรมเท็กซ์นั้น อาจจะใช้ฟอนต์ที่สร้างจากโพสต์คริปต์(postscript) หรือฟอนต์ประเภททรูไทป์ (true type) ได้แต่ต้องมีโปรแกรมที่เปลี่ยนรูปแบบแฟ้มที่ได้จากการ ทำงานของเท็กซ์ให้เป็นในรูปแบบที่ฟอนต์ชนิดนั้นรู้จัก

3. ส่วนของซอฟต์แวร์อรรถประโยชน์ (utilities software) เป็นส่วนของโปรแกรมสนับสนุนในการทำงานของเท็กซ์และเมตาฟอนต์ เช่นการแปลงรูปแบบของแฟ้มจากประเภทหนึ่งไปเป็นอีกประเภทหนึ่ง โปรแกรมในส่วนนี้จะทำให้ผู้ใช้งานเท็กซ์มีความสะดวกในการที่จะนำแฟ้มเอกสารที่จัดทำด้วยเท็กซ์หรือเมตาฟอนต์ที่ทำงานอยู่บนระบบปฏิบัติการหนึ่งสามารถเปลี่ยนไปใช้งานกับระบบปฏิบัติการอื่นๆ ได้

4. ส่วนของซอฟต์แวร์ที่ควบคุมการแสดงผลออกทางอุปกรณ์ต่างๆ (output device driver software) เป็นส่วนของโปรแกรมที่ดำเนินการนำแฟ้มที่ได้จากกระบวนการของเท็กซ์ไปแสดงผลลงบนอุปกรณ์ต่างๆ ในส่วนนี้เป็นเรื่องของอุปกรณ์แสดงผลแต่ละชนิดที่ต้องจัดหาซอฟต์แวร์ไดรเวอร์ซึ่งขึ้นอยู่กับฮาร์ดแวร์แต่ละชนิด

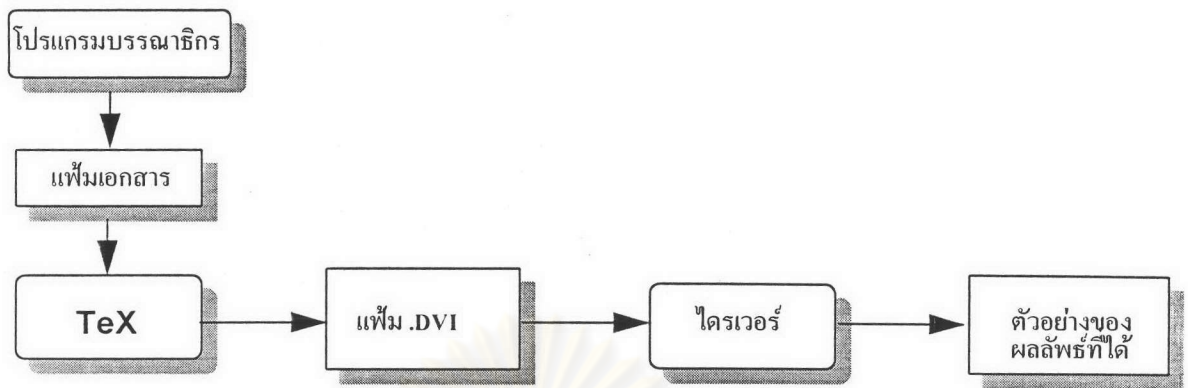
ด้วยองค์ประกอบข้างต้นซึ่งถือเป็นหัวใจในการพัฒนาระบบเท็กซ์นี้แสดงให้เห็นว่า ระบบเท็กซ์จะมีความยืดหยุ่นกับการใช้และพัฒนาค่อนข้างมาก เช่นสมมุติว่ามีอุปกรณ์แสดงผลแบบใหม่เกิดขึ้นและต้องการที่จะแสดงผลออกทางอุปกรณ์นั้นก็เพียงแค่สร้างดีไวซ์ไดรเวอร์สำหรับอุปกรณ์ชนิดนั้นขึ้นมาเท่านั้น หรือถ้าต้องการฟอนต์ชนิดใหม่ก็อาจจะใช้เมตาฟอนต์เป็นตัวสร้างฟอนต์ขึ้นมาตามความพอใจได้ หรือถ้าไม่ต้องการที่จะสร้างฟอนต์เองแต่ต้องการที่จะใช้ฟอนต์ประเภททรูไทป์ชนิดต่างๆ ก็เพียงแค่หาซอฟต์แวร์ที่สามารถแปลงรูปแบบของแฟ้มที่เกิดขึ้นในระบบเท็กซ์ให้เป็นรูปแบบที่เหมาะสมกับการใช้งานของฟอนต์ประเภทนั้นได้เช่นกัน

การพัฒนาโปรแกรมแสดงผลภาษาไทยบนระบบเท็กซ์นี้ก็อาศัยองค์ประกอบที่กล่าวข้างต้นนี้เช่นกัน โดยจะเข้าไปข้องเกี่ยวเฉพาะการพัฒนาตัวของเท็กซ์และการพัฒนาฟอนต์เท่านั้น ซึ่งวิธีการพัฒนาจะนำเสนอในบทต่อไป อย่างไรก็ตามการพัฒนาโปรแกรมเพื่อที่จะจัดเรียงพิมพ์ด้วยระบบเท็กซ์หรือการประมวลผลค่าด้วยโปรแกรมประมวลผลค่าใดๆก็ตาม สิ่งที่ใช้ต้องการนอกเหนือไปจากความสวยงามของชิ้นงานแล้ว ความสละสลวยของภาษาที่ได้ในงานพิมพ์นั้นก็เป็นที่ผู้ใช้ไม่สามารถมองข้ามได้เช่นกัน ในตัวของเท็กซ์เองจะมีอัลกอริทึมที่จัดการกับการแบ่งคำที่เป็นภาษาอังกฤษด้วยไฮเฟน โดยเฉพาะ ดังนั้นถ้าต้องการที่จะพัฒนาโปรแกรมแสดงผลตัวอักษรภาษาไทยซึ่งต่อไปจะเรียกแทนว่าโปรแกรมไทยเท็กซ์จะไม่สามารถหลีกเลี่ยงในเรื่องของการแบ่งคำหรือจัดคำในแต่ละบรรทัด ได้เช่นกัน จึงจำเป็นที่จะต้องศึกษาแนวความคิดและทฤษฎีของการตัดคำนี้รวมอยู่ด้วย

จากที่กล่าวมาสามารถแยกหัวข้อของการศึกษาเพื่อที่จะพัฒนาโปรแกรมไทยเท็กซ์ได้เป็น 3 เรื่องใหญ่ คือ การศึกษาถึงการทำงานของเท็กซ์ การทำงานของเมตาฟอนต์ และ การตัดคำ

1. การทำงานของ TeX

หน้าที่ของระบบเท็กซ์ก็คือการเรียกใช้งานของเครื่องมือที่มีอยู่ในระบบของเท็กซ์เอง เช่น ไฟล์ชนิดต่างๆ ดีไวซ์ไดรเวอร์ต่างๆ นำมาผลิตเป็นเอกสารที่มีรูปแบบซับซ้อนตามที่ต้องการ ดังนั้นการที่เราจะสร้างเอกสารให้มีรูปแบบตามที่ต้องการ ได้จำเป็นจะต้องรู้จักกับส่วนประกอบหลักของเท็กซ์ซึ่งแสดงไว้ดังรูปที่ 2.1



รูปที่ 2.1 แสดงการทำงานโดยทั่วไปของ โปรแกรมเท็กซ์

1.1 การสร้างเพิ่มเอกสาร(Document File)

เป็นขั้นตอนเริ่มแรกของการทำงานในการสร้างรูปแบบการพิมพ์ของระบบเท็กซ์ ซึ่งเราสามารถเลือกโปรแกรมบรรณาธิการที่จะใช้สร้างได้แทบจะทุกประเภทตามแต่ที่จะมีในระบบปฏิบัติการของเครื่องคอมพิวเตอร์ชนิดนั้นๆ อย่างไรก็ตามควรที่จะเลือกโปรแกรมบรรณาธิการที่สร้างเพิ่มในลักษณะของเท็กซ์ไฟล์ กล่าวคือเป็นเพิ่มข้อมูลที่เก็บข้อมูลแบบสตริงของอักขระ (Character String) ส่วนเพื่อที่จะไม่ต้องมีปัญหาของส่วนหัวของไฟล์ (File Header) ซึ่งระบบเท็กซ์อาจจะไม่รู้จัก

1.2 การประมวลผลด้วยเท็กซ์

หลังจากที่ได้จัดเตรียมเพิ่มเอกสารเรียบร้อยแล้วก็ถึงขั้นตอนของการเรียก โปรแกรมเท็กซ์เพื่อที่จะทำการประมวลผลเพิ่มเอกสารที่ได้จัดเตรียมไว้ให้เปลี่ยนสภาพเป็นเพิ่มที่ไม่ขึ้นกับอุปกรณ์ชนิดใดๆ (Device Independent File) ซึ่งหลังจากผ่านขั้นตอนนี้แล้วจะได้เพิ่มที่มีนามสกุล .dvi ในขั้นตอนนี้ดังกล่าวนี้อาจจะมีการเกิดความผิดพลาดขึ้นได้เนื่องจากการสะกดคำสั่งผิด การผิดไวยากรณ์ของเท็กซ์ หรืออาจจะเกิดจากความไม่เข้าใจ โครงสร้างของชุดคำสั่งของเท็กซ์ดีพอกก็ได้ ข้อผิดพลาดดังกล่าวนี้จะแสดงออกมาทางจอภาพและจะมีการบันทึกลงได้ในเพิ่มที่มีนามสกุล.log เพื่อให้ผู้ใช้สามารถดูและทำการแก้ไขเพิ่มเอกสารให้ถูกต้องต่อไป

โปรแกรมเท็กซ์ที่ใช้ประมวลผลนั้นมีมากมายหลายรุ่นขึ้นอยู่กับระบบคอมพิวเตอร์ที่ใช้เช่น TeX และ LaTeX สำหรับใช้บนระบบปฏิบัติการยูนิกซ์, emTeX สำหรับระบบปฏิบัติการแบบดอส เป็นต้น นอกจากนี้ยังได้มีการพัฒนาออกไปเพื่อให้ใช้ได้กับวินโดวส์และแมคอินทอชได้ด้วย แต่ที่ถือเป็นต้นแบบของระบบนี้ก็คือ TeX สำหรับโปรแกรม LaTeX นั้นเกิดจากการพัฒนาชุดคำสั่งของ TeX ให้ผู้ใช้สามารถใช้งานได้ง่ายและสะดวกยิ่งขึ้น ดังนั้นรูปแบบของคำสั่งบางคำสั่งจึงแตกต่างกันไปบ้าง

1.3 การใช้คำสั่งมาโคร (Macro Instruction)

คำสั่งควบคุมลำดับการทำงานต่างๆที่เราพิมพ์แทรกเข้าไปในแฟ้มเอกสารนี้ จะถูกกำหนดโดยชุดของมาโคร (Macro Package) ชุดของมาโครนี้ประกอบด้วยคำสั่งต่างๆที่ใช้ในเท็กซ์ซึ่งอาจจะถูกพัฒนาขึ้นโดยภาษาระดับสูง(ภาษาซี) หรืออาจจะถูกพัฒนาขึ้นจากชุดคำสั่งง่ายๆ มาเรียงต่อกัน เพื่อให้สามารถทำงานอย่างใดอย่างหนึ่งได้สะดวกยิ่งขึ้น โดยปกติชุดของมาโครจะถูกเก็บอยู่ในแฟ้มรูปแบบ (Format File) เมื่อเวลาโปรแกรมเท็กซ์เริ่มทำงานจะทำการแปลชุดคำสั่งเหล่านี้แล้วเก็บไว้ในแฟ้มรูปแบบเพื่อที่จะได้ใช้ตอนทำการประมวลผลด้วยเท็กซ์ต่อไป คำสั่งที่ถูกพัฒนาขึ้นโดยภาษาซีที่ติดมากับโปรแกรมเท็กซ์นั้นเรียกว่าคำสั่งดั้งเดิม (Primitive Instruction) คำสั่งประเภทนี้จะไม่สามารถแก้ไขเป็นอย่างอื่นได้นอกจากจะทำการพัฒนาระบบเท็กซ์ขึ้นมาใหม่ทั้งระบบ คำสั่งประเภทนี้ได้แก่ คำสั่งที่จัดการเกี่ยวกับรูปแบบของอินพุตและเอาท์พุตของตัวอักษร เป็นต้น ส่วนคำสั่งอีกประเภทหนึ่งนั้นจะถูกพัฒนาขึ้นมาจากคำสั่งดั้งเดิมเรียกว่าคำสั่งทั่วไป (Plain Instruction) ในโปรแกรม TeX จะเก็บคำสั่งเหล่านี้ไว้ในแฟ้มข้อมูลที่มีชื่อว่า plain.tex และในโปรแกรม LaTeX จะเก็บไว้ในแฟ้ม latex.tex คำสั่งทั่วไปนี้มีข้อดีคือรูปแบบการใช้งานจะง่ายกว่าแบบคำสั่งดั้งเดิมแต่มีข้อเสียอยู่ตรงที่จะสิ้นเปลืองเวลาในการประมวลผลมากกว่าเนื่องจากจะต้องทำการขยายคำสั่งที่ซ่อนอยู่ภายในออกมาให้อยู่ในรูปของคำสั่งดั้งเดิมหลายๆคำสั่งก่อนจึงจะทำการประมวลผลได้ อย่างไรก็ตามถ้าผู้ใช้งานมีความรู้ในเรื่องคำสั่งมาโครดีเพียงพอก็อาจจะเขียนชุดคำสั่งนี้ให้ทำงานตามที่ต้องการได้เช่นการสร้างโปรแกรม LaTeX นั้นจะอาศัยการเขียนคำสั่งมาโครที่โปรแกรม TeX มีมาให้แต่เพียงอย่างเดียวโดยไม่เข้าไปยุ่งเกี่ยวกับโครงสร้างหรือระบบภายในของตัวโปรแกรม TeX แม้แต่น้อย รายละเอียดของโปรแกรมดังกล่าวสามารถศึกษาได้จากแฟ้ม latex.tex

นอกจากนี้ การพัฒนาโปรแกรมไทยเท็กซ์นั้น ส่วนหนึ่งของการพัฒนาก็จำเป็นต้องอาศัยการเขียนชุดคำสั่งมาโครด้วยเช่นเดียวกัน แต่การที่ไม่สามารถพัฒนาโปรแกรมไทยเท็กซ์ทั้งหมดด้วยการเขียนชุดคำสั่งมาโคร

โครนั้นก็เนื่องจากในส่วนของการพัฒนาอรรถิทีมสำหรับตัดคำภาษาไทยนั้นจำเป็นต้องใช้ภาษาที่มีความคล่องตัวมากกว่าชุดคำสั่งมาโครของเท็กซ์ รายละเอียดดังกล่าวจะได้นำเสนอในบทต่อไป

1.4 การเรียกใช้ตัวอักษรแบบต่างๆ

จุดเด่นที่น่าสนใจของระบบเท็กซ์ก็คือลักษณะของตัวอักษรที่เรียกใช้ โดยปกติแล้วในซอฟต์แวร์ของโปรแกรมเท็กซ์เองจะมีชุดของตัวอักษรหรือฟอนต์ให้เลือกใช้มากมายหลายขนาดตามต้องการ แต่ก็ยังมีซอฟต์แวร์ให้ผู้ใช้สามารถสร้างฟอนต์ขึ้นมาใช้เองได้อีกด้วย ฟอนต์ที่จะถูกเรียกใช้งานโดยโปรแกรมเท็กซ์จะต้องอยู่ในรูปของเมตริกฟอนต์ (TeX Font Matrix) หรือแฟ้มที่มีนามสกุลเป็น .tfm ซึ่งในแฟ้มนี้จะเก็บรายละเอียดของฟอนต์แต่ละตัวที่ใช้เช่นความกว้าง ความสูง รูปร่าง เป็นต้น ซอฟต์แวร์ที่ใช้สำหรับสร้างฟอนต์นี้คือ METAFONT ซึ่งจะได้กล่าวถึงต่อไป

ตามปกติแล้วถ้าไม่ใช่คำสั่งอื่น ในชุดคำสั่งมาโครจะเลือกใช้ฟอนต์ชนิดคอมพิวเตอร์โมเดิร์นฟอนต์ (Computer Modern fonts) เสมอจนทำให้เข้าใจว่าเท็กซ์จะใช้ได้แต่เพียงฟอนต์ชนิดนี้เพียงอย่างเดียวซึ่งเป็นสิ่งที่เข้าใจผิดเป็นอย่างมากการจะใช้ฟอนต์ชนิดอื่นนั้นจะกระทำได้โดยไม่ยากนักแต่บางครั้งก็ทำให้เกิดความสับสนบ้างและก็จะอาจทำให้งานพิมพ์ที่ได้ผิดไปจากความต้องการของผู้ใช้ ขั้นตอนการทำงานของเท็กซ์ในการเรียกใช้ฟอนต์ชนิดต่างๆมีดังนี้

1.4.1 กำหนดชนิดของฟอนต์ด้วยคำสั่งมาโคร

ขั้นตอนนี้จะเป็นการบอกให้คำสั่งมาโครของเท็กซ์ทราบว่าชนิดของฟอนต์ที่จะใช้นั้นเป็นชนิดใดเช่น สมมุติว่ามีฟอนต์ตัวอักษรภาษาไทยอยู่ชุดหนึ่งซึ่งเก็บอยู่ในแฟ้มที่มีชื่อว่า san.tfm ถ้าต้องการจะเรียกใช้ฟอนต์ชุดนี้ในแฟ้มเอกสารของก็สามารถทำได้โดยสร้างคำสั่งมาโครขึ้นมาใหม่มีชื่อว่า thai จากนั้นก็กำหนดค่าให้มาโครที่มีชื่อว่า thai นี้เป็นชื่อของฟอนต์ san.tfm ดังนี้

```
\font\thai=san
```

หลังจากนั้นไม่ว่าเมื่อใดที่เราต้องการใช้ฟอนต์ชนิดนี้ก็เพียงแต่เรียกใช้คำสั่งมาโคร เท่านั้น ฟอนต์ที่ตามหลังคำสั่งดังกล่าวก็จะเป็นฟอนต์ชนิด san ทั้งหมดเว้นแต่จะมีการกำหนดฟอนต์ชนิดอื่นต่อไปอีก

1.4.2 เท็กซ์จะทำการโหลดรายละเอียดของฟอนต์ชนิดที่กำหนดไว้ในแฟ้ม san.tfm

ขั้นตอนนี้เท็กซ์จะทำการค้นหาแฟ้มดังกล่าวในไดเรกทอรีที่กำหนดไว้ ตอนเริ่มติดตั้งระบบ สิ่งที่ต้องระลึกไว้เสมอก็คือ เท็กซ์จะต้องการเฉพาะแฟ้มที่มีนามสกุล .tfm เท่านั้น ถ้าไม่เจอแฟ้มดังกล่าว เท็กซ์จะแสดงข้อผิดพลาดและหยุดการทำงาน

1.4.3 เท็กซ์จะทำการสร้างแฟ้ม .dvi ขึ้น

ในแฟ้ม .dvi จะเก็บรายละเอียดเกี่ยวกับชนิด อัตรายายของฟอนต์ที่ใช้

1.4.4 ไดรว์เวอร์จะทำหน้าที่จัดการกับแฟ้ม DVI

โดยจะกำหนดแฟ้มของฟอนต์ (font file) สำหรับฟอนต์แต่ละตัวตามขนาดของฟอนต์ที่เรียกใช้ ไดรว์เวอร์บางตัวอาจจะตรวจสอบกับรายชื่อของฟอนต์ที่มีอยู่ในตัวของเครื่องพิมพ์แต่ละชนิดด้วย ก่อนที่จะทำการกำหนดแฟ้มของฟอนต์ นอกจากนี้ไดรว์เวอร์บางตัวยังมีความสามารถที่จะเลือกฟอนต์ทดแทนได้ในกรณีที่ตรวจสอบแล้วไม่พบฟอนต์ที่อ้างถึงแฟ้มที่ไดรว์เวอร์ต้องการนั้นจะมีนามสกุลแตกต่างกันไปเช่น .pk, .pxl, .gf เป็นต้น ทั้งนี้ขึ้นอยู่กับชนิดของไดรว์เวอร์นั้นๆ สำหรับในการวิจัยครั้งนี้ใช้โปรแกรม emTeX เป็นตัวทดสอบ และไดรว์เวอร์ที่ติดมากับซอฟต์แวร์ชุดนี้จะต้องการใช้แฟ้ม .pk

1.4.5 ไดรว์เวอร์จะทำหน้าที่ผลิตเอาต์พุตออกมาให้เหมาะสมกับอุปกรณ์เอาต์พุตประเภทนั้นๆ

ข้อที่นำจดจำอีกข้อหนึ่งก็คือเท็กซ์จะไม่ใช้ตัวที่ทำหน้าที่ผลิตเอาต์พุตโดยตรง เท็กซ์จะทำหน้าที่แต่เพียงสร้างแฟ้ม .dvi เท่านั้น และแฟ้ม .dvi นี้ก็ไม่สามารถแสดงผลออกมาเป็นเอาต์พุตได้เลยในทันที จะต้องนำไปผ่านไดรว์เวอร์เสียก่อนจึงจะออกไปเป็นเอาต์พุตได้ การแสดงผลของไดรว์เวอร์จะไม่เกี่ยวข้องกันกับตัวโปรแกรมเท็กซ์แต่จะเกี่ยวข้องกับฮาร์ดแวร์ของอุปกรณ์นั้นๆ ดังนั้นถ้าเรามีแฟ้มเอกสารที่จัดเตรียมไว้ด้วยโปรแกรมเท็กซ์เราสามารถนำไปประมวลผลด้วยระบบเท็กซ์ที่มีลักษณะของไวยากรณ์หรือโครงสร้างของคำสั่งใกล้เคียงกัน หรือเหมือนกันได้ทุกระบบไม่ว่าจะเป็นระบบปฏิบัติการแบบใดก็ตาม และจะให้ผลออกมาเป็นแฟ้ม .dvi เช่นเดียวกัน

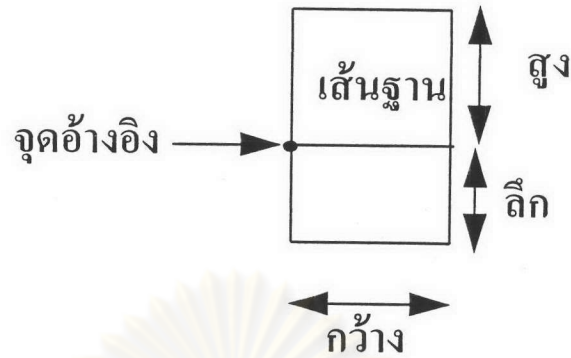
2. การทำงานของเมตาฟอนต์

เนื่องจากการออกแบบการทำงานของเท็กซ์จะใช้หลักการของกล่องและกาว กล่าวคือเท็กซ์จะไม่สนใจว่ามีอะไรอยู่ในกล่องบ้าง แต่จะสนใจว่าจะนำกล่องแต่ละกล่องมาประกอบรวมกันอย่างไรเท่านั้น ดังนั้นถ้าเท็กซ์ทำการอ้างอิงถึงกล่องดังกล่าวแต่ละกล่องได้ก็สามารถดำเนินการนำมาวางประกอบกันได้เสมอ

กล่องแต่ละกล่องอาจจะประกอบไปด้วยตัวอักษร(letter)หรือสัญลักษณ์ (symbol)คำ (word) ประโยค (sentence) หรือย่อหน้า (paragraph) ก็ได้ สำหรับหน่วยเล็กที่สุดที่เท็กซ์ถือว่าเป็นกล่องก็คือตัวอักษรแต่ละตัวหรือสัญลักษณ์

คำว่ากล่อง แท้จริงแล้วเป็นเพียงกรอบที่มีเพียง 2 มิติที่มองเห็นได้ในระนาบเดียวเท่านั้น แต่ในความหมายของเท็กซ์นี้จะมีมิติคือความกว้าง ความสูง และความลึกโดยมีการกำหนดจุดอ้างอิงหนึ่งจุด เส้นที่ลากผ่านจุดนี้ในแนวระดับจะเรียกว่าเส้นฐาน มิติความกว้างก็คือระยะทางที่ห่างจากจุดอ้างอิงไปทางขวา มิติความยาวคือ ระยะทางเหนือเส้นฐาน สำหรับความลึกนั้นเท็กซ์กำหนดให้คือระยะทางใต้เส้นฐาน องค์ประกอบต่างๆของกล่องจะเป็นดังรูปที่ 2.2 ในกล่องนี้จะบรรจุตัวอักษรหรือสัญลักษณ์ที่สร้างขึ้น ตัวอักษรหรือสัญลักษณ์ดังกล่าวอาจจะอยู่ในกล่องที่สร้างขึ้นทั้งหมดหรืออาจจะมีบางส่วนที่เกินออกมานอกกล่องบ้างก็ได้สำหรับตัวอักษรหนึ่งชุนั้นไม่มีข้อกำหนดว่ามิติของแต่ละกล่องจำเป็นต้องเท่ากัน มิติของกล่องแต่ละกล่องนั้นจะขึ้นกับรูปร่างลักษณะของตัวอักษรแต่ละตัวตลอดจนความสวยงามของตัวอักษรแต่ละชุนเป็นเกณฑ์สำคัญ การอ้างอิงถึงกล่องแต่ละกล่องจะใช้รหัสแอสกี สิ่งที่ทำให้เท็กซ์สามารถเรียกใช้กล่องเหล่านี้ได้ก็คือการอ้างอิงถึงรหัสแอสกีนั่นเอง ดังนั้นเท็กซ์จะไม่มีโอกาสรู้ว่าสิ่งที่อยู่ในกล่องจะมีรูปร่างหน้าตาเป็นอย่างไร เพียงแต่รู้จักอ้างอิงถึงรหัสสำหรับตัวอักษรหรือสัญลักษณ์ก็เพียงพอแล้ว

จุฬาลงกรณ์มหาวิทยาลัย

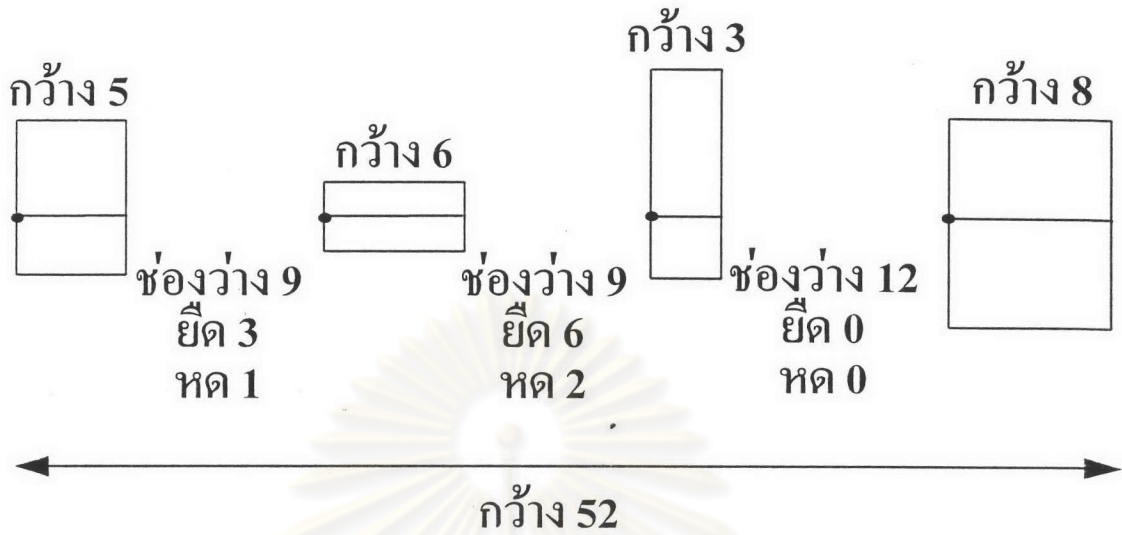


รูปที่ 2.2 แสดงภาพองค์ประกอบของกล่อง

สิ่งที่สำคัญสำหรับเท็กซ์ก็คือการนำกล่องเหล่านี้มาเรียงต่อกันเพื่อให้เป็นรูปแบบของงานพิมพ์ตามความต้องการของผู้ใช้ การจัดเรียงกล่องที่ต้องการนี้ก็อาจจะเปรียบเทียบได้กับการจัด เรียงพิมพ์หรือแสดงรูปภาพหรือข้อความบนบอร์ดตามงานนิทรรศการต่างๆกล่าวคือเริ่มต้นเราจะนำรูปภาพเหล่านั้นมากำหนดเสียก่อนว่าจะวางลงในตำแหน่งใด หลังจากที่ได้ตำแหน่งตามที่ต้องการ ทากาว แล้วจึงนำภาพหรือข้อความนั้นแปะทับลงบนกระดานว่างนั้น บางครั้งในหนึ่งรูปภาพที่นำมาแปะอาจจะมีข้อความหรือรูปภาพย่อยๆ ติดกาวอยู่บนรูปภาพนั้นก่อนแล้วก็เป็นได้ หลักการทำงานของเท็กซ์ก็เป็นเช่นเดียวกัน การกำหนดตำแหน่งที่สัมพันธ์กันของแต่ละกล่องนั้นจะอ้างอิงจากเส้นฐานเป็นเกณฑ์ ความแตกต่างอย่างหนึ่งของการติดรูปภาพบนบอร์ดและการเรียงพิมพ์ของเท็กซ์นั้นก็คือ เท็กซ์จะมีระยะยัดหรือระยะหดเพื่อให้ช่องไฟระหว่างตัวอักษรแต่ละตัวเป็นไปอย่างสวยงามตามขนาดตัวอักษรที่ใช้ในแต่ละย่อหน้า ในขณะที่รูปภาพที่ติดอยู่บนบอร์ดจะถูกกำหนดตายตัวไม่สามารถเว้นระยะยัดได้อีก

เพื่อที่จะทำความเข้าใจของหลักการของกาวในระบบเท็กซ์ได้ดียิ่งขึ้นจะพิจารณาได้ จากตัวอย่างต่อไปนี้

สมมุติว่ามีกล่องอยู่ 4 กล่องและมีช่องว่างสำหรับติดกาวอยู่ 3 ช่องดังรูปที่ 2.3



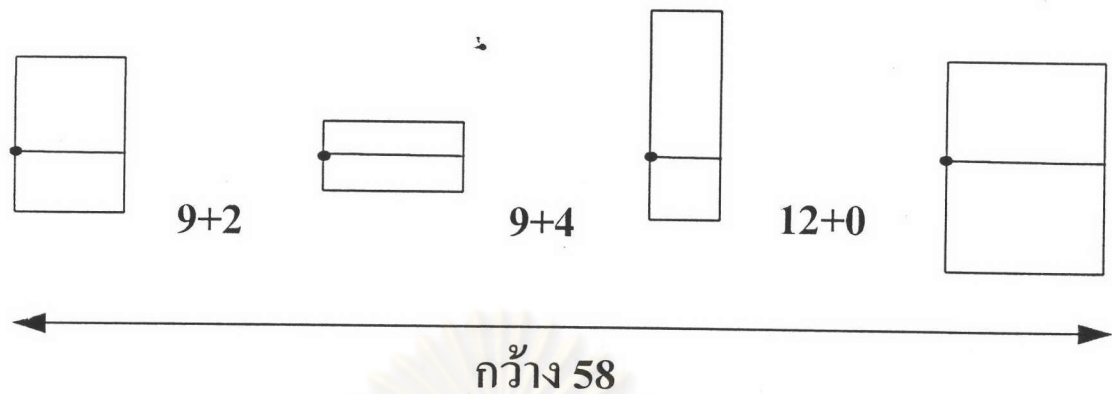
รูปที่ 2.3 แสดงภาพตัวอย่างของกล่องและกาว

ช่องแรกของกาวจะมีขนาดเปรียบเทียบเป็นช่องว่างได้ 9 หน่วยและมีระยะยึด 3 หน่วย ระยะหด 1 หน่วย ช่องที่สองจะมีขนาด 9 หน่วย ระยะยึด 6 หน่วย ระยะหด 2 หน่วย ช่องสุดท้าย จะมีขนาด 12 หน่วย ระยะยึด 0 หน่วย ระยะหด 0 หน่วยซึ่งหมายความว่าไม่สามารถยึดหรือหดได้

ความกว้างของกล่องและกาวในภาพนี้พิจารณาเฉพาะช่องว่างรวมทั้งหมดคือ

$$5+9+6++9+3+12+8 = 52 \text{ หน่วย}$$

ซึ่งจะเรียกว่าความกว้างธรรมชาติ (natural width) ในแนวราบ อย่างไรก็ตามก็อย่างที่เห็นซึ่งจะทำให้ความกว้างรวมทั้งหมดเป็น 58 หน่วยได้ถ้ามีการนำส่วนยึดที่กำหนดไว้มาบวกเพิ่มเข้าไปอีก 6 หน่วย จากตัวอย่างข้างต้นส่วนยึดที่กำหนดไว้มีขนาด $3+6+0 = 9$ หน่วย ดังนั้นเพื่อที่จะให้ผลรวมของความกว้างของทั้งบรรทัดเพิ่มขึ้นอีก 6 หน่วยก็จะต้องทำการคูณส่วนยึดของกาวแต่ละช่องด้วย $6/9$ ความกว้างของกาวช่องแรกจะเป็น $9+(6/9) \times 3 = 11$ หน่วย ความกว้างของช่องที่สองจะเป็น $9+(6/9) \times 6 = 13$ หน่วย ส่วนในช่องสุดท้ายเนื่องจากส่วนยึดมีค่าเป็นศูนย์ ค่าของความกว้างของกาวจึงเป็น 12 หน่วยเท่าเดิม ผลของการยึดขนาดของกาวที่ได้จะแสดงได้ดังรูปที่ 2.4



รูปที่ 2.4 แสดงการขยายระยะยึดของกาว

ในขณะเดียวกัน ถ้าต้องการหาค่าความกว้างของกาวแต่ละช่องให้ความกว้างของแนวระดับทั้งบรรทัดรวมกันแล้วเท่ากับ 51 หน่วย จะต้องทำการคูณส่วนหัดของกาวด้วย $1/3$ สำหรับช่องแรกและ $2/3$ สำหรับช่องที่สอง ส่วนในช่องสุดท้ายจะไม่สามารถลดลงได้เนื่องจากส่วนหัดมีค่า เป็นศูนย์

ความกว้างของกาวในแต่ละช่องนั้นจะไม่สามารถหาค่าเฉลี่ยเกินกว่าค่าส่วนหัดในแต่ละช่องนั้น ในทางกลับกัน ความกว้างของกาวในแต่ละช่องสามารถยึดได้เกินกว่าค่าของส่วนยึดเมื่อค่า ของส่วนยึดนั้นยังคงมีค่าเป็นบวก

จากที่กล่าวมาข้างต้นแล้วจะเป็นหลักการทำงานของเท็กซ์ จะเห็นว่าตัวเท็กซ์เองจะดำเนินการประมวลผลกับสิ่งที่ยอยู่นอกกล่อง เช่นการนำกล่องมาเรียงต่อกันในแนวราบหรือแนวตั้ง การยึดหรือหัดระยะระหว่างกล่อง แต่เท็กซ์จะไม่สามารถเข้าไปจัดการกับสิ่งที่ยู่ในกล่องได้ สิ่งที่จะทำหน้าที่จัดการกับสิ่งที่ยู่ในกล่องคือเมตาฟอนต์เท่านั้น

การสร้างฟอนต์นั้นจะมีวิธีหลักอยู่ 2 วิธีคือการสร้างฟอนต์แบบบิตแมพและการสร้างฟอนต์โดยใช้เวกเตอร์ แต่ละวิธีจะมีทั้งข้อดีและข้อเสียในตัว วิธีบิตแมพจะดีตรงที่ผู้ออกแบบสามารถมองเห็นภาพของฟอนต์ที่สร้างขึ้นได้ชัดเจนและกำหนดรายละเอียดต่างๆ ได้ง่าย แต่จะมีข้อเสียคือเมื่อทำการย่อหรือขยายขนาดของฟอนต์แล้ว คุณภาพของฟอนต์ที่ได้จะไม่เหมือนกับขนาดตอนเริ่มต้นออกแบบ สำหรับฟอนต์แบบเวกเตอร์นั้นจะมีข้อดีคือเมื่อทำการย่อหรือขยายแล้วจะให้คุณภาพเหมือนหรือใกล้เคียงกับขนาดที่เริ่มต้นออกแบบมากกว่า ข้อเสียของฟอนต์แบบเวกเตอร์ก็คือการออกแบบทำได้ยาก ผู้ออกแบบต้องมีความรู้เกี่ยวกับการ

ใช้เวกเตอร์และระบบแกนอ้างอิง ตลอดจนคำสั่งที่จะใช้มากพอสมควร สำหรับวิธีการของเมตาฟอนต์จะเป็นวิธีการของเวกเตอร์ ผู้ออกแบบจะต้องเข้าใจขั้นตอนการออกแบบบ้างพอสมควรดังนี้

2.1 ระบบแกนอ้างอิงและพิกัด

ในการที่จะบอกคอมพิวเตอร์ว่าต้องการจะลากเส้นอย่างไรนั้น จำเป็นที่จะต้องมียุทธศาสตร์ที่จะอธิบายว่าจุดที่จะลากเส้นผ่านนั้นอยู่ตรงไหนบ้าง เมตาฟอนต์จะใช้วิธีการอ้างอิงด้วยพิกัดคาร์ทีเซียน (Cartesian coordinates) ตำแหน่งของจุดจะถูกกำหนดโดยแกน X ซึ่งหมายถึงจำนวนหน่วยที่วิ่งไปทางขวามือจากจุดอ้างอิง และแกน Y ซึ่งหมายถึงจำนวนหน่วยที่วิ่งขึ้นข้างบนเหนือจุดอ้างอิง ในวิธีการสร้างภาพของเมตาฟอนต์นั้นสิ่งที่จำเป็นต้องกระทำคือการวาดรูปร่างของภาพที่ต้องการลงบนกระดาษกราฟอย่างคร่าวๆ จากนั้นก็ทำการกำหนดจุดที่สำคัญลงบนภาพที่วาดนั้นเพื่อความสะดวกก็อาจกำหนดให้เป็นตัวเลขจำนวนนับ หลังจากนั้นก็จะเป็นการเขียนโปรแกรมด้วยคำสั่งของเมตาฟอนต์เพื่อที่จะบรรยายว่า จุดที่กำหนดนั้นอยู่ที่ตำแหน่งใดในระบบแกนอ้างอิง และกำหนดว่าเส้นตรง หรือเส้นโค้งอะไรบ้างที่ถูกกำหนดให้ผ่านจุดเหล่านั้น

ภายในตัวเมตาฟอนต์เองจะเสมือนมีกระดาษกราฟอยู่ภายในกรอบสี่เหลี่ยมเล็กๆของกระดาษจะเรียกว่าพิกเซล (pixel) ผลลัพธ์ที่เมตาฟอนต์ให้ออกมาจะอยู่ในรูปของสถานะที่เป็นขาวหรือดำของพิกเซลเหล่านี้ เนื่องจากการใช้ระบบแกนอ้างอิงจะต้องมีหน่วยสำหรับบอกความยาว หรือระยะทางจากจุดอ้างอิงจุดหนึ่งไปยังจุดใดๆ เพื่อความสะดวกเมตาฟอนต์จะใช้ 1 หน่วยของพิกเซลเหล่านี้แทนระยะทางหนึ่งหน่วยที่อ้างถึง

การกำหนดระยะทางบนระบบแกนอ้างอิงนี้ไม่จำเป็นที่จะต้องกำหนดระยะทางเป็นเลขจำนวนเต็มเสมอไป อาจกำหนดเป็นจุดทศนิยมก็ได้ คอมพิวเตอร์จะสามารถทำงานได้ในระดับความละเอียดของจำนวนเท่าของ $1/65536$ หรือประมาณ 0.00002 ของความกว้างของพิกเซล แต่เมตาฟอนต์จะไม่สามารถแบ่งให้ในหนึ่งพิกเซลมีได้ทั้งขาวและดำในขณะเดียวกัน สิ่งที่สามารถกำหนดได้ก็คือจะต้องเป็นดำทั้งหมดหรือขาวทั้งหมดเท่านั้น

ความละเอียดของจำนวนพิกเซลต่อหนึ่งหน่วยความยาวเรียกว่าเรส โขรุขร์ชัน (resolution) ซึ่งมักจะกำหนดเป็นจำนวนพิกเซลต่อนิ้วในระบบอเมริกัน หรือพิกเซลต่อมิลลิเมตร ในระบบอื่น ในระบบที่มีเรส โขรุขร์

ชั้นสูงๆนั้นการบิดเศษของพิกเซลจะถือได้ว่าไม่มีผลต่อรูปร่างของภาพเลย แต่อย่างไรก็ตาม ถ้าเป็นในระบบที่มีเรสโซลูชันต่ำ การบิดเศษของพิกเซลทิ้งไปหรือปัดขึ้นอาจ ทำให้รูปร่างของภาพที่ต้องการนั้นเปลี่ยนไปได้

2.2 การสร้างเส้นตรงหรือเส้นโค้ง

วิธีการสร้างเส้นตรงของเมตาฟอนต์นั้นทำได้ง่ายๆ โดยการกำหนดจุดสองจุดบนระบบอ้างอิงและใช้คำสั่ง draw ระหว่างสองจุดนั้น เมตาฟอนต์จะดำเนินการลากเส้นตรงให้ทันที รูปแบบคำสั่งจะเป็นดังนี้

```
draw (x1,y1)..(x2,y2)
```

จะเป็นการลากเส้นตรงระหว่างจุด (x1,y1) ไปยังจุด (x2,y2) เป็นต้น

วิธีการสร้างเส้นโค้งต่างของเมตาฟอนต์จะมีวิธีการที่สลับซับซ้อนอยู่บ้างซึ่งสามารถอธิบายโดยยกตัวอย่างประกอบได้ดังนี้

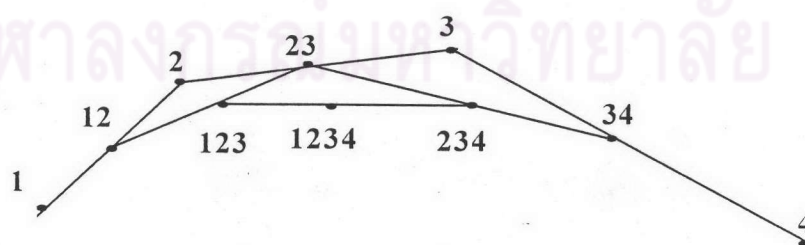
สมมุติว่าตอนเริ่มต้นมีจุดพิกัดอยู่ 4 จุดคือ z_1, z_2, z_3, z_4 จากนั้นจะมีการ กำหนดจุดที่เป็นกึ่งกลางได้ 3 จุดคือ

$$z_{12} = 1/2[z_1, z_2],$$

$$z_{23} = 1/2[z_2, z_3],$$

$$z_{34} = 1/2[z_3, z_4]$$

ดังที่แสดงในรูปที่ 2.5



รูปที่ 2.5 แสดงการจุดกึ่งกลางบนเส้นตรง

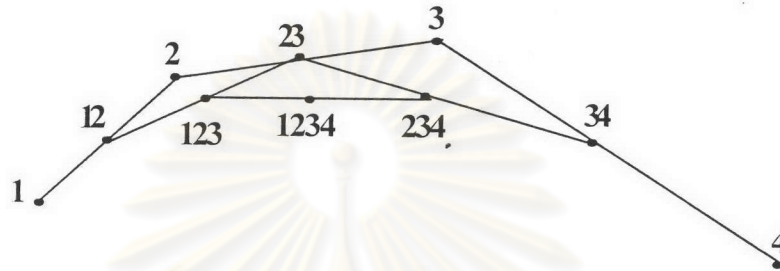
หลังจากนั้นสร้างจุดกึ่งกลางลำดับที่สองของจุดกึ่งกลางที่ได้ในขั้นแรก (z_{12}, z_{23}, z_{34}) จะได้

$$z_{123} = 1/2[z_{12}, z_{23}],$$

$$z_{234} = 1/2[z_{23}, z_{34}]$$

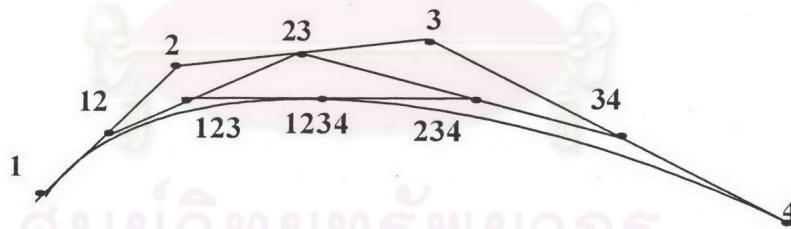
หลังจากนั้นสร้างจุดกึ่งกลางลำดับที่สาม จะได้

$$z_{1234} = 1/2[z_{123}, z_{234}] \text{ ดังรูปที่ 2.6}$$



รูปที่ 2.6 แสดงการแบ่งจุดกึ่งกลางของเส้นตรงย่อย

จุด z_{1234} ที่ได้นี้จะเป็นจุดหนึ่งที่อยู่บนเส้นโค้งที่กำหนดโดย (z_1, z_2, z_3, z_4) สำหรับการหาจุดอื่นๆที่เหลือเพื่อที่จะสามารถทำการลากเส้นโค้งนี้ได้ นั้น ก็จะใช้วิธีการเช่นเดียวกัน โดยจะเป็นการหาแนวเส้นโค้งที่กำหนดโดย $(z_1, z_{12}, z_{123}, z_{1234})$ และ $(z_{1234}, z_{234}, z_{34}, z_4)$ จะทำเช่นนี้ไปเรื่อยๆ ในที่สุดขั้นตอนการประมวลผลก็จะสิ้นสุดลงโดยรวดเร็วและก็จะได้แนวเส้นโค้งที่ต้องการดังแสดงในรูปที่ 2.7



รูปที่ 2.7 แสดงส่วนโค้งที่ได้จากการทำงานของเมตาฟอนต์

คุณสมบัติที่สำคัญของเส้นโค้งที่ได้จะมีดังนี้

- จะเริ่มต้นที่ z_1 และมีทิศทางจาก z_1 ไปยัง z_2
- จะสิ้นสุดที่ z_4 และมีทิศทางจาก z_3 ไปยัง z_4
- แนวเส้นโค้งทั้งหมดจะอยู่ภายใต้กรอบของขอบ z_1, z_2, z_3, z_4

การสร้างเส้นโค้งที่กำหนดโดยวิธีการเรียกตัวเองดังที่กล่าวข้างต้นนี้สามารถสรุปเป็นสมการได้ดังนี้

$$z(t) = (1-t)^3 z_1 + 3(1-t)^2 t z_2 + 3(1-t)t^2 z_3 + t^3 z_4$$

โดย t คือพารามิเตอร์ที่มีค่าตั้งแต่ 0 ถึง 1 สมการโพลีโนเมียลที่มีดีกรีของ t เป็น 3 นี้เรียกว่า Bernshtein polynomial

2.3 การขยายขนาดและความละเอียด (Magnification and Resolution)

จุดเด่นอย่างหนึ่งของเมตาฟอนต์ก็คือสามารถที่จะผลิตฟอนต์สำหรับอุปกรณ์ การพิมพ์หลายๆ ชนิด ถ้าผู้ออกแบบโปรแกรมได้ออกแบบให้เรสโซลูชันเปลี่ยนแปลงค่าได้

สมมติว่ามีอุปกรณ์เอาท์พุทอยู่ 2 ชนิด ชนิดแรกให้ชื่อว่าชนิด dot มีความละเอียด 200 พิกเซลต่อนิ้ว (ประมาณ 8 พิกเซลต่อมิลลิเมตร) อีกชนิดหนึ่งมีชื่อว่า laser มีความละเอียด 2000 พิกเซลต่อนิ้ว ต้องการสร้างโปรแกรมเมตาฟอนต์ที่สามารถแสดงผลได้กับอุปกรณ์ทั้งสองประเภท ดังนั้นถ้าจำเป็นต้องพิมพ์ผลลัพธ์ของแฟ้มที่มีชื่อว่า test.mf ออกทางเครื่อง dot ก็เพียงแต่ใช้คำสั่ง

```
\mode=dot; input test
```

หรือถ้าต้องการพิมพ์ออกทางเครื่องพิมพ์ที่มีคุณภาพสูงขึ้นเพียงแต่เปลี่ยนจาก dot เป็น laser เท่านั้น

```
\mode=laser; input test
```

ด้วยวิธีการดังที่กล่าวมานี้ ถ้าต้องการพิมพ์ผลลัพธ์ออกทางอุปกรณ์เอาต์พุตอื่นๆ จะสามารถทำได้ โดยเปลี่ยนค่าตัวแปรของโหมด (mode variable) เท่านั้น

วิธีการทำงานข้างต้นนี้สามารถทำได้โดยใส่คำสั่ง `mode_setup` ไว้ที่ตอนเริ่มต้นของโปรแกรม `test` และคำสั่งที่ใส่ไว้จะไปทำการกำหนดค่ากลุ่มของตัวแปรกลุ่มหนึ่งที่ทำหน้าที่ควบคุมการแสดงผลลัพธ์ซึ่งจะมีส่วนเกี่ยวข้องกับจำนวนของพิกเซลต่อหน่วยความยาว จากตัวอย่างข้างต้น เมื่อโหมดถูกกำหนดให้เป็น `dot` ค่าของ `pt` จะเท่ากับ 2.7674 พิกเซล (`pt` คือจำนวนจุดของเครื่องพิมพ์ (printer's point) ซึ่งเป็นหน่วยชนิดหนึ่งเช่นเดียวกับมิลลิเมตร $72.27 \text{ pt} = 1 \text{ นิ้ว}$) และ $\text{mm} = 7.87402$ พิกเซล แต่เมื่อโหมดถูกกำหนดให้เป็น `laser` ค่าของ `pt` = 27.674 และ $\text{mm} = 78.4017$ โปรแกรม `test.mf` ที่ทำงานภายใต้โหมด `dot` จะให้ผลลัพธ์ที่มีขนาดเล็กและสั้นกว่าที่ทำงานภายใต้โหมด `laser` ประมาณ 10 เท่า หรือในกรณีที่ต้องการเขียนเส้นตรงให้มีขนาดความยาว 3 มิลลิเมตร จากพิกัด (0,0) ไปยัง (3mm,0) ในโหมด `dot` จะใช้จำนวนพิกเซลเท่ากับ 23.6 พิกเซล และโหมด `laser` จะใช้จำนวนพิกเซลเท่ากับ 236.2 พิกเซล

ความยุ่งยากซับซ้อนที่เกิดขึ้นของการพิมพ์อีกอย่างหนึ่งก็คือการขยายขนาด หรือย่อส่วนของผลลัพธ์ซึ่งไม่ใช่ขนาดเดียวกับที่ทำการออกแบบไว้ตอนเริ่มต้นเช่น ถ้าต้องการให้ผลลัพธ์ที่ได้ภายใต้โหมด `dot` มีขนาดใหญ่ขึ้นเป็นสองเท่า เท็กซ์จะใช้คำสั่ง

```
\magnification=2000
```

ในการทำให้ขนาดของเอาต์พุตทั้งหมดใหญ่ขึ้นเป็นสองเท่า หรืออาจจะใช้คำสั่ง

```
\font\font=test scaled 2000
```

สำหรับทำให้ฟอนต์ภายใต้โปรแกรม `test.mf` เท่านั้นมีขนาดใหญ่ขึ้นเป็นสองเท่า วิธีการจัดเตรียมฟอนต์ของโปรแกรมเมตาฟอนต์เพื่อให้ใช้งานได้กับเท็กซ์สามารถทำได้โดยใช้คำสั่ง

```
\mode=dot; mag=2; input test;
```

ซึ่งจะทำให้ขนาดของ $pt = 5.5348$ พิกเซล และ $mm = 15.74803$ พิกเซล

การทำงานของรูทีน `mode_setup` จะทำการตรวจสอบว่าค่าของ `mag` เป็นเท่าใด ถ้าไม่มีการกำหนดค่าของ `mag` จะถือว่ามีความเป็น 1 ในทำนองเดียวกัน ถ้าไม่มีการกำหนดชนิดของโหมดที่ใช้ เมตาฟอนต์จะถือว่า `mode = proof`

หน่วยที่เมตาฟอนต์ใช้ในการอ้างอิงต่าง ๆ มีดังนี้

pt	printer's point	(72.27 pt = 1 in)
pc	pica	(1 pc = 12 pt)
in	inch	(1 in = 2.54 cm)
bp	big point	(72 bp = 1 in)
cm	centimeter	(100 cm = 1 meter)
mm	millimeter	(10 mm = 1 cm)
dd	didot point	(1157 dd = 1238 pt)
cc	cicero	(1 cc = 12 dd)

ในแต่ละหน่วยนั้นถ้ามีการปิดเศษของหน่วยจากการคำนวณ เมตาฟอนต์จะทำการปิดให้ใกล้เคียงกับ $1/65536$ ส่วนของพิกเซล

หน่วยของพิกเซลที่กล่าวถึงข้างต้นนั้นเรียกได้ว่าเป็นมิติทางกายภาพ (physical dimension) ซึ่งเมตาฟอนต์จะไม่ค่อยนิยมใช้เท่าใดนักเนื่องจากจะเป็นการทำให้การเขียนโปรแกรมไม่มีความยืดหยุ่น เมตาฟอนต์จะมีหน่วยสำหรับบอกมิติอีกประเภทหนึ่งซึ่งจะให้ความสะดวกในการออกแบบฟอนต์มากกว่าเช่น จะนิยมใช้ 'em' หรือ 'x_height' ในการกำหนดขนาดของตัวอักษร ซึ่งจะแปรเปลี่ยนไปตามชนิดของฟอนต์ เป็นต้น

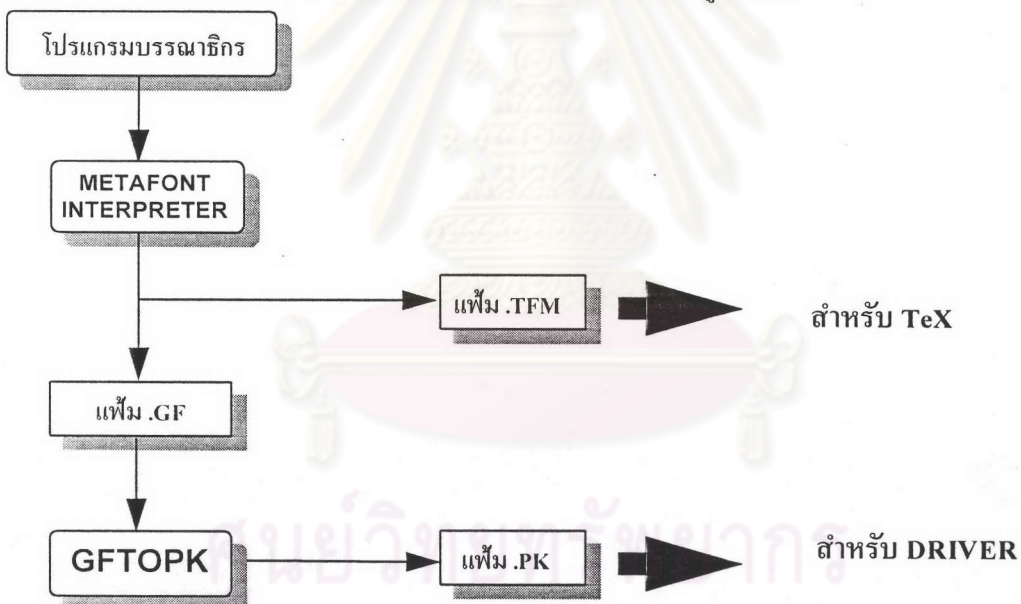
เนื่องจากการขยายขนาดหรือการกำหนดความละเอียดในโหมดต่างๆจะมีผลทำให้จำนวนของพิกเซลต่อหน่วยมิติทางกายภาพนี้เปลี่ยนไปด้วย วิธีการที่จะให้จำนวนหน่วยของพิกเซลมีค่าคงที่ไม่เปลี่ยนแปลงไปตามการกำหนดความละเอียดในโหมดต่าง ๆ นั้น เมตาฟอนต์จะกำหนดปริมาณชาร์ป (sharp quantity) ขึ้นเพื่อ

ใช้ในการนี้ โดยเติมเครื่องหมาย '#' ต่อท้ายหน่วยที่ต้องการก็จะทำให้ได้หน่วยของมิติที่ไม่เปลี่ยนแปลงค่า เมื่อมีการเปลี่ยนแปลงความละเอียด

หลักการที่กล่าวมาข้างต้นนี้เป็นหลักการที่จำเป็นในการที่จะพัฒนาโปรแกรมสร้างฟอนต์ด้วยเมตาฟอนต์ อย่างไรก็ตามขั้นตอนในการพัฒนาโปรแกรมด้วยเมตาฟอนต์สามารถสรุปได้ดังนี้

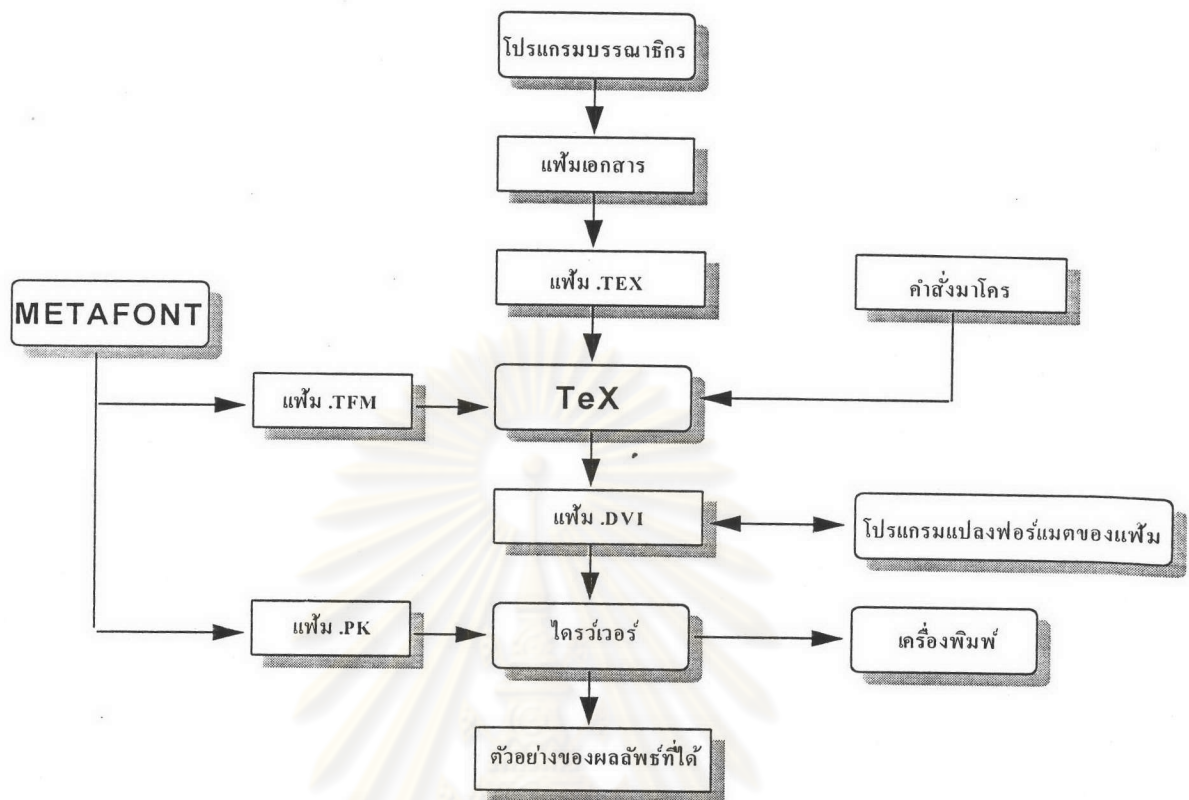
1. เขียนโปรแกรมด้วยชุดคำสั่งของเมตาฟอนต์จะได้แฟ้มที่มีนามสกุล .mf
2. ใช้เมตาฟอนต์ประมวลผลแฟ้ม .mf จะได้แฟ้มที่มีนามสกุล .tfm และ .gf (หรือ .xxxgf ขึ้นอยู่กับระบบปฏิบัติการที่ใช้)
3. เปลี่ยนแฟ้มที่มีนามสกุล .gf ให้เป็นแฟ้มที่มีนามสกุล .pk หรือ .pxl ทั้งนี้ขึ้นอยู่กับซอฟต์แวร์ไดรเวอร์ที่เท็กซ์ต้องการ

ขั้นตอนการทำงานของโปรแกรมเมตาฟอนต์สามารถแสดงได้ดังรูปที่ 2.8



รูปที่ 2.8 แสดงการทำงานโดยทั่วไปของเมตาฟอนต์

เมื่อนำมาประกอบกับขั้นตอนการทำงานของเท็กซ์แล้วจะเป็นดังรูปที่ 2.9



รูปที่ 2.9 แสดงการทำงานของเท็กซ์และเมตาฟอนต์

3. การตัดคำ

สิ่งที่มักจะพบกันอยู่เสมอในการประมวลผลคำด้วย โปรแกรมประมวลผลคำหรือเวิร์ด โพรเซสเซอร์ที่นิยมใช้กันโดยทั่วไปก็คือการตัดหรือแบ่งคำในการขึ้นบรรทัดใหม่ปัญหาที่กล่าวมานี้ เป็นตัวแปรสำคัญอย่างยิ่งที่ทำให้ความน่าใช้ของโปรแกรมประมวลผลคำลดลงอย่างมากทีเดียว ในบางครั้งจะพบว่าโปรแกรมประมวลผลคำนั้นมีฟอนต์ที่สวยงามให้เลือกใช้อย่างหลากหลายหรือมีรูทีนในการที่จะแสดงผลให้ผู้ใช้งานสามารถมองเห็นได้ก่อนที่จะทำการพิมพ์งานลงบนกระดาษก็ตาม ถ้าอัลกอริทึมในการตัดคำของโปรแกรมนั้นไม่ดีเพียงพอ หรือเมื่อทำการเว้นวรรคแล้วทำให้ช่องว่างระหว่างตัวอักษรหรือช่องว่างระหว่างคำมากเกินไป ก็จะทำให้ความน่าสนใจของโปรแกรมนั้นลดลงอย่างมากเลยทีเดียว การออกแบบอัลกอริทึมสำหรับการตัดคำนี้เป็นปัญหาที่ยากยิ่งในหลายๆภาษา รวมทั้งภาษาไทยด้วยเช่นกันในตัวโปรแกรมของเท็กซ์เองก็มีอัลกอริทึมสำหรับจัดการในส่วนนี้ด้วย แต่เนื่องจากการตัดคำของเท็กซ์นั้นทำไว้สำหรับภาษาอังกฤษ โดยโครงสร้างทางภาษาแล้วภาษาอังกฤษจะมีลักษณะเป็นคำๆ ในหนึ่งประโยคจะประกอบด้วยคำต่างๆมา

ประกอบกันและมีการแบ่งแยกระหว่างคำด้วยวรรคหรือช่องว่าง ซึ่งแตกต่างจากภาษาไทยซึ่งในหนึ่งประโยคจะนำคำมาเรียงติดกัน ไปหมดไม่มีการใช้วรรคหรือช่องว่างในการแบ่งจนกว่าจะจบประโยคหรือเริ่มต้นเรื่องใหม่ ดังนั้น การพัฒนาอัลกอริทึมที่ใช้ในโปรแกรมเท็กซ์ให้ใช้ได้กับการตัดคำภาษาไทยนั้นค่อนข้างที่จะยุ่งยากคง ต้องใช้ระยะเวลาในการศึกษากันต่อไป อย่างไรก็ตามสำหรับการตัดคำภาษาไทยเองนั้นมี อัลกอริทึมที่ใช้กันอยู่มากมาย และที่จะนำมาใช้ในการพัฒนาโปรแกรมไทยเท็กซ์นี้จะเลือกใช้อัลกอริทึมตัดคำที่ใช้อยู่ในโปรแกรม CU Writer

ถึงแม้ว่าอัลกอริทึมตัดคำที่มีอยู่ใน โปรแกรมเท็กซ์นั้นจะยังไม่สามารถนำมาใช้ได้กับโปรแกรมไทยเท็กซ์ก็ตาม แต่วิธีการที่อัลกอริทึมนี้ใช้ก็เป็นวิธีที่น่าสนใจมากซึ่งสามารถสรุปได้ดังนี้

3.1 การตัดคำด้วยวิธี *Hyphenation* ที่ใช้ใน โปรแกรมเท็กซ์

เท็กซ์จะทำการตัดคำที่นำมาโดยเริ่มต้นเปรียบเทียบกับคำที่มีใน "exception dictionary" หรือพจนานุกรมเก็บคำเฉพาะที่ต้องกำหนดวิธีแบ่งเอาไว้เป็นพิเศษ ถ้าคำที่ต้องการจะตัดตรงกับคำเฉพาะนั้นก็จะใช้วิธีการตัดคำแบบพิเศษนั้นเลย ถ้าไม่พบก็จะทำการค้นหารูปแบบ (pattern) ของคำนั้นและนำมาเปรียบเทียบกับรูปแบบ (pattern) ที่เท็กซ์มีอยู่ เพื่อที่จะให้เข้าใจง่ายขึ้นจำเป็นต้องยกตัวอย่างประกอบ โดยใช้คำว่า "hyphenation"

เมื่อเริ่มต้นเท็กซ์จะทำการใส่จุด "." เข้าไปวางไว้หน้าคำและหลังคำแห่ง ละหนึ่งจุด

.hyphenation.

จะถือว่าคำ ".hyphenation." นี้เป็นคำใหม่ที่จะดำเนินการค้นหาต่อไป ถ้านำเอาตัวอักษรที่ประกอบเป็นคำใหม่นี้มาแยกเป็นคำย่อยที่มีความยาวของตัวอักษรในแต่ละคำเท่ากับ 1 จะได้

. h y p h e n a t i o n .

ถ้าให้ความยาวของตัวอักษรในแต่ละคำเท่ากับ 2 จะได้

.h hy yp ph he en na at ti io on n.

ถ้าให้ความยาวของตัวอักษรในแต่ละคำเท่ากับ 3 จะได้

.hy hyp yph phe hen ena nat ati tio ion on.

กระทำเช่นนี้ไปเรื่อยๆจนความยาวของคำย่อยที่ได้เท่ากับคำเริ่มต้น

จากวิธีการข้างต้นนี้จะพบว่าแต่ละคำย่อยที่มีความยาวของตัวอักษรเท่ากับ k ตัว จะสามารถใส่ตัวเลขห้อยข้างหน้าตัวอักษรและข้างท้ายตัวอักษรได้ $k+1$ ตัวสำหรับตำแหน่ง ของตัวอักษรที่ติดกันเช่นคำย่อย "hen" สามารถใส่ตัวเลขได้เป็น ${}_0h{}_0e{}_0n{}_0$ ตัวเลข 2 ที่เห็นนี้จะหมายถึงตำแหน่งที่อาจจะแบ่งคำโดยใส่เครื่องหมายไฮเฟนได้ วิธีที่จะรู้ว่าคำย่อยแต่ละคำมีตัวเลขอะไรบ้างนั้นจะทำได้โดยการค้นหาใน "pattern dictionary" หรือพจนานุกรมที่เก็บ รูปแบบของคำ (ในระบบเท็กซ์พจนานุกรมดังกล่าวจะถูกเก็บไว้ในแฟ้มที่ชื่อว่า "hyphen.tex") ถ้าหาไม่พบจะถือว่าตัวเลขที่ห้อยอยู่ระหว่างตัวอักษรนั้นมีค่าเป็นศูนย์ทั้งหมด ในที่นี้เราจะนำเฉพาะคำที่หาพบออกมาเพื่อเปรียบเทียบกันต่อไป ซึ่งรูปแบบที่สามารถหาพบคือ

${}_0h{}_3p{}_0{}_0h{}_2e{}_0{}_0h{}_0e{}_0{}_4h{}_0e{}_5a{}_0{}_1n{}_0{}_0n{}_2a{}_0{}_1t{}_0{}_0{}_2i{}_0{}_0{}_2i{}_0{}_0$

หลังจากนั้นก็ทำการคำนวณหาค่าตัวเลขที่ห้อยอยู่ระหว่างตัวอักษรที่มีค่ามากที่สุดสำหรับตัวอักษรคู่เดียวกันที่อยู่ในรูปแบบที่ต่างกัน เช่นในระหว่างตัวอักษร "e" และ "n" จะมีรูปแบบที่เกี่ยวข้องอยู่ทั้งหมด 4 รูปแบบ ดังนั้นจะมีตัวเลขที่จะนำมาเปรียบเทียบอยู่ 4 ตัว (2 จาก ${}_0h{}_0e{}_0n{}_0$, 0 จาก ${}_0h{}_0e{}_0{}_4$, 0 จาก ${}_0h{}_0e{}_0{}_5a{}_0{}_1$ และ 1 จาก ${}_1n{}_0{}_0$) ตัวเลขที่มีค่ามากที่สุดสำหรับกรณีนี้คือ 2 เมื่อใช้วิธีการนี้ทำจนครบทุกรูปแบบแล้วจะได้ผลลัพธ์คือ

.h₀yp₃h₀en₂at₅ion₄o₂n₀.

ขั้นตอนสุดท้ายคือการแบ่งคำด้วยไฮเฟน วิธีการพิจารณาการแบ่งคำนี้จะเลือกเฉพาะตำแหน่งที่เป็นตัวเลขที่เท่านั้น ในตัวอย่างนี้จะสามารถแบ่งได้สองตำแหน่งคือ

hy-phen-ation

ในการทำงานของเท็กซ์ รูปแบบของคำเหล่านี้ซึ่งมีประมาณ 4447 รูปแบบจะถูกโหลดลงไว้ในหน่วยความจำของเท็กซ์ คำตัวเลขที่ห้อยอยู่จะมีค่าสูงสุดเป็น 5 ซึ่งหมายความว่าในตำแหน่งนั้นจะสามารถใช้ เป็นจุดแบ่งคำได้อย่างแน่นอน อัลกอริทึมตัดคำที่เท็กซ์ใช้นี้เรียกว่า Liang's algorithm ข้อดีของอัลกอริทึมนี้คือไม่จำเป็นต้องใช้คำทั้งหมดที่มีอยู่ในพจนานุกรมมาเป็นตัวเปรียบเทียบ จึงทำให้สามารถประหยัดเนื้อที่ของหน่วยความจำที่ใช้เป็นอันมาก และไม่จำเป็นต้องเปลี่ยนพจนานุกรมใหม่ทุกครั้งที่มีการเพิ่มคำใหม่เพิ่มขึ้น เนื่องจากพจนานุกรมรูปแบบที่มีอยู่จะบ่งบอกถึงไวยากรณ์และโครงสร้างของการผสมตัวอักษรให้เป็นคำอยู่ในตัวเองอยู่แล้ว ถ้ามีคำใหม่ซึ่งต้องการแบ่งด้วยรูปแบบพิเศษไปกว่านี้ก็จะสามารถนำไปบรรจุไว้ใน "exception dictionary" แทน

3.2 การตัดคำด้วยอัลกอริทึมของ โปรแกรม CU Writer

การทำงานของอัลกอริทึมตัดคำหรือแบ่งคำของโปรแกรม CU Writer นี้ จะใช้พจนานุกรมมาเป็นเกณฑ์ในการแบ่ง แต่วิธีการเก็บพจนานุกรมนั้นจะไม่ใช้การเก็บในลักษณะของคำต่อคำเนื่องจากจะทำให้เสียเนื้อที่หน่วยความจำเป็นจำนวนมาก และจะใช้เวลาในการค้นหา ดังนั้นจึงจำเป็นต้องมีการสร้างรูปแบบขึ้นมาเพื่อจัดการเปลี่ยนลักษณะพจนานุกรมที่เก็บคำศัพท์ชนิดคำต่อคำให้เป็นรูปแบบที่ง่ายในการค้นหาและสืบเปลี่ยนหน่วยความจำน้อยลง

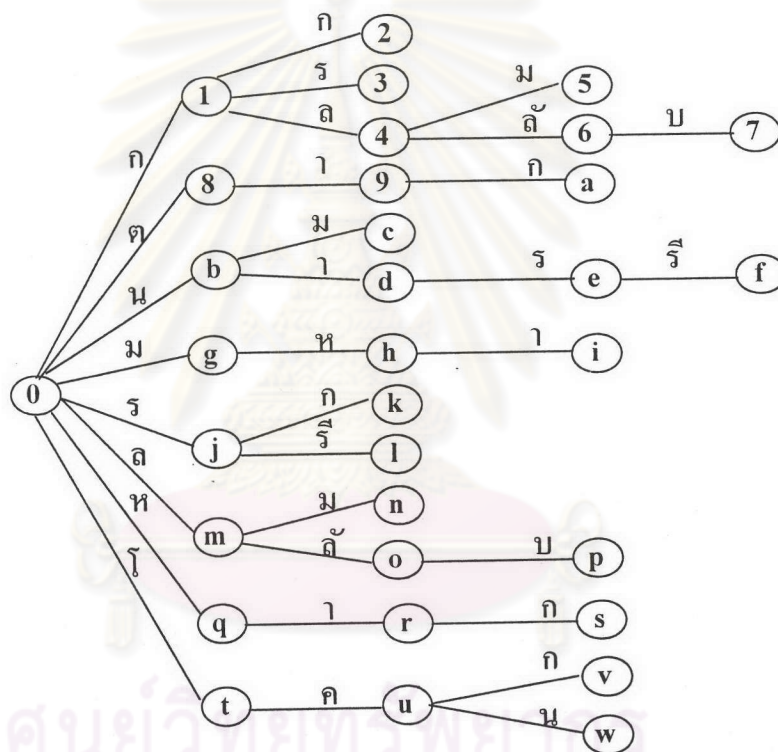
3.2.1 การแทนพจนานุกรมในหน่วยความจำด้วย โครงสร้างข้อมูลแบบทรี (Trie)

คำว่าพจนานุกรมในความหมายของอัลกอริทึมตัดคำนั้นจะหมายถึงโครงสร้างของข้อมูลแบบต้นไม้ (tree structure) โครงสร้างของข้อมูลแบบต้นไม้หรือทรี มีองค์ประกอบหลักคือ สถานะ (state) และทางเลือกที่จะไปยังสถานะใหม่เรียกว่าทรานสิชัน (transition) วิธีการที่จะทำความเข้าใจเกี่ยวกับการแทนที่พจนานุกรมด้วยทรีให้ดียิ่งขึ้นนั้นจำเป็นต้องอาศัยตัวอย่างประกอบดังนี้

สมมุติว่าในพจนานุกรมของภาษาไทยประกอบด้วยคำจำนวน 20 คำคือ

กก กร กล กลม กลับ ตา ตาก นม นา นารี มหา รก รัม ลับ หา หาก โค โลก โคน

ดังนั้นประโยคหรือวรรคในภาษาไทยก็คือสตริงของอักขระที่ประกอบด้วยคำใดๆก็ตามที่อยู่ในพจนานุกรมนี้มาเรียงต่อกันนั่นเอง สำหรับการนำคำดังกล่าวมาสร้างเป็นทรีนั้น จะสามารถทำได้ดังรูปที่ 2.10

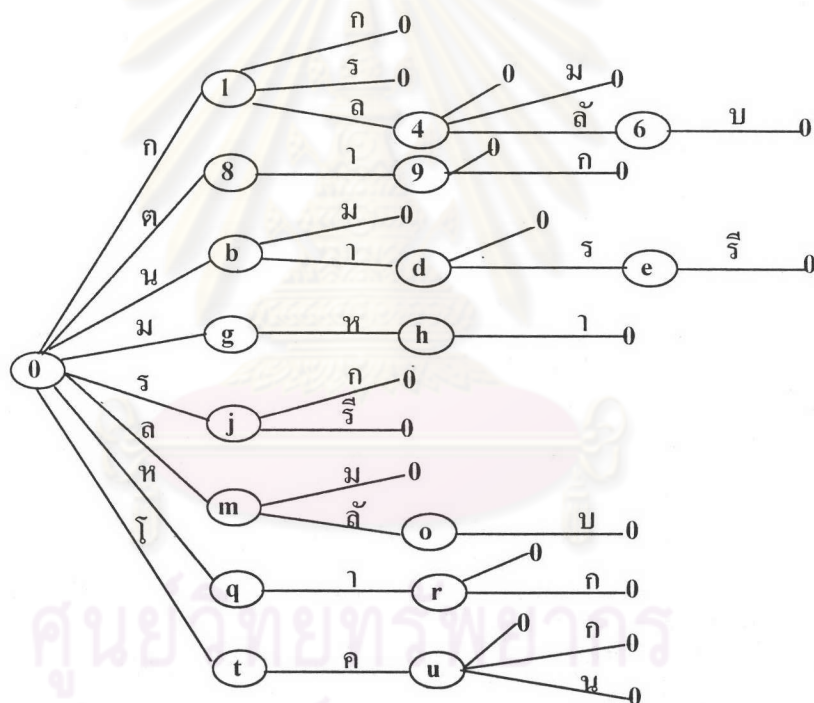


รูปที่ 2.10 แสดงแผนภาพการใช้ DFA เก็บพจนานุกรม

สิ่งที่ได้นี้เรียกว่าทรี (trie) หรือ deterministic finite state automaton (DFA) ของพจนานุกรมวิธีการทำงานของ DFA นี้จะเริ่มต้นที่สถานะ 0 จากนั้นจะรับตัวอักษรที่ป้อนเข้ามาทีละตัวและเดินไปตามเส้นทางต่างๆตามตัวอักษรที่ได้รับ ถ้าตัวอักษรที่ได้รับในแต่ละสถานะนั้นแตกต่างไปจากทางเลือกที่มีอยู่ก็แสดงว่าสตริงของอักขระนั้นไม่มีอยู่ในพจนานุกรมนี้ในทางกลับกันถ้าตัวอักษรที่ได้รับในทุกสถานะสามารถทำให้เดินไป

ถึงสถานะสุดท้ายที่เป็นสถานะสิ้นสุด (final state) ได้โดย ไม่เกิดข้อผิดพลาดขึ้นเสียก่อน แสดงว่าสตริงของอักขระนั้นเป็นคำที่มีอยู่ในพจนานุกรม (สถานะสุดท้ายหมายถึงสถานะที่ไม่มีทางเดินต่อไปอีกแล้ว) สิ่งที่กำลังกล่าวมาข้างต้นนี้ แท้จริงแล้วคือการตรวจสอบความถูกต้องของตัวสะกดของคำที่ใช้กันอยู่ในโปรแกรมประมวลผลคำต่างๆไป

ถ้าต้องการจะค้นหาจุดสำหรับแบ่งคำในประโยคหรือวรรค จะมีความแตกต่างออกไปบ้างเล็กน้อย เนื่องจากอัลกอริทึมในการตัดคำนั้นไม่ได้ต้องการแค่เพียงว่าคำนั้นอยู่ในพจนานุกรมหรือไม่ แต่ต้องการหาจุดที่เหมาะสมในการแบ่งด้วย ดังนั้นจึงต้องมีการสร้าง nondeterministic finite state automaton (NFA) ของประโยคหรือวรรคแทน ซึ่งจาก DFA ของพจนานุกรมสามารถเปลี่ยนเป็น NFA ของประโยคหรือวรรคได้ดังรูปที่ 2.11



รูปที่ 2.11 แสดงแผนภาพการสร้าง NFA

ขั้นตอนในการทำงานของการแบ่งคำนี้ก็คล้ายคลึงกับการตรวจสอบความถูกต้องของตัวสะกดของคำดังที่กำลังกล่าวมาแล้วข้างต้น จะมีข้อแตกต่างกันอยู่ตรงที่เมื่อเดินมาตามทางจนถึงสถานะสุดท้ายที่ไม่มีทางไปต่ออีกแล้วนั้น ในการแบ่งคำนี้จะไม่มีการสิ้นสุดท้ายแต่จะเป็นสถานะ 0 แทน สังเกตพบว่าเป็นสถานะเดียวกันกับสถานะเริ่มต้น ซึ่งแปลความหมายได้ว่าความยาวของคำนี้สิ้นสุดลงแล้วมีทางเลือกกว่าจะสิ้นสุดประโยคหรือ

วรรคตรงนี้ หรืออาจจะเริ่มต้นค่าใหม่ต่อไปก็ได้ แต่บางครั้งก็อาจจะเดินมาถึงสถานะสิ้นสุด (final state) ที่ไม่ใช่สถานะสุดท้าย (วงกลมล้อมรอบสถานะสิ้นสุดจะมีความเข้มมากกว่าปกติ) ถ้าเป็นเช่นนี้จะหมายความว่า อาจจะหยุดค่าลงที่ตรงตำแหน่งนี้ หรือจะสิ้นสุดประโยคเพียงเท่านี้ หรืออาจจะรอการป้อนตัวอักษรต่อไป เพื่อนำมาประกอบเป็นคำที่ยาวขึ้นก็เป็นได้

ความยุ่งยากที่ตามมาจากการมีทางเลือกมากกว่าหนึ่งทางก็คือทางไหนคือทางที่ถูกต้อง คำตอบสำหรับคำถามนี้ก็คือทางที่เดินไปแล้วสำเร็จ (successful path) เส้นทางที่เรียกได้ว่าเป็นเส้นทางที่สำเร็จก็คือเส้นทางที่เดินไปแล้วพบสถานะสิ้นสุดมีค่าเป็น 0 ตามหลักการแล้วการที่จะยอมรับว่าประโยคหรือวรรคนั้นสิ้นสุดหรือไม่จะต้องพิจารณาว่าต้องมีเส้นทางที่เดินไปแล้วสำเร็จอย่างน้อย 1 เส้นทาง นับจากสถานะเริ่มต้น (start state) ไปยังสถานะสิ้นสุด (final state)

ตัวอย่างที่สามารถยกมาอธิบายได้เช่น จากประโยค "หากโคกล" เมื่อนำมาป้อนเข้ากับ NFA จะมีเส้นทางที่เป็นไปได้อยู่ 3 เส้นทางคือ

0--h--q--a--0--k--l--โ
 0--h--q--a--r--k--0--โ--t--ค--0--ก--l--ล--0
 0--h--q--a--r--k--0--โ--t--ค--u--ก--0--ล--i

เส้นทางที่ถือว่าสำเร็จคือเส้นทางที่ 2 เพราะจบลงด้วยสถานะสิ้นสุดเป็น 0 เส้นทางที่ 1 นั้นเรียกว่าเป็นการล้มเหลวเนื่องจากไม่สามารถรับตัวอักษรต่อไปได้ สำหรับเส้นทางที่ 3 ถึงแม้ว่าจะรับตัวอักษรได้จนจบแต่ก็ไม่ถือว่าสำเร็จเพราะในสถานะนั้นยังต้องรอตัวอักษรที่จะต้องนำมาประกอบเพื่อให้สถานะนั้นสมบูรณ์ต่อไป

เมื่อได้เส้นทางที่ถือว่าสำเร็จแล้ว ก็จะนำมาพิจารณาหาจุดที่แบ่งคำ จากตัวอย่างข้างต้นเราได้เส้นทางที่สำเร็จคือ

0--h--q--a--r--k--0--โ--t--ค--0--ก--l--ล--0

ซึ่งเมื่อพิจารณาแล้วจะพบว่าค่าที่อยู่ระหว่างสถานะ 0 จะเป็นค่าที่มีอยู่ในพจนานุกรม 20 ค่าที่กำหนดไว้ตั้งแต่ต้น ดังนั้นจึงสรุปได้ว่าสถานะสิ้นสุดที่มีค่าเป็นศูนย์ก็คือจุดที่ใช้แบ่งค่านั่นเอง การแบ่งค่าที่ได้จึงเป็น (หาก)(โค)(กล)

3.2.2 การแบ่งค่าด้วยทรี

จากที่อธิบายไว้ข้างต้นจะพบว่า การแบ่งค่าจะใช้หลักการของ Finite state Automaton ซึ่งก็คือโครงสร้างของทรี ดังนั้นอัลกอริทึมแบ่งค่าที่ใช้วิธีการทำงานของทรีจะสามารถเข้ากันได้ดีกับงานประเภทนี้ สิ่งที่น่าสนใจต่อไปก็คือสถานะต่างๆที่ได้จากทรีจะมีโอกาสเป็นไปได้ 4 แบบคือ

1. สถานะที่ยังไม่สามารถแบ่งค่าได้ คือ สถานะที่ไม่ใช่สถานะสิ้นสุด
2. สถานะที่สามารถแบ่งค่าได้ คือสถานะนั้นเป็นสถานะสิ้นสุดแล้ว แต่ยังมีเส้นทางต่อไปอีก ในสถานะนี้จึงมีทางเลือกว่าจะแบ่งหรือไม่แบ่งก็ได้
3. สถานะที่บังคับให้ต้องแบ่งค่า คือสถานะสิ้นสุดที่ไม่มีเส้นทางต่อไปอีกแล้ว
4. สถานะที่ไม่สำเร็จ คือสถานะที่มีความผิดพลาดเนื่องจากตัวอักษรที่ได้รับ ไม่อยู่ในเส้นทางที่ทรีมีอยู่

ปัญหาต่อไปก็คืออะไรจะเป็นตัวที่บ่งชี้ได้ว่าสตริงของอักขระที่ป้อน เข้าไปนั้นอยู่ในสถานะใด วิธีการตรวจสอบสถานะนี้จะใช้หลักการคือ จะให้ s แทนสตริงของอักขระ current แทนค่าที่กำลังค้นหา และให้ next แทนค่าที่อยู่ในลำดับต่อจาก current (ในที่นี้มีสมมุติฐานว่าพจนานุกรมถูกเรียงตามลำดับแล้ว โดยการเรียงลำดับนี้ไม่จำเป็นต้องเป็นการเรียงตามอักขระวิธีที่ใช้ในพจนานุกรมฉบับราชบัณฑิตยสถาน) จากนั้นจะพิจารณาได้ดังนี้

1. สถานะที่ไม่สำเร็จ
ถ้าเปรียบเทียบตัวอักษรที่อยู่ใน s ด้วยจำนวนตัวอักษรที่เท่ากับที่มีอยู่ใน current แล้วพบว่าไม่เท่ากัน แสดงว่าการค้นหานี้ล้มเหลว จะเรียกว่าสถานะที่ไม่สำเร็จ
2. สถานะที่ไม่สามารถแบ่งค่าได้

ถ้าเปรียบเทียบตัวอักษรที่อยู่ใน s ด้วยจำนวนตัวอักษรเท่ากับที่มีอยู่ใน $current$ แล้วพบว่าเท่ากัน แต่ตัวอักษรที่มีอยู่ใน s มีมากกว่าตัวอักษรที่มีอยู่ใน $current$ แสดงว่า ยังไม่เป็นคำ

3. สถานะที่สามารถแบ่งคำได้

ถ้าเปรียบเทียบตัวอักษรที่อยู่ใน s ด้วยจำนวนตัวอักษรเท่ากับที่มีอยู่ใน $current$ แล้วพบว่าเท่ากัน และเมื่อเปรียบระหว่าง s กับ $next$ ด้วยจำนวนตัวอักษรเท่ากับที่มีอยู่ใน $next$ แล้วยังคงพบว่าเท่ากันอีก แสดงว่ายังมีคำที่มีความยาวมากกว่านี้ ดังนั้นจะสามารถแบ่งคำที่ตำแหน่งนี้ได้ หรืออาจจะรอสะสมตัวอักษรให้มากกว่านี้แล้วไปแบ่งในตำแหน่งต่อไปก็ได้เช่นกัน

4. สถานะที่บังคับให้ต้องแบ่ง

ถ้าเปรียบเทียบตัวอักษรที่อยู่ใน s ด้วยจำนวนตัวอักษรเท่ากับที่มีอยู่ใน $current$ แล้วพบว่าเท่ากัน และเมื่อเปรียบระหว่าง s กับ $next$ ด้วยจำนวนตัวอักษรเท่ากับที่มีอยู่ใน $next$ แล้วพบว่าไม่เท่ากัน แสดงว่าคำที่ยาวที่สุดที่มีคือคำใน $current$ ดังนั้นจะต้องแบ่งคำตรงตำแหน่งนี้ จากนั้นจึงค่อยรับตัวอักษรเข้าไปเริ่มต้นคำใหม่

3.2.3 อัลกอริทึม Backtracking Word Separation (BWS Algorithm)

จากการได้สถานะทั้งสี่ข้างต้นนี้ สามารถนำมาพัฒนาให้เป็นอัลกอริทึมสำหรับตัดคำได้อย่างสมบูรณ์ที่เรียกว่า Backtracking Word Separation หรือ BWS

เนื่องจากสถานะทั้งสี่นั้นสามารถเขียนแทนได้ด้วยทรีหรือ NFA ซึ่งก็คือรูปแบบของทรีชนิดหนึ่ง (separation tree) และจุดแบ่งคำที่ได้จากแต่ละสถานะก็ได้ มาจากตำแหน่งของโหนดที่อยู่ในทรีนั้น การค้นหาจุดแบ่งคำหรือสถานะใดๆที่ต้องการก็คือการสร้าง อัลกอริทึมที่เดินเชื่อมโหนดต่างๆนั่นเอง

อัลกอริทึมนี้จะทำการอ่านตัวอักษรจากสตริงของอักขระที่ป้อนเข้ามา ทีละตัว จากนั้นก็จะเดินไปตามเส้นทางที่อยู่ในทรีเพื่อค้นหาสถานะต่างๆ ซึ่งถ้าใช้วิธีการเขียนโปรแกรมในลักษณะของการเรียกตัวเอง (recursive) จะเขียนได้สะดวกและง่ายที่สุด แต่วิธีนี้จะทำงานได้ค่อนข้างช้าเพราะเมื่อถึงจุดที่แบ่งได้นั้น อัลกอริทึมจะมีทางเลือก (candidate) ที่จะเดินได้สองทางคือการจะแบ่งหรือไม่แบ่ง ในขั้นนี้ อัลกอริทึมจะลองเลือกทางใดทางหนึ่ง ซึ่งถ้าทำแล้วสำเร็จก็จะได้จุดแบ่งคำตามต้องการ ในทางตรงกันข้ามถ้าเลือกทางนั้นแล้ว

ไม่สำเร็จก็ต้องย้อนกลับมาอีกทางหนึ่งทำให้เกิดการเสียเวลาขึ้น ถ้าจำนวนคำมีมากขึ้นเวลาที่จะต้องเสียไปนี้จะเพิ่มมากขึ้นด้วยเช่นกัน

ในการออกแบบอัลกอริทึมตัดคำของโปรแกรม CU Writer แก้ปัญหาการเสียเวลาในการหาจุดแบ่งคำนี้โดยใช้วิธีการเขียนโปรแกรมแบบไม่เรียกตัวเอง (non-recursive) โดยจะแบ่งหน้าที่ของอัลกอริทึมออกเป็นสองส่วน ส่วนแรกจะทำหน้าที่สแกนสตริงและพิจารณาสถานะของทรีที่กำลังเดินไปจนกระทั่งไปต่อไม่ได้หรือจบสตริงนั้นหลังจากนั้น จะเข้าสู่ส่วนที่สองซึ่งจะทำหน้าที่แบ็กแทรค (back track) ในส่วนนี้จะทำหน้าที่ตรวจสอบหาคำที่ยาวที่สุดที่สามารถเป็นไปได้ เมื่อผ่านขั้นตอนนี้แล้วก็จะได้จุดแบ่งคำตามต้องการ



ศูนย์วิทยทรัพยากร
จุฬาลงกรณ์มหาวิทยาลัย