

## บทที่ 2

### แนวคิดและทฤษฎีที่ใช้

ในบทนี้จะกล่าวถึงแนวคิดและทฤษฎีที่เกี่ยวข้องที่จะนำมาใช้ในการแก้ไขปัญห โดยจะกล่าวถึง องค์ประกอบต่างๆที่เกี่ยวข้องในการค้นข้อความในเอกสารพีดีเอฟ มาตรฐานการเข้ารหัสของแบบอักษรไทย และแนวคิดการค้นข้อความโดยวิธีการเปรียบเทียบสายอักขระ

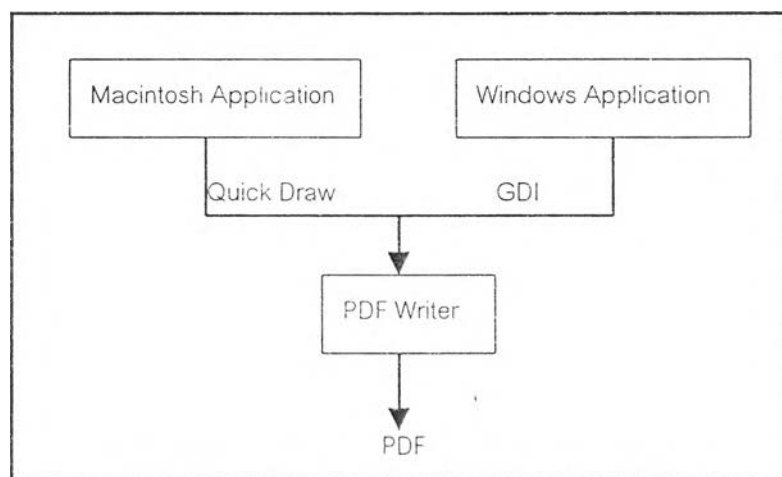
#### 2.1 เอกสารพีดีเอฟ

ก่อนที่จะสร้างกระบวนการค้นข้อความไทยในเอกสารพีดีเอฟ ควรจะทราบถึงข้อกำหนดและแนวคิดของเอกสารพีดีเอฟ เช่น วิธีในการสร้างเอกสารพีดีเอฟ โครงสร้างแฟ้มเอกสารพีดีเอฟ ข้อกำหนดประเภทแบบอักษร ข้อกำหนดการเข้ารหัสของแบบอักษรที่ใช้ในเอกสารพีดีเอฟ

##### 2.1.1 การสร้างเอกสารพีดีเอฟ

การสร้างเอกสารพีดีเอฟ<sup>1</sup> สามารถสร้างได้ 2 วิธี คือ

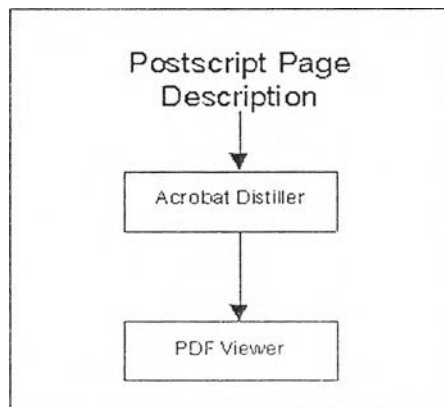
1. พีดีเอฟไรเตอร์ (PDF Writer®) สร้างเอกสารพีดีเอฟจากโปรแกรมสำเร็จประยุกต์ที่ผู้นิยมนำมาใช้โดยสั่งให้โปรแกรมสำเร็จประยุกต์นั้นพิมพ์เอกสารผ่านพีดีเอฟไรเตอร์ พีดีเอฟไรเตอร์จะทำการนำรายละเอียดการพิมพ์นั้นมาสร้างเป็นเอกสารพีดีเอฟ (ดูรูปที่ 3 ประกอบ)



รูปที่ 3 กระบวนการสร้างเอกสารพีดีเอฟโดยวิธีพีดีเอฟไรเตอร์

<sup>1</sup> Adobe Systems Incorporated.. Portable Document Format Reference Manual. Adobe Systems, 1999, pp. 19-21

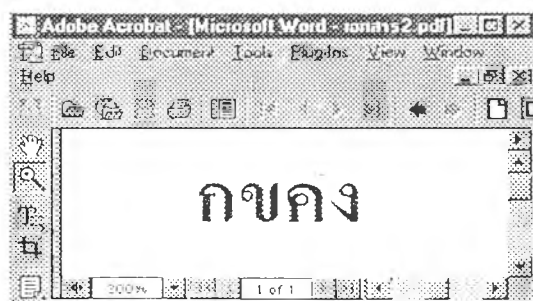
2. ดิสทิลเลอร์ (Distiller®) เป็นวิธีการสร้างเอกสารพีดีเอฟโดยจะทำการแปลงจากเอกสารโพสต์สคริปต์ (PostScript) มาเป็นเอกสารพีดีเอฟ เอกสารที่ประกอบด้วยข้อความต่างๆไปในการสร้างเป็นเอกสารพีดีเอฟสามารถทำได้ทั้ง 2 วิธี แต่ในเอกสารที่มีรูปภาพต้องการคุณภาพของรูปที่ว่าจะแสดงได้สวยงาม ควรเลือกใช้วิธีดิสทิลเลอร์ การสร้างเอกสารพีดีเอฟโดยวิธีนี้จะได้เอกสารที่มีคุณภาพดีกว่าโดยวิธีพีดีเอฟไรเตอร์ เนื่องจากผู้ใช้จะสร้างเอกสารให้อยู่ในรูปแบบเอกสารโพสต์สคริปต์ก่อน แล้วจึงใช้ดิสทิลเลอร์ทำการแปลงเอกสารนั้นเป็นเอกสารพีดีเอฟ ซึ่งในการสร้างเอกสารโพสต์สคริปต์นั้น ผู้ใช้สามารถเขียนเอกสารขึ้นเองโดยใช้คำสั่งภาษาโพสต์สคริปต์ หรือจะสร้างจากโปรแกรมสำเร็จประยุกต์ที่ผู้ใช้ใช้งานอยู่ โดยการพิมพ์ผ่านตัวขับเครื่องพิมพ์โพสต์สคริปต์ (Postscript Printer Driver) เพื่อให้ได้แฟ้มเอกสารโพสต์สคริปต์ก่อนแล้วจึงทำการสร้างเอกสารพีดีเอฟโดยดิสทิลเลอร์ (ดูรูปที่ 4 ประกอบ)



รูปที่ 4 กระบวนการสร้างเอกสารพีดีเอฟโดยวิธีดิสทิลเลอร์

### 2.1.2 โครงสร้างแฟ้มเอกสารพีดีเอฟ

หัวข้อนี้จะแสดงให้เห็นทราบถึงโครงสร้างแฟ้มเอกสารพีดีเอฟมีส่วนประกอบต่างๆอย่างไรบ้าง



รูปที่ 5 ตัวอย่างเอกสารพีดีเอฟแสดงข้อความ "กขคกง"

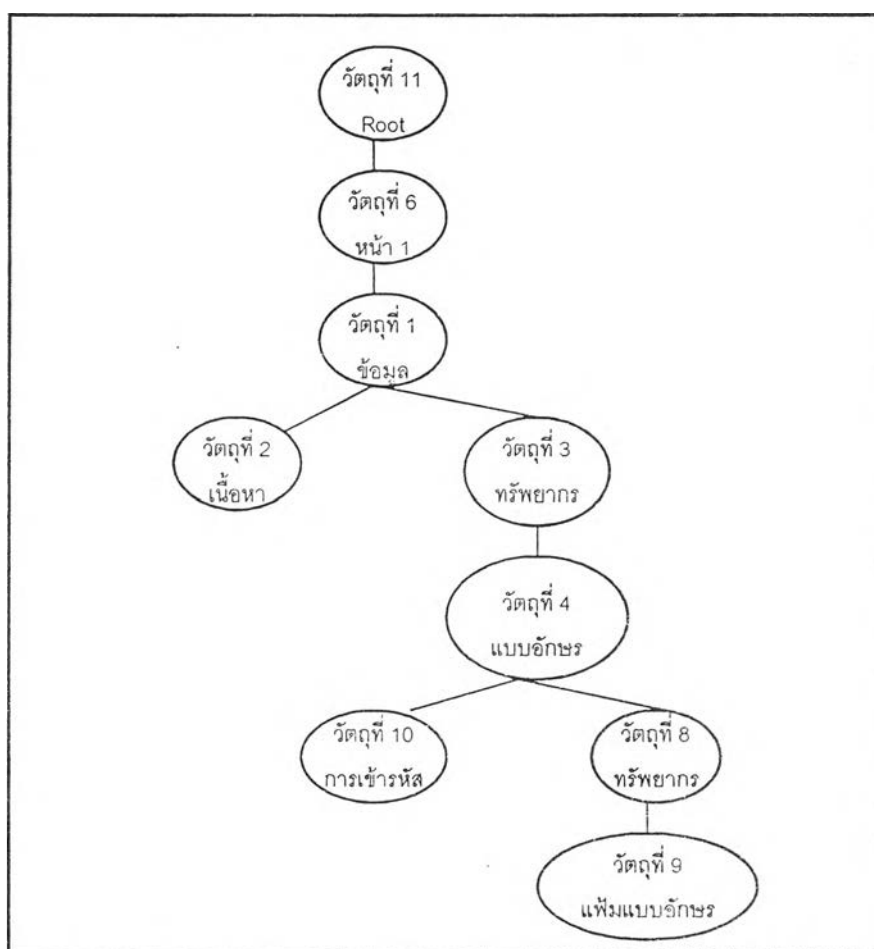


จากรูปที่ 5 และรูปที่ 6 จะแสดงตัวอย่างเอกสารพีดีเอฟและโครงสร้างภายในของแฟ้มเอกสาร ว่าแฟ้มเอกสารพีดีเอฟนี้มีโครงสร้างแฟ้มอย่างไร เอกสารจะถูกนำมาแสดงได้อย่างไร

โครงสร้างแฟ้มเอกสารพีดีเอฟ<sup>1</sup> ประกอบด้วย 4 ส่วนด้วยกัน ดังนี้

ส่วนที่หนึ่ง คือ ส่วนหัว(Header) ใช้ในการระบุว่า เอกสารพีดีเอฟนี้เป็นไปตามข้อกำหนดรุ่นที่เท่าไร ดังจะเห็นได้จากตัวอย่างในรูปที่ 6 ในบรรทัดแรกระบุ %PDF-1.2 ซึ่งหมายถึง โครงสร้างแฟ้มเอกสารพีดีเอฟนี้เป็นไปตามข้อกำหนดรุ่นที่ 1.2

ส่วนที่สอง คือ ส่วนที่ใช้ในการแสดงเอกสาร(Body) จะประกอบด้วยวัตถุต่างๆที่ใช้ในการแสดงข้อความ รูปภาพ เสียง ภาพเคลื่อนไหว และข้อมูลอื่นๆ ในส่วนนี้จะถูกนำมาใช้ในการแสดงเอกสารพีดีเอฟ วัตถุต่างๆในส่วนที่สองนี้สัมพันธ์กันในลักษณะของโครงสร้างข้อมูลแบบต้นไม้ ข้อมูลบางส่วน เช่นเนื้อหาในเอกสารหรือแฟ้มแบบอักษรที่ผนวกมากับเอกสาร โดยปกติจะถูกบีบอัดข้อมูลไว้เพื่อลดขนาดของข้อมูล ในการใช้วัตถุในเอกสารต้องทำการขยายข้อมูลที่ถูกบีบอัดไว้เสียก่อนจึงจะใช้วัตถุต่างๆในส่วนนี้ได้



รูปที่ 7 โครงสร้างต้นไม้ของวัตถุที่ใช้ในการแสดงเอกสารพีดีเอฟ

ส่วนที่สาม คือ ตารางอ้างอิง(Cross Reference Table) เป็นข้อมูลตารางที่ช่วยให้เข้าไปใช้วัตถุในเอกสารพีดีเอฟได้โดยตรง เป็นวิธีการเข้าใช้แบบ Random Access ทำให้การอ่านข้อมูลในเอกสารพีดีเอฟทำได้รวดเร็ว ข้อมูลในส่วนนี้จะต้องเริ่มต้นด้วย xref บรรทัดถัดไปจะแสดงถึงจำนวนวัตถุที่มีอยู่ในเอกสารพีดีเอฟ ตัวอักษร f จะแสดงว่าวัตถุนั้นเป็นวัตถุอิสระหรือเป็นวัตถุเริ่มต้นของ link list ตัวอักษร n จะแสดงว่าวัตถุนั้นถูกใช้ในเอกสาร ตัวเลขในแต่ละแถวจะแสดงถึงตำแหน่งเริ่มต้นของวัตถุนั้นในแฟ้มเอกสารพีดีเอฟ (ดูรูปที่ 8 ประกอบ)

```
xref
0 13
0000000000 65535 f
0000000604 00000 n
0000000816 00000 n
0000000966 00000 n
0000001098 00000 n
0000001277 00000 n
0000001720 00000 n
0000017464 00000 n
0000018545 00000 n
0000019699 00000 n
0000020179 00000 n
0000000657 00000 n
0000000796 00000 n
```

รูปที่ 8 ตัวอย่างตารางอ้างอิงในแฟ้มเอกสารพีดีเอฟ

ส่วนที่สี่ คือ ส่วนท้าย(Trailer) ในการอ่านแฟ้มเอกสารพีดีเอฟจะเริ่มอ่านที่ตอนท้ายของแฟ้มเอกสาร เพื่อให้ทราบตำแหน่งเริ่มต้นของตารางอ้างอิงเพื่อที่จะได้นำไปใช้ในการเข้าใช้วัตถุ ข้อมูลในส่วนท้ายนี้จะต้องเริ่มต้นด้วย trailer และปิดท้ายด้วย %%EOF บรรทัดก่อนที่จะระบุ %%EOF จะบอกถึงตำแหน่งเริ่มต้นของตารางอ้างอิงในแฟ้มเอกสารพีดีเอฟ นอกจากนี้ในส่วนนี้จะมีข้อมูลอื่นๆที่จะบอกให้ทราบถึง จำนวนวัตถุในแฟ้มเอกสารพีดีเอฟ วัตถุเริ่มต้น รายละเอียดในการสร้างเอกสารพีดีเอฟ เช่น ชื่อผู้สร้าง คำสำคัญ วันที่ และอื่นๆ (ดูรูปที่ 9 ประกอบ)

```
trailer
<<
Size 13
Info 12 0 P
Root 11 0 R
Prev 20458
jDj<u5304zto1477fcc8b4158fe1f5282832>]
startxref
2223
%%EOF
```

รูปที่ 9 ตัวอย่างส่วนท้ายของแฟ้มเอกสารพีดีเอฟ

### 2.1.3 แบบอักษรในเอกสารพีดีเอฟ<sup>1</sup>

โดยปกติเมื่อมีการแลกเปลี่ยนเอกสาร ผู้รับต้องมีแบบอักษรชุดเดียวกับที่ผู้ส่งใช้ในเอกสาร เพื่อให้ผู้รับจะได้เอกสารที่แสดงได้ตรงกับที่ผู้ใช้ส่งมา สำหรับในเอกสารพีดีเอฟแล้วจะใช้แบบอักษรเดียวกับที่ใช้ในเอกสารเมื่อมีแบบอักษรมันๆในเครื่องของผู้รับ แต่ถ้าไม่มีแบบอักษรมันๆจะนำแบบอักษรอื่นที่มีข้อมูลของแบบอักษรในลักษณะเดียวกันมาแสดงแทน แต่เพื่อความมั่นใจว่าเอกสารจะแสดงได้อย่างถูกต้อง ผู้สร้างเอกสารสามารถที่จะผนวกแบบอักษรที่ใช้ไปไว้ในเอกสาร วิธีนี้จะทำให้มั่นใจได้ว่าผู้รับจะได้เอกสารที่แสดงได้เหมือนกับที่ผู้สร้างสร้างเอกสาร แต่จะมีข้อเสียที่ทำให้ขนาดของเอกสารใหญ่ขึ้น

วัตถุที่ใช้ในการแสดงข้อความในเอกสารพีดีเอฟจะประกอบด้วย เนื้อหาที่ใช้ในการแสดงข้อความในเอกสารพีดีเอฟทรัพยากรที่ใช้ในการแสดงข้อความและแบบอักษร ข้อมูลแบบอักษรในเอกสารพีดีเอฟอาจมีเพียงชื่อแบบอักษรหรือมีแฟ้มชุดอักษรทั้งแฟ้มก็ได้ ขึ้นอยู่กับข้อกำหนดของผู้สร้างเอกสาร การจัดการแบบอักษรเป็นงานพื้นฐานที่สำคัญและถูกกล่าวถึงเสมอในข้อกำหนดเอกสารพีดีเอฟ

แบบอักษรเป็นข้อมูลที่สำคัญในการค้นข้อความในเอกสารพีดีเอฟ ข้อมูลแบบอักษรในเอกสารพีดีเอฟ เป็นข้อมูลชนิดพจนานุกรมที่เก็บข้อมูลหรือข้อกำหนดของแบบอักษรมันๆ เช่น ชนิดของแบบอักษร ชื่อของแบบอักษร การเข้ารหัสอักษร ข้อมูลที่ใช้ในการแสดงแบบอักษร หรือข้อมูลในการแสดงแทนเมื่อไม่มีแบบอักษรมันๆ

แบบอักษรที่ใช้ในเอกสารพีดีเอฟมีอยู่ด้วยกัน 4 ชนิด คือ

1. ประเภทที่ 0 (Type0)
2. ประเภทที่ 1 (Type1)
3. ประเภทที่ 3 (Type3)
4. ประเภททรูไทป์ (TrueType)

ประเภทที่ 0 (ในที่นี้จะขอกกล่าวถึงเพียงสังเขป ทั้งนี้เนื่องจากไม่ได้ถูกนำมาใช้กับแบบอักษรไทยที่ใช้ในเอกสารพีดีเอฟ) เป็นแบบอักษรที่ถูกออกแบบมาเพื่อสนับสนุนสำหรับอักขระที่มีตัวอักษรเป็นจำนวนมาก เช่น ภาษาจีน ภาษาญี่ปุ่น ภาษาเกาหลี ซึ่งเป็นภาษาที่มีอักขระมากกว่า 256 อักขระหรือภาษาที่เป็นอักขระภาพ การเข้ารหัสตัวอักษรของแบบอักษรมันๆจะใช้วิธีการพิเศษเพื่อที่จะสามารถเข้ารหัสอักขระจำนวนมากได้ โดยจะเก็บไว้ในข้อมูลชนิดหนึ่งที่เรียกว่าแผนที่อักษร

แบบอักษรประเภทที่ 1 มีขนาดเล็ก ให้คุณภาพตัวอักษรดี แม้ว่าจะต้องแสดงตัวอักษรขนาดเล็กบนอุปกรณ์แสดงผลที่มีรายละเอียดในการแสดงผลต่ำ เป็นแบบอักษรที่เหมาะสมที่สุดสำหรับภาษาโพสต์สคริปต์หรือเอกสารพีดีเอฟ ข้อกำหนดและรายละเอียดดังนี้

ตารางที่ 1 ข้อมูลส่วนต่างๆของแบบอักษรประเภทที่ 1

ข้อมูล	ชนิดของข้อมูล	ระบุถึง
Subtype	ชื่อ	ชนิดของแบบอักษร ต้องระบุเป็น Type1
BaseFont	ชื่อ	ชื่อของแบบอักษรโพสต์สคริปต์ที่กำหนดในแฟ้มโพสต์สคริปต์สำหรับการแสดงข้อความ เป็นชื่อที่จะถูกใช้ในการโปรแกรมภาษาโพสต์สคริปต์
FirstChar	จำนวนเต็ม	รหัสตัวอักษรตัวแรกที่ใช้ในแถวลำดับที่กำหนดความกว้างของอักขระที่ใช้ในเอกสาร
LastChar	จำนวนเต็ม	รหัสตัวอักษรตัวสุดท้ายที่ใช้ในแถวลำดับที่กำหนดความกว้างของอักขระที่ใช้ในเอกสาร
ความกว้าง	แถวลำดับ	ขนาดของแต่ละตัวอักษรที่ใช้ในเอกสารและจะระบุลำดับตัวอักษรตัวแรกและตัวอักษรตัวสุดท้ายที่ใช้ในแบบอักษร
Encoding	พจนานุกรม	การเข้ารหัสตัวอักษรของแบบอักษรที่ใช้ในเอกสาร
FontDescriptor	พจนานุกรม	ข้อมูลที่เป็นรายละเอียดต่างๆที่จะใช้ในการแสดงตัวอักษร เช่น ชื่อแบบอักษร อักขระในแบบอักษรที่ถูกใช้ในการแสดง กรอบขนาดในการแสดงแบบอักษร

```

14 0 obj <<
/Type /Font
/Subtype /Type1
/BaseFont /AGaramond-Semibold
/Encoding 25 0 R
/FontDescriptor 7 0 R
/FirstChar 0
/LastChar 255
/Widths 21 0 R
>> endobj
21 0 obj
[255 255 255 255 255 255 255 255 255 255 255 255 255 255
255 255 255 255 255 255 255 255 255 255 255 255 255 255
255 255 255 255 255 255 255 280 438 510 510 868 834
248 320 320 420 510 255 320 255 347 510 510 510 510
510 510 510 510 510 510 255 255 510 510 510 330 781
627 627 694 784 580 533 743 812 354 354 684 560 921
780 792 588 792 656 504 682 744 650 968 648 590 638
320 329 320 510 500 380 420 510 400 513 409 301 464
522 268 259 484 258 798 533 492 516 503 349 346 321
520 434 684 439 448 390 320 255 320 510 255 627 627
694 580 780 792 744 420 420 420 420 420 420 402 409
409 409 409 268 268 268 268 533 492 492 492 492 492
520 520 520 520 486 400 510 510 506 398 520 555 800
800 1044 360 380 549 846 792 713 510 549 549 510 522
494 713 823 549 274 354 387 768 615 496 330 280 510
549 510 549 612 421 421 1000 255 627 627 792 1016
730 500 1000 438 438 248 248 510 494 448 590 100 510
256 256 539 539 486 255 248 438 1174 627 580 627 580
580 354 354 354 354 792 792 790 792 744 744 744 268
380 380 380 380 380 380 380 380 380 390]
endobj

```

รูปที่ 10 ตัวอย่างแบบอักษรประเภทที่ 1 ในเอกสารพีดีเอฟ



ส่วนขยายแบบอักษรประเภทที่ 1 เป็นแบบอักษรที่เพิ่มเติมมาจากแบบอักษร ประเภทที่ 1 เพื่ออนุญาตในการสร้างรูปแบบตัวอักษรขนาดต่างๆ จากแบบอักษรแบบเดียว ในกรณีที่ไม่มีแบบอักษรนั้นในระบบที่กำลังอ่านเอกสารสามารถที่จะแสดงตัวอักษรแทนแบบอักษรนั้นได้

ตารางที่ 2 ข้อมูลส่วนต่างๆของแบบอักษรส่วนขยายประเภทที่ 1

ข้อมูล	ชนิดของข้อมูล	ระบุถึง
Subtype	ชื่อ	ชนิดของแบบอักษร ต้องระบุเป็น Multiple master Type1
BaseFont	ชื่อ	ชื่อของแบบอักษรโพสต์สคริปต์ที่ถูกกำหนดในแฟ้มโพสต์สคริปต์สำหรับใช้ในการแสดงข้อความ เป็นชื่อที่จะถูกใช้ในการโปรแกรมภาษาโพสต์สคริปต์
FirstChar	จำนวนเต็ม	รหัสตัวอักษรตัวแรกที่ใช้ในแถวลำดับที่กำหนดความกว้างของอักขระที่ใช้ในเอกสาร
LastChar	จำนวนเต็ม	รหัสตัวอักษรตัวสุดท้ายที่ใช้ในแถวลำดับที่กำหนดความกว้างของอักขระที่ใช้ในเอกสาร
Widths	แถวลำดับ	ขนาดของแต่ละตัวอักษรที่ใช้ในเอกสารและจะระบุลำดับตัวอักษรตัวแรกและตัวอักษรณตัวสุดท้ายที่ใช้ในแบบอักษร
Encoding	พจนานุกรม	การเข้ารหัสตัวอักษรของแบบอักษรที่ใช้ในเอกสาร
FontDescriptor	พจนานุกรม	ข้อมูลที่เป็นรายละเอียดต่างๆที่จะใช้ในการแสดงตัวอักษร เช่น ชื่อแบบอักษร อักขระในแบบอักษรที่ถูกใช้ในการแสดง กรอบขนาดในการแสดงแบบอักษร

```

7 0 obj <<
  /Type /Font
  /Subtype /MMType1
  /BaseFont /MinionMM_366_465_11
  /FirstChar 32
  /LastChar 255
  /Widths 19 0 R
  /Encoding 5 0 R
  /FontDescriptor 6 0 R
>>
endobj
19 0 obj
[187 235 317 430 427 717 607 168 326 326 421
619 219 317 219 282 427 427 427 427 427 427
248 320 320 420 510 255 320 255 347 510 510 510 510
510 510 510 510 510 510 255 255 510 510 510 330 781
627 627 694 784 580 533 743 812 354 354 684 560 921
780 792 588 792 656 504 682 744 650 968 648 590 638
320 329 320 510 500 380 420 510 400 513 409 301 464
522 268 259 484 258 798 533 492 516 503 349 346 321
427 427 427 427 219 219 619 619
... omitted data ...
301 301 301 569 569 0 569 607 607 607 239 400
780 792 588 792 656 504 682 744 650 968 648 590 638
320 329 320 510 500 380 420 510 400 513 409 301 464
522 268 259 484 258 798 533 492 516 503 349 346 321
400 400 400 253 400 400 400 400 400]
endobj

```

รูปที่ 11 ตัวอย่างแบบอักษรส่วนขยายประเภทที่ 1 ในเอกสารพีดีเอฟ

แบบอักษรประเภทที่ 3 เป็นแบบอักษรที่แตกต่างจากแบบอักษรอื่นๆที่ใช้ในเอกสารพีดีเอฟ เนื่องจาก แบบอักษรประเภทที่ 3 กำหนดแบบอักษรด้วยตัวมันเอง ขณะที่พจนานุกรมของแบบอักษรอื่นๆเก็บข้อมูลพื้นฐานของแบบอักษรนั้นไว้ แบบอักษรประเภทที่ 3 มีความยืดหยุ่นมากกว่าแบบอักษรประเภทที่ 1 แต่มีคุณภาพในการแสดงด้อยกว่าโดยเฉพาะ ในการแสดงตัวอักษรขนาดเล็กบนอุปกรณ์แสดงผลที่มีรายละเอียดในการแสดงผลต่ำ

ตารางที่ 3 ข้อมูลส่วนต่างๆของแบบอักษรประเภทที่ 3

ข้อมูล	ชนิดของข้อมูล	ระบุถึง
Subtype	ชื่อ	ชนิดของแบบอักษร ต้องระบุเป็น ไทป์ทรี
FirstChar	จำนวนเต็ม	รหัสตัวอักษรตัวแรกที่ใช้ในแถวลำดับที่กำหนดความกว้างของอักขระที่ใช้ในเอกสาร
LastChar	จำนวนเต็ม	รหัสตัวอักษรตัวสุดท้ายที่ใช้ในแถวลำดับที่กำหนดความกว้างของอักขระที่ใช้ในเอกสาร
Widths	แถวลำดับ	ขนาดของแต่ละตัวอักษรที่ใช้ในเอกสารและจะระบุลำดับตัวอักษรตัวแรกและตัวอักษรตัวสุดท้ายที่ใช้ในแบบอักษร
CharProcs	พจนานุกรม	ข้อมูลแต่ละตัวในพจนานุกรมคือชื่อตัวอักษร และความสัมพันธ์ที่จะนำมาใช้ในการแสดงตัวอักษร
FontBBox	พจนานุกรม	กรอบสี่เหลี่ยมที่กำหนดขนาดของแบบอักษร
FontMatrix	พจนานุกรม	ข้อมูลที่จะนำมาใช้ในการกำหนดระยะตัวอักษร
Resource	พจนานุกรม	ทรัพยากรที่จะใช้ในการแสดงอักษร
Encoding	พจนานุกรม	การเข้ารหัสตัวอักษรของแบบอักษรที่ใช้ในเอกสาร

```
6 0 obj <<
/Type /Font
/Subtype /Type3
/Name /T36
/CharProcs 10 0 R
/Resources 8 0 R
/FontBBox [3 -241 875 856]
/FontMatrix [.001 0 0 .001 0 0]
/FirstChar 1
/LastChar 4
/Widths [47 44 49 37]
/Encoding 9 0 R
>> endobj
9 0 obj
<<
/Type /Encoding
/Differences [1/G01 /G02 /G03 /G04 ]
>>
endobj
10 0 obj
<<
/G01 11 0 R
/G02 12 0 R
/G03 13 0 R
/G04 14 0 R
>>
endobj
```

รูปที่ 12 ตัวอย่างชุดแบบตัวประเภทที่ 3 ในเอกสารพีดีเอฟ

แบบอักษรประเภททรูไทป์ เป็นแบบอักษรชนิดปรับขนาดได้ ปัจจุบันถูกนำมาใช้เป็นแบบอักษรมาตรฐานระบบปฏิบัติการวินโดวส์ และเป็นมาตรฐานหนึ่งในการแสดงตัวอักษรนี้ในเอกสารพีดีเอฟ ข้อกำหนดและรายละเอียดของแบบอักษรนี้ มีรูปแบบเหมือนกับแบบอักษรประเภทที่ 1

ตารางที่ 4 ข้อมูลส่วนต่างๆของแบบอักษรประเภททรูไทป์

ข้อมูล	ชนิดของข้อมูล	ระบุถึง
Subtype	ชื่อ	ชนิดของแบบอักษร ต้องระบุเป็นทรูไทป์
BaseFont	ชื่อ	ชื่อของแบบอักษรโพสต์สคริปต์ที่กำหนดในแฟ้มโพสต์สคริปต์สำหรับการแสดงข้อความ เป็นชื่อที่จะถูกใช้ในการโปรแกรมภาษาโพสต์สคริปต์
FirstChar	จำนวนเต็ม	รหัสตัวอักษรตัวแรกที่ใช้ในแถวลำดับที่กำหนดความกว้างของอักขระที่ใช้ในเอกสาร
LastChar	จำนวนเต็ม	รหัสตัวอักษรตัวสุดท้ายที่ใช้ในแถวลำดับที่กำหนดความกว้างของอักขระที่ใช้ในเอกสาร
Widths	แถวลำดับ	ขนาดของแต่ละตัวอักษรที่ใช้ในเอกสารและจะระบุลำดับตัวอักษรตัวแรกและตัวอักษรตัวสุดท้ายที่ใช้ในแบบอักษร
Encoding	พจนานุกรม	การเข้ารหัสตัวอักษรของแบบอักษรที่ใช้ในเอกสาร
FontDescriptor	พจนานุกรม	ข้อมูลที่เป็นรายละเอียดต่างๆที่จะใช้ในการแสดงตัวอักษร เช่น ชื่อแบบอักษร อักขระในแบบอักษรที่ถูกใช้ในการแสดง กรอบขนาดในการแสดงแบบอักษร

```

17 0 obj <<
  /Type /Font
  /Subtype /TrueType
  /BaseFont /NewYork,Bold
  /FirstChar 0
  /LastChar 255
  /Widths 23 0 R
  /Encoding /WindowsEncoding
  /FontDescriptor 7 0 R
>> endobj

23 0 obj
[0 333 333 333 333 333 333 333 333 0 333 333 333 333 333
333 333 333 333 333 333 333 333 333 333 333 333 333
333 333 0 333 333 333 303 500 666 666 882 848 303
446 446 507 666 303 378 303
... omitted data ...
303 530 1280 757 605 757 605 605 355 355 355 355 803
803 790 803 780 780 780 340 636 636 636 636 636 636
636 636 636 636]
endobj

7 0 obj
<<
  /Type /FontDescriptor
  /FontName /CordialUPC
  /Flags 4
  /FontBBox [ -250 -488 950 1000 ]
  /ItalicAngle 0
  /CapHeight 898
  /Ascent 898
  /MaxWidth 792
>>
endobj

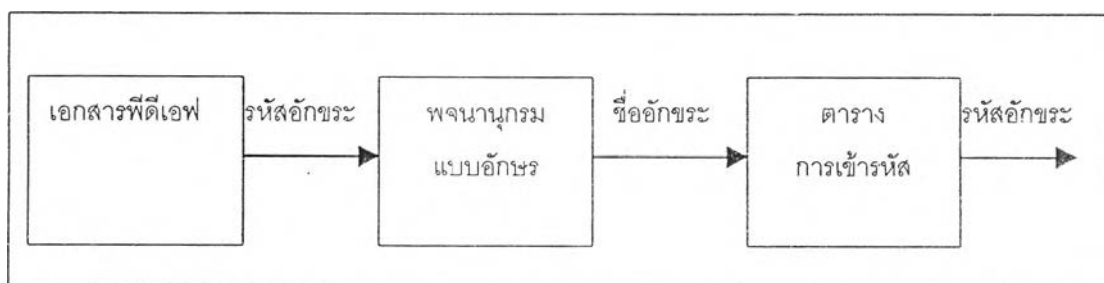
```

รูปที่ 13 ตัวอย่างแบบอักษรประเภททรูไทป์ ในเอกสารพีดีเอฟ

### 2.1.4 การเข้ารหัสตัวอักษร<sup>1</sup>

การเข้ารหัสตัวอักษรเป็นส่วนสำคัญที่จะใช้ในการระบุตัวอักษร ว่าตัวอักษรนี้คือตัวอักษรอะไรมีรูปร่างอย่างไร เพื่อนำไปใช้ในการแสดงตัวอักษรและการค้นข้อความในเอกสารพีดีเอฟ การเข้ารหัสตัวอักษรในเอกสารพีดีเอฟใช้โครงสร้างการเข้ารหัสตัวอักษรที่มีความยืดหยุ่น โดยจะเก็บข้อมูลการเข้ารหัสและข้อมูลในการอธิบายโครงสร้างตัวอักษรไว้แยกต่างหาก (ดังจะเห็นได้จากตัวอย่างแบบอักษรที่เคยได้กล่าวมาแล้ว ในหัวข้อที่ผ่านมา)

การแสดงข้อความในเอกสารโดยทั่วไป เอกสารต่างๆ ใช้รหัสอักขระในการระบุตัวอักษรในเอกสาร เช่น รหัสอักขระ 65 ก็จะเป็นตัวอักษร "A" หรือ รหัสอักขระ 161 ก็จะเป็นตัวอักษร "ก" แต่ในเอกสารพีดีเอฟ รหัสอักขระ 65 ไม่จำเป็นต้องเป็นตัวอักษร "A" ถ้ารหัสอักขระ 65 ในพจนานุกรมแบบอักษร เข้ารหัส 65 /at รหัสอักขระ 65 ก็จะเป็นตัวอักษร "@" ถ้ารหัสอักขระ 161 ในพจนานุกรมแบบอักษร เข้ารหัส 161 /zero รหัสอักขระ 161 ก็จะเป็นตัวอักษร "0" การแสดงตัวอักษรหรือถอดรหัสตัวอักษรในเอกสารพีดีเอฟ จะอ่านค่ารหัสอักขระในเอกสารพีดีเอฟ ไปดูว่าอักขระนั้นมีชื่ออักขระว่าอะไร แล้วทำการแปลงรหัสอักขระนั้นตามตารางการเข้ารหัสที่ได้กำหนดไว้ ( ดูรูปที่ 14 ประกอบ )



รูปที่ 14 ความสัมพันธ์ระหว่างรหัสตัวอักษรและชื่อตัวอักษรในเอกสารพีดีเอฟ

การเข้ารหัสตัวอักษรในเอกสารพีดีเอฟ มีอยู่ด้วยกัน 3 แบบ คือ

1. มาตรฐาน(Standard) คือการเข้ารหัสตัวอักษรจะเป็นไปตามมาตรฐานที่ข้อกำหนดเอกสารพีดีเอฟกำหนดไว้ เช่น การเข้ารหัสตัวอักษรตามมาตรฐานวินโดวส์ (WinAnsi) (ดูรูปที่ 15 ประกอบ) หรือ การเข้ารหัสตัวอักษรตามมาตรฐานอโดบี (Standard) (สามารถดูข้อกำหนดมาตรฐานการเข้ารหัสตัวอักษรตามข้อกำหนด ของเอกสารพีดีเอฟได้ที่ ภาคผนวก ก)

<sup>1</sup> Adobe Systems Incorporated, *PostScript Language Reference Manual, Third Edition*. Addison-Wesley, 1990, p.369

```
20 0 obj
<< /Type /Font
/Subtype /TrueType
/Name /F1
/BaseFont /CordiaUPC
/Encoding /WinAnsiEncoding
/FontDescriptor 22 0 R >>
endobj
```

รูปที่ 15 ตัวอย่างการเข้ารหัสตัวอักษรในเอกสารพีดีเอฟ  
ตามข้อกำหนดมาตรฐานวินโดวส์

2. ผู้ใช้กำหนดการเข้ารหัสตัวอักษรเอง(Custom) คือ ข้อมูลการเข้ารหัสจะถูกกำหนดไว้ที่  
พจนานุกรมแบบอักษรแยกออกมาเป็นส่วนข้อมูลอีกส่วนหนึ่ง

```
4 0 obj
<< /Name /F2
/Type /Font
/Subtype /Type3
/Encoding 9 0 R
/CharProcs 10 0 R
/Widths [47 44 49 37] >>
endobj
9 0 obj
<<
/Type /Encoding
/Differences [1/afii59681 /afii59682 /afii59684 /afii59687]
>>
endobj
```

รูปที่ 16 ตัวอย่างการเข้ารหัสตัวอักษรในเอกสารพีดีเอฟ  
แบบผู้ใช้กำหนดการเข้ารหัสตัวอักษรเอง



3. แบบใช้ข้อกำหนดการเข้ารหัสตัวอักษรที่มีอยู่ในแฟ้มแบบอักษร(Built-In) คือ ไม่มีข้อมูลการเข้ารหัสตัวอักษรในเอกสารพีดีเอฟ ให้ใช้การเข้ารหัสตัวอักษรที่ถูกกำหนดอยู่ภายในแฟ้มแบบอักษรนั้น

```

14 0 obj
<<
/Type /Font
/Subtype /TrueType
/FirstChar 33
/LastChar 45
/Widths [ 400 376 417 313 405 347 0 417 0 382 347 243 405 ]
/BaseFont /BHDGMA+CordiaUPC+1
/FontDescriptor 11 0 R
/ToUnicode 12 0 R
>>
obj

```

รูปที่ 17 ตัวอย่างการเข้ารหัสตัวอักษรในเอกสารพีดีเอฟ  
แบบใช้ข้อกำหนดที่มีอยู่ในแฟ้มแบบอักษร

## 2.2. มาตรฐานการเข้ารหัสอักขระไทย

### 2.2.1 มาตรฐานมอก.620

หน่วยงานมาตรฐานผลิตภัณฑ์อุตสาหกรรมไทยได้ประกาศใช้รหัสสำหรับอักขระไทยที่ใช้กับคอมพิวเตอร์เป็นครั้งแรก ที่มาตรฐานเลขที่ มอก.620-2529 ในราชกิจจานุเบกษา เล่ม 103 ตอนที่ 102 วันที่ 17 มิถุนายน พุทธศักราช 2529 แต่เนื่องจากยังขาดรายละเอียดเกี่ยวกับหมายเลขอักขระและชื่ออักขระ ซึ่งเป็นส่วนสำคัญที่จะต้องใช้ในการจดทะเบียนอักขระไทยกับองค์การระหว่างประเทศว่าด้วยการมาตรฐาน หน่วยงานมาตรฐานผลิตภัณฑ์อุตสาหกรรมไทยจึงทำการแก้ไขและกำหนดมาตรฐานที่ มอก.620-2533 ดังนี้



## 2.2.2 มาตรฐานแบบอักษรไทยในระบบปฏิบัติการวินโดวส์

บริษัทไมโครซอฟต์ (Microsoft Inc.) ได้กำหนดการใช้และแสดงตัวอักษรภาษาไทยในระบบปฏิบัติการวินโดวส์ ตามข้อกำหนดมาตรฐาน มอก.620 และมาตรฐานยูนิโค้ด (Unicode) โดยกำหนดการเข้ารหัสที่การเข้ารหัสหน้า 874 มีรายละเอียดและข้อกำหนดดังนี้

ตารางที่ 5 ตารางการเข้ารหัสแบบอักษรไทยในระบบปฏิบัติการวินโดวส์

ตัวอักษร	รหัสยูนิโค้ด	ชื่อตัวอักษร	ตัวอักษร	รหัสยูนิโค้ด	ชื่อตัวอักษร
ก	0E01	afii59681	ต	0E15	afii59701
ข	0E02	afii59682	ถ	0E16	afii59702
ช	0E03	afii59683	ท	0E17	afii59703
ค	0E04	afii59684	ธ	0E18	afii59704
ค	0E05	afii59685	น	0E19	afii59705
ฅ	0E06	afii59686	บ	0E1A	afii59706
ง	0E07	afii59687	ป	0E1B	afii59707
จ	0E08	afii59688	ผ	0E1C	afii59708
ฉ	0E09	afii59689	ฝ	0E1D	afii59709
ช	0E0A	afii59690	พ	0E1E	afii59710
ช	0E0B	afii59691	ฟ	0E1F	afii59711
ฌ	0E0C	afii59692	ภ	0E20	afii59712
ญ	0E0D	afii59693	ม	0E21	afii59713
ฎ	0E0E	afii59694	ย	0E22	afii59714
ฏ	0E0F	afii59695	ร	0E23	afii59715
ฐ	0E10	afii59696	ฤ	0E24	afii59716
ฑ	0E11	afii59697	ล	0E25	afii59717
ฒ	0E12	afii59698	ฬ	0E26	afii59718
ณ	0E13	afii59699	ว	0E27	afii59719
ด	0E14	afii59700	ศ	0E28	afii59720

ตารางที่ 5 (ต่อ) ตารางการเข้ารหัสแบบอักษรไทยในระบบปฏิบัติการวินโดวส์

ตัวอักษร	รหัสยูนิก็ัด	ชื่อตัวอักษร	ตัวอักษร	รหัสยูนิก็ัด	ชื่อตัวอักษร
ษ	0E29	afii59721	อ๋	0E48	afii59752
ส	0E2A	afii59722	อ๊	0E49	afii59753
ห	0E2B	afii59723	อึ	0E4A	afii59754
ฬ	0E2C	afii59724	อ๋	0E4B	afii59755
อ	0E2D	afii59725	อ็	0E4C	afii59756
ฮ	0E2E	afii59726	อ๋	0E4D	afii59758
ฯ	0E2F	afii59727	๑	0E4F	afii59759
๕	0E30	afii59728	๐	0E50	afii59760
อ๋	0E31	afii59729	๑	0E51	afii59761
อา	0E32	afii59730	๒	0E52	afii59762
อ๋า	0E33	afii59731	๓	0E53	afii59763
อื่อ	0E34	afii59732	๔	0E54	afii59764
อื่อ	0E35	afii59733	๕	0E55	afii59765
อื่อ	0E36	afii59734	๖	0E56	afii59766
อื่อ	0E37	afii59735	๗	0E57	afii59767
อู	0E38	afii59736	๘	0E58	afii59768
อู	0E39	afii59737	๙	0E59	afii59769
	0E3A	afii59738	๙	0E5A	afii59770
๘	0E3F	afii59743	๐	0E5B	afii59771
เ	0E40	afii59744	อ๋	f701	f702
แ	0E41	afii59745	อ๋	f702	f702
โ	0E42	afii59746	อ๋	f703	f703
ใ	0E43	afii59747	อ๋	f704	f704
ไ	0E44	afii59748	อ	f705	f705
ำ	0E45	afii59749	อ๋	f706	f706
ำ	0E46	afii59750	อ๋	f707	f707
อ๋	0E47	afii59751	อ๋	f708	f708

ตารางที่ 5 (ต่อ) ตารางการเข้ารหัสแบบอักษรไทยในระบบปฏิบัติการวินโดวส์

ตัวอักษร	รหัสยูนิโค้ด	ชื่อตัวอักษร	ตัวอักษร	รหัสยูนิโค้ด	ชื่อตัวอักษร
อ	f709	f709	อ	f712	f712
อ	f70A	f70A	อ	f713	f713
อ	f70B	f70B	อ	f714	f714
อ	f70C	f70C	อ	f715	f715
อ	f70D	f70D	อ	f716	f716
อ	f70E	f70E	อ	f717	f717
ญ	f70F	f70F	อ	f718	f718
อ	f710	f710	อ	f719	f719
อ	f711	f711	.	f71A	f71A

### 2.3 การค้นข้อความโดยวิธีการเปรียบเทียบสายอักขระ<sup>3</sup>

การค้นข้อความโดยวิธีการเปรียบเทียบสายอักขระ หลักการของวิธีนี้คือ จะทำการเปรียบเทียบกลุ่มอักขระ 2 กลุ่ม คือกลุ่มอักขระของข้อความที่ต้องการค้นกับกลุ่มอักขระของข้อความที่ต้องทำการเปรียบเทียบทีละ 1 คู่อักขระจากซ้ายไปขวา การค้นโดยวิธีนี้เป็นกรค้นหาจากข้อมูลจริงโดยตรง ไม่มีการสร้างรหัสขึ้นมาเป็นดัชนีเพื่อการเปรียบเทียบแต่อย่างใด การค้นหาข้อความเป็นการเปรียบเทียบที่ตรงกัน จะค้นพบข้อความได้ก็ต่อเมื่อมีค่าที่ต้องการค้นปรากฏในข้อมูลเหล่านั้น เช่น ถ้าค้นคำว่า "Exact" จะค้นพบคำนี้ก็ต่อเมื่อมีคำว่า "Exact" ปรากฏอยู่จริงในกลุ่มอักขระนั้น

แนวคิดวิธีการค้นข้อความโดยวิธีการเปรียบเทียบสายอักขระ

#### 2.3.1 การเปรียบเทียบสายอักขระตามแนวคิดของ Brute-force

เป็นแนวคิดเบื้องต้นในการทำการเปรียบเทียบข้อความโดยตรงไปตรงมา คือ เปรียบเทียบกันทีละตัวอักษร ถ้าตัวอักษรตัวแรกของข้อความที่ต้องการค้นตรงกับตัวอักษรของข้อความที่นำมาเปรียบเทียบก็จะขยับไป 1 ตัวอักษรแล้วเปรียบเทียบตัวอักษรถัดไปเรื่อยๆ ถ้าไม่ตรงกันเมื่อไรจะเลื่อนตัวอักษรของข้อความที่ต้องการค้นกลับมาที่ตัวอักษรตัวแรก วิธีกรนี้จะเสียเวลามากเพราะจะต้องเปรียบเทียบอักขระทุกอักขระของข้อความที่นำมาเปรียบเทียบ

<sup>3</sup> Jun-ichi Aoe, *Computer Algorithm: string pattern matching strategies*. IEEE Computer Society Press, 1994, p.1-4

### 2.3.2 การเปรียบเทียบสายอักขระตามแนวคิดของ Knuth-Morris-Pratt

เป็นแนวคิดที่พัฒนาต่อมาจาก Brute-force โดยแก้ไขส่วนที่บกพร่องในเรื่องของการขยับตัวอักษรเพื่อนำไปเปรียบเทียบ โดยในการขยับตัวอักษรจะขยับได้มากกว่าครั้งละ 1 ตัวอักษรแต่จะไม่เกินความยาวของคำที่ค้น การที่ขยับได้ครั้งละมากกว่า 1 ตัวอักษรทำให้การค้นไม่ต้องทำการเปรียบเทียบทุกตัวอักษรทำให้ค้นข้อความได้รวดเร็วกว่าวิธีการของ Brute-force

### 2.3.3 การเปรียบเทียบสายอักขระตามแนวคิดของ Boyer และ Moore

เป็นแนวคิดที่ได้รับความนิยม มีแนวคิดคล้ายกับ 2 แบบแรกที่ได้กล่าวมา แต่เพิ่มประสิทธิภาพในการเลื่อนตัวอักษรให้ดีขึ้นกว่าทั้ง 2 แบบ โดยมีแนวความคิดดังนี้

ให้  $i$  คือ ตำแหน่งตัวอักษรของข้อความที่ต้องการค้น (Pattern)

$j$  คือ ตำแหน่งตัวอักษรของข้อความทั้งหมด (String)

$n$  คือ จำนวนตัวอักษรของข้อความที่ต้องการค้น (Length of Pattern)

ในการเลื่อนตัวอักษรเปรียบเทียบ จะมีวิธีการเลื่อนตัวอักษรดังนี้

ให้  $j$  เริ่มต้นเท่ากับ  $n$

ถ้า  $\text{String}[j]$  ตรงกับตัวอักษรตัวใดๆใน  $\text{Pattern}[i]$  ให้  $j = j + (n - i)$

ถ้าไม่แล้ว ให้  $j = j + n$

ตัวอย่างเช่น	$\text{Pattern} = \text{"corn"}$	$\text{String} = \text{"Oaks of acorns"}$
ขั้นตอนที่ 1	$\text{corn}$	$\text{Oaks of acorns} \Rightarrow j = 4, \text{String}[4] \text{ ไม่ตรงกับ 'c','o','r','n'}$
ขั้นตอนที่ 2	$\text{corn}$	เมื่อ $\text{String}[4]$ ไม่ตรงกับตัวอักษรใดๆ ให้ $j = j + n \Rightarrow j = 4 + 4$
ขั้นตอนที่ 3	$\text{corn}$	$\text{Oaks of acorns} \Rightarrow j = 8, \text{String}[8] \text{ ไม่ตรงกับ 'c','o','r','n'}$
ขั้นตอนที่ 4	$\text{corn}$	เมื่อ $\text{String}[8]$ ไม่ตรงกับตัวอักษรใดๆ ให้ $j = j + n \Rightarrow j = 8 + 4$
	$\text{corn}$	$\text{Oaks of acorns} \Rightarrow j = 12, \text{String}[12] \text{ ตรงกับ Pattern}[3]$
	$\text{corn}$	เมื่อ $\text{String}[12]$ ตรงกับ $\text{Pattern}[3]$ ให้ $j = j + (n-i) \Rightarrow j = 12 + 1$
	$\text{corn}$	$\text{Oaks of acorns} \Rightarrow j = 13, \text{พบข้อความที่ต้องการค้น}$

จากตัวอย่างจะเห็นว่าพบ "corn" ในข้อความ "Oaks of acorns" แต่ถ้าหากเปรียบเทียบแล้ว ข้อความที่ต้องการค้นยังไม่ตรงกับข้อความที่ใช้เปรียบเทียบ ก็ให้ทำการเลื่อนตัวอักษรและตรวจสอบเหมือนที่ได้กล่าวมาแล้ว