

บทที่ 4

การพัฒนาระบบการค้นคืนข้อความภาษาไทยโดยใช้ดัชนีไม่แพ้

4.1 การพัฒนาระบบการค้นคืนข้อความภาษาไทยโดยใช้ดัชนีไม่แพ้

การพัฒนาระบบการค้นคืนข้อความภาษาไทยโดยใช้ดัชนีไม่แพ้ จะแสดงพร้อมอัลกอริทึมการสร้างดัชนีไม่แพ้ แพ้ต่อระยะ การค้นคืนข้อความ ทั้งรูปภาพและตัวอย่างประกอบ เพื่อเพิ่มความเข้าใจมากยิ่งขึ้น

4.2 ระบบการค้นคืนข้อความภาษาไทย

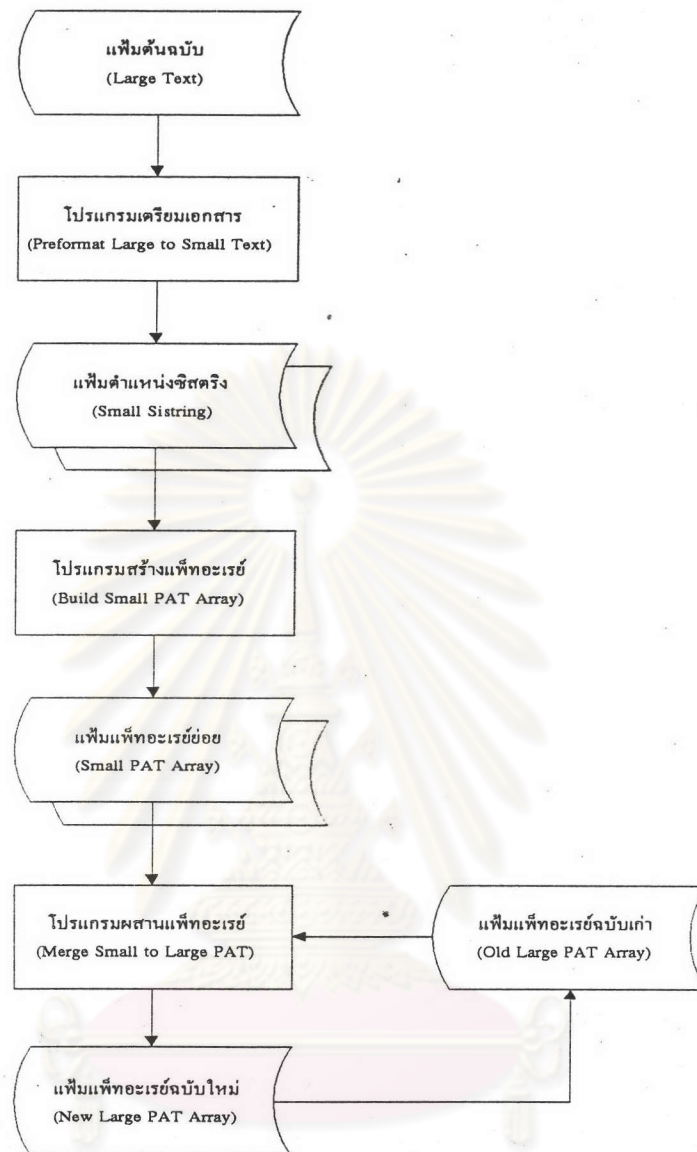
ระบบการค้นคืนข้อความภาษาไทยแบ่งออกเป็นสองส่วนหลักๆ คือ

1. การสร้างแพ็คเกจของเอกสารเพื่อใช้ค้นคืนภายหลัง
2. การค้นคืนแบบเต็มหน้ากับแบบบูล

4.3 การสร้างแพ็คเกจ

การสร้างแพ็คเกจเป็นการสร้างดัชนีให้กับเอกสารวิธีหนึ่งที่มีประสิทธิภาพ โดยแบ่งออกเป็น 3 ขั้นตอน ดังนี้

1. การเตรียมเอกสาร (Preformat Large to Small Text)
2. การสร้างแพ็คเกจ (Build Small PAT Array)
3. การผสานแพ็คเกจ (Merge Small to Large PAT Array)



รูปที่ 4.1 การสร้างแพ็ทอะเรย์

รูปที่ 4.1 แสดงขั้นตอนการสร้างแพ็ทอะเรย์ของเพิ่มต้นฉบับ (Large Text) ผ่านโปรแกรมตามขั้นตอนต่างๆ จนได้ผลลัพธ์สุดท้าย คือ เพิ่มแพ็ทอะเรย์ฉบับใหม่สุด (New Large PAT Array) โดยมีหลักการ 3 ขั้นตอน คือ

1. การตัดแบ่งเพิ่มต้นฉบับ (Large Text) ออกเป็น เพิ่มต้นฉบับย่อย (Small Text) หลายๆ ฉบับ พร้อมทั้งค้นหาตำแหน่งซิสตริงของแต่ละฉบับ และบันทึกเฉพาะ ตำแหน่งซิสตริงลงในเพิ่มตำแหน่งซิสตริง (Small Sistring) โดยอาศัยโปรแกรมเตรียมเอกสาร (Preformat Large to Small Text)

2. การสร้างแฟ้มแพ็คเกจย่อย (Small PAT Array) ซึ่งประกอบด้วยแฟ้มแพ็คเกจ (PAT Array) กับแฟ้มแพ็คเกจซ้ำ (Duplicate PAT Array) จากแต่ละแฟ้มตำแหน่ง ซิสตริง (Small Sistring) โดยอาศัยโปรแกรมสร้างแพ็คเกจ (Build Small PAT Array)

3. การผสานแต่ละแฟ้มแพ็คเกจย่อย (Small PAT Array) เข้ากับแฟ้มแพ็คเกจฉบับเก่า (Old Large PAT Array) เก็บลงแฟ้มแพ็คเกจฉบับใหม่ (New Large PAT Array) โดยอาศัย โปรแกรมผสานแพ็คเกจ (Merge Small to Large PAT Array)

หลังจากสร้างแพ็คเกจเสร็จแล้ว นำแฟ้มแพ็คเกจฉบับใหม่ไปใช้ในการค้นคืน ข้อความต่อไป

4.3.1 การเตรียมเอกสาร

ดังที่ได้กล่าวไว้แล้วเกี่ยวกับลักษณะและปริมาณของเอกสาร พบว่า ส่วนใหญ่แล้วเอกสาร ที่นำมาจัดเก็บเพื่อค้นคืนภายหลัง มักจะมีปริมาณมากเกินกว่าที่เครื่องคอมพิวเตอร์จะสามารถ ประมวลผลได้เสร็จสิ้นภายในหนเดียวได้ โดยจะทยอยสร้างแฟ้มข้อมูลเสริมเพื่อการค้นคืนไป เรื่อยๆ ดังนั้นขั้นตอนการเตรียมเอกสารจะจัดแบ่งเอกสารฉบับใหญ่ๆ หรือที่เรียกว่าแฟ้มต้นฉบับ ออกเป็นฉบับย่อยๆ เช่น 64,000 หรือ 100,000 ไบต์เป็นต้น ซึ่งมีความเหมาะสมกับการประมวล ผล ทั้งนี้ขึ้นกับจำนวนหน่วยความจำหลัก ของเครื่องคอมพิวเตอร์นั้นๆ โดยกำหนดเป็นพารามิ เตอร์ (Parameter) ตัวหนึ่งของระบบ จากนั้นจะคำนวณหาตำแหน่งซิสตริงของแฟ้มต้นฉบับย่อย แต่ละฉบับ แล้วจัดเก็บเฉพาะ ตำแหน่งซิสตริง ลงแฟ้มตำแหน่งซิสตริง เพื่อใช้ในการ สร้าง แพ็คเกจต่อไป โดยวิธีการคำนวณอาศัยหลักทางภาษาศาสตร์เบื้องต้นประกอบ คือ

1. เป็นซิสตริงภาษาอังกฤษที่ขึ้นต้นด้วย a-z, A-Z และ 0-9
2. เป็นซิสตริงภาษาไทยที่ขึ้นต้นด้วยพยัญชนะ ก-ฮ, สระ เอ, แอ, โอ, ใ, ออ และ 0-๙

อย่างไรก็ตามสำหรับกฎทางภาษาศาสตร์สามารถปรับปรุงให้มีความถูกต้องและเหมาะสม มากขึ้นในภายหลังได้

4.3.1.1 ขั้นตอนการเตรียมเอกสาร

1. โหลดข้อมูลขนาดเท่ากับพารามิเตอร์กำหนดขนาดเพิ่มต้นฉบับย่อย จากเพิ่มต้นฉบับเก็บไว้ในหน่วยความจำ
2. ตรวจสอบข้อมูลที่ละตำแหน่ง (ตัวอักษร) ว่าเป็นซีสตริงที่ถูกต้องตามกฎทางหลักภาษาศาสตร์เบื้องต้นหรือไม่ ถ้าถูกต้องบันทึกเฉพาะ ตำแหน่งซีสตริง ลงเพิ่มตำแหน่งซีสตริง จนกระทั่งหมดข้อมูลตรวจสอบ
3. ปฏิบัติตามขั้นตอนที่ 1 กับ 2 จนหมดเพิ่มต้นฉบับ ทำให้ได้เพิ่มตำแหน่งซีสตริงหลายๆ เพิ่มข้อมูล เพื่อนำไปสร้างต้นไม้แพ็คในขั้นตอนต่อไป

4.3.1.2 ข้อกำหนดเบื้องต้นของการเตรียมเอกสาร

1. พารามิเตอร์กำหนดชื่อเพิ่มต้นฉบับ จะต้องมีนามสกุลเป็น เอช ที เอ็ม (HTM) เสมอ
2. ต้องระบุ พารามิเตอร์กำหนดขนาดเพิ่มต้นฉบับย่อย
3. ต้องระบุ พารามิเตอร์กำหนดความยาวซีสตริง
4. มีทางเลือกพิมพ์ข้อมูลทางสถิติหรือซีสตริงเพื่อตรวจสอบก่อนการสร้างแพ็คเกจ ซึ่ง
 - เพิ่มข้อมูล `<input>.Brr` เก็บข้อมูลทางสถิติ โดย
 - `<input>` คือ ชื่อเพิ่มต้นฉบับ
 - B คือ ตัวระบุว่าเป็นเพิ่มข้อมูลประเภทเก็บข้อมูลสถิติ
 - rr คือ ตัวเลขต่อเนื่องกัน (Running Number)
 - เพิ่มข้อมูล `<input>.Drr` เก็บซีสตริง โดย
 - D คือ ตัวระบุว่าเป็นเพิ่มข้อมูลประเภทเก็บซีสตริง
5. เพิ่มตำแหน่งซีสตริงที่ถูกตัดแบ่งจากเพิ่มต้นฉบับ คือ `<input>.Irr` โดย
 - I คือ ตัวระบุว่าเป็นเพิ่มข้อมูลประเภทเก็บตำแหน่งซีสตริง

4.3.1.3 โครงสร้างเพิ่มตำแหน่งซีสตริง

ประกอบด้วย 2 ส่วน ดังนี้

1. ส่วนหัว (Header) มีองค์ประกอบหลักของข้อมูล คือ

ชื่อข้อมูล	รายละเอียดของข้อมูล
1. Large Text	พารามิเตอร์กำหนดชื่อเพิ่มเติมฉบับ
2. Size Small Text	พารามิเตอร์กำหนดขนาดเพิ่มเติมฉบับย่อย
3. Base Large Text	ตำแหน่งเริ่มต้นเพิ่มเติมฉบับ
4. Length Small Text	ความยาวเพิ่มเติมฉบับย่อย
5. Number of Sistrings	จำนวนซิสตริงทั้งหมด
6. Length Sistring	พารามิเตอร์กำหนดความยาวซิสตริง
7. Number of Small Text	จำนวนเพิ่มเติมฉบับย่อย
8. Flag Last Small Text	ตัวระบุเพิ่มเติมฉบับย่อยเป็นฉบับสุดท้าย

ตารางที่ 4.1 ส่วนหัวเพิ่มเติมตำแหน่งซิสตริง

2. ส่วนข้อมูล (Detail) เก็บตำแหน่งซิสตริงที่ใช้ในการอ้างอิงถึงซิสตริงในเพิ่มเติมฉบับ

สำหรับขั้นตอนการเตรียมเอกสาร อาศัยหน่วยความจำขนาดเท่ากับพารามิเตอร์กำหนดขนาดเพิ่มเติมฉบับย่อย และ ประสิทธิภาพของโปรแกรม ขึ้นกับปัจจัย ดังต่อไปนี้

1. เวลาทั้งหมดในการอ่านเพิ่มเติมฉบับจากหน่วยความจำสำรองเก็บในหน่วยความจำ
2. เวลาทั้งหมดในการตรวจสอบข้อมูลที่ตัวอักษร แต่เนื่องจากประมวลผลในหน่วยความจำ จึงใกล้เคียงศูนย์
3. เวลาทั้งหมดในการบันทึกเพิ่มเติมตำแหน่งซิสตริงทั้งหมด

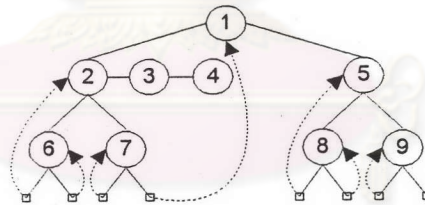
4.3.2 การสร้างแพ็คเกจ

การสร้างแพ็คเกจ คือ การจัดเรียงลำดับซิสตริงแบบหนึ่ง โดยนำ ตำแหน่งซิสตริง ที่ใช้ในการอ้างอิงถึงซิสตริงในเพิ่มเติมฉบับ ของเพิ่มเติมตำแหน่งซิสตริง ตามที่ได้แบ่งไว้ในขั้นตอน การเตรียมเอกสาร มาสร้างต้นไม้แพ็คเกจ เพื่อหาแพ็คเกจ ซึ่งแพ็คเกจที่ได้จากต้นไม้แพ็คเกจ คือ ตำแหน่งซิสตริงที่ถูกเรียงลำดับตามรหัสข้อมูลของซิสตริง แล้วจัดเก็บแพ็คเกจลงเพิ่มแพ็คเกจย่อย โดยในกรณีของการสร้างแพ็คเกจที่ไม่เกิดการซ้ำกันของซิสตริง จะจัดเก็บ

แพ็คเกจจะเรียงลงในส่วนของแพ็คเกจที่แพ็คเกจ ส่วนในกรณีเกิดการซ้ำกันของซิสตริง จะจัดเก็บแพ็คเกจที่ซ้ำลงในแพ็คเกจที่แพ็คเกจ และจัดเก็บแพ็คเกจที่ซ้ำกับตำแหน่งแพ็คเกจที่ซ้ำของแพ็คเกจที่แพ็คเกจ ลงในแพ็คเกจที่แพ็คเกจ

4.3.2.1 ขั้นตอนการสร้างแพ็คเกจ

1. โหลดแพ็คเกจต้นฉบับย่อยตามข้อมูลในส่วนหัวของแพ็คเกจตำแหน่งซิสตริงฉบับที่ต้องการสร้างแพ็คเกจเก็บในหน่วยความจำ และในกรณีเป็นฉบับสุดท้ายจะเพิ่มช่องว่างไปเป็นจำนวนความยาวของซิสตริงหักออกหนึ่ง
2. จองหน่วยความจำเท่ากับจำนวนซิสตริงทั้งหมดในส่วนหัวของแพ็คเกจตำแหน่งซิสตริงบวกหนึ่งคูณกับขนาดของโหนด (รวมโหนดรากด้วย)
3. โหลดตำแหน่งของซิสตริงทั้งหมด ในแพ็คเกจตำแหน่งซิสตริง เก็บลงแต่ละโหนดในหน่วยความจำที่จองไว้
4. สร้างต้นไม้แพ็คเกจ ถ้าเกิดซิสตริงซ้ำกันให้เก็บเรียงต่อๆ กันไป เช่น โหนดที่ 2, 3 และ 4 ซ้ำกัน ดังรูปที่ 4.2



รูปที่ 4.2 ต้นไม้แพ็คเกจที่มีซิสตริงทั้งหมด 9 ซิสตริงและซิสตริงซ้ำๆ กัน 3 ซิสตริง

5. ใช้วิธีการท่องไปในต้นไม้ เพื่อหาแพ็คเกจที่แพ็คเกจ ซึ่งคือ โหนดภายนอกที่มีกรณีซ้ำขึ้นจากทางซ้ายไปทางขวา แล้วจัดเก็บลงแพ็คเกจที่แพ็คเกจ และถ้าโหนดใดมีซิสตริงที่ซ้ำๆ กันให้จัดเก็บลงแพ็คเกจที่แพ็คเกจ

4.3.2.2 ข้อกำหนดเบื้องต้นของการสร้างแพ็คเกจ

1. ต้องระบุชื่อแพ็คเกจตำแหน่งซิสตริง

2. ต้องมีหน่วยความจำไม่น้อยกว่า ขนาดของโหนดหนึ่งโหนด คูณกับ จำนวนซิสตริงทั้งหมด รวมกับ ขนาดของเพิ่มเติมฉบับย่อ ที่กำหนดไว้ในขั้นตอนการเตรียมเอกสาร เนื่องจากต้นไม้แฟตทั้งต้นกับเพิ่มเติมฉบับย่อต้องอยู่ในหน่วยความจำ โดยมีโครงสร้างข้อมูล ดังนี้

```
typedef struct {
    short skipb;
    long left, right, displace, duplicate;
} KEYBUF;
```

skipb	left	right	displace	duplicate
-------	------	-------	----------	-----------

รูปที่ 4.3 โครงสร้างโหนดต้นไม้แฟต

โดยที่	skipb	คือ บิตข้าม (Skip Bit)
	left	คือ ครรชนชี้ทางซ้าย (Left Pointer)
	right	คือ ครรชนชี้ทางขวา (Right Pointer)
	displace	คือ ตำแหน่งซิสตริง (Position of Sistring)
	duplicate	คือ ตำแหน่งซิสตริงซ้ำ (Duplicate Position of Sistring)

หมายเหตุ อาจลดค่า skipb ลงได้อีก 1 ไบต์ (ขนาดซิสตริงยาวได้ไม่เกิน 256 ไบต์) เพื่อประหยัดหน่วยความจำในการสร้างต้นไม้แฟต

3. มีทางเลือก ให้เลือก ดังนี้

3.1 สามารถสร้างแฟตอะเรย์ เก็บลงเพิ่มข้อมูล <input>.Arr โดย A คือ ตัวระบุว่าเป็นเพิ่มข้อมูลประเภทเก็บแฟตอะเรย์

3.2 สามารถสร้างข้อมูลสถิติ เก็บลงเพิ่มข้อมูล <input>.Crr โดย C คือ ตัวระบุว่าเป็นเพิ่มข้อมูลประเภทเก็บข้อมูลสถิติ

3.3 สามารถแสดงโหนดต้นไม้แฟตทั้งหมด เก็บลงเพิ่มข้อมูล <input>.Lrr โดย L คือ ตัวระบุว่าเป็นเพิ่มข้อมูลประเภทแสดงโหนดต้นไม้แฟตทั้งหมด

3.4 สามารถแสดงรูปภาพต้น ไม้แพ็ค เก็บลงเพิ่มข้อมูล <input>.Mr โดย M คือ ตัวระบุว่าเป็นเพิ่มข้อมูลประเภทแสดงรูปภาพต้น ไม้แพ็ค

3.5 สามารถเก็บต้น ไม้แพ็ค เก็บลงเพิ่มข้อมูล <input>.Tr โดย T คือ ตัวระบุว่าเป็นเพิ่มข้อมูลประเภทเก็บต้น ไม้แพ็ค

3.6 สามารถแสดงซิสตริงที่เรียงลำดับแล้ว เก็บลงเพิ่มข้อมูล <input>.Zr โดย Z คือ ตัวระบุว่าเป็นเพิ่มข้อมูลประเภทแสดงซิสตริงที่เรียงลำดับแล้ว

4. เพิ่มแพ็คอะเรย์ย่อย คือ <input>.Pr โดย P คือ ตัวระบุว่าเป็นเพิ่มข้อมูลประเภทเก็บแพ็คอะเรย์

5. เพิ่มแพ็คอะเรย์ซ้ำ คือ <input>.Sr โดย S คือ ตัวระบุว่าเป็นเพิ่มข้อมูลประเภทเก็บแพ็คอะเรย์ซ้ำ

4.3.2.3 โครงสร้างเพิ่มแพ็คอะเรย์ย่อย

ประกอบด้วย 2 ส่วน ดังนี้

1. ส่วนหัว (Header) มีองค์ประกอบหลักของข้อมูล คือ

ชื่อข้อมูล	รายละเอียดของข้อมูล
1. Large Text	ชื่อเพิ่มต้นฉบับ
2. Position of Sistring	ชื่อเพิ่มตำแหน่งซิสตริง
3. Duplicate PAT Array	ชื่อเพิ่มแพ็คอะเรย์ซ้ำ
4. Size Small Text	ขนาดเพิ่มต้นฉบับย่อย
5. Base Large Text	ตำแหน่งเริ่มต้นเพิ่มต้นฉบับ
6. Length Small Text	ความยาวเพิ่มต้นฉบับย่อย
7. Number of Sistrings	จำนวนซิสตริงทั้งหมด
8. Number of PAT Array	จำนวนแพ็คอะเรย์ทั้งหมด
9. Length Sistring	ความยาวซิสตริง
10. Flag Last Small Text	ตัวระบุเพิ่มต้นฉบับย่อยเป็นฉบับสุดท้าย

ตารางที่ 4.2 ส่วนหัวเพิ่มแพ็คอะเรย์

2. ส่วนข้อมูล (Detail) เก็บแพ็ทอะเรย์ คือ ตำแหน่งซิสตริงที่ใช้ในการอ้างอิงถึงซิสตริงในแฟ้มต้นฉบับที่เรียงลำดับแล้ว กับ ตำแหน่งแพ็ทอะเรย์ซ้ำของแฟ้มแพ็ทอะเรย์ซ้ำ เพื่อใช้เป็นคณรรชนีในการอ้างอิงถึงซิสตริงที่ซ้ำกันในแฟ้มต้นฉบับอีกทอดหนึ่ง ดังรูปที่ 4.4

```
typedef struct {
    long dpat;
    long ddup;
} PAT_ARRAY;
```

dpat	ddup
------	------

รูปที่ 4.4 โครงสร้างแพ็ทอะเรย์

โดยที่ dpat คือ แพ็ทอะเรย์
ddup คือ ตำแหน่งในแฟ้มแพ็ทอะเรย์ซ้ำ

4.3.2.4 โครงสร้างแฟ้มแพ็ทอะเรย์ซ้ำ

ประกอบด้วย 2 ส่วน ดังนี้

1. ส่วนหัว (Header) เหมือนแฟ้มแพ็ทอะเรย์
2. ส่วนข้อมูล (Detail) เก็บแพ็ทอะเรย์ที่ซ้ำๆ กันเรียงต่อๆ กัน โดย มีค่าศูนย์เป็นตัวบ่งบอกจุดสิ้นสุดของแพ็ทอะเรย์ที่ซ้ำๆ กัน

สำหรับขั้นตอนการสร้างแพ็ทอะเรย์ต้องอาศัยหน่วยความจำเพื่อเก็บคืนไม่แพ็ท เพื่อสร้างแพ็ทอะเรย์ โดยประสิทธิภาพของโปรแกรมขึ้นกับปัจจัย ดังต่อไปนี้

1. เวลาทั้งหมดในการอ่านแฟ้มตำแหน่งซิสตริงที่นำมาสร้างแพ็ทอะเรย์เก็บในหน่วยความจำ

2. เวลาทั้งหมดในการสร้างต้นไม้เพื่อ แต่เนื่องจากประมวลผลในหน่วยความจำจึงใกล้เคียงศูนย์

3. เวลาทั้งหมดในการท่องไปในต้นไม้เพื่อหาแพ็คเกจ แต่เนื่องจากประมวลผลในหน่วยความจำจึงใกล้เคียงศูนย์

4. เวลาทั้งหมดในการบันทึกแพ็คเกจและ แพ็คเกจ

4.3.3 การผสานแพ็คเกจ

หลังจากทำการสร้างแพ็คเกจย่อยของแต่ละตำแหน่งซิสตริง เสร็จแล้ว นำแพ็คเกจเหล่านั้นมา ผสานเข้าด้วยกัน เกิดเป็นแพ็คเกจฉบับใหม่ โดยทำการผสานทีละแพ็คเกจย่อย เข้ากับแพ็คเกจฉบับเก่า โดยอาศัยหน่วยความจำที่มีอยู่อย่างจำกัดให้เกิดประสิทธิภาพมากที่สุด ซึ่งมีหลักการ คือ เก็บแพ็คเกจย่อยทั้งหมดในหน่วยความจำ แล้วนำซิสตริงของแพ็คเกจฉบับเก่า ไปค้นหาตำแหน่งที่จะแทรกเข้าไป ต่อกับซิสตริงของแพ็คเกจย่อย เพื่อได้ซิสตริงของแพ็คเกจฉบับเก่า ผสานกับแพ็คเกจย่อย เก็บต่อเนื่องกันไปได้อย่างถูกต้อง แล้วจึงจัดเก็บลงแพ็คเกจฉบับใหม่ ส่วนที่เป็นซิสตริงซ้ำๆ กันจัดเก็บลงแพ็คเกจฉบับใหม่ และเก็บตารางการผสานแพ็คเกจลงเพิ่มตารางการผสานแพ็คเกจด้วย เพื่อใช้ในการอ้างอิงถึงซิสตริงในแพ็คเกจฉบับได้อย่างถูกต้อง

4.3.3.1 ขั้นตอนการผสานแพ็คเกจ

1. กรณียังไม่มีแพ็คเกจฉบับเก่า ให้สร้างจากแพ็คเกจย่อยที่ต้องการผสาน พร้อมทั้งสร้างตารางการผสานแพ็คเกจด้วย
2. โหลดแพ็คเกจย่อย กับ แพ็คเกจฉบับย่อย ที่ต้องการผสานเก็บในหน่วยความจำ
3. จองหน่วยความจำตัวนับค่า เพื่อเก็บจำนวนซิสตริงของแพ็คเกจฉบับเก่า ขนาดเท่ากับแพ็คเกจย่อยบวกหนึ่ง
4. โหลดตารางการผสานแพ็คเกจเก็บในหน่วยความจำ
5. โหลดแพ็คเกจฉบับย่อยของแพ็คเกจฉบับเก่าขึ้นแรกตามตารางการผสานแพ็คเกจเก็บในหน่วยความจำ
6. อ่านซิสตริงของตำแหน่งซิสตริงของแพ็คเกจฉบับย่อยของแพ็คเกจฉบับเก่าตามตารางการผสานแพ็คเกจ แล้วนำไปค้นหาแบบทวิภาคกับแพ็คเกจย่อยที่ต้องการผสานใน

หน่วยความจำ ถ้า *ผลเท่ากัน* ให้เพิ่มค่าตัวนับค่าที่เก็บ *ค่าเท่ากัน* ของ *ดรรชนีตัวที่ค้นพบขึ้นอีก* หนึ่ง แต่ถ้า *ผลน้อยกว่า* ให้เพิ่มค่าตัวนับค่าที่เก็บ *ค่าน้อยกว่า* ของ *ดรรชนีตัวที่ค้นไม่พบตัวสุดท้ายขึ้นอีกหนึ่ง* และถ้า *ผลมากกว่า* ให้เพิ่มค่าตัวนับค่าที่เก็บ *ค่าน้อยกว่า* ของ *ดรรชนีตัวที่ค้นไม่พบตัวสุดท้ายถัดไปขึ้นอีกหนึ่งแทน*

7. ปฏิบัติจนครบจำนวนซิสตริงทั้งหมดของแฟ้มตำแหน่งซิสตริง แล้วปฏิบัติกับแฟ้มตำแหน่งซิสตริงถัดไปตามตารางการผสานเพื่ออะเรย์จนครบหมด ซึ่งมีผลเท่ากับปฏิบัติจนครบแฟ้มอะเรย์ฉบับเก่า

8. ผสานแฟ้มอะเรย์ย่อยเข้ากับแฟ้มอะเรย์ฉบับเก่า โดยอาศัยค่าที่เก็บในตัวนับค่าเป็นดรรชนีชี้ นำ เช่น ค่าในตัวนับค่าชุดแรกเป็น 3, 1 แสดงว่าให้อ่านค่าในแฟ้มข้อมูลที่เก็บค่าแฟ้มอะเรย์ฉบับเก่ามา 3 ค่า แล้วจึงต่อด้วยค่าของแฟ้มอะเรย์ย่อย 1 ค่า โดยมีค่าของแฟ้มอะเรย์ฉบับเก่าที่ซ้ำกับค่าของแฟ้มอะเรย์นี้อีก 1 ค่า จากนั้นปฏิบัติกับตัวนับค่าชุดถัดไปจนหมด ซึ่งผลที่ได้ คือ แฟ้มข้อมูลแฟ้มอะเรย์ฉบับใหม่

9. เพิ่มตารางการผสานเพื่ออะเรย์ขึ้นอีกหนึ่งชุด โดยนำข้อมูลจากแฟ้มอะเรย์ย่อยที่ต้องการผสานเข้าไปเพิ่มต่อท้าย

หมายเหตุ การผสานแฟ้มอะเรย์ อาจเกิดกรณีการซ้ำกันของซิสตริงได้

4.3.3.2 ข้อกำหนดเบื้องต้นของการผสานแฟ้มอะเรย์

1. ต้องระบุชื่อแฟ้มแฟ้มอะเรย์ย่อยที่ต้องการผสาน
2. ต้องมีหน่วยความจำเพียงพอที่จะเก็บ แฟ้มอะเรย์ย่อยทั้งหมด แฟ้มต้นฉบับย่อย ตารางการผสานแฟ้มอะเรย์ และ อะเรย์ตัวนับค่า ซึ่งมีโครงสร้าง ดังต่อไปนี้

```
typedef struct {
    long less;
    long equal;
} COUNT_ARRAY;
```

less	equal
------	-------

รูปที่ 4.5 โครงสร้างตัวนับค่าของแฟ้มอะเรย์ย่อย

โดยที่ `less` คือ ส่วนที่เก็บ *ค่าน้อยกว่า* ในอะเรย์ตัวนับค่า
`equal` คือ ส่วนที่เก็บ *ค่าเท่ากัน* ในอะเรย์ตัวนับค่า

3. มีทางเลือก ให้เลือก ดังนี้

3.1 สามารถสร้างข้อมูลสถิติเก็บลงเพิ่มข้อมูล Merge.sta

3.2 สามารถแสดงซิสตริงที่เรียงลำดับแล้วเก็บลงเพิ่มข้อมูล Merge.sis

4. เพิ่มแพ็คเกจอะเรย์ฉบับใหม่ คือ Merge.pat

5. เพิ่มแพ็คเกจอะเรย์ซ้ำฉบับใหม่ คือ Merge.dup

6. เพิ่มตารางการผสมแพ็คเกจอะเรย์ คือ Merge.dir โดยมีโครงสร้าง ดังนี้

```
typedef struct {
    long mdisp;
    char fname[16];
} MERGE_DIR;
```

mdisp	fname
-------	-------

รูปที่ 4.6 โครงสร้างตารางการผสมแพ็คเกจอะเรย์

โดยที่ `mdisp` คือ ตำแหน่งสะสมของเพิ่มเติมฉบับย่อย
`fname` คือ ชื่อเพิ่มเติมตำแหน่งซิสตริง

4.3.3.3 โครงสร้างเพิ่มแพ็คเกจอะเรย์ฉบับใหม่

ประกอบด้วย 2 ส่วน ดังนี้

1. ส่วนหัว (Header) มีองค์ประกอบหลักของข้อมูล คือ

ชื่อข้อมูล	รายละเอียดของข้อมูล
1. Duplicate Large PAT Array	ชื่อเพิ่มแพ็ทอะเรย์ซ้ำฉบับใหม่
2. Table of Merge PAT Array	ชื่อเพิ่มตารางการผสมแพ็ทอะเรย์
3. Number of PAT Arrays	จำนวนแพ็ทอะเรย์สะสมทั้งหมด
4. Number of Duplicate PAT Arrays	จำนวนแพ็ทอะเรย์ซ้ำสะสมทั้งหมด
5. Number of Duplicate Merge PAT Arrays	จำนวนแพ็ทอะเรย์ซ้ำในตอนผสม
6. Length Sistring	ความยาวซิสตริง
7. Number of Small PAT Array	จำนวนชิ้นที่ทำการผสมแพ็ทอะเรย์

ตารางที่ 4.3 ส่วนหัวเพิ่มแพ็ทอะเรย์ฉบับใหม่

2. ส่วนข้อมูล (Detail) เก็บค่าแพ็ทอะเรย์ใหม่ คือ ตำแหน่งซิสตริงที่เรียงลำดับแล้ว บวกกับความยาวเพิ่มเติมฉบับย่อยแบบสะสมในตารางการผสมแพ็ทอะเรย์ กับค่าตำแหน่งแพ็ทอะเรย์ซ้ำใหม่ของเพิ่มแพ็ทอะเรย์ซ้ำฉบับใหม่

4.3.3.4 โครงสร้างเพิ่มแพ็ทอะเรย์ซ้ำฉบับใหม่

ประกอบด้วย 2 ส่วน ดังนี้

1. ส่วนหัว (Header) เหมือนเพิ่มแพ็ทอะเรย์ฉบับใหม่
2. ส่วนข้อมูล (Detail) เก็บค่าแพ็ทอะเรย์ใหม่ที่ซ้ำๆ กันเรียงต่อๆ กัน โดยมีค่าศูนย์เป็นตัวบ่งบอกจุดสิ้นสุดของแพ็ทอะเรย์ที่ซ้ำๆ กัน

4.3.3.5 โครงสร้างเพิ่มตารางการผสมแพ็ทอะเรย์

ชื่อข้อมูล	รายละเอียดของข้อมูล
1. Accumulate Length Small Text	ความยาวเพิ่มเติมฉบับย่อยแบบสะสม
2. Position of Sistring	ชื่อเพิ่มตำแหน่งซิสตริงย่อย

ตารางที่ 4.4 โครงสร้างเพิ่มตารางการผสมแพ็ทอะเรย์

สำหรับขั้นตอนการผสมผสานแฟ้มต่อเรย์ พบว่าประสิทธิภาพการผสมผสานแฟ้มต่อเรย์ ขึ้นกับปัจจัย ดังต่อไปนี้

1. เวลาทั้งหมดในการอ่านซิสตริงตามตารางการผสมผสานแฟ้มต่อเรย์
2. เวลาทั้งหมดในการค้นหาแบบทวิภาคบนแฟ้มต่อเรย์ย่อย แต่เนื่องจากประมวลผลในหน่วยความจำจึงใกล้เคียงศูนย์
3. เวลาทั้งหมดในการอ่านแฟ้มต่อเรย์ฉบับเก่าตามตารางการผสมผสานแฟ้มต่อเรย์
4. เวลาทั้งหมดในการบันทึกแฟ้มต่อเรย์ฉบับใหม่ที่ได้จากแฟ้มต่อเรย์ย่อยรวมกับแฟ้มต่อเรย์ฉบับเก่าตามตารางการผสมผสานแฟ้มต่อเรย์

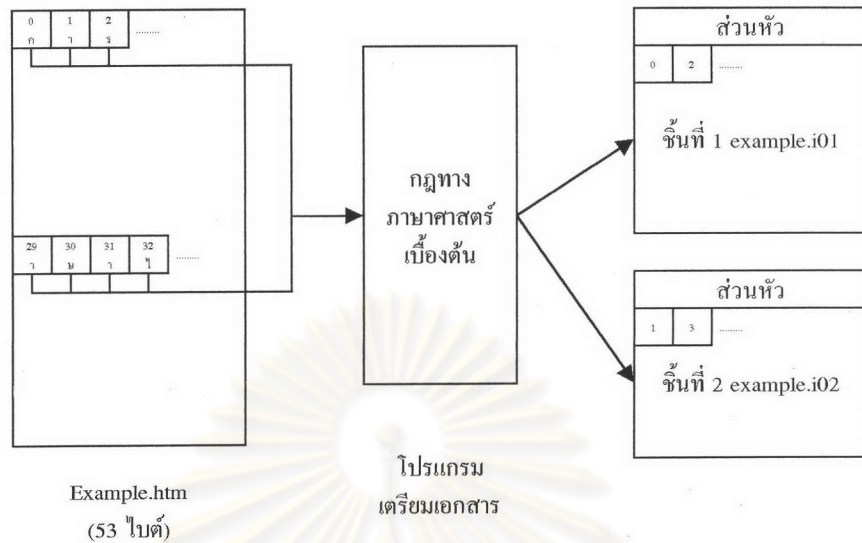
4.3.4 ตัวอย่างการสร้างแฟ้มต่อเรย์

สมมติว่า ต้องการสร้างแฟ้มต่อเรย์ ของแฟ้มต้นฉบับหนึ่ง โดยอาศัยวิธีการสร้างแฟ้มต่อเรย์ดังกล่าวข้างต้น ซึ่งมีข้อกำหนดให้ดังต่อไปนี้

1. พารามิเตอร์กำหนดชื่อแฟ้มต้นฉบับ คือ Example.htm ขนาด 53 ไบต์ ซึ่งมีข้อความในแฟ้ม คือ “ การพัฒนาระบบการค้นคืนข้อความภาษาไทยโดยใช้ต้นไม้แฟ้ม ”
2. พารามิเตอร์กำหนดขนาดแฟ้มต้นฉบับย่อย คือ 30 ไบต์
3. พารามิเตอร์กำหนดขนาดซิสตริง คือ 2 ไบต์

4.3.4.1 ตัวอย่างการเตรียมเอกสาร

การเตรียมเอกสาร คือ การจัดแบ่งแฟ้มต้นฉบับออกเป็นฉบับย่อย พร้อมทั้งค้นหาตำแหน่งซิสตริงตามกฎทางหลักภาษาศาสตร์เบื้องต้น แต่บันทึกเฉพาะตำแหน่งซิสตริงเท่านั้น เพื่อใช้เป็นครรชนีอ้างอิงถึงซิสตริงในแฟ้มต้นฉบับ ดังรูปที่ 4.7



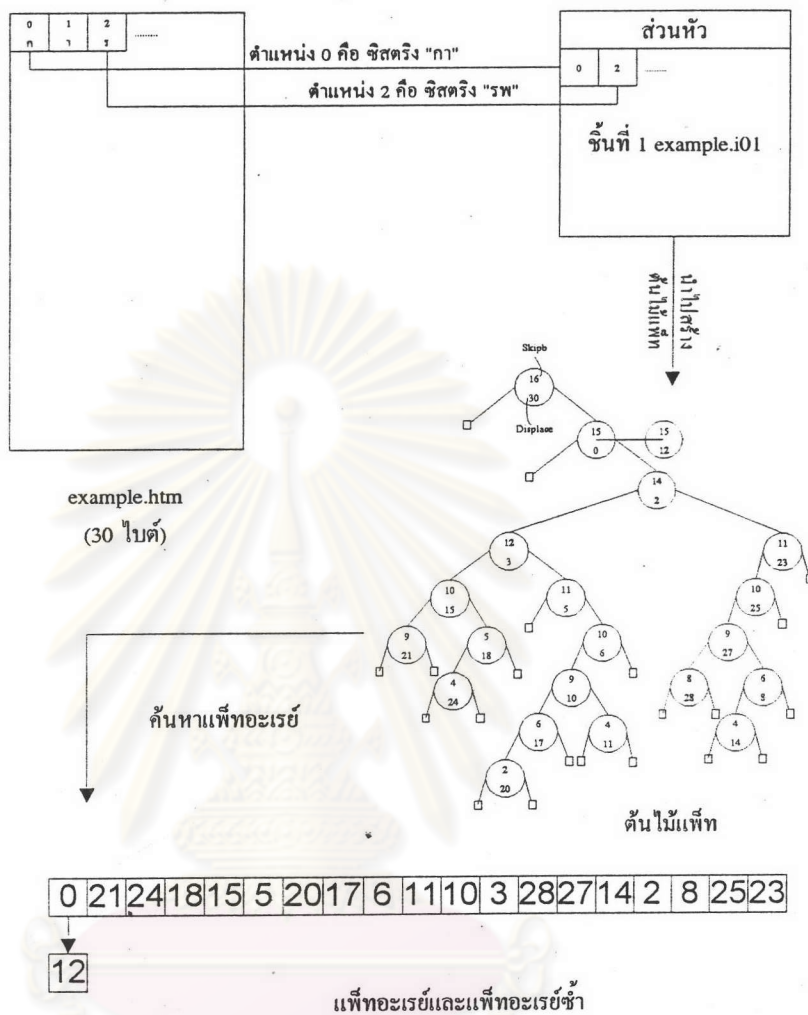
รูปที่ 4.7 ตัวอย่างการเตรียมเอกสาร

รูปที่ 4.7 เมื่อนำแฟ้มต้นฉบับ Example.htm ไปผ่านโปรแกรมเตรียมเอกสาร แล้วจะได้แฟ้มตำแหน่งซิสตริง คือ $\lceil 53 / 30 \rceil = 2$ ชั้น ชื่อ Example.i01 กับ Example.i02 โดยขนาดแต่ละแฟ้มขึ้นกับจำนวนซิสตริงของแต่ละชั้น (จำนวนซิสตริงทั้งหมด * ขนาดตำแหน่งซิสตริงในส่วนข้อมูล + ขนาดส่วนหัว) และแต่ละชั้น (แฟ้มตำแหน่งซิสตริง) เก็บตำแหน่งที่อ้างอิงถึงตำแหน่งในแฟ้มต้นฉบับ เช่น แฟ้มข้อมูล Example.i02 ข้อมูลชุดแรก คือ 1 หมายความว่า เป็นตำแหน่งที่ 1 ในแฟ้มข้อมูล Example.i02 ที่อ้างอิงถึงตำแหน่งที่ 1 บวกกับข้อมูลตำแหน่งเริ่มต้นแฟ้มต้นฉบับในส่วนหัวของแฟ้มข้อมูล Example.i02 (29) คือ $1 + 29 = 30$ คือ ตำแหน่งซิสตริง “ยา” ของแฟ้มต้นฉบับ เป็นต้น สำหรับตำแหน่งที่ไม่ถูกอ้างอิงผ่านแฟ้มตำแหน่งซิสตริงถึงเลขของแฟ้มต้นฉบับ แสดงว่าตำแหน่งนั้นไม่ใช่ซิสตริงตามกฎทางหลักภาษาศาสตร์เบื้องต้นดังกล่าวมาแล้ว

สำหรับแฟ้มตำแหน่งซิสตริง Example.i01 และ Example.i02 สามารถดูได้ โดยอาศัยทางเลือก ในข้อกำหนดเบื้องต้นของการเตรียมเอกสาร (ดูในภาคผนวก ก)

4.3.4.2 ตัวอย่างการสร้างแพ็คเกจ

การสร้างแพ็คเกจ คือ การจัดเรียงลำดับซิสตริง ตามรหัสข้อมูล โดยอาศัยต้นไม้แพ็คเกจเป็นตัวช่วย โดยแพ็คเกจ คือ โหนดภายนอกของต้นไม้แพ็คเกจที่เรียงจากทางซ้ายไปทางขวา ซึ่งโหนดภายนอกเก็บตำแหน่งซิสตริงที่ใช้อ้างอิงถึงซิสตริงในแฟ้มต้นฉบับ ดังรูปที่ 4.8



ศูนย์วิทยทรัพยากร
จุฬาลงกรณ์มหาวิทยาลัย

รูปที่ 4.8 ตัวอย่างการสร้างแฟ้มอะเรย์

รูปที่ 4.8 นำเพิ่มตำแหน่งซีสดริงชั้นแรก คือ Example.i01 ไปผ่านโปรแกรมการสร้างแฟ้มอะเรย์ เพื่อสร้างต้นไม้แฟ้ม เพื่อหาแฟ้มอะเรย์ ตามขั้นตอนการสร้างแฟ้มอะเรย์ ดังนี้

1. โหลดเพิ่มข้อมูล Example.htm ขนาด 30 ไบต์เก็บในหน่วยความจำ

2. อ่านตำแหน่งซิสตริงจากเพิ่มตำแหน่งซิสตริง Example.i01 เก็บใน `displace` ของแต่ละ โหนดในหน่วยความจำ

3. สร้างต้นไม้เพื่อติดตามค่า `displace` ในแต่ละโหนด

4. ท่องไปในต้นไม้เพื่อหาแพ็คเกจคือ 0, 12, 21, 24, 18, 15, 5, 20, 17, 6, 11, 10, 3, 28, 27, 14, 2, 8, 25, 23

เมื่อพิจารณาค่าของแพ็คเกจพบว่า ค่า 0 คือ ตำแหน่งของซิสตริง “กา” กับ ค่า 12 คือ ตำแหน่งของซิสตริง “กั” เช่นกัน ซึ่งเป็นซิสตริงที่มาก่อน ค่า 21 คือ ตำแหน่งของซิสตริง “ข” ดังนั้น การสร้างต้นไม้เพื่อหาแพ็คเกจคือ การจัดเรียงลำดับข้อมูลแบบดิจิทัลแบบหนึ่งที่มีประสิทธิภาพ

สำหรับเพิ่มตำแหน่งซิสตริงขึ้นที่สองคือ Example.i02 ทำในทำนองเดียวกับครั้งแรก ซึ่งได้แพ็คเกจคือ 10, 7, 12, 21, 4, 14, 19, 16, 5, 8, 1, 18, 6, 9, 3, 15

สำหรับเพิ่มแพ็คเกจ เพิ่มแพ็คเกจซ้ำ และ เพิ่มอื่นๆ ในขั้นตอนการสร้างแพ็คเกจสามารถดูได้ โดยอาศัยทางเลือก ในข้อกำหนดเบื้องต้นของการสร้างแพ็คเกจ (ดูในภาคผนวก ก)

4.3.4.3 ตัวอย่างการผสมผสานแพ็คเกจ

การผสมผสานแพ็คเกจ คือ การนำแพ็คเกจย่อยหลายๆ ฉบับที่ได้จากขั้นตอนการสร้างแพ็คเกจ มาผสมเข้าด้วยกันเป็นแพ็คเกจฉบับใหม่ โดยอาศัยหน่วยความจำที่มีอยู่อย่างจำกัดให้เกิดประสิทธิภาพมากที่สุด ดังรูปที่ 4.9

จุฬาลงกรณ์มหาวิทยาลัย



รูปที่ 4.9 ตัวอย่างการผสานแพ็ทอะเรย์

รูปที่ 4.9 เป็นการผสานแพ็ทอะเรย์ของขั้นแรก (Example.p01) กับขั้นที่สอง (Example.p02) เก็บลงเพิ่มแพ็ทอะเรย์ฉบับใหม่ (Merge.pat) ตามขั้นตอนการผสานแพ็ทอะเรย์ ดังนี้

1. โทลด์แพ็ทอะเรย์ย่อย (ขั้นที่ 2) กับ เพิ่มคั่นฉบับย่อย (ขั้นที่ 2) เก็บในหน่วยความจำ
2. อ่านตำแหน่งซีสตริงตามตารางการผสานแพ็ทอะเรย์ แล้วนำไปคั่นหาบนแพ็ทอะเรย์ย่อยแบบทวิภาคจนครบทุกตำแหน่ง เก็บค่านับจำนวนซีสตริงของแต่ละตัวลงในอะเรย์ตัวนับค่า ซึ่งอะเรย์ตัวนับค่า คือ อะเรย์เก็บจำนวนซีสตริงทั้งหมดที่มาก่อนและที่เท่ากันของซีสตริงของแพ็ทอะเรย์ย่อยในแต่ละตัว (Element) เช่น ค่า 6 ตัวแรกในอะเรย์ตัวนับค่า หมายความว่า มีจำนวนซีสตริงของแพ็ทอะเรย์ฉบับเก่าจำนวน 6 ซีสตริงที่มีค่าน้อยกว่า หรือเท่ากับก่อนค่า 10 ตัวแรก

ในแพ็คเกจย่อย (ชั้นที่ 2) คือ ชิสตริง 0:"กา", 12:"กา", 21:"ข", 24:"คว", 18:"ค", 15:"ค" มาก่อน ชิสตริง (40-30): "ข" หรือ 10:"ข" โดย ค่า 40 ต้องนำมาหักออกจากตำแหน่งเริ่มต้นเพิ่มเติมฉบับ คือ 30 เนื่องจากตอนผสมได้บวกค่า 30 ไว้ เป็นต้น

3. ผสานแพ็คเกจย่อย เข้ากับแพ็คเกจฉบับเก่า โดยเริ่มอ่านแพ็คเกจฉบับเก่าตามอะเรย์ตัวนับค่าทีละตัว แล้วตามด้วยแพ็คเกจย่อย (ชั้นที่ 2) ทีละตัว จนครบทุกตัว เก็บลงในแพ็คเกจฉบับใหม่

4. เพิ่มชื่อเพิ่มเติมตำแหน่งชิสตริง (fname) ที่นำมาผสม คือ Example.i02 กับค่าความยาวเพิ่มเติมฉบับย่อยแบบสะสม (mdisp) คือ ขนาดของเพิ่มเติมฉบับย่อยรวมกับค่าความยาวเพิ่มเติมฉบับย่อยแบบสะสมของชั้นสุดท้ายในตารางการผสมแพ็คเกจ กล่าวคือ $0+30$ เก็บลงในตารางการผสมแพ็คเกจ

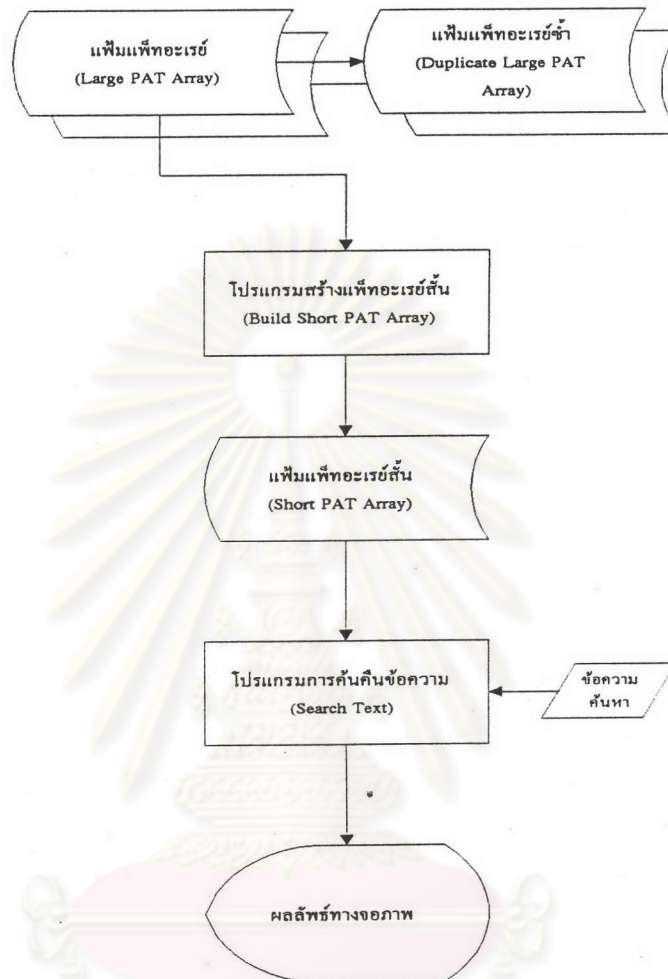
สำหรับเพิ่มแพ็คเกจ แพ็คเกจแพ็คเกจซ้ำ และ เพิ่มอื่นๆ ในขั้นตอนการผสมแพ็คเกจ สามารถดูได้ โดยอาศัยทางเลือก ในข้อกำหนดเบื้องต้นของการผสมแพ็คเกจ (ดูในภาคผนวก ก)

4.4 การคืนคืนข้อความ

การคืนคืนข้อความเป็นการค้นหาข้อมูลที่ได้จัดเตรียมไว้ล่วงหน้าในขั้นตอนการสร้างแพ็คเกจ โดยแบ่งออกเป็น 2 ขั้นตอน คือ

1. การสร้างแพ็คเกจต้น
2. การคืนคืนแบบเต็มหน้ากับแบบบูล

จุฬาลงกรณ์มหาวิทยาลัย



รูปที่ 4.10 การค้นคืนข้อความ

รูปที่ 4.10 แสดงขั้นตอนการค้นคืนข้อความ โดยอาศัยแพ็ทอะเรย์ที่เตรียมไว้ล่วงหน้าจากขั้นตอนการสร้างแพ็ทอะเรย์ แล้วนำมาสร้างแพ็ทอะเรย์สั้น เพื่อเป็นกรณีกลุ่มแพ็ทอะเรย์อีกชั้นหนึ่ง ซึ่งสามารถค้นหาข้อมูลได้ทั้งแบบเต็มหน้า หรือแบบบูล โดยมีหลักการ 2 ขั้นตอน คือ

1. สร้างแพ็ทอะเรย์สั้น เนื่องจากแพ็ทอะเรย์ทั้งหมด ไม่สามารถเก็บรวมไว้ในหน่วยความจำได้หมดในคราวเดียว ดังนั้นจึงสร้างแพ็ทอะเรย์สั้นกลุ่มอีกชั้นหนึ่ง เปรียบเสมือนเป็นกรณีของ

แพ็คเกจเก็บไว้ในหน่วยความจำทั้งหมดในคราวเดียวได้ เพื่อใช้ในการค้นหาอีกทอดหนึ่ง โดยอาศัยโปรแกรมสร้างแพ็คเกจสำเนาเป็นเครื่องมือในการสร้าง

2. ค้นคืนแบบเติมหน้า หรือแบบบูล หลังจากเตรียมแพ็คเกจสำเนาเสร็จแล้ว สามารถนำไปค้นคืนแบบเติมหน้า (Prefix) ได้อย่างรวดเร็ว เช่น ค้นหาเอกสารที่มีข้อความขึ้นต้นด้วย “กา” ได้ เช่น “การพัฒนา” หรือ “การค้นคืน” เป็นต้น และสามารถค้นคืนแบบบูล (Boolean) ได้ เช่น ค้นหาเอกสารที่มีข้อความ “กา” หรือ “ษา” ก็ได้ เป็นต้น

4.4.1 การสร้างแพ็คเกจสำเนา

ดังที่ได้กล่าวไว้แล้วเกี่ยวกับสาเหตุที่ต้องสร้างดัชนีกลุ่มแพ็คเกจอีกทอดหนึ่งเพื่อให้ดัชนี หรือ แพ็คเกจสำเนาอยู่ในหน่วยความจำทั้งหมดได้ โดยสามารถกำหนดขนาดทั้งหมดของแพ็คเกจสำเนา กับขนาดข้อมูลของแต่ละตัวของแพ็คเกจสำเนา เป็นพารามิเตอร์ตัวหนึ่งของระบบ แล้วระบบสร้างแพ็คเกจสำเนาตามพารามิเตอร์กำหนดขนาดทั้งหมดของแพ็คเกจสำเนา ซึ่งค่าของแต่ละตัวในแพ็คเกจสำเนา คือ ส่วนของข้อมูล (ซีสตริง) ของแพ็คเกจตัวสุดท้ายของกลุ่มแพ็คเกจ ที่มีความยาวตามพารามิเตอร์ขนาดข้อมูลของแต่ละตัวของแพ็คเกจสำเนา เก็บลงเพิ่มแพ็คเกจสำเนา

4.4.1.1 ขั้นตอนการสร้างแพ็คเกจสำเนา

1. กำหนดพารามิเตอร์ขนาดทั้งหมดของแพ็คเกจสำเนา (M) กับพารามิเตอร์ขนาดข้อมูลของแต่ละตัวของแพ็คเกจสำเนา (Is)

2. คำนวณหาจำนวนของแพ็คเกจสำเนา (r) จากสูตร $r = M / Is$

3. อ่านค่าของจำนวนแพ็คเกจทั้งหมด (n) จากส่วนหัวของแฟ้มแพ็คเกจ

4. คำนวณหาค่าจำนวนกลุ่มแพ็คเกจต่อหนึ่งแพ็คเกจสำเนา จากสูตร

$$b = (n / r) = (n * Is / M)$$

เพื่อต้องการแบ่งแพ็คเกจ ออกเป็นกลุ่มๆ โดยมีจำนวนกลุ่มเท่ากับจำนวนแพ็คเกจสำเนา

5. จองหน่วยความจำขนาดเท่ากับ พารามิเตอร์ขนาดทั้งหมดของแพ็คเกจสำเนา (M)

6. อ่านส่วนของข้อมูล (ชีสตรึง) ของแพ็คเกจตัวสุดท้ายของกลุ่มแพ็คเกจแพ็คเกจที่มีความยาวตามพารามิเตอร์ขนาดข้อมูลของแต่ละตัวของแพ็คเกจแพ็คเกจเก็บลงแต่ละตัวของแพ็คเกจแพ็คเกจในหน่วยความจำที่จองไว้

7. เก็บแพ็คเกจแพ็คเกจ ลงในแฟ้มแพ็คเกจแพ็คเกจ

หมายเหตุ ควรปรับค่าจำนวนกลุ่มแพ็คเกจแพ็คเกจต่อหนึ่งแพ็คเกจแพ็คเกจ (b) ให้เป็นเลขลงตัวเพื่อประสิทธิภาพในการค้นหาข้อมูลดีขึ้น

4.4.1.2 ข้อกำหนดเบื้องต้นของการสร้างแพ็คเกจแพ็คเกจ

1. ต้องระบุพารามิเตอร์ขนาดทั้งหมดของแพ็คเกจแพ็คเกจ (M) กับพารามิเตอร์ขนาดข้อมูลของแต่ละตัวของแพ็คเกจแพ็คเกจ (Is)
2. ต้องมีหน่วยความจำไม่น้อยกว่าพารามิเตอร์ขนาดทั้งหมดของแพ็คเกจแพ็คเกจ (M)
3. มีทางเลือกให้เลือก ดังนี้

3.1 สามารถสร้างข้อมูลสถิติกับ ข้อมูลในแพ็คเกจแพ็คเกจ เก็บลงแฟ้มข้อมูล Lschar.ls โดย Is คือ พารามิเตอร์ขนาดข้อมูลของแต่ละตัวของแพ็คเกจแพ็คเกจ

4. แฟ้มแพ็คเกจแพ็คเกจ คือ Spat.ls

4.4.1.3 โครงสร้างแฟ้มแพ็คเกจแพ็คเกจ

ประกอบด้วย 2 ส่วน ดังนี้

1. ส่วนหัว (Header) มีองค์ประกอบหลักของข้อมูล คือ

ชื่อข้อมูล	รายละเอียดของข้อมูล
1. Large PAT Array	ชื่อแฟ้มแพ็คเกจแพ็คเกจ
2. Duplicate Large PAT Array	ชื่อแฟ้มแพ็คเกจแพ็คเกจซ้ำ
3. Table of Merge PAT Array	ชื่อแฟ้มตารางการผสานแพ็คเกจแพ็คเกจ
4. Number of PAT Array	จำนวนแพ็คเกจแพ็คเกจสะสมทั้งหมด
5. Number of Duplicate PAT Array	จำนวนแพ็คเกจแพ็คเกจซ้ำสะสมทั้งหมด
6. Memory	จำนวนหน่วยความจำทั้งหมด
7. Number of Last PAT Block	ตัวระบุจำนวนแพ็คเกจแพ็คเกจในกลุ่มสุดท้าย

รูปที่ 4.11 เป็นการสร้างแพ็ทอะเรย์สั้น จากแพ็ทอะเรย์ที่สร้างขึ้นในขั้นตอนการสร้างแพ็ทอะเรย์ โดยมีขั้นตอนการสร้างแพ็ทอะเรย์สั้น ดังนี้

1. กำหนดค่า $M = 10$ ไบต์ กับ $ls = 2$ ไบต์
2. คำนวณค่า $r = M / ls = 10 / 2 = 5$ ตัว
3. อ่านค่า $n = 35$ ตัว
4. คำนวณค่า $b = n / r = (n * ls / M) = (35 * 2 / 10) = 70 / 10 = 7$ ตัว
5. จงหน่วยความจำขนาด $M = 10$ ไบต์
6. อ่านแพ็ทอะเรย์ตัวแทนกลุ่มแรก คือ แพ็ทอะเรย์ตัวสุดท้ายของกลุ่มแรก คือ ตัวที่ 7 มีค่าเป็น 5 นำไปตรวจสอบกับตารางการผสมแพ็ทอะเรย์ พบว่าอยู่ที่ค่าอ้างอิง 0 นำมาบวกกับ ค่า 5 ได้ $5+0 = 5$ คือ ซิสตริง “ฉน” แล้วทำงานกระทั่งครบทั้ง 5 ตัว คือ
 - 6.1 ตัวที่ 7 มีค่า 5 คือ ซิสตริง “ฉน”
 - 6.2 ตัวที่ 14 มีค่า 6 คือ ซิสตริง “นา”
 - 6.3 ตัวที่ 21 มีค่า 27 คือ ซิสตริง “มภ”
 - 6.4 ตัวที่ 28 มีค่า 25 คือ ซิสตริง “วา”
 - 6.5 ตัวที่ 35 มีค่า 45 คือ ซิสตริง “ไม”
7. เก็บแพ็ทอะเรย์สั้น ลงในแฟ้มแพ็ทอะเรย์สั้น (Spat.002)

สำหรับเพิ่มข้อมูลต่างๆ ในขั้นตอนการสร้างแพ็ทอะเรย์สั้น ดูได้ในภาคผนวก ก

จากรูปที่ 4.11 ขั้นตอน การสร้างแพ็ทอะเรย์สั้น อาศัย หน่วยความจำ 10 ไบต์ ตาม พารามิเตอร์ที่กำหนด โดยมีประสิทธิภาพ ขึ้นกับเวลาทั้งหมดในการอ่านซิสตริงตามแพ็ทอะเรย์สั้นเก็บในหน่วยความจำ และเวลาทั้งหมดในการบันทึกแพ็ทอะเรย์สั้น ลงในแฟ้มแพ็ทอะเรย์สั้น

จากขั้นตอนการสร้างแพ็ทอะเรย์กับแพ็ทอะเรย์สั้น พบว่า วัตถุประสงค์หลักของการสร้างแพ็ทอะเรย์ คือ การเรียงลำดับซิสตริงให้ถูกต้องและรวดเร็ว เพื่อเพิ่มประสิทธิภาพให้กับขบวนการค้นข้อความ ซึ่งการสร้างแพ็ทอะเรย์ คือ การเรียงลำดับซิสตริงของเอกสารที่รวดเร็ว และแพ็ทอะเรย์จะไม่เก็บข้อความของเอกสารเพื่อใช้เป็นดรรชนี แต่แพ็ทอะเรย์จะเก็บตำแหน่งที่อ้างอิงถึงซิสตริงในเอกสาร เพื่อใช้เป็นดรรชนีแทน โดยปกติแล้วปริมาณซิสตริงในเอกสาร มักมีปริมาณมากกว่าหน่วยความจำจะจัดเก็บไว้ได้หมดต่อการประมวลผลครั้งหนึ่งๆ ดังนั้นจึงจำเป็นต้อง

สร้างแพ็คเกจเรย์สั้น เพื่อเป็นกรณีคลุมแพ็คเกจเรย์อีกทอดหนึ่ง โดยแพ็คเกจเรย์สั้นจะเก็บส่วนของซิสตริงของกลุ่มแพ็คเกจเรย์ที่ถูกจัดแบ่งตามความเหมาะสม กับหน่วยความจำที่มีอยู่ และแพ็คเกจเรย์สั้นสามารถเก็บไว้ในหน่วยความจำได้ทั้งหมด ต่อการประมวลผลครั้งหนึ่งๆ ซึ่งแพ็คเกจเรย์สั้นจะมีกรณีอ้างอิงถึงแพ็คเกจเรย์ที่อยู่ในหน่วยความจำสำรองอีกทอดหนึ่ง เพื่อใช้ในการค้นหาซิสตริงที่ไม่อยู่ในแพ็คเกจเรย์สั้น แต่อาจอยู่ในแพ็คเกจเรย์ได้ ดังนั้น การกำหนดจำนวนแพ็คเกจเรย์สั้น กับการกำหนดขนาดของแพ็คเกจเรย์สั้น ให้เหมาะสมต่อการค้นหาข้อความ เป็นปัจจัยสำคัญที่ทำให้การค้นหาข้อความรวดเร็วขึ้น ทั้งนี้ขึ้นอยู่กับหน่วยความจำเป็นหลัก

4.4.2 การค้นคืนแบบเต็มหน้า

การค้นคืนแบบเต็มหน้า คือ การค้นหาข้อมูลที่มีส่วนของข้อมูล ส่วนหน้าเหมือนกัน (ข้อมูลร่วม) เช่น การค้นหาข้อความที่ขึ้นต้นด้วย “คอม” เช่น “คอมพิวเตอร์” หรือ “คอมมูนิกะชั่น” เป็นต้น ซึ่งการค้นคืนแบบเต็มหน้า โดยอาศัยแพ็คเกจเรย์มีความรวดเร็วมาก เนื่องจากแพ็คเกจเรย์ คือ ตำแหน่งซิสตริงของเอกสารที่เรียงลำดับแล้ว หรือคือ ซิสตริงของเอกสารที่เรียงลำดับแล้วต่อกัน ซิสตริงเหล่านี้มีข้อมูลร่วมที่เหมือนกันหรือใกล้เคียงกันเก็บเรียงต่อกัน ดังนั้นในการค้นคืนแบบเต็มหน้า จึงค้นหาซิสตริงที่มีข้อมูลร่วมตัวแรก กับ ซิสตริงที่มีข้อมูลร่วมตัวสุดท้าย ก็พอ (ไม่ต้องค้นหาหมดทุกตัว) และผลลัพธ์ของการค้นคืนแบบเต็มหน้า คือ ซิสตริงทุกตัว นับตั้งแต่ซิสตริงตัวแรกที่มีข้อมูลร่วมจนถึงซิสตริงตัวสุดท้ายที่มีข้อมูลร่วม

4.4.2.1 ขั้นตอนการค้นคืนแบบเต็มหน้า

1. กำหนดพารามิเตอร์ชื่อเพิ่มแพ็คเกจเรย์สั้น กับพารามิเตอร์ข้อความที่ต้องการค้นหา
2. อ่านแพ็คเกจเรย์สั้นเก็บไว้ในหน่วยความจำ (ทำหนเดียว)
3. ค้นหาแบบทวิภาคบนแพ็คเกจเรย์สั้น เพื่อหากกลุ่มแพ็คเกจเรย์สั้นตัวแรก กับตัวสุดท้าย แล้วจะได้ช่วงของกลุ่มแพ็คเกจเรย์สั้น (br) ตามความยาวของแพ็คเกจเรย์สั้นหนึ่งตัว (b)
4. ถ้าความยาวของข้อความค้นหา (lq) มากกว่า ความยาวของแพ็คเกจเรย์สั้นหนึ่งตัว (ls) ให้ค้นหาแบบทวิภาคบนกลุ่มแพ็คเกจเรย์สั้นตั้งแต่กลุ่มแพ็คเกจเรย์สั้นตัวแรกถึงกลุ่มแพ็คเกจเรย์สั้นตัวสุดท้าย (br) เพื่อหากกลุ่มแพ็คเกจเรย์สั้นตัวแรก กับตัวสุดท้าย แล้วจะได้ช่วงของกลุ่มแพ็คเกจเรย์สั้นช่วงใหม่ ตามความยาวของข้อความค้นหา (lq)
5. อ่านกลุ่มแพ็คเกจเรย์สั้นตัวแรกเก็บในหน่วยความจำ

6. ค้นหาแบบทวิภาคบนกลุ่มแพ็คเกจเริ่มต้นตัวแรกในหน่วยความจำ เพื่อหาชิสตริงตัวแรกตามข้อความค้นหา

7. อ่านกลุ่มแพ็คเกจเริ่มต้นตัวสุดท้ายเก็บในหน่วยความจำ

8. ค้นหาแบบทวิภาคบนกลุ่มแพ็คเกจเริ่มต้นตัวสุดท้ายในหน่วยความจำ เพื่อหาชิสตริงตัวสุดท้ายตามข้อความค้นหา

9. อ่านชิสตริงตั้งแต่ตัวแรกจนถึงตัวสุดท้ายตามที่ต้องการค้นหาแล้วเก็บบันทึกว่าชิสตริงผลลัพธ์อยู่ในเอกสารชิ้นใด โดยอาศัยอะเรย์เก็บจำนวนชิ้นเอกสาร ซึ่งใช้ค่า 1 บิต ต่อ เอกสารหนึ่งชิ้น

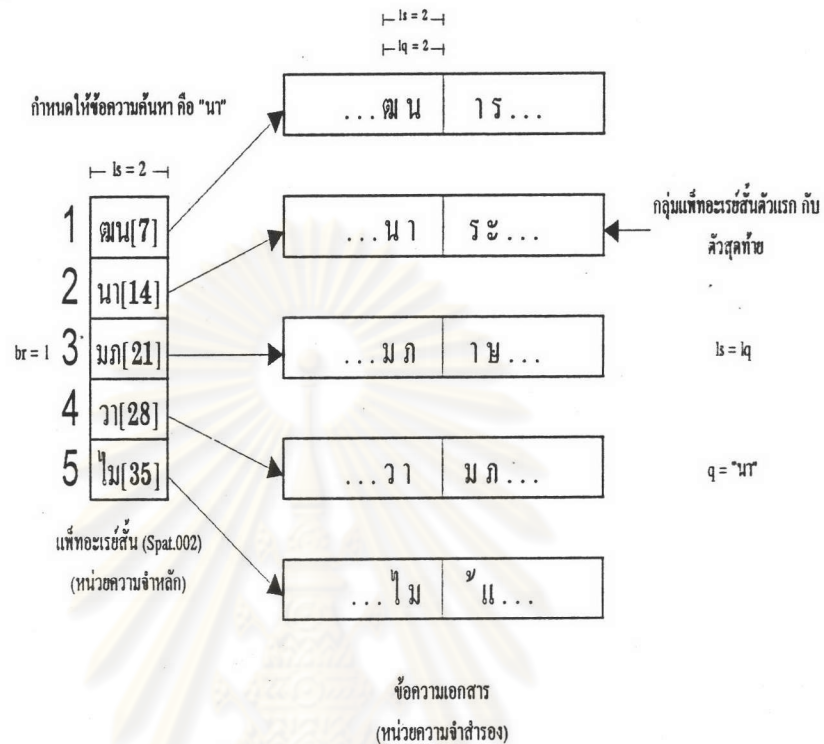
4.4.2.2 ข้อกำหนดเบื้องต้นของการค้นคืนแบบเต็มหน้า

1. ต้องระบุพารามิเตอร์ชื่อแฟ้มแพ็คเกจเริ่มต้น กับพารามิเตอร์ข้อความที่ต้องการค้นหา
2. ต้องมีหน่วยความจำไม่น้อยกว่าแพ็คเกจเริ่มต้น
3. มีทางเลือกให้เลือก ดังนี้

3.1 สามารถเก็บผลการค้นคืนแบบเต็มหน้า เก็บลงแฟ้มข้อมูล Search.log

4.4.2.3 ตัวอย่างการค้นคืนแบบเต็มหน้า

ศูนย์วิทยทรัพยากร
จุฬาลงกรณ์มหาวิทยาลัย



รูปที่ 4.12 ตัวอย่างการค้นหาแบบเดิมหน้า

รูปที่ 4.12 เป็นการค้นหาข้อความที่ขึ้นต้นด้วย "น" ตามขั้นตอนการค้นหาแบบเดิมหน้า
ดังนี้

1. กำหนดค่าพารามิเตอร์กำหนดชื่อแฟ้มแพ็ทอะเรย์สั้น คือ Spat.002 กับ ข้อความที่ต้องการค้นหาขึ้นต้นด้วย "น"
2. อ่าน Spat.002 เก็บในหน่วยความจำ (ทำหนเดียว)
3. ค้นหาแบบทวิภาคบน Spat.002 ได้ $br = 1$ และ $ls = 2$
4. ความยาวของ ls เท่ากับ lq ให้ข้ามขั้นตอนที่ 4 ไป
5. อ่านกลุ่มแพ็ทอะเรย์สั้นตัวแรกของกลุ่มแพ็ทอะเรย์ที่ค้นหาแบบทวิภาคได้ คือ กลุ่มที่ 2 เก็บในหน่วยความจำ
6. ค้นหาแบบทวิภาคบนกลุ่มแพ็ทอะเรย์สั้นตัวแรก ได้ชัดเจนคือ "น"

7. อ่านกลุ่มแพ็คเกจเรย์สั้นตัวสุดท้ายของกลุ่มแพ็คเกจเรย์ที่ค้นหาแบบทวิภาคได้ คือ กลุ่มที่ 2 เก็บในหน่วยความจำ

8. ค้นหาแบบทวิภาคบนกลุ่มแพ็คเกจเรย์สั้นตัวสุดท้าย ได้ซิสดริง คือ “นา” เหมือนกัน

9. คำตอบ คือ ซิสดริง “... นา | ระ ...” ซึ่งอยู่ในเอกสารชิ้นแรก คือ Example.i01 ของ Example.htm แล้วเก็บผลลัพธ์ลงในอะเรย์เก็บจำนวนชิ้นเอกสาร ได้ บิตแรกเป็น 1 คือ ซิสดริง “นา” อยู่ในเอกสารชิ้นแรก แล้วบันทึกผลการค้นหาลงในแฟ้มข้อมูล Search.log (ดูข้อมูลได้ใน ภาคผนวก ก)

จากรูปที่ 4.12 ขั้นตอนการค้นคืนแบบเต็มหน้า อาศัยหน่วยความจำเท่ากับ ขนาดแพ็คเกจเรย์สั้นรวมกับขนาดของกลุ่มแพ็คเกจเรย์หนึ่งกลุ่ม คือ $10+(8*3) = 34$ ไบต์ โดยมีประสิทธิภาพรวม ขึ้นกับปัจจัย ดังต่อไปนี้

1. เวลาทั้งหมดในการอ่านแพ็คเกจเรย์สั้นเก็บในหน่วยความจำ
2. เวลาทั้งหมดในการค้นหาแบบทวิภาคบนแพ็คเกจเรย์สั้น แต่เนื่องจากประมวลผลในหน่วยความจำจึงใกล้เคียงศูนย์
3. เวลาทั้งหมดในการอ่านกลุ่มแพ็คเกจเรย์สั้นกลุ่มแรกกับกลุ่มสุดท้ายเก็บในหน่วยความจำ
4. เวลาทั้งหมดในการค้นหาแบบทวิภาคของกลุ่มแพ็คเกจเรย์สั้นกลุ่มแรกกับกลุ่มสุดท้าย แต่เนื่องจากประมวลผลในหน่วยความจำจึงใกล้เคียงศูนย์

4.4.3 การค้นคืนแบบบูล

การค้นคืนแบบบูล คือ การค้นหาข้อมูลตามตรรกะ แอนด์ (And) กับ ออร์ (Or) เช่น การค้นหาข้อความที่ขึ้นต้นด้วย “คอม” แอนด์ “ระบบ” หมายความว่า ต้องการค้นหาเอกสารที่มีซิสดริง “คอม” และ “ระบบ” อยู่ในเอกสารชิ้นเดียวกัน เป็นต้น ซึ่งการค้นคืนแบบบูลโดยอาศัยแพ็คเกจเรย์จะช้ากว่าแบบเต็มหน้า เนื่องจาก จำเป็นต้องค้นหาซิสดริงของแต่ละข้อความทุกตัวก่อน (อย่างน้อยสองตัว) แล้วนำผลลัพธ์ของแต่ละซิสดริงมาประมวลผลตามตรรกะนั้น (And, Or) ที่หลัง แล้วจึงแสดงผลสุดท้ายอีกทอดหนึ่ง

4.4.3.1 ขั้นตอนการค้นคืนแบบบูล

1. กำหนดพารามิเตอร์ชื่อแฟ้มแฟ้มต่อระยะสั้น กับพารามิเตอร์ข้อความที่ต้องการค้นหา
2. อ่านแฟ้มต่อระยะสั้นเก็บไว้ในหน่วยความจำ (ทำคนเดียว)
3. ตัดแบ่งข้อความที่ต้องการค้นหา ออกเป็น ข้อความย่อยแต่ละตัว
4. นำข้อความย่อยตัวแรกไปค้นคืนแบบเต็มหน้า แล้วเก็บผลลัพธ์ลงในอะเรย์เก็บจำนวน
ชั้นเอกสาร
5. นำข้อความย่อยตัวถัดไปค้นคืนแบบเต็มหน้า แล้วเก็บผลลัพธ์ลงในอะเรย์เก็บจำนวนชั้น
เอกสารอีกตัวหนึ่ง แล้วนำอะเรย์ทั้งสองมาประมวลผลตามตรรกะของระหว่างชิสตริงทั้งสอง
6. นำข้อความย่อยตัวถัดไปประมวลผลตามขั้นตอนเดิม จนหมดทุกชิสตริง

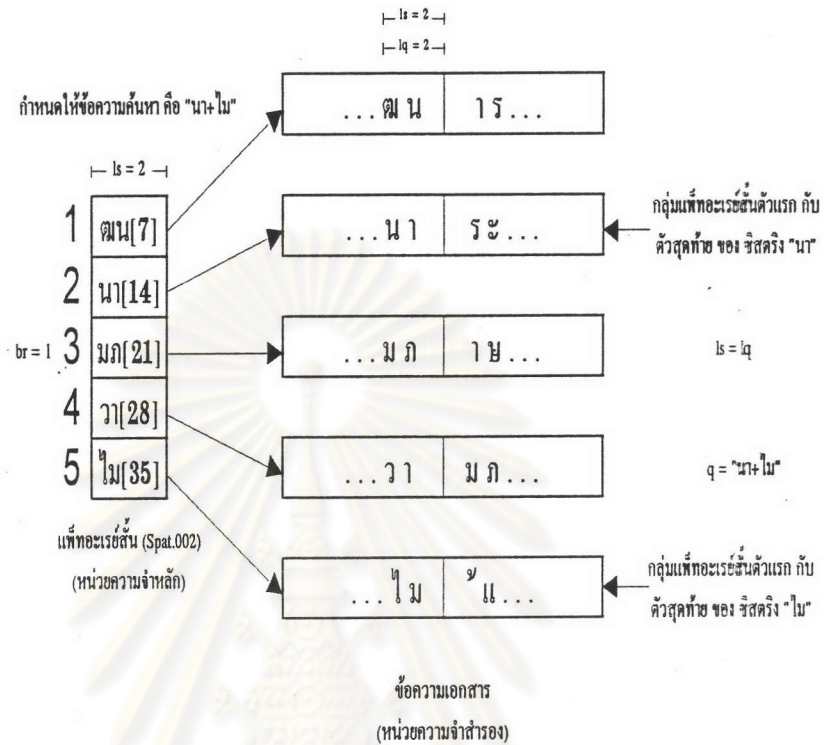
4.4.3.2 ข้อกำหนดเบื้องต้นของการค้นคืนแบบบูล

1. ต้องระบุพารามิเตอร์ชื่อแฟ้มแฟ้มต่อระยะสั้น กับพารามิเตอร์ข้อความที่ต้องการค้นหา
2. ต้องมีหน่วยความจำไม่น้อยกว่าแฟ้มต่อระยะสั้น
3. มีทางเลือกให้เลือก ดังนี้

3.1 สามารถเก็บผลการค้นคืนแบบเต็มหน้า เก็บลงแฟ้มข้อมูล Search.log

4.4.3.3 ตัวอย่างการค้นคืนแบบบูล

ศูนย์วิทยทรัพยากร
จุฬาลงกรณ์มหาวิทยาลัย



รูปที่ 4.13 ตัวอย่างการค้นหาแบบบูล

รูปที่ 4.13 เป็นการค้นหาข้อความที่ขึ้นต้นด้วย "นา" หรือ ขึ้นต้นด้วย "ไม" ตามขั้นตอนการค้นหาแบบบูล ดังนี้

1. กำหนดค่าพารามิเตอร์กำหนดชื่อแฟ้มแก๊ทอะเรย์สั้น คือ Spat.002 กับ ข้อความที่ต้องการค้นหาขึ้นต้นด้วย "นา+ไม"

2. อ่าน Spat.002 เก็บในหน่วยความจำ (ทำคนเดียว)

3. ตัดแบ่งข้อความที่ต้องการค้นหา ออกเป็น ข้อความย่อยแต่ละตัว ได้ "นา" กับ "ไม"

4. นำข้อความย่อยตัวแรก "นา" ไปค้นหาแบบเดิมหน้า แล้วเก็บผลลัพธ์ลงในอะเรย์เก็บจำนวนขึ้นเอกสาร ได้ บิตแรกเป็น 1 คือ ชิสตริงอยู่ในเอกสารชิ้นแรก

5. นำข้อความย่อยตัวถัดไป "ไม" ค้นหาแบบเดิมหน้า แล้วเก็บผลลัพธ์ลงในอะเรย์เก็บจำนวนขึ้นเอกสารอีกตัวหนึ่ง ได้ บิตสองเป็น 1 คือ ชิสตริงอยู่ในเอกสารชิ้นที่สอง แล้วนำอะเรย์ทั้งสองมาประมวลผลตามตรรกะของระหว่างชิสตริงทั้งสอง คือ ออร์ จะได้ผลลัพธ์สุดท้าย คือ บิต

แรก กับบิตที่สองเป็น 1 แสดงว่าข้อความ “นา” หรือ “ไม” มีอยู่ในเอกสารทั้งสองชิ้น คือ Example.i01 กับ Example.i02 แล้วบันทึกผลการค้นหาลงในแฟ้มข้อมูล Search.log (ดูข้อมูลได้ในภาคผนวก ก)

จากรูปที่ 4.13 ขั้นตอนการค้นหาแบบบูล อาศัยหน่วยความจำเท่ากับ ขนาดแฟ้มต่อเรย์สั้น รวมกับขนาดของกลุ่มแฟ้มต่อเรย์หนึ่งกลุ่ม คือ $10+(8*3) = 34$ ไบต์ โดยมีประสิทธิภาพรวมมากกว่าแบบเดิมหน้าเป็นสองเท่า



ศูนย์วิทยทรัพยากร
จุฬาลงกรณ์มหาวิทยาลัย