

บทที่ 4



การทดสอบ

โครงสร้างข้อมูลทรีอัดแน่นแบบจลน์

จากโครงสร้างข้อมูลทรีชนิดต่างๆ พบว่าทรีอัดแน่นเป็นโครงสร้างข้อมูลที่ดีที่สุดในแง่ของการประหยัดเนื้อที่ และอีกประการหนึ่งคือสามารถมองทรีอัดแน่นในลักษณะเช่นเดียวกับทรีแถวคู่ได้ คือมองเส้นเชื่อม (edge) ที่เก็บอยู่กับบัพ (node) เป็นฟังก์ชันการผ่าน (Transition function) ที่เรียกว่า ฟังก์ชัน goto ได้ กล่าวคือการผ่านในทรีแถวคู่มี $g(s_r, a)$ เป็น

$$t = \text{BASE}[r] + a ; \quad \text{CHECK}[t] = r$$

แต่เส้นเชื่อมของทรีอัดแน่นมี $g(s_r, a)$ เป็น

$$t = \text{BASE}[r] + a ; \quad \text{CHECK}[t] = a$$

โดย $\text{BASE}[r]$ ต้องไม่ถูกอ้างอิงด้วยเส้นเชื่อมอื่น

จากเหตุผลสองประการดังกล่าว

จึงเลือกที่จะทดสอบเปรียบเทียบกับ โครงสร้างทรีอัดแน่นเพื่อ

หาความเหมาะสมในการประยุกต์ใช้กับพจนานุกรมภาษาไทย

ทรีอัดแน่นที่สร้างขึ้นใหม่ได้ปรับปรุงให้เป็นทรีอัดแน่นแบบจลน์ คือสามารถเพิ่มลบคำได้ ทำได้โดยการเพิ่ม CHECK ซึ่งใช้เก็บสถานะที่ชี้มาโดยการผ่าน a ในลักษณะเดียวกับทรีแถวคู่ โดยแต่ละบัพในทรีอัดแน่นแบบจลน์มีโครงสร้างข้อมูลเป็นดังนี้

```
typedef struct {
    LINK next;
    LINK back;
    CHAR q;
    CHAR dummy;
```

```

} ITEM

typedef long      LINK;

typedef unsigned char  CHAR;

```

โดยเพิ่ม dummy เข้าไปอีก 1 ไบต์สำหรับการทำเขตสำหรับคำ (Word Boundary) และ back ก็คือ CHECK ในทรีแอดวู่นั้นเอง ซึ่งต้องมีไว้สำหรับการเพิ่มลบคีย์เข้าไปในทรีแอดวูแน่นอนแบบจลน์นี้ และจะกำจัดทิ้งไปเมื่อแปลงทรีแอดวูแน่นอนแบบจลน์นี้เป็นทรีแอดวูแน่นอน(แบบสถิตย์) เพราะไม่จำเป็นต้องใช้ back ในการสืบค้นคีย์ นอกจากนี้ยังได้ปรับปรุงเพิ่มส่วน TAIL ในทรีแอดวูแน่นอนแบบจลน์เหมือนกับของ ทรีแอดวู เพื่อใช้เก็บสตริงเดี่ยว (Single String) เพราะการเก็บเช่นนี้ประหยัดกว่าการเก็บเป็นบัฟและเส้นเชื่อม และยังสืบค้นได้เร็วกว่าเพราะเป็นเพียงการทำการเปรียบเทียบ (Pattern Matching)

ในส่วนของอัลกอริทึมที่ใช้เพิ่มลบคีย์ในทรีแอดวูแน่นอนแบบจลน์ก็เป็นทำนองเดียวกับทรีแอดวู แต่มีที่ต่างกันดังนี้

- 1) การหาที่ว่างสำหรับบัฟใหม่ใน BASE นอกจากการหาที่ว่างแบบพอดีก่อน (First-fit) สำหรับเส้นเชื่อมที่ชี้มาจากบัฟนั้นแล้ว ยังต้องตรวจสอบว่า BASE ที่อ้างด้วยเส้นเชื่อมนั้นว่างหรือไม่โดยใช้ BASE_USED flag
- 2) ไม่มีการเก็บ endmarker # ไว้ในทรีแอดวูแน่นอน มีแต่ IS_LAST flag ที่ใช้บอกจุดสิ้นสุดของคำ จึงแก้ไขอัลกอริทึมที่จัดการกับ endmarker เป็นจัดการกับ IS_LAST flag แทน
- 3) อัลกอริทึมการสืบค้นคีย์ตรวจสอบ $q[t] = a$ แทน $CHECK[t] = r$

ผลการทดสอบ

ทำการทดสอบเปรียบเทียบทรีแอดวูกับทรีแอดวูแน่นอนแบบจลน์

โดยใช้คำศัพท์ประเภทต่างๆดังนี้

- C คำหลักในภาษาซี
- ENGLISH คำศัพท์อังกฤษ
- THAI1K คำไทยใช้บ่อย 999 คำแรก
- THAI10K คำไทยใช้บ่อย 10,000 คำแรก
- THAI20K คำไทยใช้บ่อย 20,000 คำแรก

โดยดำเนินการทดสอบบนเครื่องไมโครคอมพิวเตอร์ที่มีหน่วยประมวลผลกลางอินเทล 486 ความเร็ว 33 เมกกะเฮิร์ตซ์ หน่วยความจำหลักขนาด 4 เมกกะไบต์ และจำลองหน่วยความจำหลัก 3 เม็ก

กะไบต์ใช้เป็นแรมดิสก์ (RAMDisk) เก็บเพิ่มข้อมูลของคำศัพท์ต่างๆ และเก็บเพิ่มข้อมูลของ โครงสร้างทรี ผลของการทดสอบเป็นดังตารางที่ 4.1

สรุปผล

1. เวลาที่ใช้ในการเพิ่มคำ

ยิ่งคำในโครงสร้างทรีทั้งสองยิ่งมาก เวลาในการเพิ่มคำก็ยิ่งมากขึ้นด้วย ทั้งนี้เป็นเพราะการหา BASE ที่ว่างเพื่อใส่สถานะและการผ่านแบบพอดีก่อน (First-fit) นั้นจะต้องเริ่มหาตั้งแต่ BASE แรกของทรี และพบว่าทรีอัดแน่นแบบจลน์ใช้เวลาในการเพิ่มคำมากกว่าทรีแถวคู่ เพราะจะต้องตรวจสอบว่า BASE ที่จะใช้นั้นถูกอ้างอิงโดยการผ่านอื่นหรือไม่ ซึ่งทำให้เสียเวลามากขึ้น

2. เวลาที่ใช้ในการสืบค้นคำ

เช่นเดียวกับการเพิ่มคำคือ คำในทรียิ่งมาก เวลาในการสืบค้นก็มากขึ้น ทั้งที่เวลาสืบค้นเฉลี่ยควรจะเท่าๆกัน ทั้งนี้เป็นเพราะทรีอยู่ในหน่วยความจำสำรองซึ่งมีโครงสร้างการเก็บข้อมูลในระบบปฏิบัติการดอสเป็นแบบรายการโยง (Linked List) โดยใช้ตารางจัดสรรเพิ่มข้อมูล (File Allocation Table — FAT) และพบว่าทรีอัดแน่นแบบจลน์และทรีแถวคู่ใช้เวลาสืบค้นเท่าๆกัน เนื่องจากการตรวจสอบการผ่านที่คล้ายกัน คือในทรีอัดแน่นแบบจลน์ใช้ $q[t]=a$ และในทรีแถวคู่ใช้ $CHECK[t]=r$

3. เนื้อที่หน่วยความจำที่ใช้

อัตราส่วนระหว่างเนื้อที่หน่วยความจำของทรีต่อเนื้อที่หน่วยความจำของเพิ่มของคำแปรผกผันกับจำนวนคำ กล่าวคือยิ่งคำมีจำนวนมาก ยิ่งมีการใช้ประสิทธิภาพของเนื้อที่หน่วยความจำของทรีดีขึ้น และพบว่าอัตราส่วนระหว่างเนื้อที่หน่วยความจำของทรีต่อเนื้อที่หน่วยความจำของเพิ่มของคำของทรีอัดแน่นแบบจลน์น้อยกว่าของทรีแถวคู่ เป็นเพราะเหตุผลเดียวกับเวลาที่ใช้ในการเพิ่มคำ คือการที่ต้องตรวจสอบ BASE ว่าว่างหรือไม่ แต่ถ้าหากแปลงทรีอัดแน่นแบบจลน์เป็นทรีอัดแน่น(แบบสถิติ) ที่ไม่ต้องมี back แล้วจะทำให้อัตราส่วนดังกล่าวนี้ลดลง

4. ความกระชับของเนื้อที่หน่วยความจำ

ความกระชับของเนื้อที่หน่วยความจำวัดจาก (เนื้อที่ที่ใช้ - เนื้อที่ที่ว่าง) / เนื้อที่ที่ใช้ พบว่าในส่วนของ TAIL มีความกระชับที่ใกล้เคียงกันไม่เปลี่ยนแปลงไปตามจำนวนของคำ แต่ส่วนของอาร์เรย์คู่และของทรีอัดแน่นแบบจลน์แปรตามจำนวนคำคือจำนวนคำยิ่งมาก ประสิทธิภาพการใช้หน่วยความจำยิ่งดี นอกจากนี้ยังพบว่าความกระชับของทรีอัดแน่นแบบจลน์ดีกว่าของทรีแถวคู่ เป็นเพราะการที่ทรีอัด

แน่นแบบจนใช้ IS_LAST เป็นการบอกการสิ้นสุดสตริงแทนการใช้ endmarker ในทรีแวก์ จึงไม่สิ้นเปลือง BASE ใหม่สำหรับการผ่าน (Transition) ของ endmarker ของทุกคำ

5. จากผลการทดสอบจึงกล่าวได้ว่าทรีแวก์เหมาะที่จะใช้เป็นโครงสร้างข้อมูลสำหรับฐานข้อมูลพจนานุกรมภาษาไทยที่จำเป็นต้องมีการเพิ่มลบคำ เพราะว่ามีประสิทธิภาพการใช้เนื้อที่หน่วยความจำได้ดี และมีเวลาที่ใช้ในการสืบค้นเพิ่มลบข้อมูลที่ค่อนข้างดี จึงยึดหยุ่นพอที่จะนำไปใช้เป็นโครงสร้างข้อมูลของพจนานุกรมสำหรับการแบ่งคำด้วย ถ้าหากผู้ใช้ต้องการที่จะทำการแก้ไขพจนานุกรมของตนเอง

	Double		Array	Trie		Dynamic		Packed	Trie	
	C	ENGLISH	THAI1K	THAI10K	THAI20K	C	ENGLISH	THAI1K	THAI10K	THAI20K
จำนวนคำ	32	19,816	999	10,000	20,000	32	19,816	999	10,000	20,000
ขนาดไฟล์คำศัพท์ (ไบต์)	230	195,914	5,951	77,274	168,904	230	195,914	5,951	77,274	168,904
ขนาดไฟล์ทรี (ไบต์)	2308	349,247	16,870	158,835	324,189	2,260	376,830	18,468	169,567	347,943
ทรี / คำศัพท์	10.03	1.78	2.83	2.06	1.92	9.82	1.92	3.10	2.19	2.06
จำนวนโหนด	271	35,893	1,784	15,896	31,837	212	31,476	2,598	13,790	27,845
จำนวนโหนดที่ใช้	50	35,705	1,535	15,431	31,245	163	31,425	2,055	13,454	27,552
ความกระชับของโหนด (%)	18.45	99.48	86.04	97.07	98.14	76.89	99.84	79.10	97.56	98.95
ขนาด TAIL (ไบต์)	140	62,103	2,598	31,667	69,493	140	62,070	1,587	31,667	69,493
ขนาด TAIL ที่ใช้ (ไบต์)	123	47,196	2,055	26,456	58,720	123	47,157	1,332	26,456	58,720
ความกระชับของ TAIL (%)	87.86	76.00	79.10	83.54	84.50	87.86	75.97	83.93	83.54	84.50
- สืบค้นเฉลี่ย (มิลลิวินาที)	1	9	3	6	11	1	9	3	6	11
- สืบค้นมากที่สุด (มิลลิวินาที)	6	54	11	22	40	6	54	11	22	40
- เพิ่มคำเฉลี่ย (มิลลิวินาที)	10	540	30	220	410	10	600	40	260	490
- เพิ่มคำมากที่สุด (มิลลิวินาที)	60	7,470	110	1,650	5,270	60	7,750	490	1,870	7,410

ตารางที่ 4.1 ผลการทดสอบกับคำประเภทต่างๆ