

นนทรีเสิร์จ : เสิร์จเอนจินสำหรับแคมปัสเน็ตเวิร์ค

NontriSearch : Search Engine For Campus Network

สุรศักดิ์ สงวนพงษ์ สุทธิพล วรางกูร
E-mail : { nguan, g4165234 }@ku.ac.th

ภาควิชาวิศวกรรมคอมพิวเตอร์
คณะวิศวกรรมศาสตร์ มหาวิทยาลัยเกษตรศาสตร์

บทคัดย่อ

นนทรีเสิร์จเป็นระบบเสิร์จเอนจินเพื่อใช้ค้นหาในเครือข่ายเว็บ การทำงานของนนทรีเสิร์จใช้หลักการจัดทำดัชนีและมีส่วนสอบถามคำที่ต้องการเพื่อสืบค้นคำในฐานข้อมูลในรูปแบบเดียวกับเว็บโรบอต ความสามารถในการสืบค้นของนนทรีเสิร์จทำได้ทั้งคำเต็มในภาษาไทยและภาษาอังกฤษ และสนับสนุนการสืบค้นโดยกำหนดเงื่อนไขบูลีน กลไกอ่านโฮมเพจเพื่อจัดทำดัชนีทำงานได้ตามข้อกำหนดของมาตรฐานข้อยกเว้นของโรบอต (Standard Robot Exclusion) ระบบได้รับการออกแบบให้มีขนาดเล็กเพื่อความเหมาะสมต่อการใช้งานเป็นเสิร์จเอนจินสำหรับแคมปัสเน็ตเวิร์คหรืออินทราเน็ตในองค์กร

1. บทนำ

อินเทอร์เน็ตในปัจจุบันมีข้อมูลข่าวสารอยู่มากมาย โดยเฉพาะข้อมูลไฮเปอร์เท็กซ์ที่มีให้บริการในระบบ เวบ เป็นที่นิยมใช้อย่างแพร่หลาย การค้นหาข้อมูลที่มีเนื้อหาที่ตรงกับความต้องการจากเว็บที่มีหลายล้านเพจในอินเทอร์เน็ตจึงเป็นเรื่องยาก หากไม่มีเครื่องมือช่วยค้นหาซึ่งรู้จักกันดีในชื่อของ เสิร์จเอนจิน (Search Engine)

บริการค้นหาเว็บเพจที่มีอยู่ในปัจจุบันเช่น AltaVista, Lycos, Yahoo, HotBot, InfoSeek [1] ฯลฯ ซึ่งมักเรียกโดยรวมทั้งหมดว่าเป็นเสิร์จเอนจินนั้น แท้ที่จริงแล้วมีบางแห่งมีกลไกทำงานแบบเสิร์จเอนจินอย่างแท้จริง และบางแห่งทำงานด้วยวิธีเก็บเว็บเพจไว้เป็น *ไครคทอรี* ซึ่งมีความแตกต่างกันโดยหลักพื้นฐานการสร้างดัชนีเว็บเพจ

ไครคทอรีมีการทำงานที่แตกต่างกับเสิร์จเอนจิน คือใช้การปรับเพิ่มข้อมูลโดยผู้ดูแลระบบเอง โดยไม่ได้ทำอย่างอัตโนมัติ เว็บไซต์ใดที่ต้องการมีรายชื่อในไครคทอรีต้องติดต่อไปยังผู้ดูแลไครคทอรี เพื่อให้ผู้ดูแลแจ้งแอดและจัดเก็บลงฐานข้อมูล ระบบการเก็บแบบไครคทอรีอาจให้ผลลัพธ์การค้นหาข้อมูลที่ตรงประเด็นมากกว่าเสิร์จเอนจินเพราะผ่านการแยกหมวดหมู่มาแล้ว ตัวอย่างของบริการค้นหาที่ใช้ระบบไครคทอรีได้แก่ Yahoo

สำหรับเสิร์จเอนจินโดยความหมายที่แท้จริงแล้วเป็นระบบซอฟต์แวร์ที่มีโปรแกรม *โรบอต* [2] ทำหน้าที่อ่านเว็บเพจจากไซต์ต่างๆ โดยอัตโนมัติ จากนั้นจึงนำเว็บเพจที่อ่านได้มาทำดัชนี เสิร์จเอนจินจะตรวจสอบลิงค์ในแต่ละหน้าของเว็บเพจเพื่อเข้าไปทำดัชนีของเว็บเพจนั้นต่อไปอีก ตัวอย่างของเสิร์จเอนจินนี้ได้แก่ HotBot หรือ Alta Vista หรือ Infoseek เป็นต้น

การค้นหาข้อมูลภาษาไทยในเว็บเพจยังคงเป็นความต้องการหลักประการหนึ่ง โดยเฉพาะอย่างยิ่งในองค์กรที่ใช้ระบบอินทราเน็ตภายในโดยมีบริการเว็บ และต้องการเสิร์จเอนจินที่ให้บริการเฉพาะการค้นหาข้อมูลภายในที่ตอบสนองได้ทั้งการสืบค้นข้อมูลภาษาไทยและอังกฤษ ในปัจจุบันเสิร์จเอนจินส่วนใหญ่ไม่สนับสนุนการทำดัชนีภาษาไทยเนื่องจากต้องอาศัยการทำงานเพิ่มเติมด้านการวิเคราะห์คำ บริการค้นหาภาษาไทยที่พบได้ในปัจจุบันอาจมีเพียง AltaVista [3] ที่ดัดแปลงให้สามารถทำดัชนีภาษาไทยได้ แต่ก็จัดอยู่ในกลุ่มของเสิร์จเอนจินขนาดใหญ่สำหรับเครือข่ายเว็บที่ครอบคลุมทั้งโลก

นนทรีเสิร์จเป็นเสิร์จเอนจินแบบโรบอต มีกลไกการทำงานแบบสร้างดัชนีคำทั้งในภาษาไทยและภาษาอังกฤษตามหลักการแบบ *ครรชนีผกผัน* (inverted index) [4] และจับเก็บคำโดยใช้วิธีแฮชชิง รูปแบบของครรชนีผกผันช่วยให้สามารถค้นหาได้อย่างรวดเร็วเพราะไม่มีกระบวนการค้นหาจากข้อมูลต้นฉบับใดๆ หากมีแต่เพียงการสืบค้นจากรางเก็บครรชนีเท่านั้น ตัวอย่างเช่นนนทรีเสิร์จสามารถใช้เวลาซีพียูโดยเฉลี่ยเพียง 2 วินาทีสำหรับสืบค้นคำในการค้นหาคำที่จากเว็บเพจประมาณ 11,000 หน้าที่มีความจุรวมราว 100 เมกะไบต์

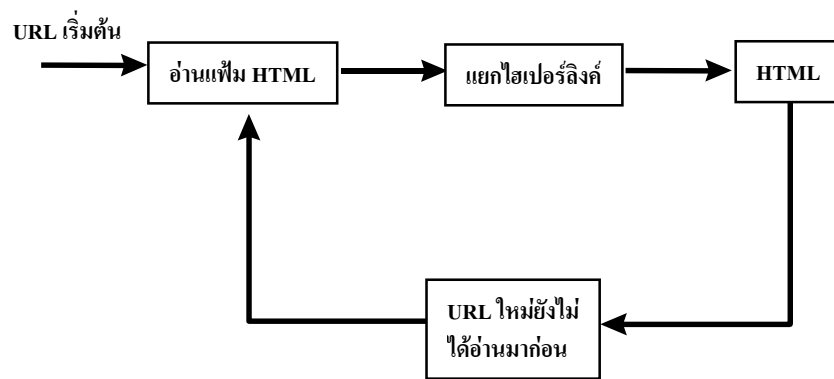
2. โครงสร้างของนนทรีเสิร์จ

โครงสร้างของนทรีเสร็จแบ่งออกเป็น 3 ส่วนหลัก คือ

- 1) นทรีบอด
- 2) อินเด็กเซอร์
- 3) นทรีเอ็นจิน

2.1 นทรีบอด

ทำหน้าที่เป็นโรบอตอ่านเว็บไซต์และจัดเก็บเพื่อส่งต่อไป อินเด็กเซอร์สร้างฐานข้อมูลต่อไป รูปที่ 1 แสดงโครงสร้างการทำงานของนทรีบอดซึ่งจะอ่านเว็บเพจและจัดเก็บโดยจำลองโครงสร้างของไคลเรทอริเช่นเดียวกับตำแหน่งที่เก็บของแฟ้มต้นฉบับในเว็บไซค์นั้น โรบอตจะสร้างคิวของเว็บเพจที่ต้องการทำคระชนนี้โดยกำหนดเพจเริ่มต้นให้ในคิว หากมองโครงสร้างของเว็บเพจเหมือนโครงสร้างต้นไม้การท่องเว็บเพจจะเป็นแบบแนวกว้าง (breadth first) โดยสามารถเลือกได้ว่าจะจำกัดระดับความลึกในการอ่านเว็บเพจหรือไม่



รูปที่ 1 การทำงานของนทรีบอด

2.2 อินเด็กเซอร์

ทำหน้าที่เป็นตัวสร้างคระชนนี้ ผลลัพธ์จากอินเด็กเซอร์ประกอบด้วยฐานข้อมูล 3 ฐานข้อมูล

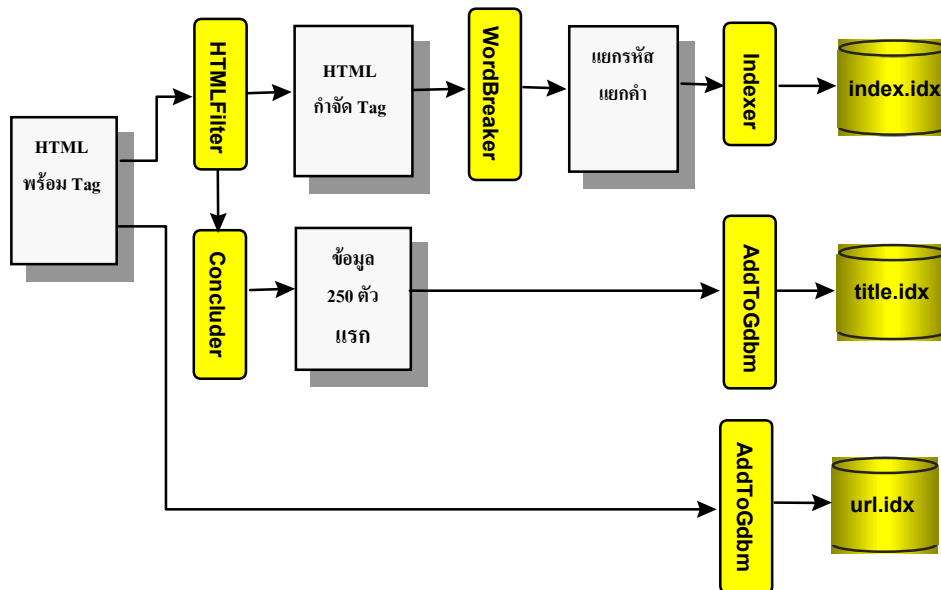
- index.idx เป็นฐานข้อมูลที่บรรจุคำศัพท์ที่พบทั้งหมด และหมายเลขกำกับ URL ที่เป็นตัวชี้ไปยังเพจที่มีคำศัพท์อยู่ โครงสร้างของฐานข้อมูลนี้แสดงได้ดังตารางที่ 1
- url.idx เป็นฐานข้อมูลที่เก็บหมายเลข URL และ ชื่อ URL โครงสร้างของฐานข้อมูลนี้แสดงได้ดังตารางที่ 2
- title.idx เป็นฐานข้อมูลที่เก็บหมายเลขกำกับ URL และคำบรรยายของเว็บเพจนั้น โดยเก็บหัวเรื่องของเว็บเพจจากแท็ก <TITLE> และอักขระ 250 ตัวแรกของเว็บเพจ ตารางที่ 3 แสดงโครงสร้างของฐานข้อมูลนี้

ตารางที่ 1 โครงสร้างภายในฐานข้อมูล index.idx

Key	Content
glance	1 2 3 4
gland	3 4 5 7 11
glare	5 7 8 9 13 120
glass	1 2 3 4 5
glaze	12 15 18 20 23 12
gleam	1 2 3 4 5 6 7 8 9
glean	2 4 5 6 8 10 12 14

โครงสร้างการทำงานของอินเด็กเซอร์ มีการองค์ประกอบย่อยดังรูปที่ 3 แต่ละส่วนทำหน้าที่ดังนี้

- 1) HtmlFilter ทำหน้าที่กรอง HTML Tag ออก และหากพบรหัสแฮชที่ใช้อักขระ จะดำเนินการแปลงไปเป็นรหัสแอสกี
- 2) Wordbreaker ทำหน้าที่แบ่งคำๆ โดยตัดคำที่ยาวที่สุดที่เป็นไปได้โดยอาศัยคลิกษณาริ
- 3) Indexer ทำหน้าที่นำคำศัพท์แต่ละคำไปเก็บในสร้างฐานข้อมูล index.idx
- 4) Concluder ทำหน้าที่หาคำบรรยายเว็บเพจโดยนำข้อความจาก <Title> xxxx </Title> มาเป็นหัวข้อ และนำอีก 250 ตัวอักษรถัดไปมาเป็นคำบรรยาย
- 5) AddToGDBM ทำหน้าที่เป็นตัวเพิ่มข้อมูลเข้าไปในฐานข้อมูล url.idx และ title.idx



รูปที่ 3 การทำงานของอินเด็กเซอร์

2.3 นนตรีเสิร์จ

เป็นส่วนที่รับคำศัพท์ที่ต้องการค้นหาจากผู้ใช้งานผ่านทางซีจีไอมาค้นหาในฐานข้อมูลและส่งผลการค้นหากลับไปแสดงบนหน้าจอ โดยนำคำศัพท์แต่ละคำที่ได้มาเปิดหาในฐานข้อมูล index.idx และนำผลลัพธ์ที่ได้จากการค้นหาของแต่ละครั้งมาเอกสารที่มีคำศัพท์ครบทุกคำ แล้วค่อยเปิดฐานข้อมูล url.idx และ title.idx เพื่อแสดงชื่อ URL และคำบรรยายออกหน้าจอ

3. ขีดความสามารถและสมรรถนะ

เทคนิคครรชนีผกผันที่ใช้ในนนตรีเสิร์จถึงแม้ว่าจะมีจุดอ่อนด้านการใช้เนื้อที่ขนาดใหญ่เพื่อสร้างครรชนี กล่าวคืออาจมีขนาดของครรชนีเป็น 50%-300% [6] ของขนาดข้อมูลต้นฉบับ และในกรณีของนนตรีเสิร์จจะมีขนาดแฟ้มครรชนีราว 185% ของขนาดข้อมูลต้นฉบับ และไม่ต้องเก็บข้อมูลต้นฉบับไว้อีกต่อไป หากแต่ในแง่ความรวดเร็วของการค้นหาแล้ววิธีนี้จะให้ผลลัพธ์ที่รวดเร็วตรงกับรูปแบบการให้บริการในระบบเครือข่ายที่ต้องตอบสนองการให้บริการแก่ผู้ใช้งานจำนวนมาก

การสร้างครรชนีสามารถเลือกให้ไม่ต้องสร้างครรชนีของคำในกลุ่ม "Stop word" เช่น a, an the, ก็, นี้ เป็นต้น เพื่อช่วยลดขนาดของแฟ้มครรชนี คำที่เหล่านี้สามารถเพิ่มหรือลดได้ตามความต้องการ

หากมีการสอบถามหาคำโดยใช้เงื่อนไขบูลีนเช่นเมื่อมีข้อความ *คำ1 AND คำ2* นนตรีเสิร์จจะค้นหารายการ URL ที่พบของแต่ละคำ จากนั้นจึงหาหมายเลข URL ร่วมกันเพื่อกำหนดเพจที่บรรจุคำทั้งสอง สัญลักษณ์ที่ใช้กำหนดฟังก์ชันบูลีนคือ

& (AND)

+ (OR)

- (NOT)

ตัวอย่างเช่น ‘ไทย & กรุงเทพฯ-อเมริกา’ หมายถึง ค้นหาเอกสารที่มีคำว่า ‘ไทย’ และ ‘กรุงเทพฯ’ แต่ไม่มีคำว่า ‘อเมริกา’ หากไม่มีการระบุเครื่องหมายจะใช้เครื่องหมาย & แทนให้โดยอัตโนมัติ สำหรับการค้นหาภาษาไทย ถ้าค้นหาเป็นประโยค นนทรีเสิร์จจะตัดประโยคที่ค้นหาออกเป็นคำๆ ให้โดยอัตโนมัติ

นนทรีเสิร์จสนับสนุนการทำงานตามข้อกำหนดของ Standard Robot Exclusion [5] ซึ่งเป็นมาตรฐานที่ช่วยให้ผู้ดูแลเซิร์ฟเวอร์กำหนดได้ว่าจะแนะนำให้โรบอตดำเนินการอ่านเพจเพื่อนำไปสร้างดัชนีหรือไม่

สำหรับสมรรถนะการค้นหาจากเว็บเพจจำนวนประมาณ 11,000 หน้า ขนาด 104 เมกะไบต์ ทำงานบนเครื่องเพนเทียม โพร 200 หน่วยความจำขนาด 32 เมกะไบต์ ฮาร์ดดิสก์แบบ IDE ทำงานภายใต้ระบบปฏิบัติการลินุกซ์ การค้นหาคำที่รับคำศัพท์ผ่านทางซีจีไอใช้เวลา 0.6-5 วินาที ตัวอย่างเช่นการค้นหาคำว่า *zurich* มีรวมกัน 10 แห่ง ใช้เวลาค้นหา 0.8 วินาที การค้นหาคำว่า *ประกาศ* มีรวมกัน 150 แห่ง ใช้เวลาค้นหา 3 วินาที แต่หากทดลองหาคำศัพท์โดยเปลี่ยนวิธีรับคำ แทนที่จะรับคำศัพท์จากซีจีไอ มาเป็นรับคำศัพท์จากไฟล์ของคำศัพท์ที่จะค้นหา และเปลี่ยนวิธีแสดงผลลัพธ์จากซีจีไอ มาเป็นหน้าจอ จะได้ผลการทดลองดังตารางที่ 4

ตารางที่ 4 ผลการทดลองเวลาที่ใช้ในการค้นหาคำศัพท์ (วินาที)

	User time	System time	Elapsed time
1 คำ	0.01	0.02	0.21
10 คำ	0.08	0.27	0.39
100 คำ	0.20	2.05	2.5
1000 คำ	3.9	19.56	38.36
10000 คำ	28.97	191.3	306.87

user time คือเวลาที่ใช้การทำงานของ โปรแกรม ซึ่งไม่นับเวลาที่ใช้ในการเรียก System call

system time คือเวลาที่ใช้ในการทำงานของ System call

elapsed time คือเวลาที่ใช้จริงทั้งหมด ซึ่งรวมเวลาที่ทำงานให้กับโปสเซอร์อื่นด้วย

4. พัฒนาการในอนาคต

การค้นหาคำในขณะนี้ยังคงจำกัดด้วยการสะกดคำที่ถูกต้อง พัฒนาการในขั้นต่อไปคือการตรวจสอบคำที่มีการสะกดใกล้เคียงและการสะกดแบบพ้องเสียง [7]

การทำดัชนีใหม่เมื่อเพจเกิดเปลี่ยนแปลงสามารถทำได้โดยไม่ต้องทำดัชนีซ้ำใหม่ทั้งหมด รูปแบบการทำดัชนีแบบเพิ่มขณะนี้กำลังอยู่ระหว่างดำเนินการปรับแต่งเพิ่มเข้าสู่นนทรีเสิร์จ นอกจากนี้ยังมีโครงการออกแบบระบบแจ้งการปรับปรุงเพจกลับมาซึ่งตัวสร้างดัชนีเพื่อให้สามารถดำเนินการปรับดัชนีใหม่ นอกเหนือไปจากการปรับปรุงตามกำหนดเวลาปกติ

คำเพื่อเลือกย่อยเพื่ออำนวยความสะดวกและให้ผู้ใช้เลือกกำหนดกลุ่มเป้าหมายของเว็บเพจเช่นวันเวลาย้อนหลัง การเลือกค้นเฉพาะเพจในโดเมนกำหนดให้ หรือการค้นหาเอกสารเฉพาะในยูสเน็ต เหล่านี้จะอยู่ในแผนการดำเนินการขั้นต่อไป

5. เอกสารอ้างอิง

- [1] Sullivan, D., *The Majors Search Engines*, <http://seachengineswatch.com/facts/major.html>.
- [2] Mauldin, M., *Lycos : Design Choices in the Internet Search Services*, IEEE Expert Vol. 15 No. 5 September 1997.
- [3] Altavista Search Network, http://altavista.digital.com/av/content/av_network.html.
- [4] Gonnet, G.H., Baeza-Yates, R., *Handbook of Algorithms and Data Structures In Pascal and C*, Addison-Wesley Publishing Company Inc., 1991.
- [5] Koster, M., *A Standard Robot Exclusion*, <http://info.webcrawler.com>.
- [6] Faloutsos, C., *Access Methods for Text*, ACM Computing Surveys, Vol. 17, 1985.
- [7] Ongroongruang, S., et.al., English to Thai Word Retrieval using Sound Index, In *the 2nd Symposium on Natural Language Processing*, Bangkok, Thailand., 1998