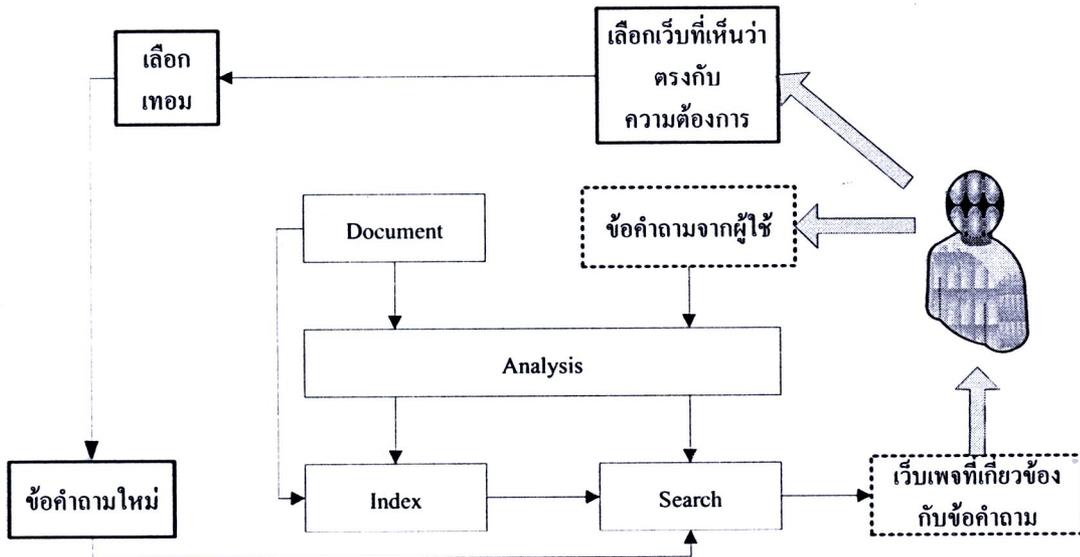


บทที่ 3 ผลการวิจัย

3.1 ส่วนประกอบของระบบ



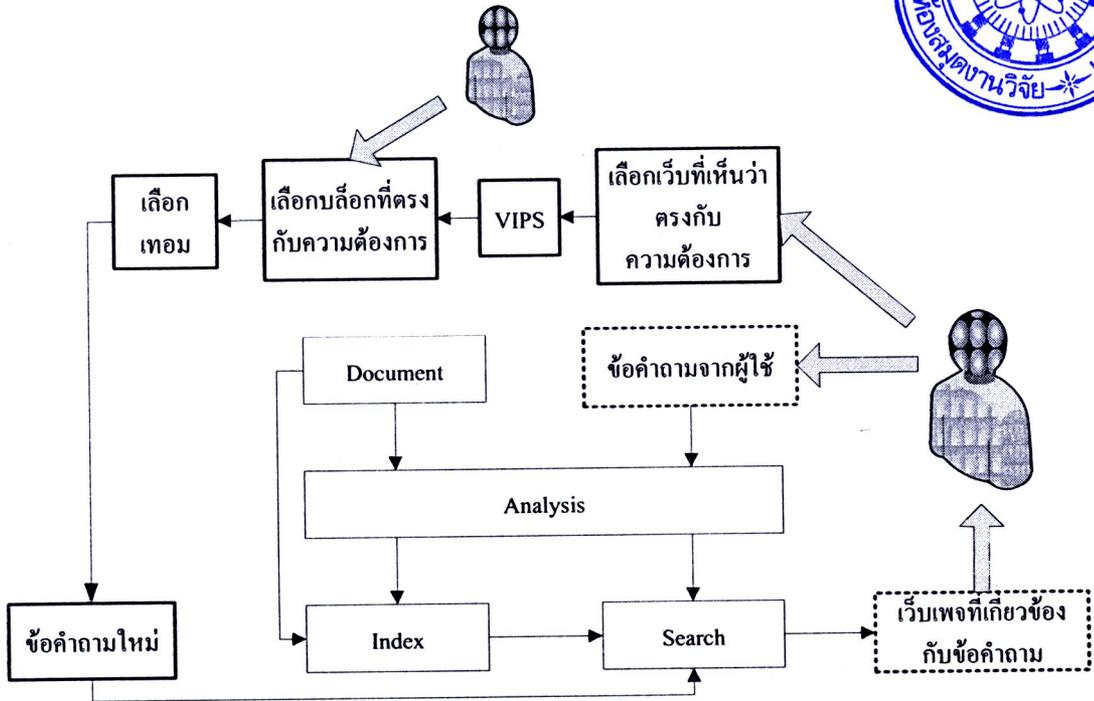
รูปที่ 3.1 แผนภาพทำงานของระบบค้นคืนสารสนเทศบนเว็บโดยใช้การค้นคืนย้อนกลับ

จากรูปที่ 3.1 การค้นคืนสารสนเทศบนเว็บโดยใช้การค้นคืนย้อนกลับจะเริ่มการทำงานจากการรับข้อความจากผู้ใช้ จากนั้นจะใช้ลูชันเพื่อหาเว็บเพจที่เกี่ยวข้องกับข้อความ เมื่อได้เว็บที่เกี่ยวข้องกับข้อความแล้ว ผู้ใช้ก็จะทำการเลือกเว็บเพจที่เกี่ยวข้องกับสิ่งที่ผู้ใช้ต้องการ และนำเว็บเพจที่เลือกไปหาข้อความใหม่ โดยการเพิ่มคำใหม่เข้าไปในข้อความเดิม ซึ่งคำใหม่ที่เพิ่มเข้าไปจะอยู่ภายในเว็บเพจที่ผู้ใช้เลือกเท่านั้น

เพื่อเป็นการเปรียบเทียบประสิทธิภาพการหาข้อความใหม่ จึงมีการออกแบบระบบค้นคืนสารสนเทศบนเว็บโดยใช้อัลกอริทึมวีไอพีเอสและการค้นคืนย้อนกลับดังรูปที่ 3.2 การทำงานจะเริ่มจากการใช้ลูชันเพื่อหาเว็บเพจที่เกี่ยวข้องกับข้อความ จากนั้นจะใช้ลูชันเพื่อหาเว็บเพจที่เกี่ยวข้องกับข้อความ เมื่อได้เว็บที่เกี่ยวข้องกับข้อความแล้วผู้ใช้ก็จะทำการเลือกเว็บเพจที่เกี่ยวข้องกับสิ่งที่ผู้ใช้ต้องการ และนำเว็บเพจที่เลือกไปทำการแบ่งออกเป็นบล็อกโดยใช้อัลกอริทึมวีไอพีเอส ผู้ใช้ก็จะทำการเลือกบล็อกที่เกี่ยวข้องกับสิ่งที่ผู้ใช้ต้องการอีกครั้งหนึ่ง ระบบก็จะหาข้อความใหม่โดยการเพิ่มคำใหม่เข้าไปในข้อความเดิม โดยคำใหม่ที่เพิ่มเข้าไปจะอยู่ในบล็อกที่ผู้ใช้เลือกเท่านั้น

เราสามารถแบ่งระบบค้นคืนออกเป็นระบบย่อยได้ 3 ส่วนด้วยกันคือ

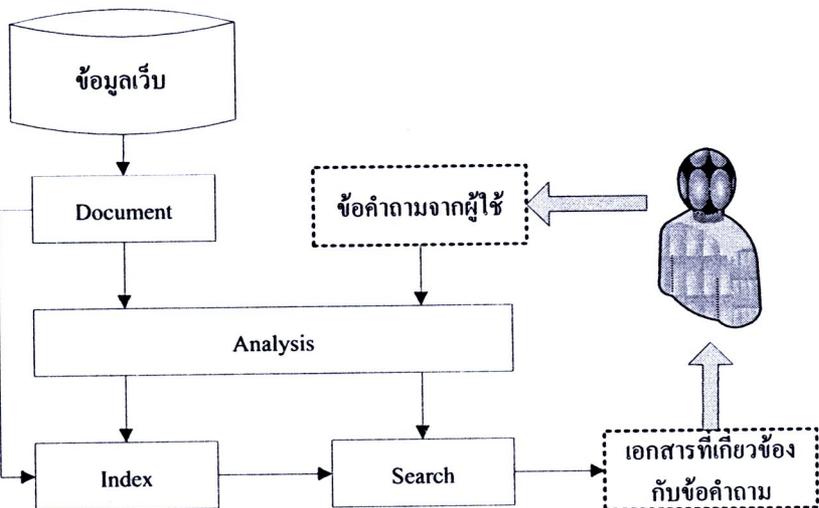
1. การค้นคืนสารสนเทศ
2. การแบ่งเว็บเพจออกเป็นบล็อก
3. การหาข้อความใหม่



รูปที่ 3.2 แผนภาพทำงานของระบบค้นคืนสารสนเทศบนเว็บโดยใช้ อัลกอริทึมวีไอพีเอสและการค้นคืนย้อนกลับ

3.1.1 การค้นคืนสารสนเทศ

การค้นคืนสารสนเทศในส่วนแรกนี้เราจะทำการศึกษาโดยใช้ลูซิน ซึ่งมีการกำหนดค่าการทำงานในลูซินดังนี้



รูปที่ 3.3 การค้นคืนสารสนเทศโดยลูซิน

1) ส่วนการจัดการเอกสาร (Document)

ดังนั้นในการค้นหาข้อมูลจะทำการค้นได้จาก 4 ส่วนนั้นคือ ไทเทิล (Title) เมตทา (Meta) บอดี (Body) และ ยูอาร์แอล (URL) โดย ฟิลด์ยูอาร์แอลและไทเทิลถือว่ามีค่าสำคัญเราจะไม่ทำการวิเคราะห์ค่าส่วนข้อมูลใน ฟิลด์บอดีที่มีจำนวนมากจะทำการบีบอัดข้อมูล เพื่อประหยัดพื้นที่ในการเก็บข้อมูล ส่วนการจัดการเอกสารที่ทำการเก็บเอกสารแต่ละหน้าเว็บเพจจะประกอบไปด้วยฟิลด์ดังต่อไปนี้

ตารางที่ 3.1 ฟิลด์ที่กำหนด

Field	ค่าของ field
title	ทำเป็นดัชนีแต่ไม่ต้องวิเคราะห์ค่า, ทำการเก็บข้อมูลดั้งเดิม
meta	ทำการวิเคราะห์ค่า ก่อนทำเป็นดัชนี, ทำการเก็บข้อมูลดั้งเดิม
body	ทำการวิเคราะห์ค่าก่อนทำเป็นดัชนี, ทำการเก็บข้อมูลดั้งเดิมและบีบอัดข้อมูล
url	ทำเป็นดัชนีแต่ไม่ต้องวิเคราะห์ค่า, ทำการเก็บข้อมูลดั้งเดิม

2) ส่วนการวิเคราะห์ค่า (Analysis)

วิเคราะห์ค่าที่ใช้ในการวิเคราะห์ข้อความและข้อความจะใช้ StandardAnalyzer ของลูชันคือ

- ทำการแบ่งเป็นคำตามช่องว่าง และอักขระพิเศษที่ไม่ใช่ตัวอักษร
- เปลี่ยนตัวอักษรภาษาอังกฤษตัวใหญ่เป็นตัวเล็กทั้งหมด
- ตัดคำที่มีอยู่ในเอกสารมากแต่ไม่ แสดงความหมายที่สำคัญ
- วิเคราะห์หลักไวยากรณ์ คือ สามารถรู้ ลักษณะอีเมล ไอพีแอดเดรส ตัวย่อและตัวอักษรที่ประกอบกับตัวเลข

ตัวอย่างการวิเคราะห์ค่าโดย StandardAnalyzer

The XY&Z Corporation – XYZ@example.com



[xy&z] [corporation] [xyz@example.com]

จากตัวอย่าง จะเห็นได้ว่าการใช้ StandardAnalyzer สามารถวิเคราะห์ข้อความโดยแบ่งเป็นคำตามช่องว่าง และอักขระพิเศษที่ไม่ใช่ตัวอักษรเปลี่ยนตัวอักษรภาษาอังกฤษตัวใหญ่เป็นตัวเล็กทั้งหมด และจะไม่ทำการแบ่งข้อความที่เป็นอีเมล

3) ส่วนดัชนี (Index)

เก็บข้อมูลดัชนีแบบเพิ่มข้อมูลผกผันตามรูปแบบของลูชัน

4) ส่วนค้นคืน (Search)

การค้นคืนข้อมูลโดยหาเอกสารที่เกี่ยวข้องกับข้อความของผู้ใช้จากการเปรียบเทียบกับ ดัชนี

ที่มี

- คิวรีเฟสเซอร์ ที่ใช้ทำการวิเคราะห์ข้อความใช้รูปแบบเดียวกับการวิเคราะห์เอกสาร (StandardAnalyzer)
- ใช้ เทอมคิวรี (TermQuery) ในการค้นหาเอกสารที่เกี่ยวข้องกับข้อความ
- การคำนวณค่าสกออร์เพื่อนำไปจัดลำดับการแสดงผลจะกำหนดให้ ทุกๆคำในข้อความ พิลด์ทุกฟิลด์และเอกสารทุกเอกสารมีค่าน้ำหนักเท่ากัน

3.1.2 การแบ่งเว็บเพจออกเป็นบล็อก

การแบ่งเว็บเพจออกเป็นบล็อกมีขั้นตอนดังนี้ คือ

1. หลังจากที่ได้ทำการค้นคืนเอกสารในครั้งแรกแล้ว ก็จะได้รายชื่อของเว็บเพจที่ระบบค้นคืนกลับมาให้ผู้ใช้ จากนั้นผู้ใช้จะทำการเลือกเว็บเพจที่คิดว่ามีความเกี่ยวข้องกับข้อความ โดยเว็บเพจที่ผู้ใช้เลือกนั้น จะนำมาทำการแบ่งเป็นบล็อกด้วย วิโอพีเอสอัลกอริทึม ซึ่งมีขั้นตอนการทำงานเริ่มจาก การตัดออกเป็นบล็อก โดยจะเป็นการนำดอมนทรีของเว็บเพจมาพิจารณาโดยใช้กฎในการแบ่งออกเป็นบล็อก จากกฎจะได้ค่า ดีโอซีของแต่ละบล็อก และเราจะกำหนดค่าพีดีโอซี ให้มีค่าเท่ากับ 4 เสมอ เพื่อกำหนดโครงสร้างเนื้อหา
2. การหาตัวแบ่ง จะเป็นการหาตัวแบ่งระหว่างแต่ละบล็อกที่ได้จากข้อหนึ่งเพื่อแยกแต่ละบล็อกให้ออกจากกัน และทำการให้ค่าน้ำหนักแต่ละตัวแบ่ง
3. การสร้างโครงสร้างเนื้อหา ในขั้นตอนนี้จะทำการรวมบล็อกโดยดูจากค่าน้ำหนักของตัวแบ่งแต่ละตัว โดยตัวแบ่งที่มีค่าน้ำหนักน้อยจะสามารถทำการรวมกันเป็นบล็อก
4. ผู้ใช้ทำการเลือกบล็อกที่เห็นว่าตรงตามความต้องการ โดยบล็อกที่ถูกเลือกจะนำไปทำการเพิ่มข้อความใหม่

3.1.3 การหาข้อความใหม่

การหาข้อความใหม่มีขั้นตอนดังนี้ คือ

1. นำข้อมูลจากเว็บที่ถูกเลือก หรือข้อมูลจากบล็อกที่ถูกเลือกที่มาจากการใช้วิโอพีเอสอัลกอริทึมมาทำการวิเคราะห์หาคำใหม่ โดยเทอมที่จะเพิ่มเข้าไปในนั้นใช้หลักการของเว็ทเตอร์โมเดล สามารถคำนวณได้จากสูตรซอสเตอร์ คือ

$$s(t_k) = \frac{1}{n_k} df_k |R^k| \quad (3.1)$$

โดยที่ n_k คือค่านอยซ์เมทซ์วีร์ของเทอมที่ k
 df_k คือความถี่ของเทอมที่ k ที่ปรากฏในหน้าเว็บที่ถูกเลือก หรือ ความถี่ของเทอมที่ k ที่ปรากฏในบล็อกที่ถูกเลือก
 R^k คือจำนวนเอกสารที่ค้นคืนได้ที่มีเทอมที่ k ปรากฏอยู่

คำนวณหาค่านอยซ์เมทซ์วีร์ได้จากสูตร

$$n_k = \sum_{i=1}^N N \times \frac{tf_{ik}}{f_k} \times \log(tf_{ik} \cdot f_k) \quad (3.2)$$

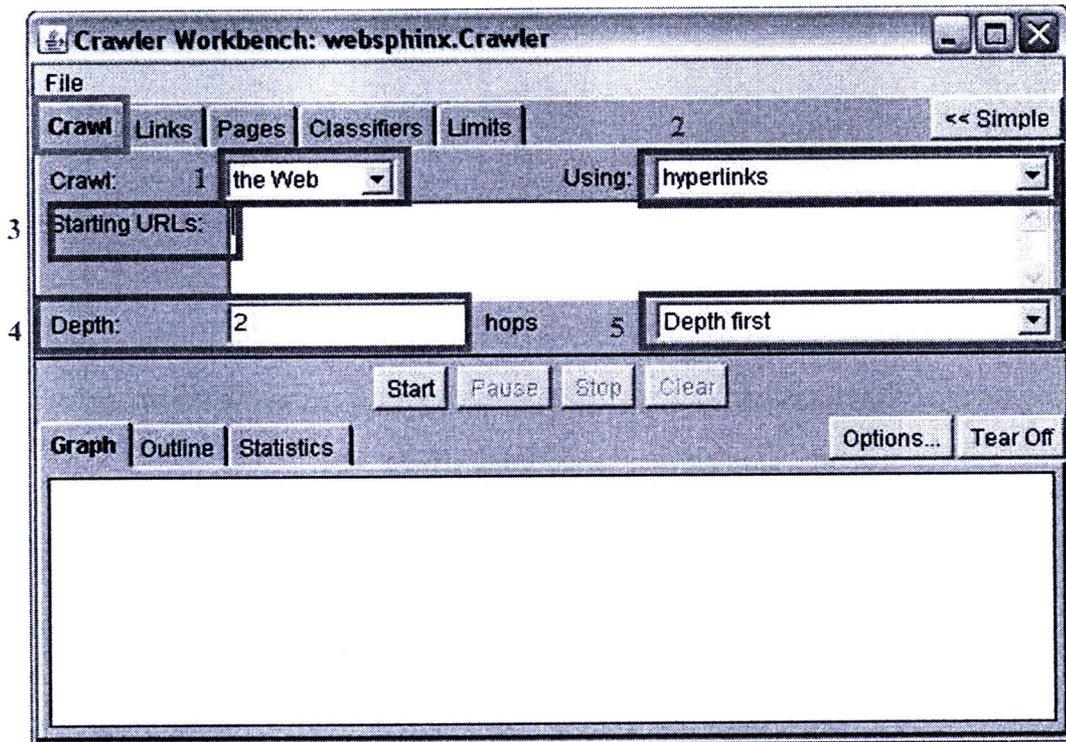
โดยที่ N คือจำนวนเทอมทั้งหมดภายในหน้าเว็บที่ถูกเลือก หรือจำนวนเทอมทั้งหมดภายในบล็อกที่ถูกเลือก

- f_{ik} คือความถี่ของเทอมที่ k ที่ปรากฏในเอกสารที่ i
- f_k คือความถี่ของเทอมที่ k ที่ปรากฏในหน้าเว็บที่ถูกเลือก หรือความถี่ของเทอมที่ k ที่ปรากฏในบล็อกที่ถูกเลือก

2. เมื่อได้ค่าซอสเตอร์ของแต่ละคำออกมาแล้ว นำมาเรียงลำดับ แล้วเลือกเพียงลำดับสูงสุดเพียง 3 คำ จากนั้นนำมาเพิ่มในข้อคำถาม ได้เป็นข้อคำถามใหม่ขึ้นมา

3.2 การเตรียมข้อมูล

การเตรียมข้อมูลจะใช้ ครอบเลอร์เว็บสฟิงซ์ (Website-Specific Processors for HTML Information Extraction: WebSPHINX) ซึ่งเป็นโอเพนซอร์สโปรแกรมครอบเลอร์มาทำการเก็บข้อมูลเว็บเพจ

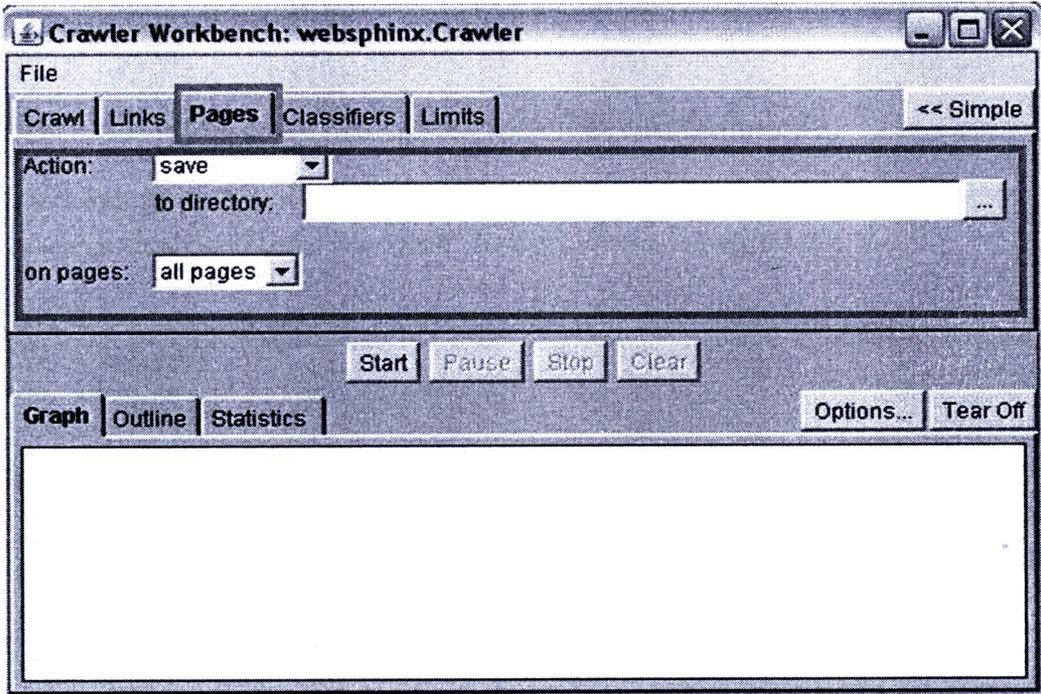


รูปที่ 3.4 แสดงโปรแกรมเว็บสฟิงซ์

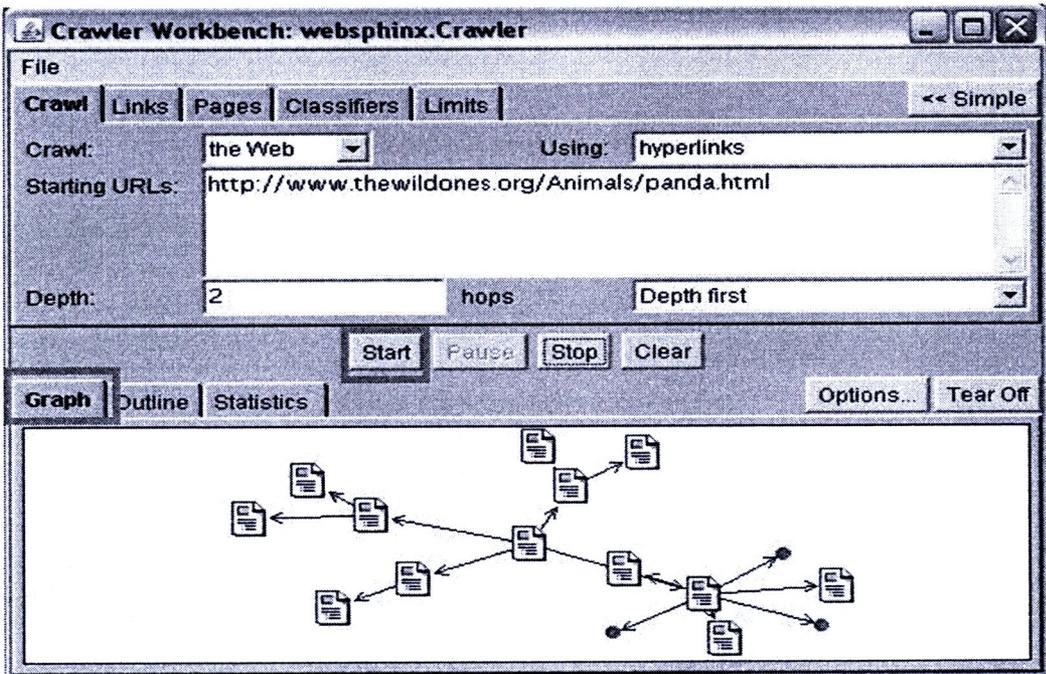
เมื่อเปิดโปรแกรมเว็บสฟิงซ์จะได้หน้าจอตั้งรูปที่ 3.4 โดยในส่วนของแท็บ crawl นี้มีการกำหนดค่าดังต่อไปนี้

1. เลือกรูปแบบการครอบเลอร์ ซึ่งจะเลือกแบบ the Web เพื่อทำการเก็บข้อมูลของหน้าเพจทั้งหมด
2. เก็บข้อมูลของลิงค์ทุกลิงค์ที่ออกจากเพจนั้น โดยใช้ hyperlink
3. สำหรับใส่ยูอาร์แอลเริ่มต้นในการเก็บข้อมูล
4. กำหนดค่าความลึกในการค้นหาให้มีค่าเป็น 2 ซึ่งเว็บสฟิงซ์จะเก็บข้อมูลเว็บเพจจากหน้าเริ่มต้นและอีกสองหน้าถัดไปจากลิงค์นั้นๆ
5. กำหนดให้ทำการเก็บข้อมูลตามแนวตั้ง คือจะเก็บข้อมูลเพจเริ่มต้นให้เรียบร้อยก่อน และไปยังลิงค์อื่นๆ ตามที่กำหนดค่าความลึก

บันทึกข้อมูลของเว็บเพจที่ทำการครอเลอร์มาได้ โดยคลิกที่แท็บ Pages เลือก Action เป็น Save และใส่ไดเรกทอรีที่ต้องการเก็บข้อมูล ไปที่ช่อง to directory ซึ่งเก็บทุกเพจที่ครอเลอร์ได้ เลือก all pages ใน on pages ดังรูปที่ 3.5



รูปที่ 3.5 กำหนดการใช้งานของแท็บ Pages

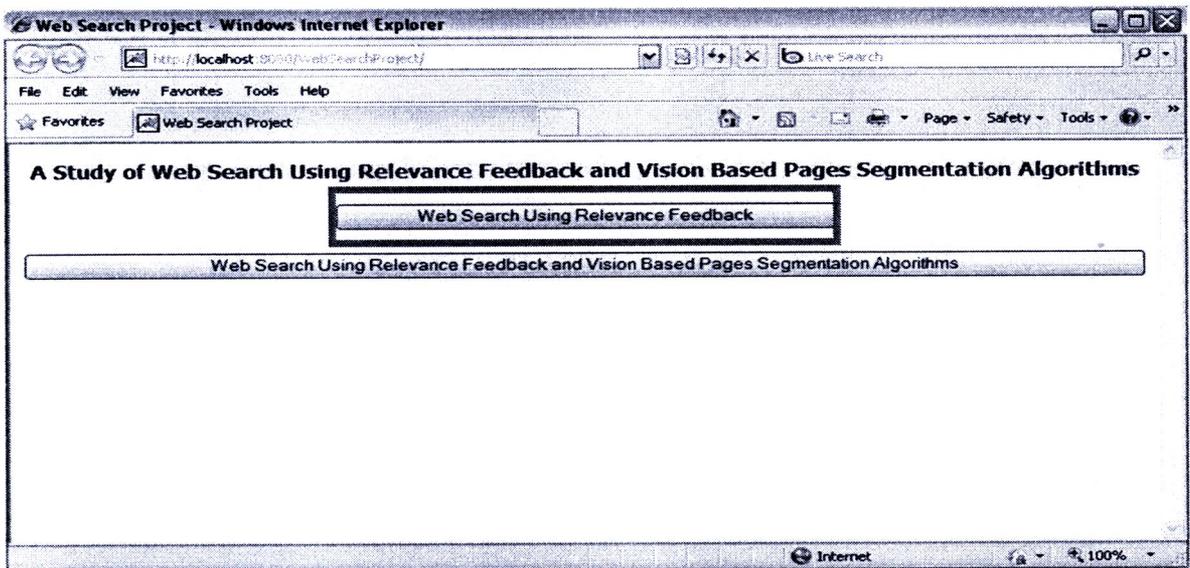


รูปที่ 3.6 แสดงการทำงานของเว็บสฟิงซ์

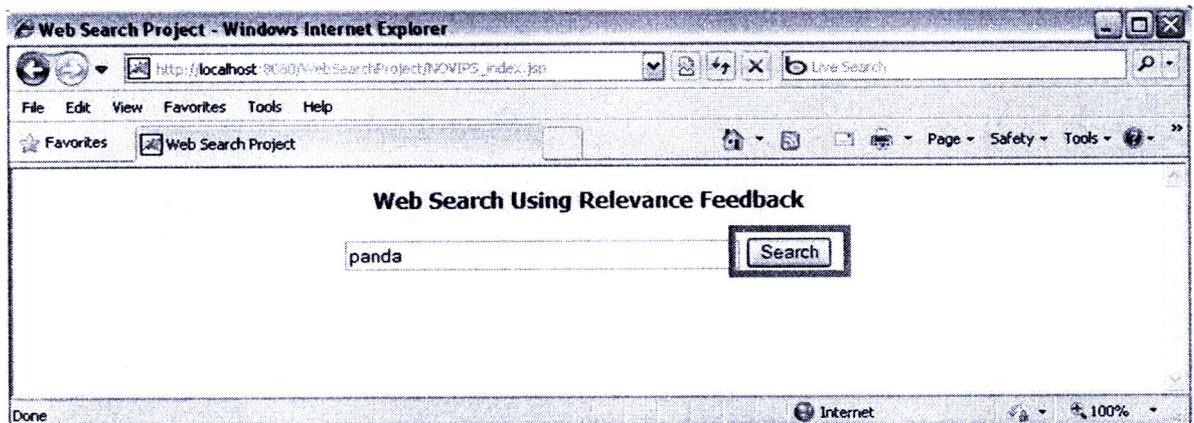
เมื่อกำหนดค่าข้างต้นแล้ว ก็สามารถ กดปุ่ม Start เพื่อเริ่มทำการเก็บข้อมูลเว็บเพจได้เลย โดยแท็บ Graph จะแสดงเส้นทางที่เว็บสฟิงซ์ได้ทำการเก็บข้อมูล ดังรูปที่ 3.6 ยูอาร์แอลของเว็บที่นำมาเก็บข้อมูลนี้ จะเป็นยูอาร์แอลของเว็บที่เกี่ยวข้องกับเรื่อง panda aids java sushi และ titanic ซึ่งเมื่อทำการครอเลอร์เรียบร้อยแล้วจะต้องนำเพจทุกเพจ มาวิเคราะห์ว่าเกี่ยวข้องกับเรื่องที่กำหนดหรือไม่ โดยให้แต่ละเรื่องมีเพจที่เกี่ยวข้อง 30 เพจจากอย่างน้อย 5 เว็บไซต์

3.3 การทำงานของระบบค้นคืนสารสนเทศบนเว็บโดยใช้การค้นคืนย้อนกลับ

การเริ่มต้นการทำงานของระบบค้นคืนสารสนเทศบนเว็บโดยใช้การค้นคืนย้อนกลับ จะต้องการเลือกปุ่ม “Web Search Using Relevance Feedback” ในหน้าเริ่มต้นของเว็บดังรูปที่ 3.7



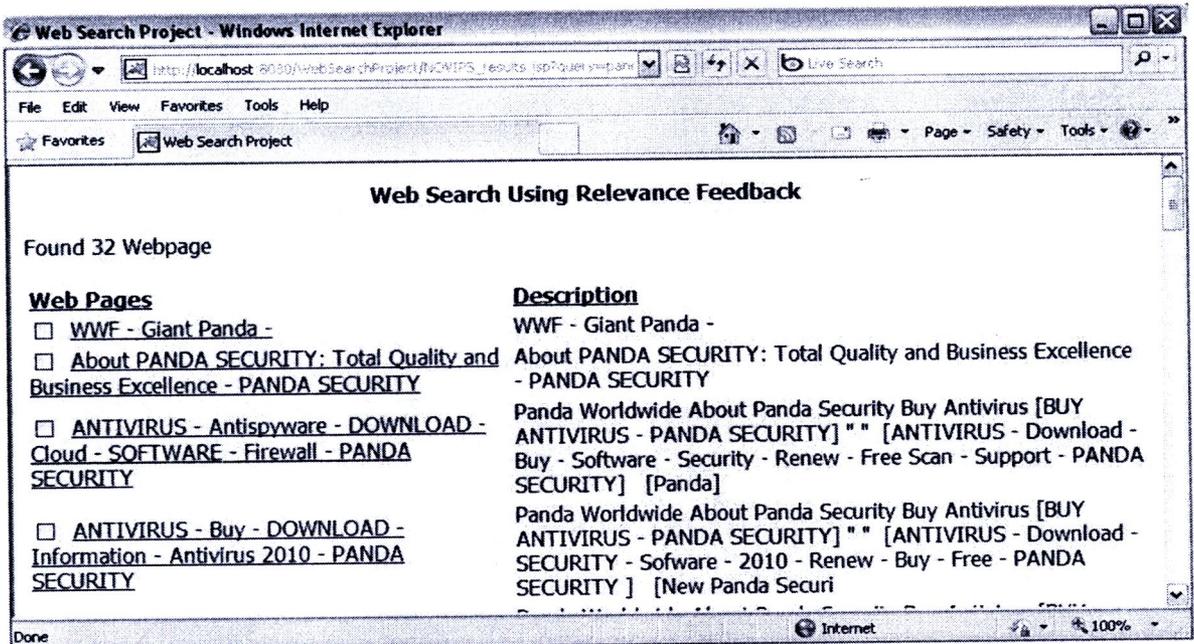
รูปที่ 3.7 หน้าเว็บเพจแรกของระบบจำลอง



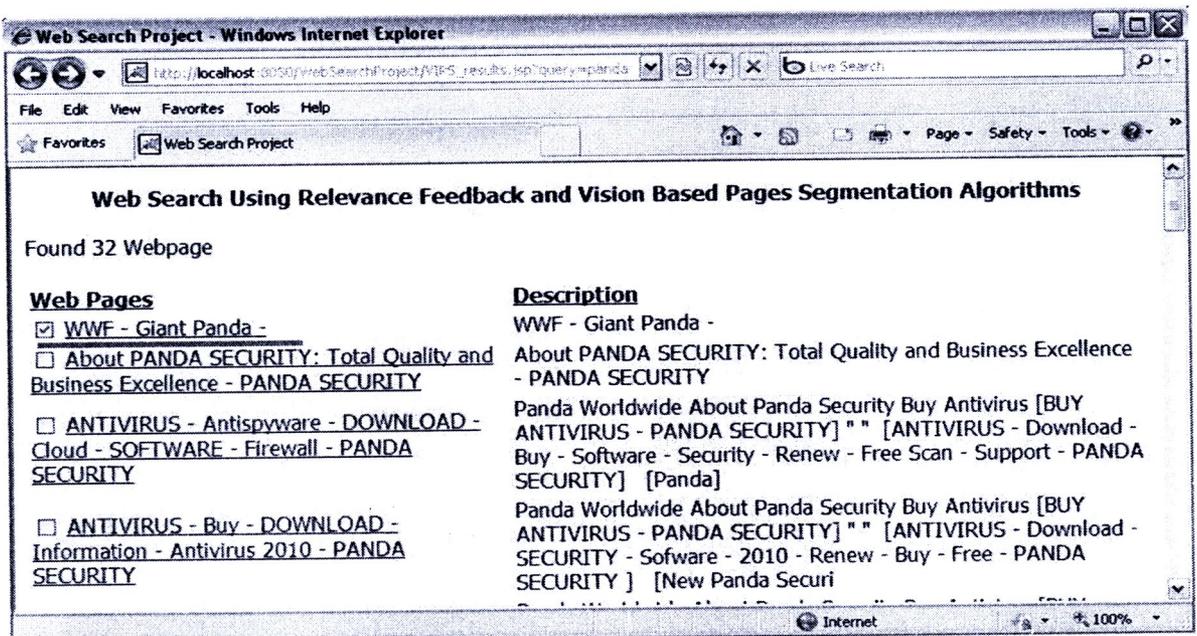
รูปที่ 3.8 หน้าเว็บเพจรับข้อความจากผู้ใช้

3.3.1 ขั้นตอนการทำงานของการค้นหาสารสนเทศบนเว็บ

ผู้ใช้ใส่ข้อความที่ต้องการทำการค้นหา ซึ่งระบบจะรองรับเฉพาะข้อความที่เป็นภาษาอังกฤษเท่านั้น แล้วกดปุ่ม Search ดังรูปที่ 3.8 ผู้ใช้ได้ใส่ข้อความคำว่า "panda" เมื่อกดปุ่ม search แล้วก็จะแสดงเว็บเพจที่เกี่ยวข้องกับข้อความออกมาโดย Web Pages คือชื่อเพจ สามารถกดลิงค์ไปยังเพจนั้นๆได้ และ Description คือคำอธิบายของเพจนั้นๆ ดังรูปที่ 3.9



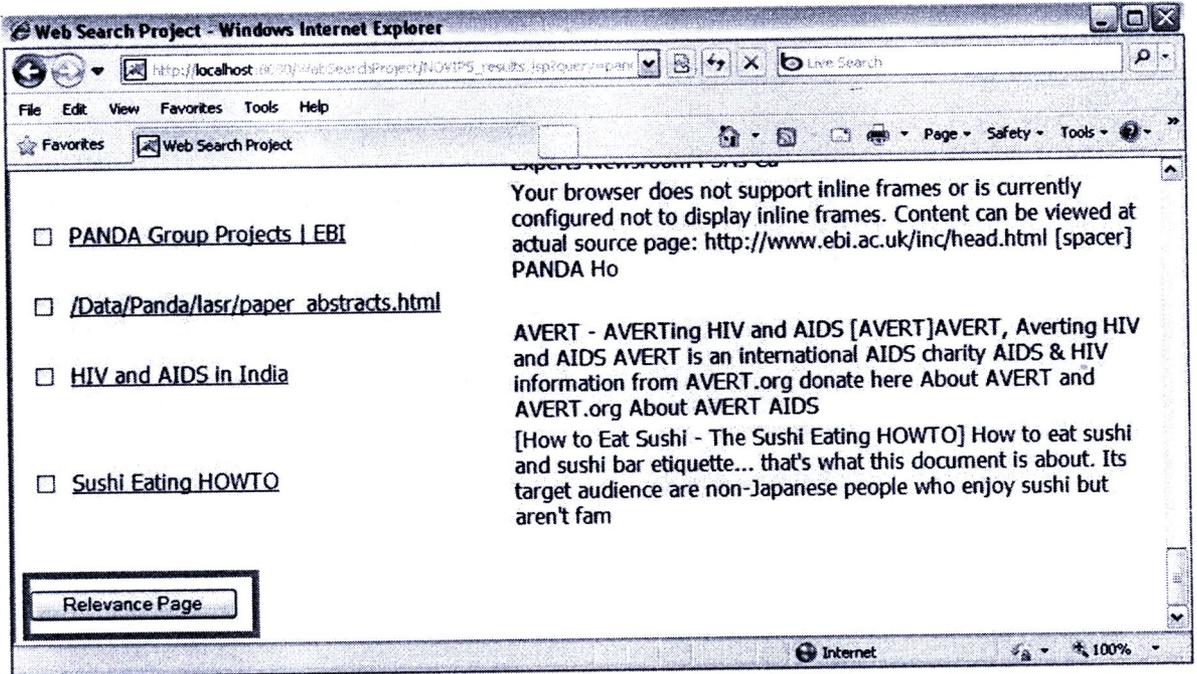
รูปที่ 3.9 ผลการค้นหาของข้อความ



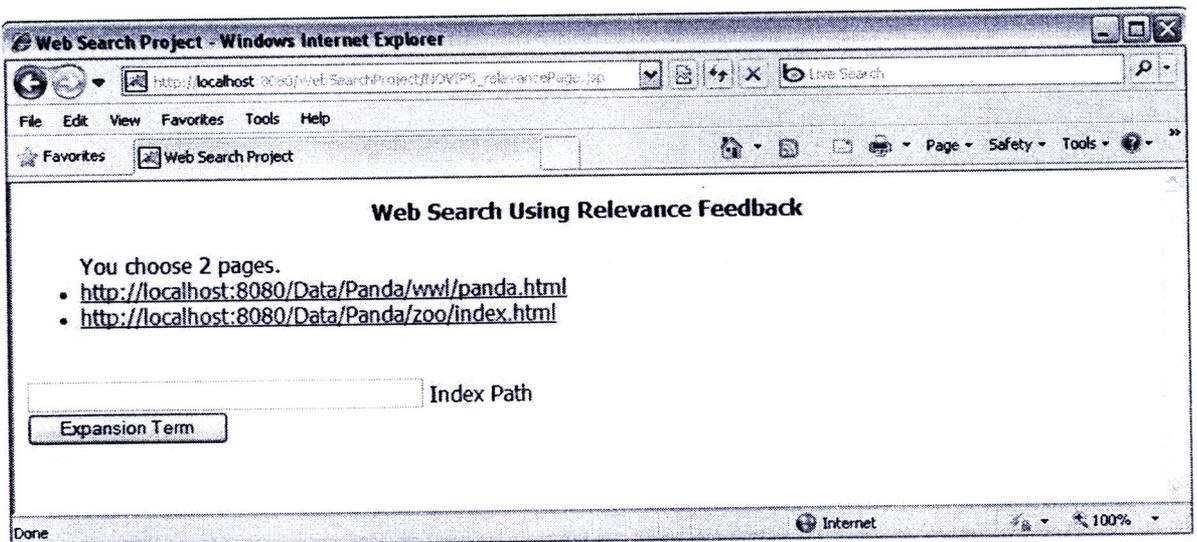
รูปที่ 3.10 หน้าเว็บเพจให้ผู้ใช้เลือกเว็บเพจที่เกี่ยวข้องกับสิ่งที่ผู้ใช้ต้องการ

3.3.2 ขั้นตอนการทำงานของการค้นหาค้นคืนย้อนกลับ

หากผู้ใช้ไม่พอใจผลการค้นคืนครั้งแรกต้องการทำการค้นคืนอีกครั้ง แล้วระบบจะทำการหาข้อความใหม่ กลับมาให้ผู้ใช้ ซึ่งในขั้นตอนการค้นคืนย้อนกลับนี้ผู้ใช้จะต้องทำการเลือกเว็บเพจที่เกี่ยวข้องกับสิ่งที่ผู้ใช้ต้องการจากผลการค้นคืนครั้งแรกโดยคลิกปุ่มที่ด้านหน้าของเพจดังรูปที่ 3.10 จากนั้นให้เลื่อนลงมายังด้านล่างของเพจแล้ว กดปุ่ม Relevance Pages ดังรูปที่ 3.11 หลังจากที่ใช้กดปุ่ม Relevance Page ระบบก็จะแสดงยูอาร์แอลของเพจที่ผู้ใช้เลือกกว่าเกี่ยวข้องกับสิ่งที่ผู้ใช้ต้องการ และช่อง "Index Path" สำหรับใส่ไอดีเรกทอรีที่เก็บไฟล์ดัชนีของเว็บที่ผู้ใช้เลือก ดังรูปที่ 3.12



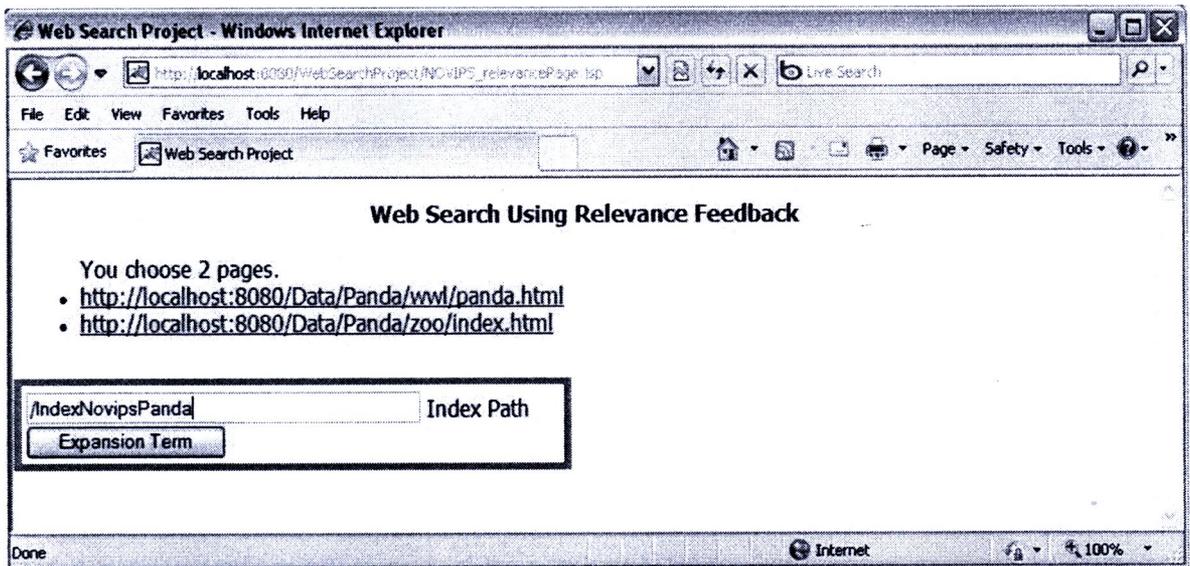
รูปที่ 3.11 หน้าเว็บด้านล่างสุดของเพจที่ทำการค้นคืนได้



รูปที่ 3.12 ผลจากการกดปุ่ม Relevance Page

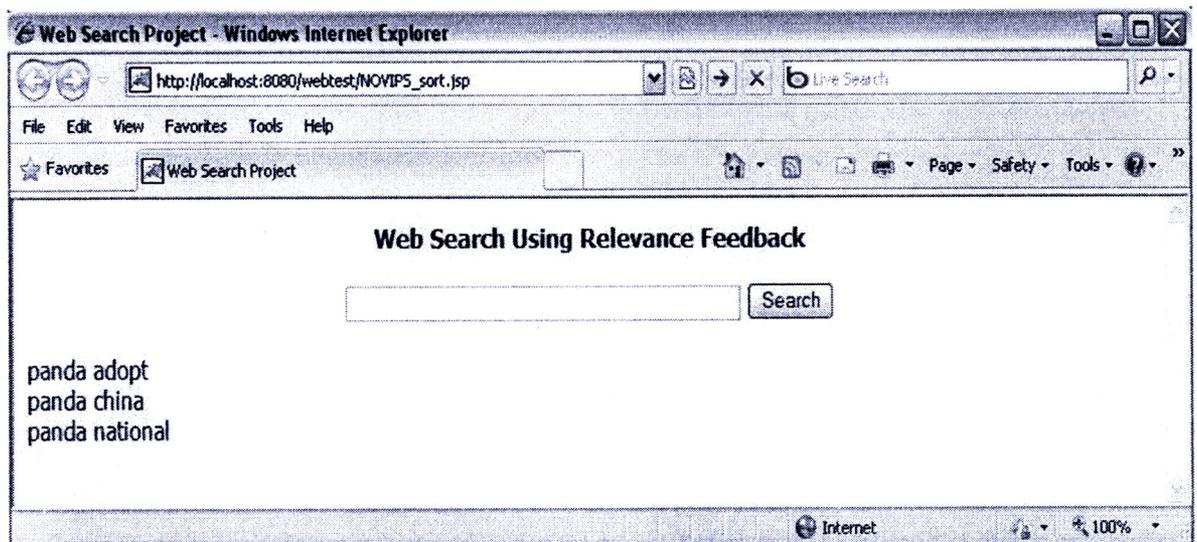
3.3.3 ขั้นตอนการทำงานของการทำงานหาคำถามใหม่

ผู้ใช้จะต้องนำเพจที่ผู้ใช้เลือกเกี่ยวข้องกับสิ่งที่ผู้ใช้ต้องการไปทำเป็นข้อมูลดัชนี และนำโดเรทอร์ที่เก็บไฟล์ดัชนีนั้นมาใส่ในช่อง “Index Path” ดังรูปที่ 3.13 ทำการเก็บไฟล์ดัชนีไว้ที่ “/IndexNovipsPanda” เมื่อใส่โดเรทอร์ของไฟล์ดัชนีเรียบร้อยแล้วก็สามารถกดปุ่ม Expansion เพื่อหาคำถามใหม่



รูปที่ 3.13 แสดงยูอาร์แอลที่ผู้ใช้เลือกและการใส่ค่าโดเรทอร์เพื่อหาคำถามใหม่

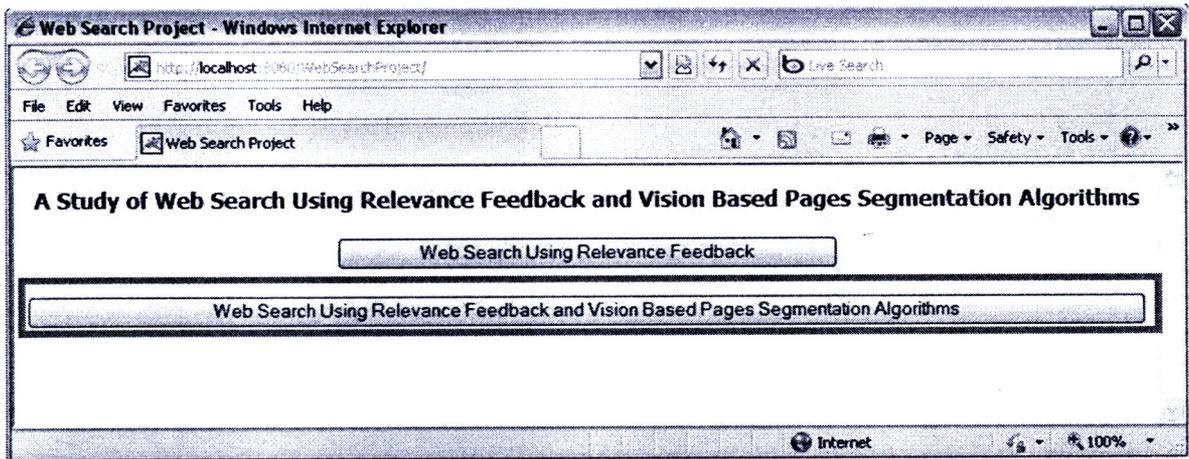
ข้อความใหม่พร้อมให้ผู้ใช้ย้อนกลับไปค้นคืนได้อีกครั้งหนึ่ง ดังรูปที่ 3.14 คำใหม่ที่นำมาเพิ่มในข้อความเดิมเพื่อให้ได้ข้อความใหม่นั้น จะมาจากการเลือกเทอมที่มีความสำคัญในเว็บเพจที่ผู้ใช้ได้ทำการเลือกมา โดยคำนวณตามอัลกอริทึมการเลือกเทอมของ ดอนน่า ฮาแมน



รูปที่ 3.14 ข้อความใหม่จากเว็บเพจที่ผู้ใช้เลือกเกี่ยวข้องกับสิ่งที่ผู้ใช้ต้องการ

3.4 การทำงานของระบบค้นคืนสารสนเทศบนเว็บโดยใช้ อัลกอริทึมวีไอพีเอสและการค้นคืนย้อนกลับ

การการทำงานของระบบค้นคืนสารสนเทศบนเว็บโดยใช้อัลกอริทึมวีไอพีเอสและการค้นคืนย้อนกลับ จะต้องการเลือกปุ่ม “Web Search Using Relevance Feedback” ในหน้าเริ่มต้นของเว็บดังรูปที่ 3.15



รูปที่ 3.15 หน้าหลักของระบบจำลอง

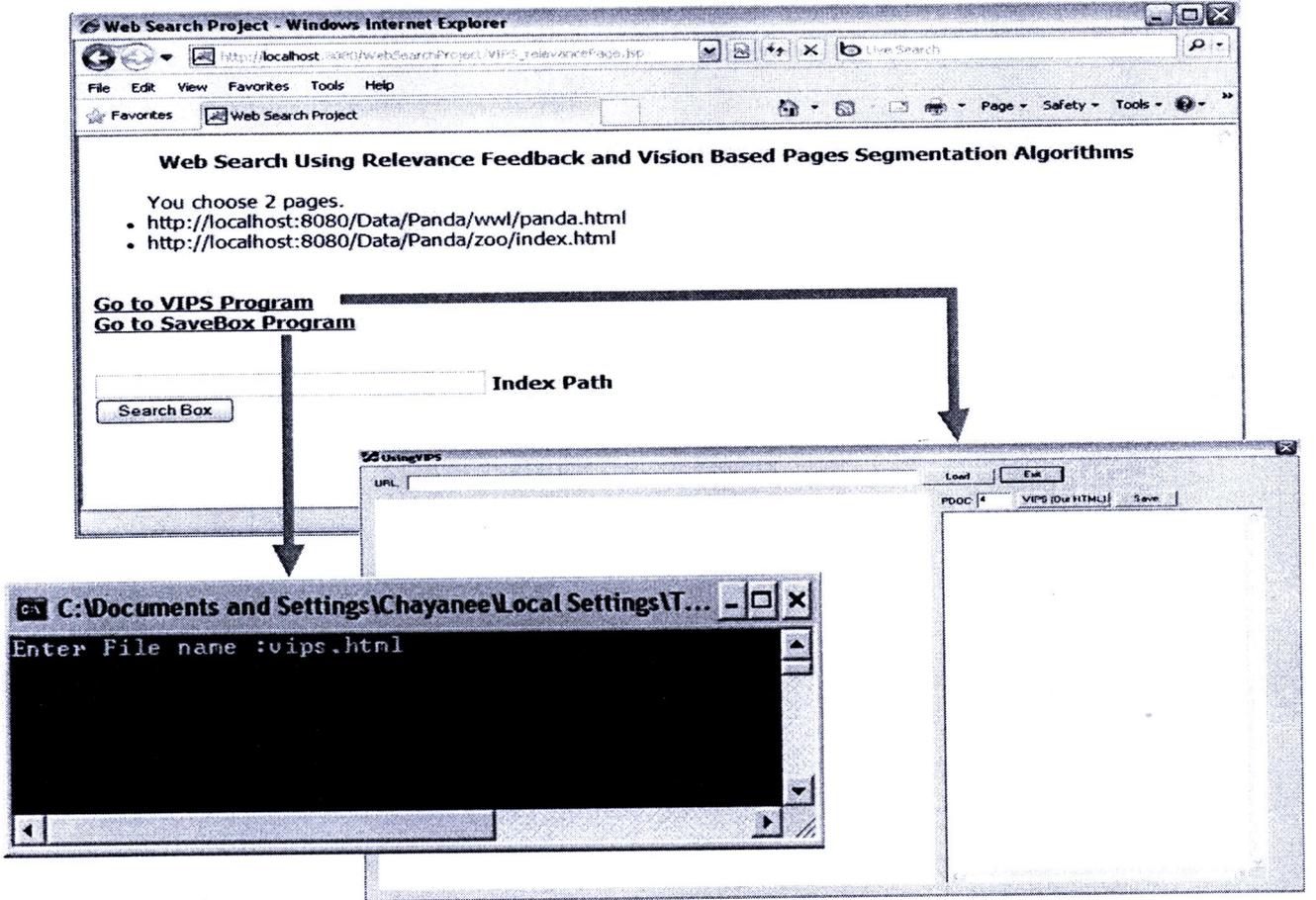
หลังจากที่ผู้ใช้ได้เลือกแล้ว การทำงานของระบบการงานของระบบค้นคืนสารสนเทศบนเว็บโดยใช้ อัลกอริทึมวีไอพีเอสและการค้นคืนย้อนกลับนี้ จะมีการทำงานเช่นเดียวกันกับระบบค้นคืนสารสนเทศบนเว็บ โดยใช้การค้นคืนย้อนกลับจนกระทั่งผู้ใช้เลือกเว็บเพจที่เกี่ยวข้องกับสิ่งที่ผู้ใช้ต้องการ และกดปุ่ม Relevance Page ระบบจะแสดงผลดังรูปที่ 3.16 โดยระบบจะแสดงยูอาร์แอลของเพจที่ผู้ใช้เลือกว่าเกี่ยวข้องกับสิ่งที่ผู้ใช้ต้องการ จะต้องนำเพจที่ผู้ใช้เลือกเหล่านั้นไปทำการแบ่งเป็นบล็อกโดยใช้โปรแกรมวีไอพีเอส และแยกไฟล์ของแต่ละบล็อกโดยใช้โปรแกรมการบันทึกไฟล์ของแต่ละบล็อก

3.4.1 ขั้นตอนการทำงานของโปรแกรมวีไอพีเอส

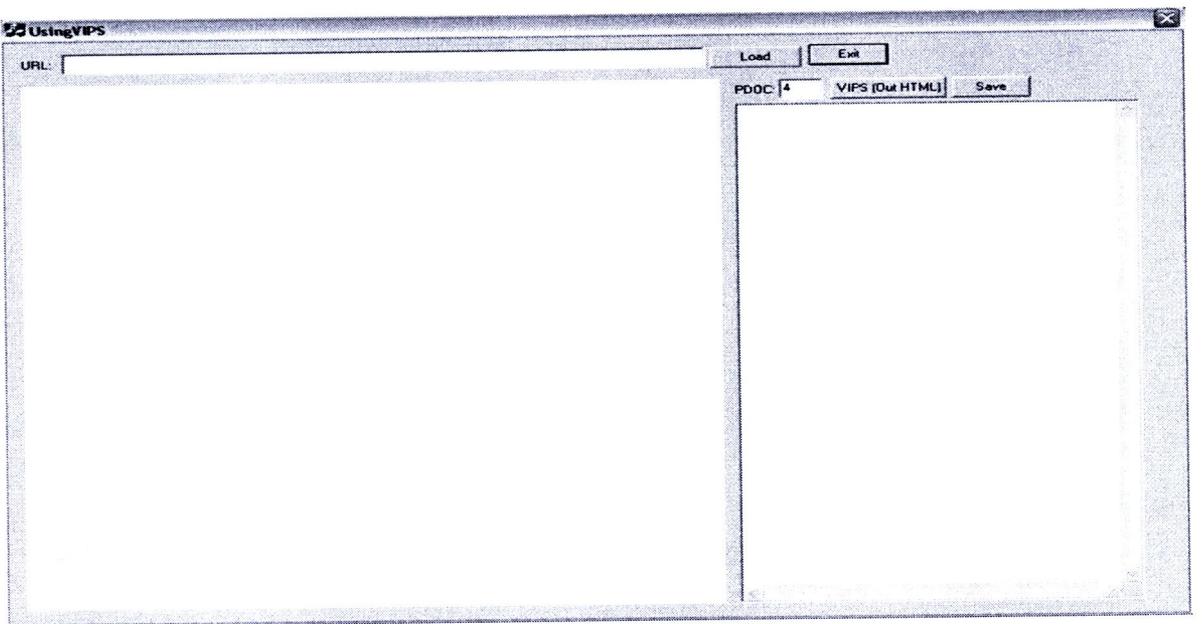
เมื่อผู้ใช้กดคลิก “Go to VIPS Program” จะเป็นการเปิดโปรแกรมขึ้นมาและพบกับส่วนติดต่อกับผู้ใช้ของโปรแกรมวีไอพีเอส แสดงดังรูปที่ 3.17 โดยมีการทำงานตามขั้นตอนดังต่อไปนี้คือ

1. ใส่ยูอาร์แอลที่ต้องการแบ่งเป็นบล็อกในช่อง URL และกดปุ่ม Load เพื่ออ่านค่าดังรูปที่ 3.18
2. กำหนดค่าความละเอียดของโครงสร้างเนื้อหาในช่อง PDOC
3. ทำการแบ่งเว็บเพจเป็นบล็อกโดยกดปุ่ม VIPS(Out HTML) จะได้โค้ดเอชทีเอ็มแอลที่มีแท็กวีไอพีเอส แสดงการแบ่งเป็นบล็อก ดังรูปที่ 3.19 จากนั้นกดปุ่ม SAVE เพื่อนำไปทำการแยกไฟล์ของแต่ละบล็อกต่อไป

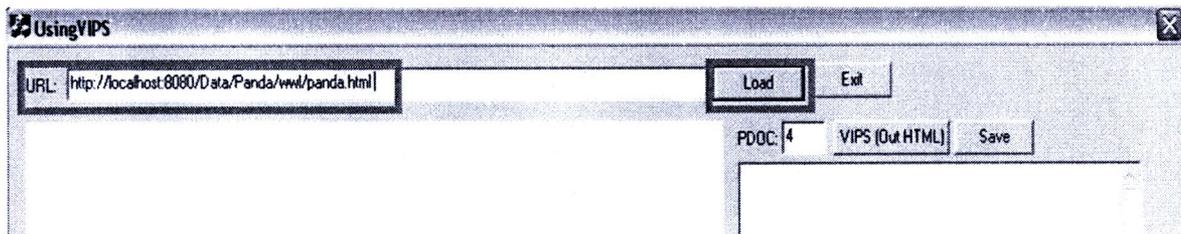




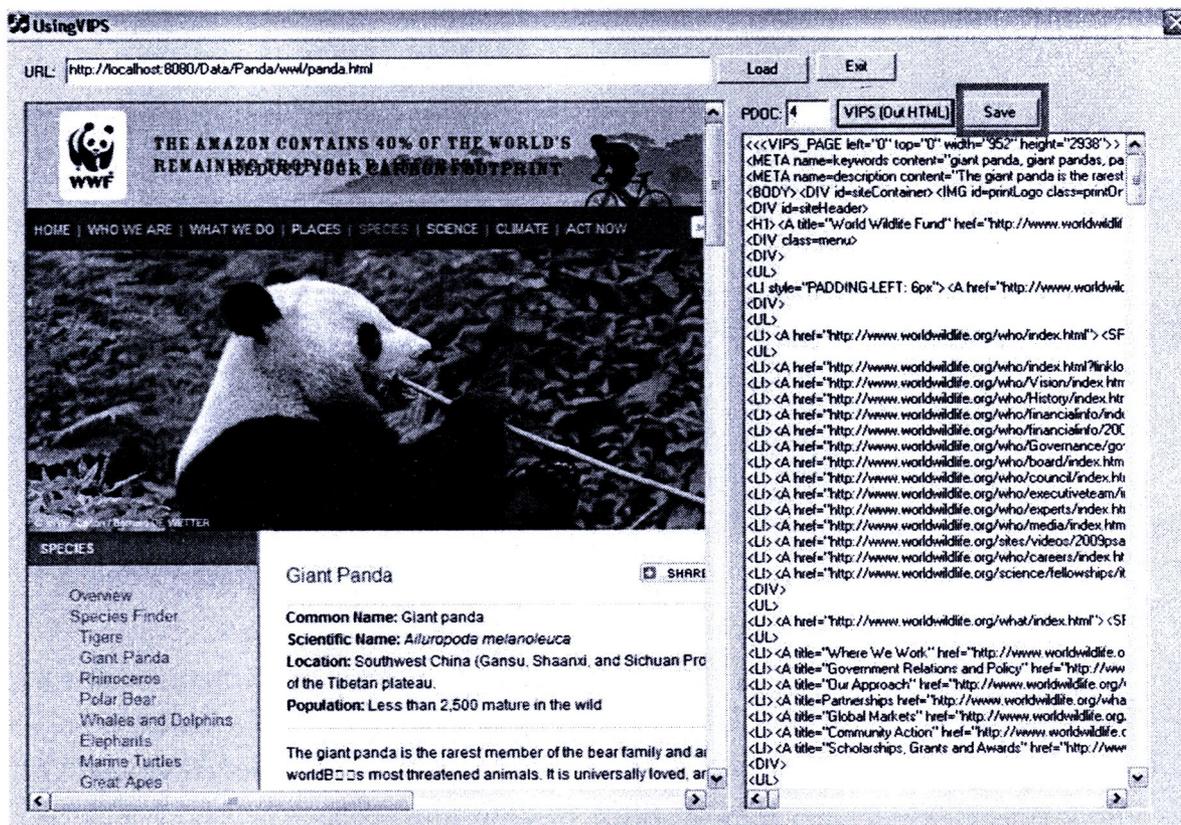
รูปที่ 3.16 แสดงการลิงก์ไปยังโปรแกรมวีไอพีเอสและโปรแกรมการบันทึกไฟล์ของแต่ละบล็อก



รูปที่ 3.17 ส่วนติดต่อกับผู้ใช้ของโปรแกรมวีไอพีเอส



รูปที่ 3.18 แสดงการใส่ค่ายูอาร์แอลของเว็บเพจที่ต้องการแบ่งเป็นบล็อก

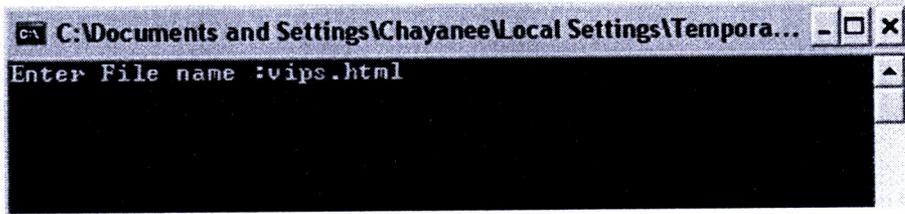


รูปที่ 3.19 ผลของการแบ่งเว็บเพจเป็นบล็อกจากโปรแกรมวีไอพีเอส

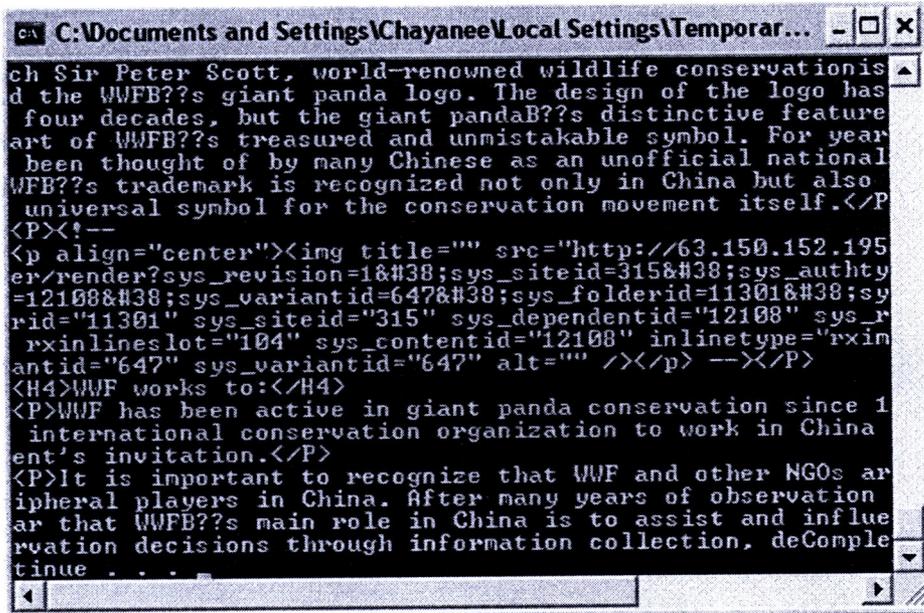
3.4.2 ขั้นตอนการทำงานของโปรแกรมบันทึกไฟล์ของแต่ละบล็อก

จากไฟล์เอชทีเอ็มแอลที่ได้จากโปรแกรมวีไอพีเอสยังไม่ได้ทำการแยกแต่ละบล็อกออกมาเป็นไฟล์จึงใช้โปรแกรมนี้ในการบันทึกไฟล์ของแต่ละบล็อก เมื่อผู้ใช้คลิก “Go to VIPS Program” ก็จะเป็นการเปิดโปรแกรมขึ้นมา โดยมีขั้นตอนการใช้งาน คือ

1. ใส่ชื่อไฟล์ที่ได้จากโปรแกรมวีไอพีเอส ดังรูปที่ 3.20
2. โปรแกรมจะทำการแยกไฟล์เป็นบล็อกตามแท็กที่โปรแกรมวีไอพีเอสกำหนด จากนั้นกดปุ่มเอนเทอร์ (Enter) เพื่อจบโปรแกรม ดังรูปที่ 3.21



รูปที่ 3.20 โปรแกรมการบันทึกไฟล์ของแต่ละบล็อกโดยใส่ชื่อไฟล์ที่ต้องการแยกไฟล์ของแต่ละบล็อก



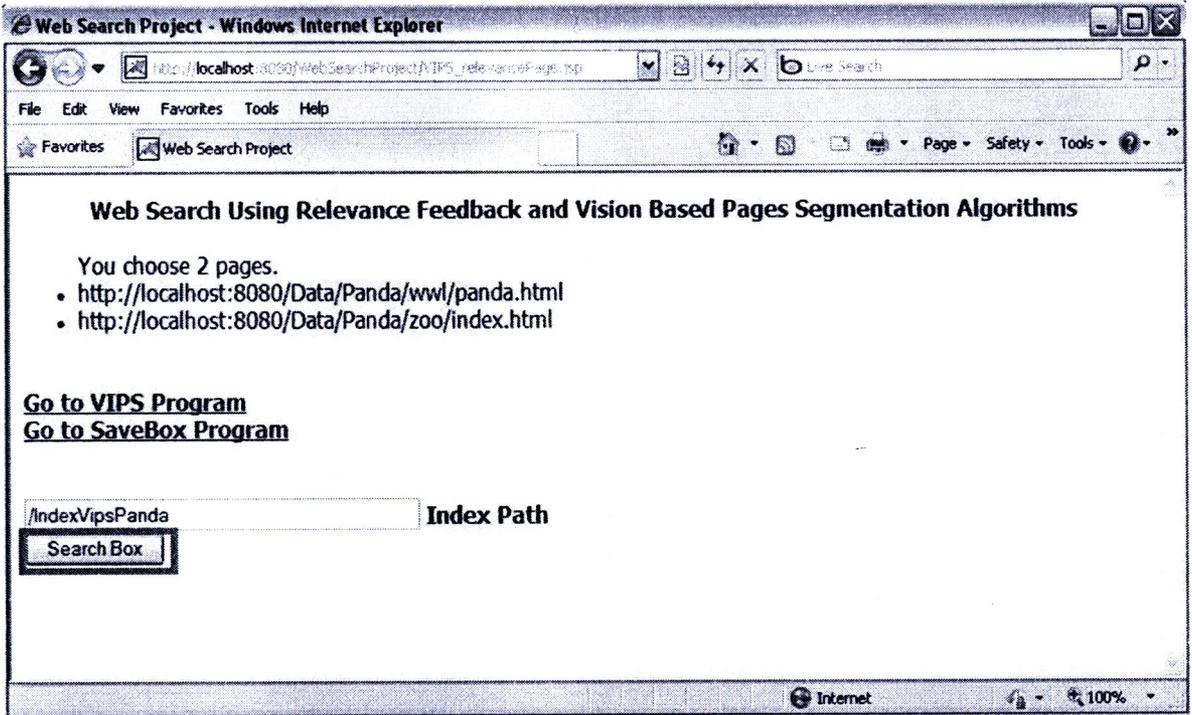
รูปที่ 3.21 การทำงานของโปรแกรมการบันทึกไฟล์ของแต่ละบล็อก

3.4.3 ขั้นตอนการทำงานของการทำงานหาคำถามใหม่

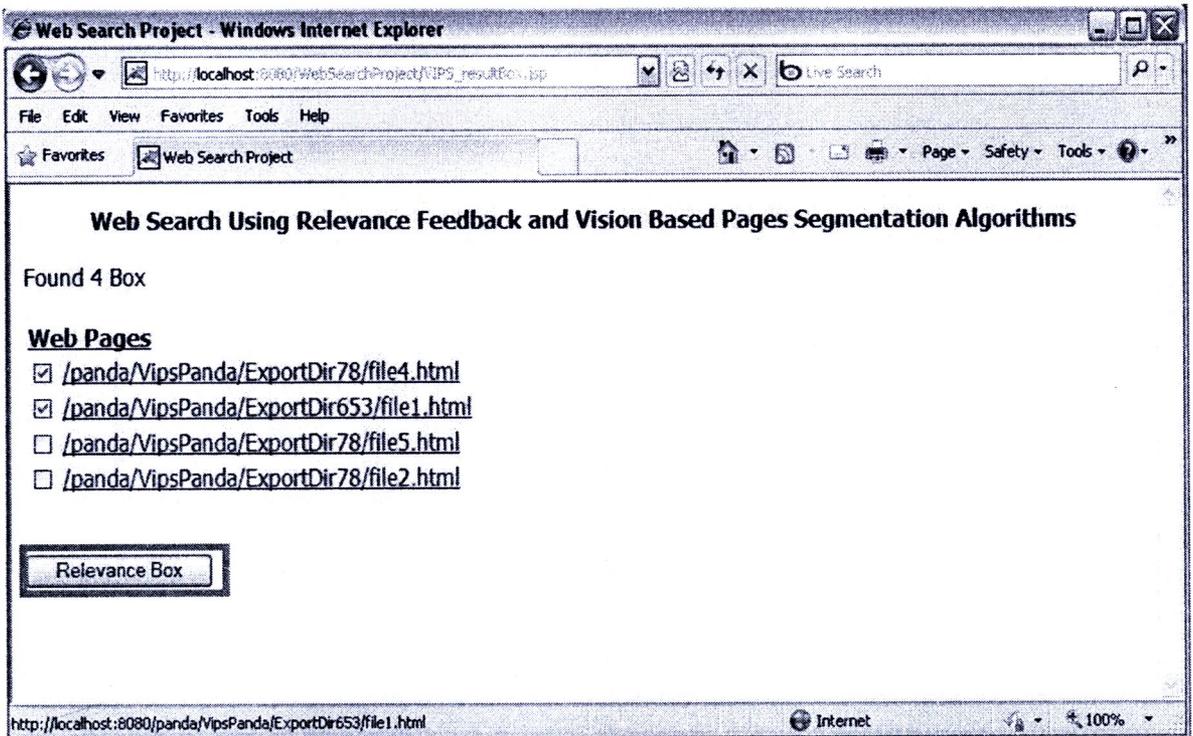
นำบล็อกที่แบ่งได้มาทำเป็นไฟล์ดัชนีและใส่โดเรทอร์ที่เก็บไฟล์นั้นในช่อง "Index Path" ดังรูปที่ 3.22 โดเรทอร์ที่เก็บไฟล์ดัชนีชื่อ "/IndexVipsPanda" จากนั้นกดปุ่ม "Search Box" ระบบจะแสดงบล็อกที่เกี่ยวข้องกับข้อความที่ทำการค้นหาตั้งแต่ต้น ดังรูปที่ 3.23 จากนั้นผู้ใช้จะต้องเลือก บล็อกที่มีข้อมูลเกี่ยวข้องกับสิ่งที่ผู้ใช้ต้องการจากนั้นกดปุ่ม "Relevance Box" ในที่นี้ได้ทำการเลือกบล็อกที่ 1 และ 2 ว่ามีความเกี่ยวข้องกับสิ่งที่ต้องการ

หลังจากที่ผู้ใช้กดปุ่ม Relevance Box ระบบก็จะแสดงยูอาร์แอลของบล็อกที่ผู้ใช้เลือกเกี่ยวข้องกับสิ่งที่ผู้ใช้ต้องการ จะต้องนำเพจที่ผู้ใช้เลือกเหล่านี้ไปทำเป็นข้อมูลดัชนี และนำโดเรทอร์ที่เก็บไฟล์ดัชนีนั้นมาใส่ในช่อง "Index Path" ในที่นี้ทำการเก็บไฟล์ดัชนีไว้ที่ "/IndexBoxPanda" ดังรูปที่ 3.24 เมื่อใส่โดเรทอร์ของไฟล์ดัชนีเรียบร้อยแล้วก็สามารถกดปุ่ม Expansion เพื่อหาคำถามใหม่ ก็จะได้ข้อความใหม่พร้อมให้ผู้นำกลับไปค้นคืนได้อีกครั้งหนึ่ง ดังรูปที่ 3.25

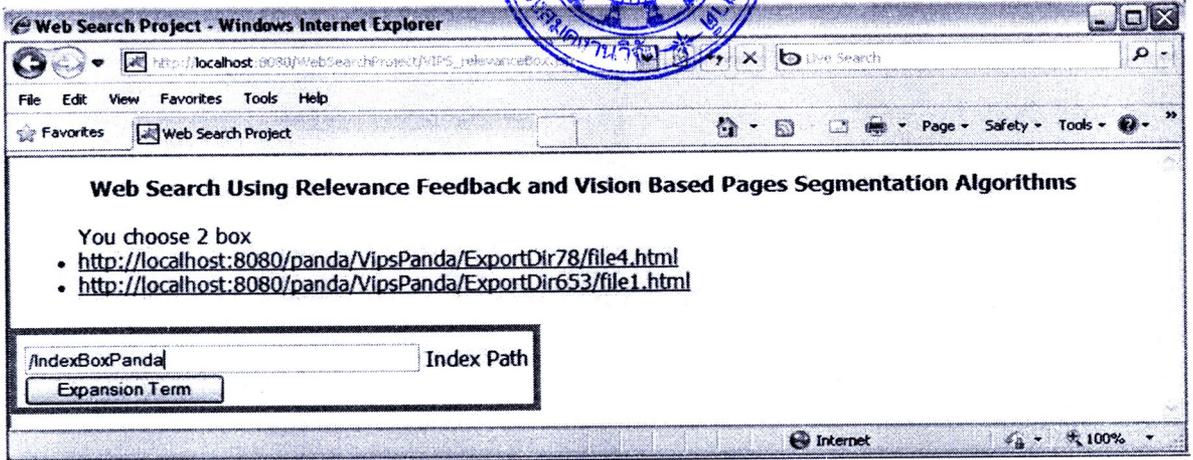
คำใหม่ที่นำมาเพิ่มในข้อความเดิมเพื่อให้ได้ข้อความใหม่นั้นจะมาจากทางเลือกที่มีความสำคัญในเว็บเพจที่ผู้ใช้ได้ทำการเลือกมา โดยคำนวณตามอัลกอริทึมการเลือกเทอมของ ดอนน่า ฮาแมน เช่นเดียวกับระบบค้นหาสารสนเทศบนเว็บโดยใช้การค้นหาค้นย้อนกลับ



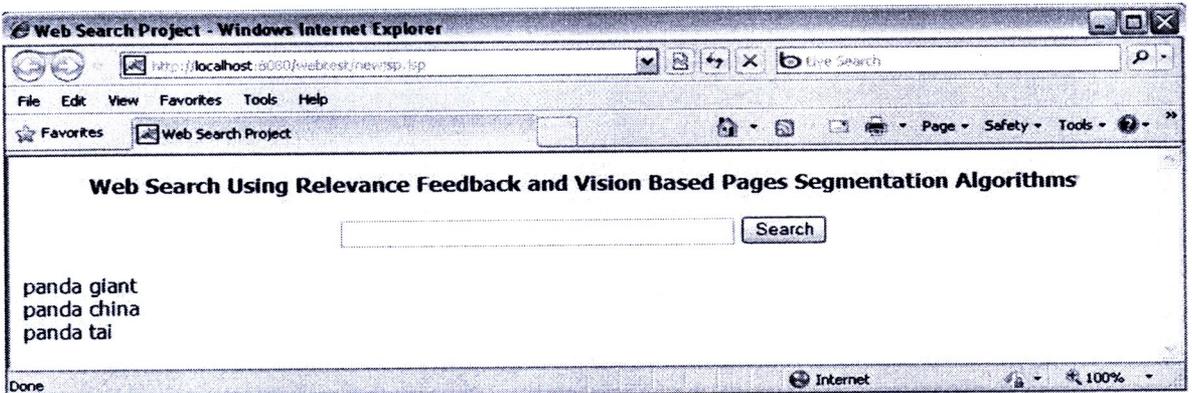
รูปที่ 3.22 การใส่ไดเรกทอรีของไฟล์ดัชนีเว็บเพจที่แบ่งเป็นบล็อกแล้ว



รูปที่ 3.23 บล็อกที่เกี่ยวข้องกับข้อความตั้งต้น



รูปที่ 3.24 หน้าเพจหลังจากกดปุ่ม Relevance Box จะแสดงยูอาร์แอลที่ผู้ใช้เลือกและการหาข้อความใหม่

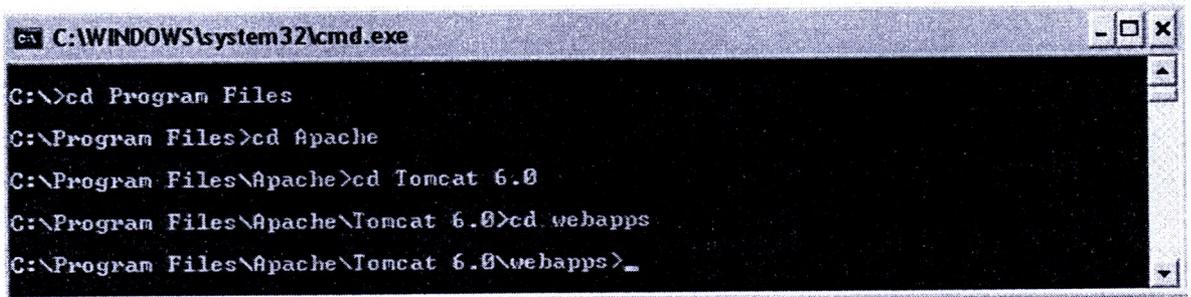


รูปที่ 3.25 ผลการหาข้อความใหม่จากบล็อกที่ผู้ใช้เลือกที่เกี่ยวข้อง

3.5 การสร้างไฟล์ดัชนี

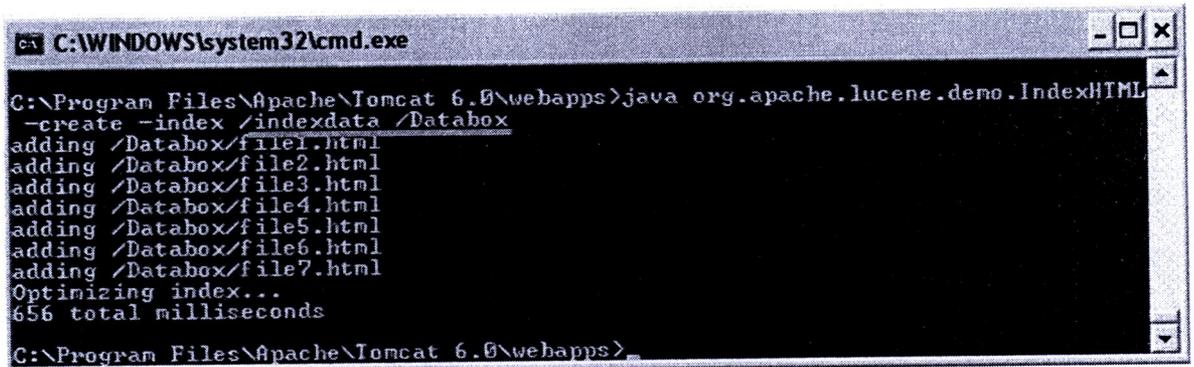
การค้นคืนของลูซินั้นจะทำการค้นคืนผ่านทางดัชนีซึ่งเป็นค่าดัชนีที่เป็นตัวแทนเอกสารโดยมีขั้นตอนการสร้างดัชนีคือ

1. เปิด command Prompt
2. เข้าไปยังไดเรกทอรี webapps ของ Server ในที่นี่ใช้ Apache Tomcat 6.0 ดังรูปที่ 3.26



รูปที่ 3.26 การเข้าไปในไดเรกทอรีของ webapps ของ Apache Tomcat 6.0

- พิมพ์คำสั่งเพื่อสั่งการสร้างดัชนี คือ "java org.apache.lucene.demo.IndexHTML -create -index ไดรกทอรี่ที่ต้องการให้เก็บค่าดัชนี ไดรกทอรี่ที่ต้องการทำไฟล์มาทำเป็นดัชนี ยกตัวอย่างการทำดัชนีของข้อมูลภายในโฟรเดอร์ Databox ให้ไปเก็บโฟรเดอร์ indexdata ดังรูปที่ 3.27



```
C:\WINDOWS\system32\cmd.exe
C:\Program Files\Apache\Tomcat 6.0\webapps>java org.apache.lucene.demo.IndexHTML
-create -index /indexdata /Databox
adding /Databox/file1.html
adding /Databox/file2.html
adding /Databox/file3.html
adding /Databox/file4.html
adding /Databox/file5.html
adding /Databox/file6.html
adding /Databox/file7.html
Optimizing index...
656 total milliseconds
C:\Program Files\Apache\Tomcat 6.0\webapps>
```

รูปที่ 3.27 การทำไฟล์ดัชนีของข้อมูลภายในโฟรเดอร์ Databox และนำไปเก็บใน โฟรเดอร์ indexdata