

**ĐẠI HỌC QUỐC GIA HÀ NỘI
TRƯỜNG ĐẠI HỌC CÔNG NGHỆ**

NGUYỄN THỊ LOAN

**NGHIÊN CỨU CÔNG NGHỆ TÌM KIẾM (MÃ
NGUỒN MỞ) LUCENE ÁP DỤNG GIẢI QUYẾT BÀI
TOÁN TÌM KIẾM TRONG HỆ THỐNG VĂN BẢN**

Ngành: Công nghệ Thông tin

Chuyên ngành: Kỹ thuật phần mềm

Mã số: 62.48.01.03

**TÓM TẮT LUẬN VĂN THẠC SĨ
NGÀNH CÔNG NGHỆ THÔNG TIN**

Hà Nội - 2017

MỤC LỤC

MỞ ĐẦU	4
CHƯƠNG 1: NGHIÊN CỨU TỔNG QUAN VỀ HỆ THỐNG TÌM KIẾM THÔNG TIN	6
1.1. Khái niệm về hệ thống tìm kiếm thông tin	6
1.2. Các bộ phận cấu thành hệ thống tìm kiếm thông tin	6
1.3. Hệ thống tìm kiếm thông tin của Google	7
1.4. Kiến trúc của hệ thống tìm kiếm thông tin	8
CHƯƠNG 2: NGHIÊN CỨU TỔNG QUAN VỀ MÃ NGUỒN MỞ LUCENE	10
2.1. Giới thiệu về thư viện Lucene.....	10
2.2. Quy trình đánh chỉ mục	11
2.3. Các toán tử đánh chỉ mục cơ bản.....	11
2.4. Tối ưu hóa việc đánh chỉ mục.....	12
2.5. Tính đồng thời, an toàn tiến trình, ngăn chặn các thực thi	12
2.6. Bộ chuyển đổi câu truy vấn của người dùng: QueryParser	12
2.7. Các biểu thức truy vấn của QueryParser.....	13
2.8. Bộ phân tích – Analyzer:	13
2.9. Sử dụng lớp IndexSearcher	13
2.10. Cú pháp truy vấn Lucene	14
2.11. Các máy tìm kiếm phát triển dựa trên Lucene.....	14
CHƯƠNG 3: XÂY DỰNG ỨNG DỤNG THỬ NGHIỆM	15

3.1. Tài mã nguồn Lucene.NET	15
3.2. Dữ liệu văn bản thử nghiệm	15
3.3. Mô hình cơ sở dữ liệu	16
3.3.1. Lược đồ cơ sở dữ liệu	16
3.4. Giao diện chính	16
3.4.1. Giao diện trang Quản lý văn bản	16
3.4.2. Giao diện trang Cập nhật văn bản.....	17
3.4.3. Giao diện trang Tìm kiếm văn bản	17
3.4.4. Giao diện trang Xem nội dung file văn bản	18
3.5. Đánh giá và thử nghiệm	19
3.5.1. Mô hình kiến trúc ứng dụng thử nghiệm	19
3.5.2. Kịch bản và kết quả.....	20
CHƯƠNG 4: KẾT LUẬN.....	23
4.1. Đánh giá kết quả nghiên cứu	23
TÀI LIỆU THAM KHẢO.....	24

MỞ ĐẦU

Với sự phát triển không ngừng của công nghệ thông tin, số lượng các tài liệu điện tử do con người tạo ra ngày càng phong phú và đa dạng, nhu cầu khai thác dữ liệu trong kho tài liệu là rất lớn, đây là một trong những nhu cầu thường ngày và thiết thực của người sử dụng. Tuy nhiên, một trong những khó khăn con người gặp phải trong việc khai thác thông tin là: Khả năng tìm kiếm chính xác thông tin cần tìm trong kho tài liệu, khả năng tìm kiếm nhanh với số lượng dữ liệu lớn. Nếu dùng các hệ quản trị cơ sở dữ liệu quan hệ để tìm kiếm dữ liệu thì sẽ gặp phải các hạn chế như: Bị giới hạn ở cú pháp của ngôn ngữ SQL, tốc độ tìm kiếm chậm khi tìm kiếm gần đúng (dùng LIKE) trong cơ sở dữ liệu lớn... Điều này đã thúc đẩy cho sự ra đời của các hệ thống tìm kiếm, điển hình nhất cho các hệ thống này là các máy tìm kiếm như Google và Yahoo... Tuy nhiên, phần lớn các công cụ tìm kiếm này đều là những sản phẩm thương mại và mã nguồn được giữ bí mật. Vì vậy, nhiều đơn vị phát triển phần mềm đã tự mình xây dựng từ đầu một công cụ tìm kiếm bằng cách sử dụng các thư viện mã nguồn mở.

Trên thế giới hiện nay có một số thư viện mã nguồn mở hỗ trợ xây dựng hệ thống tìm kiếm thông tin như: Lucene, Egothor, Xapian, MG4J, Sphinx... Trong số các mã nguồn mở này thì Lucene là thư viện mã nguồn mở được nhiều tổ chức, cá nhân sử dụng nhất, cụ thể: CNET sử dụng Lucene để tìm kiếm danh sách thể loại sản phẩm, Wikipedia dùng lucene để tìm kiếm nội dung toàn văn bản. ElasticSearch và Sorl là hai một công cụ tìm kiếm rất mạnh cũng được xây dựng và phát triển dựa trên nền tảng Lucene,... Vì vậy, trong đề tài này tôi đã lựa chọn Lucene để xây dựng thử nghiệm hệ thống tìm kiếm thông tin.

Đề tài luận văn “*Nghiên cứu công nghệ tìm kiếm (Mã nguồn mở) Lucene áp dụng giải quyết bài toán tìm kiếm trong hệ thống Văn bản*” sẽ cố gắng giải quyết các vấn đề nêu trên. Luận văn kế thừa thư viện mã nguồn mở Lucene để xây dựng hệ thống tìm kiếm với hai thành phần chính là Tạo chỉ mục và Tìm kiếm.

Luận văn tập trung nghiên cứu công nghệ mã nguồn mở Lucene áp dụng cho bài toán quản lý Văn bản, đưa ra các hướng phát triển trong tương lai. Do thời gian có hạn, việc xử lý văn bản, theo dõi tiến độ xử lý, đánh giá kết quả xử lý... là phức tạp nên luận văn chỉ tập trung hoàn thiện các chức

năng về quản lý văn bản và áp dụng công nghệ Lucene để đánh chỉ mục, tìm kiếm văn bản.

Nội dung mà luận văn nghiên cứu bao gồm: Tìm hiểu tổng quan về các hệ thống tìm kiếm thông tin. Tìm hiểu tổng quan về công nghệ tìm kiếm mã nguồn mở Lucene. Phân tích, thiết kế, xây dựng ứng dụng thử nghiệm Quản lý Văn bản.

Bố cục của luận văn như sau:

Chương 1: Nghiên cứu tổng quan về hệ thống tìm kiếm thông tin, các thành phần và nguyên lý hoạt động của hệ thống tìm kiếm thông tin.

Chương 2: Nghiên cứu các tính năng và hoạt động của mã nguồn mở Lucene, sử dụng mã nguồn mở Lucene.NET để xây dựng thử nghiệm hệ thống tìm kiếm thông tin.

Chương 3: Trên cơ sở nghiên cứu về Hệ thống tìm kiếm thông tin và mã nguồn mở Lucene, chúng tôi đề xuất xây dựng thử nghiệm hệ thống tìm kiếm Văn bản với hai thành phần chính là: Tạo chỉ mục và Tìm kiếm.

Chương 4: Trình bày các kết quả đạt được, những hạn chế của luận văn và hướng phát triển cho hệ thống quản lý Văn bản ứng dụng công nghệ Lucene trong tương lai.

CHƯƠNG 1: NGHIÊN CỨU TỔNG QUAN VỀ HỆ THỐNG TÌM KIẾM THÔNG TIN

Với những hệ thống có số lượng lớn các tài liệu thì việc tra cứu, tìm kiếm thông tin thông thường chưa đáp ứng được nhu cầu tìm kiếm của người dùng. Hệ thống chủ yếu tìm kiếm một cách chính xác dựa trên tiêu đề của tài liệu, cơ sở dữ liệu tìm kiếm đơn giản, tốc độ tìm kiếm chậm, chưa chính xác và chưa hỗ trợ các phép toán tìm kiếm. Vậy đây chính là các vấn đề cần cải thiện để cải thiện cho các hệ thống tra cứu tìm kiếm thông tin.

1.1. Khái niệm về hệ thống tìm kiếm thông tin

Theo lý thuyết, hệ thống tìm kiếm thông tin là một hệ thống thông tin. Nó được sử dụng để lưu trữ, xử lý, tra cứu, tìm kiếm và phổ biến các yếu tố thông tin đến người sử dụng. Hệ thống tìm kiếm thông tin thường thao tác với các dữ liệu dạng văn bản và không có sự giới hạn về các yếu tố thông tin trong văn bản.

1.2. Các bộ phận cấu thành hệ thống tìm kiếm thông tin

1.2.1. Bộ thu thập thông tin

Bộ phận thu thập thông tin là một chương trình chạy tự động dùng để đi thu thập, lấy dữ liệu và lưu trữ các nội dung từ các trang web trên Internet. Bộ phận này có các thành phần chính: Một thành phần để theo dõi và phát hiện các URL mới, phát hiện các URL thay đổi. Một thành phần dùng để đọc đệ quy nội dung tài liệu của tất cả các trang web từ một tập các URL đã có, phân tích tài liệu, trích xuất nội dung tài liệu dưới các định dạng như html, pdf, excel...và lưu trữ về cơ sở dữ liệu thu thập.

1.2.2. Bộ lập chỉ mục

Hệ thống lập chỉ mục là để tối ưu hóa tốc độ và hiệu suất trong việc tìm kiếm các tài liệu có liên quan cho một truy vấn tìm kiếm. Nếu không có chỉ mục, công cụ tìm kiếm sẽ quét tất cả các tài liệu trong thư viện, đòi hỏi thời gian và sức mạnh tính toán đáng kể. Chẳng hạn, trong khi một chỉ mục 10.000 tài liệu có thể được truy vấn trong vòng mili giây thì việc quét theo từng phần của mỗi từ trong 10.000 tài liệu lớn có thể mất hàng giờ.

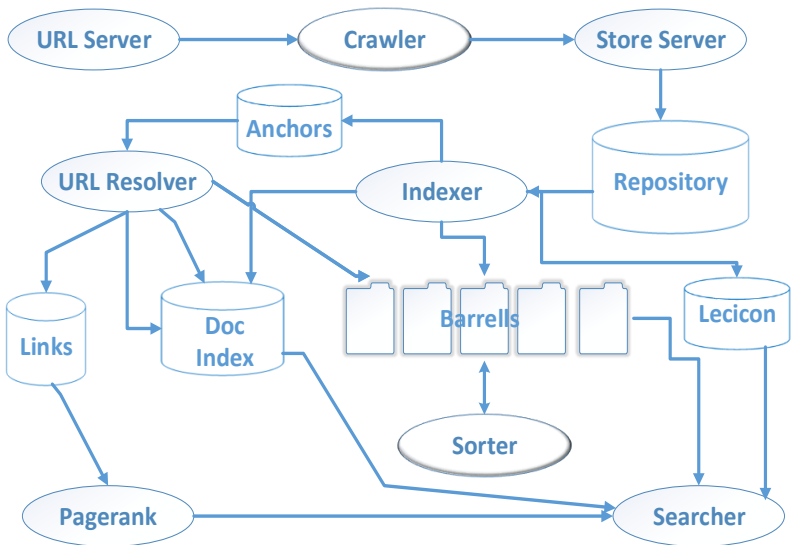
1.2.3. Bộ tìm kiếm thông tin

Bộ phận này chịu trách nhiệm tìm kiếm các tài liệu từ yêu cầu của người sử dụng, sau đó trả về danh sách các tài liệu chính xác với yêu cầu nhất. Do số lượng các trang web là rất lớn, và thông thường người dùng chỉ

đưa vào một vài từ khóa trong câu truy vấn nên tập kết quả thường rất lớn. Vì vậy bộ xếp hạng (ranking) có nhiệm vụ sắp xếp các tài liệu này theo mức độ hợp lệ với yêu cầu tìm kiếm và hiển thị kết quả cho người sử dụng.

1.3. Hệ thống tìm kiếm thông tin của Google

Google là một công ty Internet có trụ sở tại Hoa Kỳ, được thành lập vào năm 1998. Sản phẩm chính của công ty này là công cụ tìm kiếm Google, được nhiều người đánh giá là công cụ tìm kiếm hữu ích và mạnh mẽ nhất trên Internet. Trong khuôn khổ của đề tài, tôi đề xuất nghiên cứu mô hình tìm kiếm thông tin của Google để hiểu rõ hơn về kiến trúc của một Hệ thống tìm kiếm thông tin. Mô hình kiến trúc tổng thể của hệ thống tìm kiếm Google như sau:



Hình 1.3.1: Mô hình kiến trúc của hệ thống tìm kiếm Google [6]

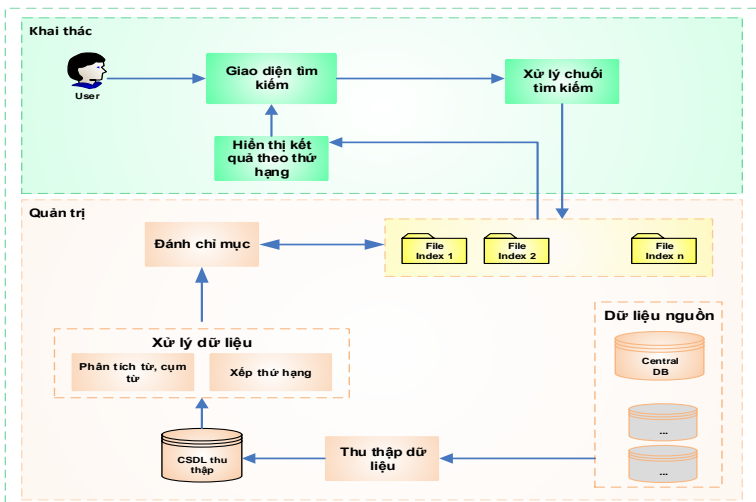
- Quy trình làm việc của hệ thống và chức năng của từng thành phần được mô tả như sau:

URL server gửi cho Crawler (được tổ chức phân tán, làm việc song song) một tập hợp các địa chỉ URLs. Các tài liệu (WebPages, hay Document) được Crawler tải xuống đưa vào Store Server, tại đây chúng được nén lại theo chuẩn Zlib (RFC 1950) và lưu trữ vào hệ thống lưu trữ tập trung Repository. Tại Repository, mỗi tài liệu

được gán cho một số number: DocID, Indexer đọc tài liệu từ Repository, giải nén và phân tích chúng. Tài liệu sau đó được chuyển đổi sang một tập các từ khóa xuất hiện bên trong nó gọi là Hits, mỗi hits là một bản ghi gồm: từ khóa, vị trí xuất hiện, font size, chữ hoa/thường. Indexer phân bổ các hits vào trong tập các kho chứa nhỏ hơn Barrels. Đồng thời nó cũng phân tích toàn bộ các đường link có trong mỗi trang và lưu trữ quan trọng vào AnchorsFile: text của link, link from, link to.

URLresolver đọc AnchorsFile rồi chuyển đổi đường dẫn tương đối về tuyệt đối và ánh xạ tương ứng các đường dẫn tuyệt đối này với DocIDs, sau đó thông tin này sẽ được đưa vào Barrels tương ứng theo DocID. Đồng thời cũng sản sinh Database link (lưu từng cặp DocIDs có mối liên kết với nhau). Sorter sắp xếp dữ liệu (hits) trong Barrels bởi DocID và sắp xếp lại bởi WordID để tạo ra Inverted Index (index nghịch đảo). Bộ phận từ điển Lexicon lấy danh sách WordID tạo ra mục từ mới. Searcher chạy bởi một WebServer sử dụng các từ điển (Lexicon) và thông tin index đảo (invert index) trong Barrels cùng với kết quả tính rank (từ PageRank) để trả về kết quả tìm kiếm.

1.4. Kiến trúc của hệ thống tìm kiếm thông tin



Hình 1.4.1.1: Mô hình kiến trúc hệ thống tìm kiếm thông tin

Dựa trên ý tưởng của Google và các hệ thống tìm kiếm thông tin khác chúng ta có thể hiểu về cơ bản một hệ thống tìm kiếm thông tin luôn có ba thành phần như sau:

- *Thành phần Thu thập dữ liệu:* thực hiện thu thập toàn bộ dữ liệu sẽ tìm kiếm đưa về một nguồn tập trung để phục vụ quá trình phân tích và đánh chỉ mục dữ liệu, thành phần này được quản lý bởi môđun thu thập dữ liệu.
- *Thành phần Đánh chỉ mục dữ liệu:* thực hiện phân tích, tiền xử lý nội dung dữ liệu, sau đó tiến hành đánh chỉ mục dữ liệu theo cách thức, cơ chế và yêu cầu của từng máy tìm kiếm cụ thể, thực hiện đánh chỉ mục dữ liệu này lưu vào các File index.
- *Thành phần Tìm kiếm dữ liệu:* thực hiện phân tích câu truy vấn và thực hiện tìm kiếm tài liệu trên các file index, sau đó kết hợp với thông tin xếp hạng (Rank) để trả lại kết quả tìm kiếm cho người dùng, thành phần này có một số chức năng chính như: Tiền xử lý khoá tìm kiếm, thực hiện phân tích từ khoá tìm kiếm, xử lý các toán tử tìm kiếm cơ bản (AND, OR, NOT,...), xử lý tìm kiếm chính xác, và xây dựng câu truy vấn dữ liệu.

Từ những nghiên cứu trên chúng ta có thể nhận thấy hệ thống tìm kiếm thông tin có những ưu điểm vượt trội hơn so với chức năng tìm kiếm trong cơ sở dữ liệu thông thường như: Hệ quản trị CSDL thông thường không thể đánh chỉ mục cho dữ liệu dạng file trong khi đó hệ thống tìm kiếm thông tin có thể đánh chỉ mục cho tất cả các tập tin dạng: pdf, html, MS Word, Excel,... Các câu truy vấn của các hệ quản trị CSDL bị giới hạn bởi cú pháp của SQL query, trong khi câu truy vấn của Hệ thống tìm kiếm gắn với yêu cầu tìm kiếm của người dùng, chúng ta có thể dùng các phép toán tìm kiếm AND, OR, NOT, tìm kiếm chính xác cụm từ, cụm từ... Ngoài ra với những dữ liệu lớn thì tốc độ tìm kiếm của Hệ thống tìm kiếm thông tin nhanh hơn nhiều so với chức năng tìm kiếm của các hệ Quản trị CSDL thông thường.

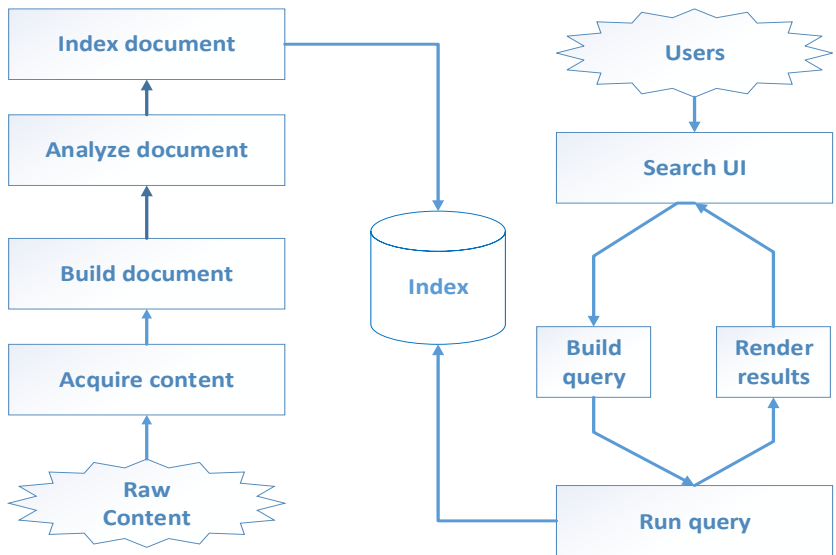
CHƯƠNG 2: NGHIÊN CỨU TỔNG QUAN VỀ MÃ NGUỒN MỞ LUCENE

Lucene là thư viện mã nguồn mở cho phép xử lý các văn bản đầu vào ở dạng văn bản (text) để tạo ra tập chỉ mục và cung cấp phương thức tìm kiếm trên tập chỉ mục đó. Nó cũng cho phép người dùng kế thừa và phát triển để phù hợp với nhiều ngôn ngữ khác nhau. Chúng tôi đề xuất nghiên cứu ứng dụng Lucene để phát triển hệ thống tìm kiếm trên các văn bản lưu trữ [2].

2.1. Giới thiệu về thư viện Lucene

Lucene là phần mềm mã nguồn mở, dùng để phân tích, đánh chỉ mục và tìm kiếm thông tin với hiệu suất cao bằng Java. Lucene được phát triển đầu tiên bởi Doug Cutting được giới thiệu đầu tiên vào tháng 8 năm 2000. Tháng 9 năm 2001 Lucene gia nhập vào tổ chức Apache và hiện tại được Apache phát triển và quản lý. Lucene không phải là một ứng dụng mà chỉ là một công cụ đặc tả API cần thiết cho việc xây dựng một search engine [10].

- **Lucene trong hệ thống tìm kiếm thông tin:**



Hình 2.1.1: Lucene trong hệ thống tìm kiếm thông tin [5]

Thành phần chức năng chính của Lucene bao gồm hai phần: Thành phần tạo chỉ mục và thành phần tìm kiếm. Đây là hai thành phần quan trọng cho một hệ thống tìm kiếm thông tin.

- **Thành phần Tạo chỉ mục:** Bao gồm các chức năng xử lý và phân tích dữ liệu để đánh chỉ mục. Lucene cho phép thiết lập các trường thông tin cần thiết để đánh chỉ mục phục vụ cho thành phần tìm kiếm, các thư viện phục vụ đánh chỉ mục mà Lucene hỗ trợ.
- **Thành phần Tìm kiếm:** bao gồm các phần chức năng xử lý tìm kiếm, trả về kết quả tìm kiếm cho người dùng, thông qua biên dịch và so khớp để lấy về kết quả tốt nhất.

2.2. Quy trình đánh chỉ mục

Để tiến hành đánh chỉ mục cho tài liệu, trước hết chúng ta phải chuyển đổi toàn bộ nội dung trong các file dữ liệu như HTML, PDF, MS WORD...sang các nội dung chỉ chứa dữ liệu dạng text. Lucene sẽ tiến hành phân tích và xử lý dữ liệu, loại bỏ những từ không có nghĩa, tách từ, cụm từ,...Sau khi dữ liệu được phân tích sẽ chuyển sang cho việc index. Lucene chứa dữ liệu này theo cấu trúc inverted index (chỉ mục có thể nghịch đảo). Cấu trúc này sẽ có hiệu quả để tiết kiệm dung lượng ổ đĩa và cho phép tìm kiếm nhanh hơn các từ khóa trong quá trình tìm kiếm. Nguyên tắc đó là thay vì phải tìm kiếm các từ nào chứa trong tài liệu đó thì với cấu trúc này sẽ tối ưu hóa việc tìm ra câu trả lời “tài liệu nào chứa từ khóa này”.

2.3. Các toán tử đánh chỉ mục cơ bản

Lucene hỗ trợ các toán tử giúp thực hiện việc đánh chỉ mục như: Thêm tài liệu mới (Document) cùng các trường (Fields): Keywords, UnIndexed, UnStored và Text. Trong mỗi tài liệu lại có thể có chứa nhiều Fields cùng tồn tại và trong mỗi Fields lại có thêm nhiều giá trị khác nhau. Xóa tài liệu ra khỏi chỉ mục (Remove Documents), lớp này sử dụng lớp IndexReader với phương thức delete() ta có thể dễ dàng xóa bỏ tài liệu được chỉ định ra khỏi chỉ mục. Lucene sẽ xem như các tài liệu này được đánh dấu như là sẽ xóa. Tuy nhiên việc này chỉ có thể thực hiện khi gọi hàm close(). Cập nhật tài liệu: Lucene không hỗ trợ thực hiện việc cập nhật tài liệu, thay vào đó sẽ xóa bỏ tài liệu và sau đó thêm lại tài liệu mới thay thế. Để đảm bảo tốc độ thực

thi thì tốt nhất việc xóa bỏ và thêm tài liệu mới nên thực hiện theo khối, không nên xen lẫn giữa việc xóa và thêm tài liệu mới.

2.4. Tối ưu hóa việc đánh chỉ mục

Việc tối ưu hóa tiến trình đánh chỉ mục là tiến trình trộn nhiều file chỉ mục lại với nhau để giảm thiểu thời gian đọc chỉ mục trong quá trình tìm kiếm. Bằng việc sử dụng API của Lucene mà cụ thể là hàm optimize() của đối tượng IndexWriter ta có thể dễ dàng tối ưu điều này. Tuy nhiên việc làm này chỉ có hiệu quả tăng tốc độ tìm kiếm trên chỉ mục đã có, mà không có tác động tới tốc độ đánh chỉ mục.

2.5. Tính đồng thời, an toàn tiến trình, ngăn chặn các thực thi

Các luật đồng thời: Lucene cung cấp cho người dùng nhiều toán tử liên quan tới việc đánh chỉ mục tài liệu như: xóa, cập nhật. Do đó trong quá trình thực hiện chúng ta phải tuân theo một số luật cụ thể để tránh việc đụng độ trong quá trình thực thi. Điều này là cần thiết khi mà có nhiều thực thi diễn ra một cách thường xuyên trước những yêu cầu gọi từ web tới ứng dụng của bạn. Sau đây là một số luật cơ bản: Bất kì toán tử chỉ đọc nào cũng có thể thực thi đồng thời, chẳng hạn là nhiều tiến trình có thể tìm kiếm cùng một chỉ mục tại một thời điểm. Bất kì toán tử chỉ đọc nào cũng có thể thực thi đồng thời trong khi một chỉ mục nào đó đang được cập nhật. Ví dụ: người dùng có thể tìm kiếm trong chỉ mục trong khi nó đang được cập nhật, thêm tài liệu mới hoặc là được xóa khỏi chỉ mục. Chỉ có duy nhất 1 toán tử cập nhật chỉ mục có thể thực thi tại một thời điểm. Một chỉ mục chỉ có thể được mở bởi chỉ một đối tượng IndexWriter hoặc là IndexReader tại một thời điểm mà thôi.

2.6. Bộ chuyển đổi câu truy vấn của người dùng: QueryParser

Hai yêu cầu quan trọng trong ứng dụng tìm kiếm đòi hỏi là: chuyển đổi câu truy vấn và truy xuất thông tin trả về. Hầu hết các phương thức Lucene đòi hỏi đối tượng Query. Việc chuyển đổi câu truy vấn là việc biểu diễn câu truy vấn của người dùng thành đối tượng Query phù hợp để sau đó truyền vào hàm tìm kiếm của lucene. Lucene có thể tìm ra kết quả chỉ khi câu truy vấn truyền vào là đúng định dạng của nó.

Để thực hiện được việc chuyển đổi câu truy vấn của người dùng, QueryParser cần thêm một đối tượng khác gọi là bộ phân tích Analyzer. Tùy

vào việc chọn lựa bộ Analyzer để phân tích chuỗi truyền vào thì kết quả sẽ khác nhau.

2.7. Các biểu thức truy vấn của QueryParser

QueryParser sử dụng nhiều toán tử luận lý để thực hiện việc chuyển đổi câu truy vấn như: OR, AND, NOT. Mặc định là OR. Chẳng hạn câu truy vấn sau: abc xyz thì sẽ được phân tích thành là abc or xyz or (abc and xyz). Để thay đổi tham số mặc định này, ta cần đặt lại toán tử cho đối tượng QueryParser.

2.8. Bộ phân tích – Analyzer:

Trong Lucene, phân tích (*analysis*) là quá trình chuyển đổi các field văn bản về dạng trình bày chỉ mục cơ bản nhất (*term*). Các terms thì được sử dụng để xác định rõ tài liệu nào sẽ phù hợp với một câu truy vấn trong quá trình tìm kiếm. Bộ phân tích (*analyzer*) là cách nói tóm lược quá trình phân tích. Analyzer phân tích trong đoạn văn bản thành *tokenizes*, đó là quá trình rút trích các từ, bỏ đi hệ thống các dấu chấm câu, chuyển toàn bộ các chữ trong văn bản về dạng chữ thường (*lowercasing* hay còn gọi là *normalizing*), loại bỏ các từ chung (*common words hay stop words*), giảm số lượng từ từ văn bản đưa vào (*root form* hay còn gọi *stemming*). Quá trình này còn được gọi là *tokenization*, chuyển đoạn văn bản thành nhiều khúc văn bản được gọi là các *token*. Tokens được kết hợp với các field name của chúng được gọi là *terms*. Sau quá trình tạo ra terms, terms sẽ là những khối dữ liệu được dùng để tìm kiếm trực tiếp. Vì vậy chọn bộ phân tích đúng đắn là cốt yếu quan trọng của quá trình phát triển phần mềm tìm kiếm. Ngôn ngữ là một yếu tố phải được nghĩ đến để chọn bộ phân tích, bởi vì đều có đặc trưng riêng và duy nhất của từng ngôn ngữ.

2.9. Sử dụng lớp IndexSearcher

Sau khi tạo ra đối tượng IndexSearcher, ta sẽ gọi phương thức search để thực hiện việc tìm kiếm. Có ba phương thức chính để tìm kiếm. Song ta chủ yếu sử dụng phương thức search(Query), tức tham số là câu truy vấn Query. Các phương thức tìm kiếm đều trả về là các Hits –chứa các thông tin đã tìm kiếm được, kết quả được sắp xếp theo thứ tự độ chính

xác. Thông qua đối tượng này ta có thể truy xuất thêm nhiều thông tin về kết quả tìm kiếm.

2.10. Cú pháp truy vấn Lucene

Lucene có một cú pháp truy vấn tùy chỉnh để truy vấn các chỉ mục của nó. Trong hầu hết các ứng dụng ta sử dụng đối tượng QueryParser để chuyển đổi câu truy vấn theo từng loại thích hợp. Lucene cung cấp bốn loại Query: QueryParse, BooleanQuery, RangeQuery và TermQuery. Sau đây ta sẽ tìm hiểu từng loại Query và lúc nào QueryParse sẽ chuyển đổi câu truy vấn thành dạng nào [7].

2.11. Các máy tìm kiếm phát triển dựa trên Lucene

- **Apache Solr:**

Solr là một máy chủ tìm kiếm mã nguồn mở phát triển dựa trên Apache Lucene có khả năng cung cấp các thư viện cho việc index (đánh chỉ mục) và search (tìm kiếm) dữ liệu. Solr nhập dữ liệu dưới dạng XML thông qua HTTP, hoặc sử dụng thư viện để nhập khối lượng lớn dữ liệu. Người dùng có thể truy vấn dữ liệu này thông qua HTTP GET và nhận về một kết quả dạng XML. Solr chạy bên trong một Java servlet container như Tomcat, Jetty hay Resin.

- **Elasticsearch:**

ElasticSearch là một máy tìm kiếm cấp doanh nghiệp (enterprise-level search engine). Mục tiêu của nó là tạo ra một công cụ, nền tảng hay kỹ thuật tìm kiếm và phân tích trong thời gian thực, có thể áp dụng hay triển khai một cách dễ dàng vào nguồn dữ liệu (data sources) khác nhau. ElasticSearch được phát triển bởi Shay Banon và dựa trên Apache Lucene, là một bản phân phối mã nguồn mở cho việc tìm kiếm dữ liệu trên máy chủ.

CHƯƠNG 3: XÂY DỰNG ỨNG DỤNG THỬ NGHIỆM

Trên cơ sở nghiên cứu, chúng tôi đề xuất xây dựng ứng dụng thử nghiệm Lucene vào trong hệ thống tìm kiếm Văn bản. Trong đó, chúng tôi kế thừa mã nguồn mở Lucene.NET để xây dựng hệ thống tìm kiếm với hai thành phần chính là: Tạo chỉ mục và Tìm kiếm. Do chỉ xây dựng hệ thống tìm kiếm nhỏ với dữ liệu văn bản ít nên hệ thống không có thành phần Thu thập dữ liệu .

3.1. Tài mã nguồn Lucene.NET

Trước tiên, chúng ta truy cập vào website <https://lucenenet.apache.org/> để tải phiên bản mới nhất của Lucene. Sau khi tải về giải nén vào thư mục làm việc mà không cần cài đặt. Sau khi giải nén, kiểm tra lại để đảm bảo trong thư mục vừa giải nén có chứa thư viện Lucene.net.dll. Thư viện này sẽ được tích hợp vào trong công cụ phát triển phần mềm Microsoft Visual Studio. Ứng dụng thử nghiệm xây dựng các giao diện tương tác với người dùng, mã nguồn của các giao diện này là các dòng lệnh gọi các phương thức do bộ thư viện Lucene.NET cung cấp để phục vụ việc đánh chỉ mục và tìm kiếm.

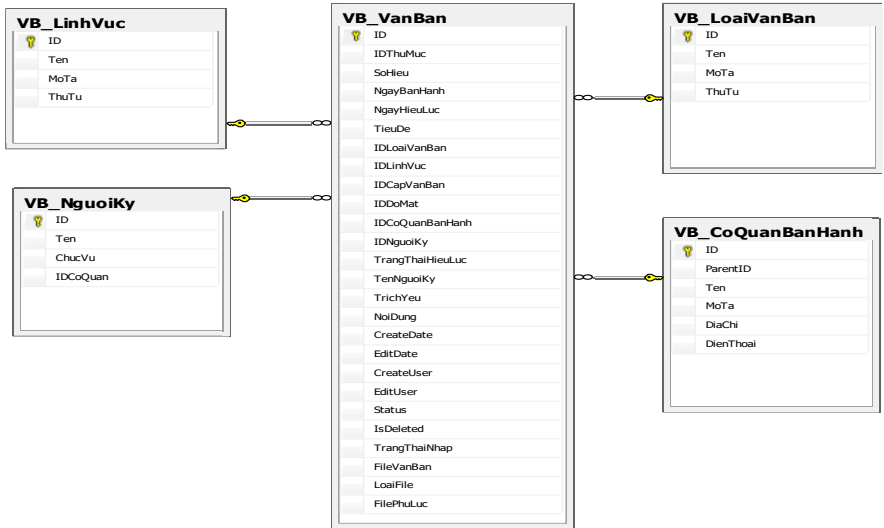
3.2. Dữ liệu văn bản thử nghiệm

Do không có thành phần thu thập dữ liệu nên dữ liệu phục vụ thử nghiệm sẽ được đưa vào hệ thống thông qua chức năng nhập dữ liệu văn bản của ứng dụng. Dữ liệu này được lưu trữ trong hệ quản trị CSDL SQL Server 2008. Sau khi có dữ liệu, chúng ta sử dụng Lucene để đánh chỉ mục văn bản phục vụ tìm kiếm.

3.3. Mô hình cơ sở dữ liệu

3.3.1. Lược đồ cơ sở dữ liệu

Lược đồ mô hình CSDL cung cấp một cách tổng quan về cấu trúc dữ liệu cũng như các mối quan hệ dữ liệu giữa các bảng trong CSDL. Mô tả chi tiết các trường dữ liệu, kiểu dữ liệu, khóa chính, khóa phụ trong các bảng. Đây là mô hình CSDL lưu trữ Văn bản phục vụ việc đánh chỉ mục và tìm kiếm của ứng dụng thử nghiệm.



Hình 3.3.1: Hình lược đồ cơ sở dữ liệu

3.4. Giao diện chính

3.4.1. Giao diện trang Quản lý văn bản

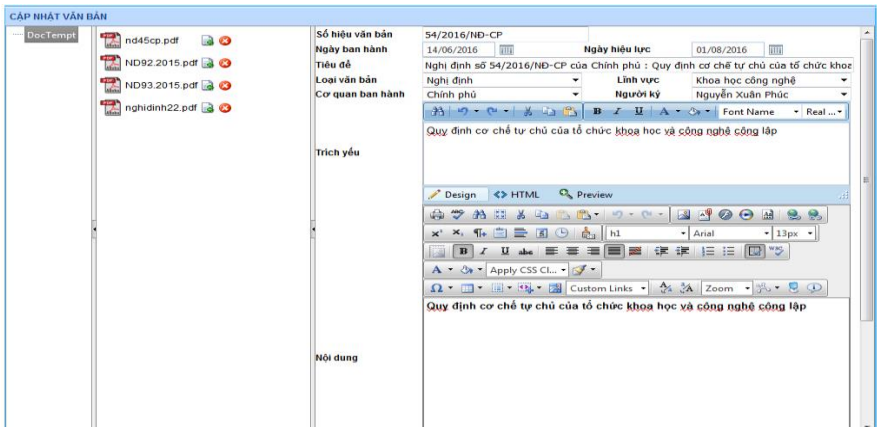
Giao diện này cho phép người dùng thêm mới, cập nhật hoặc xóa các văn bản đã có trong hệ thống. Đây chính là chức năng tạo nên Cơ sở dữ liệu văn bản phục vụ việc đánh chỉ mục để tạo ra cơ sở dữ liệu Lucene. Hệ thống cho phép đánh chỉ mục toàn bộ văn bản đã lưu trữ trong CSDL, hoặc chỉ đánh chỉ mục những văn bản vừa được thêm mới, cập nhật hoặc xóa mà không cần phải đánh chỉ mục lại từ đầu.



Hình 3.4.1: Giao diện trang Quản lý văn bản

3.4.2. Giao diện trang Cập nhật văn bản

Chức năng này cho phép người dùng thực hiện thêm mới hoặc cập nhật văn bản đã có trong hệ thống. Thông tin về văn bản bao gồm các trường thông tin như: Số hiệu văn bản, ngày ban hành, ngày hiệu lực, trích yếu, nội dung văn bản...



Hình 3.4.2: Giao diện trang Cập nhật văn bản

3.4.3. Giao diện trang Tìm kiếm văn bản

Giao diện này cho phép người dùng gõ các từ khóa tìm kiếm, hệ thống thực hiện tìm kiếm các từ khóa trong file chỉ mục, sắp xếp các kết quả và trả về danh sách các kết quả theo mức độ liên quan giữa câu truy vấn và

các tài liệu trong cơ sở dữ liệu chỉ mục. Chức năng này cho phép người dùng gõ trực tiếp từ khóa kết hợp với các phép toán mà Lucene hỗ trợ (AND, OR, NOT...), tìm kiếm chính xác, tìm kiếm gần đúng hoặc tìm với các ký tự đại diện.

TÌM KIẾM VĂN BẢN							VĂN BẢN MỚI NHẤT	
Từ khóa: <input type="text" value="quyet dinh"/> <input type="button" value="Tìm kiếm"/> Tìm nâng cao ▾							Số hiệu Trích yếu	
Tìm thấy 65 kết quả nào cho từ quyet dinh (0.02 giây)							844/QĐ-BTNMT	Thành lập ban soạn thảo và Tổ biên tập Nghị định sửa đổi, bổ sung các nghị định hướng dẫn thi hành Luật Bảo vệ môi trường năm 2014
STT	Số hiệu	Trích yếu	Ngày ban hành	Ngày hiệu lực	Tải về	Xem	789/QĐ-BTNMT	Ủy quyền phê duyệt các đề án điều tra cơ bản địa chất về khoáng sản sử dụng nguồn vốn góp của tổ chức, cá nhân: Ông Đỗ Cảnh Dương
1	1249/QĐ-BTNMT	Quyết định số 1249/QĐ-BTNMT về việc công nhận sáng kiến năm 2015	01/06/2016	01/06/2016			685/QĐ-BTNMT	Về việc tăng Kỳ niệm chương "Vì sự nghiệp Tài nguyên và Môi trường"
2	25/QĐ-BTNMT	Sửa đổi Điều 4 Quyết định số 1758/QĐ-BTNMT quy định chức năng, nhiệm vụ quyền hạn và cơ cấu tổ chức Cục Đo đạc và Bản đồ Việt Nam	06/01/2017	06/01/2017			680/QĐ-BTNMT	Thành lập Hội đồng Tư vấn chính sách tài nguyên và môi trường
3	2999/QĐ-BTNMT	Về việc Sửa đổi Điều 2 Quyết định số 1188/QĐ-BTNMT ngày 26/5/2016 của Bộ trưởng Bộ Tài nguyên và Môi trường	26/12/2016	26/12/2016			673/QĐ-BTNMT	Công bố Danh mục VBQPPL hết hiệu lực toàn bộ hoặc một phần thuộc lĩnh vực quản lý nhà nước của Bộ TN&MT đến ngày 30/01/2017
4	3110/QĐ-BTNMT	Quyết định về việc công bố 10 sự kiện nổi bật ngành Tài nguyên và Môi trường năm 2016	30/12/2016	30/12/2016			674/QĐ-BTNMT	Kiểm toán Ban chỉ đạo ứng dụng và phát triển công nghệ thông tin ngành TN&MT
5	2052/QĐ-BTNMT	Quyết định về việc tổ chức Liên hoan phim Môi trường toàn quốc lần thứ 6	07/09/2016	07/09/2016			552/QĐ-BTNMT	Thành lập Ban chỉ đạo sơ kết 5 năm thực hiện Nghị quyết số 19-NQ/TW ngày 30/10/2012 của Ban Chấp hành Trung ương Khóa XI
6	173/QĐ-BTNMT	Quyết định tặng Kỳ niệm chương Vì sự nghiệp Tài nguyên và Môi trường cho ông Chang Jae Yun, Giám đốc quốc gia của Cơ quan Hợp tác quốc tế Hàn Quốc (KOICA) tại Việt Nam.	14/02/2017	14/02/2017			1480/QĐ-BTNMT	Đánh giá việc tổ chức thi hành Luật Đất đai năm 2013, đề xuất nội dung sửa đổi, bổ sung Luật đất đai và các văn bản quy định chi tiết thi hành
7	1189/QĐ-BTNMT	Ban hành Chương trình hành động của Bộ TN&MT thực hiện Nghị quyết 19-2016/NQ về những nhiệm vụ, giải pháp chủ yếu cải thiện môi trường kinh doanh, nâng cao năng lực cạnh tranh quốc gia hai năm 2016- 2017, định hướng đến năm 2020	26/05/2016	26/05/2016			550/QĐ-BTNMT	Tăng Kỳ niệm chương "Vì sự nghiệp tài nguyên và môi trường"
8	25/2016/TT----	Căn cứ Luật Ban hành văn bản quy phạm pháp luật ngày 22 tháng 6 năm 2015; Căn cứ Nghị định số 34/2016/NĐ-CP ngày 14 tháng	15/12/2016	15/12/2016			1412/BTNMT-KH	Hướng dẫn xây dựng kế hoạch và dự toán kinh phí sự nghiệp BVMT năm 2018 của các Bộ, ngành
							464/QĐ-BTNMT	Ban hành Sổ tay hướng dẫn thực hiện dự án "Tăng cường quản lý đất đai và cơ sở dữ liệu đất đai" (Dự án VILG)
							446/QĐ-	Tăng Kỳ niệm chương "Vì sự nghiệp tài nguyên và môi trường"

Hình 3.4.3: Giao diện trang Tìm kiếm văn bản

3.4.4. Giao diện trang Xem nội dung file văn bản

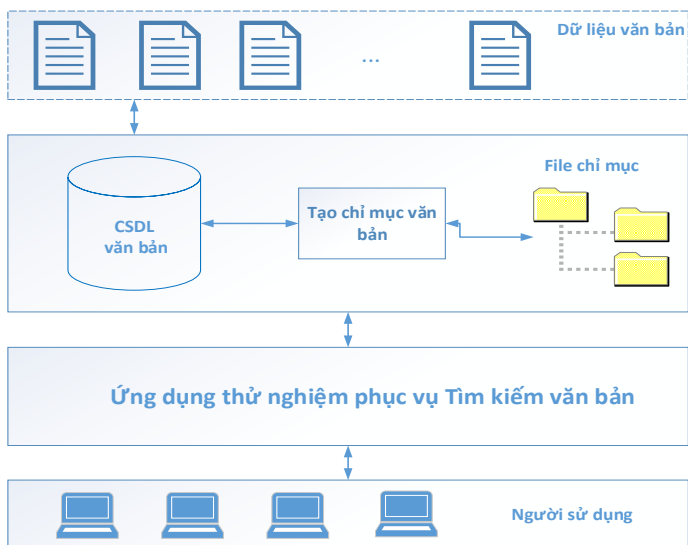
Với các văn bản đã được đưa vào hệ thống, ngoài việc xem thông tin chi tiết văn bản thì người dùng có thể xem trực tiếp nội dung của văn bản đã được lưu trữ trong file pdf. Đây chính là nội dung toàn văn bản cần được đưa vào hệ thống để đánh chỉ mục phục vụ chức năng tra cứu, tìm kiếm văn bản.



Hình 3.4.5: Giao diện trang Xem nội dung file văn bản

3.5. Đánh giá và thử nghiệm

3.5.1. Mô hình kiến trúc ứng dụng thử nghiệm



Hình 3.5.1: Hình Kiến trúc ứng dụng thử nghiệm

Hệ thống thử nghiệm có 3 thành phần chính:

- **Dữ liệu văn bản thử nghiệm:** Tất cả các văn bản được đưa vào hệ thống và lưu trữ trong hệ Quản trị CSDL SQL Server 2008 gọi là cơ sở dữ liệu Văn bản. Với mô hình như trên chúng tôi đã thử nghiệm với số lượng khoảng hơn 300 văn bản.
- **Tạo chỉ mục văn bản:** Từ những dữ liệu văn bản đã được đưa vào hệ thống, chúng tôi xây dựng chức năng đánh chỉ mục cho văn bản, chức năng này sẽ tạo ra các file chỉ mục lưu trữ trên đĩa cứng máy tính, đây chính là cơ sở dữ liệu chỉ mục hay cơ sở dữ liệu Lucene.
- **Tìm kiếm văn bản:** Người dùng truy cập vào đường dẫn phần mềm thử nghiệm và gõ từ khóa tìm kiếm văn bản. Hệ thống sẽ thực hiện tìm kiếm ở trong cơ sở dữ liệu Lucene hay tìm trong các file chỉ mục và trả về kết quả cho người dùng. Từ danh sách kết quả tìm kiếm người dùng có thể sắp xếp tăng dần, giảm dần theo số hiệu văn bản, ngày ban hành, ngày hiệu lực. Hoặc click vào một văn bản để xem thông tin chi tiết.

3.5.2. Kịch bản và kết quả

Với dữ liệu văn bản đã được đánh chỉ mục, chúng tôi đã thử nghiệm sử dụng chức năng tìm kiếm văn bản với các kịch bản như sau: Tìm kiếm theo từ khóa bất kỳ, tìm kiếm theo toán tử, tìm kiếm chính, gần đúng từ khóa,...Hệ thống đã thực hiện thành công trên máy tính cá nhân với cấu hình Intel Core i5-3210M, CPU 2.5GHz, 8GB RAM với kết quả như sau:

Bảng 3.5.2.1: Bảng Kịch bản tìm kiếm của Hệ thống tìm kiếm thông thường

Kịch bản tìm kiếm	Từ khóa tìm kiếm	Số lượng kết quả	Thời gian tìm kiếm (giây)
Tìm kiếm theo từ khóa bất kỳ	quy	67	1.25
Tìm kiếm theo từ khóa bất kỳ	Ban hành	34	1.76

Tìm kiếm theo phép toán AND	Môi trường AND	0	0.84
Tìm kiếm theo phép toán OR	Môi trường OR	0	0.74
Tìm kiếm chính xác từ khóa	“quy chế”	0	1.026
Tìm kiếm từ khóa là tiếng việt không có dấu (tìm kiếm gần đúng)	Quy che	0	0.311
Tìm kiếm với ký tự đại diện	BGD*	0	0.24

Bảng 3.5.2.2: Bảng Kịch bản tìm kiếm của Hệ thống tìm kiếm thông tin

Kịch bản tìm kiếm	Từ khóa tìm kiếm	Số lượng kết quả	Thời gian tìm kiếm (giây)
Tìm kiếm theo từ khóa bất kỳ	quy	45	0.027
Tìm kiếm theo từ khóa bất kỳ	Ban hành	76	0.016
Tìm kiếm theo phép toán AND	Môi trường AND	117	0.041
Tìm kiếm theo phép toán OR	Môi trường OR	130	0.023
Tìm kiếm chính xác từ khóa	“quy chế”	17	0.012
Tìm kiếm từ khóa là tiếng việt không có dấu (tìm kiếm gần đúng)	Quy che	48	0.017
Tìm kiếm với ký tự đại diện	BGD*	4	0.029

Qua kết quả thống kê trên cho thấy chức năng tìm kiếm của hệ thống cơ sở dữ liệu thông thường không hỗ trợ các toán tử tìm kiếm, không hỗ trợ tìm kiếm chính xác trong dấu “” và tìm kiếm tiếng việt không có dấu. Đặc biệt tốc độ tìm kiếm của hệ thống tìm kiếm thông tin nhanh hơn rất nhiều so với chức năng tìm kiếm của hệ quản trị cơ sở dữ liệu thông thường. Kết quả thử nghiệm tìm kiếm trên hệ thống tìm kiếm thông tin tương đối chính xác và hiệu quả, hỗ trợ đầy đủ các phép toán tìm kiếm với thời gian tìm kiếm rất nhanh. Như vậy hệ thống thử nghiệm tìm kiếm văn bản sử dụng mã nguồn mở Lucene đáp ứng được mục tiêu đặt ra của đề tài.

CHƯƠNG 4: KẾT LUẬN

Về lý thuyết, luận văn tìm hiểu về các thành phần cơ bản của một hệ thống tìm kiếm thông tin, tìm hiểu một cách hệ thống các tính năng và hoạt động của mã nguồn mở Lucene.

Thực nghiệm, từ cơ sở lý thuyết, luận văn đã xây dựng và cài đặt thành công ứng dụng thực nghiệm Lucene vào trong hệ thống tìm kiếm Văn bản. Trong đó, ứng dụng bộ thư viện mã nguồn mở Lucene.NET để xây dựng hệ thống tìm kiếm với hai thành phần chính là: Tạo chỉ mục và Tìm kiếm văn bản.

4.1. Đánh giá kết quả nghiên cứu

4.1.1. Kết quả đạt được:

Về cơ bản luận văn đã thực hiện tốt các nội dung đề ra và đạt được một số kết quả nhất định: Luận văn đã trình bày cơ sở lý thuyết và nguyên lý vận hành của một hệ thống tìm kiếm thông tin, trình bày một cách hệ thống các tính năng và hoạt động của mã nguồn mở Lucene. Luận văn đã ứng dụng thành công mã nguồn mở Lucene trong công tác tìm kiếm thông tin trên hệ thống Quản lý văn bản.

4.1.2. Hạn chế:

Bên cạnh những kết quả đạt được thì đề tài còn có những mặt hạn chế như: Phần thực nghiệm mới chỉ dừng lại ở phạm vi nhỏ với số lượng văn bản còn hạn chế. Đề tài chưa nghiên cứu được các bộ thư viện mã nguồn mở khác giống thư viện Lucene, chưa nghiên cứu một số tính năng liên quan đến tìm kiếm tiếng Việt.

4.2. Hướng phát triển

Hướng nghiên cứu tiếp theo của đề tài là tập trung nghiên cứu tìm hiểu các mã nguồn mở khác giống thư viện mã nguồn mở Lucene, để có thể áp dụng thử nghiệm và đưa ra được những nhận xét, đánh giá so sánh hiệu quả với thư viện Lucene.

TÀI LIỆU THAM KHẢO

Tiếng Việt

1. Đỗ Phúc, Đỗ Hoàng Cường, Nguyễn Tri Tuấn, Huỳnh Thụy Bảo Trân, Nguyễn Văn Khiết, Nguyễn Việt Hoàng, Nguyễn Việt Thành, Phạm Phú Hội, Dương Ngọc Long Nam, Nguyễn Phước Thanh Hải, “Phát triển một Hệ thống S.E” Hỗ trợ Tìm kiếm Thông tin, thuộc lãnh vực CNTT trên Internet qua từ khóa bằng tiếng Việt”, Đại học Khoa Học Tự Nhiên, TP.HCM, 2004
2. Huỳnh Đức Việt, Võ Duy Thanh, Võ Trung Hùng, ”Tập chí khoa học và công nghệ”, Đại học Đà Nẵng, 2010

Tiếng Anh

3. Michael McCandless, Erik Hatcher, Otis Gospodnetic, Lucene in action, 2010
4. Haralambos Marmanis and Dmitry Babenko, Algorithms of the Intelligent Web, 2009
5. Chris Manning and Pandu Nayak, Introduction to Information Retrieval
6. <http://infolab.stanford.edu/~backrub/google.html>
7. <http://www.lucenetutorial.com>
8. <https://www.tutorialspoint.com>
9. <https://lucenet.apache.org/>
10. <https://en.wikipedia.org/wiki/Lucene>