# CS242: Information Retrieval & Web Search

Renjie Wu (862058526), Zhihui Shao (861257080), Tong Shen (862058291)

# Part A: Collect your data and Index with Lucene

## 1. Collaboration

**a) Coding**

Renjie Wu: Overall design, implementation of class *"CrawlThread"* and class *"IndexThread"*.
Zhihui Shao: Implementation of class "*WriterThread*" and class "*RobotPolicy*".
Tong Shen: Implementation of class "*WikiCrawler*" and class "*Indexer*".

**b) Report**

Report is written by all of our three team members cooperatively.

## 2. Overview of the crawling system

### 2.1 Architecture

The crawling system generally consists of three components: crawler, writer, and database. These three components work in a pipeline: the crawler visits Wikipedia and crawls needed pages; the writer then writes those crawled pages into the database by batch; the database finally saves all necessary information associated with those pages and ensures the eventual consistency. Two intermediate layers are introduced to make those components work properly in a pipeline: a *LinkedBlockingQueue*, between the crawler and the writer, to play the queue's role in the producer-consumer model; independent database connections (SQL connection) built between the writer and the database, to perform writing transactions simultaneously. The basic architecture of our crawling system is demonstrated in Fig. 1.
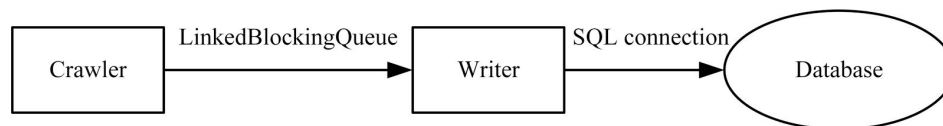


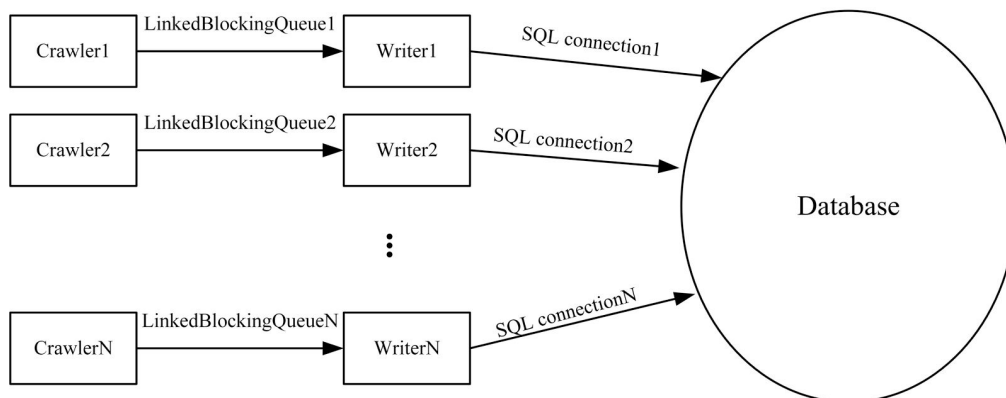Figure 1 Basic architecture of the crawling system



Figure 2 Full structure of the crawling system

To improve performance and efficiency, both the crawler and the writer would create multiple threads (*CrawlThread* for the crawler and *WriterThread* for the writer), to separate the workload and execute the crawling and the writing process parallelly. Each *CrawlThread* would work together with a *WriterThread*; they are connected by the LinkedBlockingQueue, where the *CrawlThread* pushes the crawled pages into the queue, and the *WriterThread* then removes the queue's head page and insert it into the database, through a dedicated database connection, associated with the *WriterThread*. The full structure of the crawling system is shown in Fig. 2.

## 2.2 Crawling Strategy

**a) Strategy: *CrawlThread*'s routine**
We adopted the general search algorithm in the *CrawlThread*. The default page depth limit is set to 10, since deeper pages may have already been crawled before. The following pseudo-code shows how *CrawlThread* works.

```
Start associated WriterThread;
Generate seed;
Add the random URL to Queue;
while Count<PagesToCrawl do
        if (Queue is empty) do
                Generate seed;
                Add the random URL to Queue;
        Get URL from Queue;
        Request the URL and download the web page;
        Add URL to visited set;
        Process content of the web page;
        Save processed web page into the LinkedBlockingQueue;
        if (depth of this web page exceeds limit) do
                continue;
        else do
                Get all links from the processed web page
                for URL in links do
                        if URL not in visited set and URL within Wikipedia do
                                Add URL to Queue;
    end
```

**b) Strategy: Seed generating and handling duplicate pages**
A set of seeds or a set of initial URLs of Wikipedia pages are required to start the crawler. We set the entry url as *https://en.wikipedia.org/wiki/Special:Random*, since it would randomly redirect to a page, and that is ideal for seeding. The redirected page would act as the seed for each *CrawlThread*. In addition, to respect the crawler ethics, the crawler would first check website's robots.txt and determine whether the next url to crawl is allowed, and the *CrawlThread* would sleep for 5 seconds after finished processing a given page.

Duplicate pages handling is done through a set called *visitedURL*, saving all URLs of crawled pages and it is shared among all *CrawlThread*s, to check if a given URL has already been

cralwed. *visitedURL* is implemented by *ConcurrentHashMap*, to ensure thread-safety. Page storage is done through SQLite, which provides support of eventual consistency and parallel writing, to improve the throughput of our writer.


**c) Strategy: Wikipedia article processing**
One of the important parts here is how to process the web page, and all the page process is coded in method *process()* in *CrawlThread* class. For almost all Wiki web page, there has one unique title; the main content, which is made of many paragraphs; and some words represent the category of this page. So, for every web page, we first split it into these three parts. And this is done by following three methods provided by jsoup.

> *Element elTitle = doc.getElementById("firstHeading");*
> *Element elContent = doc.selectFirst("#mw-content-text .mw-parser-output");*
> *Element elCategory = doc.getElementById("mw-normal-catlink");*

**i) Title**
As the title is unique for all Wiki web pages, the title is used for the primary key when saving to SQLite database. Using title as primary key has two advantages than using URL: one is that primary key is shorter than URL; the other is that URL can be recovered from title when necessary.
The recovery process can be this:
$$URL = "\text{https://en.wikipedia.org/wiki/}" + title$$
For instance, the URL of web page with the title Computer is:
https://en.wikipedia.org/wiki/Computer
For title with more than one work, we need change the space to "_" in title. And then follows the recovery process. For instance, the URL of of web page with the title Computer science is:
https://en.wikipedia.org/wiki/Computer_science

**ii) Content**
For the main content, we can simply convert the element to string using *text()* and *trim()* method.
*String title = elTitle.text().trim();*
But it will generate references, edit links, spaces, \r\n, tables, and captions that is irrelevant the main content. In order to make the index part more efficiency, we need remove these irrelevant parts from the main content. As is shown in the Figure 3, after processing, we obtain the clean and plain content of the webpage.

To remove all reference, we use "*elContent.select("sup[class='reference']").remove()*"
To remove the `edit` links, we use "*elContent.select("span[class='mw-editsection']").remove()*"
To remove unused tags (such as table & div), we use "*Arrays.asList("table", "div").forEach(tag -> elContent.select(tag).remove())*". To remove empty headings with no paragraphs below it, we use "*Arrays.asList("h1", "h2", "h3", "h4", "h5", "h6").forEach( tag -> elContent.select(tag + "+" + tag).stream().map(Element::previousElementSibling).forEach(Element::remove));*"

(a) Raw content of edit and unused tag part    (b) Content after process of edit and unused tag part

(c) Raw content of reference part    (d) Content after process of reference part

Figure 3 Main content process example

### iii) Category

For the words in category, we convert it to *List<String>*. And then saved into database.

### d) Strategy: Storing articles into database

In our crawling system, one crawl thread will have one write thread to write the web page generated into the SQLite database. In order to improve the efficiency and throughput, the write thread use batch. Every 50 records from one batch, and they are written at one time. In order to concurrent the crawl thread and write thread, the page queue from crawl thread to write thread is blocked queue.

In our SQLite database, there are five attributes: *title*, *content*, *categories*, *lastModify*, *outLinks*. So, in the lucene index part, three of them are used for index process: title, content, and categories. Although the other two fields, *lastModify* and *outLinks*, is not used in index process, they can be used in query part. Here is a quick view of the result of our database for crawler:

| | title | content | categories | lastModify | outLinks |
|---|---|---|---|---|---|
| 30 | Sanskrit | Sanskrit (IAST: Saṃskṛtam; IPA: [sɑ̃… | Sanskrit\|Indo-Aryan langu… | 2018-02-06 21:07 | IAST\|Sacred language\|Hinduism\|Sikhism\|Buddhism\| |
| 31 | Persian language | Persian (/ˈpɜːrʒən/ or /ˈpɜːrʃən/), … | Persian language\|Language… | 2018-02-09 11:57 | Endonym\|Western Iranian languages\|Indo-Iranian |
| 32 | Bengali literature | Bengali literature (Bengali: বাংলা সাহ… | Bengali-language literatu… | 2018-02-06 16:20 | Bengali language\|Charyapada\|Mangalkavya\|Syed Su |
| 33 | History of Bengali lite… | Ancient Age◌Charyapada◌The first wo… | Bengali-language literatu… | 2018-01-12 14:10 | Bengali language\|Buddhism\|Luipa\|Bengali people\| |
| 34 | Bangladeshi folk litera… | Bangladeshi Folk Literature (Bengali… | Bengali-language literatu… | 2017-05-19 07:59 | Bengali language\|Bengali literature\|Folklore\|Ba |
| 35 | Bengali renaissance | The Bengali renaissance or simply Be… | Bengal Renaissance\|Renais… | 2018-01-23 19:32 | Cultural movement\|Social movement\|Intellectual |
| 36 | Amar Sonar Bangla | Amar Sonar Bangla (Bengali: আমার স… | Asian anthems\|Bangladeshi… | 2018-01-30 13:15 | Bengali language\|National anthem\|Bangladesh\|Ode |
| 37 | Jana Gana Mana | "Jana Gana Mana" (Bengali: [ɟɐnɐ ɡɐɳ… | Asian anthems\|Bengali-lan… | 2018-02-09 19:07 | National anthem\|India\|Bengali language\|Rabindra |
| 38 | Sri Lanka Matha | Sri Lanka Matha (Sinhalese: ශ්‍රී ල… | Asian anthems\|1940 songs\|… | 2018-02-03 23:51 | Sinhalese language\|Tamil language\|National anthe |
| 39 | Rabindranath Tagore | Rabindranath Tagore FRAS (/rəˈbɪndrə… | Rabindranath Tagore\|1861 … | 2018-01-28 02:38 | Fellow of the Royal Asiatic Society\|Sobriquet\|Be |
| 40 | Sinhalese language | Sinhalese (/ˌsɪnəˈliːz/), known nativ… | Southern Indo-Aryan langu… | 2018-02-03 01:50 | Sinhalese people\|Sri Lanka\|Indo-Aryan languages |
| 41 | Language Movement | The Language Movement (Bengali: ভাষা … | History of Bangladesh\|His… | 2018-02-09 03:35 | Bengali language\|East Bengal\|Bangladesh\|Official |
| 42 | Dominion of Pakistan | Pakistan (Bengali: পাকিস্তান অধিরাজ্য p… | Former countries in South… | 2018-01-28 17:37 | Bengali language\|Urdu language\|Dominion\|Pakistan |
| 43 | UNESCO | The United Nations Educational, Scie… | UNESCO\|Organizations esta… | 2018-02-06 17:33 | French language\|List of specialized agencies of |
| 44 | Language Movement Day | Language Movement Day or Language Re… | 1952 protests\|February ob… | 2017-12-26 10:07 | Bengali language\|Bangladesh\|Bengali Language Mov |
| 45 | International Mother La… | International Mother Language Day (I… | February observances\|Inte… | 2018-02-06 20:26 | Linguistic diversity\|Cultural diversity\|Multilin |
| 46 | Bengali nationalism | Bengali nationalism (Bengali: বাংলা জা… | Political movements in Ba… | 2018-01-29 06:51 | Bengali language\|Constitution of Bangladesh\|Nati |
| 47 | Culture of Bengal | The culture of Bengal encompasses th… | Bengali culture\|Banglades… | 2018-02-06 00:53 | Bengal\|South Asia\|Bangladesh\|India\|West Bengal\| |
| 48 | 1st millennium BC | The 1st millennium BC encompasses th… | 1st millennium BC\|Millenn… | 2018-02-03 17:07 | Before Christ\|Iron Age\|1000 BC\|1 BC\|Neo-Assyrian |
| 49 | Gupta Empire | The Gupta Empire was an ancient Indi… | Former monarchies of Asia… | 2018-02-09 16:09 | Outline of ancient India\|Indian subcontinent\|Gu |
| 50 | Vedic and Sanskrit lite… | Vedic and Sanskrit literature compri… | Indian literature by lang… | 2017-12-14 06:07 | Oral literature\|Vedas\|Indian epic poetry\|Iron A |

Figure 4  Result in SQLite database

# 3. Overview of the Lucene indexing strategy

When building Lucene index, we also run multi-threaded batch jobs to improve the performance of our system. In order to maximize the performance, all of the threads share one connection to database. When acquiring raw content, we use batch read to fetch records from our database. Intuitively, we chose the batch factor as 50 to make a balance between reducing I/O and not consuming too much main memory.
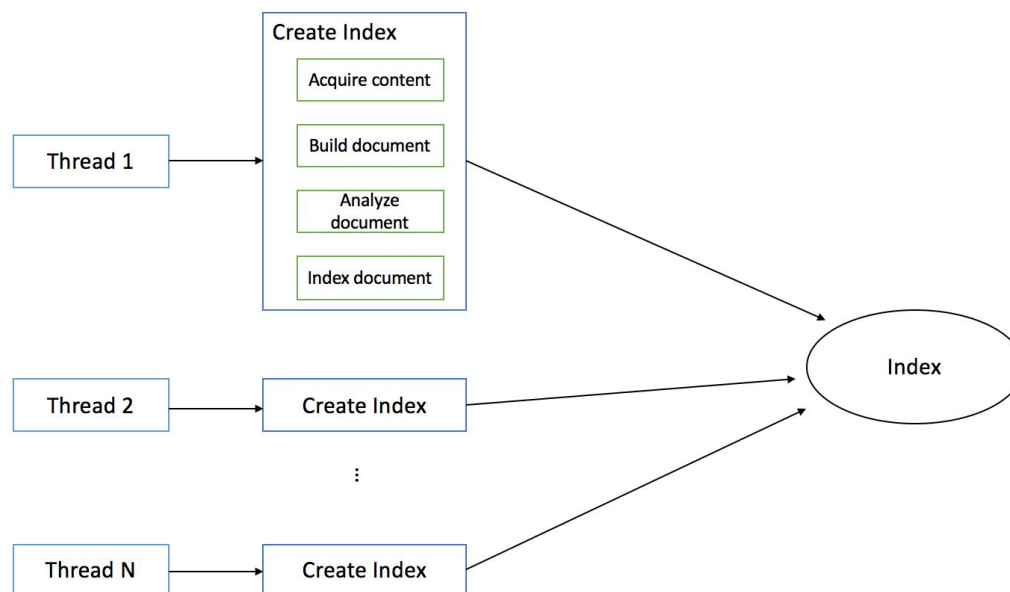


Figure 5 Architecture of Lucene Indexer

**i) Fields in the Lucene indexing**

We make one webpage as a document. To make it easier to maintain, we make each document has the three fields: title, content and categories. Since we want to search title, content and categories in the future, we build indexes for all of them, but in order to save some space, only title will be stored in the index and if user need to view the content or categories, the system will use the title as the key to fetch them from our database.

**ii) Text analyzer choice**
As for the analyzer part, since we already remove most of the irrelevant part and get the plain text while crawling the webpage as discussed before, the analyzer part is pretty simply and straightforward. We use a separate analyzer to analyze categories since when stored in the database we use the character "|" to separate different categories so that the categories field has one more special stop world "|". After that, we just use standardAnalyzer to analyze the documents.

**iii) Progress reporting**
We report the process when indexing every 1000 pages. As is shown in Figure 6, we finished indexing totally 744723 pages with in 3 hours and 20 minutes. Approximately, the running time is linear to the processed page number.
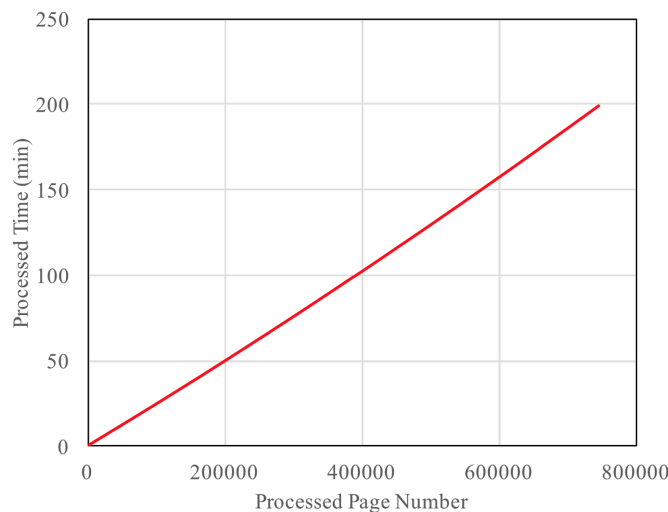


Figure 6 Run time of the Lucene index creation process

# 4. Obstacles and solutions

One big obstacle we encounter during coding is how to get only the main content of the web page. In Wiki page, there are many other parts other than the main content, such as the picture, table, edit hyperlink, reference link etc. If we do not remove these part from our main content text. There will be many irrelevant words and chars. So we spent a lot of time on this part. And it is handled using the method mentioned before.

# 5. Limitation and Achievement

There is one limitation in our crawling system: the consistency and corresponding efficiency. Although the *visitedURL* is used to avoid to search duplicate web page, the sequential consistency is still not ensured, because the *visitedURL* is not saved persistent with the database. Think about such case: Firstly, crawl the Wiki for some time, and stop; After some time, re-run the Crawler; when re-run the Crawler, the *visitedURL* is empty, and it is not consistent with the database.

To solve this, we ensure that this crawling system is eventually consistency by checking the page queue before write it database. This is done by the SQL statement: "INSERT **OR IGNORE** INTO pages". Therefore, every Wiki page in the database will be unique. But this solution is not very good. Because when we re-run the Crawler, it may crawler some pages that already saved in database before this run, and this reduce the efficiency of our crawling system.

In this project, we made the following achievements:
1. Both crawler and indexer are implemented with multi-threads and batch read & write database.
2. Respect the crawler ethic while crawling webpage, obeying robots.txt etc.
3. Wrap the code gracefully and make all parameters option.
4. Use database to store data, ensuring the eventual consistency and making reading and writing data more efficiently.
5. Crawler performs data cleaning, by removing unnecessary tags and saving paragraphs only.

# 6. Instruction

a) how to run the crawler: *./crawler.sh [options]*

Table 1 Options for crawler

| Short Name | Argument Name | Default Value |
|---|---|---|
| *t* | *numOfThreads* | 10 |
| *c* | *numOfPages* | (this value will ensure 5GB data)750000 |
| *d* | *crawlDepth* | 10 |
| *i* | *crawlInterval* | 500 |
| *u* | *entryUrl* | "https://en.wikipedia.org/wiki/Special:Random" |
| *H* | *crawlHostRegex* | "^en.wikipedia.org$" |
| *P* | *crawlPathRegex* | "^/wiki/[^:]*$" |
| | ***jdbcUrl*** | **Required parameter no default value** |
| *l* | *FILE NAME* | STDOUT |

b) how to build the Lucene index: *./indexbuilder.sh [options]*

Table 2 Options for indexer

| Short Name | Argument Name | Default Value |
|---|---|---|
| *t* | *numOfThreads* | 10 |
| *l* | *FILE NAME* | STDOUT |
| | ***jdbcUrl*** | **Required, no default value** |
| | ***indexOutputPath*** | **Required, no default value** |