

Architecture of Enterprise Applications 7

Full-text Searching

Haopeng Chen

REliable, INtelligent and Scalable Systems Group (REINS)

Shanghai Jiao Tong University

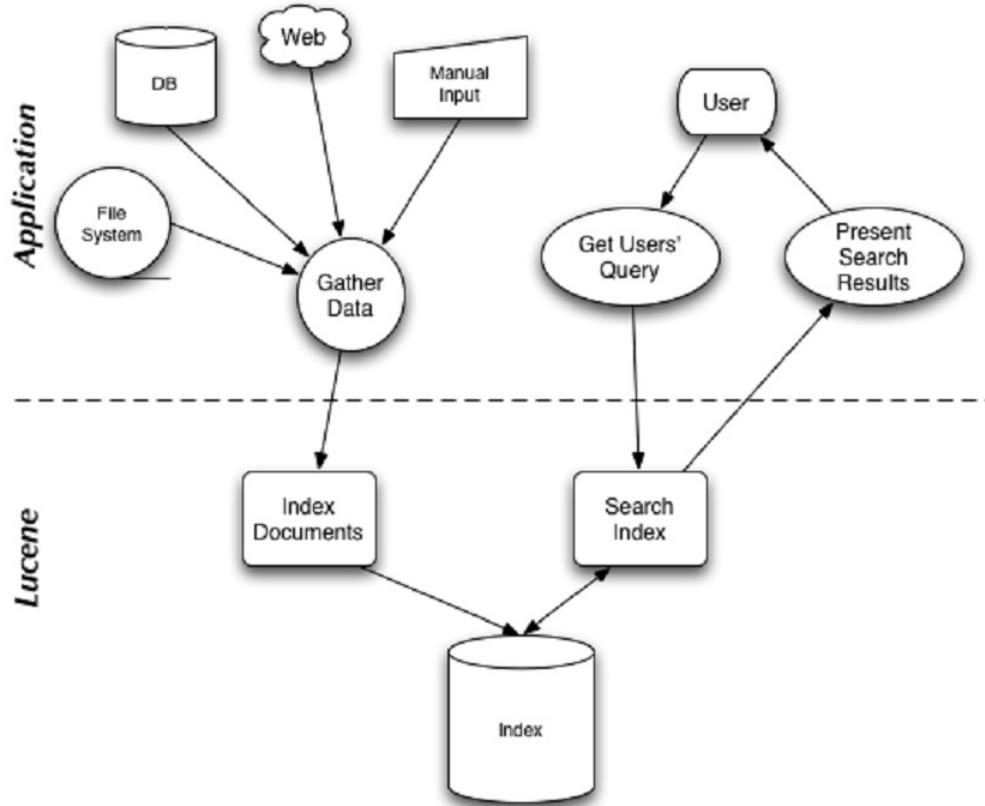
Shanghai, China

<http://reins.se.sjtu.edu.cn/~chenhp>

e-mail: chen-hp@sjtu.edu.cn

- **Contents**
 - Lucene
 - Solr
 - Elasticsearch
- **Objectives**
 - 能够根据业务需求，识别适合全文搜索引擎的场景，设计并实现基于搜索引擎的全文搜索方案

- Lucene is a high performance, scalable Information Retrieval (IR) library.
 - It lets you add indexing and searching capabilities to your applications.
 - Lucene is a mature, free, open-source project implemented in Java.
 - it's a member of the popular Apache Jakarta family of projects, licensed under the liberal Apache Software License.
- Lucene provides a simple yet powerful core API
 - that requires minimal understanding of full-text indexing and searching.



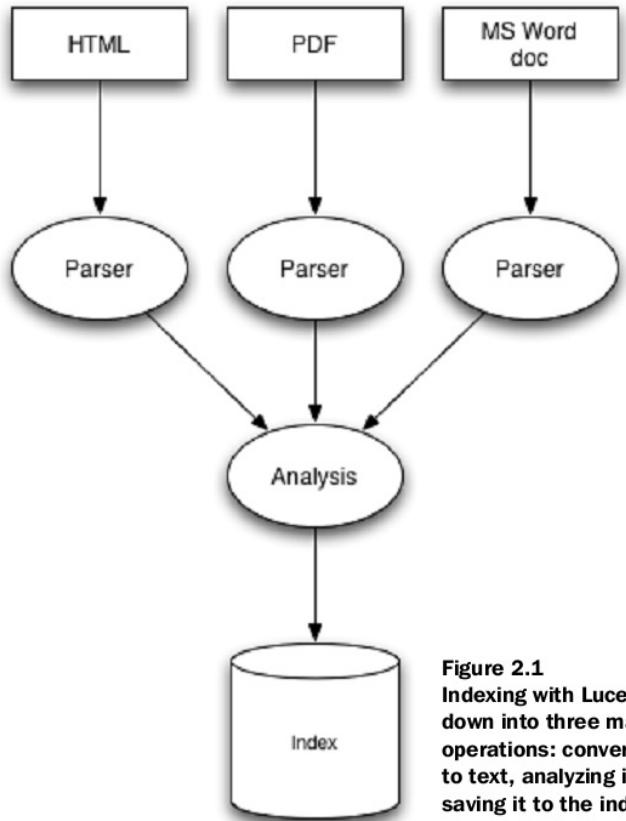
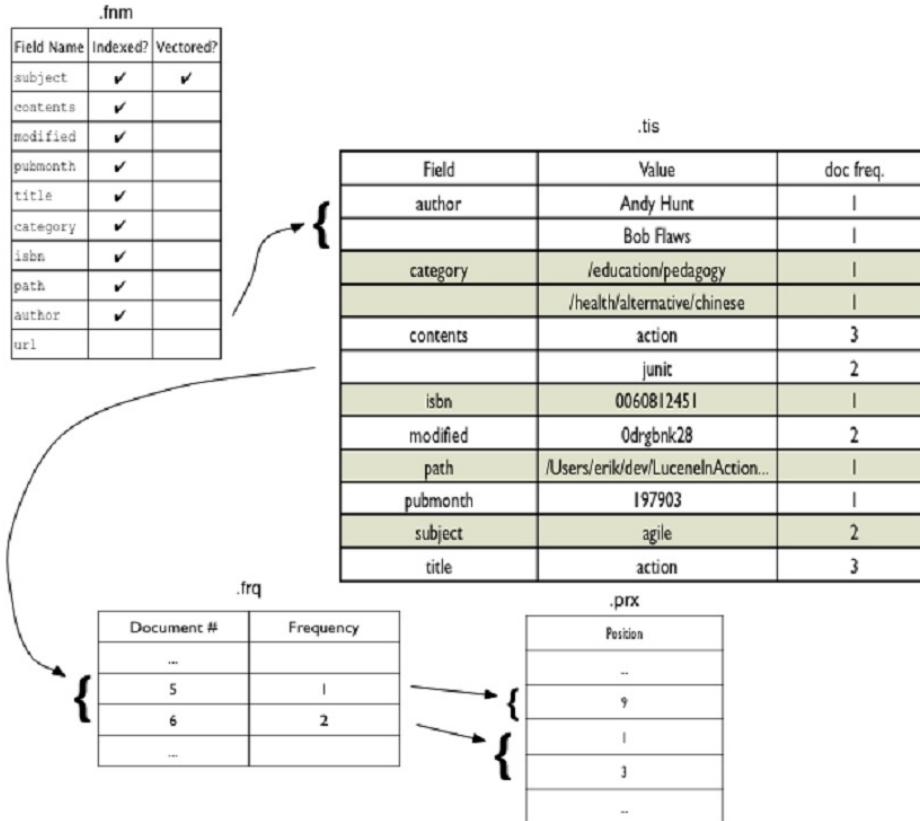


Figure 2.1
Indexing with Lucene breaks down into three main operations: converting data to text, analyzing it, and saving it to the index.

- At the heart of all search engines is the concept of indexing:
 - processing the original data into a **highly efficient cross-reference lookup** in order to facilitate rapid searching.
- Suppose you needed to search a large number of files, and you wanted to be able to find files that contained a certain word or a phrase
 - A naïve approach would be to sequentially scan each file for the given word or phrase.
 - This approach has a number of flaws, the most obvious of which is that it doesn't scale to larger file sets or cases where files are very large.

- This is where indexing comes in:
 - To search large amounts of text quickly, you must first index that text and convert it into a format that will let you search it rapidly, eliminating the slow sequential scanning process.
 - This conversion process is called indexing, and its output is called an index.
 - You can think of an index as a data structure that allows fast random access to words stored inside it.

Inverting index



- Searching is the process of looking up words in an index to find documents where they appear.
- The quality of a search is typically described using precision and recall metrics.
 - Recall measures how well the search system finds relevant documents, whereas precision measures how well the system filters out the irrelevant documents.
- A number of other factors
 - speed and the ability to quickly search large quantities of text.
 - Support for single and multi term queries, phrase queries, wildcards, result ranking, and sorting are also important, as is a friendly syntax for entering those queries.

A sample application

- Suppose you need to index and search files stored in a directory tree, not just in a single directory
- These example applications will familiarize you with Lucene's API, its ease of use, and its power.
- The code listings are complete, ready-to-use command-line programs.

Creating an Index

```
/**  
 * This code was originally written for  
 * Erik's Lucene intro java.net article  
 */  
public class Indexer {  
  
    public static void main(String[] args) throws Exception {  
        if (args.length != 2) {  
            throw new Exception("Usage: java " + Indexer.class.getName()  
                + " <index dir> <data dir>");  
        }  
        File indexDir = new File(args[0]);  
        File dataDir = new File(args[1]);  
  
        long start = new Date().getTime();  
        int numIndexed = index(indexDir, dataDir);  
        long end = new Date().getTime();  
  
        System.out.println("Indexing " + numIndexed + " files took " + (end - start) + " milliseconds");  
    }  
}
```



Create Lucene index in this directory
Index files in this directory

Creating an Index

```
// open an index and start file directory traversal
public static int index(File indexDir, File dataDir) throws IOException {
    if (!dataDir.exists() || !dataDir.isDirectory()) {
        throw new IOException(dataDir
            + " does not exist or is not a directory");
    }

    IndexWriter writer = new IndexWriter(indexDir,
        new StandardAnalyzer(), true);
    writer.setUseCompoundFile(false);                                ← Create Lucene index
    indexDirectory(writer, dataDir);

    int numIndexed = writer.docCount();
    writer.optimize();
    writer.close();                                                 ← Close index
    return numIndexed;
}
```

Creating an Index

```
// recursive method that calls itself when it finds a directory
private static void indexDirectory(IndexWriter writer, File dir)
throws IOException {

    File[] files = dir.listFiles();

    for (int i = 0; i < files.length; i++) {
        File f = files[i];
        if (f.isDirectory()) {
            indexDirectory(writer, f);           ← recurse
        } else if (f.getName().endsWith(".txt")) {
            indexFile(writer, f);             ← Index .txt files only
        }
    }
}
```

Creating an Index

```
// method to actually index a file using Lucene
private static void indexFile(IndexWriter writer, File f)
throws IOException {

if (f.isHidden() || !f.exists() || !f.canRead()) {
    return;
}

System.out.println("Indexing " + f.getCanonicalPath());

Document doc = new Document();
doc.add(Field.Text("contents", new FileReader(f)));           ← Index file content

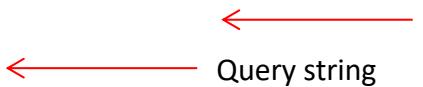
doc.add(Field.Keyword("filename", f.getCanonicalPath()));      ← Index file name
writer.addDocument(doc);                                     ← Add document to Lucene index
}
}
```

Running Indexer

```
% java lia.meetlucene.Indexer build/index/lucene
Indexing /lucene/build/test/TestDoc/test.txt
Indexing /lucene/build/test/TestDoc/test2.txt
Indexing /lucene/BUILD.txt
Indexing /lucene/CHANGES.txt
Indexing /lucene/LICENSE.txt
Indexing /lucene/README.txt
Indexing /lucene/src/jsp/README.txt
Indexing /lucene/src/test/org/apache/lucene/analysis/ru/stemsUnicode.txt
Indexing /lucene/src/test/org/apache/lucene/analysis/ru/test1251.txt
Indexing /lucene/src/test/org/apache/lucene/analysis/ru/testKOI8.txt
Indexing /lucene/src/test/org/apache/lucene/analysis/ru/testUnicode.txt
Indexing /lucene/src/test/org/apache/lucene/analysis/ru/wordsUnicode.txt
Indexing /lucene/todo.txt
Indexing 13 files took 2205 milliseconds
```

Searching an index

```
/**  
 * This code was originally written for  
 * Erik's Lucene intro java.net article  
 */  
public class Searcher {  
  
    public static void main(String[] args) throws Exception {  
        if (args.length != 2) {  
            throw new Exception("Usage: java " + Searcher.class.getName()  
                + " <index dir> <query>");  
        }  
        File indexDir = new File(args[0]);  
        String q = args[1];  
        if (!indexDir.exists() || !indexDir.isDirectory()) {  
            throw new Exception(indexDir +  
                " does not exist or is not a directory.");  
        }  
        search(indexDir, q);  
    }  
}
```



← Index directory created by Indexer
← Query string

Searching an index

```
public static void search(File indexDir, String q)
throws Exception {
Directory fsDir = FSDirectory.getDirectory(indexDir, false);
IndexSearcher is = new IndexSearcher(fsDir);           ← Open Index

Query query = QueryParser.parse(q, "contents", new StandardAnalyzer());    ← Parse query
long start = new Date().getTime();
Hits hits = is.search(query);          ← Search Index
long end = new Date().getTime();

System.err.println("Found " + hits.length() + " document(s) (in " + (end - start) +
" milliseconds) that matched query '" + q + "'");   ← Write search stats

for (int i = 0; i < hits.length(); i++) {           ← Retrieve matching document
    Document doc = hits.doc(i);
    System.out.println(doc.get("filename"));           ← Display filename
}
}
```

Running Searcher

```
%java lia.meetlucene.Searcher build/index 'lucene'  
Found 6 document(s) (in 66 milliseconds) that matched query 'lucene':  
/lucene/README.txt  
/lucene/src/jsp/README.txt  
/lucene/BUILD.txt  
/lucene/todo.txt  
/lucene/LICENSE.txt  
/lucene/CHANGES.txt
```

- **IndexWriter**
 - This class creates a new index and adds documents to an existing index.
- **Directory**
 - The Directory class represents the location of a Lucene index.
- **Analyzer**
 - The Analyzer, specified in the IndexWriter constructor, is in charge of extracting tokens out of text to be indexed and eliminating the rest.
- **Document**
 - A Document represents a collection of fields.
- **Field**
 - Each field corresponds to a piece of data that is either queried against or retrieved from the index during search.

- **IndexSearcher**
 - IndexSearcher is to searching what IndexWriter is to indexing
- **Term**
 - A Term is the basic unit for searching.
- **Query**
 - Query is the common, abstract parent class. It contains several utility methods
- **TermQuery**
 - TermQuery is the most basic type of query supported by Lucene, and it's one of the primitive query types.
- **Hits**
 - The Hits class is a simple container of pointers to ranked search results

Adding documents to an index

```
public abstract class BaseIndexingTestCase extends TestCase {  
    protected String[] keywords = {"1", "2"};  
    protected String[] unindexed = {"Netherlands", "Italy"};  
    protected String[] unstored = {"Amsterdam has lots of bridges",  
                                  "Venice has lots of canals"};  
    protected String[] text = {"Amsterdam", "Venice"};  
    protected Directory dir;  
  
    protected void setUp() throws IOException {  
        String indexDir =  
            System.getProperty("java.io.tmpdir", "tmp") +  
            System.getProperty("file.separator") + "index-dir";  
        dir = FSDirectory.getDirectory(indexDir, true);  
        addDocuments(dir);  
    }  
}
```

Adding documents to an index

```
protected void addDocuments(Directory dir) throws IOException {  
    IndexWriter writer = new IndexWriter(dir, getAnalyzer(), true);  
    writer.setUseCompoundFile(isCompound());  
    for (int i = 0; i < keywords.length; i++) {  
        Document doc = new Document();  
        doc.add(Field.Keyword("id", keywords[i]));  
        doc.add(Field.UnIndexed("country", unindexed[i]));  
        doc.add(Field.UnStored("contents", unstored[i]));  
        doc.add(Field.Text("city", text[i]));  
        writer.addDocument(doc);  
    }  
    writer.optimize();  
    writer.close();  
}  
protected Analyzer getAnalyzer() { return new SimpleAnalyzer();}  
protected boolean isCompound() { return true; }  
}
```

- All fields consist of a name and value pair.
 - **Keyword**—Isn't analyzed, but is indexed and stored in the index verbatim.
 - **UnIndexed**—Is neither analyzed nor indexed, but its value is stored in the index as is.
 - **UnStored**—The opposite of UnIndexed. This field type is analyzed and indexed but isn't stored in the index.
 - **Text**—Is analyzed, and is indexed. This implies that fields of this type can be searched against, but be cautious about the field size.

- **Heterogeneous Documents**

- One handy feature of Lucene is that it allows Documents with different sets of Fields to coexist in the same index.
- This means you can use a single index to hold Documents that represent different entities.
- For instance, you could have Documents that represent retail products with Fields such as **name** and **price**, and Documents that represent people with Fields such as **name, age, and gender**.

- **Appendable Fields**
- Suppose you have an application that generates an array of synonyms for a given word, and you want to use Lucene to index the base word plus all its synonyms.
- like this:

```
String baseWord = "fast";
String synonyms[] = String {"quick", "rapid", "speedy"};
Document doc = new Document();
doc.add(Field.Text("word", baseWord));
for (int i = 0; i < synonyms.length; i++) {
    doc.add(Field.Text("word", synonyms[i]));
}
```

- Internally, Lucene appends all the words together and index them in a single Field called **word**, allowing you to use any of the given words when searching.

Removing Documents from an index

```
public class DocumentDeleteTest extends BaseIndexingTestCase {  
    public void testDeleteBeforeIndexMerge() throws IOException {  
        IndexReader reader = IndexReader.open(dir);  
        assertEquals(2, reader.maxDoc());  
        assertEquals(2, reader.numDocs());  
        reader.delete(1);  
        assertTrue(reader.isDeleted(1));  
        assertTrue(reader.hasDeletions());  
        assertEquals(2, reader.maxDoc());  
        assertEquals(1, reader.numDocs());  
        reader.close();  
        reader = IndexReader.open(dir);  
        assertEquals(2, reader.maxDoc());  
        assertEquals(1, reader.numDocs());  
        reader.close();  
    }  
}
```

Removing Documents from an index

```
public void testDeleteAfterIndexMerge() throws IOException {  
    IndexReader reader = IndexReader.open(dir);  
    assertEquals(2, reader.maxDoc());  
    assertEquals(2, reader.numDocs());  
    reader.delete(1);  
    reader.close();  
    IndexWriter writer = new IndexWriter(dir, getAnalyzer(), false);  
    writer.optimize();  
    writer.close();  
    reader = IndexReader.open(dir);  
    assertFalse(reader.isDeleted(1));  
    assertFalse(reader.hasDeletions());  
    assertEquals(1, reader.maxDoc());  
    assertEquals(1, reader.numDocs());  
    reader.close();  
}
```

Undeleting Documents

- Because Document deletion is deferred until the closing of the IndexReader instance,
 - Lucene allows an application to change its mind and undelete Documents that have been marked as deleted.
- A call to IndexReader's `undeleteAll()` method undeltes all deleted Documents
 - by removing all `.del files` from the index directory.
- Subsequently closing the IndexReader instance therefore leaves all Documents in the index.
 - Documents can be undeleted only if the call to `undeleteAll()` was done using the same instance of IndexReader that was used to delete the Documents in the first place.

Updating Documents in an index

- “How do I update a document in an index?”
 - is a frequently asked question on the Lucene user mailing list.
- Lucene doesn’t offer an update(Document)method;
 - instead, a Document must first be deleted from an index and then re-added to it

```
IndexReader reader = IndexReader.open(dir);
reader.delete(new Term("city", "Amsterdam"));
reader.close();
```

```
IndexWriter writer = new IndexWriter(dir, getAnalyzer(), false);
Document doc = new Document();
doc.add(Field.Keyword("id", "1"));
doc.add(Field.UnIndexed("country", "Netherlands"));
doc.add(Field.UnStored("contents", "Amsterdam has lots of bridges"));
doc.add(Field.Text("city", "Haag"));
writer.addDocument(doc);
writer.optimize();
writer.close();
```

- Not all Documents and Fields are created equal
 - Document boosting is a feature that makes such a requirement simple to implement.
 - By default, all Documents have no boost—or, rather, they all have the same boost factor of 1.0.
 - By changing a Document's boost factor, you can instruct Lucene to consider it more or less important with respect to other Documents in the index.
 - The API for doing this consists of a single method, `setBoost(float)`

Boosting Documents and Fields

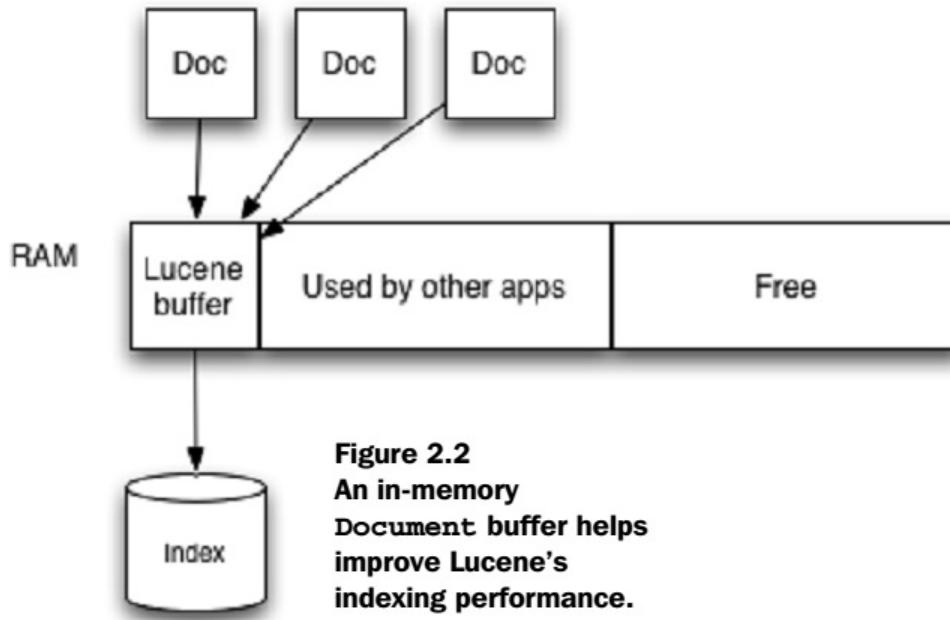
```
public static final String COMPANY_DOMAIN = "example.com";
public static final String BAD_DOMAIN = "yucky-domain.com";

Document doc = new Document();
String senderEmail = getSenderEmail();
String senderName = getSenderName();
String subject = getSubject();
String body =getBody();

doc.add(Field.Keyword("senderEmail", senderEmail));
doc.add(Field.Text("senderName", senderName));
doc.add(Field.Text("subject", subject));
doc.add(Field.UnStored("body", body));

if (getSenderDomain().endsWithIgnoreCase(COMPANY_DOMAIN)) {
    doc.setBoost(1.5);
}
else if (getSenderDomain().endsWithIgnoreCase(BAD_DOMAIN)) {
    doc.setBoost(0.1);
}
writer.addDocument(doc);
```

Tuning indexing performance



Parallelizing indexing by working with multiple indexes

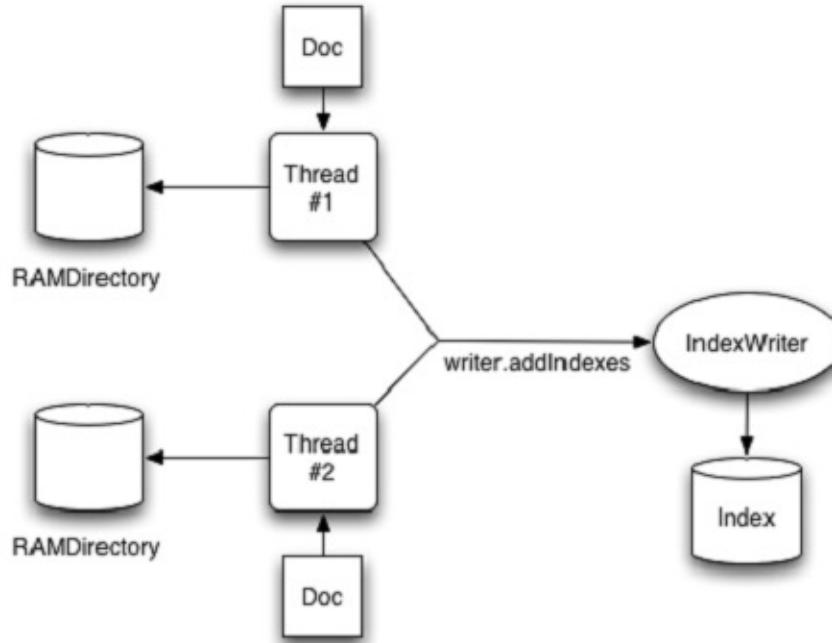


Figure 2.3 A multithreaded application that uses multiple `RAMDirectory` instances for parallel indexing.

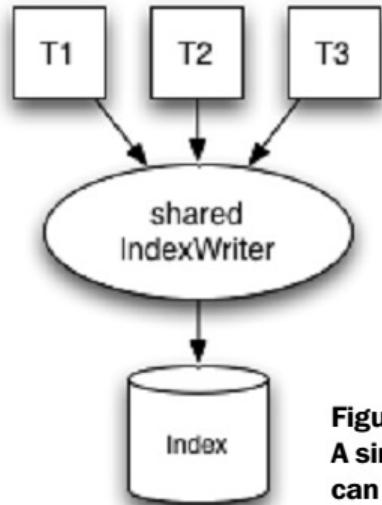


Figure 2.7
A single `IndexWriter` or `IndexReader`
can be shared by multiple threads.

Searching for a specific term

- A term is a value that is paired with its containing field.

```
IndexSearcher searcher = new IndexSearcher(directory);
```

```
Term t = new Term("subject", "ant");
```

```
Query query = new TermQuery(t);
```

```
Hits hits = searcher.search(query);
```

```
t = new Term("subject", "junit");
```

```
hits = searcher.search(new TermQuery(t));
```

```
searcher.close();
```

Parsing a user-entered query expression: QueryParser

```
IndexSearcher searcher = new IndexSearcher(directory);
Query query = QueryParser.parse("+JUNIT +ANT -MOCK", "contents",
                                new SimpleAnalyzer());
Hits hits = searcher.search(query);
```

```
Document d = hits.doc(0);
```

```
query = QueryParser.parse("mock OR junit", "contents",
                           new SimpleAnalyzer());
hits = searcher.search(query);
```

Understanding Lucene scoring

- The score is computed for each document (d) matching a specific.
 - This score is the raw score.
 - Scores returned from Hits aren't necessarily the raw score, however.
 - If the top-scoring document scores greater than 1.0, all scores are normalized from that score, such that all scores from Hits are guaranteed to be 1.0 or less.

Understanding Lucene scoring

$$\sum_{t \text{ in } q} tf(t \text{ in } d) \cdot idf(t) \cdot boost(t.\text{field in } d) \cdot lengthNorm(t.\text{field in } d)$$

Table 3.5 Factors in the scoring formula

Factor	Description
$tf(t \text{ in } d)$	Term frequency factor for the term (t) in the document (d).
$idf(t)$	Inverse document frequency of the term.
$boost(t.\text{field in } d)$	Field boost, as set during indexing.
$lengthNorm(t.\text{field in } d)$	Normalization value of a field, given the number of terms within the field. This value is computed during indexing and stored in the index.
$coord(q, d)$	Coordination factor, based on the number of query terms the document contains.
$queryNorm(q)$	Normalization value for a query, given the sum of the squared weights of each of the query terms.

- **APACHE SOLR™**

- is the popular, blazing-fast, open source enterprise search platform built on Apache Lucene™.
- is highly reliable, scalable and fault tolerant, providing distributed indexing, replication and load-balanced querying, automated failover and recovery, centralized configuration and more.
- Solr powers the search and navigation features of many of the world's largest internet sites.



- Launch Solr in SolrCloud Mode
 - `$./bin/solr start -e cloud`

The screenshot shows the Solr Admin interface with the following details:

- Request-Handler (qt):** /select
- Common Query Parameters:** q: *:
fq:
sort:
start, rows: 0, 10
fl:
df:
- Raw Query Parameters:** key1=val1&key2=val2
- WT Options:** -----,
- Checkboxes:** debugQuery, dismax, edismax, hl, facet, spatial, spellcheck
- Buttons:**
- Result Preview:** The right panel shows the JSON response from the Solr server, listing several book documents with their details like id, category, name, price, author, etc.

- SolrJ
 - is an API that makes it easy for applications written in Java (or any language based on the JVM) to talk to Solr.
- For projects built with Maven, place the following in your `pom.xml`:

```
<dependency>
  <groupId>org.apache.solr</groupId>
  <artifactId>solr-solrj</artifactId>
  <version>8.4.1</version>
</dependency>
```

Using SolrJ

- Querying.java

```
public class Querying {  
    public static void main(String[] args) throws IOException, SolrServerException {  
        final SolrClient client = getSolrClient();  
  
        final Map<String, String> queryParamMap = new HashMap<String, String>();  
        queryParamMap.put("q", "*:*");  
        queryParamMap.put("fl", "id, name");  
        queryParamMap.put("sort", "id asc");  
        MapSolrParams queryParams = new MapSolrParams(queryParamMap);  
  
        final QueryResponse response = client.query("techproducts", queryParams);  
        final SolrDocumentList documents = response.getResults();  
  
        System.out.println("Found " + documents.getNumFound() + " documents");  
        for (SolrDocument document : documents) {  
            final String id = (String) document.getFirstValue("id");  
            final String name = (String) document.getFirstValue("name");  
  
            System.out.println("id: " + id + "; name: " + name);  
        }  
    }  
}
```

Using SolrJ

- Querying.java

```
public static SolrClient getSolrClient() {  
    final String solrUrl = "http://localhost:8983/solr";  
    return new HttpSolrClient.Builder(solrUrl)  
        .withConnectionTimeout(10000)  
        .withSocketTimeout(60000)  
        .build();  
}  
  
}
```

Using SolrJ

The screenshot shows the Solr admin interface for the 'techproducts' collection. On the left, there's a sidebar with various navigation links like Dashboard, Logging, Cloud, Collections, Java Properties, Thread Dump, Suggestions, and a dropdown for 'techproducts'. The 'Query' link under 'techproducts' is currently selected. The main area has a 'Request-Handler (qt)' dropdown set to '/select'. Below it are several input fields: 'common' (with 'fq' containing '*:*' highlighted with a red box), 'sort' (set to 'id asc'), 'start, rows' (set to '0' and '10'), and 'fl' (set to 'id, name'). Under 'Raw Query Parameters', there's a key-value pair 'key1=val1&key2=val2'. The 'wt' field is set to '-----'. There are also checkboxes for 'indent off', 'debugQuery', 'dismax', 'edismax', 'hl', 'facet', 'spatial', and 'spellcheck'. At the bottom is a blue 'Execute Query' button.

The right side of the interface shows the URL `http://localhost:8983/solr/techproducts/select?fl=id%2Cname&q=%3A*&sort=id%20asc` and the raw JSON response:

```
{
  "responseHeader": {
    "zkConnected": true,
    "status": 0,
    "qTime": 12,
    "params": {
      "q": "*:*",
      "fl": "id, name",
      "sort": "id asc",
      "_": "1583898721030"
    }
  },
  "response": {
    "numFound": 52,
    "start": 0,
    "docs": [
      {
        "id": "/Users/haopengchen/solr-8.4.1/example/exemplodocs/post.jar",
        "name": null
      },
      {
        "id": "/Users/haopengchen/solr-8.4.1/example/exemplodocs/sample.html",
        "name": null
      },
      {
        "id": "/Users/haopengchen/solr-8.4.1/example/exemplodocs/solr-word.pdf",
        "name": null
      },
      {
        "id": "/Users/haopengchen/solr-8.4.1/example/exemplodocs/test_utf8.sh",
        "name": null
      },
      {
        "id": "0060248025",
        "name": null
      },
      {
        "id": "0380014300",
        "name": "Nine Princes In Amber"
      },
      {
        "id": "0441385532",
        "name": "Jhereg"
      },
      {
        "id": "0553293354",
        "name": "Foundation"
      },
      {
        "id": "0553573403",
        "name": "A Game of Thrones"
      },
      {
        "id": "055357342X",
        "name": "A Storm of Swords"
      }
    ]
  }
}
```

Below the response is a 'Run' section titled 'Querying' with a list of log messages:

- SLF4J: Failed to load class "org.slf4j.impl.StaticLoggerBinder".
- SLF4J: Defaulting to no-operation (NOP) logger implementation
- SLF4J: See <http://www.slf4j.org/codes.html#StaticLoggerBinder> for further details.

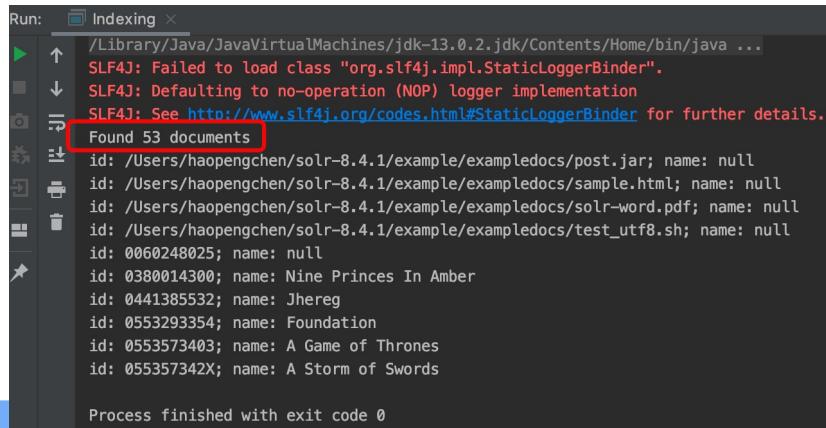
The 'Run' section also lists the 52 documents found, each with its ID and name. At the bottom, it says 'Process finished with exit code 0'.

Indexing

- Indexing.java

```
final SolrInputDocument doc = new SolrInputDocument();
doc.addField("id", UUID.randomUUID().toString());
doc.addField("name", "Amazon Kindle Paperwhite");

final UpdateResponse updateResponse = client.add("techproducts", doc);
// Indexed documents must be committed
client.commit("techproducts");
```



```
Run: Indexing x
/Library/Java/JavaVirtualMachines/jdk-13.0.2.jdk/Contents/Home/bin/java ...
SLF4J: Failed to load class "org.slf4j.impl.StaticLoggerBinder".
SLF4J: Defaulting to no-operation (NOP) logger implementation
SLF4J: See http://www.slf4j.org/codes.html#StaticLoggerBinder for further details.
Found 53 documents
id: /Users/haopengchen/solr-8.4.1/example/exemplardocs/post.jar; name: null
id: /Users/haopengchen/solr-8.4.1/example/exemplardocs/sample.html; name: null
id: /Users/haopengchen/solr-8.4.1/example/exemplardocs/solr-word.pdf; name: null
id: /Users/haopengchen/solr-8.4.1/example/exemplardocs/test_utf8.sh; name: null
id: 0060248025; name: null
id: 0380014300; name: Nine Princes In Amber
id: 0441385532; name: Jhereg
id: 0553293354; name: Foundation
id: 0553573403; name: A Game of Thrones
id: 055357342X; name: A Storm of Swords

Process finished with exit code 0
```

Java Object Binding

- TechProduct.java

```
public class TechProduct {  
    @Field  
    public String id;  
    @Field  
    public String name;  
  
    public TechProduct(String id, String name) {  
        this.id = id;  
        this.name = name;  
    }  
  
    public TechProduct() {}  
  
    public String getId() {  
        return this.id;  
    }  
  
    public String getName() {  
        return this.name;  
    }  
}
```

Java Object Binding

- ObjectBinding.java

```
public class ObjectBinding {  
    public static void main(String[] args) throws IOException, SolrServerException {  
        final SolrClient client = getSolrClient();  
        final TechProduct kindle = new TechProduct("kindle-id-4", "Amazon Kindle Paperwhite");  
        final UpdateResponse response = client.addBean("techproducts", kindle);  
        client.commit("techproducts");  
  
        final SolrQuery query = new SolrQuery("*:*");  
        query.addField("id");  
        query.addField("name");  
        query.setSort("id", SolrQuery.ORDER.asc);  
  
        final QueryResponse responseOne = client.query("techproducts", query);  
        final List<TechProduct> products = responseOne.getBeans(TechProduct.class);  
        for (TechProduct product : products) {  
            final String id = product.getId();  
            final String name = product.getName();  
            System.out.println("id: " + id + "; name: " + name);  
        }  
    }  
}
```

Java Object Binding

- ObjectBinding.java

```
public static SolrClient getSolrClient() {  
    final String solrUrl = "http://localhost:8983/solr";  
    return new HttpSolrClient.Builder(solrUrl)  
        .withConnectionTimeout(10000)  
        .withSocketTimeout(60000)  
        .build();  
}
```

```
Run: ObjectBinding  
/Library/Java/JavaVirtualMachines/jdk-13.0.2.jdk/Contents/Home/bin/java ...  
SLF4J: Failed to load class "org.slf4j.impl.StaticLoggerBinder".  
SLF4J: Defaulting to no-operation (NOP) logger implementation  
SLF4J: See http://www.slf4j.org/codes.html#StaticLoggerBinder for further details.  
id: /Users/haopengchen/solr-8.4.1/example/exampledocs/post.jar; name: null  
id: /Users/haopengchen/solr-8.4.1/example/exampledocs/sample.html; name: null  
id: /Users/haopengchen/solr-8.4.1/example/exampledocs/solr-word.pdf; name: null  
id: /Users/haopengchen/solr-8.4.1/example/exampledocs/test_utf8.sh; name: null  
id: 0060248025; name: null  
id: 0380014300; name: Nine Princes In Amber  
id: 0441385532; name: Jhereg  
id: 0553293354; name: Foundation  
id: 0553573403; name: A Game of Thrones  
id: 055357342X; name: A Storm of Swords  
Process finished with exit code 0
```

- Elasticsearch is a highly scalable open-source full-text search and analytics engine.
 - It allows you to store, search, and analyze big volumes of data quickly and in near real time.
 - It is generally used as the underlying engine/technology that powers applications that have complex search features and requirements.
 - <https://www.elastic.co>



- Near Real Time
 - Elasticsearch is a near real time search platform. What this means is there is a slight latency (normally one second) from the time you index a document until the time it becomes searchable.
- Document
 - A document is a basic unit of information that can be indexed.
 - This document is expressed in JSON (JavaScript Object Notation) which is a ubiquitous internet data interchange format.
- Index
 - An index is a collection of documents that have somewhat similar characteristics.
 - An index is identified by a name (that must be all lowercase) and this name is used to refer to the index when performing indexing, search, update, and delete operations against the documents in it.

- Shards & Replicas
 - An index can potentially store a large amount of data that can exceed the hardware limits of a single node.
 - To solve this problem, Elasticsearch provides the ability to subdivide your index into multiple pieces called shards.
- Sharding is important for two primary reasons:
 - It allows you to horizontally split/scale your content volume
 - It allows you to distribute and parallelize operations across shards (potentially on multiple nodes) thus increasing performance/throughput
- Replication is important for two primary reasons:
 - It provides high availability in case a shard/node fails. For this reason, it is important to note that a replica shard is never allocated on the same node as the original/primary shard that it was copied from.
 - It allows you to scale out your search volume/throughput since searches can be executed on all replicas in parallel.

- Index some documents

PUT /customer/_doc/1 (http://localhost:9200/customer/_doc/1)

The screenshot shows the Postman application interface. The top navigation bar includes 'New', 'Import', 'Runner', 'My Workspace', 'Invite', 'Upgrade', 'Launchpad', 'http://localhost:9200/customer/_doc/1', and environment settings. The main workspace shows a 'PUT' request to 'http://localhost:9200/customer/_doc/1'. The 'Body' tab is selected, containing the following JSON payload:

```
1 [ {  
2   "name": "John Doe"  
3 }]
```

The bottom section displays the response details: Status: 201 Created, Time: 764ms, Size: 284 B, and a 'Save Response' button. The 'Pretty' view of the response JSON is shown below:

```
1 {  
2   "_index": "customer",  
3   "_type": "_doc",  
4   "_id": "1",  
5   "_version": 1,  
6   "result": "created",  
7   "_shards": {  
8     "total": 2,  
9     "successful": 1,
```

At the bottom, there are icons for 'Bootcamp', 'Build', 'Browse', and help.

- Index some documents

GET /customer/_doc/1(http://localhost:9200/customer/_doc/1)

The screenshot shows the Postman application interface. The top navigation bar includes 'New', 'Import', 'Runner', 'My Workspace', 'Invite', 'Upgrade', and various status icons. The main header shows a 'GET' request to 'http://localhost:9200/customer/_doc/1'. Below the header, tabs for 'Params', 'Authorization', 'Headers (9)', 'Body', 'Pre-request Script', 'Tests', 'Settings', 'Cookies', and 'Code' are visible, with 'Params' being the active tab. Under 'Query Params', there is a table with columns 'KEY', 'VALUE', and 'DESCRIPTION'. A single row is present with 'Key' and 'Value' both empty and 'Description' set to 'Description'. The 'Body' tab is selected, showing the JSON response from the Elasticsearch query. The response is a multi-line JSON object:

```
1 {  
2   "_index": "customer",  
3   "_type": "_doc",  
4   "_id": "1",  
5   "_version": 1,  
6   "_seq_no": 0,  
7   "_primary_term": 1,  
8   "found": true,  
9   "_source": {  
10     "name": "John Doe"  
11   }  
12 }
```

Below the JSON, the status bar indicates 'Status: 200 OK Time: 36ms Size: 246 B' and a 'Save Response' button. The bottom navigation bar includes 'Pretty', 'Raw', 'Preview', 'Visualize', 'JSON', and search/filter icons. The footer features links for 'Bootcamp', 'Build', 'Browse', and help.

- Indexing documents in bulk

http://localhost:9200/bank/_bulk?pretty&refresh

The screenshot shows the Postman application interface. The request URL is `http://localhost:9200/bank/_bulk?pretty&refresh`. The response status is `200 OK`, time `504ms`, size `7.62 KB`. The response body is a JSON object indicating the operation took 471 milliseconds, had no errors, and processed 2 items. Each item was indexed into the 'bank' index as a '_doc' type document with ID '1' and version '2', resulting in an 'updated' status and forcing a refresh.

```
1 {  
2   "took": 471,  
3   "errors": false,  
4   "items": [  
5     {  
6       "index": {  
7         "_index": "bank",  
8         "_type": "_doc",  
9         "_id": "1",  
10        "_version": 2,  
11        "result": "updated",  
12        "forced_refresh": true,  
13        "_shards": {  
14          "total": 2  
15        }  
16      }  
17    }  
18  ]  
19}
```

- Indexing documents in bulk
 - http://localhost:9200/_cat/indices?v

The screenshot shows the Postman application interface. At the top, there are navigation tabs: New, Import, Runner, My Workspace, Invite, and Upgrade. Below the header, there's a search bar with the URL `GET http://localhost:9200/_cat/indic...` and a status indicator. To the right of the search bar are buttons for '+' and '***'. The main workspace is titled 'Untitled Request' and contains a 'Params' tab selected, showing a single parameter 'v' with a value of 'v'. Other tabs include Authorization, Headers (10), Body, Pre-request Script, Tests, Settings, Cookies, and Code. The 'Body' tab is active, displaying the response from the API call. The response is a table with columns: _id, health, status, index, uuid, pri, rep, docs.count, docs.deleted, store.size, pri.store.size. The data rows are:

_id	health	status	index	uuid	pri	rep	docs.count	docs.deleted	store.size	pri.store.size
1	yellow	open	bank	4qqz83NPPrpCjX0tCVvN5Jw	1	1	1000	1000	831kb	831kb
2	yellow	open	customer	jz4_75FQS0GHT4AwahjFfw	1	1	1	0	3.4kb	3.4kb
3										
4										

At the bottom, there are buttons for Pretty, Raw, Preview, Visualize, Text, and a copy icon. The status bar at the bottom right shows Status: 200 OK, Time: 75ms, Size: 311 B, and Save Response.

- Delete Index

- Now let's delete the index that we just created and then list all the indexes again:

```
DELETE /customer?pretty
GET /_cat/indices?v
```

- And the response:

```
health status index uuid pri rep docs.count docs.deleted store.size pri.store.size
```

- Before we move on, let's take a closer look again at some of the API commands that we have learned so far:

```
PUT /customer
PUT /customer/_doc/1
{
  "name": "John Doe"
}
GET /customer/_doc/1
DELETE /customer
```

- Indexing and Replacing Documents
 - Let's recall that command again:

```
PUT /customer/_doc/1?pretty
{
    "name": "John Doe"
}
```

- If we then executed the above command again with a different (or same) document, Elasticsearch will replace (i.e. reindex) a new document on top of the existing one with the ID of 1:

```
PUT /customer/_doc/1?pretty
{
    "name": "Jane Doe"
}
```

- Indexing and Replacing Documents

- If, on the other hand, we use a different ID, a new document will be indexed and the existing document(s) already in the index remains untouched.

```
PUT /customer/_doc/2?pretty
{
  "name": "Jane Doe"
}
```

- This example shows how to index a document without an explicit ID:

```
POST /customer/_doc?pretty
{
  "name": "Jane Doe"
}
```

- Deleting Documents
 - Deleting a document is fairly straightforward.

```
DELETE /customer/_doc/2?pretty
```

- Batch Processing

```
POST /customer/_doc/_bulk?pretty
{"index": {"_id": "1"}}
{"name": "John Doe" }
{"index": {"_id": "2"}}
{"name": "Jane Doe" }
```

```
POST /customer/_doc/_bulk?pretty
{"update": {"_id": "1"}}
{"doc": { "name": "John Doe becomes Jane Doe" } }
{"delete": {"_id": "2"}}
```

Elasticsearch

- Start searching

- http://localhost:9200/bank/_search

```
{  
  "query": { "match_all": {} },  
  "sort": [  
    { "account_number": "asc" }  
  ]  
}
```

The screenshot shows the Postman application interface. At the top, there's a navigation bar with 'Postman' and various icons. Below it is a toolbar with buttons for 'New', 'Import', 'Runner', and 'Invite'. The main area has tabs for 'Launchpad' and 'My Workspace'. A search bar shows 'GET http://localhost:9200/bank/_search'. On the right, there are buttons for 'Send', 'Save', and 'Comments (0)'. The central workspace is titled 'Untitled Request' and contains a 'Body' tab selected. The body content is a JSON document:

```
1 {  
2   "query": { "match_all": {} },  
3   "sort": [  
4     { "account_number": "asc" }  
5   ]  
6 }
```

Below the body, there's a large empty text area for the response. At the bottom, there are tabs for 'Body', 'Cookies', 'Headers (3)', and 'Test Results'. The status bar at the bottom right indicates 'Status: 200 OK' and 'Time: 179ms'.

Elasticsearch

- Start searching

- http://localhost:9200/bank/_search

```
{  
  "query": {  
    "bool": {  
      "must": [  
        {"match": {"age": "40"} }]  
      ],  
      "must_not": [  
        {"match": {"state": "ID"} }]  
    ]  
  }  
}
```

The screenshot shows a Postman window with the following details:

- Header Bar:** Postman, My Workspace, Invite, Upgrade.
- Request URL:** GET http://localhost:9200/bank/_search
- Body Tab:** Contains the Elasticsearch search query in JSON format.
- Code Tab:** Shows the raw JSON code with line numbers.
- Response Tab:** Status: 200 OK, Time: 15ms, Size: 1.05 KB. The response body is partially visible, showing the search took 9 units of time, had no timed out shards, and one successful hit.

- 请你在大二开发的E-Book系统的基础上，完成下列任务：
 1. 在你的项目中增加基于Solr或Lucene的针对书籍简介的全文搜索功能，用户可以在搜索界面输入搜索关键词，你可以通过全文搜索引擎找到书籍简介中包含该关键词的书籍列表。为了实现起来方便，你可以自己设计文本文件格式来存储书籍简介信息。例如，你可以将所有书籍的简介信息存储成为JSON对象，包含书的ID和简介文本，每行存储一本书的JSON对象。
 - 请将你编写的相关代码整体压缩后上传，请勿压缩整个工程提交。
- 评分标准：
 1. 能够正确实现上述搜索功能，并且集成到工程中(3分)

- Apache Lucene
 - <http://lucene.apache.org/>
- Lucene in Action
 - By Otis Gospodnetic & Erik Hatcher
 - MANNING Publishing
- Lucene 8.4.1 demo API
 - https://lucene.apache.org/core/8_4_1/demo/overview-summary.html#overview.description
- Apache Solr
 - <https://lucene.apache.org/solr/>
- Solr Tutorial
 - https://lucene.apache.org/solr/guide/8_4/solr-tutorial.html#launch-solr-in-solrcloud-mode
- Using SolrJ
 - https://lucene.apache.org/solr/guide/8_4/using-solrj.html
- Getting started with Elasticsearch
 - <https://www.elastic.co/guide/en/elasticsearch/reference/current/getting-started.html>
- Java API [7.6] » Document APIs
 - <https://www.elastic.co/guide/en/elasticsearch/client/java-api/current/java-docs.html>



Thank You!