**NLP and IR**

# Building your first Search Engine with Lucene

## Aliaksei Severyn

University of Trento, Italy

**March 07, 2013**

# Plan for the lab

- **Introduction to Lucene Search Engine**

- **Lucene concepts**

- **Hands-on experience with indexing and searching**

  - HelloWorld example

- **Using search engine to retrieve answer passages for Question Answering system**

  - Indexing and searching 50k of QA corpus
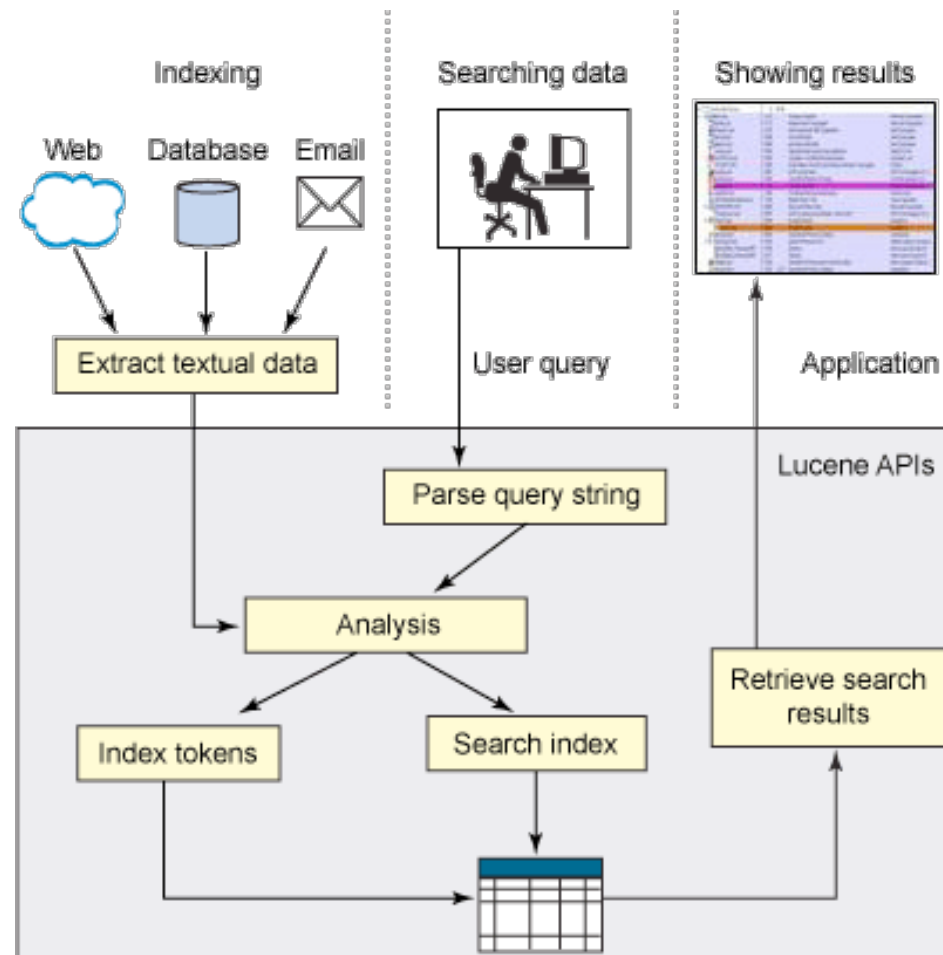
# What is Lucene

- software library for search

- open source

- not a complete application

- set of java classes

- active user and developer communities

- widely used in research and also in production

# High level overview

- Lucene is a full-text search library

- Designed to add search to your application

- Maintains a full-text index.

- Searches the index and returns results ranked by either the relevance to the query (or by an arbitrary field such as a document's last modified date.)

# Typical architecture of a Lucene search app

# Our first HelloWorld app with Lucene

- Create an in-memory index

- Add a few documents

- Construct a query

- Search an index

- Display results

# Setting up our first example

Download Lucene sources and binary from:

http://www.apache.org/dist/lucene/java/3.5.0/

Or download everything from :

https://github.com/aseveryn/NLPIR-2013

E.g. try the following in your terminal:

```
$ git clone https://github.com/aseveryn/NLPIR-2013.git
```

# Create a new project

# Name your project

# Drag HelloLucene.java to the src folder

# Add lucene-core-3.5.0.jar

# Run your first Lucene app!



```
        w.addDocument(doc);
    }
```

Problems  @ Javadoc  Declaration  Console

&lt;terminated&gt; HelloLucene [Java Application] /System/Library/Frameworks/JavaVM.framework/Versions/1.4/Home/bin/java (Apr 4, 2012 2:02:07 PM)
Found 2 hits.
1. Lucene in Action
2. Lucene for Dummies

# Adding Lucene documentation to the project

Go to Project properties->Libraries

Select lucene-core-3.5.0.jar

Select javadoc location

Locate lucene-core-3.5.0.jar

# Adding javadoc

# Setting the working environment

To be able to look at the Lucene internals:

Go to Project properties->Libraries

Select lucene-core-3.5.0.jar

Select source attachment

Locate src folder

# Adding sources to the JAR

# Now we can examine the sources and documentation

# Basic concepts in Lucene Search Engine

- Indexing

- Documents

- Fields

- Searching

- Queries

# Indexing

- Instead of searching the text directly it searches the index

- Uses **inverted index** - inverts a document-centric data structure (document->words) to a keyword-centric data structure (word->documents)

# Documents

- In Lucene, a **Document** is the unit of search and index.

- An index consists of one or more Documents.

- **Indexing –** adding Documents to an IndexWriter

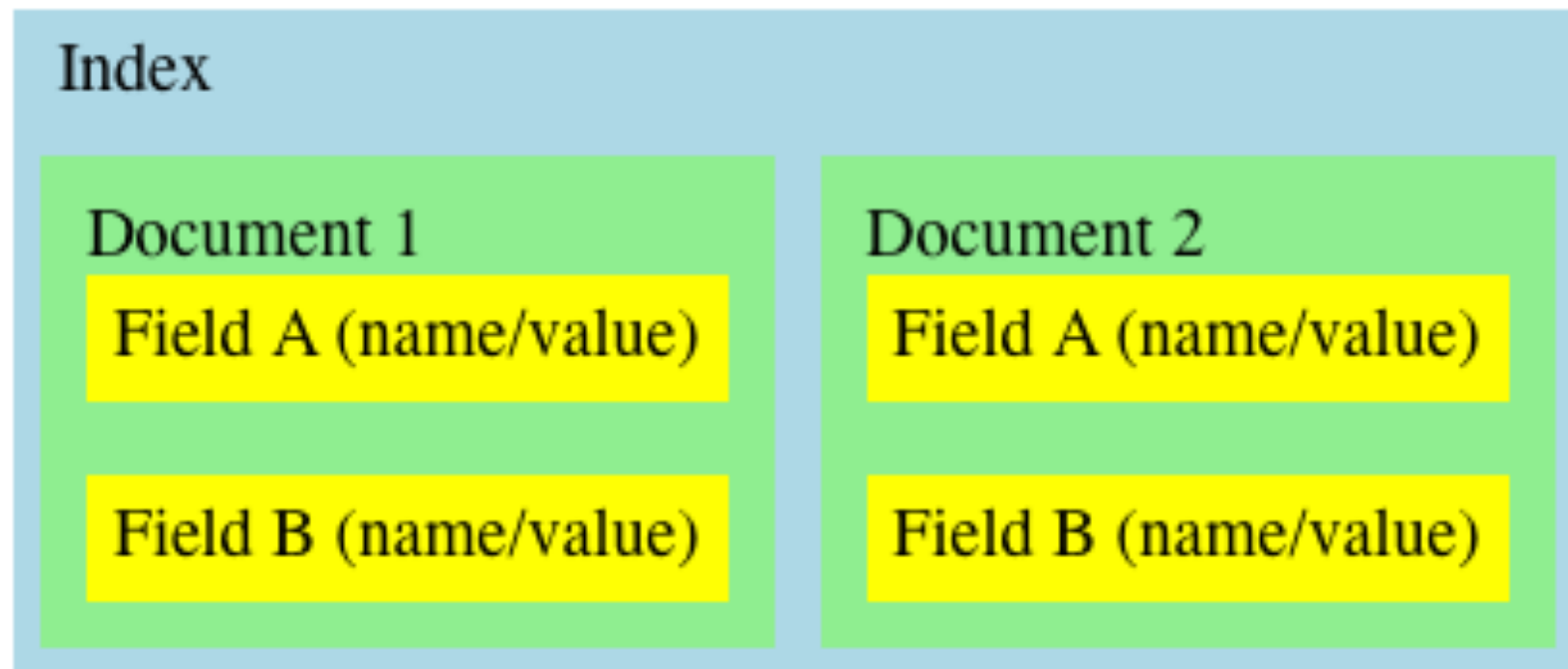- **Searching -** retrieving Documents from an index via an IndexSearcher.

# Fields

- A Document consists of one or more Fields.

- A Field is simply a name-value pair.

- For example, a Field commonly found in applications is *title*.

- Indexing in Lucene thus involves creating Documents of one or more Fields, and adding these Documents to an IndexWriter.

# Documents and fields in Lucene

Index

Document 1
Field A (name/value)

Field B (name/value)

Document 2
Field A (name/value)

Field B (name/value)

# Adding documents to the index

```java
Directory index = new RAMDirectory();

IndexWriterConfig config = new IndexWriterConfig(Version.LUCENE_35, analyzer);

IndexWriter w = new IndexWriter(index, config);
addDoc(w, "Lucene in Action");
addDoc(w, "Lucene for Dummies");
addDoc(w, "Managing Gigabytes");
addDoc(w, "The Art of Computer Science");
w.close();
```

```java
private static void addDoc(IndexWriter w, String value) throws IOException {
  Document doc = new Document();
  doc.add(new Field("title", value, Field.Store.YES, Field.Index.ANALYZED));
  w.addDocument(doc);
}
```

# Queries

Lucene has its own mini-language for performing searches.

Allows the user to specify which field(s) to search on, which fields to give more weight to (boosting), the ability to perform boolean queries (AND, OR, NOT) and other functionality.

# Query

We read the query from stdin, parse it and
build a lucene Query out of it.

```java
String querystr = args.length > 0 ? args[0] : "lucene";

// the "title" arg specifies the default field to use
// when no field is explicitly specified in the query.
Query q = new QueryParser(Version.LUCENE_35, "title", analyzer).parse(querystr);
```

# Searching

Searching requires an index to have already been built.

Very simple process:

- Create a **Query** (usually via a QueryParser)

- Handle this Query to an **IndexSearcher**

- Process a list of results

# Searching

- Using the Query we create a Searcher to search the index.

- Then instantiate a TopScoreDocCollector to collect the top 10 scoring hits.

```java
int hitsPerPage = 10;
IndexSearcher searcher = new IndexSearcher(index, true);
TopScoreDocCollector collector = TopScoreDocCollector.create(hitsPerPage, true);
searcher.search(q, collector);
ScoreDoc[] hits = collector.topDocs().scoreDocs;
```

# Display of results

Now that we have results from our search, we display the results to the user.

```
System.out.println("Found " + hits.length + " hits.");
for(int i=0;i<hits.length;++i) {
  int docId = hits[i].doc;
  Document d = searcher.doc(docId);
  System.out.println((i + 1) + ". " + d.get("title"));
}
```

# Let's get more practical

Build a Search Engine for answer passage retrieval in the Question Answering system

Use community QA site: Answerbag*

Use ~180k of automatically scraped question/answer pairs from over 20 categories

To reduce the amount of junk content focus only on professionally answered questions

http://www.answerbag.com/

# QA system with AnswerBag data

# Build high-quality QA corpus

# Professionaly researched answers

# Setting up QA example

**Index:**

$ sh run_qa_index.sh indexdir data/answers.50k.txt

**Search:**

$ sh run_qa_interactive.sh index 10

# Exercise

- Try out constructing boolean queries
- Experiment with a more involved StemAnalyzer for indexing and searching
- Answer: How old is too old for tires?