



MASTER IN  
COMPUTER  
SCIENCE

# Citation Search Engine

Academic paper search engine

## Master Thesis

Faculty of Natural Sciences  
University of Bern

January 2015

Prof. Dr. Oscar Nierstrasz

Mr. Haidar Osman, Mr. Boris Spasojevic

Software Composition Group

Institut für Informatik und angewandte Mathematik

University of Bern, Switzerland

*u<sup>b</sup>*

---

<sup>b</sup>  
UNIVERSITÄT  
BERN

**unine**  
UNIVERSITÉ DE  
NEUCHÂTEL

**UNI  
FR**  
■

UNIVERSITÉ DE Fribourg  
UNIVERSITÄT FREIBURG

# Abstract

Nowadays the number of documents in the World Wide Web grows at extremely fast rate<sup>1</sup>. Tools that can facilitate information retrieval (IR) present a particular interest in the modern world. We believe that considering meta information helps to build enhanced search systems that can facilitate IR. Particularly, we target an IR task for scientific articles. We consider citations in scientific articles as important text blocks summarizing or judging previous scientific findings, assisting in creating new scientific work.

We propose *Citation Search* – a software system that automatically extracts, indexes and aggregates citations from collections of scientific articles in a PDF format. Our system proposes a search interface to discover relevant scientific results based on a statement query. Also search interface of Citation Search allows to search for citations based on meta-information queries.

We evaluated searching capabilities of our system by conducting user evaluation experiments that compare it with alternative approaches. In the first set of experiments, we measured the efficiency of our system, i.e. how fast users can find relevant results in comparison with *Google Scholar*. We found that Citation Search performs equally well as Google Scholar. Secondly, we developed a citation aggregation feature to create automatic summaries of scientific articles and asked domain experts to evaluate summaries created by Citation Search and TextRank algorithms. We found that Citation Search outperforms TextRank algorithm in generating article summaries.

---

<sup>1</sup><http://googleblog.blogspot.ch/2008/07/we-knew-web-was-big.html>

# Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
1.1	Thesis statement . . . . .	4
1.2	Contributions . . . . .	5
1.3	Outline . . . . .	5
1.4	Glossary of Terms . . . . .	6
1.5	Typical Web Search Engine . . . . .	7
1.6	Inverted Index . . . . .	7
1.7	Dynamic Indexing . . . . .	11
1.8	Retrieving Search Results . . . . .	11
<b>2</b>	<b>Related Work</b>	<b>13</b>
2.1	Citations In Scientific Publications . . . . .	13
2.2	Popular Academic Search Engines . . . . .	14
<b>3</b>	<b>Citation Search Engine</b>	<b>16</b>
3.1	System Overview . . . . .	16
3.2	Parser . . . . .	17
3.2.1	PDF Processing . . . . .	19
3.2.2	Document Publishing . . . . .	22
3.2.3	Challenges . . . . .	25
3.3	Indexer . . . . .	27
3.3.1	Solr Ranking Model . . . . .	27
3.4	Meta Data Storage . . . . .	29
3.5	Web Search Interface . . . . .	29
3.5.1	Citation Search Main Page . . . . .	29
3.5.2	Search by Bibliography Page . . . . .	29

<i>CONTENTS</i>	3
<b>4 Evaluation Experiments</b>	<b>34</b>
4.1 Experiment Setup . . . . .	34
4.2 Questionnaires before and after each experiment . . . . .	36
4.3 Evaluation Results . . . . .	36
4.3.1 Results for Task 1 . . . . .	37
4.3.2 Results for Task 1' . . . . .	39
4.3.3 Results for Task 2 . . . . .	39
4.3.4 Final Questionnaire . . . . .	40
4.3.5 Results Summary . . . . .	40
<b>5 Conclusion</b>	<b>42</b>
<b>A User Guide for Citation Search Deployment</b>	<b>47</b>
A.1 Solr Installation . . . . .	47
A.1.1 Solr Configuration . . . . .	48
A.1.2 Enhanced Solr Search Features . . . . .	50
A.2 MongoDB Installation . . . . .	51
A.2.1 MongoDB configuration . . . . .	51
A.3 Running the parser . . . . .	52
A.4 Search Interface Deployment . . . . .	53

# 1

## Introduction

### 1.1 Thesis statement

Searching large collections of documents manually is not time efficient. The common way of making search effective in terms of time is indexing documents in advance. Before indexing procedure, documents should be parsed and analyzed. In this work we propose a search system over collections of scientific articles in PDF format. Scientific articles are different from ordinary documents in following sense. First, scientific articles in PDF format need to be converted into textual representation. Second, scientific articles have a certain layout format, where each section needs to be considered separately. In Citation Search Engine, we do not aim to index entire contents of articles. Instead, we aim to take advantage of indexing extracted meta-information. We argue that considering meta-information enhances a search system.

In general, a citation is composed of two parts: a citation text in a body of a document and a bibliographic link in a reference section. Citation Search Engine indexes both parts separately thus allowing user to make a choice between different indexed collections. We also store a context of a citation, i.e., text surrounding a citation, to provide better understanding of a cited text. Bibliographic links identify the cited origin paper. By finding bibliographic links pointing to the same original source we can merge citations related to the same paper. Such functionality is especially useful to automatically find citations of a given article in other articles.

## 1.2 Contributions

Following are the main contributions of this work:

- A novel IR system for scientific articles based on citations,
- A search interface to discover relevant scientific results based on a statement query,
- A search interface to discover citations based on meta-information queries,
- A new method of summary generation by means of citation aggregation,
- An empirical evaluation of the system by means of user study experiments.

## 1.3 Outline

The rest of the paper structured as follows:

**Chapter 1** gives a high overview of the architecture of a typical web search engine. It describes the main steps to construct an inverted index.

**Chapter 2** surveys the research related to citations in scientific publications. It overviews two popular academic search engines: Google Scholar and CiteSeer.

**Chapter 3** describes the design of Citation Search Engine. It first shows overall architecture of the proposed system and then shows details of implementation of each component.

**Chapter 4** describes user evaluation experiments and analysis the results.

**Chapter 5** concludes the work and describes potential future work.

## 1.4 Glossary of Terms

**Citation** is “a reference to a published or unpublished source”<sup>1</sup>. Consist of two parts: a cited text in the body of the article and a bibliographic link in the references section of the article. If it is not specified, by citation in this work we mean a cited text in the body of the article.

**Bibliographic link or bibliographic reference** is a bibliographic entry in the references section of the article associated with the cited text in the body of an article.

**Document** a broader term, having multiple meanings. In this work we can use a term *document* to refer to a single file, i.e a PDF article. A *document* term can be used to refer to a basic storage unit, i.e basic storage unit of a Indexes storage or a MongoDB database.

---

<sup>1</sup><http://en.wikipedia.org/wiki/Citation>

## 1.5 Typical Web Search Engine

Figure 1.1 illustrates a high level architecture of a standard web engine. It consists of three main components:

- Crawler
- Indexer
- Index Storage
- Search interface

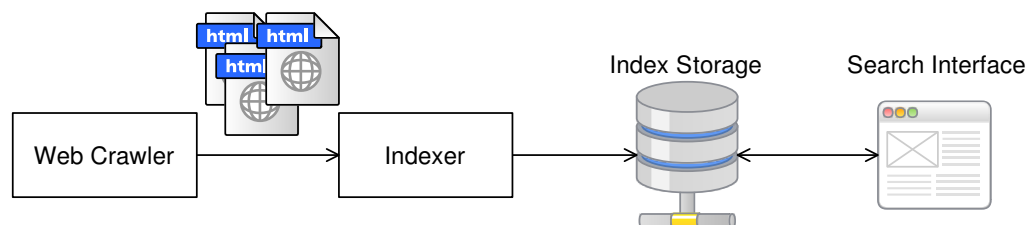


Figure 1.1: A high-level architecture of a typical web search engine

Web Crawler is a program that browses the World Wide Web reading the content of web pages in order to provide up-to-date data to Indexer. Indexer decides how a page content should be stored in an index storage. Indices help to quickly query documents from the index storage. Users can search and view query results through the Search Interface. When a user makes a query the search engine analyzes its index and returns best matched web pages according to specific criteria.

Web crawlers that fetch web pages with the content in the same domain are called focused or topical crawlers. An example of a focused crawler is an academic-focused crawler that crawls scientific articles. Such crawlers become components of focused search engines. Examples of popular academic search engines are Google Scholar and CiteSeer. Chapter 2 gives an overview of these search engines.

## 1.6 Inverted Index

Search engines like CiteSeer or Google Scholar deal with a large collection of documents. The way to avoid linear scanning the text of all documents for each query is to *index* them in advance. Thereby we are coming to the concept of *inverted index*, which is a major concept in IR. The



term *inverted index* comes from the data structure storing a mapping from content, such as words or numbers, to the parts of a document where it occurs. Figure 1.2 shows the basic idea of an inverted index. We have a dictionary of terms appearing in the documents. Each term maps to a list that records which documents the term occurs in. Each item in the list, conventionally named as *posting*, records that a term appears in a document, often it records the position of the term in the document as well. The dictionary on Figure 1.2 has been sorted alphabetically and each postings list is sorted by document ID. Document ID is a unique number that can be assigned to a document when it is first encountered.

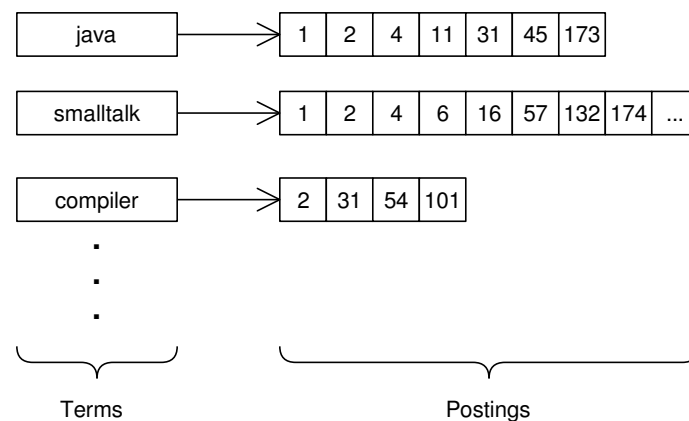


Figure 1.2: Example of an inverted index. Each term in a dictionary maps to a posting list consisting of document IDs, where this term occurs. Dictionary terms are sorted alphabetically and posting lists are sorted by document IDs

The construction of the inverted index has following steps:

1. Documents collection
2. Breaking each document into tokens, turning a document into a list of tokens
3. Linguistic preprocessing of a list of tokens into normalised list of tokens
4. Index documents by creating an inverted index, consisting of a dictionary with terms and postings

First step of the index construction is documents collection. It aims at obtaining a set of documents containing a sequence of characters. The input to the indexing process could be digital documents that are bytes in a file or a web server. Consider, for example, a collection of PDF

files. First, we need to correctly decode out characters of the binary representation. Next we should determine a document unit. For example, it might be a chapter in a book, or a paragraph in a scientific article.

The next step after getting the sequences of characters in the document units, is to break up documents into *tokens*. Tokens can be thought of as the semantical units for processing. For example, it might be a word or a number. During tokenisation some characters like punctuations can be thrown out.

Here is a tokenisation example:

Input: Sometimes, I forget things.

Output: Sometimes I forget things

The third step is normalization. Consider an example of querying the word *co-operation*. A user might also be interested in getting documents containing *cooperaion*. *Token normalisation* is a process of turning a token into a canonical form so matches can occur despite lexical differences in the character sequences. One way of token normalisation is keeping relations between unnormalized tokens, which can be extended to manual constructed synonym lists, such as *car* and *automobile*. The most standard way of token normalization however is creating *equivalence classes*. If tokens become identical after applying a set of rules then they are in the equivalence classes. Consider the following common normalization rules that are often used:

**Stemming and Lemmatization** Words can be used in different grammatical forms. For instance, *organize*, *organizes*, *organizing*. However in many cases it sounds reasonable for one of these words to return documents contain other forms of the word. The goal of stemming and lemmatization is to reduce the form of the word to a common base form.

Here is an example:

am, are, is =>be

car, cars, car's, cars' =>car

The result of applying the rule to the sentence:

three frogs are flying =>three frog be fly

Stemming and lemmatization are closely related concepts however there is a difference. *Lemmatization* usually refers to finding a *lemma*, common base of a word, with the help of a vocabulary, morphological analysis of a word and requires understanding the context of a word and language grammar. *Stemming* however operates with a word without knowing its context and thereby can't distinguish that the same words have different meanings depending on the context.

Here is an example:

*better* => *good*, can only be matched by lemmatization since *it* requires dictionary look-up

*writing* => *write*, can be matched by both lemmatization and stemming

*meeting* => *meeting(noun)* or *to meet(verb)*, can be matched only by lemmatization since it requires the word context

In general, stemmers are easier to implement and run faster. The most common algorithm for stemming is *Porter's* algorithm [22].

**Capitalization/Case-Folding** A simple strategy is to reduce all letters to a lower case, so that sentences with *Automobile* will match to queries with *automobile*. However this approach would not be appropriate in some contexts like identifying company names, such as *General Motors*. Case-folding can be done more accurately by a machine learning model using more features to identify whether a word should be lowercased.

**Accents and Diacritics** Diacritics in English language play an insignificant role and simply can be removed. For instance *cliché* can be substituted by *cliche*. In other languages diacritics can be part of the writing system and distinguish different sounds. However, in many cases, users can enter queries for words without diacritics.

The last step of building the inverted index is sorting. The input to indexing is a list of pairs of normalized tokens and documents IDs for each document. Consider an example of three documents with their contents:

- Document 1: Follow the rules.
- Document 2: This is our town.
- Document 3: The gates are open.

After applying tokenisation and normalisation steps of the listed documents the input to the indexing is shown in Table 1.1. The indexing algorithm sorts the input list so that the terms are in alphabetical order as in Table 1.2. Then it merges the same terms from the same document by folding two identical adjacent items in the list. And finally instances of the same term are grouped and the result is split into a dictionary with postings, as shown in Table 1.3.

Term	DocumentID
follow	1
the	1
rule	1
this	2
be	2
our	2
town	2
the	3
gate	3
be	3
open	3

Table 1.1: Input to the indexing algorithm is a list of pairs of a term and document ID, where this term occurs.

## 1.7 Dynamic Indexing

So far we assumed that document collection is static. However there are many cases when the collection can be updated, for example, by adding new documents, deleting or updating existing documents. Simple way to deal with dynamic collection is to reconstruct the inverted index from scratch. This might be acceptable if the changes made in the collection are small over time and the delay in making new documents searchable is not critical. However if there is one of the aforementioned conditions is violated, one might be interested in another more dynamic solution like keeping an auxiliary index. Thus we have a large main index and we keep auxiliary index for changes. The auxiliary index is kept in memory. Every time a user makes a query the search runs over both indexes and results are merged. When the auxiliary index becomes too large it can be merged with the main index.

## 1.8 Retrieving Search Results

When a user makes a query it would be good to give her back a result document containing all terms in the query, so that terms are located close to each other in the document. Consider an example of querying a phrase containing 4 terms. The part of the document that contains all terms is named a *window*. The size of the window is measured in number of words. For instance the smallest window for 4-term query will be 4. Intuitively, smaller windows represent better results for users. Such window can become one of the parameters ranking a document in the search result. If there is no document containing all 4 terms, a 3-term phrase can be queried.

Term	DocumentID
be	2
be	3
follow	1
gate	3
open	3
our	2
rule	1
the	1
the	3
this	2
town	2

Table 1.2: Indexing algorithm sorts all terms in a alphabetical order. The result is a list of sorted terms with document IDs

Term	Postings
be	2 3
follow	1
gate	3
open	3
our	2
rule	1
the	1 3
this	2
town	2

Table 1.3: Indexing algorithm groups the same terms with creating postings. The result is a dictionary with terms as keywords and values as postings.

Usually search systems hide the complexity querying from the user by introducing *free text query parsers* so a user can make only one query.

# 2

## Related Work

### 2.1 Citations In Scientific Publications

Citations are the subject of many interesting scientific studies.

Bradshaw et al. [3] showed that citations provide many different perspectives on the same article. They believe that citation provide means to measure the relative impact of articles in a collection of scientific literature. In their work the authors improved the relevance of documents in the search engine results with a method called Reference Directed Indexing. Reference Directed Indexing (RDI) is based on a comparison of the terms authors use in reference to documents.

Bertin and Atanassova [1] [14] [2] automatically extract citations and annotate them using a set of semantic categories. In [14] and [1] they used linguistic approach, which used the contextual exploration method, to annotate automatically the text. In [2] they proposed a hybrid method for the extraction and characterization of citations in scientific papers using machine learning combined with rule-based approaches.

There are several studies that used citations to evaluate science by introducing map of science. Map of science graphically reflects the structure, evolution and main contributors of a given scientific field [6] [11] [13] [26].

Kessler [10] first used the concept of bibliographic coupling for document clustering. To build a cluster of similar documents Kessler used a similarity function based on the degree of

bibliographic coupling. Bibliographic coupling is a number of citations two documents have in common. The idea was developed further by Small in co-citation analysis [25]. Later co-citation analysis and bibliographic coupling was used by Larson [12] for measuring similarity of www pages.

Another approach is to use citations to build summaries of scientific publications. There are three categories of summaries proposed based on citations: overview of a research area (multi-document summarization) [20], impact summary (single document summary with citations from the scientific article itself) [16] and citation summary (multi- and single document summarization, in which citations from other papers are considered) [23]. In work by Nakov et al. citations have been used to support automatic paraphrasing [19].

An expert literature survey on citation analysis was made by Smith [27], she reviewed hundred of scientific articles on this topic.

## 2.2 Popular Academic Search Engines

**CiteSeer<sup>x</sup>** CiteSeer<sup>x</sup> is built on the concept of citation index. The concept of citation index was first introduced by Eugene Garfield in [7]. In terms of Eugene Garfield citations are bibliographic links or referrers linking scientific documents. In his work Eugene Garfield proposed an approach where citations between documents were manually cataloged and maintained so that a researcher can search through listings of citations traversing citation links either back through supporting literature or forward through the work of later researchers [4].

Lawrence et al. automated this process in CiteSeer<sup>x</sup> <sup>1</sup> [8], a Web-based information system that permits users to browse the citation links between documents as hyperlinks. CiteSeer<sup>x</sup> automatically parses and indexes publicly available scientific articles found on the World Wide Web.

CiteSeer<sup>x</sup> is built on top of the the open source infrastructure SeerSuite <sup>2</sup> and uses Apache Solr <sup>3</sup> search platform for indexing documents. It can extract meta information from papers such as title, authors, abstract, citations. The extraction methods are based on machine learning approaches such as ParseCit [5]. CiteSeer<sup>x</sup> currently has over 4 million documents with nearly 4 million unique authors and 80 million citations.

CiteSeer<sup>x</sup> indexes citations more precisely bibliographic links or referrers while in Citation Search Engine we intend to index not only bibliographic links but also cited text in a body of a document. If by indexing bibliographic links CiteSeer<sup>x</sup> mainly aims to simplify navigation

---

<sup>1</sup>CiteSeer, <http://citeseerx.ist.psu.edu/>

<sup>2</sup>SeerSuite, <http://citeseerx.sourceforge.net/>

<sup>3</sup>Apache Solr, <http://lucene.apache.org/solr/>

between linked documents, in Citation Search we focus on simplifying retrieval of documents containing a text of interest.

**Google Scholar** Google Scholar is a freely accessible web search engine that makes full-text or metadata indexing of scientific literature <sup>4</sup>. Besides the simple search Google Scholar proposes following features. Unique ranking algorithm, an algorithm that ranks documents “the way researchers do, weighing the full text of each document, where it was published, who it was written by, as well as how often and how recently it has been cited in other scholarly literature” <sup>5</sup>. “Cited by” feature allows to view abstracts of articles citing the given article. “Related articles” feature shows the list of closely related articles. It is also possible to filter articles by author name or published date. Google Scholar contains roughly 160 million of documents by May 2014 [21].

---

<sup>4</sup>Google Scholar, <http://scholar.google.ch/>

<sup>5</sup><https://scholar.google.com/scholar/about.html>



# 3

## Citation Search Engine

### 3.1 System Overview

The components of Citation Search is shown on Figure 3.1. There are three main operations performed by Citation Search Engine: parsing PDF files, indexing documents and querying the resulted indexes. Correspondingly, there are three major components responsible for accomplishment of these operations: *Parser*, *Indexer* and *Search Web App*. The system has two more components for storing data: *Indexes Storage* and *Meta Data Storage*. We use *Indexes Storage* for storing indexes built on citations. This storage is very simple and was not designed to represent any relations in data structures. Moreover, it does not allow to perform any sophisticated operations over the stored data. Therefore, we use *Meta Data Storage* to represent complex data structures and perform sophisticated queries, like aggregating citations referred to the same article.

The workflow of the system is shown on Figure 3.2. The first operation performed by the system is parsing. The *Parser* converts a PDF file into the text. Then it extracts meta information, like citations and references, from the textual representation of the file. Next it packages extracted information into data units corresponding to formats acceptable by *Indexer* and *Meta data storage*. A data unit publishable to the *Indexer* consist of a citation, that should be indexed and an additional information related to this citation (citation context, a file URI, bibliographic references) that

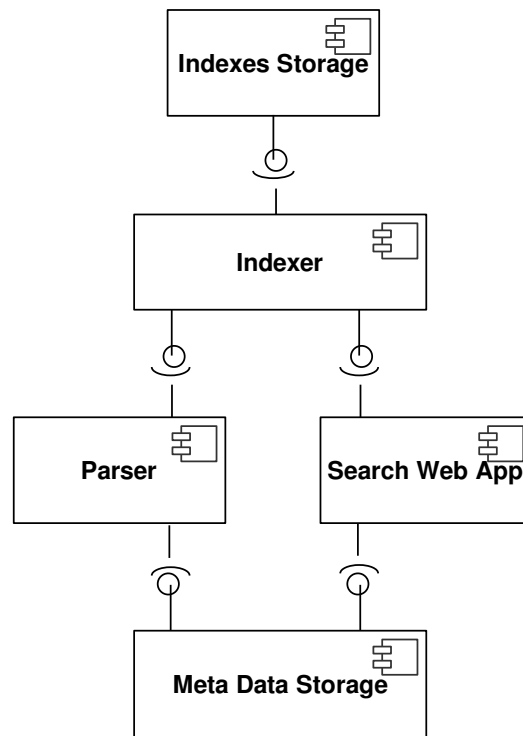


Figure 3.1: Component diagram of Citation Search.

should be stored. A data unit publishable to the *Meta data storage* consist of a citation, a source paper identifier and bibliographic references. We use *Meta data storage* for aggregating citations referred to the same source. Once the *Parser* processed the PDF file it can proceed to the next paper if there is any left. When all papers are processed, user can make queries over *Search Web App*.

The next sections of this chapter describe the implementation of each component in detail and show the reasons behind choosing a particular solution.

## 3.2 Parser

It is practical to divide the work of the *Parser* into two phases: *PDF processing* and *Document publishing*, as in Figure 3.3. The output of the *PDF processing* phase is the input to the *Documents publishing* phase.

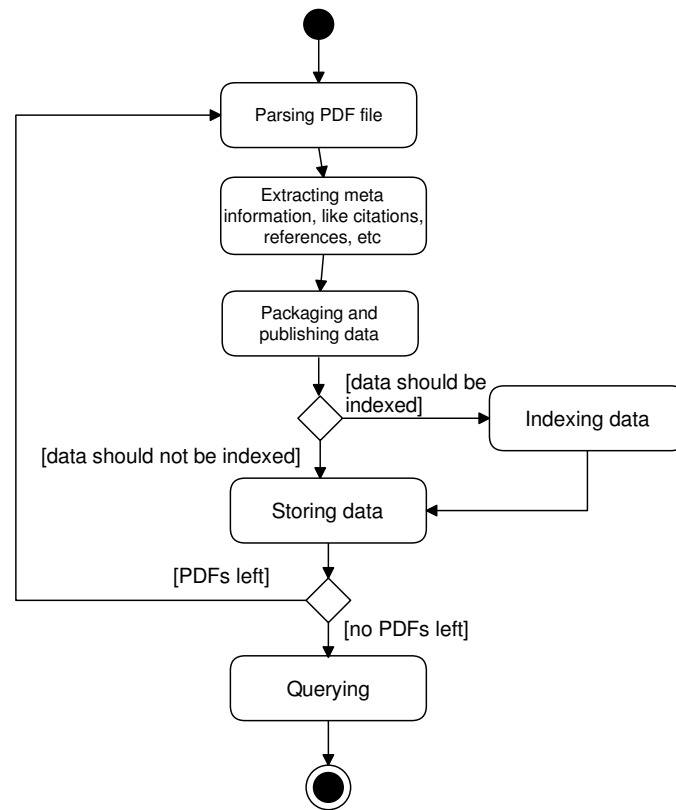


Figure 3.2: Activity diagram of Citation Search.

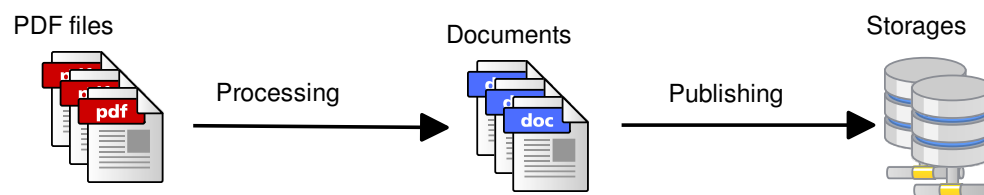


Figure 3.3: Parser workflow

### 3.2.1 PDF Processing

The main role of *PDF processing* phase is to parse scientific articles into text and extract citations and bibliographic references to create documents for publishing. Parsing PDF-files from different sources is a very challenging task as there is a large number of scientific articles in different formats. Besides that, there exist scientific articles written decades ago which can use different encoding algorithms from articles created recently. Thereby, building a universal parser is very hard in practice. In our case, we try to identify common patterns covering the majority of the scientific article formats or at least the formats found in our collection of articles.

*PDF processing* phase starts with recursively walking through the directory tree of the collection library. While walking through the directory, the *Parser* filters non-pdf files and parses and processes each PDF-file separately. We use Apache PDFBox library <sup>1</sup> to parse PDF-files into text. The library extracts full text from PDF-files, but without any hints on the initial structure of the article. In our case, to find citations and bibliographic references in text, we need to look them up in different parts of the article. We designed an algorithm to break the PDF-text into sections.

Generally, we are interested in identifying a body of the document where we can find citations and reference blocks where we can find bibliographic links. One way of finding these sections can be using keywords that might signify the beginning or the end of some sections. Based on those keywords, one can extract different sections of a document. Figure 3.4 shows a sample text of a parsed pdf-document with keywords.

One can notice the following characteristics of scientific articles:

- The body of a document comes before the references section.
- The appendix or author's biography sections can come after the references section.
- Each document contains the “Abstract” and the “References” words and might contain the “Appendix” word. We call these words keywords.

The keywords can be written in different formats, like using upper or lower cases. Table 3.1 illustrates some variations of the keywords.

body	references	appendix
Abstract	References	Appendix
ABSTRACT	References: REFERENCES	APPENDIX

Table 3.1: Keywords identifying different sections in a document

<sup>1</sup>Apache PDFBox, <https://pdfbox.apache.org/>

After breaking a document down into sections as shown in Figure 3.4, the text is presented in one-column format. There are two aspects regarding this format. First, sentences can be split by new line symbols at the end of a line. Second, words can be split by dash symbol at the end of a line. We introduce a normalization step where new lines are substituted by white spaces and dashes are removed in the end of a line to obtain continuous text.

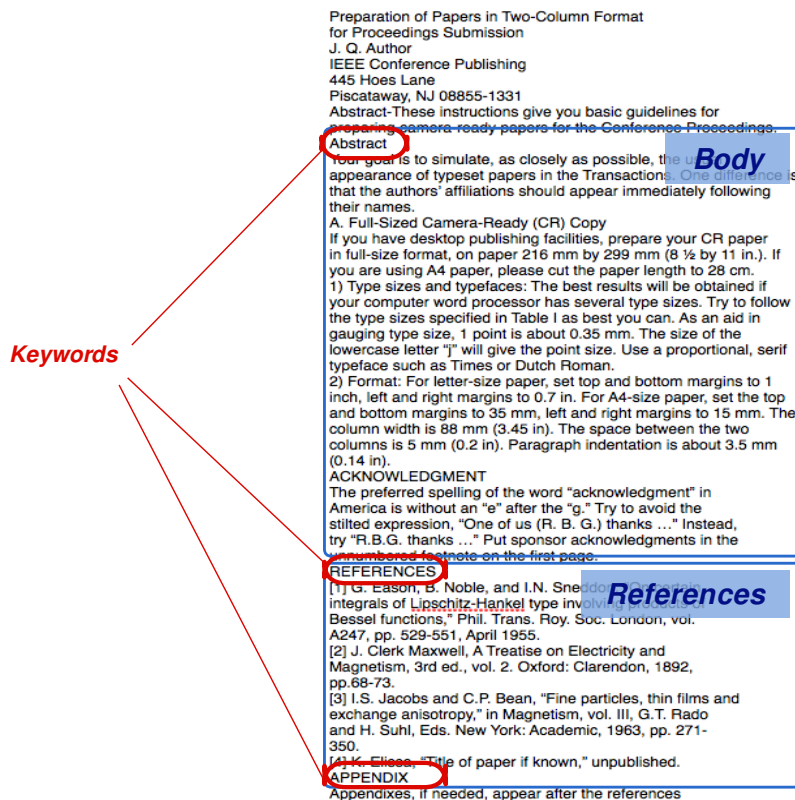


Figure 3.4: Sample text of the parsed scientific article. Keywords help to break the document into sections.

As a result of the normalization step, we have a document divided into body and references sections. Before searching citations in the body of a document, we break the body into sentences. In general, breaking text into sentences is not an easy task. Consider a simple example with a period. Period not only indicates the end of a sentence but also can be encountered inside the sentence itself, like an item of a numbered list or a name of a scientist. Besides, not all sentences end with a period, like the title of a section or an item of a list. We used Stanford CoreNLP library

that employs natural language and machine learning techniques to extract sentences [9].

Next, we search for the citations in the body and for the bibliography links in the references section. When an author makes a citation she puts a link to a bibliographic reference in the sentence. It is common to use square brackets ([ and ]) to make a link to a bibliographic reference in the sentence. Thus, we can identify citations by detecting square brackets in the text. After analyzing some set of articles we found multiple patterns in using square brackets for citations, as shown in Table 3.2.

Patterns of using [ ]	Example in text
[21]	Our conclusion is that, contrary to prior pessimism [21], [22], data mining static code attributes to learn defect predictors is useful.
[20, 3, 11, 17]	In the nineties, researchers focused on specialized multivariate models, i.e., models based on sets of metrics selected for specific application areas and particular development environments [20, 3, 11, 17].
[24, Sections 6.3 and 6.4 ]	Details on the life-cycle of a bug can be found in the BUGZILLA documentation [24, Sections 6.3 and 6.4].
[PJe02]	In a lazy language like Haskell [PJe02] this is not an issue - which is one key reason Haskell is very good at defining domain specific languages.

Table 3.2: Frequent patterns in using square brackets ([ and ])for citing

Further we need to extract bibliographic links from the references section. For that we studied most common variants of composing the references sections. Table 3.3 summarises these findings. To extract bibliographic links we make a list of regular expressions matching one of those patterns listing in Table 3.3. Extracting bibliographic links allows us to match citations with bibliographic links.

The pipeline of PDF processing stages described above is depicted on Figure 3.5. As seen from Figure 3.5 to complete the processing stage we need to perform one more step: extracting titles from bibliographic links. The objective point of extracting titles from bibliographic links is to collect citations referred to the the same source (scientific article). In general case, different formats of bibliographic links can identify the same source or scientific article. For example, an article may have different editions, published in different journals in different years or simply different authors may use different style formatting. What we consider to be identical for all

References section templates
[1] J. Bach. Useful features of a test automation system (partiii) ...
[2] B. Beizer. Black-Box Testing. John Wiley and Sons, ...
...
1. J. R. Hobbs, Granularity, Ninth International Joint Conference ...
2. W. Woods, What's in a Link: Foundations for Semantic Networks, ...
...
[1]. Arnold, R.S., Software Reengineering, ed. ...
[2]. Larman, C., Applying UML and Patterns. 1998, ...
...
[ASU86] A. V. Aho, R. Sethi, and J. D. Ullman. Compilers: ...
[AU73] A.V. Aho and J.D. Ullman. The theory of parsing, translation ...
...

Table 3.3: Frequent patterns of writing bibliographic links in a references block

bibliographic links citing the same paper is the paper's title.

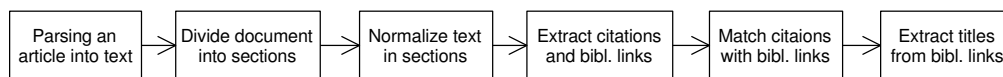


Figure 3.5: Pipeline of the PDF processing stage

**Processing bibliographic links** As usual we try to recognize common patterns covering the majority of bibliographic links. Table 3.4 shows some examples of bibliographic links. First, we noticed that if a bibliographic link contains some sort of quotes, for example, double quotes (“”) or single quotes (‘’), then it is highly probable that a title is enclosed by these quotes. Then, we made some observations for bibliographic links without quotes. Very often, a bibliographic link is structured as follows: it begins by listing the paper's authors, then the title, and then comes the rest of the link (see Figure 3.6). It turned out that Core NLP library used for breaking text into sentences is good in dividing a link into parts according to our view. In most of cases it is enough to take the second part of the bibliographic link to be a title.

### 3.2.2 Document Publishing

There are two systems where documents are published to: Solr and MongoDB. Solr is used for indexing citations and MongoDB for storing relations in meta-data. We aim to use MongoDB for aggregating citations referred to the same source paper.

Conradi, R., Dyba, T., Sjoberg, D.I.K., and Ulsund, T., "Lessons learned and recommendations from two large norwegian SPI programmes." Lecture notes in computer science, 2003, pp. 32-45."
P. Molin, L. Ohlsson, 'Points & Deviations - A pattern language for fire alarm systems,' to be published in Pattern Languages of Program Design 3, Addison-Wesley.
R. P. Wilson and M. S. Lam. Effective context sensitive pointer analysis for C programs. In PLDI, pages 112, June 1995. 289
Allen, Thomas B. Vanishing Wildlife of North America. Washington, D.C.: National Geographic Society, 1974.
I. Herraiz, J. M. Gonzalez-Barahona, and G. Robles. Towards a Theoretical Model for Software Growth. In Proceedings of the 4th International Workshop on Mining Software Repositories, Minnesota, USA, May 2007.

Table 3.4: Some examples of bibliographic links

Authors	Title	Rest
R. P. Wilson and M. S. Lam.	Effective context sensitive pointer analysis for C programs	In PLDI, pages 1–12, June 1995. 289

Figure 3.6: Common structure of a bibliographic link

**Publishing documents to Solr** The common way to interact with Solr is using REST API. Solr provides client libraries for many programming languages to handle interactions with Solr's REST API. In our project we used SolrJ client library for Java language. This library abstract interaction with Solr into java objects instead of using typical XML request/response format. Basic Solr storage unit is called document and SolrJ has document abstraction implementation. For every detected citation we compose a document to publish. Figure 3.7 represents a structure of documents we publish to Solr.

Document
<ul style="list-style-type: none"> <li>- id: int</li> <li>- text: String</li> <li>- context: String</li> <li>- path: String</li> <li>- references : List</li> </ul>

Figure 3.7: Document structure publishing to Solr



Every document representing one citation consist of following fields:

- id: document unique id, mandatory field for publishing to Solr
- text: text of the citation that we want to index
- context: citation with a text framing it, we take 1 sentence before and 1 after the citation
- path: URL of a document where citation was found
- references: list of bibliographic links from references section matching this citation

Solr uses NoSQL-like data storage, but actually it is even more limited than traditional NoSQL databases. The data stored in Solr is very ‘flat’, which means that Solr cannot store hierarchical data [28], [24]. In our case, along with the references, we intend to store a title of the scientific article parsed from the reference string, so we can aggregate citations referred to the same scientific article. It is worth to say that citation itself can refer to multiple scientific articles. We are also interested in a solution that does not require reviewing all Solr documents to find citations referred to the same scientific article as it will be too slow and will decrease the quality of user experience. Thus we use an external storage solution that can keep the titles of scientific articles and all the citations referred to a specific article. As there are few relations in our data and we would like to have a scalable solution we decided to use MongoDB as an external storage.

**Publishing documents to MongoDB** MongoDB is a document-oriented NoSQL database that stores data in JSON-like documents with dynamic schema <sup>2</sup>. To connect to the database we used a Java driver provided by MongoDB. Although MongoDB is a ‘schemaless’ database we adhere to the JSON structure of the document shown in Listing 1. The json document consist of following fields:

- id: document id, field automatically assigned by MongoDB
- title: title of a scientific article
- citations: citations with its references of the scientific article identifying by title field

Every time we send a new citation with a paper title to MongoDB, we check if a document with the same title already exists. If so, we add a new citation to document, otherwise create a new document.

---

<sup>2</sup>MongoDB database, <http://www.mongodb.org/>

```

{
  "_id" : ObjectId("547ef1b219795f049d6a0ad0"),
  "title" : "Re-examining the Fault Density-Component Size Connection",
  "citations" : [
    {
      "citation" : "Hatton, [19], claims that there is compelling empirical
                    evidence from disparate sources to suggest that in any
                    software system, larger components are proportionally more
                    reliable than smaller components.",
      "references" : [
        "[19] L. Hatton, Re-examining the Fault Density-Component Size ..."
      ]
    },
    {
      "citation" : "Hatton examined a number of data sets, [15], [18] and
                    concluded that there was evidence of macroscopic behavior
                    common to all data sets despite the massive internal
                    complexity of each system studied, [19].",
      "references" : [
        "[15] K.H. Moeller and D. Paulish, An Empirical Investigation of ...",
        "[18] T. Keller, Measurements Role in Providing Error-Free Onboard ...",
        "[19] L. Hatton, Re-examining the Fault Density-Component Size ..."
      ]
    }
  ]
}

```

Listing 1: Sample document stored in MongoDB

### 3.2.3 Challenges

One of the challenges we encounter during parsing is finding cited sentences that do not contain any specific identifiers. Indeed the task is straightforward when a sentence contains square brackets, for example, ‘[34]’ or ‘[Ali86]’. However, sometimes a link to bibliography can be composed only from the authors’ names. In this case it becomes difficult to distinguish a citation from any other sentence (see Figure 3.8).

Another challenge in using square brackets as citation identifiers arises when square brackets are used not as links to bibliography. For example, a parser might mix up an array in a code snippet with a link to bibliography (see Figure 3.9). Usage of code snippets are common in computer science literature so it would be nice to have a method to distinguishing a snippet of source code from the natural text.

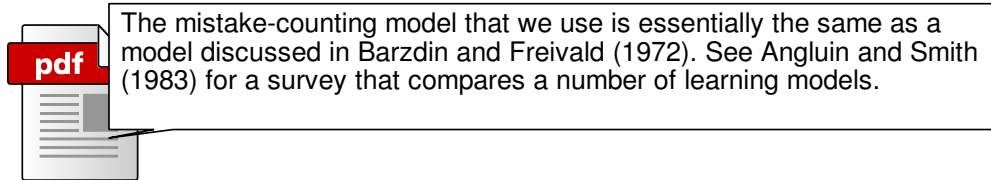


Figure 3.8: An example of citations from the article “Learning Abound: Quickly When Irrelevant Attributes A New Linear-threshold Algorithm” authored by Nick Littlestone. In both sentences a link to bibliography composed from authors’ names making it hard to distinguish a citation from any other sentences.

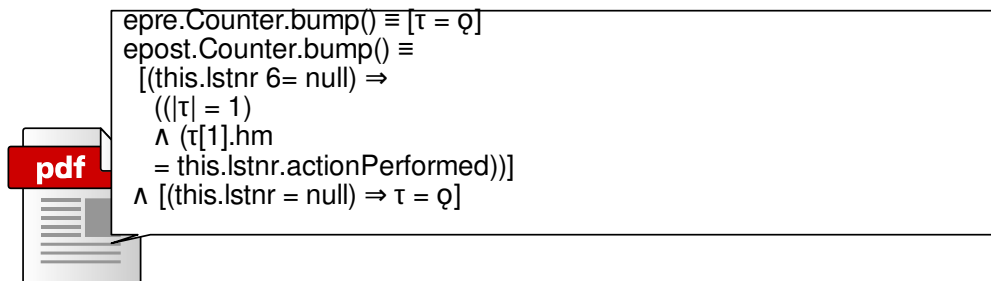


Figure 3.9: An example of the code snippet from the article “Modular Verification of Higher-Order Methods with Mandatory Calls Specified by Model Programs” authored by Steve M. Shaner et al. The code snippet is wrongly considered as a citation and is matched to the first bibliographic entry.

### 3.3 Indexer

We use Solr as an Indexer: a solution from Apache Software Foundation built on Apache Lucene. Apache Lucene is an open source, IR library that provides indexing and full text search capabilities<sup>3</sup>. While web search engines focus on searching content on the Web, Solr is designed to search content on corporate networks of any form. Some of the public service that use Solr as a server are Instagram (photo and video sharing social network), Netflix (movie hosting service) and StubHub.com (public entertainment events ticket reseller).

Figure 3.10 illustrates a high level architecture of Solr. Solr is distributed as a Java web application that runs in any servlet container, for example, Tomcat or Jetty. It provides REST-like web services so external applications can make queries to Solr or index documents. Once the data is uploaded, it goes through text analysis pipeline. In this stage, different preprocessing phases can be applied to remove duplicates in the data or some document-level operations prior to indexing, or to create multiple documents from a single one. Solr comes with a variety of query parser implementations responsible for parsing the queries passed by the end user as search strings. For example, *TermQuery*, *BooleanQuery*, *PhraseQuery*, *PrefixQuery*, *RangeQuery*, *MultiTermQuery*, *FilteredQuery*, *SpanQuery* and others. Solr has xml configuration files (schema.xml and solrconfig.xml) to define the structure of the index and how fields will be represented and analyzed (see Appendix A.1 for Solr installation and configuration).

#### 3.3.1 Solr Ranking Model

Solr ranking model is based on the Lucene scoring algorithm, also known as a TF-IDF model [15]. This model takes into consideration following factors:

- **tf** - term frequency, a frequency of the term in a document. The higher the term frequency, the higher a document score.
- **idf** - inverse document frequency, an inverse frequency of the term in all documents. The rare the term occurs in all documents, the higher its contribution to the document's score.
- **coord** - coordination factor, takes into account the number of query terms in a document. The more query terms in a document, the higher score it has.

The exact scoring formula can be found on the official web page of the Lucene documentation

4.

---

<sup>3</sup>Apache Lucene, <http://lucene.apache.org/core/>

<sup>4</sup>Apache Lucene, scoring formula

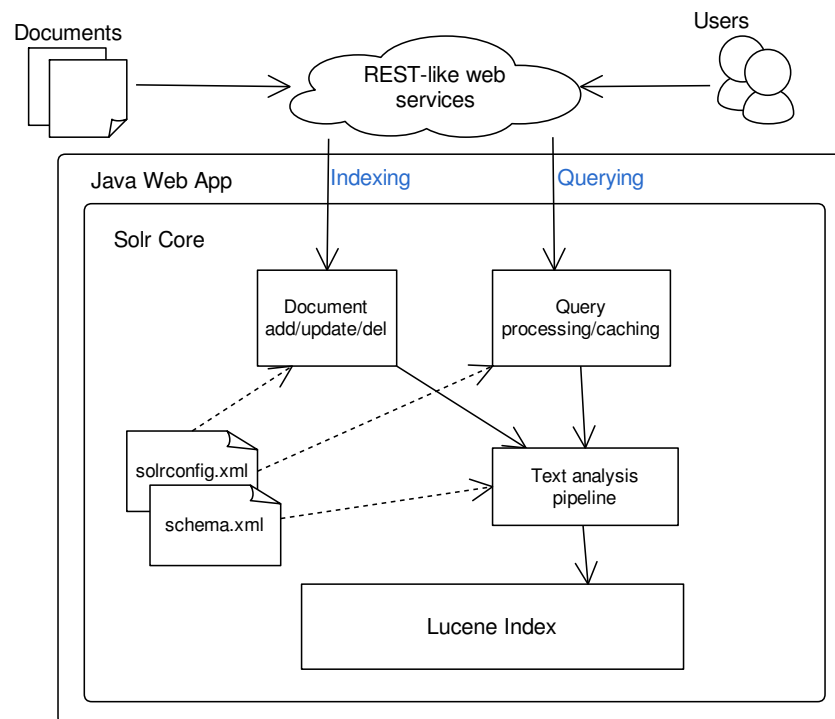


Figure 3.10: High level architecture of a Solr system

## 3.4 Meta Data Storage

As a meta information storage we used MongoDB. MongoDB is a NoSQL document-oriented database. Data in MongoDB is stored in JSON-like documents with a dynamic schema. The format of stored data is called BSON, which stands for Binary JSON. BSON is an open standard developed for human readable data exchange <sup>5</sup>. MongoDB requires a little amount of configuration to start to work with (see Appendix A.2).

## 3.5 Web Search Interface

Web Search interface is a Java web application running in a servlet container. Figure 3.11 shows an architecture of a web search application. The application is based on MVC (model-view-controller) architectural pattern implemented with Struts framework <sup>6</sup>. The application communicates with Solr via Solr REST API and with Mongo database via Java database connector.

### 3.5.1 Citation Search Main Page

The main page of Citation Search presents a simple search interface allowing user to search for citations. Figure 3.12 shows a sample response to the user query “software testing is time-consuming”. As a result a user see a list of documents matching the query. Each document has a citation with a list of bibliographic links supporting this citation. A user can click to “Show context” link to see a text surrounding the citation in the original paper. If the source paper is available online then user can open it using a link “See pdf on SCG resources”.

If a reference has a title recognizable by Citation Search then a user can see all citations referred to the paper from this reference by clicking on the button next to the reference. Figure 3.13 demonstrates this feature. A user can see all citations of the paper “Software maintenance and evolution: a roadmap” in a popover dialog. User can get more information on each citation by following a “View details” link.

User can take advantage of using enhanced search query syntax. The query syntax is explained on the help page of the Citation Search interface and in the Appendix of this article.

### 3.5.2 Search by Bibliography Page

Another feature provided by citation Search is the possibility to search by bibliography entries. For example, a user can search by authors, title or publication venue. An example of a search by author is shown in Figure 3.14. A user see a list of bibliographic entries with matched author’s

---

<sup>5</sup>BSON specification, <http://bsonspec.org/>

<sup>6</sup>Struts framework, <https://struts.apache.org/>

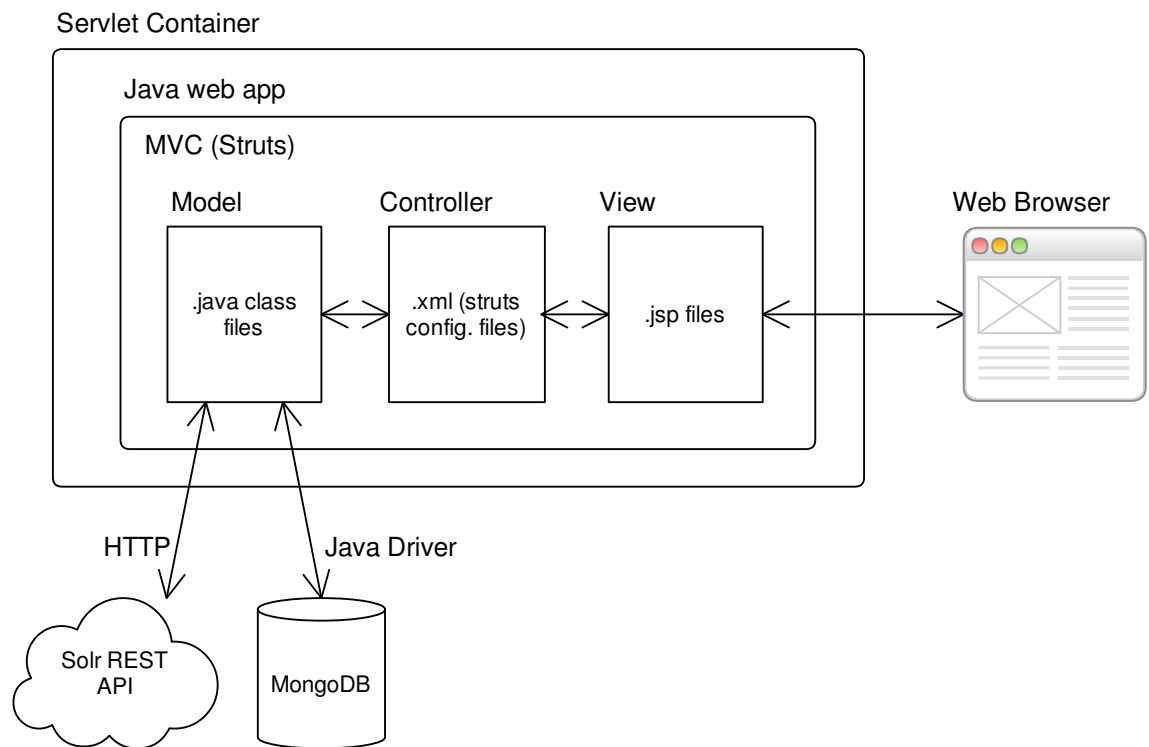


Figure 3.11: Architecture overview of a web search application

## Citation Search Engine ?

Type a statement to search citations or [search by bibliography...](#)

software testing is time-consuming

Search

Page 1 of about 4342.0 pages:

"Generally, testing is the most time-consuming activity during software maintenance [1]."

[Show context](#)

[See pdf on SCG resources](#)

References:

- [1] K. H. Bennett and V. T. Rajlich. Software Maintenance and Evolution: a Roadmap. In Proceedings of the IEEE International Conference on Software Engineering - The Future of Software Engineering, pages 73–87, 2000.

[Google](#)

"It is well accepted that test case generation is one of the most time consuming aspect of software testing [1]."

[Show context](#)

[See pdf on SCG resources](#)

References:

- [1] S. Rapps and E.J. Weyuker, "Selecting Software Test Data Using Data Flow Information," IEEE Transactions on Software Engineering, Vol. SE-11, No. 4, pages 367-375, April 1985. 357 Proceedings of the International Conference on Software Maintenance (ICSM '95) 1063-6773 /95 \$10.00 © 1995 IEEE

[Google](#)

Figure 3.12: A screenshot of the main page of the Citation Search interface showing results for a statement query “software testing is time-consuming”



## Citation Search Engine ?

Type a statement to search citations or [search by bibliography...](#)

software te

Citations of "Software maintenance and evolution: a roadmap"

Search

Page 1 of

"Gener

[Show conte](#)

References:

- [1] K. H
- Future

"It is w

[Show conte](#)

References:

- [1] S. R
- 375, A

"Regres

- Introduction It is widely acknowledged that software maintenance occupies a large fraction of the software development cycle [2, 24], where maintenance includes modifying, extending, debugging, testing, and documenting the application.
- The arguments raised in [4] might point to a solution.
- Hence, whether a system is kept in evolution or downgraded to the servicing stage [4] is a managerial decision.
- For comprehending this aspect of software evolution, empirical knowledge about process and human aspects is still needed [5].
- Introduction and rationale Over the years studies of industrially evolved software have contributed to a body of knowledge of the software evolution phenomenon [1, 3, 14, 15, 16].
- These sheets are subject to similar evolutionary patterns of possibly staged maintenance [3] as conventional software of comparable strategic importance.
- For instance, the object-orientation, often considered the "solution for the software maintenance" [1], created new problems for the maintenance [4] and should be used with very care to assure that the maintenance will not be more problematic than in the traditional legacy systems.
- I. INTRODUCTION Understanding the rationale and design knowledge of a software system is essential when maintaining or evolving the system [1].
- Therefore, as Bohner and Arnold [3] and Bennett and Rajlich [2] point out, a good software architecture should be evolvable, flexible.
- This makes it more and more difficult to evolve the program, and will most likely result in code decay [2].
- Generally, testing is the most time-consuming activity during software maintenance [1].
- Furthermore, it is reported that corrective changes account to only 21% of changes [2], so we cannot conclude that defects are the main reason for change.

[View Details](#)

Figure 3.13: A screenshot of the main page of a Citation Search interface showing citations referred to the article with the recognizable title "Software maintenance and evolution: a roadmap."

name. If an entry has an extractable title then user can see citations from other papers referred to the entry.

Type a phrase to search bibliography...

**Page 1 of about 11 pages:**

Mircea F. Lungu. Reverse Engineering Software Ecosystems. PhD thesis, University of Lugano, 2009.

Mircea Lungu. Reverse Engineering Software Ecosystems. PhD thesis, University of Lugano, November 2009.

Mircea Lungu, Michele Lanza, and Oscar Nierstrasz. Evolutionary and collaborative software architecture recovery with Softwrenaut. Science of Computer Programming (SCP), 2012.

Mircea Lungu, Michele Lanza, Tudor Gîrba, and Romain Robbes. The Small Project Observatory: Visualizing software ecosystems. Science of Computer Programming, Elsevier, 75(4):264–275, April 2010.

Oscar Nierstrasz and Mircea Lungu. Agile software assessment. In Proceedings of International Conference on Program Comprehension (ICPC 2012), pages 3–10, 2012. doi: 10.1109/ICPC.2012.6240507.

Figure 3.14: A screenshot of the ‘search by bibliography’ page of a Citation Search interface illustrating a search for citations based on a meta-information query. Here a user searches for citations of scientific articles authored by Mircea Lungu.

# 4

## Evaluation Experiments

To measure effectiveness of Citation Search Engine we conducted evaluation experiments comparing it with other search engines. We had two main candidates to compare with Citation Search Engine: CiteSeerX and Google Scholar. Preparatory tests showed that CiteSeerX is too slow in showing results. Moreover, users complained that resulting documents are not relevant. Too many results were from different than Computer Science domain, for example, Biology or Physics. Thus all experiments were conducted with Google Scholar and Citation Search Engine.

There are many aspects on how two search engines might be compared. In our experiments we focused on efficiency and usability. By efficiency we imply how quickly users can find documents and by usability we imply simplicity of search interfaces and personal impression.

### 4.1 Experiment Setup

For evaluation experiments we used dataset of scientific articles collected by SCG members over decades. Collection contains for about 16000 scientific articles and covers various topics in Software Engineering and Programming Languages. Dataset of Google Scholar is very large compared to our dataset, so we reduced searching space to the domain of Software Engineering and Programming Languages. Besides we intentionally looked for experts in these domains that can participate in experiments. Nine experts with different experiences (7 PhD candidates, 1

postdoctoral researcher, 1 professor) voluntarily agreed to participate in experiments. Participants were split into two groups. Both groups were asked to perform the same tasks however the second group was asked to do one more additional task. Time given to complete each task was limited to 5 minutes. All tasks will be described below. Each experiment was setup to last for approximately 45 minutes. During the experiment all participants were provided with a laptop and their actions was recorded with a screen casting application.

**Task1** As a first task an expert was asked to find a reference to a citation from one of his articles written in the past using the particular search engine. Figure illustrates one of a citation from the paper given to an expert. An expert can read the cited sentence as well as the context of this sentence however he is not aware of referred source paper. The task is to find a referred paper that proves given citation.

In overall, an expert was given 4 different citations from one of his papers and asked to find references using Google Scholar or Citation Search alternately. Every time the order of using search engines was changed. Before the experiment, an article with a citation was removed from our dataset so an expert could not find an exact match using Citation Search.

During the execution of tasks following observations were recorded:

- Citation and search engine name
- Search time
- Number of queries made
- Number of words in a query
- Expert comments

**Task1'** This task was given to the second group as an extra task. An expert was also given a citation to find a reference to it as in Task 1. But this time a search engine was not specified and a citation was taken from the paper not authored by expert.

Every expert received only one citation for this task. As for the previous task the paper containing citation was removed from Citation Search dataset.

During Task1 we observed which search engine was used to find a reference.

**Task2** For this task we asked every expert in advance to provide with her favorite paper that she thinks is important in her research field. We then built two automatic summaries of this paper. One summary was built using state of the art TextRank algorithm [18]. TextRank algorithm is a graph-based ranking algorithm for Natural Language Processing (NLP) [17]. It extracts sentences from the text based on their importance. We used Python implementation of this algorithm that can be found on GitHub. Another summary was built using citations to the paper collected by

Citation Search. If TextRank provides with the extractive summary, Citation Search collects external opinions about the paper. We limited the size of summaries to the abstract size that is approximately 9-10 sentences.

As a task an expert was given two summaries and asked to give a score from 0 to 10. The score should reflect how summaries are good in describing a paper.

## 4.2 Questionnaires before and after each experiment

Before the beginning of the experiment each expert was asked to fill a form with following questions:

- What search engines do you usually use to find scientific literature?
- How often do you use listed above search engines? (Possible answers: everyday, a few times per week, once a week or less)
- Have you ever used following search engines?(Google Scholar, Citation Search )

After filling the form each expert was introduced to the simple query syntax and to interfaces of search engines. Then experiment begins with Task1 and Task1. After Task1 and Task1 experts were asked some general questions about their impressions on using Google Scholar and Citation Search, for example:

- What did you like, what you did not like about using particular search engine?
- What difficulties did you have?

Right after the experiment experts were asked to fill a form with following questions:

- Has the experiments changed your opinion on the two search engines?
- Would you consider using one of these search engines?

## 4.3 Evaluation Results

First survey, taken by experts before the experiment, showed that all experts use Google Scholar to find scientific literature. Some experts mentioned that they use IEEE Explorer, ACM digital library and DBLP as well. Half of the respondents use Google Scholar every day (Figure 4.1).

Four respondents answered positively on the question if they have ever used Citation Search(Figure 4.2). However all of them mentioned that they used Citaion Search only a few times or rare.

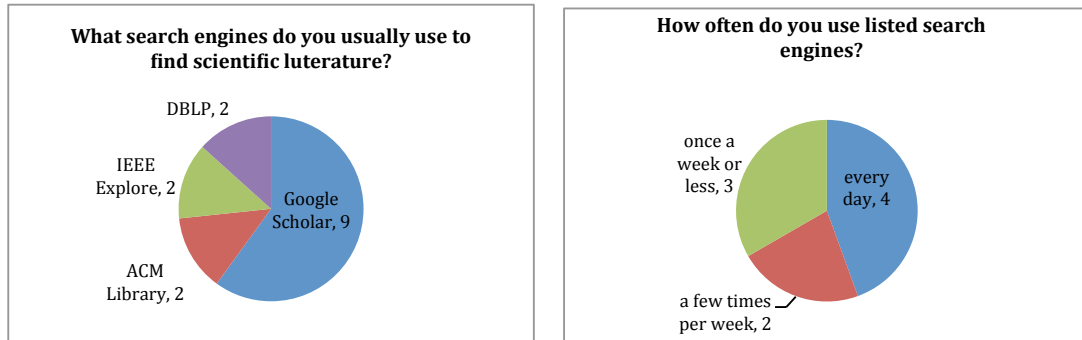


Figure 4.1: Results for questions 1 and 2

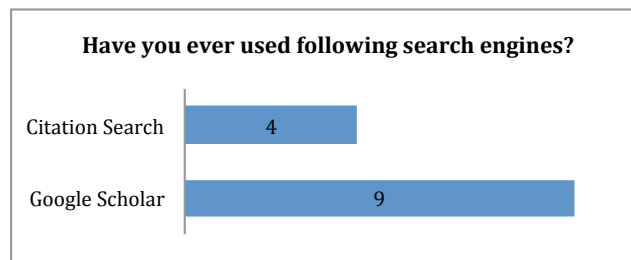


Figure 4.2: results for question 3

### 4.3.1 Results for Task 1

Table 4.1 and Figure 4.3 give results for search time in Task 1. Table 4.1 shows mean and standard deviation values for search time for both search engines. Figure 4.3 illustrates a normal distribution of search time for both search engines. From Table 4.1 and Figure 4.3 we observed that average time to find a reference for a given citation is approximately 2.5 minutes. Experts were slightly faster with finding results using Citation Search. However there is no statistically significant difference between search times of Citation Search and Google Scholar according to ttest with significant level  $p = 5\%$ .

	Mean(sec)	Std(sec)
Citation Search	150	97
Google Scholar	160	78

Table 4.1: Mean and standard deviation for search time in Task 1.

Table 4.2 shows average values and standard deviations for number of queries made by experts to find references. Table 4.3 shows average values and standard deviations for average number of words in queries. We did not see any significant differences in a number of queries and

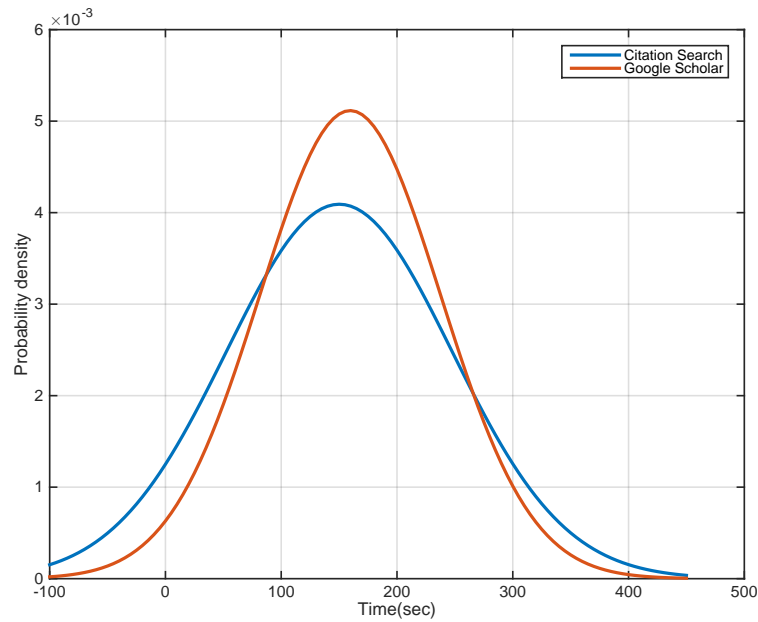


Figure 4.3: Normal distribution fit for a search time for both search engines

average number of words in a query between two search engines. From Table 4.2 we concluded that in average experts made 2-3 queries before finding a referred paper. And from Table 4.3 we concluded that average number of words in a query was 4.

	Mean	Std
Citation Search	2.0	1.5
Google Scholar	2.5	1.5

Table 4.2: Mean and standard deviation values for a number of queries made in Task 1.

During the experiments we noticed that experts were more familiar with Google Scholar search interface so participants spent some time exploring Citation Search interface. This could affect search time for Citation Search making it longer.

We also noted that the way search engines present results is important factor of the search engine usability. For example, most of experts admitted that they like that Citation Search engine shows the exact place from the article where match was found. In contrast, Google Scholar shows a title of the article and a beginning of the abstract, so it is not clear where match was found. In this case experts had to open the article and make a manual search over the text.

	Mean	Std
Citation Search	4.3	1.7
Google Scholar	4.2	1.2

Table 4.3: Mean and standard deviation values for an average number of words in a query in Task 1.

### 4.3.2 Results for Task 1'

Results for Task 1' are shown on Table 4.4. Task 1' was given to four experts. Table 4.4 shows in what search engine a referred paper was found and if an expert tried to use both search engines. From Table 4.4 we concluded that all experts found a referred paper using Citation Search. Meantime three of experts tried both search engines and one expert did not use Goggle Scholar at all.

	Search Engine	Tried to use both SEs
Expert1	Citation Search	Yes
Expert2	Citation Search	Yes
Expert3	Citation Search	Yes
Expert4	Citation Search	No

Table 4.4: Results for Task 1'. Table shows in what search engine a result was retrieved and if an expert tried to used both search engines during the task.

### 4.3.3 Results for Task 2

Results for Task 2 are shown in Figure 4.4. Figure 4.4 illustrates scores from 0 to 10 given by expert to summaries generated with TextRank and citation from Citation Search. Almost all experts except one gave better score to the summary composed by citations from Citation Search. According to ttest with a significance level  $p = 5\%$  there is a significant difference between scores given to summaries generated with TextRank and Citation Search. Experts noted that a summary generated with TextRank tends to consist of too general or not important for understanding a paper sentences. According to experts, summary composed with citations tends to contain sentences more relevant for understanding a paper. However, sometimes sentences rephrase each other not adding new meaning to the summary. Also, compared to TextRank there is no a natural flow in a summary from citations, saying in other words sentences in a summary are not ordered in a story manner.



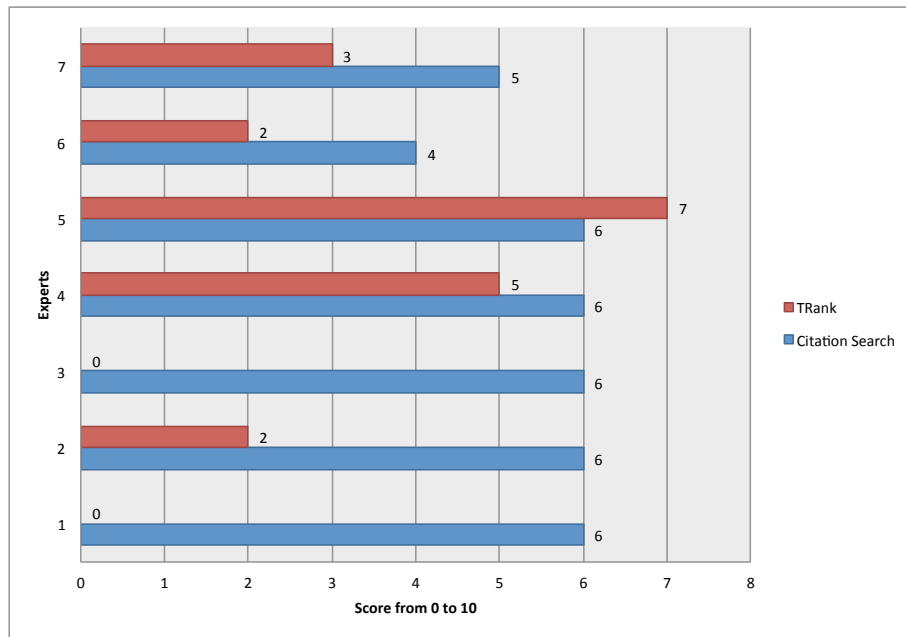


Figure 4.4: Scores from 0 to 10 given by experts in Task 2 to summaries generated with TextRank and citations from Citation Search.

#### 4.3.4 Final Questionnaire

Final questionnaire show that some experts were pleasantly surprised by the capabilities of another type of search engine. All experts answered positively if they are willing to continue to use Citation Search. Some experts specified that they will use Google Scholar and Citation Search for different purposes. According to some experts Citation search is more appropriate to search for a related work on the given topic. Others think that Citation Search is good to prove claims while writing a scientific paper. One expert opinion states that Citation Search is useful for discovering new works in the given domain. Experts highlighted following feature of Citation Search: possibility to see citations with a context, possibility to search by bibliography entries.

#### 4.3.5 Results Summary

In our evaluation experiments we compared Citation Search with Google Scholar. During experiments we collected statistics on search time, number of queries and average number of words in queries. The results show that Citation Search performs slightly better for mean value of search time, however there is no statistically significant difference among search engines. Overall, given that Google Scholar is one of the most popular academic search engines in WWW, Citation Search might be a good alternative to Google Scholar. Indeed, when experts have a possibility

to choose between two search engines, all experts succeeded in the task accomplishment using Citation Search.

The results for summaries comparison show that summary generated with citations give better description of a paper. Automatic citation aggregation feature of Citation Search could be used to generate summaries or even judge the importance of a paper, for example, by counting number of citations.

# 5

## Conclusion

In this work we address the problem of IR for scientific articles. We believe that considering meta-information helps to build enhanced search systems. We particularly focused on citations considering them as important text blocks. We designed and implemented Citation Search engine which automatically extracts and indexes citations from scientific articles in a PDF format. In general, a citation composed of two parts: a cited text in the body of a document and a bibliographic link in a reference section of a document. In Citation Search engine we index both parts separately allowing user to make a choice between different indexed collections. One more feature of Citation Search is the ability to aggregate citations referred to the same source article. This feature could be used to build automatic summaries of scientific articles.

We evaluated our system by conducting user evaluation experiments. In the first part of our experiments, we compared our system with a popular academic search engine Google Scholar. We observed how quickly can users find results using both search engines. Our results showed that Citation Search performs equal or better than Google Scholar. In the second part of our experiments, we used an aggregation feature of Citation Search to build summaries of scientific articles. We compared those summaries with summaries built using a TextRank algorithm. Our results showed that Citation Search gives better description of scientific articles according to expert's opinion.

Nowadays when the amount of information grows fast, efficient IR systems present particular

interest. Our system is different from typical IR systems that performs full-text indexing. Instead we index citations with the intent to reduce search space at the same time making it more effective.

## Bibliography

- [1] Marc Bertin and Iana Atanassova. Semantic enrichment of scientific publications and metadata : Citation analysis through contextual and cognitive analysis. *D-Lib Magazine*, 18:8, 2012.
- [2] Marc Bertin and Iana Atanassova. Extraction and characterization of citations in scientific papers. In Valentina Presutti, Milan Stankovic, Erik Cambria, Ivn Cantador, Angelo Di Iorio, Tommaso Di Noia, Christoph Lange, Diego Reforgiato Recupero, and Anna Tordai, editors, *Semantic Web Evaluation Challenge*, volume 475 of *Communications in Computer and Information Science*, pages 120–126. Springer International Publishing, 2014.
- [3] Shannon Bradshaw. Reference directed indexing: Redeeming relevance for subject search in citation indexes. In Traugott Koch and IngeborgTorvik Slvberg, editors, *Research and Advanced Technology for Digital Libraries*, volume 2769 of *Lecture Notes in Computer Science*, pages 499–510. Springer Berlin Heidelberg, 2003.
- [4] Shannon Glenn Bradshaw. *Reference directed indexing: Indexing scientific literature in the context of its use*. Northwestern University, 2002.
- [5] Isaac G. Councill, C. Lee Giles, and Min yen Kan. Parscit: An open-source crf reference string parsing package. In *International Language Resources and Evaluation*. European Language Resources Association, 2008.
- [6] Gaizka Garechana, Rosa Rio, Ernesto Cilleruelo, and Javier Gavilanes. Visualizing the scientific landscape using maps of science. In Suresh P. Sethi, Marija Bogataj, and Lorenzo Ros-McDonnell, editors, *Industrial Engineering: Innovative Networks*, pages 103–112. Springer London, 2012.
- [7] Eugene Garfield et al. Science citation index-a new dimension in indexing. *Science*, 144(3619):649–654, 1964.

- [8] C. Lee Giles, Kurt D. Bollacker, and Steve Lawrence. Citeseer: An automatic citation indexing system. In *Proceedings of the Third ACM Conference on Digital Libraries*, DL '98, pages 89–98, New York, NY, USA, 1998. ACM.
- [9] The Stanford Natural Language Processing Group. Stanford CoreNLP API. <http://nlp.stanford.edu/software/corenlp.shtml>.
- [10] M. M. Kessler. Bibliographic coupling between scientific papers. *American Documentation*, 14(1):10–25, 1963.
- [11] Richard Klavans and Kevin W. Boyack. Toward a consensus map of science. *J. Am. Soc. Inf. Sci. Technol.*, 60(3):455–476, March 2009.
- [12] Ray R Larson. Bibliometrics of the world wide web: An exploratory analysis of the intellectual structure of cyberspace. In *PROCEEDINGS OF THE ANNUAL MEETING-AMERICAN SOCIETY FOR INFORMATION SCIENCE*, volume 33, pages 71–78, 1996.
- [13] Loet Leydesdorff, Stephen Carley, and Ismael Rafols. Global maps of science based on the new web-of-science categories. *CoRR*, abs/1202.1914, 2012.
- [14] Bertin M, Descls J. P, Djioua B, Krushkov Y, and Lalicc Umr. Automatic annotation in text for bibliometrics use.
- [15] Christopher D Manning, Prabhakar Raghavan, and Hinrich Schütze. *Introduction to information retrieval*, volume 1. Cambridge university press Cambridge, 2008.
- [16] Qiaozhu Mei and ChengXiang Zhai. Generating impact-based summaries for scientific literature. In *Proceedings of ACL-08: HLT*, pages 816–824, Columbus, Ohio, June 2008. Association for Computational Linguistics.
- [17] Rada Mihalcea. Graph-based ranking algorithms for sentence extraction, applied to text summarization. In *Proceedings of the ACL 2004 on Interactive poster and demonstration sessions*, page 20. Association for Computational Linguistics, 2004.
- [18] Rada Mihalcea and Paul Tarau. Textrank: Bringing order into texts. Association for Computational Linguistics, 2004.
- [19] Preslav I. Nakov, Ariel S. Schwartz, and Marti A. Hearst. Citances: Citation sentences for semantic analysis of bioscience text. In *In Proceedings of the SIGIR04 workshop on Search and Discovery in Bioinformatics*, 2004.

- [20] Hidetsugu Nanba and Manabu Okumura. Towards multi-paper summarization reference information. In *Proceedings of the 16th International Joint Conference on Artificial Intelligence - Volume 2*, IJCAI'99, pages 926–931, San Francisco, CA, USA, 1999. Morgan Kaufmann Publishers Inc.
- [21] Enrique Orduña Malea, Juan M. Ayllón, Alberto Martín-Martín, and Emilio Delgado López-Cózar. About the size of google scholar: playing the numbers, July 2014.
- [22] M. F. Porter. Readings in information retrieval. chapter An Algorithm for Suffix Stripping, pages 313–316. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1997.
- [23] Vahed Qazvinian and Dragomir R. Radev. Scientific paper summarization using citation summary networks. In *Proceedings of the 22Nd International Conference on Computational Linguistics - Volume 1*, COLING '08, pages 689–696, Stroudsburg, PA, USA, 2008. Association for Computational Linguistics.
- [24] Jonathan Rochkind. Thinking like solr its not an rdbms. <https://bibwild.wordpress.com/2011/01/24/thinking-like-solr-its-not-an-rdbms/>.
- [25] Henry Small. Co-citation in the scientific literature: A new measure of the relationship between two documents. *Journal of the American Society for Information Science*, 24(4):265–269, 1973.
- [26] Henry Small. Visualizing science by citation mapping. *J. Am. Soc. Inf. Sci.*, 50(9):799–813, July 1999.
- [27] Linda C. Smith. Citation analysis. [https://www.ideals.illinois.edu/bitstream/handle/2142/7190/librarytrendsv30ili\\_%20opt.pdf?sequence=1](https://www.ideals.illinois.edu/bitstream/handle/2142/7190/librarytrendsv30ili_%20opt.pdf?sequence=1), 1981.
- [28] Solr Wiki. Why use solr? <http://wiki.apache.org/solr/WhyUseSolr>.



# User Guide for Citation Search Deployment

## A.1 Solr Installation

Solr installation requires JDK and any servlet container to be installed on the server machine. Here we describe the configuration of Solr for Apache Tomcat container. We need to download Solr distribution that can be found on the official Solr home page <sup>1</sup>. Solr is distributed as an archive. After unzipping the archive, the extracts have following directories:

- **contrib/** - directory containing extra libraries to Solr, such as Data Import Handler, MapReduce, Apache UIMA, Velocity Template, and so on.
- **dist/** - directory providing distributions of Solr and some useful libraries such as SolrJ.
- **docs/** - directory with documentation for Solr.
- **example/** - Jetty based web application that can be used directly.
- **Licenses/** - directory containing all the licenses of the underlying libraries used by Solr.

Copy the dist/solr.war file from the unzipped folder to \$CATALINA\_HOME/webapps/solr.war. Then point out to Solr location of home directory describing a collection:

---

<sup>1</sup>Apache Solr, <http://lucene.apache.org/solr/>



- **Java options:** one can use following command so that the container picks up Solr collection information from the appropriate location:

```
$export JAVA_OPTS="$JAVA_OPTS -Dsolr.solr.home=/opt/solr/example"
```

By a collection in Apache Solr one indicates a collection of Solr documents that represents one complete index.

The Solr home directory contains configuration files and index-related data. It should consist of three directories:

- **conf/** - directory containing configuration files, such as solrconfig.xml and schema.xml
- **data/** - default location for storing data related to index generated by Solr
- **lib/** - optional directory for additional libraries, used by Solr to resolve any plugins

### A.1.1 Solr Configuration

Configuring Solr instance requires defining a Solr schema and configuring Solr parameters.

**Defining Solr schema** Solr schema is defined in the schema.xml file placed in the conf/ directory of the Solr home directory. Solr distribution comes with a sample schema file that can be changed for the needs of the project. The schema file defines the structure of the index, including fields and field types. The basic overall structure of the schema file is:

```
<schema>
  <types>
  <fields>
  <uniqueKey>
  <copyField>
</schema>
```

The basic unit of data in Solr is document. Each document in Solr consists of fields that are described in the schema.xml file. By describing data in the schema.xml, Solr understands the structure of the data and what actions should be performed to handle this data. Here is an example of a field in the schema file:

```
<field name="id" type="integer" indexed="true" stored="true" required="true"/>
```

Table A.1 lists and explains major attributes of field element.

Here is a fragment of schema file defining fields of a document in Citation Search Engine collection:

Name	Description
default	default value if it is not read while importing a document
indexed	true if field should be indexed
stored	when true a field is stored in index store and is accessible while displaying results
compressed	when true a field will be zipped, applicable for text-type fields
multiValued	if true, field can contain multiple values in the same document.

Table A.1: Major attributes of field element in a schema.xml file

```

<fields>
  <field name="_version_" type="long" indexed="true" stored="true"
    multiValued="false"/>
  <field name="id" type="string" multiValued="false"/>
  <field name="text" type="text_en" indexed="true" multiValued="false"/>
  <field name="context" type="string" indexed="false" multiValued="false"/>
  <field name="path" type="string" indexed="false" multiValued="false"/>
  <field name="reference" type="string" indexed="false" stored="true"
    multiValued="true" />
</fields>

```

Every document represents a citation with matching bibliographic references. In the schema file we indicate that we want to index a text field which is the citation text. We store an id of a citation, that is a generated value, calculated from the hash of the citation string. Specifying the id is particularly useful for updating documents. We also store a context for a citation and a path to the scientific article where the citation was found. As a citation can refer to multiple sources, we make the reference field multivalued.

In the schema configuration file, one can define the field type, like string, date or integer and map them to Java classes. This can be handy when we define custom types. A field type includes the following information:

- Name
- Implementation class name
- If the field type is a TextField, it will include a description of the field analysis
- Field attributes

A sample field type description:

```
<fieldType name="text_ws" class="solr.TextField" positionIncrementGap="100">
  <analyzer>
    <tokenizer class="solr.WhitespaceTokenizerFactory"/>
  </analyzer>
</fieldType>
```

Other elements in the Solr schema file listed in Table A.2:

Name	Description
uniqueKey	specifies which field in documents is a unique identifier of a document, should be used if you ever update a document in the index
copyField	used to copy field's value from one field to another

Table A.2: Description of some elements in schema.xml

**Configuring Solr Parameters** To configure a Solr instance we need to describe the solrconfig.xml and solr.xml files.

**solr.xml** The solr.xml configuration is located in solr home directory and used for configuration of logging and advanced options to run Solr in a cloud mode.

**solrconfig.xml** The solrconfig.xml configuration file primarily provides you with an access to index-management settings, RequestHandlers, listeners, and request dispatchers. The file has a number of complex sections and mainly is changed when a specific need is encountered.

### A.1.2 Enhanced Solr Search Features

Solr provides a number of additional features that can enhance the search system. One of the features we use is synonyms. To use this feature you need to specify synonyms.txt file with listed synonyms. This file is used by synonym filter to replace words with their synonyms. For example, a search for "DVD" may expand to "DVD", "DVDs", "Digital Versatile Disk" depending on the mapping in this file. This file can be also used for spelling corrections. Here is an example of synonyms.txt file:

```
GB, gib, gigabyte, gigabytes
MB, mib, megabyte, megabytes
Television, Televisions, TV, TVs
Incident_error, error
```

Additionally, there are other configuration files that appear in the configuration directory. We are listing them in Table A.3 with the description of each configuration:

Name	Description
protwords.txt	file where you can specify protected words that you do not wish to get stemmed. So, for example, a stemmer might stem the word "catfish" to "cat" or "fish".
spellings.txt	file where you can provide spelling suggestions to the end user.
elevate.txt	file where you can change the search results by making your own results among the top-ranked results. This overrides standard ranking scheme, taking into account elevations from this file.
stopwords.txt	Stopwords are those that will not be indexed and used by Solr in the applications. This is particularly helpful when you really wish to get rid of certain words. For example, in the string, "Jamie and joseph," the word "and" can be marked as a stopwords.

Table A.3: Additional configuration files in Solr

## A.2 MongoDB Installation

To install MongoDB follow instruction on the official web site <http://docs.mongodb.org/manual/installation/>.

### A.2.1 MongoDB configuration

Once MongoDB is downloaded, it is very easy to set up a database server. All we need to start the MongoDB server is to type *mongod* command. In our case we would like to specify database location with *--dbpath* parameter and default listening port:

```
> mongod --dbpath /home/aliya/mongodb2 --port 27272
```

MongoDB provides REST API, [t](#)→[.T](#) To enable REST API use parameter *--rest*:

```
> mongod --dbpath /home/aliya/mongodb2 --port 27272 --rest true
```

The simple way to communicate with the MongoDB server is to use the MongoDB shell, in our case we specify *--port* parameter to connect to our instance of MongoDB:

```
> mongo --port 27272
```

Compared to relational databases MongoDB operates with terms collection, which is equivalent to table, and document, which is equivalent to record in relational databases. MongoDB does

not require creating databases and collections explicitly. Databases and collections can be created while starting to use MongoDB. To see list of databases or collections, type `show dbs` in mongo shell:

```
> show dbs
```

MongoDB shell allows to make queries, updates, deletes on collections, get various statistics on data and server usage, and manipulate with data with map-reduce interface, full documentation can be found on the official web site <sup>2</sup>.

### A.3 Running the parser

Before running the parser Solr web application should be deployed on the Tomcat web server and MongoDB instance should be run. One should use Java version 7 or above to run the parser. Get the parser distribution:

```
> git clone git@scg.unibe.ch:citation-search-engine
```

The cloned directory consists of three modules:

- **solr** - Solr related configuration files,
- **citation\_search** - a parser of scientific articles, that extracts meta-information and publish documents to Solr and MongoDB,
- **citation\_search\_web** - a web application for searching citations.

All files related to the parser are located in the *citation\_search* directory. The *citation\_search* directory has a standard Maven project layout <sup>3</sup>. Go to the resources *parser.properties* according to your development environment. Table A.4 describes properties of a *parser.properties* file with sample values.

One can change default logging properties for Log4j <sup>4</sup> library in a *log4j.properties* file. Once property files are configured, build a jar file executing following command from the directory containing a *pom.xml* file:

```
> mvn assembly:assembly -DdescriptorId=jar-with-dependencies -DskipTests
```

Maven will generate a jar file *citation\_search-1.0-jar-with-dependencies.jar* in a *target* folder. To execute the jar file run following command:

```
> java -jar citation_search-1.0-jar-with-dependencies.jar
```

---

<sup>2</sup>MongoDB database, <http://www.mongodb.org/>

<sup>3</sup>Apache Maven, <https://maven.apache.org>

<sup>4</sup>Apache Log4j, <http://logging.apache.org/log4j/2.x/>

Property	Description	Sample value
solr.url.citations	Endpoint for publishing citations.	http://localhost:8088/solr/collection1/
solr.url.bibliography	Endpoint for publishing bibliographic links.	http://localhost:8088/solr/collection2/
db.host	MongoDB host server IP address.	127.0.0.1
db.port	MongoDB listening port.	27272
db.name	MongoDB database name.	CS
db.collection	MongoDB database collection name.	papers
pdfs.path	Location of pdf files.	/home/aliya/Library

Table A.4: Explanation of properties of a *parser.properties* file.

## A.4 Search Interface Deployment

All web application related web files are located in a *citation\_search\_web* directory. The directory has a standard Maven project layout <sup>5</sup>. Change a *search.properties* file in a *resources* folder. Table A.5 describes properties of a *search.properties* file with sample values .

Property	Description	Sample value
solr.url.citations	Endpoint for querying citations.	http://localhost:8088/solr/collection1/
solr.url.bibliography	Endpoint for querying bibliographic links.	http://localhost:8088/solr/collection2/

Table A.5: Explanation of properties of a *search.properties* file.

One can change default logging properties for Log4j <sup>6</sup> library in a *log4j.properties* file. Once property files are configured, build a war file executing following command from the directory containing a *pom.xml* file:

```
> mvn package -DskipTests
```

Maven will generate a war file in a *target* folder. Deploy to the Tomcat web server by putting a warfile in a Tomcat *webapp* directory or use a Tomcat web interface to deploy through Tomcat's manager.

<sup>5</sup>Apache Maven, <https://maven.apache.org>

<sup>6</sup>Apache Log4j, <http://logging.apache.org/log4j/2.x/>