



Jakarta Lucene

Eine Java-Bibliothek zur Suchindex-Erstellung

Seminararbeit
Tilman Schneider





Agenda

- Definition: Suchmaschine
- Vorstellung von Jakarta Lucene
- Erstellung eines Suchindex
- Suchen mit dem Suchindex
- Fazit

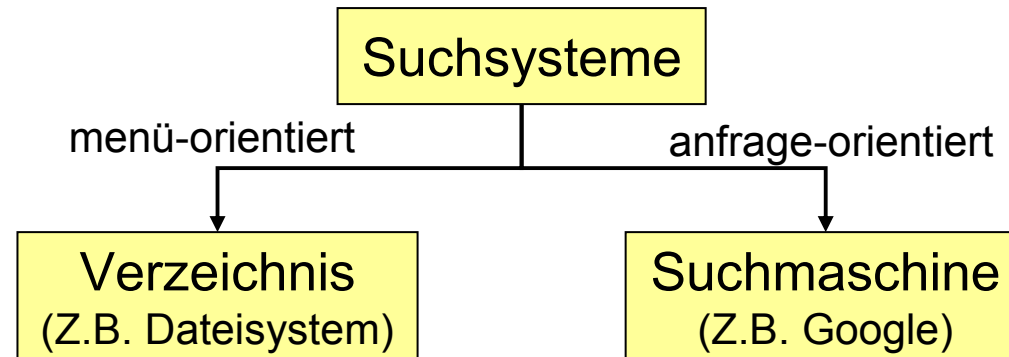


Agenda

- **Definition: Suchmaschine**
- Vorstellung von Jakarta Lucene
- Erstellung eines Suchindex
- Suchen mit dem Suchindex
- Fazit



Definition: Suchmaschine



„Suchmaschine: Ein System zur Volltextsuche in großen Datenmengen“



Agenda

- Definition: Suchmaschine
- **Vorstellung von Jakarta Lucene**
- Erstellung eines Suchindex
- Suchen mit dem Suchindex
- Fazit



Was ist Jakarta Lucene?

- Projekt der Apache Jakarta Gruppe
- 100% pure Java.
Mittlerweile gibt es auch Implementierungen in C++, .NET, Python und Perl.
- Vielseitig einsetzbar
- Keine Applikation, sondern API



Agenda

- Definition: Suchmaschine
- Vorstellung von Jakarta Lucene
- **Erstellung eines Suchindex**
- Suchen mit dem Suchindex
- Fazit

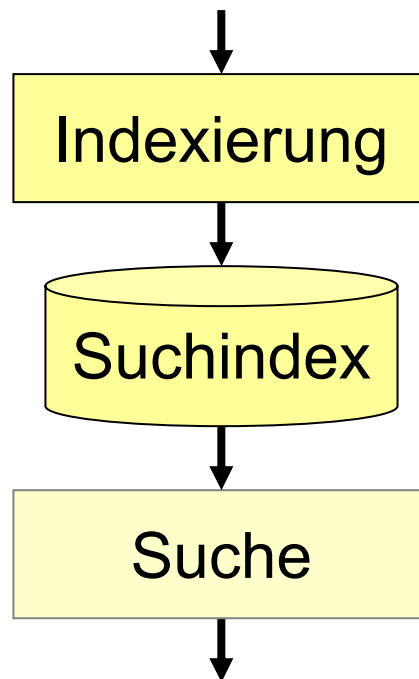


Wozu ein Suchindex?

- Herkömmliche Volltextsuche sehr langsam: $O(n)$, n = Größe des Datenbestands
- Daher: Aufbereitung des Suchraums in einer geeigneten Datenstruktur



Gesamtablauf



- Durchforsten des Suchraums nach Texten
- Aufbereitung des Textes
- Eigentliche Indexierung

- Interpretation der Anfrage
- Eigentliche Suche
- Ergebnisaufbereitung



Umgekehrter Wortindex

Originaltext

Blaukraut bleibt Blaukraut und
Brautkleid bleibt Brautkleid



blaukraut	0, 17
bleibt	10, 42
brautkleid	31, 49
und	27

Suchindex



Textaufbereitung

Durchforsten des Suchraums nach Texten



Textaufbereitung

Nicht übel, sprach der Dübel!

Nicht

übel

sprach

der

Dübel

Zerlegung

übel

sprach

Dübel

Stopword-Filter

ueb

sprach

dueb

Stemming



Indexierung



Lucene: Documents

- Für jeden zu indizierenden Text wird ein Document-Objekt erzeugt.
- Ein Document ist eine Menge von Feldern
- Ein Feld ist ein Key/Value-Paar
- Beispiele für mögliche Felder:
 - Der Text selbst
 - Die URL des Textes
 - Das Datum der letzten Änderung
 - Die Größe des Textes



Lucene: Fields

- Für jedes Feld wird angegeben:
 - Ob es gespeichert wird (store-Flag)
 - Ob es indexiert wird (index-Flag)
 - Ob es vor der Indexierung zerlegt wird (token-Flag)

Fabrikmethode	store	index	token
<code>Field.Text(String name, String value)</code>	x	x	x
<code>Field.Text(String name, Reader value)</code>		x	x
<code>Field.Keyword(String name, String value)</code>	x	x	
<code>Field.UnIndexed(String name, String value)</code>	x		
<code>Field.UnStored(String name, String value)</code>		x	x



Codebeispiel: Datei-Indexierung 1

```
private static void addToIndex(IndexWriter index, File file) {
    InputStream is = new FileInputStream(file);

    // Wir erzeugen ein Document mit zwei Feldern:
    // Eines mit dem Pfad der Datei und eines mit seinem Inhalt
    Document doc = new Document();
    String fileName = file.getAbsolutePath();
    doc.add(Field.UnIndexed("path", fileName));
    doc.add(Field.Text("body", new InputStreamReader(is)));

    // Document zu Index hinzufügen
    index.addDocument(doc);
    is.close();
}
```



Codebeispiel: Datei-Indexierung 2

```
public static void main(String[] args) throws Exception {
    // Neuen Index erzeugen
    String indexPath = args[0];
    IndexWriter index
        = new IndexWriter(indexPath, new GermanAnalyzer(), true);

    // Gegebenes Verzeichnis durchsuchen
    File dir = new File(args[1]);
    crawlDirectory(index, dir);

    // Index optimieren und schließen
    index.optimize();
    index.close();
}

private static void crawlDirectory(IndexWriter index, File dir) {
    // Alle Dateien und Unterverzeichnisse zum Index hinzufügen
    File[] childArr = dir.listFiles();
    for (int i = 0; i < childArr.length; i++) {
        if (childArr[i].isDirectory()) {
            crawlDirectory(index, childArr[i]);
        } else {
            addToIndex(index, childArr[i]);
        }
    }
}
```

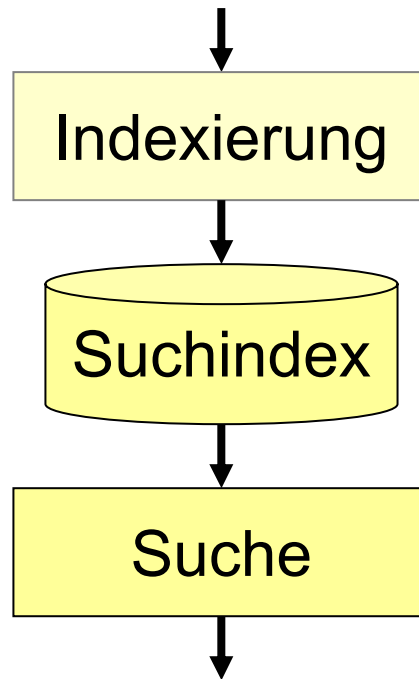


Agenda

- Definition: Suchmaschine
- Vorstellung von Jakarta Lucene
- Erstellung eines Suchindex
- **Suchen mit dem Suchindex**
- Fazit



Gesamtablauf



- Durchforsten des Suchraums nach Texten
- Aufbereitung des Textes
- Eigentliche Indexierung

- Interpretation der Anfrage
- Eigentliche Suche
- Ergebnisaufbereitung

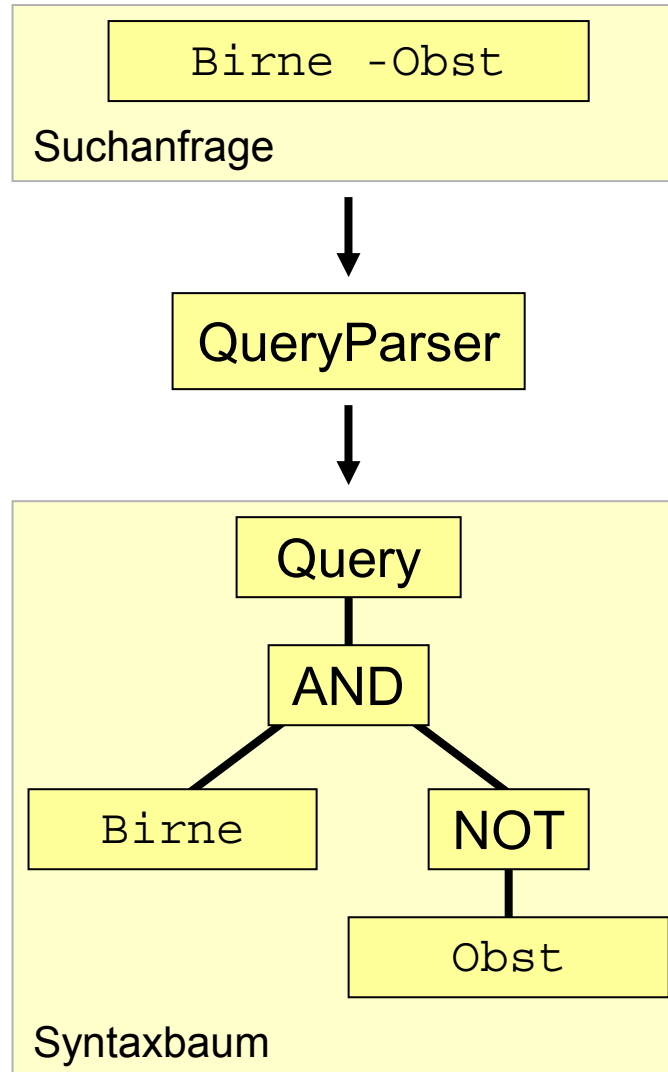


Lucene Querysyntax

- Lucene bietet mächtige QuerySyntax:
 - Bool'sche Operatoren:
„Lucene and Seminar“ oder „+Lucene -Seminar“
 - Suche in Feldern: „Lucene title:Seminararbeit“
 - Gruppierung: „(Jakarta and Lucene) or Seminar“
 - Wildcards: „te?t“ oder „text*“
 - Fuzzy-Suche: „Maier~“
 - ...

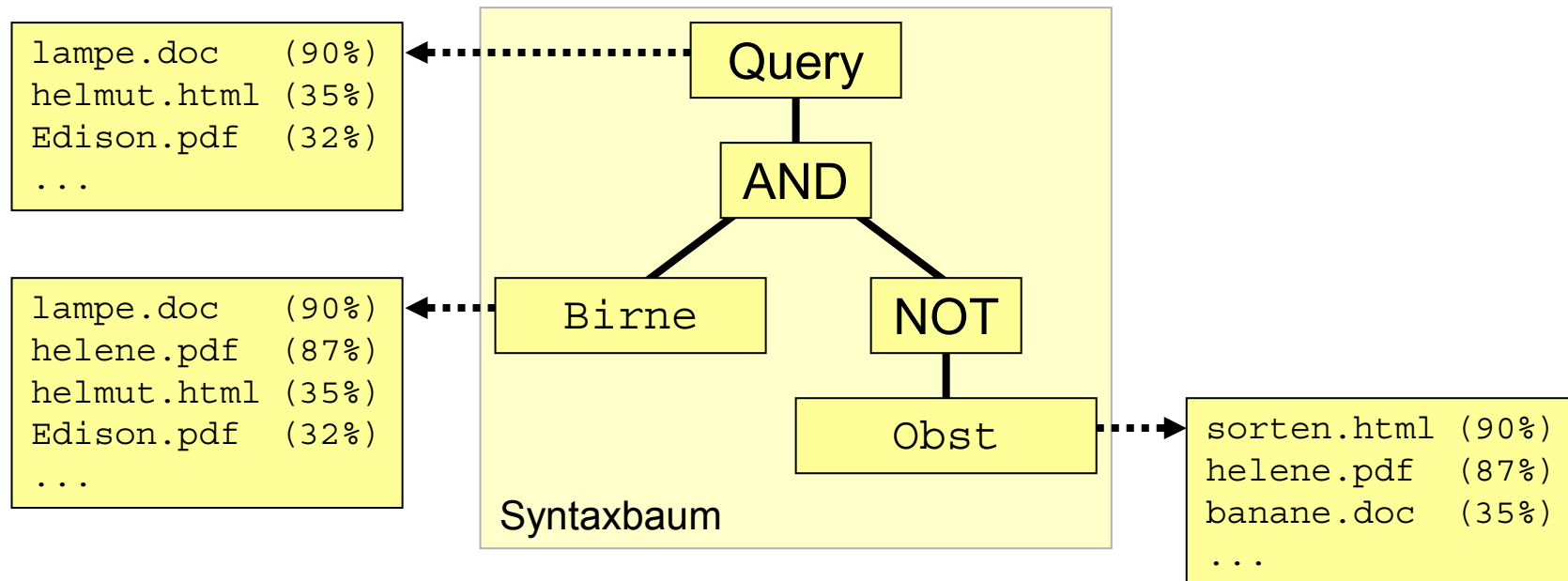


Interpretation der Suchanfrage



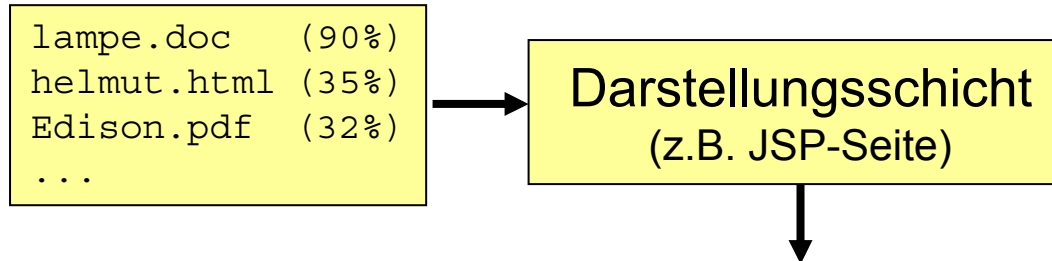


Eigentliche Suche





Ergebnisaufbereitung



Suche nach java tutorial - Mozilla Firefox

File Edit View Go Bookmarks Tools Help

http://localhost:8080/LuceneSearch/SearchOutput.jsp?index=main&query=java+tutorial&maxres

Suchen nach: Search

Die dm-Seiten wurden nach **java tutorial** durchsucht. Ergebnisse 1-5 von insgesamt **1886**. Suchdauer: **0.03** Sekunden.

[Java IDL: Tutorial](#) (Relevanz: 100%)
Java IDL: Tutorial Tutorial Hello World Tutorial Home Concepts Programming Reference Tutorial Copyright © 1996-98 Sun Microsystems, Inc., 2550 Garcia Ave., Mtn. View, CA. 94043-1100 USA., All...
file://p:/www_data/cheestah_docs/java_apis/jdk131/guide/idl/examples.html - 807 Byte

[Java IDL: Tutorial](#) (Relevanz: 100%)
Java IDL: Tutorial Tutorial Hello World Tutorial Home Concepts Programming Reference Tutorial Copyright © 1996-98 Sun Microsystems, Inc., 2550 Garcia Ave., Mtn. View, CA. 94043-1100 USA., All...
file://p:/www_data/cheestah_docs/java_apis/jdk141/guide/idl/examples.html - 807 Byte

[Search The Java Tutorial](#) (Relevanz: 94%)
Search The Java Tutorial The Java TM Tutorial Search the Java TM Tutorial What are you looking for? Help on searching is available. All of the material in The Java Tutorial is copyright...
file://p:/www_data/cheestah_docs/Java-Tutorial/search.html - 2,44 kB

[The Java Tutorial Books](#) (Relevanz: 91%)
The Java Tutorial Books The Java Tutorial About the Java Tutorial Books > The Java Tutorial Continued CD-ROM Contents index | TOC | CD-ROM | Errata The CD-ROM that accompanies The...
file://p:/www_data/cheestah_docs/Java-Tutorial/books/continued/cdinfo.html - 17,11 kB

[The Java Tutorial Books](#) (Relevanz: 91%)
The Java Tutorial Books The Java Tutorial About the Java Tutorial Books > The JFC Swing Tutorial CD-ROM Contents Index | TOC | CD-ROM | Errata The Java Tutorial CD-ROM is...
file://p:/www_data/cheestah_docs/Java-Tutorial/books/swing/cdinfo.html - 10,38 kB

Ergebnisseite: [1](#) [2](#) [3](#) [4](#) [5](#) [6](#) [7](#) [8](#) [9](#) [10](#) [11](#) [12](#) [13](#) [14](#) [15](#) [Weiter >>](#)



Codebeispiel: Suche

```
public class SimpleSearcher {  
  
    public static void main(String[] args) throws Exception {  
        String indexPath = args[0];  
        String queryString = args[1];  
  
        Query query = QueryParser.parse(queryString, "body",  
                                        new GermanAnalyzer());  
  
        Searcher searcher = new IndexSearcher(indexPath);  
        Hits hits = searcher.search(query);  
  
        for (int i = 0; i < hits.length(); i++) {  
            System.out.println(hits.doc(i).get("path")  
                               + " (score: " + hits.score(i) + ")");  
        };  
    }  
}
```



Agenda

- Definition: Suchmaschine
- Vorstellung von Jakarta Lucene
- Erstellung eines Suchindex
- Suchen mit dem Suchindex
- **Fazit**



Fazit

- **Vorteile von Lucene:**
 - Einfache API
 - Trotzdem hohe Funktionalität
 - Sehr ressourcenschonend (Kleiner Index, Wenig Ansprüche an CPU und Arbeitsspeicher)
 - Alle Module vom Analyzer über den QueryParser bis zur Persistenz des Index sind austauschbar.
 - Vielseitig: Lässt sich in jedes Projekt integrieren, das eine Suchfunktion benötigt.
- **Nachteil:**
 - Lucene bietet nur den Kern einer Suchmaschine, die Darstellung des Suchergebnis und die Aufbereitung des Suchraums muss man selbst übernehmen.



Regain

- Open Source Projekt
- Bietet Crawler für Dateisystem und HTTP
- Kann viele geläufige Dateiformate lesen: HTML, XML, Excel, Powerpoint, Word, PDF and RTF
- Bietet TagLibrary für einfache Einbindung der Suchergebnisse in die eigene JSP-basierte Website.



Referenzen

- **Lucene:** <http://jakarta.apache.org/lucene>
- **Regain:** <http://regain.sf.net>
- **Java World-Artikel über Lucene:**
<http://www.javaworld.com/javaworld/jw-09-2000/jw-0915-lucene.html>
- **Lucene-Tutorial:** <http://darksleep.com/lucene>