

Die Suche in OpenCms

Die Suche in OpenCms basiert auf dem Lucene Framework. Die Jakarta Lucene ist eine leistungsstarke, voll funktionsfähige Text-Such-Engine, geschrieben in Java. Sie ist eine Technologie, die für fast jede mögliche Anwendung verwendbar ist, die Volltextsuche erfordert, besonders eine Cross-Plattform. Lucene ist allerdings keine fertige Suchmaschine, sondern stellt Klassen und Funktionen zur Verfügung, um für beliebige Projekte eine eigene Suchmaschine zu bauen.

Lucene unterstützt standardmäßig reichhaltige Suchoptionen. Allerdings besteht auch die Möglichkeit, eine eigene Suchsyntax zu entwickeln.

1. Operatoren

Welche Suchmöglichkeiten und Operatoren bietet die Lucene Engine in OpenCms.

Operatoren	Lucene Search Engine in OpenCms
<p>Fuzzy Suche</p> <p>undeutliche Suche zu formulieren mit Tilde-Symbol ("~")</p> <p>Beispiel: seit~liefert: Seite, seinen, weitere, Leiter...</p>	Ja
<p>Boolesche Suche</p> <p>Verknüpfen von Wörtern mit UND, OR, AND</p> <p>Beispiel: <i>Guten AND Tag</i> + für soll vorkommen</p> <p>Beispiel: oder aber <i>+Tag + guten</i> NOT nicht vorkommen</p>	Ja
<p>Proximity-Suche/Distanzsuche</p> <p><i>"Guten wiedersehen"~10</i> (zeigt Ergebnisse, indem die beiden Wörter maximal 10 Wörter zwischen sich haben).</p>	Ja
<p>Wildcard</p> <p>Platzhalter (?) z. B. <i>te?t</i></p>	Ja, aber nicht am Anfang des Wortes: z. B. ?wort, ?gang geht nicht z. B. Haus?, Baum? geht

<p>Feldsuche</p> <p>Lucene unterstützt Felddaten.</p> <p>Beispiel: <i>"title:index AND text: Tag"</i>.</p>	<p>Ja, z. B. Suche kann über Volltext oder Seitentitel gesteuert werden.</p> <p>Dazu muss das Feld indiziert sein.</p>
<p>Bereichssuche</p> <p>Ein Bereich kann die Zeit zwischen Erstelldatum und Bearbeitungsdatum sein oder zwischen 1. Januar 2010 und 8. Januar 2010 inklusive Start- und Enddatum.</p> <p>oder Beispiel: Findet alle Dokumente ohne den Begriff Wissensmanagement.</p>	<p>Ja, ist möglich.</p> <p>created: [20100101 TO 20100108]</p> <p>lastmodified: [20100101 TO 20100108]</p> <p>title:{Wissensmanagement}</p>
<p>Verstärkungsfaktor</p> <p>um einen Verstärkungsfaktor zuzuordnen, wird das "^"-Symbol verwendet</p> <p>Tag^4 du = Tag ist 4 x stärker gewichtet</p>	<p>Ja</p>
<p>Weitere Funktionen</p>	
<p>Tokenisierung</p> <p>Zerlegen eines oder mehrerer Texte in kleine Einheiten, so genannte Tokens. Diese werden anschließend gespeichert in der Token-Dokument-Beziehung.</p>	<p>Ja</p>
<p>Indexierung</p> <p>Mit Indexierungsverfahren wird durch die Ermittlung von Worthäufigkeiten und Wortrelevanz eine Auswahl getroffen und somit Wörter in den Index aufgenommen.</p> <p>Anhand des erstellten Index kann dann die Volltextsuche, wie später beschrieben, effizient durchgeführt werden.</p>	<p>Ja</p> <p>Bei der Installation von OpenCms erstellt die Lucene Engine automatisch aus den Inhalten in der Datenbank einen Index.</p> <p>Neue Dokumente werden sofort indiziert.</p> <p>Änderungen werden auch sofort berücksichtigt.</p> <p>Lucene untersucht den von OpenCms übergebenen Text und schreibt Informationen über die Häufigkeit der im Text enthaltenen Worte sowie deren Relevanz in den Index.</p> <p>Das hat den Vorteil, dass bei einer Suche nicht das gesamte Dokument durchsucht werden muss (bei großen Dokumenten würde das sehr lange dauern), sondern nur nach dem Wort in dem Index geschaut werden muss und sofort die Trefferstelle und Häufigkeit feststeht.</p>

<p>Anfrageanalyse oder effiziente Suche</p> <p>Analyse von Text-, PDF-, Word-, Excel- und OpenOffice-Dokumenten.</p>	<p>Für jedes Dokumentformat existiert ein „Analyzer“. Dieser durchsucht und indexiert in optimierter Weise verschiedene Dokumentformate. OpenCms unterstützt das.</p>
<p>Rekursives Crawling</p> <p>Suche durch mehrere Verzeichnisse und externe Laufwerke</p>	<p>Nein, dazu muss die Lucene Engine erweitert werden.</p>
<p>Paging</p> <p>Blättern durch Ergebnislisten.</p>	<p>Ja</p>
<p>TagCloud</p> <p>Eine Schlagwortwolke –ist eine Methode zur Informationsvisualisierung, bei der eine Liste aus Schlagworten, oft alphabetisch sortiert, flächig angezeigt wird, wobei Wörter, die häufiger gesucht werden größer angezeigt werden als weniger gesuchte.</p>	<p>Standardmäßig gibt es sowas aber in OpenCms noch nicht.</p>
<p>Semantische Suche – Beispiel</p> <p>Eine Suchanfrage nach ‚Baum‘ an eine bedeutungerschließende Suchmaschine verwendet zur Suche auch Begriffe, die im Zusammenhang mit ‚Baum‘ genannt werden, auch wenn sie in der Anfrage selbst nicht genannt werden. Es werden Ergebnisse, die die Wörter (z. B. Ahorn, Eiche, Linde) beinhalten, angezeigt.</p>	<p>Nein</p>

2. Verhalten von OpenCms und Lucene

1. OpenCms erkennt, dass eine neue Ressource in die Suche aufgenommen werden soll.
2. OpenCms schaut, um welche Ressourcenart es sich handelt und extrahiert suchrelevante Daten - bei XML-Content z. B. Titel, Body, Dateipfad des Anhangs etc., bei Bildern z. B. Bildname, Bildpfad etc. - als Text.
3. OpenCms ruft Lucene auf und übergibt die aus der neuen Ressource suchrelevanten extrahierten Daten; dabei kann OpenCms noch angeben, ob diese eine besondere Gewichtung - "boost" - haben sollen. Außerdem wird Text, welcher aus wenigen Einzelworten besteht, höher gewichtet, so z. B. beim Titel, der oft nur aus einem Wort besteht, im Gegensatz zum Content, der aus Tausenden Wörtern bestehen kann.
4. Lucene untersucht den von OpenCms übergebenen Text und schreibt Informationen über die Häufigkeit der im Text enthaltenen Worte sowie deren Relevanz in den Index.

3. Indexierung

Die Indexierung dient der Aufbereitung des Suchraums und wird typischerweise zu Beginn und bei Änderungen des Suchraums durchgeführt. Anhand des erstellten Index kann dann die Volltextsuche effizient durchgeführt werden.

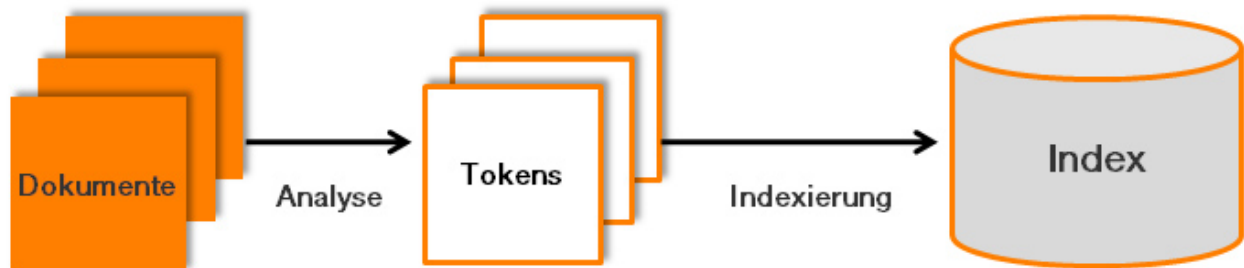


Abbildung: Allgemeiner Aufbau einer Volltextindexierung

Zunächst zum Begrifflichen: Ein Index mehrerer Dokumente bezeichnet im Allgemeinen ein Abbild von Elementen (Tokens oder Terme) dieser Dokumente auf eine Liste der Dokumente, in denen das jeweilige Token vorkommt. Unter der Indexierung versteht man die Generierung eines Index aus einem oder mehreren Texten. Der Indexierungsprozess in Lucene benötigt die folgenden Grundkonzepte mit entsprechenden Java-Klassen:

- **Dokument:** Die von Lucene zu indexierenden Dokumente sind in der Klasse `org.apache.lucene.document.Document` abstrahiert. Ein Dokument kann hierbei alles von einer Buchseite über eine Text- oder HTML-Datei bis zu einem Datenbankeintrag umfassen.
- **Feld:** Ein Dokument besteht aus mehreren (benannten) Feldern (`org.apache.lucene.document.Field`), die auf Wunsch indexiert werden können. Diese Aufteilung ermöglicht eine feiner strukturierte Suche. Beispielsweise könnte man bei einem Word-Dokument neben dem eigentlichen Text als zusätzliche Felder Autor und Beschreibung indexieren.

Textvorverarbeitung (Analyse): Der zu indexierende Text aus den Feldern der Dokumente wird nicht direkt in den Index geschrieben, sondern zunächst von einem sogenannten Analyzer, einem Objekt vom Typ `org.apache.lucene.analysis.Analyzer`, geparkt. Dadurch wird der in den Index einzutragende Text vorverarbeitet - das kann von so trivialen Dingen wie Normalisierung auf Kleinbuchstaben über das Weglassen bestimmter trivialer Wörter bis hin zum Eintrag nur des Wortstammes reichen. Alle Suchanfragen werden dann auch vom entsprechenden Analyzer vorverarbeitet.

Indexerstellung: Der eigentliche Index wird von der Klasse `org.apache.lucene.index.IndexWriter` geschrieben. Dazu muss man angeben, wo der Index gespeichert werden soll - sei es direkt im Dateisystem, in einer Datenbank oder nur im Arbeitsspeicher - und welcher Analyzer verwendet werden soll. In den Index lassen sich dann mehrere Dokumente eintragen. Auch können mehrere Indizes zusammengefügt und optimiert werden.

Indexierung hat den Vorteil, dass bei einer Suche nicht das gesamte Dokument durchsucht werden muss, sondern nur nach dem Wort in dem Index geschaut werden muss und sofort die Trefferstelle und Häufigkeit feststeht. Bei großen Dokumenten bzw. Dokumentenmengen würde das sehr lange dauern und die Performance verschlechtern.

4. Die eigentliche Suche

1. Benutzer gibt in OpenCms einen Suchbegriff ein.
2. OpenCms gibt den Begriff an Lucene weiter.
3. Lucene bereitet den Begriff ggf. auf (Operatoren wie AND, OR, Wildcards etc.).
4. Lucene holt sich Informationen über den Begriff aus dem Index und gibt eine Liste von Treffern zurück.
5. Die Liste ist dann so geordnet, dass relevante Treffer am Anfang stehen.

4.1. Darstellung der Suchergebnisse

- Standardmäßig werden in OpenCms die Suchergebnisse nach relevanten Treffern gelistet.
- Standardmäßig ist das Suchwort markiert oder hervorgehoben.
- Quelle ist ersichtlich, wenn die Suche über mehrere Websites geht.
- Standardmäßig wird die Trefferanzahl angezeigt.

4.2. Suchvorschläge

Standardmäßig macht OpenCms keine Suchvorschläge. Dazu muss der Suchindex oder die Metadatenverwaltung an die Sucheingabemaske angebunden werden. Nur so wird während des Schreibens in das Suchfeld ein Wort aus dem Index oder der Metadatenverwaltung vorgeschlagen, das dort indexiert bzw. aufgenommen wurde. Den Vorschlägen liegt eine Technologie zugrunde, die sich daran orientiert, was der Nutzer mit seiner Suchanfrage am ehesten gemeint haben könnte. Die bereitgestellten Vorschläge laufen dabei im Hintergrund mit und werden erst eingeblendet, wenn der Nutzer das Drop-Down-Menü direkt unter dem Suchfenster öffnet. Die angezeigten weiteren Begriffe und Themenvorschläge sollen dabei helfen, das gesuchte Thema noch präziser einzugrenzen. Ziel wäre es, den User bei seiner Suche schneller und bequemer ans Ziel zu führen.

5. Option Solr

Für eine hochskalierbare Suche kommt die führende Open Source Software für Enterprise Search Applications Solr zum Einsatz. Solr ist ein Enterprise Search Server auf der Basis der Lucene Java-Bibliothek und kann somit unabhängig vom Portal betrieben werden. Dies führt zu einer besseren Performance, Skalierbarkeit, Lastverteilung und Wartbarkeit von Solr. Zu den Schnittstellen von Solr gehört beispielsweise eine HTTP-API, mit der Dokumente hinzugefügt, geändert oder gelöscht werden können. Zu den weiteren Funktionen gehören XML/HTTP und JSON APIs, Hit-High-Lighting, facetierte Suche, Caching, Replikation sowie eine Web-Administrations-Oberfläche.

Im Endeffekt steht Solr zwischen Client-Anwendung (z. B. OpenCms) und Lucene. Der Vorteil ist, dass man dadurch die Suche vereinheitlichen kann. Man kann OpenCms an Solr anschließen und dadurch einen Suchindex erzeugen lassen, welcher Dokumente beider Systeme enthält.

5.1. Vorteile von Solr

- **Individualisierung:** Flexible Anpassung der Suchalgorithmen.
- **Transparenz:** offene API, Protokolle, Formate und Suchalgorithmen.
- **Portabilität:** Lucene/Solr läuft auf allen Plattform-Systemen, welche Java unterstützen; die erstellten Indizes sind unabhängig vom Plattform-System und können somit ohne Probleme zwischen verschiedenen Plattformen ohne Anpassungen portiert werden.
- **Sicherheit:** Solr wird bereits in geschäftskritischen Anwendungen bei über 4.000 Unternehmen weltweit eingesetzt, darunter Branchengrößen wie MySpace, AOL, Nike, LinkedIn oder Monster.com.
- **Performance:** Schnelle Antwortzeiten (intern oft unter 50 ms), da keine Datenbankzugriffe nötig sind; verbessert die Geschwindigkeit Ihrer Seite und gleichzeitig das Ranking in Suchmaschinen.
- **Skalierbarkeit:** Auch bei wachsenden Datenbeständen sind keine überproportionalen Investitionen in Hardware nötig - dies schont Ihre IT-Budgets; große Anwender können von den Replikationsmöglichkeiten und Load-Balancer-Systemen für Solr profitieren.

comundus realisiert auf Basis von OpenCms Mitarbeiterportale und Internetauftritte. In beiden Unternehmenslösungen spielt die Suche eine zentrale Rolle. Eine intelligente Suchfunktion kann darüber entscheiden, ob aus Besuchern Kunden werden und Mitarbeiter effizient und effektiv ihr Wissen teilen können.



comundus GmbH

Heerstraße 111

D-71332 Waiblingen

Telefon: +49 7151 96528-0

E-Mail: info@comundus.com

Internet: www.comundus.com