

Thema: **Web-Page Crawler auf der Basis von Konzepten von Cho, Molina, Page**  
 Veranstaltung: Google tanzt. - Zur Anatomie einer Suchmaschine  
 Dozent: Professor Stock  
 Referent: Daniel Ritter  
 Semester: Sommersemester 2004

## Was ist ein Crawler ?

<p>Ein Crawler ist ein <b>Softwareprogramm</b></p> <p>Ein Crawler ruft <b>automatisiert</b> Webseiten auf und findet auf ihnen <b>Hypertext-Links zu weiteren Webseiten</b>, welche danach ebenfalls in einer <b>festgelegten Reihenfolge aufgerufen</b> werden.</p> <p>Dies findet nach vorher festgeschriebenen <b>Regeln</b> statt.</p> <p>Ein Crawler ist an einen <b>Index</b> gekoppelt, in welchen er die gefundenen Seiten einspeist.</p> <p>Ein Crawler kann auch zur <b>Index-Pflege</b> benutzt werden, indem er nicht mehr existierende Seiten aus dem Index <b>entfernt</b> oder geänderte und umgezogene Seiten im Index <b>aktualisiert</b>.</p>	<p>Schritt 1 Übergabe einer Start-URL an den Crawler</p> <p>Schritt 2 Einlesen des Inhaltes der Start-URL. Schreiben der gefundenen Links in einen Speicher. Weitere individuelle Schritte für den spezifischen Crawler</p> <p>Schritt 3 <b>Einlesen des Inhaltes einer der URL's aus dem URL-Speicher.</b> Schreiben der gefundenen Links in den Speicher. Weitere individuelle Schritte für den spezifischen Crawler</p> <p>Schritt 4 Überprüfen, ob der Speicher leer ist oder eine andere Abbruchbedingung erfüllt wurde. Falls weiter gecrawlt werden soll, Rücksprung zu Schritt 3. Ansonsten Sprung zu Schritt 5.</p> <p>Schritt 5 Ende des Crawlvorgangs</p>
---	--

## Web-Page Crawler auf der Basis von Konzepten von Cho, Molina, Page

„Wenn man nicht zufällig IRGENDEINE URL als nächstes besucht, sondern versucht die relevantesten Seiten zuerst zu besuchen, wird der Index der Suchmaschine qualitativ besser werden. Ausserdem können Ressourcen gespart werden“

<p>Kriterien nach denen die nächste zu crawlende URL ausgesucht werden kann:</p> <p><b>1. Ähnlichkeit zu einem Suchwort für den Crawler (Similarity to driving query)</b> Seiten werden zuerst gecrawlt, auf die Links zeigen, die dem Suchwort ähnlich sind. Beim aktualisieren von Seiten werden solche zuerst gecrawlt, die das Suchwort (oft) enthalten.</p> <p><b>2. Backlink Anzahl</b> Die Anzahl der Links die auf eine URL zeigen generieren ihren Backlinkwert</p> <p><b>3. Pagerank</b> Einsatz des Pageranks, besonders interessant für bereits erfasste Seiten die aktualisiert werden sollen, bei nicht bekannten Seiten muss der PR geraten werden.</p> <p><b>4. Location Metric</b> Bestimmte Eigenschaften der URL des Dokuments, z.B. wenige Unterverzeichnisse, Endung .com, bestimmtes Suchwort in URL enthalten</p>	<p>Verschiedene Abbruchbedingungen für Crawler:</p> <p><b>1. Crawl &amp; Stop</b> In diesem Konzept besucht der Crawler eine festgelegte Anzahl von K Seiten mit einer beliebigen Start-URL. Ein ideal funktionierender Crawler hätte nach dem Besuch von K Seiten aus allen verfügbaren Seiten nur die K relevantesten extrahiert. Wären also N Seiten durch Links (<math>K &lt; N</math>) verfügbar gewesen, hätte der Crawler aus diesen N verfügbaren Seiten die K relevantesten Seiten extrahiert und den Vorgang beendet. Einfach ausgedrückt könnte man sagen, dass der Crawler eine gewissen Menge von Seiten besuchen darf, die Menge der verfügbaren Seiten jedoch grösser ist und er nur die relevantesten Seiten aus der Gesamtmenge besuchen sollte.</p> <p><b>2. Limited Buffer Crawl</b> Auch bei diesem Vorgang besucht der Crawler eine festgelegte Anzahl von K Seiten. Sein Speicherplatz in den er erfasste Seiten ablegen kann ist jedoch nur S gross (<math>S &lt; K</math>). Der Crawler hat also nicht genug Platz um alle gecrawlten Dokumente unterzubringen. Hier entfernt der Crawler während des Crawlingvorgangs Dokumente aus seinem Speicher um Platz für neue zu schaffen. Bei einem ideal funktionierenden Crawler wären nach dem Erfassen von S Dokumenten nur noch die relevantesten aus der Gesamtmenge K im Speicher. Da dem Crawler die Relevanz von neue erfassten Dokumenten nicht bekannt ist muss er versuchen sie auf der Basis der bereits vorhandenen Dokumente zu „erraten“.</p> <p><b>3. Crawl &amp; Stop with Threshold (Threshold = Schwelle)</b> Hier wird vorgegangen wie bei der ersten Methode. Jedoch wird ein Ziel-Qualitätskriterium für das Crawlten vergeben. Nach dem Crawlvorgang von K Dokumenten muss ein perfekter Crawler alle Dokumente mit dem höchsten Qualitätswert erfasst haben.</p>
--	---

# Eigenschaften des Googlebots

<p>Googlebot ist auf mehrere Rechner <b>verteilt</b>.</p> <p>Es gibt 2 verschiedene <b>Googlebot-Varianten</b>          64.*.* (Freshbot läuft ständig)          209.*.* (Deepbot läuft einmal im Monat)</p> <p>Seiten mit <b>höherem Pagerank</b> werden <b>häufiger besucht</b></p> <p>Googlebot <b>crawl langsamer als er könnte</b> um Server nicht voll auszulasten.</p> <p>Für lokale Google-Sites wie z.B. google.de werden <b>lokalsprachige Dokumente anders gewichtet</b>.</p> <p>Googlebot hält sich an die <b>robots.txt</b></p>	<p>Googlebot <b>berücksichtigt Backlinks nur bei Seiten mit einem PR &gt; 3</b></p> <p>Freshbot <b>besucht eine Seite sehr schnell</b>, nachdem man sich für <b>Ad-Sense</b> angemeldet hat.</p> <p>Googlebot kann folgende <b>Dateiformate</b> besuchen und speichern: html, pdf, asp, jsp, hhtml, shtml, xml, cfm, doc, xls, ppt, rtf, wks, lwp, wri.</p> <p>Das Schalten von <b>Ad-Words</b> für eine Seite <b>beeinträchtigt NICHT ihr Ranking</b> und die Besuche vom Bot (offizielles Google Statement)</p> <p>Das <b>Ändern der URL</b> einer Seite wird das Ranking der Seite und damit auch die <b>Besuchsfrequenz des Bots verändern</b>, auch wenn die Inhalte gleich bleiben.</p>
--	---

## Literatur:

Chu Molina, Page - Efficient crawling through URL ordering  
<http://www7.scu.edu.au/programme/fulpapers/1919/com1919.htm>

Definitions for crawler  
 Google Suche: define:crawler

Google's FAQ for the Googlebot  
<http://www.google.com/bot.html>

The GNU wget crawler  
<http://www.gnu.org/software/wget/wget.html>

The robots exclusion standard  
<http://www.robotstxt.org/>

Google Watch - A look at how Google's monopoly, algorithms, and privacy policies are undermining the Web  
<http://www.google-watch.com>

## Download:

Download des Thesenpapiers und der Präsentation unter  
<http://www.daniel-ritter.de/arbeiten/>