

## Kurzfassung der 2. Studienarbeit

<b>Abteilung</b>	<b>Informatik</b>
<b>Name der Studentin / des Studenten</b>	<b>Sebastian de Castelberg Markus Kinzler</b>
<b>Studienjahr</b>	<b>SS04</b>
<b>Titel der 2. Studienarbeit</b>	<b>Projekt Geometabot – Ein fachspezifischer WebCrawler</b>
<b>Examinatorin / Examinator</b>	<b>Stefan F. Keller</b>
<b>Kurzfassung der 2. Studienarbeit</b>	
<p>Im Projekt GeometaBot0.1 geht es darum einen fokussierten WebCrawler auf Basis von Java zu Implementieren. Der Crawler sammelt automatisch URL's von georelevanten Seiten und bereitet diese auf, so dass sie in den Index vom Suchportal <a href="http://www.geometa.info">www.geometa.info</a> aufgenommen werden können.</p> <p>Als Basis des Crawlers dient der WebCrawler Heritrix. Heritrix bietet sehr gute Schnittstellen um dem Crawler zu konfigurieren und zu jedem Zeitpunkt des Crawls seine eigenen Module einzubinden.</p> <p>Um GeometaBot auf den Geobereich zu fokussieren wurden zwei Komponenten eingefügt:</p> <ul style="list-style-type: none"> <li>• LinkDistiller: Priorisiert das Besuchen von Links die georelevante Wörter beinhalten</li> <li>• Classifier: Der Inhalt der heruntergeladenen Seiten wird analysiert und mittels Mathematischer Methoden auf die Georelevanz geprüft.</li> </ul> <p>Der Crawler ist zurzeit noch etwas langsam, was aber in zukünftigen Heritrix Releases verbessert wird.</p> <p>Der Classifier, der ganze Texte klassifizieren soll, ist nicht so gut wie erwartet. Dies funktionierte mit den WEKA Implementationen von J48 und NaiveBayes nicht wie gewünscht. Zu viele (auch grobe) Fehler machen den Classifier in diesem Entwicklungsstadium nicht brauchbar .</p>	



HSR  
HOCHSCHULE FÜR TECHNIK  
RAPPERSWIL

