

# Abschlussbericht

## Studienprojekt: TaxoSearch

Lehrstuhl für Computerlinguistik  
Ruprecht-Karls-Universität Heidelberg

vorge stellt von

Thorsten Beinhorn, Vesna Cvoro, Khaled Dhaoui und  
Christian Pretzsch

# Was ist TaxoSearch?

- TaxoSearch ist ein Tool zur Unterstützung bei der Internetsuche. Unser Programm hilft dem Benutzer Suchbegriffe präziser zu wählen und unterstützt ihn bei der Sichtung der Suchergebnisse durch individuelle Darstellung des Rankings

# Ziele vom Projekt

- Steigerung der Trefferquote durch automatische semantische Erweiterung der Suchanfrage durch Wordnet
- thematische Eingrenzung der Suchergebnisse mittels Ontologie
- Schnittstellen zu mehreren Suchmaschinen
- Information Retrieval zur Aufbereitung der Ergebnisse

# Realisierung I

- Erweiterung der Suchanfrage
  - Wordnet als einzige Alternative
- Manuell vs. Automatisch
  - Automatische Disambiguierung der Suchbegriffe nicht möglich

# Realisierung II

- Schnittstellen zu verschiedenen Suchmaschinen
  - Praktische Gründe und strukturelle Ähnlichkeiten der Suchmaschinen => Google
  - Nachträgliche Erweiterung problemlos möglich

# Realisierung III

- Information Retrieval
  - Ranking und individuelle Darstellung durch Selektion relevanter Begriffe
  - Vektormodell

# Arbeiten mit TaxoSearch I

- Eingabe der Suchbegriffe
  - Unterstützung des Benutzers durch Thesaurus Funktion (WordNet)
  - Erweiterte GoogleSuche in GUI vereinfacht dargestellt

# Arbeiten mit TaxoSearch II

- Information Retrieval:  
Ranking der Trefferseiten
  - Anpassung des Rankings durch individuelle  
Zusammenstellung relevanter Begriffe



# Umsetzung

- Programmiersprache: Python
- Plattform: Unix, Mac OS, Windows
- GUI: TkInter
- Information Retrieval: Vektormodell

# Vektormodell

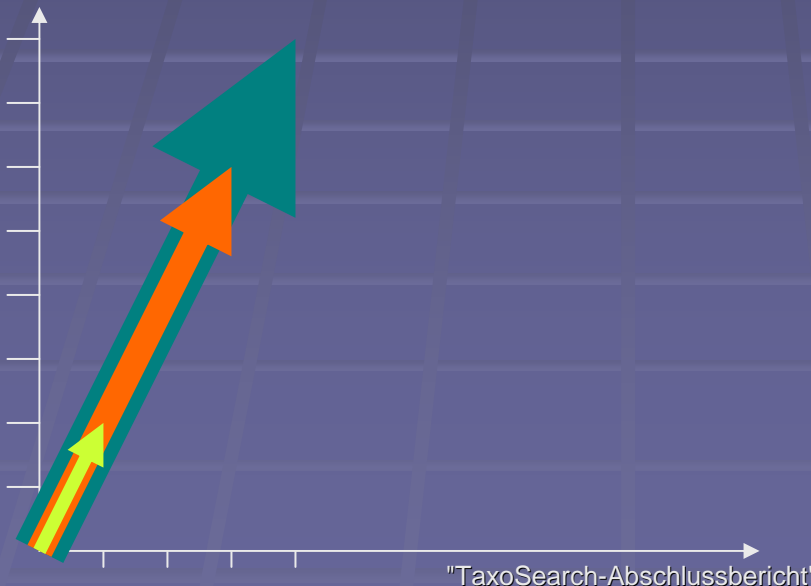
- Berücksichtigung partieller Übereinstimmung zwischen Anfragetermen und Dokumenttermen durch nicht - binäre Werte für Termgewichtung
- Berechnung der Ähnlichkeit zwischen Anfragetermen und Dokumenttermen
- Sortierung von Dokumenten nach Grad der Ähnlichkeit

(Baeza-Yates/Ribeiro-Neto, 1999,27)

# Vektormodell - Beispiel

Dokumentvektor: (4,8,0)  
Vektor Query<sub>1</sub>: (1,2,0)  
Vektor Query<sub>2</sub>: (3,6,0)

	Term 1 <i>Öl</i>	Term 2 <i>Preis</i>	Term 3 <i>Alaska</i>
	4	8	0
	1	2	0
	3	6	0



(vgl. Kowalski, 1997,153)

# Externe Komponenten

- WordNet 1.7.1
- pywordnet 1.4
- pyGoogle 0.5.3
- MontyTagger 1.2

# WordNet

- WordNet ist eine semantische (allgemeinorientierte) Wortdatenbank, die an der University of Princeton erstellt wurde
  - <http://www.cogsci.princeton.edu/~wn/>
- WordNet ist eine hierarchische Datenbank, die erfassten Konzepte werden in Baumform konzeptualisiert.
- WordNet ist ein erweitertes Lexikon, dessen Design an aktuelle psycholinguistische Theorien des menschlichen Wortgedächtnisses angelehnt ist.
- Englische Nomen, Verben und Adjektive werden in sogenannten *Synonym Sets* organisiert, von denen jeder ein zugrundeliegendes lexikalisches Konzept repräsentiert. Verschiedene Beziehungen verbinden diese Synonym Sets.

# WordNet II

- WordNet speichert seine Informationen nicht nach alphabetischen, sondern nach konzeptuellen Gesichtspunkten.
  - Begriffe gleicher oder ähnlicher Bedeutung werden zusammen abgelegt.
  - die Datenbank von WordNet umfaßt im Gegensatz zu früheren Projekten ungefähr 95600 verschiedene Wortformen, welche in ca. 70100 Bedeutungsklassen (Synonym Sets) eingeordnet sind.
  - WordNet unterteilt das Lexikon in 5 Kategorien : Nomen, Verben, Adjektive, Adverbien und Funktionswörter (nicht implementiert). Diese Einteilung beruht auf Untersuchungen über Wortassoziationen von Fillenbaum und Jones (1965)
    - Der Nachteil dieser Vorgehensweise besteht darin, daß manche Wörter in mehr als eine syntaktische Kategorie hineinfallen können und somit mehrmals gespeichert werden müssen.

# Monty Tagger I

- Regelbasierter POS – Tagger
- Basiert auf Brills 1994 entwickelten Transformational – Based Learning POS Tagger
- In plattformunabhängigem Python und Java
- Benutzt als Grundlage die Penn Treebank



<http://web.media.mit.edu/~hugo/research/montytagger.html>

# Monty Tagger II

- Tokenizer:
  - „Tokenization“ des Eingabetextes
    - children's --> children 's
    - parents' --> parents '
    - won't --> wo n't
    - I'm --> I 'm
  - Trennung von Wörtern und Interpunktion durch Leerzeichen
  - Ausnahme: Abkürzungen und Akronyme



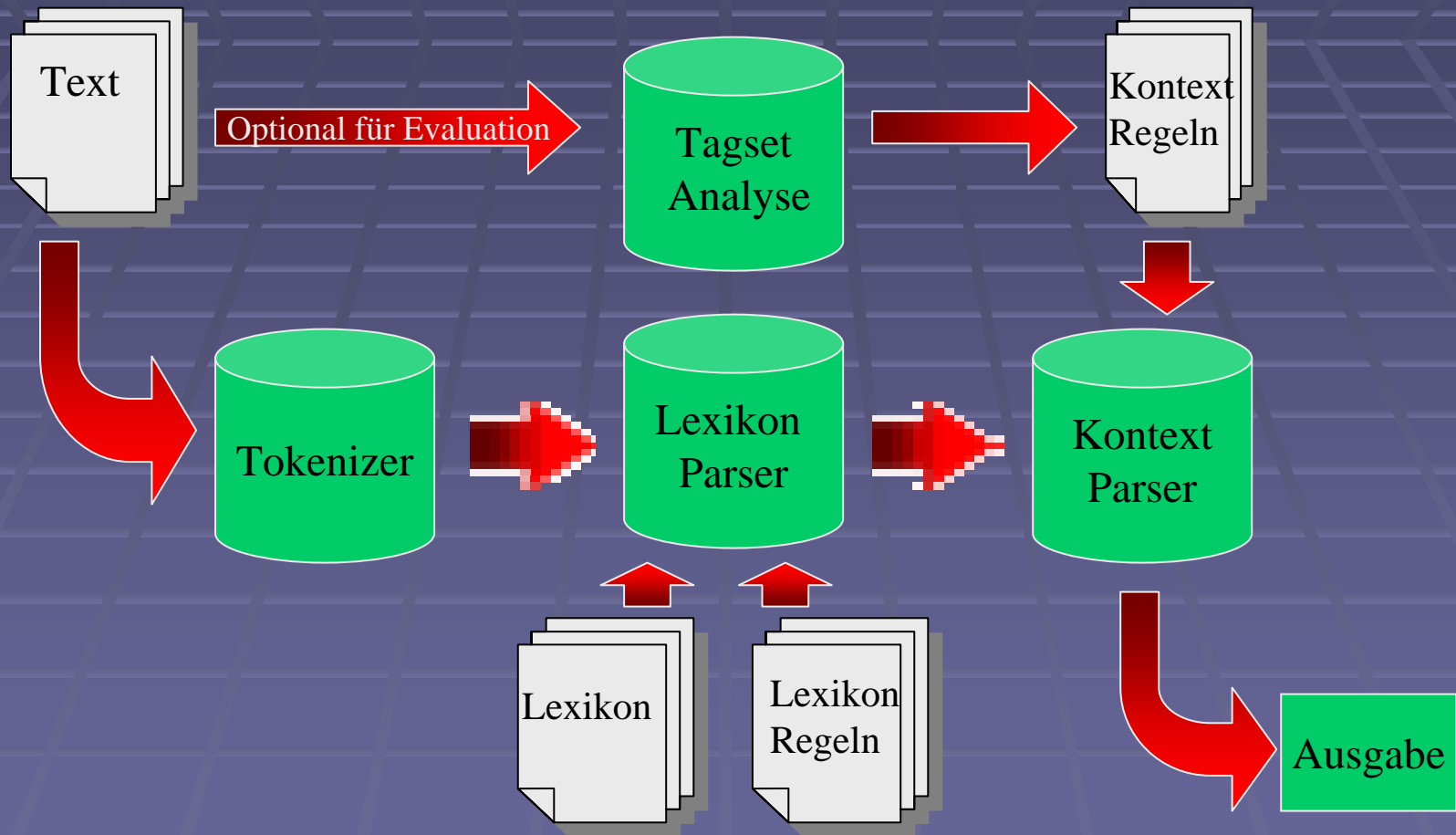
# Monty Tagger III

- Lexikon und lexikalische Regeln:
  - Einbindung eines Lexikons und eines Regelsets („Brill94 lexical rule files“)
  - Morphosyntaktische Analyse
  - Zuordnung des „wahrscheinlichsten“ Tags
    - „golden gate“ -> /NNP
    - „race“ -> /NNS oder /VB ?

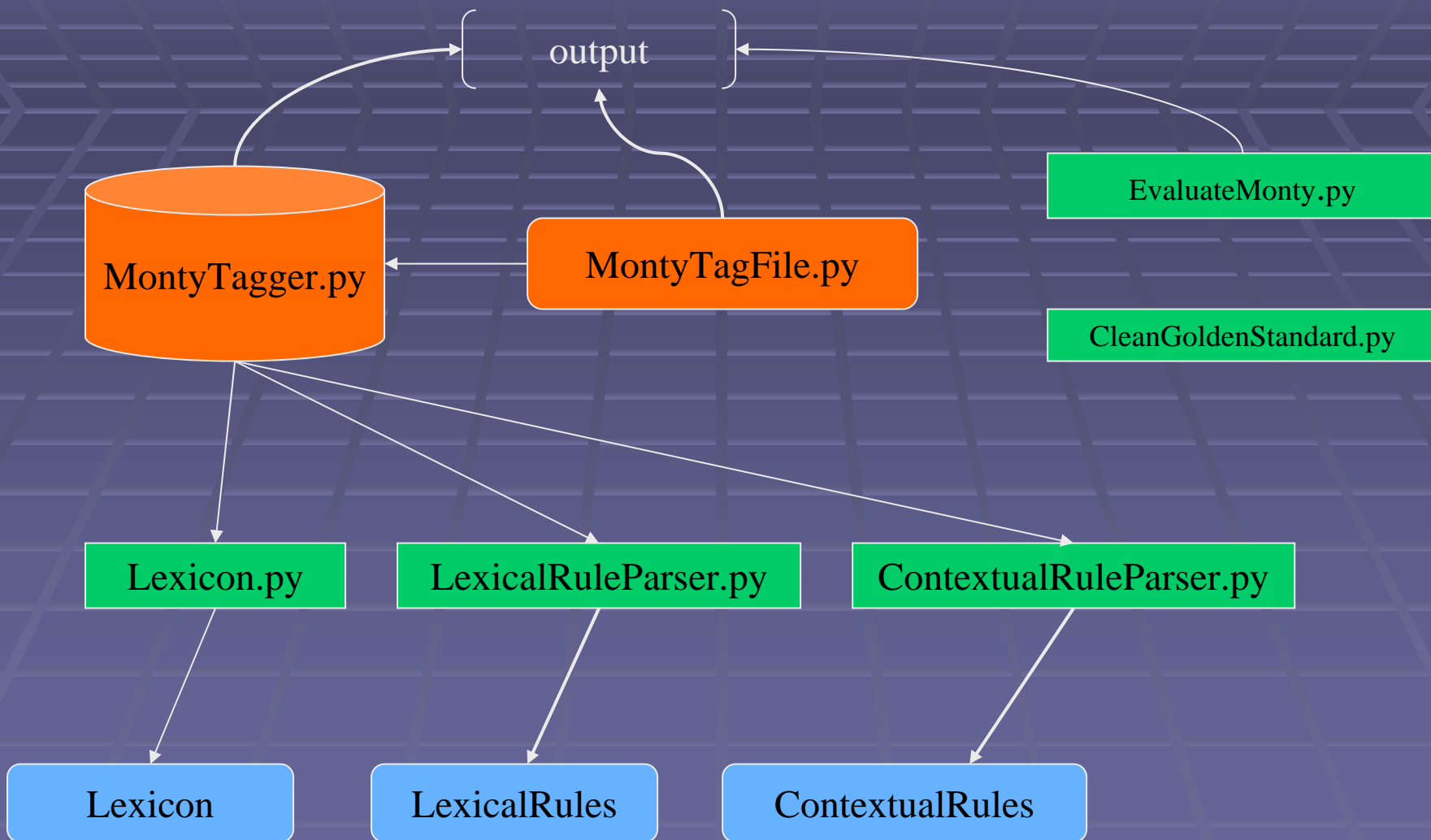
# Monty Tagger IV

- Kontextregeln und Syntaxanalyse:
  - Einbindung der Kontextregeln („Brill94 context rule files“)
  - Syntaktische (Kontext-) Analyse: jede Regel wird für alle Wörter geprüft
  - Überprüfung und anschließende Zuordnung bzw. Auswahl der „wahrscheinlichsten“ Tags
    - „golden gate“ -> /NNP
    - „race“ -> /NNS oder /VB ? -> Entscheidung: race /NNS

# Beispiel (Monty Tagger)



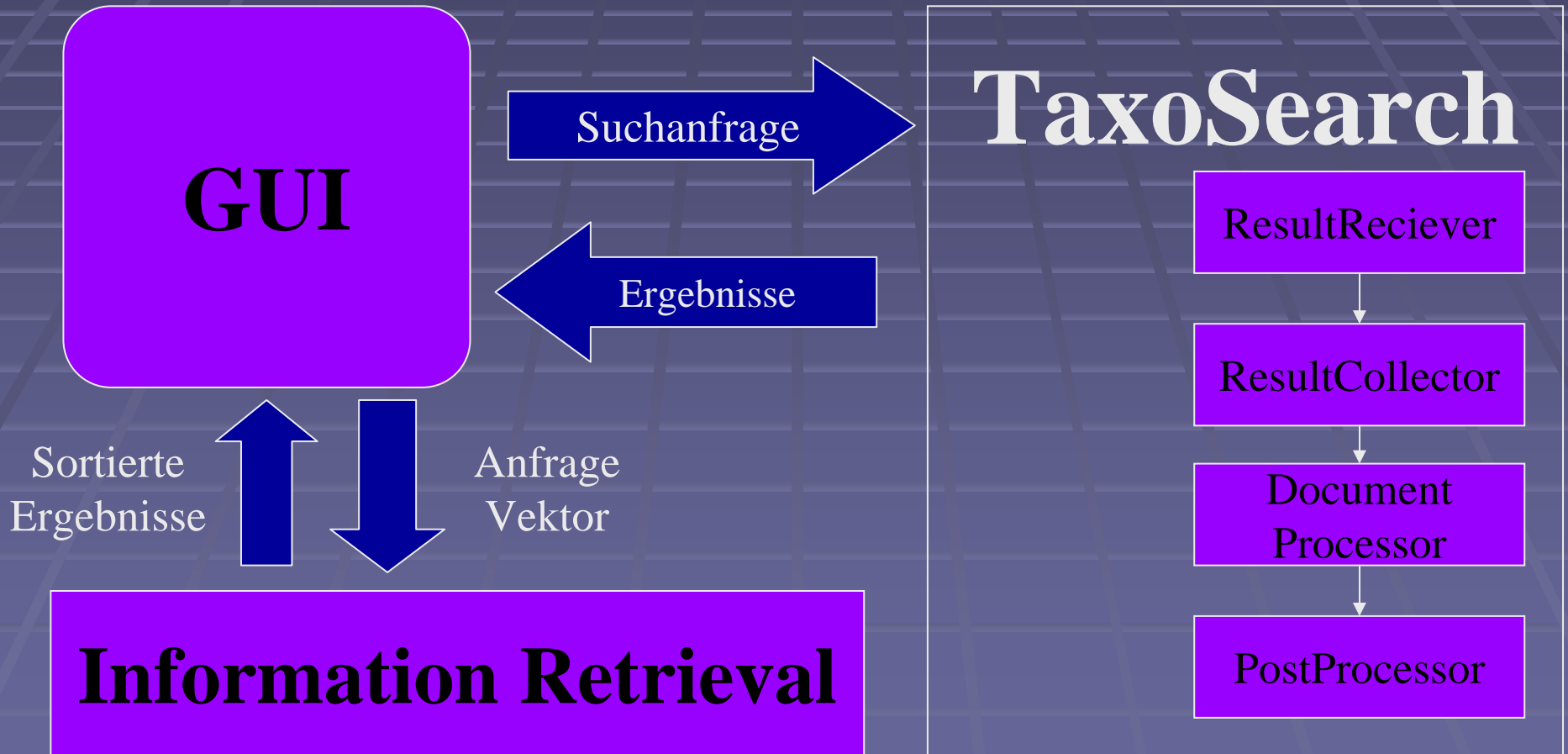
# Architektur



# Evaluation: Monty Tagger

- Performance
  - 500 Wörter /s
- Genauigkeit
  - ca. 96 – 97 %

# TaxoSearch - Systemarchitektur



# Modul: mod\_DocumentObjects

- **Class DocumentCollection**
  - Speichert eine Sammlung an Dokumenten in einem Dictionary
- **Eigenschaften**
  - *dicDocumentCollection*
    - Dictionary zum Speichern der DocumentObjects (Key=URL, Value=DocumentObject)
- **Methoden**
  - *AddDocument(self, PageURL, Title)*
    - Fügt ein neues DocumentObject in die Auflistung ein.
  - *UpdateDocument(self, DocumentObject)*
    - Aktualisiert ein DocumentObject in der Auflistung mit einem extern veränderten DocumentObject

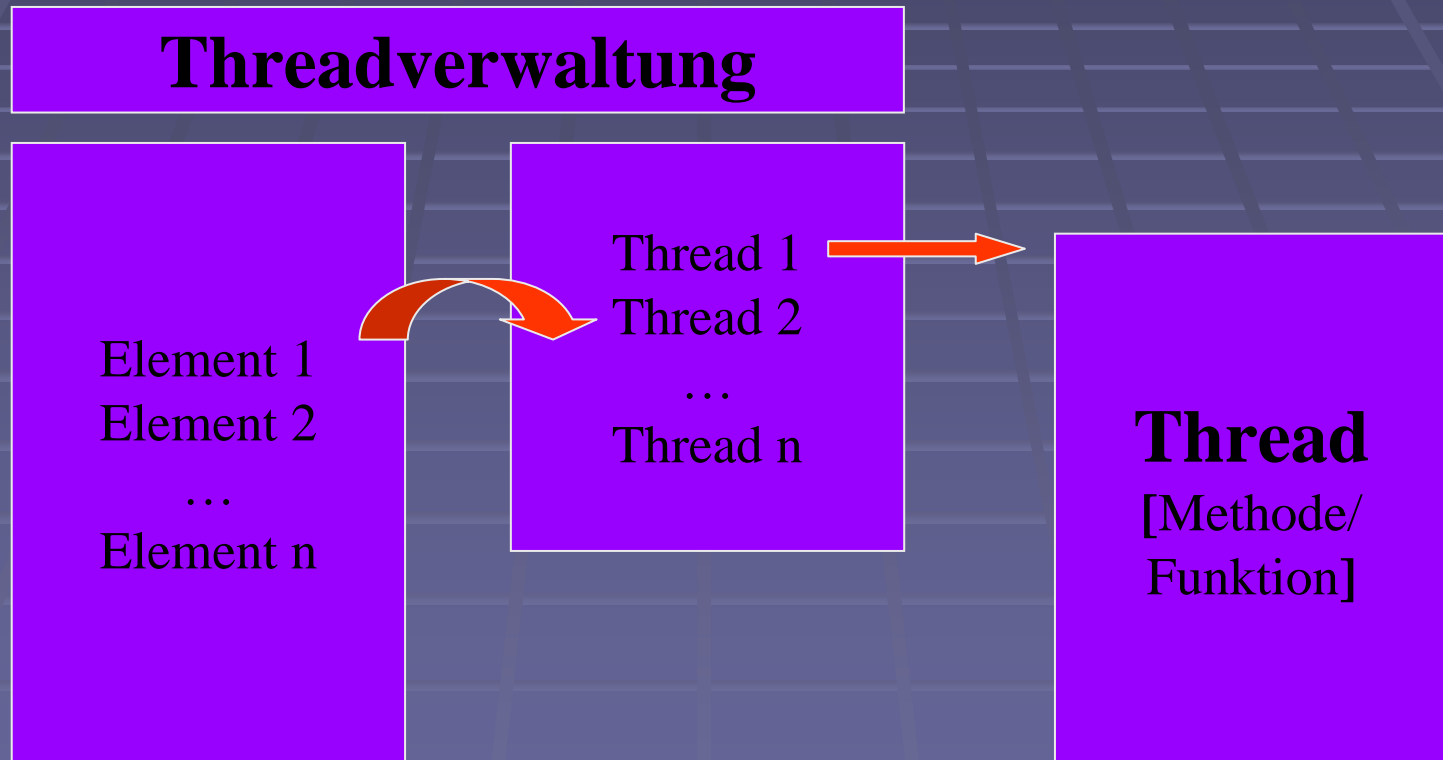
# Modul: mod\_DocumentObjects

## Class DocumentObject

- Repräsentiert ein eine Trefferseite mit all ihren wichtigen Eigenschaften
- **Eigenschaften**
  - *URL*
    - Internet-Adresse der Trefferseite
- *Title*
  - Titel der Trefferseite (aus HTML-Header)
- *Page*
  - Kompletter, unveränderter Inhalt einer Trefferseite
- *PageContent*
  - Inhalt einer Trefferseite nach DeHTML, also ohne Tags
- *dicDocumentVector*
  - Dokumentvector einer Trefferseite, enthält alle Nomen und Verben  
Key=Wort, Value=Anzahl
- *numWordCount*
  - Anzahl aller Wörter im Dokumentvector
- *numDocumentRating*
  - PageRankingWert nach Information Retrieval



# Multithreading



# TaxoSearch in Aktion

**TaxoSearch**

Please enter your query words here:

Autogenerate synonyms

AND  
  OR  
  NOT  
  EXACT / FORCE

Nouns	Verbs	Adjectives
<p><b>submarine</b>  <i>a submersible warship usually armed with torpedoes</i></p> <p><b>Synonyms</b>  <a href="#">submarine</a> <a href="#">pigboat</a> <a href="#">sub</a> <a href="#">U-boat</a></p> <p><b>Hypernyms</b>  <a href="#">submersible</a> <a href="#">submersible warship</a></p> <p><b>Hyponyms</b>  <a href="#">attack submarine</a> <a href="#">auxiliary research submarine</a>  <a href="#">fleet ballistic missile submarine</a>  <a href="#">nautilus</a> <a href="#">nuclear submarine</a> <a href="#">nuclear-powered subr</a></p> <p><b>submarine</b>  <i>a large sandwich made of a long crusty roll split lengthwise and filled with meats and cheese (and tomato and onion and lettuce and condiments)</i></p> <p><b>Synonyms</b>  <a href="#">bomber grinder</a> <a href="#">hero</a> <a href="#">hero sandwich</a>  <a href="#">hoagie</a> <a href="#">hoagy</a> <a href="#">Cuban sandwich</a> <a href="#">italian sandwich</a>  <a href="#">poor boy</a> <a href="#">sub</a> <a href="#">submarine</a> <a href="#">submarine sandwich</a>  <a href="#">torpedo</a> <a href="#">wedge zep</a></p> <p><b>Hypernyms</b>  <a href="#">sandwich</a></p> <p><b>Hyponyms</b>            -none-</p>	<p><b>submarine</b>  <i>move forward or under in a sliding motion</i></p> <p><b>Synonyms</b>  <a href="#">submarine</a></p> <p><b>Hypernyms</b>  <a href="#">skid</a> <a href="#">slip</a> <a href="#">slue</a> <a href="#">slew</a> <a href="#">slide</a></p> <p><b>Hyponyms</b>            -none-</p> <p><b>submarine</b>  <i>throw with an underhand motion, as of a baseball</i></p> <p><b>Synonyms</b>  <a href="#">submarine</a></p> <p><b>Hypernyms</b>  <a href="#">flip</a> <a href="#">toss</a> <a href="#">skv</a> <a href="#">pitch</a></p> <p><b>Hyponyms</b>            -none-</p> <p><b>submarine</b>  <i>bring down with a blow to the legs, in sports</i></p> <p><b>Synonyms</b>  <a href="#">submarine</a></p> <p><b>Hypernyms</b>  <a href="#">down</a> <a href="#">knock down</a> <a href="#">cut down</a> <a href="#">push down</a>  <a href="#">pull down</a></p> <p><b>Hyponyms</b>            -none-</p> <p><b>submarine</b>  <i>control a submarine</i></p> <p><b>Synonyms</b>  <a href="#">submarine</a></p> <p><b>Hypernyms</b>  <a href="#">operate</a> <a href="#">control</a></p> <p>..</p>	<p><b>submarine</b>  <i>beneath the surface of the sea</i></p> <p><b>Synonyms</b>  <a href="#">submarine</a> <a href="#">undersea</a></p> <p><b>Hypernyms</b>            -none-</p> <p><b>Hyponyms</b>            -none-</p>



Very important ▲

Important ▲

Less important ▲

Back to search

**CNN.com - Russian submarine rescue bid under way - August 14, ...**

<http://www.cnn.com/2000/WORLD/europe/08/14/russia.submarine.06/>

 70%

**CNN.com - Bad weather hampers Russian submarine rescue - August ...**

<http://www.cnn.com/2000/WORLD/europe/08/14/russia.submarine.07/>

 70%

**A Russian Submarine (17-Dec-2000)**

[http://www.sandelman.ottawa.on.ca/People/Michael\\_Richardson/russian-submarine.html](http://www.sandelman.ottawa.on.ca/People/Michael_Richardson/russian-submarine.html)

 63%

**Guardian Unlimited | Special reports | Minister blames ...**

<http://www.guardian.co.uk/submarine/story/0,7369,1033266,00.html>

 62%

**Russian Submarine (captioned photo)**

<http://familyofmann.tripod.com/rph021.htm>

 61%

**scorpion**

<http://www.beachcalifornia.com/scorpion.html>

 60%

**Russian Submarine Disaster - Latest - Jane&#39;s Naval Forces**

[http://www.janes.com/defence/naval\\_forces/news/jfs/jfs000815\\_1\\_n.shtml](http://www.janes.com/defence/naval_forces/news/jfs/jfs000815_1_n.shtml)

 59%

**Russian Submarine Emergency - Jane&#39;s Naval Forces**

[http://www.janes.com/defence/naval\\_forces/news/jfs/jfs000814\\_1\\_n.shtml](http://www.janes.com/defence/naval_forces/news/jfs/jfs000814_1_n.shtml)

 59%

**Pravda.RU Search for sailors from Russian submarine sunk in ...**

<http://newsfromrussia.com/accidents/2003/08/30/49676.html>

 57%

**A Visit to a Russian Submarine**

<http://cosmo.pasadena.ca.us/adventures/submarine/>

 57%

**Pravda.RU Swedish divers discover hull of Russian submarine which ...**

<http://newsfromrussia.com/science/2003/08/01/49012.html>

 57%

**Russian Submarine**

<http://www.queenmary.com/QMweb/html/sub.html>

 50%

Very important ▲  
kursk ▼

Important ▲  
▼

Less important ▲  
▼

Back to search

<b>Guardian Unlimited   Special reports   Minister blames ...</b> <a href="http://www.guardian.co.uk/submarine/story/0,7369,1033266,00.html">http://www.guardian.co.uk/submarine/story/0,7369,1033266,00.html</a>		96%
<b>Russian Submarine Disaster - Latest - Jane's Naval Forces</b> <a href="http://www.janes.com/defence/naval_forces/news/jfs/jfs000815_1_n.shtml">http://www.janes.com/defence/naval_forces/news/jfs/jfs000815_1_n.shtml</a>		90%
<b>Russian Submarine Emergency - Jane's Naval Forces</b> <a href="http://www.janes.com/defence/naval_forces/news/jfs/jfs000814_1_n.shtml">http://www.janes.com/defence/naval_forces/news/jfs/jfs000814_1_n.shtml</a>		83%
<b>CNN.com - Bad weather hampers Russian submarine rescue - August ...</b> <a href="http://www.cnn.com/2000/WORLD/europe/08/14/russia.submarine.07/">http://www.cnn.com/2000/WORLD/europe/08/14/russia.submarine.07/</a>		58%
<b>CNN.com - Russian submarine rescue bid under way - August 14, ...</b> <a href="http://www.cnn.com/2000/WORLD/europe/08/14/russia.submarine.06/">http://www.cnn.com/2000/WORLD/europe/08/14/russia.submarine.06/</a>		52%
<b>A Russian Submarine (17-Dec-2000)</b> <a href="http://www.sandelman.ottawa.on.ca/People/Michael_Richardson/russian-submarine.html">http://www.sandelman.ottawa.on.ca/People/Michael_Richardson/russian-submarine.html</a>		13%
<b>Russian Submarine (captioned photo)</b> <a href="http://familyofmann.tripod.com/rph021.htm">http://familyofmann.tripod.com/rph021.htm</a>		13%
<b>scorpion</b> <a href="http://www.beachcalifornia.com/scorpion.html">http://www.beachcalifornia.com/scorpion.html</a>		13%
<b>Pravda.RU Search for sailors from Russian submarine sunk in ...</b> <a href="http://newsfromrussia.com/accidents/2003/08/30/49676.html">http://newsfromrussia.com/accidents/2003/08/30/49676.html</a>		12%
<b>A Visit to a Russian Submarine</b> <a href="http://cosmo.pasadena.ca.us/adventures/submarine/">http://cosmo.pasadena.ca.us/adventures/submarine/</a>		12%
<b>Pravda.RU Swedish divers discover hull of Russian submarine which ...</b> <a href="http://newsfromrussia.com/science/2003/08/01/49012.html">http://newsfromrussia.com/science/2003/08/01/49012.html</a>		12%
<b>Russian Submarine</b> <a href="http://www.queenmary.com/QMweb/html/sub.html">http://www.queenmary.com/QMweb/html/sub.html</a>		11%

**TaxoSearch**

Very important: [kursk](#) | Important: [explosion](#) | Less important: [rescuer](#) | [Back to search](#)

Very important  
 Important  
 Less important

<b>Russian Submarine Disaster - Latest - Jane's Naval Forces</b> <a href="http://www.janes.com/defence/naval_forces/news/fts/fts000815_1_n.shtml">http://www.janes.com/defence/naval_forces/news/fts/fts000815_1_n.shtml</a>	83%	<a href="#">rescue</a>
<b>Guardian Unlimited   Special reports   Minister blames ...</b> <a href="http://www.guardian.co.uk/submarine/story/0,7369,1033266,00.html">http://www.guardian.co.uk/submarine/story/0,7369,1033266,00.html</a>	72%	<a href="#">rescuer</a>
<b>Russian Submarine Emergency - Jane's Naval Forces</b> <a href="http://www.janes.com/defence/naval_forces/news/fts/fts000814_1_n.shtml">http://www.janes.com/defence/naval_forces/news/fts/fts000814_1_n.shtml</a>	71%	<a href="#">researcher</a>
<b>CNN.com - Bad weather hampers Russian submarine rescue - August ...</b> <a href="http://www.cnn.com/2000/WORLD/europe/08/14/russia.submarine.07/">http://www.cnn.com/2000/WORLD/europe/08/14/russia.submarine.07/</a>	58%	<a href="#">reserve</a>
<b>Pravda.RU Search for sailors from Russian submarine sunk in ...</b> <a href="http://newsfromrussia.com/accidents/2003/08/30/49676.html">http://newsfromrussia.com/accidents/2003/08/30/49676.html</a>	45%	<a href="#">resilience</a>
<b>CNN.com - Russian submarine rescue bid under way - August 14, ...</b> <a href="http://www.cnn.com/2000/WORLD/europe/08/14/russia.submarine.06/">http://www.cnn.com/2000/WORLD/europe/08/14/russia.submarine.06/</a>	43%	<a href="#">resort</a>
<b>Pravda.RU Swedish divers discover hull of Russian submarine which ...</b> <a href="http://newsfromrussia.com/science/2003/08/01/49012.html">http://newsfromrussia.com/science/2003/08/01/49012.html</a>	24%	<a href="#">resource</a>
<b>A Russian Submarine (17-Dec-2000)</b> <a href="http://www.sandelman.ottawa.on.ca/People/Michael_Richardson/russian-submarine.html">http://www.sandelman.ottawa.on.ca/People/Michael_Richardson/russian-submarine.html</a>	10%	<a href="#">respect</a>
<b>Russian Submarine (captioned photo)</b> <a href="http://familyofmann.tripod.com/rsh021.htm">http://familyofmann.tripod.com/rsh021.htm</a>	10%	<a href="#">restrict</a>
<b>scorpion</b> <a href="http://www.beachcalifornia.com/scorpion.html">http://www.beachcalifornia.com/scorpion.html</a>	9%	<a href="#">result</a>
<b>A Visit to a Russian Submarine</b> <a href="http://cosmo.pasadena.ca.us/adventures/submarine/">http://cosmo.pasadena.ca.us/adventures/submarine/</a>	9%	<a href="#">resume</a>
<b>Russian Submarine</b> <a href="http://www.queenmary.com/QMweb/html/sub.html">http://www.queenmary.com/QMweb/html/sub.html</a>	8%	
<b>The Russian Submarine, Located in Fokestone Harbour, UK.</b> <a href="http://www.sovietsub.co.uk/">http://www.sovietsub.co.uk/</a>	7%	
<b>you'll enjoy this! russian submarine</b> <a href="http://book-data.saveurbois.com/russian-submarine.html">http://book-data.saveurbois.com/russian-submarine.html</a>	7%	
<b>Armchair Travel for Virtual Travel (R) - panorama, travel, ...</b>	0%	

red

# Literatur

- Karin Haenelt
- **Alexander Valet / Christian Pretzsch / Vanessa Micelli: Ansätze des Tagging – Ein Seminarreferat**, „Hauptseminar: Parsing“, Universität Heidelberg, SS 2003  
[kontext.fraunhofer.de/haenelt/kurs/Referate/Micelli\\_Pretzsch\\_Valet\\_SS03/MontyAnsatzRef.pdf](http://kontext.fraunhofer.de/haenelt/kurs/Referate/Micelli_Pretzsch_Valet_SS03/MontyAnsatzRef.pdf)
- MontyTagger Version 1.2:  
<http://web.media.mit.edu/~hugo/research/montytagger.html>
- WordNet 1.7.1
- pyWordNet
- pyGoogle 0.5.3