#### **Developing Intelligent Search Engines**

#### Isabel 'MaineC.' Drost



#### The world wide web





• Basic architecture of search engines.

Properties of the WWW

• One machine learning task in search engines.

• How can YOU contribute?

#### Search engines from scratch



#### Search engines from scratch



### Crawling the WWW



# Crawling the WWW

- Choose starting page (e.g. DMOZ).
- Get and save page content.
- Extract and follow links recursively.
  - Parallelize crawling.
  - Prioritize links to follow.
  - Avoid gathering too much content at once DOS.

#### Search engines from scratch



# Indexing web pages

- Parse web pages.
  - Extract keywords.
  - Extract keywords from anchor texts of inlinks.
- Index content of pages.
  - Each keyword part of a very long list.
  - URLs of matching pages associated with list entries.

# Indexing web pages

AbiBycambew 27 take (5 uhel 22 CBrythet fikst place in Berlin.

- 2005 http://www.ccc.de/congress/2005
- 22C3 http://www.fukami.de http://www.ccc.de/congress/2005
  - 27<sup>th</sup> http://www.ccc.de/congress/2005
    - at http://www.ccc.de/congress/2005
- Berlin http://www.berlin.de http://www.ccc.de/congress/2005

december http://www.ccc.de/congress/2005

- in http://www.ccc.de/congress/2005
- place http://www.ccc.de/congress/2005
  - take http://www.ccc.de/congress/2005
  - the http://www.ccc.de/congress/2005
  - will http://www.ccc.de/congress/2005

#### Search engines from scratch



## Handling user queries

- Retrieve queries from user.
  - Parallelize query retrieval.
- Retrieve matching pages from index.
- Rank found pages.
- Create easy to understand representation.

#### Handling user queries



#### Handling user queries

#### 22C3: Private Investigations

The 22nd Chaos Communication Congress (**22C3**) is a four-day conference on ... More information on **22C3** is coming up. For now, study and distribute our Call ... www.ccc.de/congress/2005/ - 3k - <u>Cached</u> - <u>Similar pages</u>

#### CCC | Chaos Communication Congress - [ Translate this page ]

**22C3**: Private Investigations. Der 22. Chaos Communication Congress findet vom 27. bis 30. ... All Informationen zum **22C3** finden sich unter [Interner Link] ... www.ccc.de/congress/ - 13k - <u>Cached</u> - <u>Similar pages</u> [<u>More results from www.ccc.de</u>]

#### Upcoming.org: 22C3: Private Investigations at Berliner Congress ...

The 21nd Chaos Communication Congress (**22C3**) is a conference on technology, society and utopia ... Chaos Communication Congress **22C3**: Private Investigations ... upcoming.org/event/23579/ - 9k - <u>Cached</u> - <u>Similar pages</u>

#### 22c3 h07 - NerdPedia

Darum: META-REFRESH http://hackerhippie.org/newiki/**22c3**. [old content deleted...use diff feature of this wiki] ... www.nerdpedia.org/index.php/**22c3** h07 - 8k - <u>Cached</u> - <u>Similar pages</u>

#### 22C3 (Private Investigations) - Call for Papers | Uwe Hermann

The first signs of the upcoming **22C3** congress ("Private Investigations") are ... The 22nd Chaos Communication Congress (**22C3**) is a four-day conference on ... www.hermann-uwe.de/blog/ **22c3**-private-investigations-call-for-papers - 44k - <u>Cached</u> - <u>Similar pages</u>

#### The Lunatic Fringe » Blog Archive » 22C3 Updates

**22C3** Updates. Sunday October 23rd 2005, 18:41 Filed under: General. The last months have been quite exhausting in terms of getting people on track for the ... tim.geekheim.de/2005/10/23/**22c3**-updates/ - 17k - <u>Cached</u> - <u>Similar pages</u>

#### 22C3 (Plan 9 wiki)

**22C3**. Chaos Communication Congress Where: Berlin, Germany When: 27-30/12/2005 URL: http://www.ccc.de/congress/. Who will be there? 20h; MTG; garbeam; uriel ... www.cs.bell-labs.com/wiki/plan9/22C3/ - 2k - <u>Cached</u> - <u>Similar pages</u>

#### maha's blog » 22c3

**22c3**. Sunday, July 10th, 2005. Gestern fand die erste Sitzung der ... You are currently browsing the archives for the **22c3** category. ... www.maha-online.de/blog/category/ccc/**22c3**/ - 14k - <u>Cached</u> - <u>Similar pages</u>

symlink.ch | **22C3**: Fnord-Jahresrückblick oder nicht? - [<u>Translate this page</u>] **'22C3**: Fnord-Jahresrückblick oder nicht?' | Einloggen/Account erstellen | 9 Kommentar(e) ... Re: Bin gespannt auf **22c3** by Anonymer Feigling Wednesda

# Handling user queries in depth

• What makes ranking difficult?

- Enhanced query handling.
  - Presentation of search results.
  - Enriching queries by further information.
  - Building search engines that answer questions.

### **Ranking search results**

- Unknown what the user actually is looking for.
- Rank pages according to:
  - Text based relevance for query.
  - Overall importance.
- Measure for page importance: Page Rank.
- Problems when ranking search results:
  - Users tweaking pages to spam the ranking.
  - Evaluation of rankings.

### Influencing the ranking of pages



# Ranking by page rank

- Page rank expresses popularity.
- Based on in-links.



KOraanizer

#### HowTo inflate page rank?



# A Google bomb

	] http://www.google.de/search?hs=Bot&hl=de&client=opera&rls=de&q=%22miserable+failure%22&btnG=Suche&meta=	
Google <sup>Web</sup> "mise Suche	) <u>Bilder Groups Verzeichnis News Froogle<sup>Neu!</sup> Mehr »</u> erable failure" Suche <u>Enweiterte Suche</u> <u>Einstellungen</u> e: • Das Web • Seiten auf Deutsch • Seiten aus Deutschland	
Web	Erg	ebnis
Tipp: Suchen nur nach Ergebniss Biography of President Geo Biography of the president from the www.whitehouse.gov/president/gw Past Presidents - Kids On Weitere Ergebnisse von we BBC NEWS   Americas   'M Web users manipulate a popular is to the president's page. news.bbc.co.uk/2/hi/americas/329 Welcome to MichaelMoore.co Official site of the gadfly of corpor and the television show The Awful www.michaelmoore.com/ - 39k - 1 :: ringfahndung :: Google ver :: ringfahndung :: ist ein netzine fu der kunst und des neuen digitaler www.ringfahndung.de/index.shtml	<ul> <li>an auf Deutsch. Sie können Ihre bevorzugten Spracheinstellungen in Einstellungen angeben.</li> <li>Drge W. Bush - [Diese Seite übersetzen] he official White House web site.</li> <li>wbbio.html - 33k - 15. Sept. 2005 - Im Cache - Ähnliche Seiten hly - Current News - President www.whitehouse.gov.»</li> <li>Iiserable failure' links to Bush - [Diese Seite übersetzen] search engine so an unflattering description leads</li> <li>298443.stm - 31k - 14. Sept. 2005 - Im Cache - Ähnliche Seiten rations, creator of the film Roger and Me ul Truth. Includes mailing list, message board, 14. Sept. 2005 - Im Cache - Ähnliche Seiten</li> <li>rbindet Bush mit 'Miserable failure fuer freunde des schraegen humors, der satire, n stils :). kurz :: wir schreiben hier all nl - 11k - 14. Sept. 2005 - Im Cache - Ähnliche Seiten</li> </ul>	

#### How was the bomb created?



# What is link spam?



#### Alternative Medicine Links - Add Your URL

Alternative **Medicine** Links section of the health directory at PillWatch.com. ... Order All Heralife Weight Loss products here. www.katsherbs.com/Herbalife ... www.pillwatch.com/directory/alternative-**medicine**/ - 65k - Supplemental Result -<u>Cached</u> - <u>Similar pages</u>

#### Express-Scripts.com

Track Your **Order** · Get Helpful Answers. Get details about your coverage, ... Follow our 7 Steps to Safety to avoid common **medicine** mistakes. ... www.express-scripts.com/ - 23k - 24 Sep 2005 - <u>Cached</u> - <u>Similar pages</u>

#### Entrez PubMed

National Library of **Medicine's** search service provides access to over 10 million citations in Medline, PreMedline, and other related databases, ... www.ncbi.nlm.nih.gov/entrez/guery.fcgi - Similar pages

#### Experimental Biology and Medicine Order Form

Experimental Biology and **Medicine Order** Form. Please print this form, complete and send via fax to: 201-291-2988 or by mail to: Society for Experimental ... www.sebm.org/order.htm - 4k - <u>Cached</u> - <u>Similar pages</u>





News

Frooale

Local

more »

#### HgH Human Growin normone

Web

<u>Images</u>

... Drug Information by RxList TOP 200 DRUGS OF 2002 Drug Information for the most prescribed products with **links** to Side Effects, Drug ... http://www.pillwatch.com/. ... www.mdhealthline.com/links/online-pharmacies.htm - 26k - Supplemental Result - <u>Cached</u> - <u>Similar pages</u>

Groups

#### Transportation Links

**links**, bonuses section, search engine and directory. ... http://www.pillwatch.com DESCRIPTION : An up to date price calculator that finds the lowest price ... www.rosshealthcareclinic.com/links.html - 27k - <u>Cached</u> - <u>Similar pages</u>

## **Classifying link spam**

- Page rank assigns weight to pages.
- Spam disturbs coupling of links/ relevance.
- Identification of two spam types:
  - Link exchange/ link farms.
  - Guestbook spam.
- Evaluation of performance of classifier.

#### Crawling examples with nutch



## **Context similarity features**

 In natural networks neighboring nodes tend to be linked.



### **Context similarity features**

 Clustering coefficient measures ratio of existing links between neighbors.



## **Common features in linked pages**

- Link farms tend to be automatically created.
- Linking pages are mostly similar.
- We look for common length, IP, MD5 hash of content, link target and link sources.

#### **Common features in linked pages**

• Lets look at the average number of pages linking to our example with common length:



 $common(webpage) = \frac{3}{5}$ 

# Similarity of example and links

- Some trivial link farms tend to recite on the same server as their target.
- We look for linked/-ing pages with the same IP, length, MD5 hash as our example.

## Similarity of example and links

• Lets look at the average number of inlinking with same length as our example:



$$similar(webpage) = \frac{2}{5}$$

#### **Evaluating spam detection**

- Compare to gold standard manual labeling.
- Use implicit ranking evaluation measures:
  - Compare the average rank of the link a user clicked on before and after changing the ranking.
  - In case you are a provider: Compare the ranking of a page vs. the actual usage of this site.

#### **Performance of classifier**



- Best features:
  - Link based features.
  - Clustering Coefficient.
- Content useless for adversarial classification.



#### • Ranking search results as adversarial game.



- Spammers will react on any classifier.
  - Is there any kind of equilibrium in this game?
  - Are there "unspammable" ranking criteria?
  - Can internet usage data help win this game?

# **Enhanced Query Handling**

- Cluster search results by topic.
- Enrich queries by further information.
- Do not ask for queries terms but for questions.
## **Clustering Search Results**

- Most queries are ambiguous: panther, cluster.
- Handling ambiguity:
  - Cluster search results.
  - Add additional context to the query.
  - Categorize search results thematically.

## **Clustering Search Results**

Vivísimo*	company   products   solutions   customers   demos   press									
Clustered Results	NEW search the <u>Wikipedia</u> at <u>Clusty.com</u> Top 231 results of at least 6,571,543 retrieved for the query ccc ( <u>Details</u> )									
<ul> <li>CCC (231)</li> <li>College (55)</li> <li>Community (24)</li> <li>Colleges of Chicago (10)</li> <li>Programs (12)</li> <li>Faculty, Daley (4)</li> <li>CCC Student (3)</li> <li>Check (2)</li> <li>Wilbur Wright College (2)</li> <li>Other Topics (3)</li> <li>Codes, Cheats (42)</li> <li>Conservation Corps (13)</li> <li>Canadian (11)</li> <li>Solution (7)</li> <li>Club (9)</li> <li>Chess (2)</li> <li>Other Topics (2)</li> </ul>	<ol> <li>CCC Information Services Inc. [new window] [frame] [cache] [preview] [clusters] CCC Information Services Inc supplies the automotive claims and collision repair industries with advanced sof and Internet and wireless-enabled technology. Based in WWW.cccis.com - MSN Search 1, Looksmart 3, Lycos 4, MSN 6, Ask Jeeves 23</li> <li>Copyright.com - Copyright Licensing and Compliance Solutions from [new window] [frame] [cache] [preview] [clusters]  Do you want to register a copyright? What's New CCC Announces Improvements to Copyright.com. Learn r Permissions Building Block for Blackboard. Read the release . Learn Www.copyright.com - Wisenut 1, MSN 3, MSN Search 8</li> <li>Bombich Software: Carbon Copy Cloner [new window] [frame] [cache] [preview] [clusters]  computer? Then CCC is the tool for you! CCC makes these tasks simple by harnessing the Unix power buil features that CCC has provided www.bombich.com/software/ccc.html - MSN 1, Wisenut 6, MSN Search 7</li> <li>CCC - The Center for Computational Chemistry [new window] [frame] [cache] [preview] [clusters]  Center for Computational Chemistry. The University of Georgia, Athens, Georgia Welcome to the new CCC undergraduate student with an interest in physical chemistry, organic www.ccc.uga.edu - Open Directory 1, MSN Search 5, MSN 30</li> <li>CCC [new window] [frame] [cache] [preview] [clusters] CCC Group A world-class wireless systems provider and a forerunner in software solutions to the Industry and WWw.cccft - MSN Search 4, MSN 8, Wisenut 8</li> </ol>									

## Categorizing Queries – KDD Cup 05

- Queries in common only consist of few words.
- Categorize queries into topic hierarchies.
- KDD Cup 2005 task:
  - 1Mio queries
  - 80 categories
  - Only few examples about 70.
  - 20 international research groups took part.
  - HU Berlin has won runner up award ;)

## **Answering Questions**

<b>(</b>	Brair When was	the CCC fou	Duestion even	CRYTHING		Ask
						Ads by God
	Ads by Google	<u>Sea Turtles</u> CC was founde	Conservation	Leatherback	<u>Sea Turtle Inn</u> vo - big groups	Live it Di Marine co 2 weeks 1
	one in Berlin, th http://www.ccc.de/	ne other one in club/?language=e	Hamburg. In I I I I I I I I I I I I I I I I I I	<u>ore]</u> 12 1981 by visic	voorv Way Holland	www.reefci.
	and others in an influence the was http://www.masterli	nticipation that ay people live a ness.com/a/Chaos	information tech and communicat .Computer.Club.htm	nology will come e on this planet. 。	up strong and	Sea Turt Many spe turtles Ac
	<ul> <li>The CCC was function others in anticipy strong and influchttp://en.wikipedia.</li> </ul>	ounded in Berli pation of the fac ence the way p org/wiki/Chaos_Co	n on September ct that informatio beople live and c omputer_Club	12 , 1981 by Wa n technology wo ommunicate on t 碹 <u>[Read More]</u>	u Holland and uld come up his planet.	Marine c
	<ul> <li>Production - En was founded in effective, and rishttp://www.chicher</li> </ul>	vironment - Sa 1986 to provide sk protective er n.com/ 🍉 🖻 🏾	fety DO IT RIGH e innovative, tech ngineering servic Read Morel	T - RIGHT NOW mologically soph es.	with CCC! . CCC isticated, cost	GVI Expe unique cc project www.gvi.co
,	<ul> <li>CCC was found Group). http://www.comem</li> </ul>	ed in 1972 and berships.com/07-a	is now owned b	y the Good Sam .htmi 🔌 🖻 <u>[Reac</u>	Club (Affinity I More]	Voluntee Biospher Help with

## Building your own search engine

- Don't start out to solve this problem alone.
  - Many obstacles to take.
  - Many interesting projects seeking help.
- Examples for interesting projects:
  - Google API an API to Google search.
  - Nutch/ AspSeek standalone search engines.
  - YaCy Peer to peer search engine.

## Google API

- From Google you need to download for free:
  - The example source code.
  - Your personal API key.
- Now you can write applications, that automate querying Google Web search.
- The number of daily free queries per key is restricted to keep API traffic under control.

## Google API

- Advantages:
  - Use the power of Google for free.
  - Very easy to understand and use.
  - Quickly see the result of your work without having to buy expensive hardware or crawl web pages.
- Disadvantages:
  - Restricted number of queries per day.
  - Your application will rely on Google's service.

# Google API



#### Painting Marx

Quote: "Karl Marx and Frederick Engels had an excellent knowledge of world art and truly For this reason, he said, the appearance and poetic glorification of". Thus, From Hal Draper, Karl Marx's Theory of Revolution, Vol.1: State and Bureaucracy, He glorifies his rich Jew by painting him in pleasing colors. And, the ideas of Marx, on the contrary, took forward Blake's fourfold vision', which combined He also studied drawing and painting, combined all these. We must think of this as proven.

#### References

- 1. Marx and Engels On Literature and Art Preface.
- 2. Hal Draper: Marx and the Economic-Jew Stereotype (1977)
- 3. Marx and the Fourfold Vision of William Blake | Libertarian

This was one of Theorybot's 8 million constantly updated theories.

Show Random Theory





i.				Ag	es	ha	re	<u></u>	"ŀ	(äs	se	foi	nd	ue	"	50 - 10	0			
Share	**					144								ALC: N						
A																				
Age	8	10	12	14	16	18	20	22	24	26	28	30	32	34	36	38	40	42	44	46
Keywords (Phrase):				käse fondue																
Level of Detail:					Normal (every 2nd year)															
								Са	lcul	ate										

#### What is shown here?

Your guess (1 of 3): 0 seconds to make up your mind.



#### Nutch/ Aspseek

- Open source search engine implementations.
- For nutch online installations exist.
- From the project homepage you need to get:
  - Source code of the search engine you prefer.
  - Maybe some documentation on its internals.
- Now you can extend it according to your needs or work on tasks on the issue tracker.

## Nutch/ Aspseek

- Advantages:
  - At least Nutch is open source and fun to read.
  - There are many interesting tasks in issue tracker.
  - Active project members who know the obstacles.
  - Requests on mailing lists are answered quickly.
- Disadvantages:
  - Need expensive hardware and fast internet.
  - Same problems with spam as any other SE has.

#### Nutch in action



## YaCy – Search going peer to peer

- Distributes crawling and indexing on clients:
  - After installation YaCy can be configured to crawl the WWW and publicly share the created index.
  - Each query answered based on distributed index.
- From the project homepage you get:
  - The

## YaCy – peer to peer search

#### • Cool:

- Open source search engine already running.
- Basic idea: Distribute crawling, indexing and searching to those that want to search.
- Again: many interesting tasks, ambitious project members :)
- Uncool:
  - Will need some time to be really useful.

## YaCy – in action

Jilj	YaCv - Distributed Web Indexing - Administration
Giobai Index	
Search Page	
Help	
Distributed Crawler	
🔒 Index Create	
🔒 Index Control	
Index Monitor	
Local Proxy	
🔒 Blacklist	
🔒 Proxy Indexing	PZP WEB SEARCH
🔒 Cache Monitor	
Scookie Monitor	
Communication /	chaos
Publication	
Home Page	Max. number of results: 10 order by: Quality-Date
File Share	Resource: global 🗾 🛛 Max. search time (seconds): 10 🗾
Wiki	LIBL mask:
Messages	OTE MOST -
Personal	10 results from 23 ordered links of a total number of 1314 known. $\rightarrow$ Catch up more links from
Profile	'late' peers.
Peer Control	
Status	lopwords (to refine search): amazon 2005 heise obidos exec books ccc debian bentham
Network	archives puecher ispanari archiv politik spieger
News O Log	heise online - News-Archiv
Continue	http://www.heise.de/newsticker/archiv/2005/36/
Settings	Sa, 17 Sep 2005
Performance	
Canyuaye	BulkCarrier's WebLog : April 2005
SKIIIS	mehr beim chaos computer club posted by holger appel at 22-19 categories aktuelles 08.
The Project	nttp://www.puikCarrier.de/arcnives/04-01-2005_04-30-2005.html
Project Home	50, 20 36p 2003
Fruject News	Amazon de: Bezensionen English Books: Stations of the Tide
Deutsches Forum	http://www.amazon.de/exec/obidos/tg/stores/detail/-/books-intl-de/0380817616/reviews/
Download VaCu	ref=cm_rev_more_2/028-0663114-7383759
Dominuau Tacy	Ma 26 Cap 2006

#### Conclusion

• Search not only interesting for big companies.

• Many challenging and important tasks.

 Please take part in existing projects – its a very broad topic with many obstacles to take.

## Search – only for big companies?

- To run traditional search engines we need:
  - Large bandwidth.
  - Large hard disks.
  - Fast CPUs.
  - A lot of intelligent algorithms to perform all tasks.
- Never the less, every one can contribute.
- A selection of examples:

#### Properties, that can be faked

- Content easy to mask :)
   Key word stuffing, unreadable text.
- Content format equally easy to adjust. Enlarge query terms, put them in head lines...
- Link popularity once upon a time was hard.
   Link farms, spam guest books, exchange/ buy links.
- Anchor texts remember Google bombs?

#### WWW – less obvious properties

- Link distribution power law governed. Outliers often are spam.
- Link graph has small world properties.
  - Small average shortest path.
  - Comparably small clustering coefficient.
- Degree of nodes correlated.
   Low in-degree nodes link high in-degree nodes

and vice versa.

## Used features – page itself

- Length of URL
- Length of domain.
- # of sub domains.
- Tilde in URL?
- TLD=edu/org?
- TLD=com/biz?

- Number of inlinks.
- Number of outlinks.
- Is x a redirection?
- Size of x.
- Length of description
- Length of title.
- Number of keywords.



#### **Power law**

- Rich get richer.
  - Pages tend to link to known pages.
  - Assumption holds, if no links get removed.
- Outliers often spam pages.



### We live in a small world.

- Small clustering coefficient. People know only comparably few people.
- Small shortest path between nodes.
- Link farms often without these properties.





Properties of the WWW

• Basic architecture of search engines.

• One machine learning task in search engines.

• How can YOU contribute?



Properties of the WWW

• Basic architecture of search engines.

• One machine learning task in search engines.

• How can YOU contribute?

#### **People to contact – examples**

- Watch out for open source projects like:
  - Nutch.
  - Asp Seek.
  - YaCy.
- Google Sommer of Code.
- Faculty of computer science of your choice, e.g. Knowledge Management at HU Berlin.

#### Research

• Many open questions research topics.

- Each year:
  - Challenges e.g. KDD Cup to demonstrate state of the art technology.
  - New interesting publications on machine learning, information retrieval, link mining, ...

## Rankingbewertung

- Google Proxy User Feedback
- Explicit User Feedback
- Implicit User Feedback
- Incorporate User Location

## Click spam





Properties of the WWW

• Basic architecture of search engines.

• One machine learning task in search engines.

• How can YOU contribute?

#### Conclusion

• Search not only interesting for big companies.

• Many challenging and important tasks.

 Please take part in existing projects – its a very broad topic with many obstacles to take.

#### Violations of observations

- Outliers in power law plot often link spam.
- Networks of pages without small world properties often artificial link farms.

#### **Correlation of node degree**

- High in-degree nodes link low in-degree ones.
- Comparable to natural networks.


## Data mining the NSA^h^h/WWW

- What to do with all these data?
  - Search for relevant information what's relevant?
  - Find similar web pages.
  - Track evolving topics or find new ones.
- Are there any major problems?
  Spam, unreliable or false information.
  New types of media (Blogs, Wikis).
  - Presentation of huuuuge amount of search hits.

## Search engines from scratch

- What we need today:
  - Laaaarge hard disk, for good coverage.
  - Looots of bandwidth, to improve crawling speed.
  - Intelligent algorithms to find most important links.
- Do we have that much money? ;)

## Search engines from scratch

- What we need today:
  - Fast CPU to increase parsing/analysis speed.
  - Large hard disk to save index.
  - Intelligent algorithms for parsing pages.
  - Intelligent algorithms for scoring pages.
- Can be done with common PC.

## Search engines from scratch

- What we need today:
  - Bandwidth for query retrieval/ result delivery.
  - CPU cycles for analysis.
- Intelligent algorithms for:
  - Finding similar results.
  - Categorizing web pages.
  - Ranking matching pages.
  - Disambiguation and extension of queries.
  - You can think of many more features here :)