# Search Engines

# Exercise 2: Crawling

Dustin Lange & Saeedeh Momtazi

28 April 2011

# Googlewhacking Evaluation

My favourite terms with perfect scores:

| Potsdam | Search | HPI |
| --- | --- | --- |
| Archäologiefund | Bartkrätze | genussfreudig |
| Ericssonmotor | gartenbauwettbewerb | Hinterlegungsverfahren |
| gartenbaukurs | Gärtnerlohn | hühnerkäfig |
| Haveldampfschifffahrtsgesellschaft | Sekundärstandarddosimetrielabor | Prozessorerweiterungen |
| neutronenbeschleunigung | Wegtragsel | räuberhöhlen |
| Vergrauungsinhibitor | | Ziegenweide |

# Googlewhacking Evaluation

- 28 submissions – thanks to all!

- 10 submissions with perfect results

- First submission = first submission with perfect results

- Winner submission at 10:21 am

- Congratulations to Marika Marszalkowski and Peter Retzlaff!

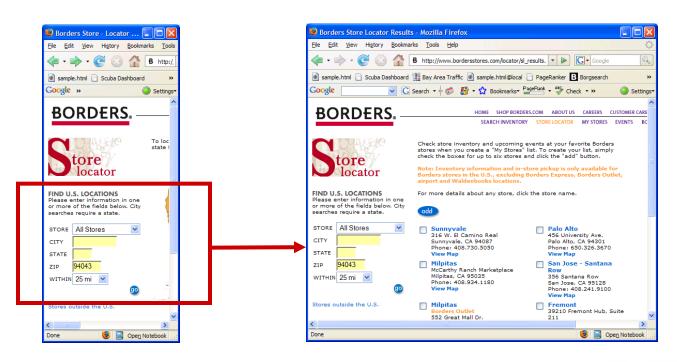  – kreuchend, fraternisierendes, druckendes

# **Task: Journal Club**

- Read one of the presented papers

- Present the key ideas (not the entire paper)

  - 5-10 minutes

  - Focus on answering the questions

  - Your fellow students should get the idea

# Paper 1: Google's Deep-Web Crawl

http://www.cs.cornell.edu/~lucja/Publications/i03.pdf



- VLDB 2008

- Cited 87 times

- Which attributes to query?

- Which values to use?

Slide contents by Jayant Madhaven, Google Inc., 2008

# Paper 1: Google's Deep-Web Crawl
# Questions

1. Introduction
   a. What is the **deep web**? Give an example. Which kind of deep web categories (from the lecture) is addressed in the paper?
2. The surfacing problem
   a. Find an **example form** in the web that is not mentioned in the paper. Give examples for the following terms using your example form: web form, inputs, selection inputs, wild card value, presentation inputs, database.
   b. What is a query template? Reuse your example.
3. Selecting query templates
   a. What are characteristics of good query templates?
   b. How is informativeness determined? (brief)
   c. Briefly describe the algorithm for **incremental template search**.
4. Generating input values
   a. Why is it difficult to generate appropriate input values for text boxes?
   b. Briefly describe the **iterative probing algorithm**. (How are seed words determined? What happens during one iteration? How are final keywords selected?)

# Paper 2: Do Not Crawl in the DUST

http://www2007.org/papers/paper194.pdf

- DUST – Different URLs with Similar Text
- Examples:
  - "http://domain.name/index.html" →
    "http://domain.name"
  - "http://news.google.com" →
    "http://google.com/news"
- How to find URL transformation rules from a list of URLs?

- WWW 2007
- Cited 37 times



Slide contents by Uri Schonfeld, Technion, 2007

# Paper 2: Do Not Crawl in the DUST
# Questions

1. Introduction
   a. What is **DUST**? Find an example that is not mentioned in the paper.
2. Problem Definition
   a. What is the **definition** of DUST rules?
   b. What is the **definition** of valid DUST rules?
3. Basic heuristics (briefly describe the three **basic heuristics**)
   a. Why are rules with large support sufficient?
   b. Why are small buckets more interesting?
   c. How can the similarity of two pages help?
4. DustBuster
   a. Briefly describe the **main algorithm for discovering likely DUST rules**. (Do not discuss details.)
   b. What are redundant rules? How can they be detected (what is the key idea)?
   c. Why is validation of DUST rules necessary? How can rules be validated?

# Selection Procedure

- Who wants to read which paper?

# Submissions & Next Exercise

- Selection:
  - Select a paper **today**: http://goo.gl/jv2ED
- Submissions:
  - Create slides to present your selected paper.
  - Send us your presentation
    - as PDF or PPT(X) or ODP:
      *SearchEngines**2**[Name1][Name2].[pdf|ppt|pptx|odp]*
    - via e-mail with subject: *Search Engines 2 [Paper 1|Paper2]*
    - to *dustin (dot) lange (at) hpi (dot) …*
    - until **4 May 2011, 5:00 pm**
- On **5 May 2011**: Be prepared to present the paper
  - English (or German)
  - Absent: Send me an e-mail in advance

# Feedback

# Thanks for Listening

- Updates
  - See website
  - Mailing list: tbd (very soon)
- Questions
  - Via e-mail:
    - dustin (dot) lange (at) hpi (…)
    - saeedeh (dot) momtazi (at) hpi (…)
  - Office: A-1.6 / A-1.7