Computer Science Presentations                                      Computer Science

10-2-2009

# (Re-) Discovering Lost Web Pages

Martin Klein
*Old Dominion University*

Michael L. Nelson
*Old Dominion University*, mnelson@odu.edu

Recommended Citation

# (Re-) Discovering Lost Web Pages

Mathematics & Computer Science Seminar
Emory University
October 2, 2009

Martin Klein & Michael L. Nelson
Department of Computer Science
Old Dominion University
Norfolk VA
www.cs.odu.edu/~{mklein,mln}

# The Problem

- Web links "break"
    - 404 http status code -- "not found"
    - "soft 404" -- http server returns "200 OK", but the resource isn't really there

- Is the content really gone?
    - Did it just move somewhere else in the web?
    - Is there a copy in search engine caches or web archives?

- To find new or different copies, we need to augment *digital preservation* with *information retrieval* techniques

# The Actors



Put a human -- lots of humans -- in the loop
for preservation purposes

# The Environment

## Web Infrastructure (WI) [McCown07]

- Web search engines (Google, Yahoo, MSN Live) and their caches
- Research Projects (CiteSeer, NSDL)
- Web archives (Internet Archive, Web Base)

NASA Technical Memorandum 109025 (Revision 1)

# A Comparison of Queueing, Cluster and Distributed Computing Systems

ftp://techreports.larc.nasa.gov/pub/techreports/larc/93/tm109025.ps.Z
http://techreports.larc.nasa.gov/ltrs/PDF/tm109025.pdf

Joseph A. Kaplan and Michael L. Nelson
Langley Research Center, Hampton, Virginia

June 1994

# Web Infrastructure: Refreshing & Migrating



12 versions found

3 remote & 4 cached versions

2 cached PDF versions

3 versions (2 nasa.gov & 1 mpg.de)

# Lapsed Website

http://web.archive.org/web/*/http://www.dl00.org/    Q~ Google

**INTERNET ARCHIVE**
**WayBack Machine**

Enter Web Address: http://   [ All ⬍ ]   ( Take Me Back )   Adv. Search  Compare Archive Pages

Searched for http://www.dl00.org/     **56** Results

Note some duplicates are not shown. See all.
* denotes when site was updated.

Gambling

Search Engine Portal

## Search Results for Jan 01, 1996 - May 03, 2005

| 1996 | 1997 | 1998 | 1999 | 2000 | 2001 | 2002 | 2003 | 2004 | 2005 |
|------|------|------|------|------|------|------|------|------|------|
| 0 pages | 0 pages | 0 pages | 1 pages | 12 pages | 11 pages | 5 pages | 13 pages | 11 pages | 0 pages |
|  |  |  | Oct 01, 1999 * | Jan 05, 2000 | Jan 24, 2001 | Jan 18, 2002 * | Feb 05, 2003 | Jan 01, 2004 * |  |
|  |  |  |  | Mar 06, 2000 | Feb 02, 2001 | May 31, 2002 * | May 02, 2003 * | Jan 03, 2004 * |  |
|  |  |  |  | Mar 10, 2000 | Feb 04, 2001 | Jun 06, 2002 | Jun 10, 2003 * | Jan 30, 2004 * |  |
|  |  |  |  | May 11, 2000 | Feb 13, 2001 | Nov 25, 2002 * | Jul 30, 2003 * | Apr 03, 2004 * |  |
|  |  |  |  | May 20, 2000 | Mar 01, 2001 | Nov 27, 2002 | Aug 05, 2003 | Apr 11, 2004 * |  |
|  |  |  |  | Jun 20, 2000 | Apr 01, 2001 |  | Aug 08, 2003 * | May 25, 2004 * |  |
|  |  |  |  | Jun 21, 2000 | Apr 05, 2001 |  | Sep 27, 2003 * | Jun 08, 2004 * |  |
|  |  |  |  | Aug 15, 2000 | Apr 14, 2001 * |  | Oct 05, 2003 * | Jun 09, 2004 |  |
|  |  |  |  | Aug 16, 2000 | Apr 21, 2001 |  | Oct 07, 2003 * | Jun 10, 2004 * |  |
|  |  |  |  | Oct 09, 2000 | Aug 31, 2001 * |  | Dec 13, 2003 * | Jun 12, 2004 |  |
|  |  |  |  | Oct 18, 2000 | Nov 27, 2001 * |  | Dec 17, 2003 | Jun 16, 2004 * |  |
|  |  |  |  | Dec 04, 2000 |  |  | Dec 27, 2003 * |  |  |
|  |  |  |  |  |  |  | Dec 28, 2003 * |  |  |

ACM DL
Conference

Porn

# URI Content Mapping Problem

| | | | |
|---|---|---|---|
| **1** | U1 → C1, U1 → C1 / A — time — B | **same** URI maps to **same** or very similar content at a later time | **2** | U1 → C1, U1 → C2 / A — time — B | **same** URI maps to **different** content at a later time |

| | | | |
|---|---|---|---|
| **3** | U1 → 404, U2 → C1 / U1 → C1 / A — time — B | **different** URI maps to **same** or very similar content at the same or at a later time | **4** | U1 → C1, U1 → ?? / A — time — B | the content can **not be found at any URI** |

# Scenario 1: Same URI, Same Content

## JCDL 2008

http://www.jcdl2008.org/
**July 2008**

http://www.jcdl2008.org/
**Today**

# Scenario 2: Same URI, Different Content

## Hypertext 2006

http://www.ht06.org/
**August 2006**

http://www.ht06.org/
**Today**

# Scenario 3a: Same Content, Different URI

## PSP 2003

http://www.pspcentral.org/events/annual_meeting_2003.html
### August 2003

http://www.pspcentral.org/events/archive/annual_meeting_2003.html
### Today



February 3-5, 2003

The Association of American Publishers, PSP Division

*invites you to join us for*

SMART CONTENT: NEW WAYS TO ADD VALUE

*2003 PSP Annual Conference*

Renaissance Mayflower Hotel
Washington, DC

- Download the brochure (.doc)
- Click here to register (.doc)
- Exhibitor Information

Program of Events

MONDAY, FEBRUARY 3, 2003

8:30am-
3:00pm
*Pre-Conference Session (separate registration fee)*
**Where�s the Customer for Smart Content?**
*(pre-conference full-day seminar/separate registration fee)*
*Moderator:* **Eileen Dolan**, Vice President, Wiley InterScience, John Wiley & Sons, Inc.

This seminar will discuss:

- Why the online customer is important
- Identifying the needs, desires and priorities of the online customer
- Creating value for the online customer

The landmark *Usage Statistics* White Paper, published by the PSP Electronic Information Committee, will be available.

4:00pm **Conference Opens**

4:00pm-
6:00pm
***Public Policy is Everyone's Concern:* Copyright -- A Perennial Rallying Point**
*Moderator:* **Marc Brodsky**, Executive Director & CEO, American Institute of Physics



February 3-5, 2003

The Association of American Publishers, PSP Division

*invites you to join us for*

SMART CONTENT: NEW WAYS TO ADD VALUE

*2003 PSP Annual Conference*

Renaissance Mayflower Hotel
Washington, DC

- Download the brochure (.doc)
- Click here to register online / Postal mail (.doc)
- Exhibitor Information

PSP Í03 Annual Conference Planning Committee
Pieter Bolman (PSP ExCo Chair), Patrick Bernuth, Donald Burden, Nigel Fletcher-Jones, Andrew Grabois, Doug LaFrenier, Eric Massant, Ted Nardin, Hill Slowinski *AAP:* Barbara Meredith, Sara FIrestone

Program of Events

MONDAY, FEBRUARY 3, 2003

8:00am-
9:00am
**Continental Breakfast**

3:00pm-
4:00pm
**Your chance to visit the New Technologies/Services Exhibitors**

9:00am-
3:00pm
*Pre-Conference Session (separate registration fee)*
**WHERE IS THE USER FOR YOUR SMART CONTENT?**
*Produced by the AAP/PSP Electronic Information Committee*

9:00am-9:15am
*Moderator:* **Eileen Dolan**, Vice President, Wiley InterScience, John Wiley & Sons, Inc.

# Scenario 3b: Similar Content, Different URI

## ECDL 1999

http://www-rocq.inria.fr/EuroDL99/
**October 1999**

http://www.informatik.uni-trier.de/~ley/db/conf/ercimdl/ercimdl99.html
**Today**

# Scenario 4: Content Not Findable At Any URI

## Greynet 1999

http://www.konbib.nl/infolev/greynet/2.5.htm

**1999**

**Today**

?                    ?

**Miller**: A lot o' people don't realize what's really going on. They view life as a bunch o' unconnected incidents 'n things. They don't realize that there's this, like, lattice o' coincidence that lays on top o' everything. Give you an example; show you what I mean: suppose you're thinkin' about a plate o' shrimp. Suddenly someone'll say, like, *plate, or shrimp, or plate o' shrimp* out of the blue, no explanation. No point in lookin' for one, either. It's all part of a cosmic unconsciousness.

**Otto**: You eat a lot of acid, Miller, back in the hippie days?

# Synchronicity



- Experience of causally unrelated events occurring together in a meaningful manner

- Events reveal underlying pattern, framework bigger than any of the synchronous systems

- Carl Gustav Jung (1875-1961)

  - "meaningful coincidence"

# Synchronicity Architecture



- Firefox extension catches 404 error (or initiated by user if a "soft" 404 is suspected)

- Discovers copy of missing page in WI (1) and provides to user (2)

- Generates a search engine query based on what the missing page is "about" (3)

- Finds old content at new URI or provides a "good enough" alternative page (4,5,6)

# What Was That Web Page About?

- If an "old" copy can be found:
  - Lexical Signatures
  - \<title>…\</title>
- If no archived/cached copy:
  - Tags
  - Link Neighborhoods; LSs, anchor tags

```
GET https://user:pass@api.del.icio.us/v1/posts/suggest?url=http://yahoo.com/

<?xml version="1.0" encoding="UTF-8"?>
<suggest>
  <popular>web</popular>
  <popular>tools</popular>
  <popular>searchengines</popular>
  <recommended>yahoo!</recommended>
  <recommended>yahoo</recommended>
  <recommended>web</recommended>
  <recommended>tools</recommended>
  <recommended>search</recommended>
  <recommended>reference</recommended>
  <recommended>portal</recommended>
  <recommended>news</recommended>
</suggest>
```

A → ?  ← C

B → ?

http://nicnichols.com/

Google

HOME    PHOTOGRAPHY ▶    ABOUT

Search ...

nicnichols.c

FEATURED ITEM

DOCUMENTA
PHOTOGRA

LINKS

Blog : Four

Online Shop

RECE

| LS | NICNICHOLS NICHOLS NIC STUFF SHOOT COMMAND PENITENTIARY |
|---|---|
| Title | NICNICHOLS.COM : DOCUMENTARY TOY CAMERA PHOTOGRAPHY OF NIC NICHOLS : HOLGA, LOMO AND OTHER LO-FI CAMERAS! |
| Tags | PHOTOGRAPHY BLOG PHOTOGRAPHER PORTIFOLIO PORTFOLIO INSPIRATION PHOTOGRAPHERS |
| LNLS | NICNICHOLS PHOTO SPACER VIEW PHIREBRUSH SUBMISSION BOONIKA |

Table 1: Data Obtained from www.nicnichols.com

VIEW ALL A

RECENT TWITTER POSTS

Find:    Next    Previous    Highlight all

Done

# What is a Signature?
**(aka "message digest", examples include "md5" and "sha-1")**



image from Eddie Kohler http://www.cs.ucla.edu/~kohler/

# What is a Lexical Signature?

- First introduced by Phelps and Wilensky

[Phelps00]

- Small set of terms capturing the "aboutness" of a document
  - Phelps and Wilensky assumed 5
- "lightweight metadata"

**Resource**

**Abstract**

**LS**

REMOVAL
HIT
RATE
PROXY
CACHE

Query
Google

"Removal Policies in Network Caches for World-Wide Web Documents"

# LSs as Proposed by Phelps and Wilensky

- "Robust Hyperlink Cost Five Words Each"
- Append LS to URL:

http://www.cs.berkeley.edu/~wilensky/NLP.html

becomes:

http://www.cs.berkeley.edu/~wilensky/NLP.html?lexical-signature=texttiling+wilensky+disambiguation+subtopic+iago

- Limitations:

  1. Applications (browsers) need to be modified to exploit LSs

  2. LSs need to be computed a priori

  3. Works well with most URLs but not with all of them

# Lexical Signatures -- Examples

| Rank/Results | URL | LS |
|:---:|:---|:---|
| **1/1** | **http://www.cs.berkeley.edu/˜wilensky/NLP.html** | **texttiling wilensky disambiguation subtopic iago** <br> http://www.google.com/search?q=texttiling+wilensky+disambiguation+subtopic+iago |
| **1/221,000** (1/174,000 in 01/2008) | **http://www.loc.gov** | **library collections congress thomas american** <br> **http://www.google.com/search?q=library+collections+congress+thomas+american** |
| **1/51** (2/77 in 01/2008) | **http://www.jcdl2008.org** | **libraries jcdl digital conference pst** <br> **http://www.google.com/search?q=libraries+jcdl+digital+conference+pst** |
| **0/10** | **http://www.dli2.nsf.gov** | **nsdl multiagency imls testbeds extramural** <br> **http://www.google.com/search?q=nsdl+multiagency+imls+testbeds+extramural** |

A "Googlewhack" (http://en.wikipedia.org/wiki/Googlewhack) can be thought of as a two-term LS that produces a 1/1 ranking.

# Generating LSs

- ## Term Frequency (TF)
  - "How often does this term occur in this document?"

- ## Inverse Document Frequency (IDF)
  - "In how many documents does this term appear?"

$$TF_{ij} = \frac{f_{ij}}{m_i}$$
$$f_{ij} = freq\ of\ j\ in\ i$$
$$m_i = max\ freq\ in\ i$$

$$IDF_j = \log\left(\frac{N}{n_j}\right) + 1$$
$$N = total\ number\ of\ documents$$
$$n_j = number\ of\ documents\ j\ occurs\ in$$

# Generating LSs

• Park et al. [Park03] investigated performance of various LS generation algorithms

• Evaluated "tunability" of TF and IDF
- • Weight on TF increases recall (completeness, ex. "photography, blog")
- • Weight on IDF improves precision (exactness, ex. "nicnichols, penitentiary")

• Computed IDF on closed system (not live web)

• Also assumed "5" to be a good number

• Compared results after 6 months, but did not do an in-depth analysis of LSs over time

# Theoretical Underpinnings of Synchronicity

• Estimating IDF values for the Web (WIDM 2008, ECIR 2009)

• Investigated how lexical signatures change over time (ECDL 2008)

• Compared retrieval performance of lexical signatures with titles, tags and lexical signatures generated from link neighborhoods (submitted)

• Investigated how titles change over time (InDP 2009, in preparation)

# Hacks for Estimating IDF



1. everyone knows this value is flaky
2. get N from: http://www.worldwidewebsize.com/

# For LS purposes, it doesn't matter much...

URL: http://www.perfect10wines.com
Year: 2007
Union: 12 unique terms

| Rank | Local Universe | | Screen Scraping | | N-grams | |
|------|------|--------|------|--------|------|--------|
| | Term | TF-IDF | Term | TF-IDF | Term | TF-IDF |
| 1 | perfect | 7.77 | wines | 5.97 | wines | 7.56 |
| 2 | wines | 6.95 | robles | 5.3 | perfect | 7.25 |
| 3 | 10 | 6.57 | perfect | 4.35 | robles | 7.18 |
| 4 | paso | 6.29 | paso | 4.27 | paso | 6.93 |
| 5 | wine | 6.18 | wine | 3.26 | wine | 4.86 |
| 6 | robles | 5.4 | sauvignon | 3.16 | 10 | 4.52 |
| 7 | sauvignon | 3.54 | chardonnay | 3.15 | chardonnay | 3.99 |
| 8 | cabernet | 3.54 | robles84 | 3.11 | sauvignon | 3.93 |
| 9 | monterey | 3.36 | cabernet | 3.09 | cabernet | 3.89 |
| 10 | chardonnay | 3.36 | enthusiast85 | 2.91 | monterey | 3.49 |

**Comparing LSs**

Top 5, 10 and 15 terms

LC – local universe
SC – screen scraping
NG – N-Grams

~4 of 5 LS terms are the same

# How Does Google N-grams TC Relate to DF?

- Google N-grams has only Term Count (TC), not Document Frequency
  - where TC >= DF
  - http://googleresearch.blogspot.com/2006/08/all-our-n-gram-are-belong-to-you.html

| $d_1 = Please\ Please\ Me$ | | | | $d_3 = All\ You\ Need\ Is\ Love$ | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| $d_2 = Can't\ Buy\ Me\ Love$ | | | | $d_4 = Long,\ Long,\ Long$ | | | | | |

| Term | All | Buy | Can't | Is | Love | Me | Need | Please | You | Long |
|---|---|---|---|---|---|---|---|---|---|---|
| TC | 1 | 1 | 1 | 1 | 2 | 2 | 1 | 2 | 1 | 3 |
| DF | 1 | 1 | 1 | 1 | 2 | 2 | 1 | 1 | 1 | 1 |

- Idea: compare TC & DF in a known collection, then compare that collection's TC to the Google N-grams TC
  - we used ukWaC, from WaCKy: http://wacky.sslmit.unibo.it/

# TC Ranks vs DF Ranks Within ukWaC

# Rank Correlation Within ukWaC



semi-log scale

# TC Frequencies in ukWaC and N-Grams

# LS Evolution Over Time

Copies of web pages from the IA (1996-2007)



300 Random URLs, winnowed to 98, 10493 observations over 12 years

# Evolution Over Time -- Example

## 10-term LSs generated for
http://www.perfect10wines.com

|    | 2005 | | 2006 | | 2007 | |
|----|------|-------|------|-------|------|-------|
|    | **Term** | **Score** | **Term** | **Score** | **Term** | **Score** |
| 1  | wines | 8.56 | wines | 6.52 | wines | 5.25 |
| 2  | perfect | 5.00 | wine | 4.80 | wine | 4.50 |
| 3  | wine | 3.03 | perfect | 4.70 | paso | 4.50 |
| 4  | 10 | 2.60 | 10 | 3.45 | perfect | 4.10 |
| 5  | monterey | 2.24 | paso | 3.01 | robles | 3.75 |
| 6  | chardonnay | 2.24 | robles | 2.89 | 10 | 3.40 |
| 7  | merlot | 2.20 | monterey | 2.79 | monterey | 2.25 |
| 8  | robles | 1.99 | chardonnay | 2.79 | cabernet | 2.25 |
| 9  | paso | 1.99 | ripe | 1.86 | chardonnay | 2.25 |
| 10 | blonde | 1.38 | vanilla | 1.86 | sauvignon | 2.25 |

$$for\ all\ terms: \quad |\cup| = 14\ and\ |\cap| = 8$$

# Two Methods for Measuring Evolution

## Idea

- Generate LSs from copies of URLs
- Conduct overlap analysis

# Evolution Over Time - Rooted

| compare to | \multicolumn{11}{c}{Year of First Observation} | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1996 | 1997 | 1998 | 1999 | 2000 | 2001 | 2002 | 2003 | 2004 | 2005 | 2006 |
| 1997 | 0.33 | | | | | | | | | | |
| 1998 | 0.13 | 0.33 | | | | | | | | | |
| 1999 | 0.13 | 0.20 | 0.56 | | | | | | | | |
| 2000 | 0.13 | 0.33 | 0.49 | 0.51 | | | | | | | |
| 2001 | 0.20 | 0.27 | 0.31 | 0.46 | 0.58 | | | | | | |
| 2002 | 0.13 | 0.33 | 0.33 | 0.32 | 0.48 | 0.64 | | | | | |
| 2003 | 0.13 | 0.13 | 0.40 | 0.40 | 0.47 | 0.54 | 0.66 | | | | |
| 2004 | 0.13 | 0.13 | 0.36 | 0.35 | 0.40 | 0.53 | 0.60 | 0.66 | | | |
| 2005 | 0.13 | 0.07 | 0.38 | 0.37 | 0.37 | 0.42 | 0.50 | 0.63 | 0.58 | | |
| 2006 | 0.13 | 0.20 | 0.31 | 0.35 | 0.38 | 0.48 | 0.51 | 0.46 | 0.62 | 0.80 | |
| 2007 | 0.20 | 0.20 | 0.27 | 0.29 | 0.37 | 0.44 | 0.50 | 0.37 | 0.52 | 0.60 | 0.90 |

- Little overlap between the early years and more recent ones

- Highest overlap in the first 1-2 years after creation of the LSs

- Rarely peaks after that – once terms are gone they do not return

# Evolution Over Time - Sliding

| comparison | Year of First Observation | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1996 | 1997 | 1998 | 1999 | 2000 | 2001 | 2002 | 2003 | 2004 | 2005 | 2006 |
| **1996-1997** | 0.33 | | | | | | | | | | |
| **1997-1998** | 0.40 | 0.33 | | | | | | | | | |
| **1998-1999** | 0.73 | 0.27 | 0.56 | | | | | | | | |
| **1999-2000** | 0.53 | 0.40 | 0.49 | 0.51 | | | | | | | |
| **2000-2001** | 0.47 | 0.87 | 0.56 | 0.62 | 0.58 | | | | | | |
| **2001-2002** | 0.53 | 0.73 | 0.51 | 0.52 | 0.63 | 0.64 | | | | | |
| **2002-2003** | 0.60 | 0.73 | 0.67 | 0.55 | 0.67 | 0.64 | 0.66 | | | | |
| **2003-2004** | 0.93 | 0.80 | 0.76 | 0.69 | 0.80 | 0.83 | 0.73 | 0.66 | | | |
| **2004-2005** | 0.87 | 0.80 | 0.73 | 0.66 | 0.82 | 0.68 | 0.83 | 0.74 | 0.58 | | |
| **2005-2006** | 0.93 | 0.47 | 0.71 | 0.72 | 0.77 | 0.72 | 0.84 | 0.51 | 0.76 | 0.80 | |
| **2006-2007** | 0.87 | 0.53 | 0.80 | 0.68 | 0.83 | 0.76 | 0.81 | 0.49 | 0.68 | 0.80 | 0.90 |

- Overlap increases over time

- Seem to reach steady state around 2003

# Performance of LSs

**Idea**

• Measure performance in respect to age of LS and number of terms it contains

• Query Google search API with LSs

• Identify URL in result set:

1. Top ranked
2. Ranked between 2-10
3. Ranked between 11-100
4. Ranked beyond 100 (considered undiscovered)

# Performance – Number of Terms

| | 1 | 2-10 | 11-100 | ≥101 | MR |
|---|---|---|---|---|---|
| **2-term** | 24.3 | 14.9 | 13.2 | 47.6 | 53.1 |
| **3-term** | 40.2 | 15.0 | 15.0 | 29.8 | 36.5 |
| **4-term** | 43.9 | 15.7 | 11.4 | 29.0 | 33.8 |
| **5-term** | 47.0 | 19.4 | 3.4 | 30.2 | 32.7 |
| **6-term** | 51.2 | 11.4 | 3.4 | 34.1 | 36.0 |
| **7-term** | 54.9 | 9.4 | 1.5 | 34.2 | 35.5 |
| **8-term** | 49.8 | 7.7 | 2.2 | 40.4 | 41.9 |
| **9-term** | 47.0 | 6.6 | 0.9 | 45.5 | 46.4 |
| **10-term** | 46.1 | 4.0 | 0.9 | 49.0 | 49.8 |
| **15-term** | 39.8 | 0.8 | 0.6 | 58.9 | 59.5 |

- 2-, 3- and 4-term LSs perform poorly
- 5-, 6- and 7-term LSs seem best
    - Top mean rank (MR) value with 5 terms
    - Most top ranked with 7 terms
    - Binary pattern: either top 10 or undiscovered
- 8+ terms -- decreased performance

# Performance – Age

## Score of LSs consisting of 2, 5, 7 and 10 terms

## Fair

## Optimistic



•Example, scores for the position of an URL in a list of 10:
• fair: 10/10, 9/10, 8/10 ... 1/10, 0
• optimistic: 1/1, 1/2, 1/3 ... 1/10, 0

# Titles (TI), 5- & 7-term Lexical Signatures (LS5, LS7), Tags (TA)

| | Google | | | | Yahoo | | | | MSN | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Top | Top10 | Top100 | Undis | Top | Top10 | Top100 | Undis | Top | Top10 | Top100 | Undis |
| LS5 | 50.8 | 12.6 | 4.2 | 32.4 | **67.6** | 7.8 | 2.3 | 22.3 | 63.1 | 8.1 | 1.6 | 27.2 |
| LS7 | 57.3 | 9.1 | 2.6 | 31.1 | **66.7** | 4.5 | 1.9 | 26.9 | 62.8 | 5.8 | 1.6 | 29.8 |
| TI | **69.3** | 8.1 | 2.9 | 19.7 | 63.8 | 8.1 | 0.6 | 27.5 | 61.5 | 6.8 | 1.0 | 30.7 |
| TA | 2.1 | 10.6 | 12.8 | 75.5 | **6.4** | 17.0 | 12.8 | 63.8 | 0 | 8.5 | 10.6 | 80.9 |

Table 2: Relative Number of URLs Retrieved with one Single Method from Google, Yahoo and MSN

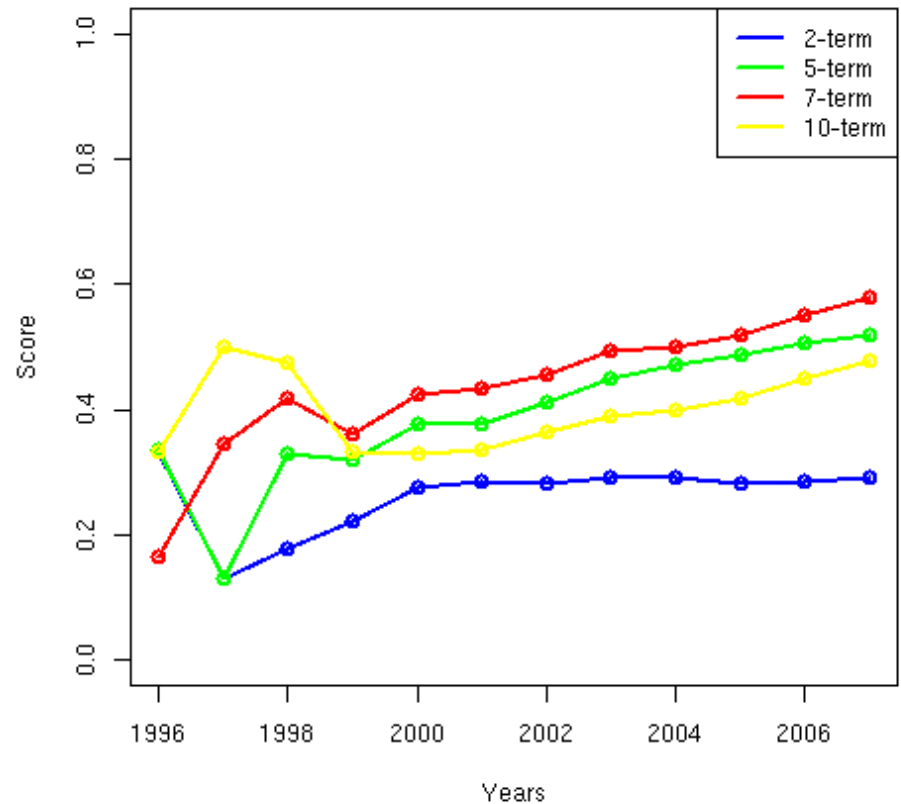| | Google | | | | Yahoo | | | | MSN | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Top | T10 | T100 | Undis | Top | T10 | T100 | Undis | Top | T10 | T100 | Undis |
| LS5-TI | 65.0 | 15.2 | 6.1 | 13.6 | **73.8** | 10.0 | 2.3 | 14.0 | 71.5 | 10.0 | 1.9 | 16.5 |
| LS7-TI | 70.9 | 11.7 | 4.2 | 13.3 | **75.7** | 7.4 | 1.9 | 14.9 | 73.8 | 9.1 | 1.9 | 15.2 |
| TI-LS5 | 73.5 | 9.1 | 3.9 | 13.6 | **75.7** | 9.1 | 1.3 | 13.9 | 73.1 | 9.1 | 1.3 | 16.5 |
| TI-LS7 | 74.1 | 9.4 | 3.2 | 13.3 | **75.1** | 8.7 | 1.3 | 14.9 | 74.1 | 9.1 | 1.6 | 15.2 |
| LS5-TI-LS7 | 65.4 | 15.2 | 6.5 | 12.9 | **73.8** | 10.0 | 2.6 | 13.6 | 72.5 | 10.4 | 2.6 | 14.6 |
| LS7-TI-LS5 | 71.2 | 11.7 | 4.2 | 12.9 | **76.4** | 7.8 | 2.3 | 13.6 | 74.4 | 9.1 | 1.9 | 14.6 |
| TI-LS5-LS7 | 73.8 | 9.1 | 4.2 | 12.9 | **75.7** | 9.1 | 1.6 | 13.6 | 74.1 | 9.4 | 1.9 | 14.6 |
| TI-LS7-LS5 | 74.4 | 9.4 | 3.2 | 12.9 | **75.7** | 9.1 | 1.6 | 13.6 | 74.8 | 9.1 | 1.6 | 14.6 |
| LS5-LS7 | 52.8 | 12.9 | 6.5 | 27.8 | **68.0** | 7.8 | 2.9 | 21.4 | 64.4 | 8.4 | 2.6 | 24.6 |
| LS7-LS5 | 59.9 | 9.7 | 2.6 | 27.8 | **71.5** | 4.9 | 2.3 | 21.4 | 66.7 | 7.1 | 1.6 | 24.6 |

Table 3: Relative Number of URLs Retrieved with Two or More Methods Combined

500 random URLs from dmoz.org winnowed to 309 (only 47 of 309 had tags in delicious.com).
Due to query restrictions, link neighborhood only run on Yahoo -- results were similar to tags.

# Number of Title Changes and Observations in the IA

ordered in increasing order by:
1) observations
2) changes



- generally low number of changes

- max changes: 25

- number of observations does not impact the number of changes

6000 random URLs from dmoz.org, winnowed to 1090 URLs and 100k+ observations

# Mean Time Delta Between Changes
# Time Span Between First and
# Last Observation in the IA

ordered in
increasing order
by:
1) observations
2) changes



- time span between
  observations
  decreases with
  increasing number of
  observations

- overall time span just
  slightly increases

- URLs with many
  observations are being
  crawled frequently in a
  short period of time

# Mean Levenshtein Scores of all Titles - Sliding



- 5 URLs with score = 0

- 85% of URLs with score >=0.8

- titles rarely change drastically

# Mean Levenshtein Scores of all Titles - Rooted



- 9 URLs with score = 0

- 56% of URLs with score >=0.8

- titles more likely to change compared to their first observation

mean Levenshtein score
sliding: 0.84 rooted: 0.29

1998-01-27
Sun Software Products Selector Guides -Solutions Tree

1999-02-20
Sun Software Solutions

2002-02-01
Sun Microsystems Products

2002-06-01
Sun Microsystems - Business & Industry Solutions

2003-08-01
Sun Microsystems - Industry & Infrastructure Solutions

2004-02-02
Sun Microsystems - Solutions

2004-06-10
Gateway Page - Sun Solutions

2006-01-09
Sun Microsystems Solutions & Services

2007-01-03
Services & Solutions

2007-02-07
Sun Services & Solutions

2008-01-19
Sun Solutions

mean Levenshtein score
sliding: 0.68 rooted: 0.15

2000-06-19
DataCity of Manassas Park Main Page

2000-10-12
DataCity of Manassas Park sells Custom
Built Computers & Removable Hard Drives

2001-08-21
DataCity a computer company in Manassas
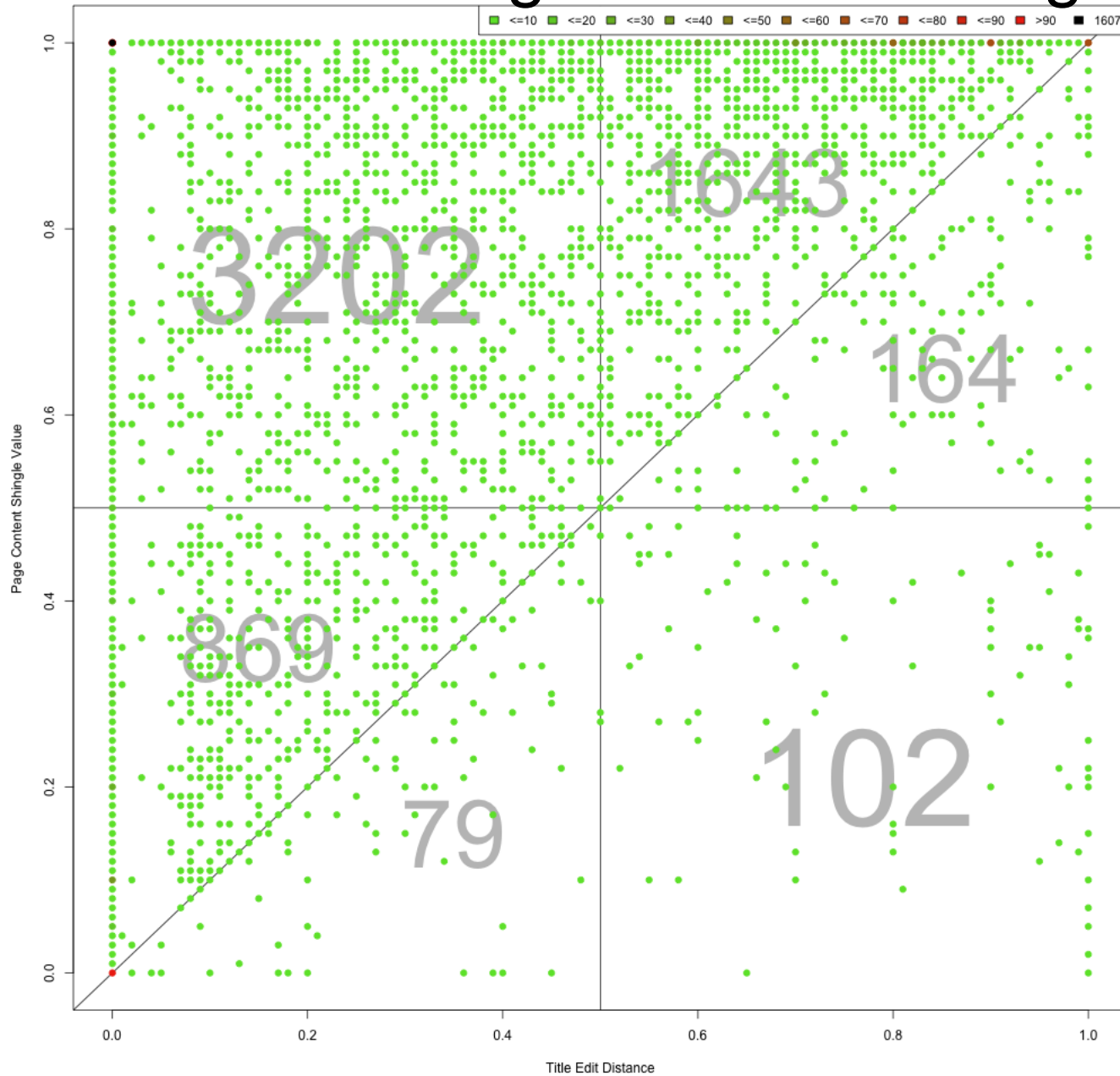Park sells Custom Built Computers & Removable
Hard Drives

2002-10-16
computer company in Manassas Virginia sells
Custom Built Computers with Removable Hard
Drives Kits and Iomega 2GB Jaz Drives
(jazz drives) October 2002 DataCity
800-326-5051 toll free

2006-03-14
Est 1989 Computer company in Stafford
Virginia sells Custom Built Secure
Computers with DoD 5200.1-R Approved
Removable Hard Drives, Hard Drive Kits
and Iomega 2GB Jaz Drives (jazz drives),
introduces the IllumiNite&reg; lighted
keyboard DataCity 800-326-5051 Service
Disabled Veteran Owned Business SDVOB

# Content Change vs. Title Change

# Conclusions & Future Work

- LSs decay over time, Titles decay less
  - Rooted: quickly after generation
  - Sliding: seem to stabilize
- Titles give comparable performance to LSs
- Titles + LSs give better performance

- Future work:
  - can we know in advance if a title is "good"? (i.e., not "welcome to my home page")
  - can we use tags to augment titles / LS?
  - how big should a link neighborhood be?

- Contact us to get a beta version of the Firefox extension (real soon now!)

# Necronomicon