

# Agreeing to Disagree: Search Engines and their Public Interfaces

Frank McCown  
Old Dominion University  
Computer Science Department  
Norfolk, Virginia, USA 23529  
fmccown@cs.odu.edu

Michael L. Nelson  
Old Dominion University  
Computer Science Department  
Norfolk, Virginia, USA 23529  
mln@cs.odu.edu

## ABSTRACT

Many studies of popular commercial search engines have used either the web user interface (WUI) or the application programming interface (API) to obtain their results, but anecdotal evidence suggests that the interfaces produce different results. We provide the first in depth quantitative analysis of the results produced by the Google, MSN and Yahoo API and WUI interfaces. We submitted 3500 identical queries to Google, MSN and Yahoo every day for five months (from late May to Oct 2006) using both interfaces. We measured the disagreements between the interfaces and found significant discrepancies in several categories. Google's API failed to produce a single top 100 result that was identical to the Google WUI results, and MSN and Yahoo only produced identical results 0.2% of the time. In general, we found MSN to produce the most consistent results between their two interfaces. Our findings suggest that the API indexes are *not* older, but they are *probably* smaller for Google and Yahoo.

We also examined how search results decay over time and built predictive models based on the observed decay rates. We found the decay rates for popular term results differ significantly from computer science term results for Google and MSN, and the top 10 results decay at a much slower rate than do the top 100 results. Based on our observations, it can take over a year for half of the top 10 results to a popular query to be replaced with other results in Google and Yahoo; for MSN it can take only 2-3 months.

**Categories and Subject Descriptors:** H.3.5 [Information Storage and Retrieval] Online Information Services: Web-based services

**General Terms:** Measurement, Experimentation, Design

**Keywords:** search engine results, search engine interfaces, distance measurements

## 1. INTRODUCTION

Commercial search engines have long been used in academic studies. Sometimes the search engines themselves are being studied, and sometimes they are used to study the Web or web phenomena. In the past, researchers have either manually submitted queries to the search engine's web user interface (WUI), or they have created programs that automate the task of submitting queries. The returned re-

sults have been processed manually or by programs that rely on brittle screen-scraping techniques.

But data collection mechanisms for search engines have changed in the past few years. Three of the most popular commercial search engines (Google, MSN and Yahoo) have developed freely available APIs for accessing their index. Researchers now can use these APIs in their automated data collection processes in place of the screen-scraping techniques they used in the past.

Unfortunately, the APIs do not always give the same results as the WUI. The list serves and news groups that cater to the API communities are full of questions regarding the perceived differences in results between the two. None of the search engines publicly disclose the inner workings of their APIs, so users are left wondering if the APIs are giving second-rate data. This anecdotal evidence has led some researchers to question the validity of their findings. For example, Bar-Yossef and Gurevich [8] state that "due to legal limitations on automatic queries, we used the Google, MSN, and Yahoo! web search APIs, which are, reportedly, served from older and smaller indices than the indices used to serve human users." Other researchers may be totally unaware of the differences. When writing about a 2004 experiment using the Google search engine, Thelwall [45] stated that "the Google API... could have been used to automate the initial collection of the results pages, which should have given the same outcome" as using the WUI.

The purpose of this study is to examine the differences between what is reported between the WUI and API when queried with a variety of queries that researchers and users frequently use. This is the first study to provide quantitative analysis comparing the WUI and API of the top three commercial search engines. We have queried Google, MSN and Yahoo every day for 5 months (late May to Oct 2006) using both the WUI and API. We issued queries for 100 search terms, queries to test if 100 URLs were indexed and cached, queries ask for backlinks to 100 URLs and queries to ask for the total number of indexed pages for 100 websites.

Our experiments allow us to address the question of whether the APIs are pulling from older and smaller indexes. We also examine how search results decay over time and provide predictive models for determining the half-lives of search results.

## 2. BACKGROUND AND RELATED WORK

### 2.1 Search Engine APIs

In early 2002, Google generated a great deal of fanfare

Copyright is held by the author/owner(s).

WWW2007, May 8-12, 2007, Banff, Canada.

when it became the first search engine to release a free-to-use SOAP-based API for accessing its index. The mechanics of the API have changed little over the years; in fact, Google still maintains that the service is in beta [20]. Google requires users to register for a license key and allows 1000 requests per key to be issued per day. Google at one time allowed researchers to request additional queries per day, but they are no longer granting such requests, even if researchers are willing to pay [21].

Yahoo and MSN released public APIs in early and late 2005, respectively, but in an effort to out-do Google, they greatly reduced the usage restrictions of their APIs. Yahoo’s REST-based API allows 5000 daily requests per IP address, per application ID [52]. MSN’s SOAP API allows 10,000 daily requests per IP address, per application ID [38]. By allowing more requests per day and only enforcing the quota to a particular IP address, Yahoo and MSN allow developers to create applications for third parties which do not require users to register for a license key.

All three search engines have expanded their APIs to include facilities like maps, spelling, news, etc. Despite their popularity (or because of it), none of the search engines give direct technical support for their APIs; there are on-line forums [22, 49] and mailing lists [51] that developers can use, and representatives from the search engines keep an eye on the discussions and occasionally respond.

Researchers have been attracted to using the APIs not only because it eases data collection activities, but because of the legal restrictions against automating data collection using the WUIs of several search engines. According to Google’s terms of service [19], users “may not send automated queries of any sort to Google’s system without express permission in advance from Google.” MSN and Ask.com have similar policies [37, 1]. Only Yahoo does not explicitly restrict automated queries in their terms of service [50].

## 2.2 Using APIs in Academic Studies

There are numerous studies which have evaluated search engines exclusively through either the WUI or API. For example, early studies on the size of the indexable web and search engine overlap automated thousands of queries against the WUI when no API was available [9, 28]. The 2005 size and overlap experiment performed by Gulli and Signorini [24] also used the WUI exclusively, although recent measurements taken by Bar-Yossef and Gurevich [8] used only the APIs. In some of our recent work [35], we have used a combination of APIs and WUIs in the same experiment because the Google API’s limit of 10 results per query.

Because the Yahoo and MSN APIs have only been accessible for a little over a year, most academic studies up to now have used only the Google API. A number of studies have used the Google API to perform general search queries where ranking of the items is important [11, 27, 32, 25, 41], search for backlinks [39] and determine if particular content was indexed or not [8, 34, 44]. In all these studies, the findings would be altered if the Google API reported significantly different results than the WUI.

Evaluation of API search results has received little attention in the literature. Bar-Ilan [5] provided an excellent summary of current search engine deficiencies in 2005, but she did not evaluate any search engine APIs. Mayr and Tosques [30] performed a limited study of the Google API

in a 2004-2005 experiment. They found a large discrepancy between the results obtained between the Google API and WUI; the WUI produced as much as six times the number of hits for the word *webometrics* over the period of a year. They also found more differences in the search results obtained from the API than the WUI, but they do not provide any measurements.

A number of studies compare the results obtained from one search engine with others when given the same query [4, 6, 14, 15, 43], and/or they examine how the results change over time [2, 3, 6, 7, 46]. The studies may use a handful of queries (e.g., [4]) or a very large data set (e.g., [43]), and typically just the top 10 to 50 results are examined since these are the results most frequently examined by human searchers. While a few studies employ humans to judge the relevance or ranking of the results (e.g., [46]), most compare the results by applying a number of difference measures like overlap, Kendall tau, Spearman’s footrule, and other measures that take into account rank ordering. We have employed these distance measures (more detail in Section 2.3) in our study and do not attempt to examine relevance or technical precision [3]. All these studies have concluded that most search engines produce very different results when given the same query, and that the results also change at different rates over time. This has prompted active research into meta-searching [17, 43]. All of the studies cited have used results obtained using the WUI interfaces, and none of them have examined how the results would be different if using the search engine APIs. We are unaware of any studies that have used queries to measure website indexing or backlinks over time.

## 2.3 Comparing Search Engine Results

Most of our analysis focuses on comparing the top 10 and 100 search results obtained using a variety of query terms. In order to compare the results over time and the differences between the two interfaces, we review three distance measures which have been used in similar studies to compare top  $k$  search results.

As discussed in Section 2.2, computing overlap is a common way to compare search results. The formal definition we will use for the normalized overlap  $P$  of two top  $k$  lists (lists of size  $k$ )  $\tau_1$  and  $\tau_2$  is:

$$P = |\tau_1 \cap \tau_2|/k \quad (1)$$

For example, consider the following top 4 search results where each letter represents a unique URL:

1. ABCD
2. EDAF

In this case URL A was returned first in list 1 and third in list 2, B was returned second in list 1 but is not present in list 2, etc. The overlap  $P$  for these two lists is  $2/4 = 0.5$ .

Overlap does not convey changes in ranking but merely set membership. A well-known measure like Kendall tau distance is a stronger distance measure that takes into account changes of rank between two lists. Kendall tau distance counts the number of pairwise disagreements (discordant pairs) between two lists and is equivalent to the number of swaps needed to transform one list to another using the bubblesort algorithm. Kendall tau distance  $K_\tau$  for two lists

$\tau_1$  and  $\tau_2$  is defined:

$$K_\tau(\tau_1, \tau_2) = \sum_{\{i,j\} \in P} \bar{K}_{i,j}(\tau_1, \tau_2) \quad (2)$$

where  $P$  is the set of unordered pairs of distinct elements in  $\tau_1$  and  $\tau_2$ , and  $\bar{K}_{i,j}(\tau_1, \tau_2) = 1$  if  $i$  and  $j$  are in the opposite order in  $\tau_1$  and  $\tau_2$  and 0 if they are in the same order.

As in our example, there are often results that are not shared between two lists of search results. Since Kendall tau distance assumes both lists have all elements in common, we must modify it for use on top  $k$  lists. Fagin et al. [15] have developed an equivalence class of distance measurements based on Kendall’s tau and Spearman’s footrule that can be applied to top  $k$  lists that do not necessarily share all elements. Fagin et al. show that all the distance measures they developed are essentially the same, and therefore it is not important which one is used. We use the Kendall tau distance  $K^{(p)}$  with penalty  $p = 0$ , what Fagin et al. call the “optimistic approach”.

Assume  $P$  is the set of all unordered pairs of distinct elements in  $\tau_1 \cup \tau_2$  where  $\tau_1$  and  $\tau_2$  are two top  $k$  lists.  $K^{(0)}$  is defined formally as

$$K^{(0)}(\tau_1, \tau_2) = \sum_{\{i,j\} \in P} \bar{K}_{i,j}(\tau_1, \tau_2) \quad (3)$$

where

- $\bar{K}_{i,j}(\tau_1, \tau_2) = 1$  if any of the following conditions hold:
  1.  $i$  and  $j$  are in both lists and in the opposite order
  2.  $i$  and  $j$  are in the same list with  $j$  ranked ahead of  $i$ , and only  $i$  is in the other list
  3.  $i$  appears only in one list and  $j$  is only in the other
- $\bar{K}_{i,j}(\tau_1, \tau_2) = 0$  for all other conditions

In order to compare  $K^{(0)}$  with our overlap distance measure, we need to normalize  $K^{(0)}$  so two identical lists have a value of 1, and two lists that share no elements have a value of 0. If two top  $k$  lists have no shared elements, then  $K^{(0)}$  would be equal to  $k^2$ . Therefore our normalized version of  $K^{(0)}$ , which we denote  $K$ , becomes:

$$K = 1 - \frac{K^{(0)}(\tau_1, \tau_2)}{k^2} \quad (4)$$

To illustrate how  $K$  can be intuitively computed using  $K_\tau$ , consider again the top 4 search results from the previous example. Since B and C do not appear in list 2, we assume optimistically that they have slid down to positions 5 and 6 in list 2, just beyond the bounds of our top 4 list. Since E and F do not appear in list 1, we again assume that they are at positions 5 and 6 in list 1. So to compare these lists with  $K_\tau$ , we take the unshared elements of list 1 and append them to list 2 and vice versa. The two lists become:

1. ABCDEF
2. EDAFBC

We can now compute  $K_\tau(\tau'_1, \tau'_2)$  on the altered lists since they now have all elements in common. In this case there are 9 pairwise disagreements. Since  $K_\tau(\tau'_1, \tau'_2)$  and  $K^{(0)}(\tau_1, \tau_2)$  are equivalent,  $K$  is  $1 - 9/16 = 0.4375$ .

The third distance measure we will use in our experiment is the  $M$  measure developed by Bar-Ilan et al. [7]. The  $M$  measure is based on findings of an eye-tracking study [13] that showed participants where much more likely to view

the top set of search results on a page than the bottom. The  $M$  measure encodes the intuition that changes in the top positions of the search results should be weighed more heavily than changes at the bottom of the result set. The unnormalized value  $M'$  for two top  $k$  lists  $\tau_1$  and  $\tau_2$  is defined:

$$M'(\tau_1, \tau_2) = \sum_{\tau_1 \cap \tau_2} \left| \frac{1}{rank_1(i)} - \frac{1}{rank_2(i)} \right| + \sum_{\tau_1 \setminus \tau_2} \left( \frac{1}{rank_1(i)} - \frac{1}{k+1} \right) + \sum_{\tau_2 \setminus \tau_1} \left( \frac{1}{rank_2(i)} - \frac{1}{k+1} \right) \quad (5)$$

where  $rank_j(i)$  is the rank of element  $i$  in  $\tau_j$ .

To normalize  $M'$  so the value is 1 when the two lists are identical and 0 when the two lists share no common elements, we compute the normalization factor  $n$  to be 4.03975 for  $k = 10$  and 8.39456 for  $k = 100$ :

$$M = 1 - \frac{M'(\tau_1, \tau_2)}{n} \quad (6)$$

Returning again to our previous example,  $M$  would be  $1 - 2.2/2.56667 = 0.1429$  where  $n = 2.56667$  is the value of  $M'$  when two top 4 lists share no elements. Here we see  $M$  assigning a much lower distance to the two lists than did  $P$  (0.5) or  $K$  (0.4375) because of large changes to the first couple of elements. If A had remained the first element in list 2 (AEDF),  $P$  would remain 0.5, but  $K$  would now be 0.5625, and  $M$  would be 0.6623.

### 3. EXPERIMENT SETUP

We devised an experiment to compare the WUI and API results based on four types of queries:

1. **General search terms** - query for the top 100 results and the total number of results
2. **URL backlinks** - query for the number of backlinks to a particular URL
3. **Total URLs indexed for a website** - query for the number of pages indexed for a website
4. **URL indexing and caching** - query if a particular URL is indexed and/or cached

We obtained 50 popular search terms from Lycos Top 50 [29] and 50 computer science (CS) terms obtained from a list at Wikipedia [48]. We assumed the CS terms would produce more stable results than the popular search terms. For each search term we obtained the first 100 results produced by each search engine. We ignored all sponsored results from the WUI queries (API queries do not give sponsored results). We did not use any advanced syntax in our searches. Care was taken to ensure that both the WUI and API were asked to return the same sets of results (no filtering, similar results or restrictions). Although Google often returns ‘supplemental results’ [47] in their WUI results but not in their API results, there is no way to control for it.

At the outset of our experiment, we randomly selected 100 URLs from the top 100 results generated by all the search terms. We used these URLs to ask 1) if the given

URL was indexed by the search engine (using ‘info:’ for Google and ‘url:’ for MSN and Yahoo), and 2) if there was a cached URL for the given URL. We also asked each search engine to report the number of backlinks it had recorded for each of these URLs using the ‘link:’ query. We used the website domain name of the same 100 randomly obtained URLs to ask each search engine the number of web pages it had indexed for the website. For example, if the URL was `http://foo.org/abc.html`, we would query for ‘site:foo.org’. The same URLs were used throughout the experiment.

We used a single server (beatitude.cs.odu.edu with IP address 128.82.4.22) to automate our data collection activities. Each of our queries were issued serially, first to the WUI and then to the API. The queries were executed beginning at 2 am in the morning (Eastern Daylight Time) when search engine traffic is typically light [40]. We delayed a random 15-60 seconds per query to each WUI to avoid generating too much traffic in a small amount of time and being detected as an automated script. Google and Yahoo have both taken steps in the past few years to deny access to IP addresses where they have detected what they believe to be automated queries [18, 36]. These efforts are likely in response to aggressive querying from the SEO community and from viruses that use search engines to find new targets [16].

We issued a total of 3500 queries each day to the three search engines. The WUI queries were aimed at the US version of each search engine: `www.google.com`, `search.msn.com` and `search.yahoo.com`. We did not target specific datacenters since this would be atypical behavior of general users [12]. All WUI web page responses were archived in case changes to the HTML format broke any of our screen-scraping regular expressions (this happened a couple of times with MSN and Yahoo). The WUI queries produced roughly 50 MB of data (5.5 MB compressed) daily. We started making daily queries in late May 2006.

## 4. EXPERIMENT RESULTS

### 4.1 Query Errors

We encountered numerous transitory errors throughout the experiment. The Yahoo API typically generated a few “Bad Request: service temporarily unavailable” responses each morning. On occasion the error would be received 20-30 times on the same day. Towards the end of the experiment, Google’s API frequently returned 502 “Bad Gateway” errors. These errors became so prevalent that we modified our scripts to delay for 15 seconds and re-submitted our query each time the 502 error was returned. On rare occasions, the WUI of all three search engines would respond with an http 500 response.

Throughout our experiment, the Google API frequently returned back an error response (“Exception from service object: For input string XYZ”) for the queries *list* and *database*. This error was caused on Google’s back end when they attempted to return a total result value like 3,450,000,000 that was too large for the 32-bit integer used by the API.

The most troubling error was incurred on Aug. 29 when MSN invalidated our license key without warning. For a period of 17 days we continued to receive this “Invalid value for AppID in request” error before replacing the key with a new one. We received this error several more times throughout the experiment, but our key appeared to function properly

for the vast majority of queries.

## 4.2 Search Term Results

### 4.2.1 Intra Interface Comparisons

When examining the top  $k$  search results produced each day from popular and CS terms, we used the three distance measures discussed in Section 2.3: overlap  $P$ , Fagin et al.’s Kendall tau for top  $k$  lists  $K$ , and Bar-Ilan et al.’s measure  $M$ . We applied these measures using intra-interface comparisons (WUI vs. WUI results and API vs. API results) and comparisons between interfaces (WUI vs. API).

First we examine how the top 100 WUI and API results change over time. We averaged the distances for all search results, separating the CS and popular term results. Table 1 shows the descriptive statistics when comparing the WUI results on day  $n$  to day  $n - 1$  and the API results on day  $n$  to day  $n - 1$ . In Figure 1 we plot the distance measures for Google, MSN and Yahoo using the  $K$  distance (the API gap for MSN was due to the invalid key error discussed in Section 4.1). Graphs using the other distance measures looked very similar in terms of how the WUI and API distance lines closely followed each other. Excluding days 15-25 for Google, the distance lines for all three search engines appear to move in synchronization, even when very large changes take place. For example, we see a spike on day 50 for Yahoo where the WUI and API results both report less than a 0.7  $K$  distance with day 49 results. We see very similar results when examining the top 10 results, and we refer the reader to [31] that contains detailed findings.

For the most part, the CS and popular term results averaged the same degree of change for both Google and Yahoo. Only MSN reported higher average distance measures for CS results. The averages reported by both interfaces were largely the same; MSN reported the greatest average disagreement between the interfaces using the  $M$  measure for popular term results, a difference of 0.05.

When examining each term’s results individually, we found that several of the popular term results did not have synchronized interface updates in Google. Google’s API results for *carmen electra* and *jessica simpson*, for example, appeared to change independently from the WUI results. Other search terms like *angelina jolie* and *anna nicole smith* showed a high degree of synchronization. All 50 of the CS terms also showed close change synchronization between the Google WUI and API results as did all the Yahoo and MSN results for all terms.

### 4.2.2 WUI vs. API Comparisons

We have seen that the daily results produced by the WUI and API changed at nearly the same rate. Now we examine how similar the WUI and API results are to each other on any given day. Figure 2 plots the distance between the API and WUI results on each day using all three distance measures. Descriptive statistics are given in Table 2. When we examined just the top 10 results, we found the measurements to be nearly the same (less than a change of 0.05) except in the case of Google whose popular and CS  $P$  averages dropped by 0.13 and 0.15, respectively, and whose  $K$  average dropped by 0.6.

The  $M$  measure appears significantly lower than the other measures when examining Google’s results for both popular and CS terms (Table 2). MSN also appears to have lower

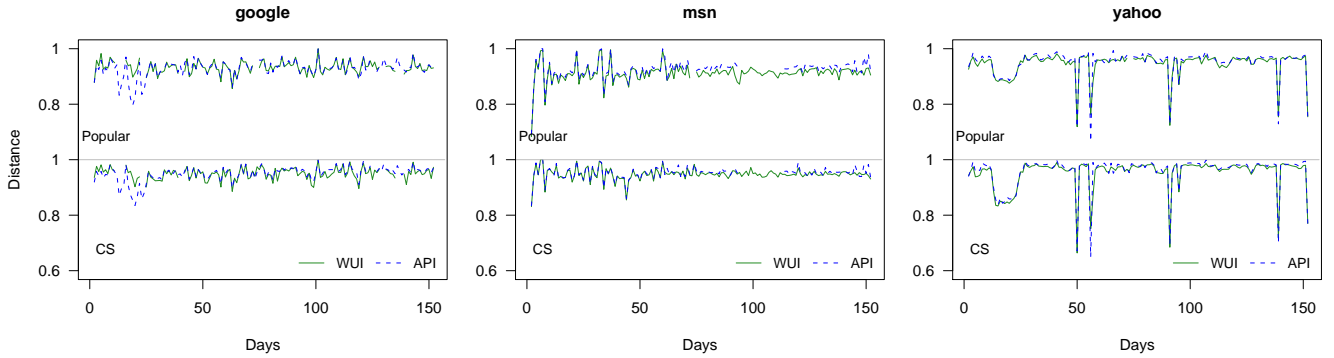


Figure 1:  $K$  distance between top 100 search results when comparing day  $n$  to day  $n - 1$

Table 1: Statistics for Top 100 Search Results Comparing WUI to WUI and API to API

Type	Dist	Interface	Mean	Google		MSN			Yahoo		
				Min	Max	Mean	Min	Max	Mean	Min	Max
Popular	P	WUI	0.92	0.49	1	0.89	0	1	0.93	0.31	1
		API	0.91	0.30	1	0.90	0.19	1	0.94	0.32	1
	K	WUI	0.92	0.53	1	0.91	0.28	1	0.95	0.44	1
		API	0.93	0.39	1	0.93	0.28	1	0.95	0.35	1
	M	WUI	0.94	0.45	1	0.84	0.16	1	0.95	0.34	1
		API	0.93	0.22	1	0.89	0.05	1	0.95	0.34	1
CS	P	WUI	0.94	0.53	1	0.93	0.24	1	0.94	0.29	1
		API	0.93	0.47	1	0.93	0.24	1	0.95	0.31	1
	K	WUI	0.95	0.62	1	0.95	0.31	1	0.95	0.36	1
		API	0.95	0.59	1	0.95	0.32	1	0.96	0.39	1
	M	WUI	0.94	0.34	1	0.93	0.26	1	0.95	0.35	1
		API	0.95	0.43	1	0.94	0.05	1	0.95	0.35	1

$M$  measures for popular terms. This suggests that the top results are the most significantly different between the interfaces for both Google and MSN.

We also see from Table 2 that Google’s overall distance measures are lower for CS terms than popular terms, but the reverse is true for Yahoo. A two-sample Wilcoxon signed rank test confirms the significance at the  $p < 0.001$  level. Although the mean differences between the term type results are also significantly different for MSN ( $p < 0.001$ ), the values are much closer.

In an effort to see if the API results were older (or newer) than the WUI results for any given day, we compared each set of top 100 results on day  $n$  with day  $m$  where  $m$  ranged from  $n - 10$  to  $n + 10$ . The results are graphed in Figure 3. For all three search engines, the highest degree of similarity occurs on the same day (offset 0) which means that the API results are *not* pulling from an older (or newer) index. This is true when examining the top 10 results as well.

The number of times the WUI and API produced identical top 100 results was nearly 0. In Table 3 we show the percentage of times that the WUI and API interfaces produced top  $k$  results that were either identical in rank ( $K = 1$ ) or identical in set membership ( $P = 1$ ). Not one time in 5 months (15,200 total queries to each interface) did the Google API ever produce a single top 100 result that was identical to the Google WUI. MSN and Yahoo did only slightly better (0.2%). When we examine only the top 10, Google’s interfaces only give an identical result set 4% of the time although MSN and Yahoo improve considerably. There is only modest improvement by all of the search engines when we

Table 2: Statistics for Top 100 Search Results Comparing WUI to API

Type	SE	Dist	Mean	Min	Max	
Popular	Google	P	0.83	0.05	1	
		K	0.86	0.07	0.99	
		M	0.77	0.01	0.99	
		P	0.93	0	1	
		K	0.93	0.52	1	
		M	0.87	0.06	1	
	Yahoo	P	0.89	0.27	1	
		K	0.92	0.40	1	
		M	0.88	0.24	1	
	CS	Google	P	0.94	0.52	1
			K	0.93	0.58	0.99
			M	0.86	0.42	0.99
MSN		P	0.95	0.43	1	
		K	0.95	0.65	1	
		M	0.93	0.06	1	
Yahoo		P	0.81	0.35	0.99	
		K	0.84	0.41	0.99	
		M	0.83	0.37	0.99	

Table 3: Identical WUI and API Search Results

	Identical in Rank		Identical in Membership	
	Top 100	Top 10	Top 100	Top 10
Google	0%	4.0%	3.1%	4.2%
MSN	0.2%	38.2%	6.4%	39.7%
Yahoo	0.2%	31.6%	0.2%	33.0%

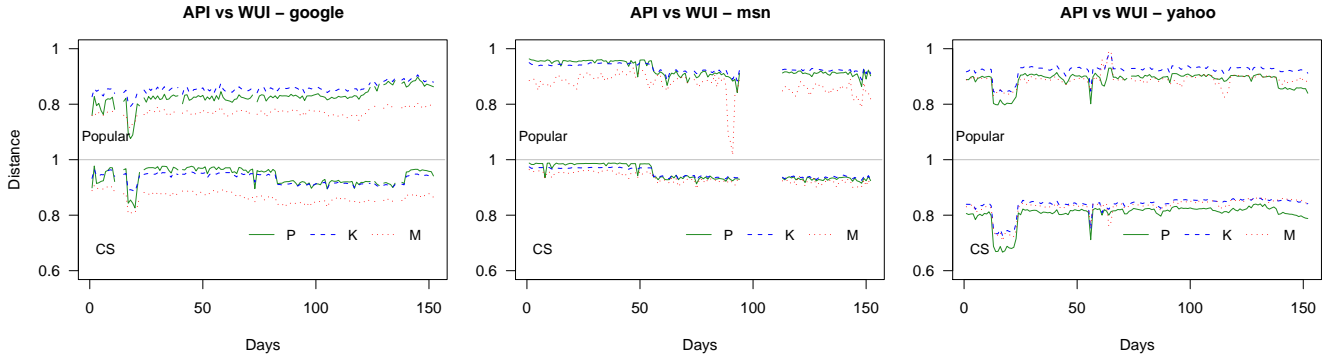


Figure 2: Distance between WUI and API top 100 search results

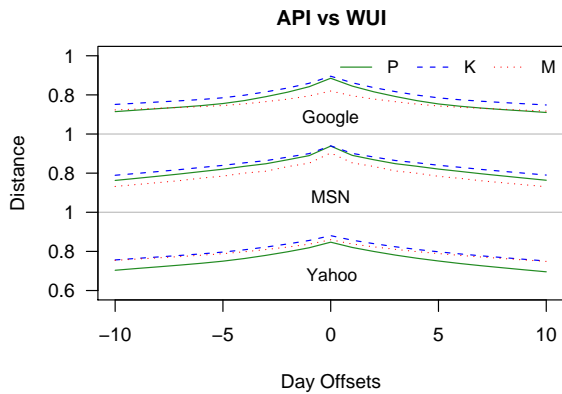


Figure 3: Distance between WUI and API results on day  $n$  compared to day  $n+offset$

disregard ranking. In light of these findings, researchers and developers should expect that the results obtained through any search engine’s API will rarely be identical to what the general WUI user sees.

#### 4.2.3 Decay of Search Results

A number of studies (e.g., [26, 33, 42]) have examined the decay (or linkrot) of the Web over time. In a similar vein, we now examine how the search results for the same query decays over time. Search result decay is a previously unexamined phenomena. Although researchers have noted changes in search results over time, they have not reported change measures on a daily basis over a significant amount of time.

The decay of search results are influenced by a number of factors: 1) decay of the Web (many web pages are here today and gone tomorrow), 2) changes to the content of web pages that make them more or less relevant to a given search term, 3) changes in page popularity (often measured by search engines in terms of inlinks) 4) indexing of new content by search engines, 5) de-duping and de-spamming efforts, and 6) changes to the search engine’s ranking strategy. Although there are many factors influencing decay, we would expect each search engine’s WUI and API to report similar decay rates. We would also expect the results for popular search

terms to decay faster than those of common CS terms since we would expect new Web content in blogs, news websites, commercial sites, etc. to focus on popular subjects of the day.

To compute the decay of search results over time, we used the overlap distance measure  $P$  since it accounts for the presence of items and ignores their rank. We have also computed the **search result half-life** for each search engine and each interface, that is the amount of time it takes for half of the top  $k$  search results obtained on day  $n$  to completely drop from the result set.

There are two methods that may be used to compute the decay of search results over time. The first method is to compare the results on day 1 with every other result obtained in our experiment. While this captures the specific changes for the results obtained on day 1, we are likely to see different changes over time if we had arbitrarily started with the results on day 10 instead, since there could have been a huge discrepancy between results on days 1 and 2 and a very small discrepancy between days 10 and 11. In order to remove this bias and calculate a more generalized set of changes over time, we averaged the offset results over the set of all results obtained in the experiment. For example, to calculate the decay for results that are 20 days apart, we averaged the overlap between days  $(1,21)$ ,  $(2,22)$ , ...,  $(n-20,n)$  where  $n$  is the total number of days in the experiment. More formally, the decay  $D$  for results  $d$  days apart is defined:

$$D(d) = \frac{\sum_{i=1}^{n-d} P(\tau_i, \tau_{i+d})}{n-d} \quad (7)$$

When we examined the top 100 results for decay, each individual term produced results that decayed at different rates for each search engine. The interfaces reported nearly identical decay rates for each search engine except in a few cases like *jessica simpson* for Google and *programming language semantics* for Yahoo. When we examined only the top 10 results, the decay lines became somewhat less synchronized for the interfaces. We also saw some results that “de-decayed” over time. For example, Google’s WUI results for *mother’s day* had decayed by almost 50% two months into the experiment but then rebounded to about 80% by the end of the experiment. On the other hand, Google’s API reported a more linear decay line for *mother’s day*.

**Table 4: Half-Life of Top  $k$  Search Results**

Type	Top	Google		MSN		Yahoo	
		WUI	API	WUI	API	WUI	API
Popular	100	94	71	69	73	47	48
	10	376*	529*	74	96	783*	433*
CS	100	215*	235*	327*	338*	44	48
	10	672*	1480*	228*	264*	67	71

\* Indicates values predicted by our model.

To get a better composite picture of the decay for CS and popular term results, we averaged the decay rates for each category. Figure 4 plots the decay of the top 100 search results over time for each of the three search engines. An **X** marks the half-life for each interface. For top 100 results using CS terms, we did not collect enough data to observe a decay of 0.5 for Google and MSN. Nor did we observe a half-life for top 10 results except MSN using popular terms and Yahoo using CS terms. To predict these missing half-lives, we fitted the decay lines using the linear model:

$$f(\text{day}) = a - b \cdot \log(\text{day}) \quad (8)$$

for which  $a$  ranged from 0.912 to 1.167 (mean of 1.034),  $b$  ranged from 0.161 to 0.393 (mean of 0.253) and R-squared ranged from 0.951 to 0.994 (mean of 0.968). We refer the reader to [31] for specific values. We then calculated the half-life by solving for  $\text{day}$  when  $f(\text{day}) = 0.5$ . The observed and predicted half-lives for each search engine is presented in Table 4. For top 100 results, the predicted half-lives for both interfaces lie close to each other. For top 10 results, the predicted half-lives vary greatly; although we predict Google’s top 10 half-life for CS terms at around 2 years for the WUI, the half-life is approximately 4 years for the API. This large discrepancy can be accounted for by the different decay rates we observed for the top 10 CS results from Google’s interfaces (for  $D(150)$  the WUI reported 0.56, and the API reported 0.62) and because of the flatness of the curve produced by our predictive models over long periods of time.

The CS results appear to be the most stable in Google and MSN. The most unstable results are from Yahoo whose popular and CS top 100 results have half-lives ranging from 44-48. In most cases the stability of the search results increased when examining only the top 10 results (MSN’s CS results are the lone exception). This suggests that it typically takes much longer for a relevant search result to break into the first page of a search engine’s results than it takes to break into the top 100.

### 4.3 Total Results per Search Term

The total results produced by a search term is a very rough estimate of the number of indexed pages containing the term since providing an accurate number is resource-intensive. Large numbers are usually impossible to verify since search engines will only return the first 1000 or so results whether using the API or WUI. Therefore we do not attempt to prove the accuracy of this number, but we compare the numbers given by the interfaces for each search term.

Very rarely did the API and WUI report the exact same total number of results. Google reported the exact same number about 2% of the time and Yahoo about 0.5% of the time. MSN was consistent at providing the exact same

**Table 5: Loose Disagreements (Means)**

		Total results	Total backlinks	Pages indexed
Google	API > WUI	7.9%	0.6%	4.9%
	WUI > API	46.5%	1.5%	46.0%
MSN	API > WUI	0.9%	2.2%	5.4%
	WUI > API	0.6%	21.4%	7.3%
Yahoo	API > WUI	1.0%	24.8%	14.1%
	WUI > API	37.5%	28.1%	8.6%

number (over 99% of the time) until 2006-07-21 when their API and WUI began to report slightly different numbers. Since 2006-07-21, MSN’s interfaces have agreed only 6% of the time.

Since the search engine interfaces frequently disagree, we define a more general sense of disagreement. *Loose disagreement* occurs when the API value is greater than or less than  $\pm 10\%$  of the WUI value. In Table 5 we summarize the percentage of loose disagreements for three of the query types we generated. In Figure 5 we have plotted the percentage of loose disagreements per day between the search engine interfaces. We divide loose disagreements into two types, where the WUI reports higher or lower values than the API. The graphs show that both Google and Yahoo typically report higher values for the WUI than the API. We also see the total percentage of disagreements between Google’s interfaces growing from 40% to 70% on day 96 and then to 90% three weeks later. MSN’s interfaces produce very little disagreement, and Yahoo’s disagree consistently about 40% of the time.

### 4.4 Total Backlinks

In Figure 6 we have plotted the percentage of loose disagreements each day when examining the number of backlinks reported by each interface. Google’s WUI and API typically agreed on every backlink request. For the most part, Google reported the same number of backlinks every-day for all 100 URLs, updating their values only occasionally. MSN and Yahoo reported daily fluctuations. Yahoo’s disagreements held steady at about 35% until they began providing WUI responses from their Yahoo Site Explorer on day 72. Google also reported values that were unusually low when compared to the values reported by MSN and Yahoo. This is not a surprising finding since Google has publicly acknowledged that they do not disclose all backlinks [23].

### 4.5 Total Pages Indexed per Website

Figure 7 plots the percentage of loose disagreements each day when examining the number of pages indexed as reported by each interface. Google appears to have had a major internal change on day 65 which caused more disagreement between the interfaces. Google’s WUI regularly reports much larger values than the API. MSN appears to be the most consistent at returning identical values. Yahoo’s increase in disagreements can again be blamed on the redirection to Yahoo Site Explorer beginning on day 72.

### 4.6 Indexed and Cached URLs

All three search engines appeared to give more consistent responses between their interfaces when asked about the indexed and cached status of URLs. Table 6 shows the percentage of times the WUI responded ‘yes’ to the queries

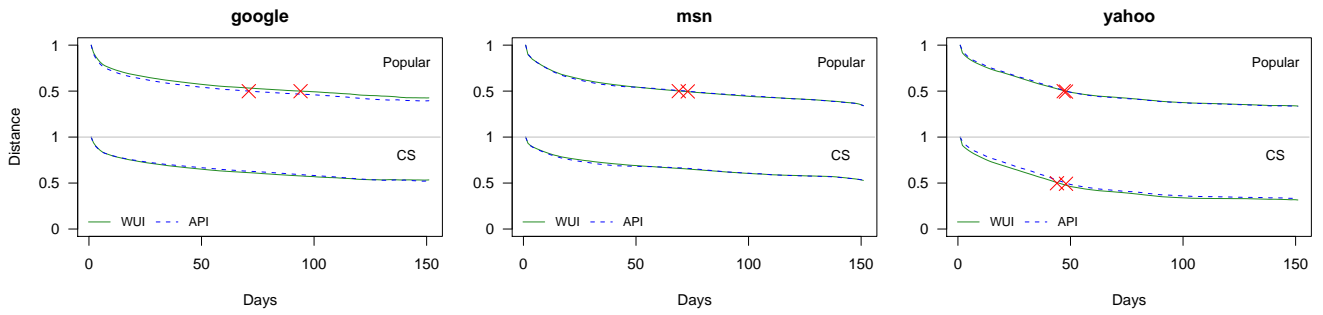


Figure 4: Decay of search results over time (WUI vs. WUI, API vs. API) and half-lives

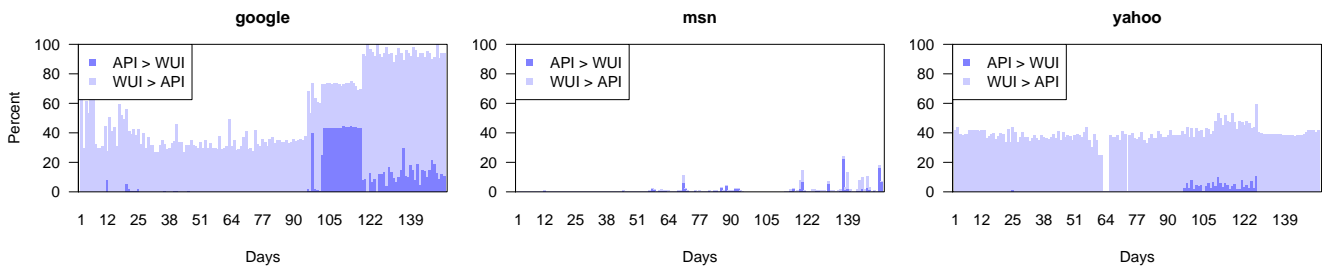


Figure 5: Daily percentage of loose disagreements between interfaces for total search results

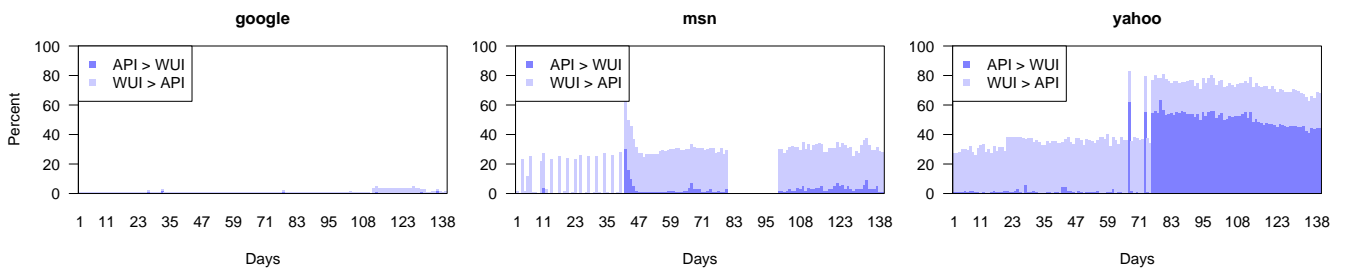


Figure 6: Daily percentage of loose disagreements between interfaces for total backlinks

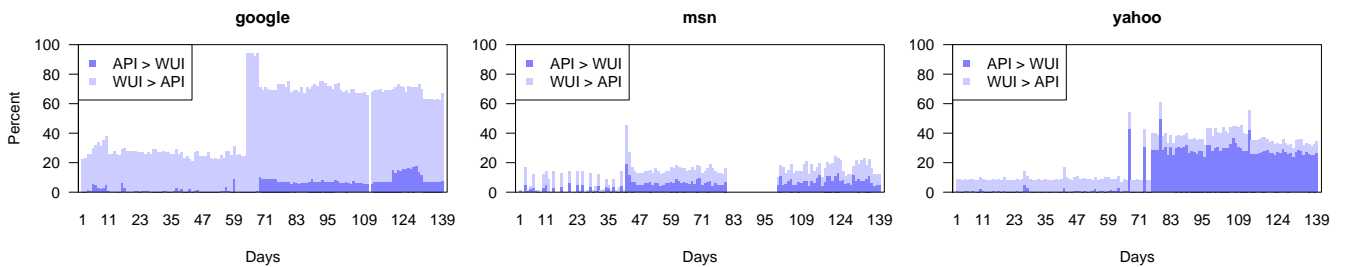


Figure 7: Daily percentage of loose disagreements between interfaces for total indexed pages

**Table 6: Indexed and Cached Status Disagreements**

		Indexed		Cached	
		Disag	Total URLs	Disag	Total URLs
Google	WUI is yes	0.96%	4	1.08%	13
	API is yes	0.04%	6	0.69%	15
MSN	WUI is yes	0.24%	21	0.24%	21
	API is yes	0.91%	23	0.91%	23
Yahoo	WUI is yes	6.30%	28	6.52%	34
	API is yes	0.50%	29	0.49%	27

**Table 7: Synchronized Interfaces**

Type	Most synchronized	Least synchronized
Searching for popular terms	MSN	Google
Searching for CS terms	MSN	Yahoo
Total results	MSN	Google
Total backlinks	Google	Yahoo
Pages indexed per website	MSN	Google
Indexed/cached status	Google/MSN	Yahoo

‘Is URL  $x$  indexed?’ and ‘Is URL  $x$  cached?’. The table also reports the total number of URLs for which a disagreement was found.

Yahoo appears to have the greatest inconsistency between interfaces, averaging 7% disagreements when asked if a URL is indexed and cached. MSN and Google average fewer than 2% disagreements each day. Google only recorded disagreements for a handful of URLs and consistently reported <http://www.koolcelebrities.com/actress/angelina/> as being indexed by the WUI but not the API. MSN and Yahoo had trouble with a broader range of URLs.

## 5. DISCUSSION

Our analysis of search engine results over a five month span has shown that there are significant differences between the reported values of the API and WUI interfaces. We can categorize those search engines which we found to have the least synchronized and most synchronized interfaces by examining the total loose disagreements and average distances between interfaces. In Table 7 we have classified the three search engines across all the query types we have examined. MSN appears to have the most synchronized interfaces overall. Although Yahoo provided the least synchronized results in terms of total backlinks, we remind the reader that Yahoo has designed Site Explorer APIs that would likely match the results from Yahoo’s Site Explorer WUI. We did not test these APIs because we were attempting to measure Yahoo’s standard WUI to their standard API. Our findings show that Yahoo’s standard API is working off an entirely different set of data than the Yahoo Site Explorer WUI.

As noted in Section 4.2.2, a surprising result from our study was that the WUI and API interfaces typically give less similar results depending on the type of search term used. Popular terms revealed significantly higher distances between the interfaces for Google, and CS terms revealed significantly higher distances for Yahoo. Because the WUI and API results appear to be changing in a synchronized fashion, it appears that there are slight differences in ranking algorithms between the two which likely consider web spam differently.

When we ask the question, “Are the APIs serving results from an *older* index?”, the answer is *no*. We have seen that all three search engines provide nearly synchronized changes in their results each day. Although the WUI and API results are consistently less than unity with regards to their distance measures, we found the highest distance measures always occurred on the same day (Figure 3).

When we ask the question, “Are the APIs serving results from a *smaller* index?”, the answer is *probably* for Google and Yahoo. Google and Yahoo’s WUIs consistently report total results that are higher than the APIs, but MSN’s interfaces regularly agree (of course we cannot directly verify these estimates, and Google adds ‘supplemental results’ to their WUI results). The backlink counts from MSN’s and Yahoo’s WUI were also consistently larger than the API counts. And Google’s WUI consistently produced larger website page counts than did their API. To give a more definitive answer to this question, we suggest a future experiment that randomly samples from each corpus and compares the overlap [8] or that uses newer methods proposed by Broder et al. [10].

## 6. CONCLUSIONS

Our five month experiment has uncovered a variety of disagreements between the interfaces of Google, MSN and Yahoo. Our findings suggest that the API indexes are not older, but they are *probably* smaller. Although the indexes used by the WUI and API appear to be updated at the same rate for all three search engines, the top 100 WUI and API results are rarely identical. When examining just the top 10 results, MSN and Yahoo only give identical results about one third of the time. In general, we found MSN to produce the most consistent results between their two interfaces.

We have also examined how search results decay over time. We have built predictive models based on the observed decay rates of the popular and CS search term results used in our experiment. In general, we have found that the decay rates for popular results differ significantly for Google and MSN, and the top 10 results decay at a much slower rate than do the top 100 results. It can take over a year for half of the results to a popular query to be replaced with other results in Google and Yahoo; for MSN it can take only 2-3 months.

We hope that our findings will allow other researchers to better understand the differences between results obtained through the search engine interfaces. Researchers may need to use caution when generalizing their results obtained from the API to those results that the common user sees using the WUI. It is our hope that commercial search engines will make a committed effort to provide more synchronized interfaces for the academic community in the future.

## 7. REFERENCES

- [1] Ask terms of service, 2006. [http://sp.ask.com/en/docs/about/terms\\_of\\_service.shtml](http://sp.ask.com/en/docs/about/terms_of_service.shtml).
- [2] J. Bar-Ilan. Search engine results over time - A case study on search engine stability. *Cybermetrics*, 2/3(1), 1998/99.
- [3] J. Bar-Ilan. Methods for measuring search engine performance over time. *Journal of the American Society for Information Science and Technology*, 53(4):308–319, 2002.
- [4] J. Bar-Ilan. Comparing rankings of search results on the Web. *Information Processing & Management*, 41(6):1511–1519, Dec 2005.
- [5] J. Bar-Ilan. Expectations versus reality - search engine features needed for web research at mid 2005. *Cybermetrics*, 9(1), 2005.

- [6] J. Bar-Ilan, M. Levene, and M. Mat-Hassan. Dynamics of search engine rankings - A case study. In *Proceedings of the 3rd International Workshop on Web Dynamics*, May 2004.
- [7] J. Bar-Ilan, M. Mat-Hassan, and M. Levene. Methods for comparing rankings of search engine results. *Compututer Networks*, 50(10):1448–1463, July 2006.
- [8] Z. Bar-Yossef and M. Gurevich. Random sampling from a search engine’s index. In *Proceedings of WWW ’06*, pages 367–376, 2006.
- [9] K. Bharat and A. Broder. A technique for measuring the relative size and overlap of public web search engines. In *Proceedings of WWW7*, pages 379–388, 1998.
- [10] A. Broder, M. Fontura, V. Josifovski, R. Kumar, R. Motwani, S. Nabar, R. Panigrahy, A. Tomkins, and Y. Xu. Estimating corpus size via queries. In *Proceedings of CIKM 2006*, 2006.
- [11] K. Curran and A. Doherty. Automated broadcast media monitoring using the Google API. In *IEEE Consumer Communications and Networking Conference (CCNC 2006)*, volume 2, pages 1098–1102, 2006.
- [12] M. Cutts. Google datacenters. Video, July 31 2006. <http://video.google.com/videoplay?docid=8726665066825965913>.
- [13] Did-it, Enquiro, and Eyetools uncover search’s Golden Triangle. 2005. <http://www.enquiro.com/eye-tracking-pr.asp>.
- [14] W. Ding and G. Marchionini. A comparative study of web search service performance. In *Proceedings of the ASIS Annual Meeting*, volume 33, pages 136–142, 1996.
- [15] R. Fagin, R. Kumar, and D. Sivakumar. Comparing top  $k$  lists. *SIAM Journal on Discrete Mathematics*, 17(1):134–160, 2003.
- [16] P. Festa. Google worm targets AOL, Yahoo. Dec 28 2004. <http://news.com.com/Google+worm+targets+AOL%2C+Yahoo/2100-7349.3-5504769.html>.
- [17] S. Gauch, G. Wang, and M. Gomez. Profusion: Intelligent fusion from multiple, distributed search engines. *The Journal of Universal Computer Science*, 2(9):637–649, 1996.
- [18] B. Gillette. Google blacklisting researchers? Dec 14 2004. [http://www.emailbattles.com/2005/12/14/virus.aacdehdic\\_ei/](http://www.emailbattles.com/2005/12/14/virus.aacdehdic_ei/).
- [19] Google privacy center: Terms of service, 2006. [http://www.google.com/terms\\_of\\_service.html](http://www.google.com/terms_of_service.html).
- [20] Google SOAP Search API. <http://code.google.com/apis/soapsearch/>.
- [21] Google Web API - Frequently Asked Questions. [http://code.google.com/apis/soapsearch/api\\_faq.html](http://code.google.com/apis/soapsearch/api_faq.html).
- [22] Google Web APIs newsgroup, 2006. <http://groups.google.com/group/google.public.web-apis>.
- [23] GoogleGuy. GoogleGuy’s posts. 2005. <http://www.webmasterworld.com/forum30/29720.htm>.
- [24] A. Gulli and A. Signorini. The indexable web is more than 11.5 billion pages. In *Proceedings of WWW ’05*, pages 902–903, May 2005.
- [25] N. Jain, M. Dahlin, and R. Tewari. Using Bloom filters to refine web search results. In *Proceedings of WebDB 2005*, 2005.
- [26] W. Koehler. A longitudinal study of web pages continued: A consideration of document persistence. *Information Research*, 9(2), 2004.
- [27] M. Koo and H. Skinner. Improving web searches: Case study of quit-smoking web sites for teenagers. *Journal of Medical Internet Research*, 5(4), November 2003.
- [28] S. Lawrence and C. L. Giles. Accessibility of information on the web. *Intelligence*, 11(1):32–39, 2000.
- [29] The Lycos 50, 2006. <http://50.lycos.com/>.
- [30] P. Mayr and F. Tosques. Google Web APIs - an instrument for webometric analyses? In *Proceedings of ISSI 2005*, 2005.
- [31] F. McCown. Comparison of search engine APIs, 2006. [http://www.cs.odu.edu/~fmccown/research/se\\_apis/](http://www.cs.odu.edu/~fmccown/research/se_apis/).
- [32] F. McCown, J. Bollen, and M. L. Nelson. Evaluation of the NSDL and Google for obtaining pedagogical resources. In *Proceedings of ECDL ’05*, pages 344–355, 2005.
- [33] F. McCown, S. Chan, M. L. Nelson, and J. Bollen. The availability and persistence of web references in D-Lib Magazine. In *Proceedings of the 5th International Web Archiving Workshop (IWA’05)*, Sept 2005.
- [34] F. McCown, X. Liu, M. L. Nelson, and M. Zubair. Search engine coverage of the OAI-PMH corpus. *IEEE Internet Computing*, 10(2):66–73, Mar/Apr 2006.
- [35] F. McCown and M. L. Nelson. Evaluation of crawling policies for a web-repository crawler. In *Proceedings of HYPERTEXT ’06*, pages 145–156, 2006.
- [36] M. Moffatt. Yahoo error: Unable to process request at this time – error 999. Feb 14 2005. <http://murraymoffatt.com/software-problem-0011.html>.
- [37] MSN terms of service, 2006. <http://tou.live.com/en-us/default.aspx>.
- [38] MSN Web Search API. <http://msdn.microsoft.com/msn/msnsearch/>.
- [39] G. Pant. Deriving link-context from HTML tag tree. In *Proceedings of DMKD ’03*, pages 49–55, 2003.
- [40] R. Pike, S. Dorward, R. Griesemer, and S. Quinlan. Interpreting the data: Parallel analysis with sawzall. *Dynamic Grids and Worldwide Computing*, 13(4):277–298, November 2005.
- [41] C. Snelson. Sampling the Web: The development of a custom search tool for research. *Library and Information Science Research Electronic Journal*, 16(1), Dec 2005.
- [42] D. Spinellis. The decay and failures of web references. *Commun. ACM*, 46(1):71–77, 2003.
- [43] A. Spink, B. J. Jansen, C. Blakely, and S. Koshman. A study of results overlap and uniqueness among major web search engines. *Information Processing & Management*, 42(5):1379–1391, September 2006.
- [44] K. C. Sua, S. E. Waldren, and T. B. Patrick. Differences in the effects of filters on health information retrieval from the internet in three languages from three countries: A comparative study. In *Proceedings of MEDINFO 2004*, 2004.
- [45] M. Thelwall. Can the Web give useful information about commercial uses of scientific research? *Online Information Review*, 28:120–130, 2004.
- [46] L. Vaughan. New measurements for search engine evaluation proposed and tested. *Information Processing & Management*, 40(4):677–691, May 2004.
- [47] What’s a “supplemental result?”. *Webmaster Help Center*, 2006. <http://www.google.com/support/webmasters/bin/answer.py?answer=34473&topic=8523>.
- [48] Wikipedia: List of basic computer science topics, 2006. [http://en.wikipedia.org/wiki/List\\_of\\_basic\\_computer\\_science\\_topics](http://en.wikipedia.org/wiki/List_of_basic_computer_science_topics).
- [49] Windows Live Search: Development (forum for search API web services), 2006. <http://forums.microsoft.com/MSDN/ShowForum.aspx?ForumID=111&SiteID=1>.
- [50] Yahoo terms of service, 2006. <http://docs.yahoo.com/info/terms/>.
- [51] Yahoo! Web Search developer community discussion list, 2006. <http://finance.groups.yahoo.com/group/yws-search-web/>.
- [52] Yahoo! Web Search APIs. <http://developer.yahoo.net/search/web/>.