



SEARCH ENGINES AND THEIR PUBLIC INTERFACES: WHICH APIs ARE THE MOST SYNCHRONIZED?



FRANK MCCOWN AND MICHAEL L. NELSON

DEPARTMENT OF COMPUTER SCIENCE, OLD DOMINION UNIVERSITY, NORFOLK, VIRGINIA, UNITED STATES

RESEARCHERS: SCREEN SCRAPE OR USE THE APIs?

WEB USER INTERFACE (WUI)

This screenshot shows the Google search results page for the query "search engines api". The results include links to various developer documentation and forums related to Google's API services.

THIS IS THE 3RD RESULT OUT OF ABOUT 24,100,000.

APPLICATION PROGRAMMING INT. (API)

This screenshot shows the Google Web APIs developer documentation page. It provides instructions for developing applications using Google's web APIs, including how to create a Google Account and use the Google API Client Library.

NO, IT'S THE 10TH RESULT OUT OF ABOUT 16,300,000!

This screenshot shows the Microsoft Live Search results page for the query "www2007 index size". The results mention the International Web Index Conference and its index size.

THERE ARE 2,911 PAGES INDEXED.

This screenshot shows the Microsoft Live Search results page for the query "www2007 index size". The results mention the International Web Index Conference and its index size.

I SEE ONLY 1,740.

This screenshot shows the Yahoo! search results page for the query "www2007 index size". The results mention the International Web Index Conference and its index size.

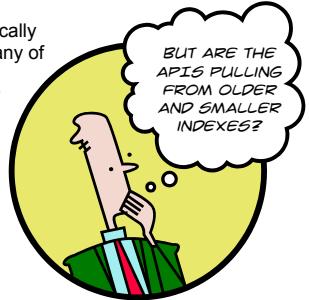
THE URL IS INDEXED AND CACHED.

This screenshot shows the Yahoo! developer network search results page for the query "www2007 index size". The results mention the International Web Index Conference and its index size.

IT'S MISSING ENTIRELY FROM MY INDEX.

Google Terms of Service: "You specifically agree not to access (or attempt to access) any of the Services through any automated means (including use of scripts or web crawlers...)"

Windows Live Terms of Service: "You may not... use any automated process or service to access and/or use the service (such as a BOT, a spider, ...)"



5 MONTH EXPERIMENT

Late May to Oct 2006:

1. **General search terms.** Queried for the top 100 results and total results using 50 popular search terms and 50 computer science (CS) terms.
2. **URL backlinks.** Queried for the number of backlinks (inlinks) to 100 randomly selected URLs.
3. **Pages indexed for a website.** Asked how many pages were indexed for 100 randomly selected websites.
4. **URL indexing and caching.** Queried to see if 100 randomly selected URLs were indexed and/or cached.

COMPARING SEARCH RESULTS

1. Overlap (P)
2. Kendall tau for top k results (K)¹
3. Penalize changes at the top more heavily (M)²

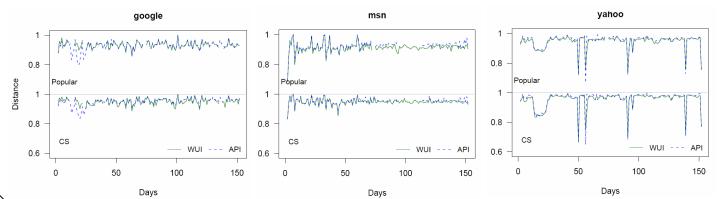
0 1
More similar

¹R. Fagin, R. Kumar, and D. Sivakumar. Comparing top k lists. *SIAM Journal on Discrete Mathematics*, 17(1):134–160, 2003.
²J. Bar-Ilan, M. Mat-Hassan, and M. Levene. Methods for comparing rankings of search engine results. *Computer Networks*, 50(10):1448–1463, July 2006.

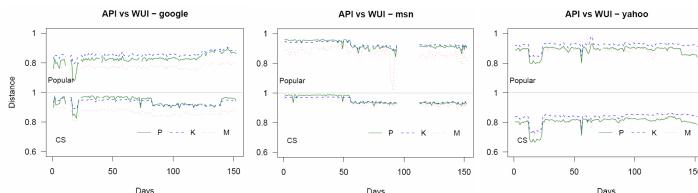
1. ABCD
 2. EDAF
- Examples**

$$P = 0.50 \quad K = 0.44 \quad M = 0.14 \quad P = 0.50 \quad K = 0.56 \quad M = 0.66$$

COMPARING WUI TO WUI & API TO API ON SUCCESSIVE DAYS

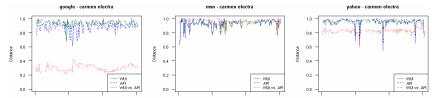


COMPARING WUI TO API



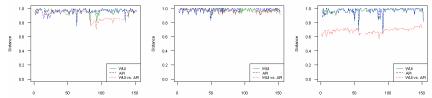
FOR ALL 3 SEARCH ENGINES, THE WUI & API ARE MOST SYNCHRONIZED ON THE SAME DAY.

EXAMPLES



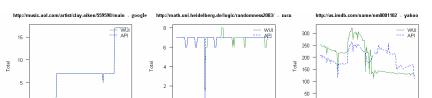
MSN IS MOSTLY SYNCHRONIZED FOR "ALGORITHM".

GOOGLE IS LESS SYNCHRONIZED FOR POPULAR TERMS.



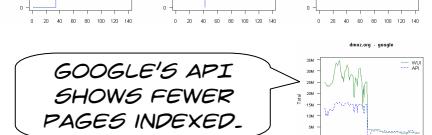
YAHOO IS LESS SYNCHRONIZED FOR CS TERMS.

HOW MANY TOTAL RESULTS ARE THERE?

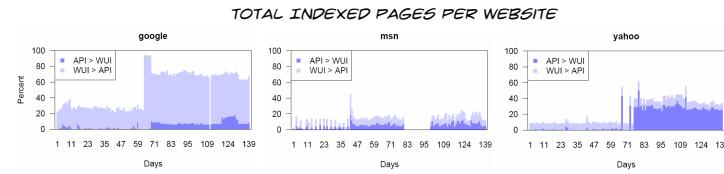
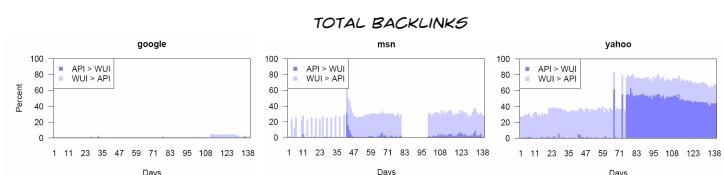
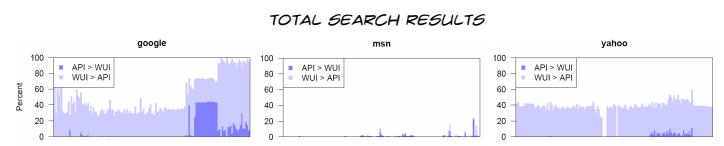


WHOSE BACKLINK COUNTS ARE CORRECT?

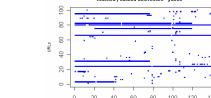
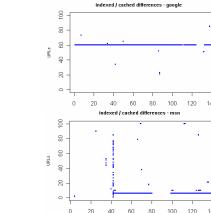
GOOGLE'S API SHOWS FEWER PAGES INDEXED.



LOOSE DISAGREEMENTS



INDEXED / CACHED DISAGREEMENTS



GOOGLE & YAHOO MIGHT BE PULLING FROM SMALLER INDEXES.

Table 1: Loose Disagreements (Means)

	Total results	Total backlinks	Pages indexed
Google	API > WUI (7.0%)	WUI > API (46.5%)	0.6% (1.5%)
MSN	API > WUI (0.9%)	WUI > API (0.6%)	2.2% (21.4%)
Yahoo	API > WUI (1.0%)	WUI > API (37.5%)	24.8% (28.1%)

YAHOO SEEMS TO BE CONFUSED.

KAB POW!!!

Table 2: Synchronized Interfaces

Type	Most synched	Least synched
Search for popular terms	MSN	Google
Search 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

SEE ALSO

Frank McConell and Michael L. Nelson. Agreeing to Disagree: Search Engines and their Public Interfaces. ACM IEEE Joint Conference on Digital Libraries (JCDL 2007). To appear.

All graphs: http://www.cs.odu.edu/~fmccown/research/se_apis/
Complete data set available upon request.



Lazy Preservation: Reconstructing Websites from the Web Infrastructure
<http://www.cs.odu.edu/~fmccown/warrick/>



Other research projects at Old Dominion University:

mod_oai: An Apache Module for Efficient, Automatic Web Harvesting
<http://www.modoai.org/>

