

1 ALBERT GIDARI, JR., *pro hac vice*
 (AGidari@perkinscoie.com)
 2 PERKINS COIE LLP
 1201 Third Avenue, Suite 4800
 3 Seattle, WA 98101
 Telephone: (206) 359-8000
 4 Facsimile: (206) 359-9000

5 LISA A. DELEHUNT, Bar No. 228551
 (LDelehunt@perkinscoie.com)
 6 PERKINS COIE LLP
 180 Townsend Street, 3rd Floor
 7 San Francisco, California 94107-1909
 Telephone: (415) 344-7000
 8 Facsimile: (415) 344-7050

9 Attorneys for Respondent
 GOOGLE INC.

11 **UNITED STATES DISTRICT COURT**
 12 **NORTHERN DISTRICT OF CALIFORNIA**
 13 **SAN JOSE DIVISION**

15 ALBERTO R. GONZALEZ, in his official
 capacity as ATTORNEY GENERAL OF THE
 16 UNITED STATES,
 17 **Movant,**
 18 **v.**
 19 **GOOGLE INC.,**
 20 **Respondent.**

CASE NO. 5:06-mc-80006-JW
 DECLARATION OF MATT CUTTS

22 MATT CUTTS states as follows:
 23
 24
 25
 26
 27
 28

I. INTRODUCTION

1
2 1. I am a Senior Staff Software Engineer and head of the Webspam group at Google.
3 I have held this position since November 2004. I have worked with Google since February 2000
4 as a technical engineer in one capacity or another. I have personal knowledge of the facts set forth
5 below and, if called as a witness about those facts, could testify competently thereto.

6 2. Currently, I manage the Webspam group. "Webspam" does not refer to spam
7 emails, but rather to pages from the World Wide Web ("Web") that violate Google's quality
8 guidelines in an attempt to rank higher on Google. For example, a page containing pornographic
9 material that attempted to show up as a result to a search for "Disney cartoons" would be
10 considered webspam. Accordingly, I am knowledgeable about how Google crawls, indexes, and
11 ranks web pages. I have written a variety of programs for Google, including programs having to
12 do with filtering pornography and detecting web pages that attempt to bypass Google's quality
13 guidelines. I have a Master of Science in Computer Science, *summa cum laude*, from University
14 of North Carolina at Chapel Hill. I am recipient of a National Science Foundation Graduate
15 Research Fellowship, a Link Foundation Fellowship, and a Gaines Fellowship in the Humanities.
16 I also have a Bachelor of Science in Computer Science and Mathematics, *summa cum laude*, from
17 the University of Kentucky. I have co-authored and published a variety of articles relating to
18 computer graphics and computer vision.

19 3. I have reviewed the Motion to Compel Compliance with Subpoena Duces Tecum
20 brought by Alberto R. Gonzales (the "Government") against Google, Inc. and the supporting
21 declarations and exhibits, including the Declaration of Professor Philip B. Stark, Ph.D. The
22 Government has requested that Google produce a multistage random sample of one million URLs
23 from Google's databases. URL stands for "Uniform Resource Locator" and acts as the global
24 address of documents and other resources on the World Wide Web ("Web"). URLs consist of the
25 protocol (e.g. HTTP) and the Internet Protocol ("IP") address or domain name. The Government
26 has suggested that Google select at random 100 of its data centers containing URLs, and then
27

1 select at random 10,000 URLs from each of those data centers. In his declaration, Professor Stark
2 has stated that he could be involved in the random selection process.

3 4. Additionally, the Government has requested that Google produce every single
4 search query entered at www.google.com over a given week.

5 **II. SUMMARY OF DECLARATION**

6 5. In my declaration, I testify as to the following:

7 A. In Section III, I summarize how Google's search engine works.

8 B. In Section IV, I explain how examining the search queries Google's users
9 have typed into Google's search engine at www.google.com does not reveal accurate
10 information about the likelihood of minors' exposure to material the Government deems
11 Harmful to Minors ("HTM") or about the search patterns of users.

12 C. In Section V, I explain that there is little to no correlation between a URL's
13 presence in Google's index and if or when that URL would be returned as a result
14 responsive to a specific set of query words.

15 D. In Section VI, I discuss the fact that, in order to protect the privacy of its
16 users, Google does not disclose to third parties the raw records of searches entered into its
17 search engine.

18 E. In Section VII, I explain that Google does not disclose the queries it
19 receives from its users because it offers information about the nature of Google's users and
20 of Google's market share in both the United States and other countries and languages,
21 which could allow Google's competitors to compete more effectively with Google.
22 Additionally, Google does not disclose its index of URLs or any such large samples from
23 its index, because one could glean from a sample of the URLs from Google's index
24 confidential and proprietary information about how Google's search engine works.

25 F. In Section VIII, I explain those measures of which I am personally aware
26 that Google undertakes to protect the confidentiality of its proprietary information.

1 G. In Section IX, I discuss reasonable, viable alternative public sources of the
2 information the Government is seeking from Google.

3 H. In Section X, I explain the burden that would be imposed on Google if it
4 were required to respond to the Government's subpoena.

5 **III. GOOGLE'S SEARCH ENGINE**

6 6. Google provides the world's most-used search engine at www.google.com. In
7 general terms, search engines allow users of the Web to search a considerable portion of the vast
8 amount of material located on the Web by entering text queries into the engine. The Web is
9 composed of billions of publicly accessible web sites from around the world and other information
10 sources that web browsers can access. Search engines offer this search capability by creating an
11 index of certain Web content through the use of programs that "crawl" the Web and automatically
12 fetch Web pages. The search engine then allows users to search that index of certain Web content
13 through an interface such as that found at www.google.com. When a user enters a query, the
14 engine searches its Web index and retrieves results relevant to the query. Google's search engine
15 is open to the public. Google is the leader among the major search engines (including Yahoo!,
16 Ask Jeeves, and MSN Search) to a large degree because of its sophisticated and proprietary
17 technology that returns the most relevant and useful results in response to user queries.

18 7. Some of Google's services, like Google News, allow users to store or establish
19 repeat search queries. Google users may set up recurring searches with results sent to email
20 accounts at defined intervals. Google then processes these requests for its users and sends the
21 results to the email accounts, wireless phones, or other designated mobile devices.

22 **IV. SEARCH QUERIES**

23 8. Examining the search queries Google's users have typed into Google's search
24 engine at www.google.com does not reveal accurate information about how likely it is that minors
25 might be exposed to material the Government deems Harmful to Minors ("HTM") or about the
26 search patterns of users.

1 programs that detect whether a computer is connected to the Web by sending periodic queries to
2 Google; such programs can send thousands of queries to Google each day. Google also receives a
3 fraction of its queries from malicious programs, such as the Santy worm.

4 14. In addition to queries sent by "bots," test programs, and worms, an individual web
5 page owner (also known as a "Webmaster") may send dozens of queries by hand to check on how
6 his or her websites rank in Google. Many website owners check their rankings every day, which
7 can cause further skew in the query log. I am also aware of efforts by some Google users to
8 deliberately send pornography queries to Google in reaction to the Government's subpoena. One
9 individual wrote a feature for the popular Firefox web browser that will send a random
10 pornography query to Google whenever a user performs a normal query, as if the pornography
11 query had also been entered by the user. Attached as Exhibit A is a true and correct copy of this
12 blog entry which can be found at <http://www.hughes-family.org/wordpress/2006/01/23/help-the-justice-department-out-via-greasemonkey/>, and which was printed on February 16, 2006.

14 15. Without removing these artificial and automatic queries, one week of raw query
15 logs will be skewed beyond usability for many purposes. It would be difficult—perhaps
16 impossible—for most researchers to distinguish between such artificial queries and real queries
17 without considerable effort and the use of proprietary techniques.

18 16. Google's proprietary techniques for returning search results are not static, and
19 Google's algorithms change regularly. Thus, the identical search query run in Google's search
20 engine today is likely to yield different search results than an identical search conducted yesterday,
21 last week, or last month.

22 V. GOOGLE'S INDEX OF URLS

23 17. The presence of a URL in Google's index is not in any way representative of the
24 frequency that the URL will be shown to a user. There is little to no correlation between a URL's
25 presence in the index Google maintains of copies of documents it has collected from crawling the
26 web and if or when that URL would be returned as a result responsive to a specific set of query
27 words. Using its proprietary and confidential technology, Google employs more than 100 factors

1 in scoring documents for relevancy besides just the document's URL. Some of the factors are
2 straightforward. For example, Google considers whether query words are present in the title of the
3 document, how many times the query words appear in the document, the proximity of the query
4 words to each other, as well as a document's PageRank, which is Google's patented method of
5 measuring the reputation of a page. Most of Google's scoring factors and how they are combined
6 is confidential.

7 18. Beyond scoring factors, some documents in Google's index will have additional
8 demotions or will be blocked for webspam or legal reasons. For example, Google's index contains
9 documents for which we have received valid complaints under the Digital Millennium Copyright
10 Act, and Google does not return those documents in our search results.

11 19. For all these reasons, it does not follow that a sample of URLs from Google's index
12 will indicate how often searches will return HTM material. Simply put, the fact that a URL in
13 Google's index may appear to contain content the Government considers HTM is not
14 representative of whether and how often HTM content appears as a result of a query. For
15 example, if 5 percent of the URLs listed in Google's index contain HTM content, that fact alone
16 does not mean that any given search would yield 5 percent HTM content.

17 20. The adage "you can't judge a book by its cover" applies on the Web: URLs alone
18 do not indicate whether content at that URL is HTM. A URL may or may not be logically
19 connected to its content. Sites may have names that suggest explicit sexual material, but actually
20 do not contain HTM content. A URL such as
21 <http://www.pbs.org/wgbh/pages/prontline/shows/porn/etc/links.html> contains the word "porn" but
22 provides links to anti-pornography organizations such as the American Family Association and the
23 Family Research Council. Likewise, a URL such as <http://www.porn-free.org/> may appear
24 pornographic from the URL, but is a faith-based site that is against pornography. The reverse is
25 also true: URLs may have name that seem innocent but actually contain material the Government
26 might consider HTM. The classic example was that www.whitehouse.gov was a harmless website
27 about the White House, while, until recently, www.whitehouse.com contained pornographic

1 material. Furthermore, the content of web pages at URLs is fluid and dynamic, and the actual
2 content can be difficult to ascertain. For example, until recently the URL
3 <http://www.crisiscentersyr.org> was evidently owned by a rape crisis center. It appears that the
4 domain expired in late 2005 and was instead registered by someone else. The new owner kept the
5 appearance of the site, making it appear it still belongs to a rape crisis center, but added links to
6 pornography at the bottom of the document. In September 2005, that URL contained non-
7 pornographic material, while in January 2006, it contained links to pornographic material. In
8 short, a URL alone is insufficient to determine whether a document is harmful to minors. This
9 example also shows the dynamic nature of Web content. The content of a web page can change at
10 anytime, and some Web pages update their content daily or more frequently.

11 21. Additionally, in an attempt to rank higher in Google's search results, some
12 Webmasters show different results to Google than they do to individual users. That is, Google
13 uses a program to crawl the Web to find Web pages for its index. Google's program comes from a
14 specific set of IP addresses and identifies itself as "Googlebot." Some Webmasters program their
15 Web pages to show certain content to Google, content that may be innocent and may seem
16 relevant to a search query. The Webmaster, however, shows different content to individual users,
17 and the content may be pornographic or material otherwise prohibited by Google's guidelines.
18 Such "bait and switch" tactics are called "cloaking." Ultimately, through its diligence, Google
19 finds and removes those documents, but the nature of these Web pages is another example of how
20 a URL does not indicate what content will be displayed in response to a search query.

21 22. The presence of URLs in Google's index is not reflective of the entire Web.
22 Google uses its technology to crawl and index only the best documents from the Web. The Web
23 can be viewed as having an infinite number of pages. For example, a single web server running a
24 calendar application could generate dynamic web pages with dates going forward forever.
25 Therefore, Google's mission must be to retrieve and index the most useful pages that it can from
26 the infinite number of potential pages that can be retrieved.

1 23. The only way to get a current and accurate snapshot of the search results that would
2 be returned by a query into Google's search engine would be to run the query on Google's search
3 engine. Notably, running millions of queries on Google above and beyond the normal use of
4 Google would burden Google's system and possibly shut it down. Alternatively, one would need
5 to understand how Google crawls the web, collects URLs, sorts them, indexes them, ranks them,
6 and returns them as search results—in other words, exactly how Google's crown-jewel trade
7 secrets function.

8 **VI. SEARCH QUERIES AND PERSONALLY IDENTIFYING INFORMATION**

9 24. Google does not publicly disclose the searches queries entered into its search
10 engine. If users believe that the text of their search queries into Google's search engine could
11 become public knowledge, they may be less likely to use the search engine for fear of disclosure
12 of their sensitive or private searches for information or websites.

13 25. There are ways in which a search query alone may reveal personally identifying
14 information. For example, many internet users have experienced the mistake of trying to copy-
15 and-paste text into the search query box, only to find that they have pasted something that they did
16 not intend. Because Google allows very long queries, it is possible that a user may paste a
17 fragment of an email or a document that would tie the query to a specific person. Users could also
18 enter information such as a credit card, a social security number, an unlisted phone number or
19 some other information that can only be tied to one person. Some people search for their credit
20 card or social security number deliberately in order to check for identity theft or to see if any of
21 their personal information is findable on the Web.

22 **VII. GOOGLE'S COMPETITIVE AND CONFIDENTIAL INFORMATION**

23 26. Another reason that Google does not disclose the queries it receives from its users
24 is because it offers information about the nature and effectiveness of Google's search engine and
25 Google's market share in the United States and other countries, which could allow Google's
26 competitors to compete more effectively with Google. Disclosing an entire week's worth of
27 queries would give an estimate of the number of queries that Google processes, which could be
28

1 used to deduce market share among search engines. Google protects this information not only by
2 not disclosing search queries but by not disclosing even the amount of computers Google
3 maintains to run its search engine. A week's worth of queries would also indicate the percentage
4 of queries Google tends to receive in each language, which would allow competitors to estimate
5 Google's relative market share in a given language or country, to learn opportunities for market
6 growth, and to decide where to allocate resources for each language or country. The data would
7 also indicate the average length of query that Google's search engine typically receives and how
8 Google's users search, such as what percentage of queries use special search operators or
9 punctuation. This information could allow competitors to better understand the type of users who
10 seek out Google or to allocate resources on specific search operators. Similarly, the queries could
11 relate information regarding whether users tend to use Google for navigational help, for research,
12 or for shopping.

13 27. Google also does not disclose its index or such a large sample from its index. One
14 could estimate from a sample of the URLs from Google's index information such as (i) the size of
15 Google's index; (ii) how "deeply" Google crawls in different countries or languages (i.e., how
16 many URLs are crawled from each website on average); and (iii) the ability of Google's crawl
17 metrics to measure the reputation of pages or domains. For example, the amount of crawling
18 between different top-level domains such as .com and .uk compared to .pl and .jp would disclose
19 much about how Google's crawling works. The depth of the crawl, the languages of the URLs
20 crawled, and the number of distinct sites crawled could all reveal confidential information about
21 Google's technology. As described above, Google uses multiple methods for crawling the Web,
22 collecting URLs, indexing them, ranking them, and providing relevant results to Google users.
23 Google developed these methods over a number of years and at considerable expense. As far as I
24 know, these methods are not known by Google's competitors. These methods are critical to
25 Google's success as the world's leading search engine and popularity as the world's most-used
26 search engine.

1 28. While Professor Stark has not disclosed what he intends to do with Google's data or
2 how he intends to do it, it seems logical that Professor Stark will need personal knowledge of how
3 the URLs are collected and maintained. Given my experience in this area, I believe that to obtain
4 that level of knowledge, Professor Stark will need to understand how Google crawls the Web and
5 indexes information. In addition, anyone wishing to attack the randomness of the sample,
6 Professor Stark's methodology (which he has not disclosed), or his ultimate opinions might want
7 to probe into everything that goes into maintaining Google's URL database. For example, it might
8 be important to defend or attack Professor Stark's methodology to know what percentage of
9 queries are initiated within the United States, what percentage of the Web is found on Google's
10 URL index, and how a URL is returned as a result to a search query.

11 **VIII. GOOGLE'S MEASURES TO MAINTAIN THE CONFIDENTIALITY**
12 **OF ITS INFORMATION**

13 29. I am personally aware of some of the measures Google takes to maintain the
14 confidentiality of its query log, index, and proprietary crawling, indexing, and retrieval
15 technology. The methods of which I am aware are described below. Additional methods are
16 described in the Declaration of Marty Lev.

17 30. The Google buildings that I have visited all have secured doors and entrances; I
18 have no reason to think that any Google building that houses the information I describe above
19 would be unsecured. Employees and visitors are required to carry appropriate security badges that
20 security personnel routinely check. Certain entrances are restricted by access devices. There are
21 also security guards on duty in the building.

22 31. Google's information is compartmentalized so that only employees with a need to
23 know have access to certain information. Thus, the information that a sales employee can access
24 is different from the information an engineer can access. For example, sales employees do not
25 have access to Google's source code or to certain parts of Google's technical information. Access
26 to Google's computers is controlled by employee logins and passwords, and employee logins
27 allow them to access only the information they need and blocks them from sensitive information

1 they do not need. Interns, independent contractors, and temporary employees require special
2 clearance before they may access parts of Google's technical information.

3 32. Access to Google's index is restricted to engineers. I believe that most employees
4 do not even know the specifics of where the index is stored. The index cannot be viewed in
5 decipherable form unless one uses a source code compiled with specific Google libraries.
6 Additionally, Google is implementing a system where its most sensitive data—such as
7 documentation reflecting the methods it uses to measure the reputation of web pages—will be
8 taken off line and kept in hard copy in a locked location, and an employee wanting access will
9 have to request such access and sign out the file.

10 33. Likewise, access to Google's query log is particularly restricted to a much smaller
11 group of employees with special clearance who need access to perform their job duties. These
12 employees must reaffirm their need for access to the query log periodically, and anyone who does
13 not reaffirm that need loses their access automatically.

14 34. If a Google employee does not have the correct security clearance for the query log,
15 he or she must file a request to get it, explaining in writing his or her need to know and promising
16 to keep the information confidential. I am aware that any abuse of Google's internal security and
17 privacy policies is grounds for immediate termination.

18 35. I am aware that Google has refused requests for search queries in the past. For
19 example, a professor by the name of Amanda Spink at the University of Pittsburgh, requested
20 access to query data from Google and her request was denied. Additionally, Google generally
21 does not share its index.

22 IX. ALTERNATIVE SOURCES

23 36. Both query data and web documents are available from several other sources.
24 Metasearch engine Dogpile provides a service called "SearchSpy" to see queries done on that
25 engine at <http://www.dogpile.com/info.dogpl/searchspy/>. Metasearch engine MetaCrawler
26 provides a service called "Metaspy" which provides a similar list of queries at
27 <http://www.metacrawler.com/info.metac/searchspy>. The search engines Infospace and

1 WebCrawler also provide similar information. For under \$250 per year, Wordtracker.com sells a
2 database of over 330 million queries which is updated on a weekly basis. Ask Jeeves recently
3 offered a service called "Ask Jeeves Take a Peek" at <http://www.ask.com/docs/peek/> which
4 showed questions being typed at Ask Jeeves, and the page would refresh automatically twice a
5 minute. According to her website, Prof. Spink successfully obtained query data from Excite,
6 AltaVista, Ask Jeeves, Fast/AllTheWeb, Vivisimo, Dogpile, Metacrawler, Webcrawler, and
7 Infospace. Prof. Spink has also written a book discussing trends and characteristics of user
8 interaction with search engines.

9 37. An index of web documents is also available from other sources. One approach
10 would be to crawl the public web, and several research groups have done that. In addition, an
11 Amazon.com subsidiary, Alexa.com ("Alexa"), recently released a service that allows anyone to
12 access Alexa crawl data. The system allows users to process over 4 billion URLs and over 1.7
13 billion full-text documents. The system is specifically intended for "[r]esearchers who wish to
14 tackle problems related to Web content," according to
15 http://pages.alexa.com/awsp/docs/WebHelp/Introduction/Who_Should_Use_the_Platform.htm.
16 Alexa.com also provides tools to search, process, and publish one's own custom subset of data. A
17 researcher or developer could use this system to test code, including pornography filters, over the
18 full text of documents (not just URLs) and to test code with much more than one million URLs.
19 Attached as Exhibit B is a copy of a web log entry about Alexa which can be found at
20 <http://battellemedia.com/archives/002116.php>, and which was printed from the internet on
21 February 16, 2006.

22 38. I believe that Alexa offers Professor Stark a reasonable and viable option to
23 accomplish his goals as set forth in his declaration. I believe that Alexa data would allow for more
24 thorough testing than a sampled list of URLs from Google, as, according to Alexa, their system
25 provides over 300 terabytes of data for researching and testing.

X. THE DEMAND ON GOOGLE TO COMPILE THE INFORMATION REQUESTED

1 39. Google does not index URLs in the form the Government has requested. Rather,
2 Google maintains an index of copies of the billions of web pages that it has crawled. The index
3 requires vast amounts of computer space and is maintained on multiple computers. There is no
4 existing method to simply copy and hand over in electronic or paper form a selection of URLs
5 from Google's index. The type of information sought by the Government is not created or used
6 internally.
7

8 40. Accordingly, an engineer would need to write a computer program capable of
9 gathering the URLs that the Government seeks. The engineer would need to develop the code by
10 determining a scheme to map the data logs and pull random URLs from it. The code would then
11 need to be implemented, which involves a potentially lengthy debugging process in which the
12 engineer tests portions of the code (and eventually the entire program) to identify and fix
13 problems. It is common for debugging to take longer than the initial coding. Finally, the selected
14 documents must be extracted somehow and copied in a form that can be provided to the
15 Government. I estimate that this will take two to five full-time days of an engineer's time. This
16 procedure could be further complicated if Professor Stark needs to participate to verify the
17 randomness of the sample. Google does not dedicate any engineers for this type of task, and so an
18 engineer would have to be diverted from his or her normal job responsibilities.

19 41. Similarly, Google also does not maintain search queries in the form requested by
20 the Government. The query logs maintained by Google will have to be scrubbed of any personally
21 identifying information. Thus, as in searching the URLs, responding to this request will require
22 writing and implementing new code. I estimate that this will take one to three full-time days of an
23 engineer's time.

24 42. In addition, implementing the code to search and pull random URLs from Google's
25 index and to search and copy the search queries without identifying information will require
26 extended hours of processing time on Google's computers. Because these are not routine functions
27 for Google, the computers will be processing an additional and burdensome program. Running
28

1 such programs above and beyond the normal demand on Google's computers will interfere to
2 some unknown degree with the day-to-day operations of Google. The computers that maintain the
3 index and search queries also process Google's spam filters and process advertising reports.
4 Accordingly, these functions could be slowed down or completely interrupted by processing the
5 Government's request, resulting in lower quality of service to users of Google's search engine and
6 to Google's advertisers.

7
8 I DECLARE UNDER PENALTY OF PERJURY that the foregoing is true and correct.

9 DATED at Mountain View, California, this 16th day of February, 2006.

10
11
12 
13 _____
14 MATT CUTTS
15
16
17
18
19
20
21
22
23
24
25
26
27
28