



GREG CONTI

# Googling

SECURITY



HOW MUCH DOES GOOGLE KNOW ABOUT YOU?

Many of the designations used by manufacturers and sellers to distinguish their products are claimed as trademarks. Where those designations appear in this book, and the publisher was aware of a trademark claim, the designations have been printed with initial capital letters or in all capitals.

The author and publisher have taken care in the preparation of this book, but make no expressed or implied warranty of any kind and assume no responsibility for errors or omissions. No liability is assumed for incidental or consequential damages in connection with or arising out of the use of the information or programs contained herein.

The views expressed in this article are those of the author and do not reflect the official policy or position of the U.S. Military Academy, the Department of the Army, the Department of Defense, or the U.S. government.

The publisher offers excellent discounts on this book when ordered in quantity for bulk purchases or special sales, which may include electronic versions and/or custom covers and content particular to your business, training goals, marketing focus, and branding interests. For more information, please contact:

U.S. Corporate and Government Sales  
(800) 382-3419  
corpsales@pearsontechgroup.com

For sales outside the United States, please contact:

International Sales  
international@pearson.com

Visit us on the web: [www.informit.com/aw](http://www.informit.com/aw)

Library of Congress Cataloging-in-Publication Data available on request.

Copyright © 2009 Pearson Education, Inc.

All rights reserved. Printed in the United States of America. This publication is protected by copyright, and permission must be obtained from the publisher prior to any prohibited reproduction, storage in a retrieval system, or transmission in any form or by any means, electronic, mechanical, photocopying, recording, or likewise. For information regarding permissions, write to:

Pearson Education, Inc.  
Rights and Contracts Department  
501 Boylston Street, Suite 900  
Boston, MA 02116  
Fax (617) 671 3447

Google, Gmail, AdSense, AdWords, Google Maps, Google Analytics, Google Groups, and Google Mobile are all registered trademarks of Google Inc.

ISBN-13: 978-0-321-51866-8

ISBN-10: 0-321-51866-7

Text printed in the United States on recycled paper at RR Donnelley in Crawfordsville, Indiana.  
First printing October 2008

**Editor in Chief**

Karen Gettman

**Acquisitions Editor**

Jessica Goldstein

**Development Editor**

Sheri Cain

**Managing Editor**

Kristy Hart

**Project Editor**

Jovana San Nicolas-  
Shirley

**Copy Editor**

Krista Hansing Editorial  
Services, Inc.

**Indexer**

Erika Millen

**Proofreader**

Kathy Ruiz

**Publishing Coordinator**

Romny French

**Cover Designer**

Gary Adair

**Compositor**

Bronkella Publishing

---

# Preface

---

The rise of the information economy dramatically shifted how we conduct business and live our lives. In the competitive world of business, the organization that has the best access to information can make more timely and effective decisions than its rivals, creating a distinct advantage. In our day-to-day lives, easier access to information and improved methods of communication enrich virtually every facet of our existence. At the heart of this revolution is the Internet, particularly the World Wide Web.

Shortly after the formation of the web in the early 1990s, its commercialization began in earnest. Online companies struggled to find business models that worked in this brave new world, where many of the traditional rules of business no longer applied. The combination of free, easy-to-use tools, along with targeted advertising, emerged as one of the most viable approaches. Customized advertising, by definition, requires insight into the needs of the individual user, which, in turn, requires logging and data mining to be most effective. As a result, online companies have logged virtually every conceivable type of data associated with our use of web-based tools. The existence of this data enables online companies to constantly improve our user experience and support their goal of selling customized advertising.

The value of this data is unprecedented in the history of mankind. If you consider the sum of your online searching, mapping, communicating, blogging, news reading, shopping, and browsing, you should realize that you've revealed a very complete picture of yourself and placed it on the servers of a select few online companies. The thin veneer of anonymity on the web is insufficient to protect you from revealing your identity. If you aren't even a little concerned, you should be. The value of this information is staggering

and ripe for misuse. The threat is even worse when you consider the sum of the disclosures of your company. Everything from the dark little secrets of the corporate executives to the strategic plans of the company exists on someone else's servers. Like water rising behind a dam, this is an issue we need to address sooner or later. Today there are *only* one billion Internet users; however, this value represents just 18% of the world population. Web-based information disclosure will certainly grow as billions more users join us online. Although there is no miracle cure on the horizon, this book is a first step toward a solution. This book clearly illustrates and analyzes the problem of web-based information disclosure and provides you with countermeasures that you can employ now to minimize the threat.

## WHO SHOULD READ THIS BOOK

Stated simply, if you use the web, you should read this book—unless, of course, you have nothing to hide. Not everyone is a pedophile or a terrorist, right? However, I argue that you wouldn't want everything from your health to your politics to your social network stored on someone's server, even if that someone is Google or one of its competitors. Perhaps today this won't be a problem, but history has shown that information leaks, privacy policies change, companies merge, data spills, and attacks occur. The mere existence of this information, and the power it proffers, ensures that it will be coveted by many, including business competitors, insurance companies, law enforcement officials, and governments.

As with any book, there are limits on what can fit between two covers. I've chosen to address a broad range of topics, to make the contents as accessible as possible (and increase the positive impact of the book), but at the same time give enough technical detail to provide insight into the technical challenges that exist behind the scenes. My intent is to raise awareness of the privacy implications of using the tools of Google and other online tools. The threat of web-based information disclosure is an open problem; I've attempted to outline the problem in detail, but complete solutions do not exist. For many, this book will be eye opening. However, I believe some IT and security professionals will have considered some of the points the book brings up. That being said, there will be insights and ah-ha moments for even the most security-savvy readers.

## WHY GOOGLE?

This book studies the security implications of using Google's products and services. Why am I picking on Google when other companies offer similar products and services? Frankly, Google's success and marked innovation makes it the subject of this book. True, many other online companies offer similar products and services, but none comes close to the innovative and comprehensive offerings of Google. Much imitated but rarely surpassed, Google is truly the market leader and *the* company to beat when it comes to online services. Because of its success, Google has been the first to encounter many unanticipated challenges. What do you do when the U.S. Department of Justice subpoenas your search query logs? What do you do when China refuses you access to its two billion potential customers unless you place your servers inside China and comply with the government's requests to censor search results? How does Google find the ethical, and profitable, path when even a simple change in its search-ranking algorithm could make one person a millionaire or destroy a thriving business? Because of Google's sheer size, the tiniest changes have tremendous impact. Even a simple change in Google's default background color from white to an energy-conserving black could save 750 Megawatt hours per year.<sup>1</sup> You can easily see that ethical challenges abound for Google. By its very oft-quoted motto "Don't be evil," Google has set a very high ethical bar for itself. Any perceived deviation draws intense criticism.

Google's mission to "organize the world's information and make it universally accessible and useful" is very admirable and has already made many aspects of our lives easier. Google pursues this mission by offering high-quality tools, often for free. However, by using the tools offered by Google and its competitors, we are actually paying for their use with micropayments of personal information. By studying our interactions with Google, we are studying virtually the entire range of available offerings from *any* competing company. Throughout this book, I analyze the implications of using Google's tools and services, but when appropriate, I include other similar offerings from different vendors. I also include detailed discussion of the threat posed by Internet service providers (ISPs). ISPs have similar visibility on your online activities, although from a different network vantage point. Online companies see information streams from across the planet from hundreds of millions of users. ISPs see all activities from their customers but lack the global reach of Google.

Personally, I am a big fan of Google and use many of its services on a daily basis. By no means do I want Google to fail. However, we have to recognize and address the problem of web-based information disclosure before we reach a point of crisis—a point that I believe is rapidly approaching. We felt a tremor of this problem when AOL inadvertently released the search activity of 658,000 users in 2006.<sup>2</sup> My goal with this book is to outline the problem so we can all start making more informed decisions regarding our use

of “free” web tools, as well as jointly seek solutions that will allow companies such as Google to innovate and thrive, while still meeting the privacy requirements of individuals and organizations. In short, by having such a high-grade product, Google makes itself a high-profile target, and that’s a problem. It is certainly fair game to consider what we provide to Google and how Google protects that product.

It is important to note that this book is based entirely on publicly available information. I did not seek out any proprietary or internal-use information. As you will see, publicly available information is more than enough to understand the problem.

## **A MAP OF THE BOOK**

The book covers many facets of the problem of web-based information disclosure as seen through the lens of Google’s tools and services. The first chapter, “Googling,” is an analysis of Google, its capabilities, its motivations, and its reach; it provides an overview of the types of information individuals and organizations reveal when using the wide variety of tools Google makes available. Chapter 2, “Information Flows and Leakage,” places these disclosures into big-picture context by studying the same information flows, seen from the vantage point of network service providers and individual workstations.

Chapter 3, “Footprints, Fingerprints, and Connections,” studies the information we leave behind as we use the web and how this information can be used to profile our behaviors, be tied to our real-world identities, and be connected with other users, businesses, and groups. Chapters 4–6 deeply examine the risks associated with major classes of online tools, including search, communication, and mapping. Chapter 7, “Advertising and Embedded Content,” illustrates the increasing number of ways users can be tracked as they browse hundreds of thousands of (non-Google) web sites, thanks to embedded advertising, YouTube videos, and similar content.

Chapter 8, “Googlebot,” describes how Google and other large online companies collect and process information around the clock using an army of automated web crawlers. Chapter 9, “Countermeasures,” presents techniques for reducing the impact of web-based information disclosure. Finally, Chapter 10, “Conclusions and a Look to the Future,” analyzes current trends and illustrates what future risks could lie ahead.

The web is an ever-growing and continuously evolving space. Although I have carefully chosen the most relevant topics to include here, a single book is not sufficient to document and analyze the full range of current and future possibilities. With this in mind, I encourage you to visit this book’s companion web site at [www.informit.com/title/9780321518668](http://www.informit.com/title/9780321518668) for additional information.

## ENDNOTES

1. "Change Google's Background Color Background Color to Save Energy?" Slashdot.org, 27 July 2007. <http://hardware.slashdot.org/hardware/07/07/27/054249.shtml>, last accessed 25 September 2007.
2. Ryan Singel, "FAQ: AOL's Search Gaffe and You," Wired.com, 11 August 2006. [www.wired.com/politics/security/news/2006/08/71579](http://www.wired.com/politics/security/news/2006/08/71579), last accessed 25 May 2008.

---

# Advertising and Embedded Content

---

*[My web history is] mine—you can't have it. If you want to use it for something, then you have to negotiate with me. I have to agree, I have to understand what I'm getting in return.<sup>1</sup>*  
—Sir Tim Berners-Lee

Publishing information is the backbone of web content, and bloggers and webmasters frequently rely on embedded content from companies such as Google to enhance the quality of their sites. Unfortunately, embedding third-party content is the equivalent of planting a web bug in web pages,<sup>2</sup> alerting the source of the embedded content to a user's presence on a given site and facilitating logging, profiling, and fingerprinting. More important, the source of the third-party content can aggregate these single instances and track users as they browse the web. This notion deserves restating: The simple act of web browsing across many disparate sites has the potential to generate a continuous stream of information back to the providers of third-party content. The more popular a given third-party service is, the more sites will deploy their content, and the greater the window of visibility on users' web surfing activity becomes. In the case of advertising

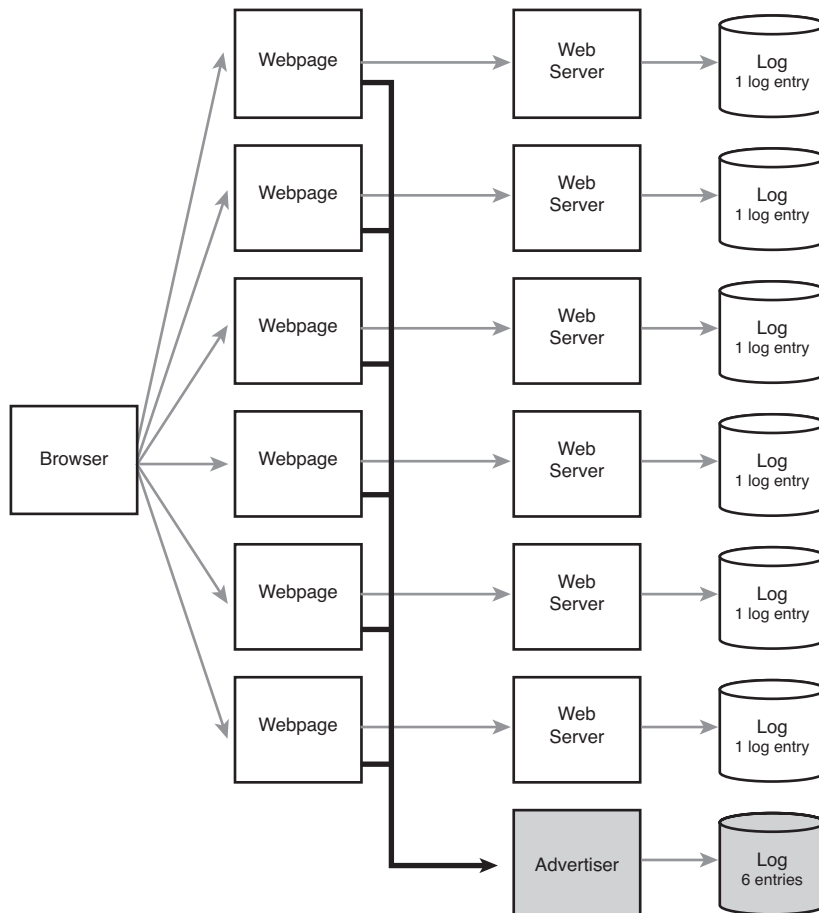


networks such as Google/DoubleClick and web-analytics services such as Google Analytics, the risk is large indeed.

This chapter explores the risks associated with embedded content by focusing on Google's advertising network and Google Analytics, but it also provides an overview of other forms of embedded content that present related risks, such as embedded YouTube videos, maps, and Google's Chat Back Service.

## **CROSS-SITE TRACKING**

As mentioned in Chapter 3, "Footprints, Fingerprints, and Connections," many web sites embed third-party content in their sites. Third-party content can take the form of legitimate images and video clips, among other forms of content, but it can also be used to track users as they surf the web. Advertisers and web-analytic services give webmasters enticing analysis tools and advertising profit, simply requiring that, in exchange, webmasters add small snippets of HTML and JavaScript to their pages. Unfortunately, such third-party content is a severe privacy and web-based information-disclosure risk because the user's web browser automatically visits these third-party servers,<sup>3</sup> where their visit is presumably logged and their browser tagged with cookies. More important, the larger the advertising network is, the larger the window a given company has on a user's online activity. For example, if a user visits 100 different web sites, each containing advertisements from a single advertising service, that service can observe the user as he or she visits each site. Figure 7-1 depicts cross-site tracking via an advertising network. In this figure, a user visits six distinct web sites, each hosting content from a single advertiser. In turn, the user's visits create one set of log entries on each of the six legitimate servers. However, because each visit contained an advertisement from a single advertising network, the advertiser is able to log all six visits.



**Figure 7-1** Example of cross-site tracking by an advertising network. When the user visits six distinct sites, he or she generates one set of log entries at each site. However, if each site contains advertisements from a single advertising network, the advertiser is able to record all six visits.

Let's look at a real-world example by visiting a popular web site, MSNBC (see Figure 7-2). As it turns out, the MSNBC web site is laden with third-party content.



**Figure 7-2** Analyzing the MSNBC web site demonstrates that it contains a great deal of third-party content. The problem is rampant among other web sites, large and small.

In a world without third-party content, the user should simply receive content from MSNBC's domain, `msnbc.msn.com`. In the real world, however, the user visits 16 additional domains from 10 different companies. Two of these domains, DoubleClick and GoogleAnalytics, are owned by Google. The web browser provides the user with little assistance in detecting third-party content. The user simply sees the browser's status bar rapidly flicker as the browser contacts each new site. To provide a clearer picture, I captured the raw network activity using the Wireshark protocol-analysis tool and created Table 7-1 to detail each of the third party domains visited.

**Table 7-1** Third-Party Sites Visited When Browsing the MSNBC Web Site

Domain	Notes
a365.ms.akamai.net a509.cd.akamai.net	Domain owned by Akamai.com, a mirroring service for media content
ad.3ad.doubleclick.net	Digital marketing service, acquired by Google
amch.questionmarket.com	Hosting web site where online surveys are posted
c.live.com.nsac.net c.msn.com.nsac.net rad.msn.com.nsac.net	Registered to Savvis Communications, a networking and hosting provider
context3.kanoodle.com	Search-targeted sponsored links service
global.msads.net.c.footprint.net hm.sc.msn.com.c.footprint.net	Registered to Level 3 Communications, a large network provider
msnbcom.112.2o7.net	Registered to Omniture, a web analytics and online business optimization provider
prpx.service.mirror-image.net wrpx.service.mirror-image.net	Registered to Mirror Image Internet, a content delivery, streaming media, and web computing service
switch.atdmt.com view.atdmt.com	Registered to aQuantive, parent company to a family of digital marketing companies
www-google-analytics.l.google.com	Traffic measurement and interactive reporting service offered by Google

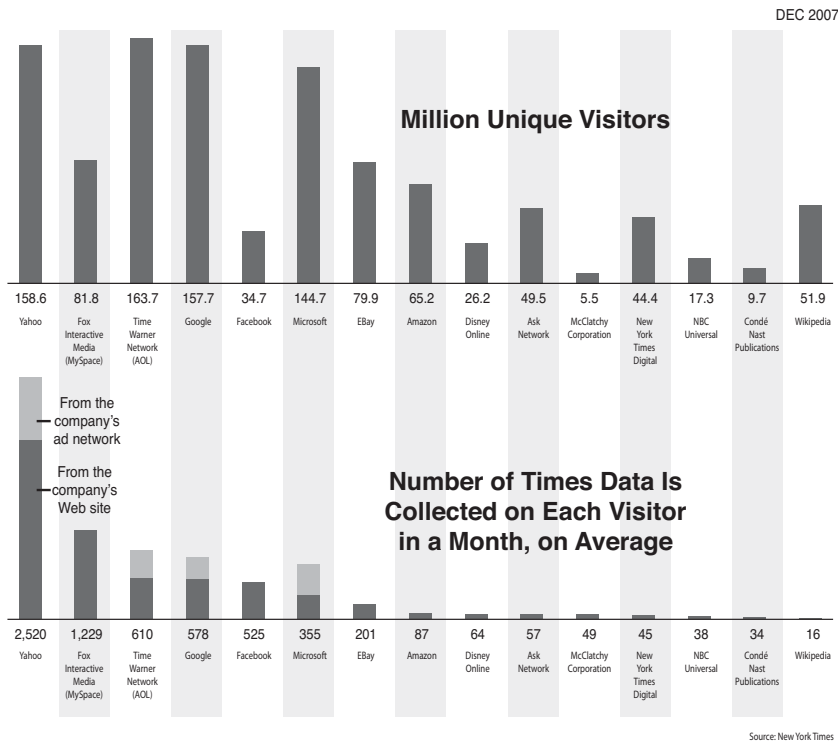
Think about it. Simply visiting a single web page from a popular news service informs 16 third-party servers of the visit, a 16-fold magnification of logging. This is not a manufactured example, but it is representative of a common practice. Embedding third-party content in web sites is ubiquitous, and so is the problem. The end result is that web surfers are frequently tracked by companies they've never even heard of. It is also worth considering that information sharing via embedded content doesn't occur only with "third parties"—sharing can also occur between ostensibly separate entities that are actually owned by the same parent company. For example, the A9 search engine (an Amazon.com company) inserts search term-related Amazon book advertisements adjacent to search results. These advertisements allow Amazon to track what A9 users search for, click through, and possibly buy online. If the user does make a purchase on Amazon.com, Amazon knows that user's real-world identity, including billing and shipping information. In the case of A9, Amazon makes clear on the A9 site that A9 is an Amazon.com company, but the important idea is that corporate ownership—and, hence, implicit information sharing—might not be obvious as users browse the web.

## ADVERTISING

*A famous New Yorker cartoon from 1993 showed two dogs at a computer, with one saying to the other, “On the Internet, nobody knows you’re a dog.” That may no longer be true.<sup>4</sup>*

—Louise Story

Advertising is the fuel behind virtually all free online tools. Advertising is also the means for tracking your web surfing across the Internet. Anytime you visit sites that serve advertisements from a common advertising network, your activities can be logged. These logs can then be used to create precise profiles, facilitating tailored advertising. Importantly, log analysis isn’t a static process. For example, advertising companies are actively developing technologies to anticipate people’s next steps.<sup>5</sup> Based on the popularity of the largest web companies, the amount of information they can collect is staggering (see Figure 7-3).



**Figure 7-3** New York Times analysis of the data collection conducted by some of the largest online companies.

Web advertisements are big business, and some of the largest services are offerings by Google, including AdSense, AdWords, and DoubleClick. However, advertisements are appearing in many other forms. Microsoft is quietly offering an ad-funded version of its Works office suite.<sup>6</sup> Pudding Media, a San Jose–based startup, is offering advertisement-supported phone service.<sup>7</sup>

Google understands the growth potential of advertising. In 2007, the company announced that it would pay \$3.1 billion to acquire DoubleClick, a leading online advertiser, whose current estimated revenues at the time were \$150 million.<sup>8</sup> A year ago, Yahoo! made a similar, smaller-scale move by acquiring the online global ad network BlueLithium for approximately \$300 million.<sup>9</sup> At the time of this writing Yahoo! was reportedly considering an agreement to carry search advertisements from Google, amid a potential hostile takeover attempt by Microsoft.<sup>10</sup> Although it is impossible to know for sure at this time, such an allegiance carries the very real potential of allowing Google to gather search terms, issue cookies, and conduct other activities based on Yahoo!’s extremely large user base. Make no mistake, a key component of such acquisitions and alliances is access to user data and the power it provides.

#### **NOTE**

Online advertising is a growth industry that is eating into traditional media markets. For example, the media-analysis company Simmons found that Internet video advertisements are 47% more engaging—and, hence, more effective—than traditional television advertisements.<sup>11</sup>

Chapter 3 should have convinced you that users scatter significant, and often personally identifiable, information behind as they surf the web. When *New York Times* reporter Louise Story decided to determine how personally identifiable the information maintained by large web companies is, she asked four large online companies a single question: “Can you show me an advertisement with my name in it?” Story provided the following summary of the responses.<sup>12</sup>

Microsoft says it could use only a person’s first name. AOL and Yahoo! could use a full name, but only on their sites, not the other sites on which they place ads.

Google isn’t sure—it probably could, but it doesn’t know the names of most of its users.

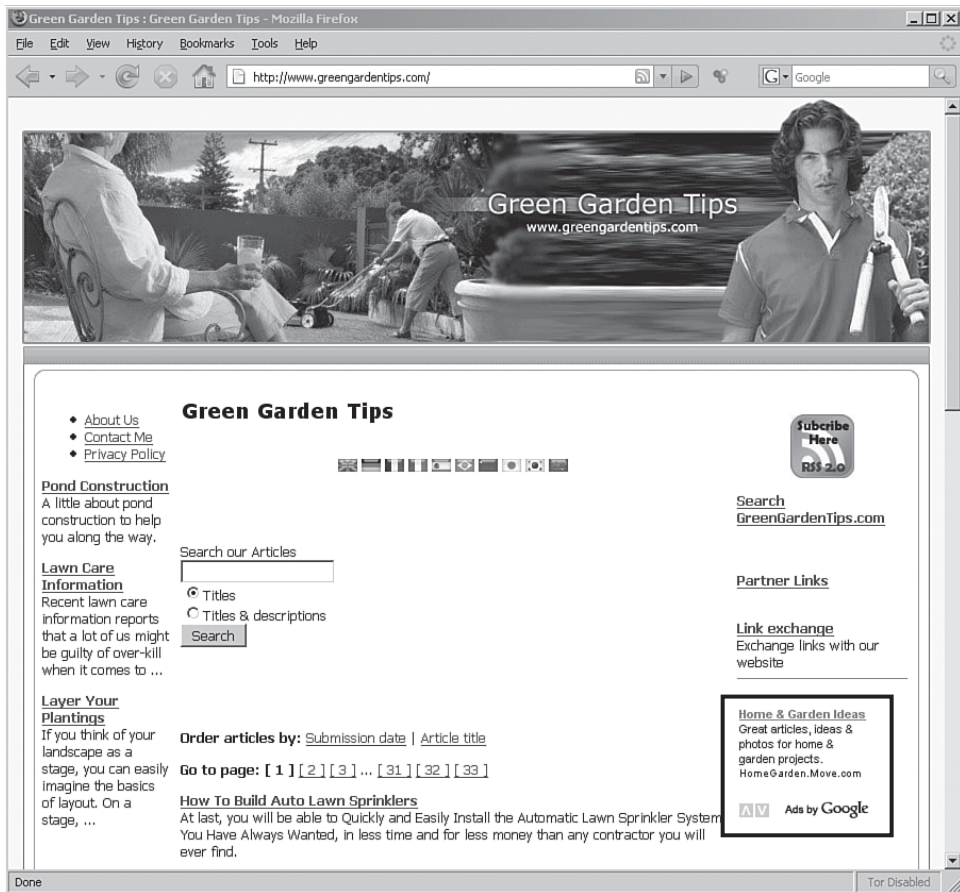
Although these results are telling in their own right, keep in mind that these are official responses provided to a *New York Times* reporter. In other words, these are legal

opinions on the subject, as opposed to technical capability. What the four companies have the *capability* to do is an entirely different matter. Advertising networks are the net that allows large online companies to gather this precise information and use it for user profiling, data mining, and targeted advertising. Because of this capability to aggregate and analyze user information, advertising campaigns can follow users as they switch from independent sites. You might have encountered this technique when shopping online. For example, if you were searching for flat panel monitors on site X, when you hopped over to visit site Y, low and behold, there were advertisements for flat panel monitors.

Microsoft is openly touting an “Engagement Mapping” approach that seeks to move beyond the “outdated ‘last ad clicked’” model by understanding “how each ad exposure—whether display, rich media or search, seen multiple times on multiple sites and across many channels—influenced an eventual purchase.”<sup>13</sup> The article also states that Microsoft intends to use data on user behavior *before* clicking an ad to be able to say that an unclicked ad still made an impression. If the article is to be believed, Microsoft intends to collate that information with search queries and sites visited within the period of a day or a couple of days. Two questions immediately arise. Where does that other data come from (and is it a mere coincidence that Microsoft is seeking to acquire Yahoo! and other search companies)? Now that we are talking about keeping tabs of long-term user behavior, and hypothesizing why a user did this or that, what other kinds of data mining will this floodgate open? They are essentially saying, “We are going to dedicate a lot of resources to watching where you go and what you see on the web.” The fact that Microsoft Windows is the most popular operating system on Earth<sup>14</sup> just magnifies the concern. The law surrounding online advertising and the collection of user data is still immature, so advertisers have a very wide lane in which to operate. For example, a U.S. Federal Court ruled that ads displayed by search engines are protected as free speech when deciding what advertisements to display.<sup>15</sup>

## AdSense

Google AdSense<sup>16</sup>, sometimes called Google Syndication, is an advertising service Google provides that allows webmasters to earn advertising revenue by hosting AdSense ads (see Figure 7-4). These revenues aren’t trivial, commonly ranging from a few hundred dollars a month to \$50,000 or more per year, making the service extremely popular.<sup>17</sup> AdSense advertisements are context-sensitive ads served by Google based on the hosting site’s content. Unfortunately, merely visiting a web site hosting these advertisements informs Google of the user’s IP address and gives Google the opportunity to log the user’s visit and tag the user’s browser with a cookie.



**Figure 7-4** Screenshot of Google AdSense. Notice the (debatably) unobtrusive advertisement in the bottom-right corner. Embedded advertisements such as these alert the advertising network of your presence on a site.

AdSense isn't limited to textual ads on traditional web pages. Google is experimenting with AdSense for other forms of content, including RSS Feeds,<sup>18</sup> web site search boxes,<sup>19</sup> mobile content,<sup>20</sup> video, and Cost Per Action AdSense.<sup>21,22</sup> Nor are AdSense and similar services limited to minor sites. Major online retailers also participate. For example, eBay signed deals to run ads from Google and Yahoo!.<sup>23</sup> Figure 7-5 shows an example with



eBay and Yahoo!. When searching for an item on eBay, Yahoo! servers provide contextual advertisements, leaving open the likelihood of Yahoo!'s logging of eBay visitors.<sup>24</sup>



**Figure 7-5** Screenshot of Yahoo! advertisements embedded in an eBay web page. Because these ads are pulled directly from Yahoo! servers, user information such as cookies and IP addresses is disclosed directly to Yahoo!.

**WARNING**

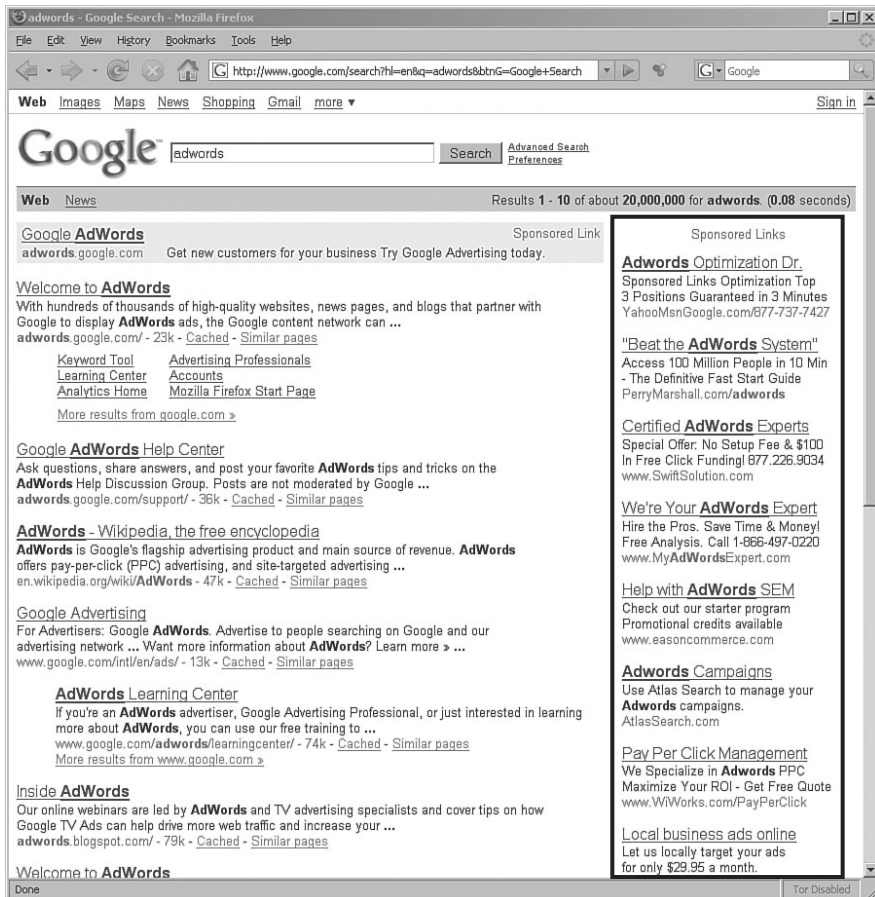
It is also important to consider the identity-disclosure risks from the perspective of the webmaster, who must use the user's registered Google account to log on to the service and use the web interface to administer his or her accounts.<sup>25</sup>

The future of AdSense is difficult to determine. Some analysts believe that Google is acquiring sites that will provide traffic itself instead of paying advertising fees to third-party sites.<sup>26</sup>

**AdWords**

AdWords is a fundamental part of Google's business model. According to the BBC, every time a user conducts a search on Google, the company makes 12¢ in revenue.<sup>27</sup> When you consider that Google receives more than 60 billion searches per year in the United States alone, you can see that the program generates huge profit. Google believes that AdWords "is the largest program of its kind."<sup>28</sup> Using AdWords, would-be advertisers bid on search terms that are displayed as part of the user's search results; the better the placement, the higher the cost. AdWords are relatively unobtrusive, but quite effective, advertisements (see the right side of Figure 7-6).<sup>29</sup>

AdWords poses both information-disclosure risk and other security risks. Attackers have used AdWords and similar services to misdirect users to malicious sites; see the "Malicious Ad Serving" section later in the chapter.<sup>30</sup> However, Google's AdWord Partners are a significant information-disclosure risk because searches from these sites can be sent to Google. According to Google those who have already joined their "growing advertising network" include AOL, Ask.com, Ask Jeeves, AT&T Worldnet, CompuServe, EarthLink, Excite, and Netscape.<sup>31</sup> Even third-party search engines that delete their logs locally are still at risk. Take, for example, Ask.com, who took an industry-leading position by offering AskEraser, a function that deletes search activity from Ask.com servers.<sup>32</sup> However, Google delivers the bulk of Ask's advertisements, so user information, including the search query and IP address, are passed back to Google each time a page is served to a visitor.<sup>33</sup>



**Figure 7-6** Screenshot of Google AdWords. Notice the advertisements on the right side of the image.

## GOOGLE DOUBLECLICK

DoubleClick is a major online advertising service, long criticized for using cookies and IP addresses to track users as they surf the web.<sup>34</sup> DoubleClick is an extremely popular advertising service and counts a large number of Fortune 500 companies as clients. In 2007, Google announced a definitive agreement with DoubleClick for \$3.1 billion in cash to acquire the company. The acquisition drew the attention of the U.S. Federal Trade Commission and European Regulators, who investigated antitrust and privacy

implications but eventually acquiesced.<sup>35, 36, 37, 38</sup> Google closed the acquisition of DoubleClick shortly thereafter.<sup>39</sup>

**WARNING**

DoubleClick allows users to opt out of its ad-serving and search products by issuing a special cookie. This is an elegant solution from the perspective of DoubleClick because users concerned enough to opt out might delete this cookie by accident as they remove other traditional tracking cookies. One possible solution is a browser plug-in that allows easy-to-use and fine-grained cookie control, although this approach still assumes that companies employing opt-out cookies will honor the request and not log user activity anyway.<sup>40, 41</sup>

The implications of a combined Google–DoubleClick dreadnaught are significant. Google excels in search advertising (AdWords) and simple textual advertisements (AdSense). On the other hand, DoubleClick excels in “display advertising,” such as flashy banner ads and video advertisements, which reach between 80% and 85% of the web population.<sup>42, 43</sup> The end result is a broad net that permits Google to track a user’s web searches and web site visits, with the potential to impact the privacy interests of more than 1.1 billion Internet users worldwide.<sup>44</sup> This acquisition underscores the fact that mergers and acquisitions are about data, including both existing data stockpiles and access to continued data streams.

**NOTE**

Google’s acquisition of DoubleClick places significant competitive pressure on Microsoft.<sup>45</sup> As a response, Microsoft actively sought to acquire Yahoo! to increase its advertising reach.<sup>46</sup>

**ADVERTISING RISKS**

Advertising poses more risks than those already discussed regarding AdWords, AdSense, and DoubleClick. Attackers can exploit advertising networks to compromise end-user machines, unethical interface techniques can trick users into disclosing sensitive information, and historically unbiased network providers can insert advertisements as the web pages make their way to the user’s browser.

## MALICIOUS AD SERVING

Advertising networks are more than just information-disclosure risks. They also serve as a malware attack vector. Advertising services pay web site owners for publishing advertisements on their web sites. A very common technique is the banner ad we've all seen at the top of web pages. Such ads usually take the form of animated GIF files, but they now include many image and video formats. Individuals and organizations that want to advertise using such a service create a media file and pay an advertiser a fee, and the advertiser serves the image to thousands of visitors of sites that belong to its advertising network. The risks here are twofold. Attackers have created misleading advertisements as a means to draw traffic to a malware serving or other malicious web site.<sup>47</sup> The users' trust of the advertisement company and the hosting web site increases their trust of the advertisements, leaving web surfers more vulnerable to such an attack. Virus writers have used the Google Adwords service to serve text ads that appeared to link to legitimate destination sites, but silently infected vulnerable web surfers by routing users through an intermediate, malicious site. Attackers also have used a vulnerability in Internet Explorer to compromise visitors as they passed through the intermediate site, before ultimately arriving at the legitimate site.<sup>48</sup>

The ads themselves have also been used to attack the web user directly. Malformed graphical images are one common technique. For example, a banner advertisement displayed on MySpace served spyware to more than one million visitors. In this case, attackers exploited a flaw in the way Windows processed Windows Meta File (WMF) images to install a Trojan horse.<sup>49</sup> Because the attack occurs when the browser displays the image, the user needn't click the advertisement to be infected. Another attack, served by DoubleClick, used rich media advertisements created in Adobe Flash to exploit a similar vulnerability, with the malicious advertisements appearing on extremely popular sites, including those of *The Economist* and Major League Baseball.<sup>50</sup> Rich media advertisements are highly interactive and are becoming increasingly popular. Even the seemingly simple task of securing browsers against malicious images is proving difficult, and complex, rich media tools such as Flash are proving to be an even greater challenge.<sup>51</sup> Other complex environments, such as ads embedded in Adobe PDF files, are now being explored as ways to reach potential customers, and I expect that similar issues will arise.<sup>52</sup>

### NOTE

Google has taken an active role in countering malicious advertisements and web sites. For example, Google warns users of potentially malicious sites by clearly labeling suspect sites in their search result listings. As another example, Google's I'm Feeling Lucky search button no longer automatically redirects a user to a suspicious site; instead, the user is presented with a list of search results instead.<sup>53</sup>

## MALICIOUS INTERFACES

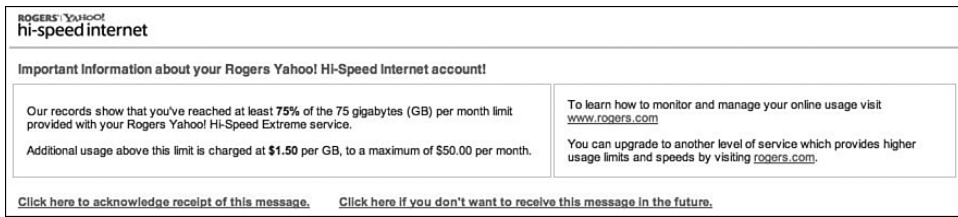
Beyond malicious ad serving, advertisers employ another concerning strategy that I call *malicious interfaces*. In the idealistic world of interface design theory, interface designers always operate in the best interests of their users. Designers carefully study user tasks and painstakingly craft interfaces and applications to help users accomplish them. However, in the world of online advertising, the exact opposite is true. Designers frequently violate design best practices to coerce or mislead users into viewing advertising. Examples abound on the web: fake hyperlinks that pop up advertisements, giant advertisements that cover the text of articles, distracting advertising videos that begin playing the moment a page is viewed, banner ads with fake buttons that appear to be a part of the interface, advertisements embedded in video clips ... the list goes on. Malicious interface designers are creative; new “innovations” are coming out regularly. The only constraining factor appears to be the tolerance of the user. The invasiveness of advertisements is getting worse. You may have heard the term “banner blindness.” Users quickly learned that banner ads are of little value and ignore the advertisements, to the point that they barely perceive banner ads anymore. This defense mechanism has forced advertisers to become more aggressive in capturing user attention.<sup>54, 55</sup> Although, I don’t claim that Google employs malicious interface design, the trend is concerning; it seems that malicious interface designers are carefully seeking the sweet spot between making advertising profit and annoying the user so much that they abandon a given site altogether.

### WARNING

In many cases the link displayed by an embedded advertisement is not the actual link. Nor will hovering over the link display the destination URL in the user’s browser status bar. The actual link goes first to the ad server so the click can be logged. The user’s browser is then redirected to the page chosen by the advertiser.

## HOSTILE NETWORKS

As carriers of key components of network infrastructure, ISPs and web hosting services are flexing their muscle to place advertisements in front of users. For example, domain registrar and web hosting provider Network Solutions hijacked customers’ unused sub-domains to resort to ad-laden “parking” pages.<sup>56</sup> As another example, bloggers Lauren Weinstein and Sarah Lai Stirland reported Canadian ISP Rogers modification of web pages en route (see Figure 7-7).



**Figure 7-7** Screenshot of an ISP altering a Google web page

Inserting advertisements into web pages as they transit an ISP is potentially very big business. ISPs already have access to a tremendous amount of personal information about users' online activity. This data is a veritable gold mine if used to target online advertising. Some ISPs are reportedly selling significant amounts of user data to online marketers.<sup>57</sup> In the United Kingdom, three major ISPs have announced plans to use user clickstream data to insert relevant advertisements as they surf, through a new startup called Phorm.<sup>58</sup> ISP data contains some of the most sensitive information disclosures made by online users. If this advertising technique becomes widespread, virtually every web surfer's activities will be passed on to advertisers in some form. Fighting the issue will be difficult, and users might find that they are faced with little alternative than to accept this new status quo.

## AFFILIATE SERVICES


Affiliate advertising isn't just confined to Google via its AdSense and DoubleClick programs. It is a popular marketing practice in which, in its common form, web authors embed advertisements that contain unique tracking data to identify the correct affiliate to compensate. Such advertisements can be static—that is, the advertisement exists entirely on the server of the web author. In this case, the user must click the advertisement before the advertiser is aware of the user. However, many affiliate services provide dynamic content that is pulled directly from the online company without any action by the user, immediately linking the user to the visited web site (see Figure 7-8 for an example). In addition, the online company might tag the user with a cookie or retrieve an existing cookie, opening the possibility of identifying the user by name, billing address, and shipping address.

### Build Your Amazon Widgets

To build a widget or a link, start by selecting one of the options below. You'll then be able to customize your widget or link, and we'll provide the HTML for you to use on your site. [Learn more](#)

Showing 1 - 21 of 21 Sort by

#### MP3 Clips Widget



Add music to your web site with the MP3 Clips widget. Search through Amazon's catalog of DRM-free MP3 music and add entire albums or select specific MP3 tracks to add to your widget.

[Add to your Web page](#)

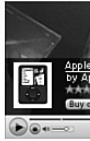
#### Carousel Widget



Take your products for a spin in the Carousel widget. Hand-pick your favorite products or choose from several lists including Amazon Bestsellers and Hot New Releases.

[Add to your Web page](#)


#### Your Video Widget



Add product links to the videos you already post on your web page!

[Add to your Web page](#)

#### Deals Widget



Showcase the hottest deals from Amazon on your web page. Delight your viewers with the Deal of the Day, Lightning Deals or Our Best Deals from across Amazon.

[Add to your Web page](#)


#### Unbox Videos



Show studio previews for a huge range of recent and classic movies or tv shows available from the Unbox video download service.

[Add to your Web page](#)

#### Slideshow



Make elegant slideshow widgets out of images chosen from products across the entire Amazon catalog - CD cover art, DVD's, books, anything!

[Add to your Web page](#)

#### Product Cloud



Take tag clouds one step further by showing Product Clouds - clusters of product titles relevant to your page

[Add to your Web page](#)


#### My Favorites



Express yourself by recommending and commenting on products from Amazon.com. Let everyone know how you feel and what you like!

[Add to your Web page](#)

#### Wish List



Show your Amazon Wish List in a widget and let everyone know what you like. Or show your friend or relatives Wish List to help plan a special occasion.

[Add to your Web page](#)

**Figure 7-8** Example of Amazon affiliate network dynamic advertisements. Some affiliate networks, such as Amazon, encourage web authors to embed these advertisements into their web content, allowing the user to be logged, tagged with a cookie, and perhaps identified by name, merely by visiting the web page containing the advertisement.<sup>59</sup>

## FACEBOOK BEACON

Facebook's Beacon illustrates both the detailed insight that large online companies have into the activities of their user populations and the lengths some companies will go to increase profit. Facebook's Beacon service allows Facebook users to share their purchases from affiliated online companies, such as books, movies, and gifts, with their Facebook



friends.<sup>60</sup> However, when the service was first offered, participation was “opt out.” The end result is that many users were infuriated and the civil action group MoveOn.org initiated a very effective online petition that rapidly gained more than 50,000 supporters.<sup>61</sup> MoveOn even blamed Facebook for “ruining Christmas” because the Beacon advertising system allowed users to see holiday present purchases made by friends and family.<sup>62</sup>

Shortly thereafter, Facebook changed the service to “opt in” by requiring explicit user permission before publishing purchases to the user’s Facebook friends.

## **OTHER CROSS-SITE RISKS**

The preceding sections illustrated the risks associated with AdWords, AdSense, and DoubleClick, but cross-site information-disclosure risks do not end with advertising. The web functions thanks to hyperlinks and embedded content. Through these vectors, Google and other large companies can gather tremendous amounts of user information and help link clusters of information to individual users, companies, and other organizations. The following examples share one common characteristic: Each relies upon third-party web masters to embed tracking (or trackable) content into their web sites. Most web masters would not add such content arbitrarily; instead, they are enticed by at least a nominal incentive for cooperation.

### **WARNING**

Embedding Google or other third-party content into a web page isn’t the only concern. Posting your own content on a Google-hosted service, such as Blogger, YouTube, or Orkut, allows Google to track all users to the site. In addition, you become dependent on Google and its ideas about censorship, privacy, and quality of service.

## **GOOGLE ANALYTICS**

Google Analytics is a free tool for webmasters that provides a powerful and intuitive interface for analyzing web log data (see Figure 7-9).<sup>63</sup> Google Analytics is part of a class of applications that provide statistical and graphical analyses of web visitor activity based on web server log data and (optionally) on data gained via cookies placed on users’ computers, web bugs, and JavaScript code. Such tools display site visitor reports (for example, geographic locations of visitors, most active visitors, and browsers used), page view

reports (for example, entry/exit pages, most popular time of day, and number of requests for each page), server reports (for example, amount of bandwidth consumed and which files were requested), and referrer reports (for example, search queries and referring URLs). Other popular web analytics software includes Webalizer ([www.mrunix.net/webalizer/](http://www.mrunix.net/webalizer/)) and WebTrends ([www.webtrends.com/](http://www.webtrends.com/)).



**Figure 7-9** Google Analytics gives Google the capability to track users as they visit any Google Analytics member site.

Google Analytics is easy to install. Webmasters need only paste code similar to the following into web pages they want the service to track.

```
<script src="http://www.google-analytics.com/urchin.js"
type="text/javascript">
</script>

<script type="text/javascript">
_uacct = "UA-994065-1";
urchinTracker();
</script>
```

This code is straightforward JavaScript. It serves as a hook in each web page to contact Google whenever a page is loaded and download a JavaScript file called `urchin.js`. The `_uacct` variable stores a unique tracking code assigned to the webmaster. The script then launches the `urchinTracker()` function in the newly downloaded `urchin.js` file. Unfortunately, the code within `urchin.js` is far more complex and apparently obfuscated.<sup>64,65</sup> The following is a short snippet:<sup>66</sup>

```
function urchinTracker(page) {
  if (_udl.protocol=="file:") return;
  if (_uff && (!page || page=="")) return;
  var a,b,c,xx,v,z,k,x="",s="",f=0;
  var nx=" expires="+_uNx()+"";
  var dc=_ubd.cookie;
  _udh=_uDomain();
  if (!_uVG()) return;
  _uu=Math.round(Math.random()*2147483647);
  _udt=new Date();
  _ust=Math.round(_udt.getTime()/1000);
  a=dc.indexOf("__utma="+_udh);
  b=dc.indexOf("__utmb="+_udh);
  c=dc.indexOf("__utmc="+_udh);
  if (_udn && _udn!="") { _udo=" domain="+_udn+""; }
  if (_utimeout && _utimeout!="") {
    x=new Date(_udt.getTime()+(_utimeout*1000));
    x=" expires="+x.toGMTString()+"";
  }
}
```

Although many webmasters sing the praise of Google Analytics, the tool also poses a significant privacy concern for web surfers. Each time they visit a web page that contains the request to download `urchin.js`, the user's web browser contacts a Google server and downloads and then executes the script, leaving behind all the typical web-browsing footprints described in Chapter 3. The `urchin.js` script presumably discloses additional information, but Google does not provide specific details. The primary risk of Google Analytics is that it gives Google the capability to track users as they browse from web site to web site, including the use of cookies.<sup>67</sup> There is no official count of the number of participating web sites, but several years ago, an analyst estimated the number to be about 237,000.<sup>68</sup> The number now is presumably far greater. Some of the most popular sites on the web employ Google Analytics, such as `Slashdot.org`, which downloads Google's newer `ga.js` script, an equally difficult script to interpret. You can see a snippet in Figure 7-10.

```

var _gat=new Object({c:"length",p:"cookie",b:undefined,bb:function(d,a){this.wb=d;this.Hb=a},o:"__utma=",
0;h--)(o=d.charCodeAt(h);a=(a<<6268435455)+o+(o<<14);c=a&266338304;a=c!0?a^c>>21:a))return a),B:function
(c,d):escape(d),z:function(d,a){var c=decodeURIComponent,h=d.split("+").join("")};if(c instanceof Func
v:function(d,a){return d.indexOf(a),D:function(d,a,c){c=_gat.b=c?d[_gat.c]:c;return d.substring(a,c)},m
c("alltheweb","q"),c("gigablast","q"),c("voila","rdata"),c("virgilio","qs"),c("live","q"),c("baidu","wd"),
"/";d.ha=100;d.Da="/__utm.gif";d.ta=1;d.ua=1;d.F="";d.sa=1;d.qa=1;d.nb=1;d.f="auto";d.C=1;d.Ga=100;d.Mc
i[s][0]+b,e)};f.Eb=function(){return n.b=g|g=f.t();f.Ba=function(){return m?m:"-"};f.Qb=function(k){m
for(var b=0;b<j[y];b++)if(b<4&&1.n.Ea(j[b]))j[b]="-";f.tc=function(){return p};f.Ic=function(k){p=k};f.Ic
b)(var e=f.U,i=B,l,s,f.Ha(k);B.l=b;for(s=0;s<c[y];s++)if(!t(e[s][1]{}))e[s][3]{};B.l=1;f.Yb=function(){l
new Image(1,1);u.src=h.Da+r;u.onload=function(){j()};if(1==B|2==B){var x=new Image(1,1);x.src=("https:"=
m)p.g(j.ca,new p.h.ab(r,d,a,c,h,o));else(m.Xb=r;m.Oa=d;m.K=a;m.qb=c;m.Ub=h;m.Kb=o);_gat.h$.prototype.Bb
n)(n=new f.h.$(d,a,c,h,o,j,m,r);f.g(p.la,n))else(n.mb=a;n.Ub=c;n.Vb=h;n.Sb=o;n.sb=j;n.Ub=m;n.vb=r);return
new XMLHttpRequest("6");p=p+WIN 6,0,21,0";f.AllowScriptAccess="always";p=f.GetVariable(t)}catch(B){}if(!
a.n&&a.n.javaEnabled()?!0:a.yb=o?j{}:c;a.rb=h.d(a.a.characterSet?a.a.characterSet:(a.a.charset?a.a.chars
v[z]:if(f(q,t(1.wb))){x=y(x,"?").join("&");if(f(x,"&"+1.Hb+"="))(u=y(x,"&"+1.Hb+"=")[1];if(f(u,"&"))u=y(u
q=y(q,"/")[])}if(O==m.v(q,"www."))q=m.D(q,4);return new m.k.q(p,q,p,"referral","referral",p,u);j.kc=fu
";x=j.kc(j.a.location);if(j.r.H&&q.Eb()){z=q.Ca();if(!r(z)&&f(z,"");)(q.Ra();return"");z=n(e,m.X+1,"");
F?1:F;q.Rb([1,j.ja,F,k,l.Ka()).join("").);q.Ra();return"utmcr=1"}else return"utmcr=1"};_gat.k.q=functi
a(d.ra));_gat.k.q.prototype.zb=function(d){var a=this,c=_gat,h=function(o){return c.z(c.B(d,o,""))};a.u
l[j[b]];if(a.b!=e){if(k)v+=j[b];v+=x(e);k=false}else k=true)return v}function x(i){var v=[];k,b;for(b=0;b
[1.M(i)];k;for(k in c)if(a.b!=c[k]&&1.yc(k))a.g(v,k.toString()+c[k]);return v.join("");d._setKey=funct
o));_gat.cc=function(d,a){var c=this;Cc=a;c.Dc=d;c._trackEvent=function(h,o,j){return a._trackEvent(c.j
i){if(o(b)|o(e)|o(i))return"-";var s=r(b,c.o+a.e,e),w;if(!o(s)){w=f(s,"");w[5]=w[5]?w[5]*1+1:w[3]=w[
]"==g.f||"none"==g.f){g.f="";return 1}q(i);if(g.nb)return c.t(g.f);else return 1};a.lc=function(b,e){if(o(
return b)a.Oc=function(b){if(a.P()){var e="";if(a.j!=h&&a.j.M().length>0)e="sume="+c.d(a.j.M())+e+a.I
a.a.createElement("script");b.type="text/javascript";b.id="gasajs";b.src="https://www.google.com/analyti
o(1.za())(E=x(G,"&"),e);a.L=true}else(C=f(i.I(),"");s=C[0])else if(F){if(!I){D}(E=x(b,"");e);a.L=true)el
function(i){var b;if(!B){a.zc(i);a.ea.rc(i);a.s=new c.Y(a,a,g)}if(z(i)).xc(i);if(!B){if(z(i))(a.vaa.lc(a.a.r
r(b,"gasos","&");r(a.a[c.p],c.Sa,"");if(e[t]>10){a.A=e;if(a.V.addEventListener)a.V.addEventListener("lo
r(g.J&&g.J[t]>0)a.Jc(i);a.Oc(b);a.L=false});a._trackTrans=function(i){var b=a.e,e=[i,s,w,A];a._initData().
l(w[e]);i=f(w[e],g.F);for(s=0;s<i[t];s++)i[s]=l(i[s]);if("T"==i[0])a._addTrans(i[1],i[2],i[3],i[4],i[5],i
a.G,a.a,e)});a._link=function(b,e){if(g.H&&b){a._initData();a.a[n].href=a._getLinkerUrl(b,e)});a._linkB
function(i){a._initData();return new c.Z};a._sendXEvent=function(b){var e="";a._initData();if(a.P()){e="&
false;return w};a._trackOutboundUrl=function(b){a._initData();if(a.P()){var e=new c.Z;e._setKey(6,1,b);y.i
b);a._clearIgnoredRef=function(){g.ga=[]};a.Tc=function(){return g.ga};a._setAllowHash=function(b){g.nb=
function(i){return g.ta};a._setLocalGifPath=function(b){g.Da=b};a._getLocalGifPath=function(){return g.Da
a._setCampContentKey=function(b){g.db=b};a._setCampIdKey=function(b){g.eb=b};a._setCampMediumKey=function

```

**Figure 7-10** Screenshot of the Google Analytics script `ga.js` downloaded from Google when users visit Slashdot.org

The risk of being tracked across 250,000 or more web sites is concerning enough, but the true risk of Google Analytics is that the user data can be combined with web sites participating in Google's AdSense and AdWords programs, enabling the company to track users across a broad swath of the most popular portions of the web. Users see only a brief flicker in their browser's status bar as their browser contacts Google's servers. The potential of "free" web-analytics software is not lost on Google's competitors; both Yahoo! and Microsoft recently released free web-analytics tools.<sup>69,70</sup>

## CHAT BACK

Google's Chatback service enables web authors to embed a status indicator, a "badge," directly into their web pages. When the page is loaded, the badge (see Figure 7-11) indicates whether the user is available for communication via Google Talk. Merely visiting the page causes the user's browser to pull the Chatback badge from Google's servers, leaving behind footprints in their logs. Clicking the link can start an online conversation,

leaving open the eavesdropping and logging risks discussed in Chapter 5, “Communications.” Although this is a text-based service, similar risks exist via VoIP “call-me buttons” offered by companies such as Jajah, Jangle, Jaxtr, Tringme, and Grand Central.<sup>71,72</sup>



---

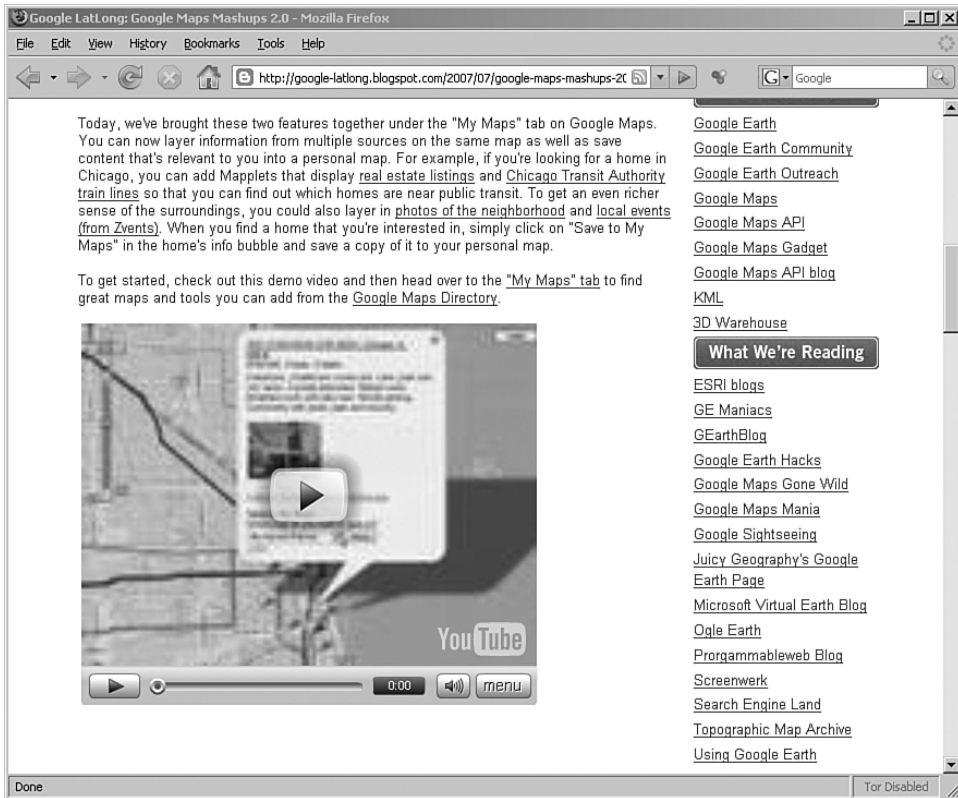
**Figure 7-11** Sample Google Chatback badge. Web authors place small snippets of Google-provided code in their web pages, and visitors to the page can see whether the author is available to chat via Google Talk.

## YOU TUBE VIDEOS

Embedding YouTube videos is an extremely popular practice by web authors (see Figure 7-12). When doing so, authors place code similar to the following in their web pages.<sup>73</sup>

```
<object height="350" width="425">
<param name="movie" value="http://www.youtube.com/v/KJukKpQDVLQ">
<param name="wmode" value="transparent">
<embed src="http://www.youtube.com/v/KJukKpQDVLQ"
  type="application/x-shockwave-flash" wmode="transparent"
  height="350" width="425">
</embed>
</object>
```

Notice that the code embeds a movie object pulled from Google’s servers. Again, users need only visit a page containing an embedded YouTube video to leave themselves open to tracking by Google, even if the page is run by a third party and there are no DoubleClick or AdSense advertisements.



**Figure 7-12** Example of a YouTube video embedded in a web page. When the image is merely displayed in the user's browser, that user can be immediately logged by YouTube.

## SEARCH ON YOUR WEB PAGE

Another common practice is for web authors to include a Google Search box on their site (see Figure 7-13). Although some visitors find this useful, it also facilitates the disclosure of search queries, as well as the user's IP address and the site he or she is visiting, to

Google. In some implementations, the disclosure takes place only when the user clicks Submit, as in the following code:<sup>74</sup>

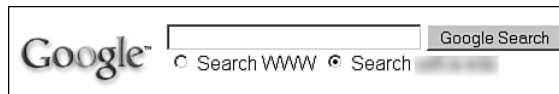
```
<form method="get" action="http://www.google.com/search">

<input type="text" name="q" size="31"
  maxlength="255" value="" />
<input type="submit" value="Google Search" />
<input type="radio" name="sitesearch" value="" />
  The Web
<input type="radio" name="sitesearch"
  value="askdave.taylor.com" checked /> Ask Dave Taylor<br />

</form>
```

However, note the Google logo in the image. If the webmaster includes the logo on the page, he or she can choose to download the image directly from Google; this immediately informs Google when someone visits a given site.

Google also offers AdSense for search, which helps webmasters earn revenue by creating a custom search engine for a site.<sup>75</sup> Along with customized search results, users see targeted advertisements.



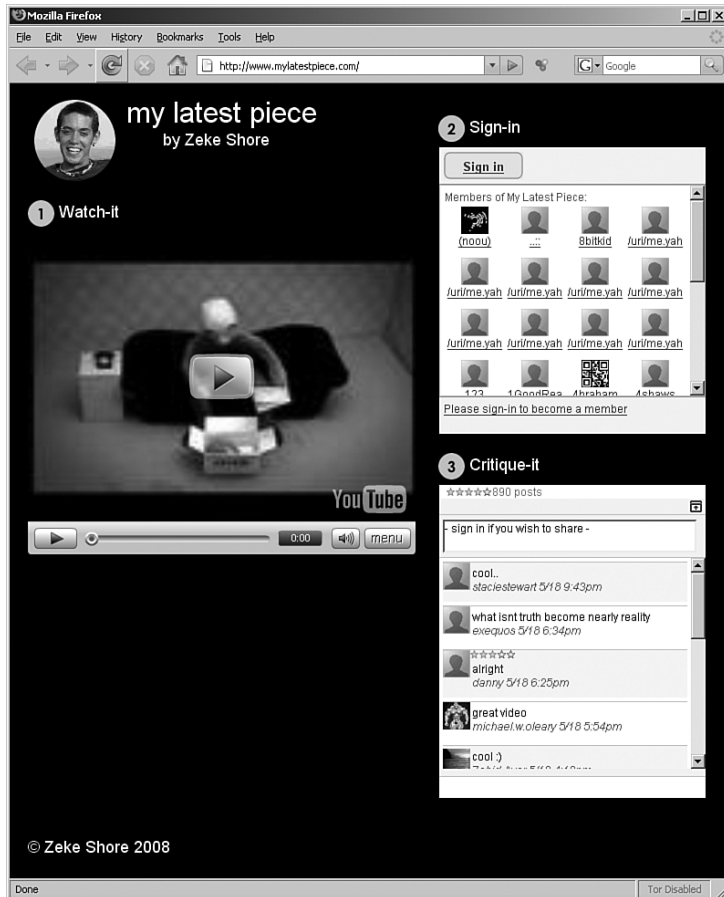
---

**Figure 7-13** Many web sites include a Google search field, which encourages users to disclose search terms and the site they are visiting.

## FRIEND CONNECT

Friend Connect is a new service offered by Google that enables web authors to add social networking facets to their sites by embedding small snippets of code. Google Friends Connect “offers a core set of social gadgets such as member management, message board, reviews, and picture sharing.”<sup>76</sup> Figure 7-14 shows a sample site provided by Google and illustrates several concerns with Friends Connect. Visitors to the site are offered the opportunity to sign in using their existing credentials, which uniquely identifies them to Google or one of several other participating services, including Yahoo! and AOL. The sites’ members, photos, and comments can be disclosed. In addition, because

this site includes an embedded YouTube video and two Friend Connect widgets, the user can be logged three times by Google's servers by merely visiting the site. Friend Connect is an interesting service that will likely be very popular. Therein lies the risk: Friends Connect and future generations of social networking applications will amplify user disclosure and facilitate uniquely identifying users.



**Figure 7-14** Google's new Friends Connect service enables web authors to add social networking functions to their sites.

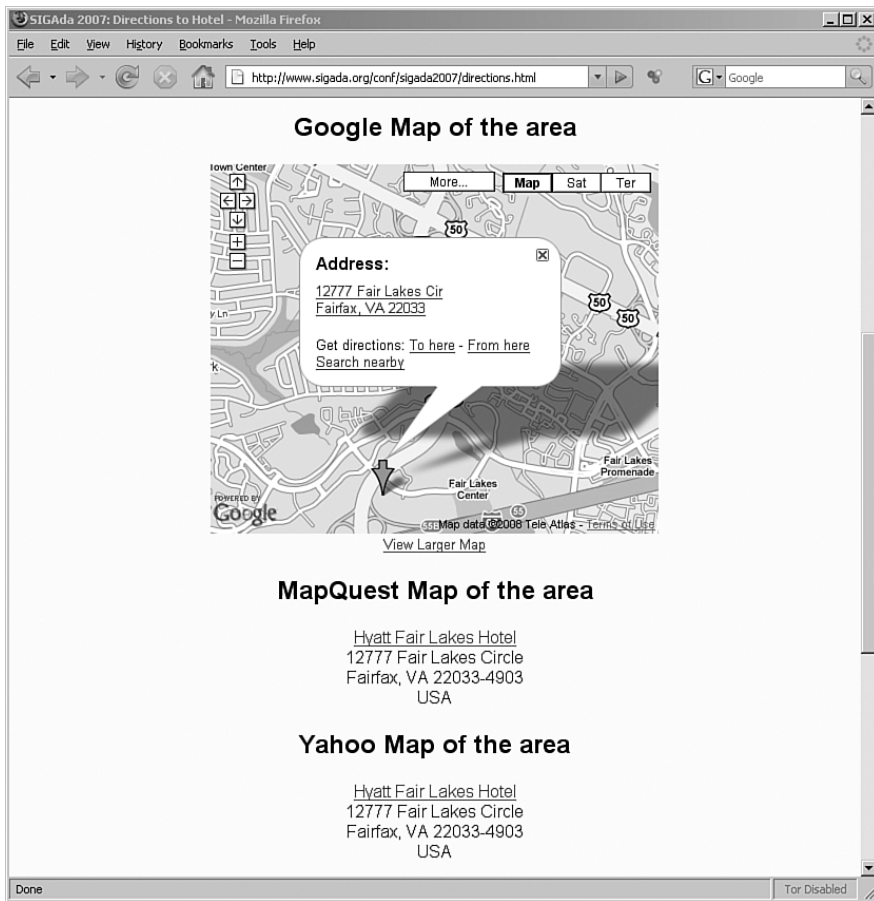


## EMBEDDED MAPS

Another common practice, and subsequent cross-site disclosure risk, is embedding maps within web pages. From hotels to tourist attractions, to business and social events, web authors rely upon third-party mapping services such as Google Maps and MapQuest to provide easy-to-use, interactive maps for their site's visitors. Unfortunately, the practice also informs the mapping service of the IP address of the visitor, HTTP cookies, the site the user came from, and a location he or she is interested in. For example, in Figure 7-15, a web author for an academic conference directly embedded a Google map into the conference web site.<sup>77</sup> Thus, every potential conference attendee who browses the conference's directions page immediately informs Google, and possibly Yahoo! and MapQuest, of his or her interest in the conference and probable attendance. With thousands, perhaps millions, of embedded maps in sites across the web, this practice greatly extends the cross-site visibility of large online companies such as Google and Yahoo!. The future of information-disclosure risks associated with embedded mapping is likely to worsen. Simple mapping is giving way to *mapplets* (or *mashups*), which combine mapping with virtually any type of location-based data (think homes for sale, local coffee shops, or driving ranges). The end result is a growth in the type and quantity of information disclosed via embedded maps and their progeny.

### NOTE

The term *mashup* applies to far more than just mapping. Mashups are a core Web 2.0 tenet and apply to web applications that combine more than one data source into single integrated tools. Beyond mapping, examples include combinations of images, videos, news, search, and shopping data.<sup>78</sup> Mashups increase information-disclosure risks because their use can share user information with many disparate mashup data source providers.



**Figure 7-15** Embedded maps in web pages immediately inform the mapping service of the user's visit to a given web page, as well as that user's interest in a specific area, as is the case for this academic conference.

## SUMMARY

Web browsing isn't a one-to-one conversation with a single web site. Instead, embedded content such as maps, images, videos, advertisements, web-analytics code, and social networking widgets immediately disclose each user's visit to a third party when that user merely views a page in his or her web browser. Web authors and webmasters gain a great

deal of value by embedding these “small snippets of code” in their web sites, such as gaining access to advertising revenue, free web-analytics reports, improved customer contact, and richer, more compelling web content. The true benefit is to the online companies, which gain a greatly increased field of view that isn’t constrained to their own properties, but instead encompasses a major swath of the Internet. As these companies innovate and field compelling new services, expect their field of view to increase further, as webmasters and web authors across the Internet embed new and better content. A key conclusion is that embedded third-party content forces the user to accept many different privacy policies from many different companies, most likely without even being aware of it. This creates a lowest common denominator effect of privacy policies; your real privacy in terms of visiting a web site is the equivalent of the worst policy of all the sites embedded there. This is a huge issue. Consider the MSNBC example earlier in the chapter. Most users might be aware that they fall under the MSNBC privacy policy, but they likely are not aware of the information being collected by the ten other companies providing embedded content, let alone the finer points of each of these companies’ privacy policies, if they even exist.

## ENDNOTES

1. Rory Cellan-Jones, “Web Creator Rejects Net Tracking,” BBC, 17 March 2008. <http://news.bbc.co.uk/2/hi/technology/7299875.stm>, last accessed 21 April 2008.
2. Sometimes third parties actually employ web bugs, typically in the form of a 1×1 transparent GIF.
3. For advertisements, these third-party servers are often called *central ad servers*.
4. Louise Story, “To Aim Ads, Web Is Keeping Closer Eye on You,” *The New York Times*, 10 March 2008. [www.nytimes.com/2008/03/10/technology/10privacy.html?\\_r=1&oref=slogin](http://www.nytimes.com/2008/03/10/technology/10privacy.html?_r=1&oref=slogin), last accessed 2 May 2008.
5. *Ibid.*
6. Ina Fried, “Microsoft Quietly Offering Ad-Funded Works,” CNET, 18 April 2008. [www.news.com/8301-13860\\_3-9922750-56.html?tag=nefd.top](http://www.news.com/8301-13860_3-9922750-56.html?tag=nefd.top), last accessed 10 May 2008.
7. Louise Story, “Company Will Monitor Phone Calls to Tailor Ads,” *The New York Times*, 24 September 2007. [www.nytimes.com/2007/09/24/business/media/24adcol.html?\\_r=3&ei=5065&en=5822f6a12e575488&ex=1191297600&partner=MYWAY&pagewanted=print&oref=slogin&oref=slogin&oref=slogin](http://www.nytimes.com/2007/09/24/business/media/24adcol.html?_r=3&ei=5065&en=5822f6a12e575488&ex=1191297600&partner=MYWAY&pagewanted=print&oref=slogin&oref=slogin&oref=slogin), last accessed 10 May 2008.

8. Catherine Holahan, "Google's DoubleClick Strategic Move," *BusinessWeek*, 14 April 2007. [www.businessweek.com/technology/content/apr2007/tc20070414\\_675511.htm](http://www.businessweek.com/technology/content/apr2007/tc20070414_675511.htm), last accessed 2 May 2008.
9. Mark Hendrickson, "Yahoo! Acquires Ad Network Blue Lithium," *TechCrunch*, 4 September 2007. [www.techcrunch.com/2007/09/04/yahoo-acquires-ad-network-bluelithium/](http://www.techcrunch.com/2007/09/04/yahoo-acquires-ad-network-bluelithium/), last accessed 10 May 2008.
10. Kevin Delaney, "Yahoo!–Google Pact May Be Close," *The Wall Street Journal*, 2 May 2008. [http://online.wsj.com/article/SB120968562237161201.html?mod=rss\\_whats\\_news\\_technology](http://online.wsj.com/article/SB120968562237161201.html?mod=rss_whats_news_technology), last accessed 2 May 2008.
11. Nate Anderson, "Study: Ads in Online Shows Work Better Than Ads on TV," *Ars Technica*, 26 December 2007. <http://arstechnica.com/news.ars/post/20071226-study-ads-in-online-video-work-better-than-ads-on-tv.html>, last accessed 17 May 2008.
12. Louise Story, "Where Every Ad Knows Your Name," *The New York Times Bits Blog*, 10 March 2008. <http://bits.blogs.nytimes.com/2008/03/10/where-every-ad-knows-your-name/?hp>, last accessed 10 May 2008.
13. "Microsoft to Test New Measure of Web Ads," Reuters, 25 February 2008. <http://www.reuters.com/article/internetNews/idUSWNAS219120080225?feedType=RSS&feedName=internetNews&pageNumber=1&virtualBrandChannel=0>, last accessed 11 May 2008.
14. David Legard, "IDC: Consolidation to Windows Won't Happen," *LinuxWorld*, 27 April 2004. [www.linuxworld.com.au/index.php/id;940707233;fp;2;fpid;1](http://www.linuxworld.com.au/index.php/id;940707233;fp;2;fpid;1), last accessed 2 July 2008.
15. "Google Ads Are a Free Speech Issue," *Slashdot*, 28 February 2007. <http://yro.slashdot.org/yro/07/02/28/0139222.shtml>, last accessed 10 May 2008.
16. This section focuses on information disclosure risks and does not cover Click Fraud, the illegitimate clicking of advertisements to generate revenue or charge an advertiser. See [www.businessweek.com/magazine/content/06\\_40/b4003001.htm](http://www.businessweek.com/magazine/content/06_40/b4003001.htm) for more information.
17. Jeffrey Graham, "Google's AdSense a Bonanza for Some Web Sites," *USA Today*, 10 March 2008. [www.usatoday.com/tech/news/2005-03-10-google-ads-usat\\_x.htm](http://www.usatoday.com/tech/news/2005-03-10-google-ads-usat_x.htm), last accessed 5 May 2008.
18. "What Is AdSense for Feeds?" Google AdSense Help Center. [www.google.com/adsense/support/bin/answer.py?hl=en&answer=20012](http://www.google.com/adsense/support/bin/answer.py?hl=en&answer=20012), last accessed 5 May 2008.
19. "Google AdSense for Search," Google AdSense, 2008. [www.google.com/adsense/static/en\\_US/WsOverview.html?hl=en\\_US](http://www.google.com/adsense/static/en_US/WsOverview.html?hl=en_US), last accessed 5 May 2008.

20. "AdSense for Mobile Content," GoogleAdSense, 2008. [www.google.com/adsense/www/en\\_US/mobile/](http://www.google.com/adsense/www/en_US/mobile/), last accessed 5 May 2008.
21. Shamim Samadi, "AdSense for Video Now in Beta," Google Blog, 21 February 2008. <http://googleblog.blogspot.com/2008/02/adsense-for-video-now-in-beta.html>, last accessed 5 May 2008.
22. "Google Testing AdSense Cost Per Action (CPA)," Search Engine Roundtable, 21 June 2006. [www.seroundtable.com/archives/003988.html](http://www.seroundtable.com/archives/003988.html), last accessed 17 May 2008.
23. Juan Carlos Perez, "Merchants Say eBay Ad Programs Drive Buyers Away," *ITWorld*, 8 October 2007. [www.itworld.com/Tech/2403/071008ebayads/](http://www.itworld.com/Tech/2403/071008ebayads/), last accessed 11 May 2008.
24. eBay also has its own contextual advertising network; see <http://slashdot.org/article.pl?sid=06/06/11/013216>.
25. Google AdSense sign-in page, Google AdSense, 2008. [www.google.com/adsense/login/en\\_US/?gsessionid=JjtK8kJKBuw](http://www.google.com/adsense/login/en_US/?gsessionid=JjtK8kJKBuw), last accessed 5 May 2008.
26. Mitch Ratcliffe, "How Google Falls: Unprofitable in 2009," ZDNet, Rational Rants Blog, 15 February 2007. <http://blogs.zdnet.com/Ratcliffe/?p=265>, last accessed 10 May 2008.
27. Mark Ward, "Searching for the Net's Big Thing," BBC News, 13 March 2006. <http://news.bbc.co.uk/1/hi/technology/4780648.stm>, last accessed 6 May 2008.
28. "Company Overview," Corporate Information, Google, 2008. [www.google.com/intl/en/corporate/index.html](http://www.google.com/intl/en/corporate/index.html), last accessed 6 May 2008.
29. Google is also experimenting with video ads on search result pages; see <http://bits.blogs.nytimes.com/2008/02/14/google-tests-video-ads-on-search-results-pages/>.
30. For an insightful overview of AdWords attacks, see StankDawg's Defcon 13 talk, available at [www.defcon.org/html/links/defcon-media-archives.html#dc\\_13](http://www.defcon.org/html/links/defcon-media-archives.html#dc_13).
31. "Google Partners Put You in Front of More Potential Customers," Google, 2008. <https://adwords.google.com/select/partner.html>, last accessed 6 May 2008.
32. "About AskEraser," Ask.com, 2008. <http://sp.ask.com/en/docs/about/askeraser.shtml>, last accessed 5 May 2008.
33. Thomas Claburn, "Google Keeps What Ask.com Erases," *Information Week*, 13 December 2007. [www.informationweek.com/news/internet/ebusiness/showArticle.jhtml;jsessionid=BQP3K401V1M1QQSNDLRSKH0CJUNN2JVN?articleID=204802233&\\_requestid=54316](http://www.informationweek.com/news/internet/ebusiness/showArticle.jhtml;jsessionid=BQP3K401V1M1QQSNDLRSKH0CJUNN2JVN?articleID=204802233&_requestid=54316), last accessed 5 May 2008.
34. Gwendolyn Mariano, "DoubleClick Able to Settle Privacy Suits," CNET News, 21 May 2002. [www.news.com/DoubleClick-able-to-settle-privacy-suits/2100-1023\\_3-919895.html](http://www.news.com/DoubleClick-able-to-settle-privacy-suits/2100-1023_3-919895.html), last accessed 16 May 2008.

35. Steve Lohr, "Google Deal Said to Bring U.S. Scrutiny," *The New York Times*, 29 May 2007. [www.nytimes.com/2007/05/29/technology/29antitrust.html?\\_r=2&ref=business&oref=slogin&oref=slogin](http://www.nytimes.com/2007/05/29/technology/29antitrust.html?_r=2&ref=business&oref=slogin&oref=slogin), last accessed 16 May 2008.
36. "E.U. Lobby Says Google, DoubleClick Merger Hurts Privacy," Reuters, 20 December 2007. [www.reuters.com/article/internetNews/idUSL2051059120071220](http://www.reuters.com/article/internetNews/idUSL2051059120071220), last accessed 16 May 2008.
37. Grant Gross, "FTC Approves Google/DoubleClick Deal," *Macworld*, 20 December 2007. [www.macworld.com/article/131204/2007/12/doubleclick.html](http://www.macworld.com/article/131204/2007/12/doubleclick.html), last accessed 16 May 2008.
38. "European Regulators Approve Google–DoubleClick Deal," CBC News, 11 March 2008. [www.cbc.ca/technology/story/2008/03/11/tech-google-doubleclick.html](http://www.cbc.ca/technology/story/2008/03/11/tech-google-doubleclick.html), last accessed 16 May 2008.
39. "Google Closes Acquisition of DoubleClick," Google Press Release, 11 March 2008. [www.google.com/intl/en/press/pressrel/20080311\\_doubleclick.html](http://www.google.com/intl/en/press/pressrel/20080311_doubleclick.html), last accessed 15 May 2008.
40. "DART Ad-Serving and Search Cookie Opt-Out," DoubleClick, 2008. [www.doubleclick.com/privacy/dart\\_adserving.aspx](http://www.doubleclick.com/privacy/dart_adserving.aspx), last accessed 16 May 2008.
41. Note that DoubleClick's opt-out cookie does not prevent ad targeting based on the user's operating system, Windows version, local time, or IP address.
42. Louise Story and Miguel Helft, "Google Buys an Online Ad Firm for \$3.1 Billion," *The New York Times*, 14 April 2007. [www.nytimes.com/2007/04/14/technology/14deal.html](http://www.nytimes.com/2007/04/14/technology/14deal.html), last accessed 16 May 2008.
43. Stefanie Olsen, "Privacy Concerns Dog Google–DoubleClick Deal," CNET News, 17 April 2007. [www.news.com/Privacy-concerns-dog-Google-DoubleClick-deal/2100-1024\\_3-6177029.html](http://www.news.com/Privacy-concerns-dog-Google-DoubleClick-deal/2100-1024_3-6177029.html), last accessed 16 May 2008.
44. Electronic Privacy Information Center Complaint with the Federal Trade Commission, 20 April 2007. [http://epic.org/privacy/ftc/google/epic\\_complaint.pdf](http://epic.org/privacy/ftc/google/epic_complaint.pdf), last accessed 16 May 2008.
45. Clint Boulton, "Meet Google: Search Giant, Monopolist Extraordinaire," Google Watch Blog, Eweek.com, 24 December 2007. [http://googlewatch.eweek.com/content/google\\_vs\\_microsoft/meet\\_google\\_search\\_giant\\_monopolist\\_extraordinaire.html](http://googlewatch.eweek.com/content/google_vs_microsoft/meet_google_search_giant_monopolist_extraordinaire.html), last accessed 16 May 2008.
46. Frank Rose, "Microsoft's Bid for Yahoo! Is All About Big-Budget Brand Advertising," *Wired Magazine*, 24 March 2008. [www.wired.com/techbiz/it/magazine/16-04/bz\\_microsoft\\_yahoo](http://www.wired.com/techbiz/it/magazine/16-04/bz_microsoft_yahoo), last accessed 15 May 2008.

47. StankDawg, "Hacking Google AdWords," Defcon 13, 2005. [www.defcon.org/images/defcon-13/dc13-presentations/DC\\_13-Stankdawg.pdf](http://www.defcon.org/images/defcon-13/dc13-presentations/DC_13-Stankdawg.pdf), last accessed 2 July 2008.
48. Brian Krebs, "Virus Writers Taint Google Ad Links," Security Fix Blog, WashingtonPost.com, 25 April 2007. [http://blog.washingtonpost.com/securityfix/2007/04/virus\\_writers\\_taint\\_google\\_ad.html](http://blog.washingtonpost.com/securityfix/2007/04/virus_writers_taint_google_ad.html), last accessed 4 May 2008.
49. Brian Krebs, "Hacked Ad Seen on MySpace Served Spyware to a Million," Security Fix Blog, WashingtonPost.com, 19 July 2006.
50. Betsy Schiffman, "Hackers Use Banner Ads on Major Sites to Hijack Your PC," Wired.com, 15 November 2007. <http://www.wired.com/techbiz/media/news/2007/11/doubleclick>, last accessed 4 May 2008.
51. Nonmalicious Flash applets hosted on web pages have also been shown to be vulnerable to attack. See Dan Goodin's "Serious Flash Vulns Menace at Least 10,000 Websites," [www.theregister.co.uk/2007/12/21/flash\\_vulnerability\\_menace/](http://www.theregister.co.uk/2007/12/21/flash_vulnerability_menace/)
52. Eric Auchard, "Adobe, Yahoo! test running ads inside PDF documents." Reuters, 28 November 2007. [www.reuters.com/article/marketsNews/idUKN2754715120071129?rpc=44&spr=true](http://www.reuters.com/article/marketsNews/idUKN2754715120071129?rpc=44&spr=true), last accessed 4 May 2008.
53. Robert Freeman, "I'm Feeling Lucky," Frequency X Blog, IBM Internet Security Systems, 29 April 2008. <http://blogs.iss.net/archive/FeelingLucky.html>, last accessed 17 May 2008.
54. Jakob Nielsen, "Banner Blindness: Old and New Findings," Alertbox blog, 20 August 2007. [www.useit.com/alertbox/banner-blindness.html](http://www.useit.com/alertbox/banner-blindness.html), last accessed 7 May 2008.
55. An interesting exception is NCSoft's game *City of Heroes*, which makes viewing in-game advertisements optional; see <http://yro.slashdot.org/article.pl?sid=08/04/06/0554230>.
56. Cade Metz, "Network Solutions Hijacks Customer Subdomains for Ad Fest," *The Register*, 11 April 2008. [www.theregister.co.uk/2008/04/11/network\\_solutions\\_sub\\_domain\\_parking/](http://www.theregister.co.uk/2008/04/11/network_solutions_sub_domain_parking/), last accessed 8 May 2008.
57. Henry Blodget, "Compete CEO: ISPs Sell Clickstreams for \$5 a Month," Seeking Alpha, 13 March 2007. <http://seekingalpha.com/article/29449-competite-ceo-isps-sell-clickstreams-for-5-a-month>, last accessed 8 May 2008.
58. "U.K. ISPs to Start Tracking to Serve You Ads," *TechDirt*, 18 February 2008. <http://techdirt.com/articles/20080218/024203278.shtml>, last accessed 8 May 2008.
59. For example, an Amazon cookie disclosed by the user's browser might have been associated with previous Amazon purchases and, hence, tied to shipping and billing information.

60. "Facebook Beacon," Facebook Business Solutions, 2008. [www.facebook.com/business/?beacon](http://www.facebook.com/business/?beacon), last accessed 18 May 2008.
61. Henry Blodget, "Facebook's 'Beacon' Infuriate Users, MoveOn," Silicon Alley Insider, 21 November 2007. [www.alleyinsider.com/2007/11/facebooks-beacon.html](http://www.alleyinsider.com/2007/11/facebooks-beacon.html), last accessed 18 May 2008.
62. Josh Catone, "Is Facebook Really Ruining Christmas?" ReadWriteWeb, 21 November 2007. [www.readwriteweb.com/archives/facebook\\_moveon\\_beacon\\_privacy.php](http://www.readwriteweb.com/archives/facebook_moveon_beacon_privacy.php), last accessed 2 July 2008.
63. Kristen Nicole, "Google Analytics Gets a Beautiful New Interface," Mashable.com, 8 May 2007. <http://mashable.com/2007/05/08/google-analytics/>, last accessed 15 May 2008.
64. The density of the code could also be seen as an attempt to decrease the size of the file, to improve response time.
65. Blogger Garrett Rogers provides an introductory analysis to the operation of `urchin.js` at <http://blogs.zdnet.com/Google/?p=39>.
66. The full script is more 13 pages long. You can view it at [www.google-analytics.com/urchin.js](http://www.google-analytics.com/urchin.js)
67. For additional information on Google Analytics cookies, see [www.customizegoogle.com/block-google-analytics-cookies.html](http://www.customizegoogle.com/block-google-analytics-cookies.html).
68. Garrett Rogers, "How Do I Know the Number of Google Analytics Accounts?" Googling Google ZDNet Blog, 28 November 2005. <http://blogs.zdnet.com/Google/?p=42>, last accessed 15 May 2008.
69. "Microsoft adCenter Analytics Registration," Microsoft Digital Advertising Solutions, 2008. <http://advertising.microsoft.com/advertising/adcenter-analytics-registration>, last accessed 15 May 2008.
70. Aurelie Pols, "Yahoo! Buys Indextools: 80% of the Functionality of Omniture for Free!" OX2 Web Analytics Blog, 8 April 2008. <http://webanalytics.ox2.eu/2008/04/15/yahoo-buys-indextools-80-of-the-functionality-of-omniture-for-free/>, last accessed 15 May 2008.
71. Nick Gonzalez, "TringMe: Phone Free Click to Call," *Tech Crunch*, 2 October 2007. [www.techcrunch.com/2007/10/02/tringme-phone-free-click-to-call/](http://www.techcrunch.com/2007/10/02/tringme-phone-free-click-to-call/), last accessed 17 May 2008.
72. Erik Schonfeld, "Google Talk Adds a Chatback Widget," *TechCrunch*, 26 February 2008. [www.techcrunch.com/2008/02/26/google-talk-adds-a-chatback-widget/](http://www.techcrunch.com/2008/02/26/google-talk-adds-a-chatback-widget/), last accessed 17 May 2008.



73. Thai Tran, “Google Maps Mashups 2.0,” Google Lat Long Blog, 11 July 2007. <http://google-latlong.blogspot.com/2007/07/google-maps-mashups-20.html>, last accessed 17 May 2008.
74. Dave Taylor, “How Can I Add a Google Search Box to My Web Site?” Ask Dave Taylor, 10 December 2004. [www.askdave.taylor.com/how\\_can\\_i\\_add\\_a\\_google\\_search\\_box\\_to\\_my\\_web\\_site.html](http://www.askdave.taylor.com/how_can_i_add_a_google_search_box_to_my_web_site.html), last accessed 19 May 2008.
75. “The Power of Google Search on Your Site,” Google AdSense, 2008. [www.google.com/adsense/www/en\\_US/afs/index.html](http://www.google.com/adsense/www/en_US/afs/index.html), last accessed 19 May 2008.
76. “More Info About Google Friend Connect,” Google Friend Connect BETA, 2008. [www.google.com/friendconnect/home/moreinfo](http://www.google.com/friendconnect/home/moreinfo), last accessed 18 May 2008.
77. SIGAda2007, Association for Computing Machinery Special Interest Group on Ada Conference web site, 10 October 2007. [www.sigada.org/conf/sigada2007/directions.html](http://www.sigada.org/conf/sigada2007/directions.html), last accessed 23 May 2008.
78. “Mashup (Web Application Hybrid),” Wikipedia, 2 July 2008. [http://en.wikipedia.org/wiki/Mashup\\_\(web\\_application\\_hybrid\)](http://en.wikipedia.org/wiki/Mashup_(web_application_hybrid)), last accessed 2 July 2008.

---

# Index

---

## A

Abdullah, Kulsoom, 271

accessing web server logs, 62-63

accidental data spills, 18

Accoona, 108

accounts

- GrandCentral accounts, 156-157

- registered user accounts, 78-80, 285

accuracy of mapping services, 196

Acquisti, Alessandro, 292

addresses

- IP addresses

  - city-level geolocation, 68-69

  - dynamic IP, 66

  - identifying end users with, 67-69

  - IPv4, 65-66

  - IPv6, 66

  - NAT (Network Address Translation), 66

  - purpose of, 65

  - static IP, 66

  - zip code-based geolocation, 69

network. *See* network addresses

AdSense, 212-215

advanced search operators, 113-114

advertising, 301

- AdSense, 212-215

- AdWords, 215

- cross-site tracking, 206-209

- DoubleClick, 216-217

- overview, 205-206, 210-212

- risks

  - affiliate services, 220

  - Chatback service, 225-226

  - Facebook Beacon, 221-222

  - Google Analytics, 222-225

  - hostile networks, 219-220

  - malicious ad serving, 218

  - malicious interfaces, 219

  - YouTube videos, 226

AdWords, 215

Adwords Traffic Estimator, 123

affiliate services, 220

Ahn, Luis von, 167

AI (artificial intelligence), future of, 301-302

Alerts, 115-117

algorithms, genetic, 301  
allintitle googling technique, 113  
allinurl googling technique, 113  
alternate surfing locations, 285  
Ames, Aldrich, 20  
analysis  
    computer analysis of communications, 166-167  
    imagery analysis, 192-194  
    traffic analysis, 162  
Analytics (Google), 222-225  
anonymization  
    anonymizing proxies, 278-280  
    data anonymization, 292-293  
anonymizing proxies, 278-280  
AOL Psycho, 260  
AOL search dataset  
    AOL User 2649647, 103-105  
    AOL User 3558174, 105-107  
    AOL User 1963201, 101-102  
    AOL User 789586, 102-103  
    AOL User 98280, 99-101  
AOL Stalker, 260  
AOL Searchlogs, 260  
API (Google), 117  
artificial intelligence (AI), future of, 301-302  
Asimov, Isaac, 301  
attacks. *See* data breaches  
audio, information flow and leakage, 40-41  
Australia, control of Internet access, 304  
avoiding registered accounts, 285

## B

Baidu, 109  
Baiduspider, 247  
Bank of America, loss of customer data, 18  
Bankston, Kevin, 197  
Basic Input/Output System (BIOS), 34

Battelle, John, 22  
Beacon (facebook), 221-222  
behavioral targeting, 88-89  
Berners-Lee, Tim, 97, 205  
BIOS (Basic Input/Output System), 34  
boot process, chain of trust, 34-36  
boot sector, 35  
botnets, 252  
breaches of data security  
    accidents, 18  
    challenges of controlling electronic data, 17-18  
    deliberate sharing with third parties, 18  
    legal requests for information, 21  
    malware and software vulnerabilities, 19  
    targeted attacks, 19-20  
    tension between shareholder needs and individual privacy, 21-22  
Brin, Sergey, 128  
browsers  
    header fields, 70-71  
    loading web pages, 60  
    requests, 60  
BrowserSpy web site, 71-72

## C

cache technique, 113  
CALEA (Communications Assistance for Law Enforcement Act), 45  
Callas, Jon, 276  
CAPTCHA (Completely Automated Public Turing Test to Tell Computers and Humans Apart), 167  
censorship, self-censorship, 127-128  
ChaCha, 108  
chaffing, 274-275  
Chatback service, 225-226  
Chernobyl virus (CIH), 35  
Chilling Effects web site, 291

- China
  - control of Internet access, 304
  - laws and regulatory requirements, 21
- chips, trust in, 34
- Chirac, Jacques, 108
- CIH (Chernobyl virus), 35
- Citigroup, accidental loss of back-up tapes, 18
- city-level geolocation, 68-69
- Cleveland Clinic, partnership with Google, 305
- click-through tracking, 84-85
- Clinton, Bill, 109
- Clusty, 305
- Code Green Networks, 270
- Comcast, 309
- communications
  - e-mail (Gmail), 140
    - forwarding, 143-144
    - history of, 139
    - HTML e-mail, 145
    - labeling, searching, and filtering, 144
    - online e-mail services, 145
    - out-of-office messages, 145
    - popularity of, 139-141
    - privacy policy, 141-143
    - spam, 145-146
    - table of security and privacy risks, 146-148
    - typographic errors, 146
  - groups, 151-155
  - mobile devices
    - archiving of messages, 162-163
    - computer analysis of communications, 166-167
    - convergence, 164
    - dependency, 161-162
    - eavesdropping, 163
    - emergent social networks, 165-166
    - filtering, 163
    - GrandCentral accounts, 156-157
    - language translation disclosures, 163
    - mobile and location-based searches, 159-161
    - overview, 155-156
    - text messaging, 157-158
    - traffic analysis, 162
  - overview, 139
  - voice, video, and instant messaging, 149-151
- Communications Assistance for Law Enforcement Act (CALEA), 45
- Completely Automated Public Turing Test to Tell Computers and Humans Apart (CAPTCHA), 167
- computer analysis of communications, 166-167
- connections across multiple platforms/networks, identifying, 91-93
- content filtering, 268-270
- controlling cookies, 266-267
- convergence, 164
- cookies, 72-76
  - controlling, 266-267
  - definition, 73
  - elements of, 74
  - example, 74
  - lifespan, 76
  - minimum browser support, 76
  - name/value pairs, 75
  - persistent cookies, 73
  - session cookies, 73
- Corona photo reconnaissance system, 193
- countermeasures. *See* protecting against information disclosure
- Cranor, Lorrie, 292
- cross-site tracking with advertising and embedded content, 206-209
- "Crowds: Anonymity for Web Transactions" (Reiter and Rubin), 283
- cryptography, 275-277
- Cryptome.org, 188
- Cuill, 247

## D

dangers of information disclosure, xv

data, definition of, 32

data anonymization, 292-293

data breaches

- accidents, 18
- challenges of controlling electronic data, 17-18
- deliberate sharing with third parties, 18
- legal requests for information, 21
- malware and software vulnerabilities, 19
- targeted attacks, 19-20
- tension between shareholder needs and individual privacy, 21-22

data generation, reducing/eliminating, 291

data permanence, 16-17

data retention, 16-17

- minimizing, 286
- reducing/eliminating, 291-292

death of Google, 306-307

deceit by mapping services, 196

dependency

- on Google, 22
- mobile devices and, 161-162

desktop search, 119-120

DHCP (Dynamic Host Configuration Protocol), 66

direction services. *See* mapping services

Disallow rule (Googlebot), 245

disclosure of information. *See* information disclosure

disclosure vectors, future of, 307-308

Discovery Engine, 247

DNS cache poisoning, 46

DNS (Domain Name System), 63-64

do-not-track lists, 286

- individual and organization policies, 287
- online company policies, 287-290

DogPile, 108

Domain Name System (DNS), 63-64

DoubleClick, 216-217

- Google acquisition of, 310

Dranor, Lorrie, 263

driving traffic, 122

Dynamic Host Configuration Protocol (DHCP), 66

dynamic IP, 66

## E

e-mail (Gmail), 140

- forwarding, 143-144

- Google Alerts, 115-117

- Google API, 117

- history of, 139

- HTML e-mail, 145

- labeling, searching, and filtering, 144

- online e-mail services, 145

- out-of-office messages, 145

- popularity of, 139-141

- privacy policy, 141-143

- spam, 145-146

- table of security and privacy risks, 146-148

- typographic errors, 146

eavesdropping, mobile devices and, 163

educating users about information disclosure

- overview, 260

- raising awareness, 260

- what users need to know, 262

Egelman, Serge, 292

electromagnetic (EM) radiation, information flow and leakage, 39-40

Electronic Communications Privacy Act, 17

Electronic Frontier Foundation, 290

Electronic Privacy Information Center, 290

eliminating

- data generation, 291

- data retention, 291-292

EM (electromagnetic) radiation, information flow and leakage, 39-40

**embedded content**

- AdSense, 212-215

- AdWords, 215

- cross-site tracking, 206-209

- DoubleClick, 216-217

- overview, 205-206, 210-212

**risks**

- affiliate services, 220

- Chatback service, 225-226

- embedded maps, 230

- Facebook Beacon, 221-222

- Friend Connect, 228-229

- Google Analytics, 222-225

- Google Search boxes on web pages, 227-228

- hostile networks, 219-220

- malicious ad serving, 218

- malicious interfaces, 219

- YouTube videos, 226

**emergent social networks, 165-166****encryption, 275-277****enterprise search, 120-122****errors, typographic errors in e-mail, 146****Exalead, 247****Expanding Your Horizons, 118****experts in Google's employ, 4-5****external media, information flow and**

- leakage, 41-42

**Eyeball-series.org, 188**

---

**F****Facebook Beacon, 221-222****Field Programmable Gate Arrays (FPGAs), 35****fields, browser header fields, 70-71****files, robots.txt, 245-246****filetype technique, 113****filtering**

- content, 268-270

- e-mail, 144

- mobile devices and, 163

**fingerprinting, 128-129****firewalls (NAT), 284-285****fitness function, 301****flow of information**

- boot process chain of trust, 34-36

- overview, 31-33

- on networks, 43

- Google versus ISPs, 46-49

- Internet data communication, 43-46

- on PCs

- comparison of information disclosure

- vectors, 42

- EM(electromagnetic) radiation, 39-40

- external media, 41-42

- holistic approach, 36-37

- networks, 37-38

- peripherals, 38-39

- power lines, 41

- sound, 40-41

- trust in hardware/software, 33-34

**forwarding e-mail, 143-144****Foundation series (Asimov), 301****FoxyProxy, 280****FPGAs (Field Programmable Gate Arrays), 35****Friend Connect, 228-229****future**

- of Google

- death of Google, 306-307

- disclosure vectors, 307-308

- government collaboration, 306

- services, 305-306

- of web applications

- office services, 308-309

- sensors and RFID, 309

- Web 2.0, 309

- of web-based information disclosure

- advertising, 301

- artificial intelligence, 301-302

- network control, 302-305

- overview, 299-300

## G

- Garfinkel, Simson, 263
- Garlinghouse, Brad, 166
- genetic algorithms, 301
- geographic relationships, linking user classes
  - via, 184-186
- geolocation
  - city-level geolocation, 68-69
  - zip code-based geolocation, 69
- geotagging, 191
- Gizmo, 149
- global presence of Google, 5-6
- Gmail
  - e-mail forwarding, 143-144
  - e-mail labeling, searching, and filtering, 144
  - HTML e-mail, 145
  - out-of-office messages, 145
  - popularity of, 140-141
  - privacy policy, 141-143
  - spam, 145-146
  - table of security and privacy risks, 146-148
  - typographic errors, 146
- Google
  - acquisition of DoubleClick, 310
  - AdSense, 212-215
  - AdWords, 215
  - Alerts, 115-117
  - Analytics, 222-225
  - API, 117
  - Chatback service, 225-226
  - click-through tracking, 84-85
  - compared to ISPs, 46-49
  - dependence on, 22
  - DoubleClick, 216-217
  - Earth. *See* mapping services
  - experts in Google's employ, 4-5
  - Friend Connect, 228-229
  - future of
    - death of Google, 306-307
    - government collaboration, 306
    - services, 305-306
  - global and multicultural presence, 5-6
  - Groups, 151-155
  - growth of, 3-4
  - hacking, 251-252
  - information disclosure. *See*
    - information disclosure
  - Maps. *See* mapping services
  - market leadership, xv
  - popularity of, 7-8
  - privacy policy
    - data sharing, 18
    - overview, 2
    - tension between shareholder needs and
      - individual privacy, 21-22
  - resources and power, 4-7
  - revenue, 4
  - Search boxes, 227-228
  - searches. *See* searches
  - Sightseeing site, 188
  - Talk, 149-151
  - as target for attach, 19
  - Toolbar, 116
- Googlebot
  - definition, 239
  - Disallow rule, 245
  - entries in web server logs, 242-244
  - how it works, 240-242
  - overview, 239-240
  - risks
    - Google hacking, 251-252
    - overview, 247-248
    - placing sensitive information on web, 250-251
    - potential attacks with, 252-253
    - spoofing, 248-249
  - robots.txt file, 245-246
  - User-Agent rule, 245

googling  
  definition of, 1  
  technique, 113

government collaboration, 306

GPG, 276

GrandCentral accounts, 156-157

GreenMaven, 118

groups, 151-155

growth

  of Google, 3-4

  of Internet, 3

## H

hacking, 113-114, 251-252

Hanssen, Robert, 20

headers, browser header fields, 70-71

hostile networks, 219-220

Howe, Daniel, 275

HTML (Hypertext Markup Language), 60

  e-mail, 145

  META tag, 245

HTTP (HyperText Transfer Protocol), 60

  referer data, 76-78

## I

iChat, 149

imagery analysis, 192-194

imagery. *See* mapping services

importance of web-based information

  disclosure, 309-311

information, definition of, 32

information disclosure

  with advertising and embedded content

    AdSense, 212-215

    AdWords, 215

    affiliate services, 220

    Chatback service, 225-226

  cross-site tracking, 206-209

  DoubleClick, 216-217

  embedded maps, 230

  Facebook Beacon, 221-222

  Friend Connect, 228-229

  Google Analytics, 222-225

  Google Search boxes on web pages, 227-228

  hostile networks, 219-220

  malicious ad serving, 218

  malicious interfaces, 219

  overview, 205-206, 210-212

  YouTube videos, 226

behavioral targeting, 88-89

browser header fields, 70-71

connections across multiple platforms/networks,

  identifying, 91-93

cookies, 72-76

  definition, 73

  elements of, 74

  example, 74

  lifespan, 76

  minimum browser support, 76

  name/value pairs, 75

  persistent cookies, 73

  session cookies, 73

dangers of, xv

data retention and permanence, 16-17

with e-mail (Gmail)

  e-mail forwarding, 143-144

  HTML e-mail, 145

  labeling, searching, and filtering, 144

  online e-mail services, 145

  out-of-office messages, 145

  overview, 139-141

  privacy policy, 141-143

  spam, 145-146

  table of security and privacy risks, 146-148

  typographic errors, 146

future of

  advertising, 301

  artificial intelligence, 301-302

  disclosure vectors, 307-308



- Google, 305-307
  - network control, 302-305
  - overview, 299-300
  - Web applications, 308-309
- Google tools/services and information
  - disclosure examples, 9-12
- with groups, 151-155
- HTTPReferer data, 76-78
- importance of threat, 309-311
- IP addresses
  - city-level geolocation, 68-69
  - dynamic IP, 66
  - identifying end users with, 67-69
  - IPv4, 65-66
  - IPv6, 66
  - NAT (Network Address Translation), 66
  - purpose of, 65
  - static IP, 66
  - zip code-based geolocation, 69
- with mapping services
  - accuracy and deceit, 196
  - benefits of, 177
  - common actions and information
    - disclosed, 198-199
  - imagery analysis, 192-194
  - interacting with, 179-182
  - linking user classes via geographic relationships, 184-186
  - list of services, 198
  - location disclosure, 179
  - mashups, 186-187
  - obscuring sensitive locations, 194-195
  - overview, 177-178
  - personalized maps, 184
  - risks of combining search and mapping, 183-184
  - satellite imagery threats, 188-192
  - street-level views, 196-197
- with mobile devices
  - archiving of messages, 162-163
  - computer analysis of communications, 166-167
  - convergence, 164
  - dependency, 161-162
  - eavesdropping, 163
  - emergent social networks, 165-166
  - filtering, 163
  - GrandCentral accounts, 156-157
  - language translation disclosures, 163
  - mobile and location-based searches, 159-161
  - overview, 155-156
  - text messaging, 157-158
  - traffic analysis, 162
- protecting against. *See* protecting against
  - information disclosure
- registered user accounts, 78-80
- risk assessment, 12-13
- scenarios, 13-15
- searches
  - advanced search operators, 113-114
  - AOL User 2649647 dataset, 103-105
  - AOL User 3558174 dataset, 105-107
  - AOL User 1963201 dataset, 101-102
  - AOL User 789586 dataset, 102-103
  - AOL User 98280 dataset, 99-101
  - desktop search, 119-120
  - driving traffic, 122
  - enterprise search, 120-122
  - exploitation by malicious software, 124-126
  - fingerprinting, 128-129
  - Google Alerts, 115-117
  - Google API, 117
  - Google hacking, 113-114
  - Google search services, 109-111
  - overview, 97-98
  - search box and related applications, 112-113
  - search engine optimization (SEO), 123-124
  - search engines, 108-109
  - search queries, 98-99
  - search queries of others, 126-127
  - self-censorship, 127-128
  - site-based search, 118
- Sobiesk Information Disclosure Metric (SIDM), 8

table of common web-based information  
     disclosure, 87-88  
 tension between shareholder needs and  
     individual privacy, 21-22  
 trend toward web-based applications, 15-16  
 trust, 17  
 uniqueness, 89-91  
 value of, xiii-xiv  
 with voice, video, and instant  
     messaging, 149-151  
 web server logs  
     access to, 62-63  
     contents of, 64-65  
 web site navigation  
     inter-web site navigation, 81-85  
     intra-web site navigation, 85-87  
 information flow. *See* flow of information  
 information leakage  
     on networks, 43  
         Google versus ISPs, 46-49  
         Internet data communication, 43-46  
     on PCs  
         comparison of information disclosure  
             vectors, 42  
         EM (electromagnetic) radiation, 39-40  
         external media, 41-42  
         holistic approach, 36-37  
         networks, 37-38  
         peripherals, 38-39  
         power lines, 41  
         sound, 40-41  
 insider threats, 20  
 instant messaging, 149-151  
 inter-web site navigation, 81-85  
 interacting with mapping services, 179-182  
 interfaces, malicious, 219  
 Internet  
     amount of data of, 3  
     growth of, 3  
     information flow and leakage, 43-46

Internet backbone providers, 46-49  
 Internet Protocol. *See* IP addresses  
 Internet service providers (ISPs), xv  
     compared to Google, 46-49  
     and network control, 303-304  
 intitle googling technique, 113  
 intra-web site navigation, 85-87  
 inurl googling technique, 113  
 IP addresses  
     city-level geolocation, 68-69  
     dynamic IP, 66  
     identifying end users with, 67-69  
     IPv4, 65-66  
     IPv6, 66  
     NAT (Network Address Translation), 66  
     purpose of, 65  
     static IP, 66  
     zip code-based geolocation, 69  
 IPv4, 65-66  
 IPv6, 66  
 ISPs (Internet service providers), xv  
     compared to Google, 46-49  
     and network control, 303-304

## J-K

Japan, control of Internet access, 304  
 Java Anonymous Proxy, 284  
  
 kernels, 35-36  
 Kessinger, Kevin, 18  
 -keyword googling technique, 113  
 knowledge, definition of, 32

## L

labeling e-mail, 144  
 language translation disclosures, 163  
 lawmakers, petitioning, 290

- leakage of information
  - on networks, 43-49
  - on PCs
    - comparison of information disclosure vectors, 42
    - EM (electromagnetic) radiation, 39-40
    - external media, 41-42
    - holistic approach, 36-37
    - networks, 37-38
    - peripherals, 38-39
    - power lines, 41
    - sound, 40-41
- legal requests for information, 21
- legislation
  - CALEA (Communications Assistance for Law Enforcement Act), 45
  - Electronic Communications Privacy Act, 17
- Lexxe, 305
- lifespan of cookies, 76
- link googling technique, 113
- linking user classes via geographic relationships, 184-186
- lists, do-not-track lists, 286
  - individual and organization policies, 287
  - online company policies, 287-290
- LivePlasma, 108
- loading web pages, 60
- location, disclosing to mapping services, 179
- logs, web server logs
  - access to, 62-63
  - contents of, 64-65
  - Googlebot entries, 242-244
- Lyman, Peter, 3
- M**
- MacWorld, 118
- Malaysia, control of Internet access, 304
- malicious ad serving, 218
- malicious interfaces, 219
- malware, 19, 124-126
- mapping services
  - accuracy and deceit, 196
  - benefits of, 177
  - common actions and information disclosed, 198-199
  - embedded maps, 230
  - imagery analysis, 192-194
  - interacting with, 179-182
  - linking user classes via geographic relationships, 184-186
  - list of, 198
  - location disclosure, 179
  - mashups, 186-187
  - obscuring sensitive locations, 194-195
  - overview, 177-178
  - personalized maps, 184
  - risks of combining search and mapping, 183-184
  - satellite imagery threats, 188-192
  - street-level views, 196-197
- MapQuest. *See* mapping services
- mashups, 186-187
- MBR (Master Boot Record), 35
- media, information flow and leakage, 41-42
- message externals, 162
- messages
  - archiving, 162-163
  - message externals, 162
  - text messaging, 157-158
  - voice, video, and instant messaging, 149-151
- META tag, 245
- minimizing data retention, 286
- Mixmaster, 284
- Mixminion, 284
- mobile devices
  - GrandCentral accounts, 156-157
  - mobile and location-based searches, 159-161
  - overview, 155-156

risks  
 archiving of messages, 162-163  
 computer analysis of communications,  
   166-167  
 convergence, 164  
 dependency, 161-162  
 eavesdropping, 163  
 emergent social networks, 165-166  
 filtering, 163  
 language translation disclosures, 163  
 traffic analysis, 162  
 text messaging, 157-158  
 monitoring, self-monitoring, 270-273  
 motherboards, trust in, 34  
 Motion Picture Association of America  
   (MPAA), 304  
 Mozilla Online, 109  
 MPAA (Motion Picture Association of  
   America), 304  
 Ms. Dewey, 108  
 msnbot, 247  
 multicultural presence of Google, 5-6

## N

n..m googling technique, 113  
 NAT (Network Address Translation), 66, 284-285  
 navigation between web sites  
   inter-web site navigation, 81-85  
   intra-web site navigation, 85-87  
 Network Address Translation (NAT), 66, 284-285  
 network addresses, protecting  
   anonymizing proxies, 278-280  
   NATfirewalls, 284-285  
   Tor, 280-284  
 network control, 302-305  
 networks  
   emergent social networks, 165-166  
   hostile networks, 219-220

information flow and leakage, 37-38, 43  
   Google versus ISPs, 46-49  
   Internet data communication, 43-46  
 network addresses, protecting  
   anonymizing proxies, 278-280  
   NATfirewalls, 284-285  
   Tor, 280-284  
 Nissenbaum, Helen, 275

## O

obscuring sensitive locations in mapping  
   services, 194-195  
 office services, Web-based, 308-309  
 online e-mail services, 145  
 operators  
   advanced search operators, 113-114  
   Google hacking, 113-114  
 out-of-office messages, 145  
 Overture, 123, 301

## P

Page, Larry, 301  
 “The PageRank Citation Ranking: Bringing Order  
   to the Web” (paper), 241  
 paying for privacy, 292  
 PCs, information flow and leakage  
   boot process, 34-36  
   comparison of information disclosure vectors, 42  
   EM (electromagnetic) radiation, 39-40  
   external media, 41-42  
   holistic approach, 36-37  
   networks, 37-38  
   peripherals, 38-39  
   power lines, 41  
   sound, 40-41  
 PDAs (personal digital assistants). *See*  
   mobile devices

- peripherals, information flow and leakage, 38-39
- permanence of data, 16-17
- persistent cookies, 73
- personal digital assistants (PDAs). *See*
  - mobile devices
- personalization, 80
- personalized maps, 184
- petitioning law and policy makers, 290
- PGP, 276
- phone numbers, GrandCentral accounts, 156-157
- policies, Gmail privacy policy, 141-143
- polymakers, petitioning, 290
- popularity of Google, 7-8
- power analysis attacks, 41
- power of Google, 4-7
- power lines, information flow and leakage, 41
- prefetching, 116
- privacy organizations, 290
- privacy policy
  - Gmail, 141-143
  - Google
    - data sharing, 18
    - overview, 2
  - tension between shareholder needs and individual privacy, 21-22
- privacy. *See* information disclosure
- Privoxy, 268
- Proofpoint, 270
- protecting against information disclosure
  - avoiding registered accounts, 285
  - content filtering, 268-270
  - controlling cookies, 266-267
  - data anonymization, 292-293
  - do-not-track lists, 286
    - individual and organization policies, 287
    - online company policies, 287-290
  - educating users
    - overview, 260
    - raising awareness, 260
    - what users need to know, 262
- encryption, 275-277
- minimizing data retention, 286
- network address protection
  - alternate surfing locations, 285
  - anonymizing proxies, 278-280
  - NAT firewalls, 284-285
  - Tor, 280-284
- overview, 259-260
- paying for privacy, 292
- petitioning law and policy makers, 290
- reducing/eliminating data generation, 291
- reducing/eliminating data retention, 291-292
- search term chaffing, 274-275
- self-monitoring, 270-273
- supporting privacy organizations, 290
- usable security, 263-265
- protocols
  - DHCP (Dynamic Host Configuration Protocol), 66
  - DNS (Domain Name System), 63-64
  - HTML (Hypertext Markup Language), 60
  - HTTP (HyperText Transfer Protocol), 60
  - IP addresses
    - city-level geolocation, 68-69
    - dynamic IP, 66
    - identifying end users with, 67-69
    - IPv4, 65-66
    - IPv6, 66
    - NAT (Network Address Translation), 66, 284-285
    - purpose of, 65
    - static IP, 66
    - zip code-based geolocation, 69
- proxies, anonymizing, 278-280
- PSTN (Public Switched Telephone Network), 149

## Q

quantum computing, 302  
queries, search, 98-99

- AOL User 2649647 dataset, 103-105
- AOL User 3558174 dataset, 105-107
- AOL User 1963201 dataset, 101-102
- AOL User 789586 dataset, 102-103
- AOL User 98280 dataset, 99-101
- search engines, 108-109
- search queries of others, 126-127

Quintura, 108

## R

raising awareness of information disclosure, 260  
Readings in Database Systems (Hellerstein and Stonebraker), 122  
Recording Industry Association of American (RIAA), 304  
reducing

- data generation, 291
- data retention, 291-292

referrer data (HTTP), 76-78  
registered user accounts, 78-80, 285  
Reiter, Michael, 283  
related googling technique, 113  
resources at Google's disposal, 4-7  
retention of data, 16-17  
RFID tags, 36, 309  
RIAA (Recording Industry Association of American), 304  
risk assessment, 12-13  
robots.txt file, 245-246  
Rollyo, 305  
root servers, 46  
Rubin, Avi, 283

## S

satellite imagery, 188-192  
Schmidt, Eric, 4, 92  
*The Search* (Battelle), 22  
search box, 112-113  
search engine optimization (SEO), 123-124  
search engines, 108-109  
search queries, 98-99

- AOL User 2649647 dataset, 103-105
- AOL User 3558174 dataset, 105-107
- AOL User 1963201 dataset, 101-102
- AOL User 789586 dataset, 102-103
- AOL User 98280 dataset, 99-101
- search engines, 108-109
- search queries of others, 126-127

search term chaffing, 274-275  
Search Wikia, 108  
searches

- advanced search operators, 113-114
- desktop search, 119-120
- e-mail, 144
- enterprise search, 120-122
- Google Alerts, 115-117
- Google API, 117
- Google hacking, 113-114
- Google search services, 109-111
- mapping services. *See* mapping services
- mobile and location-based searches, 159-161
- overview, 97-98

risks

- driving traffic, 122
- exploitation by malicious software, 124-126
- fingerprinting, 128-129
- search engine optimization (SEO), 123-124
- search queries of others, 126-127
- self-censorship, 127-128

search box and related applications, 112-113

- search queries, 98-99
  - AOL User 2649647 dataset, 103-105
  - AOL User 3558174 dataset, 105-107
  - AOL User 1963201 dataset, 101-102
  - AOL User 789586 dataset, 102-103
  - AOL User 98280 dataset, 99-101
  - search engines, 108-109
  - search queries of others, 126-127
- search term chaffing, 274-275
- site-based search, 118
- Second Life*, 302
- security. *See also* information disclosure
  - data breaches
    - accidents, 18
    - challenges of controlling electronic data, 17-18
    - deliberate sharing with third parties, 18
    - legal requests for information, 21
    - malware and software vulnerabilities, 19
    - targeted attacks, 19-20
    - tension between shareholder needs and individual privacy, 21-22
  - overview, 1-3
- Security and Usability: Designing Secure Systems That People Can Use* (Cranor and Garfinkel), 263
- self-censorship, 127-128
- self-monitoring, 270-273
- semantic disclosures
  - overview, 78
  - registered user accounts, 78-80
  - web site navigation
    - inter-web site navigation, 81-85
    - intra-web site navigation, 85-87
- sensitive information, posting on web via webbots, 250-251
- sensors, 309
- SEO (search engine optimization), 123-124
- servers
  - root servers, 46
  - web. *See* web servers
- services
  - affiliate services, 220
  - Google search services, 109-111
- session cookies, 73
- sharing data with third parties, 18
- SIDM (Sobiesk Information Disclosure Metric, 8
- site-based search, 118
- site googling technique, 113
- Skype, 149
- Sobiesk Information Disclosure Metric (SIDM), 8
- social networks, emergent, 165-166
- software vulnerabilities, 19
- sound, information flow and leakage, 40-41
- SOUPS (Symposium on Usable Privacy and Security), 265
- spam, 145-146
- Speedy Spider, 247
- spoofing Googlebot, 248-249
- static IP, 66
- Sterling, Bruce, 309
- Story, Louise, 210
- street-level views (mapping services), 196-197
- supporting privacy organizations, 290
- Swicki, 305
- SwitchProxy, 280
- Symposium on Usable Privacy and Security (SOUPS), 265
- system calls, 35

## T

- Tacoda, 88
- tags, META, 245
- targeted attacks, 19-20
- Teoma, 247
- text messaging, 157-158
- third parties, sharing data with, 18

Time Warner, loss of customer data, 18

Tomlinson, Ray, 139

Toolbar, 116

Tor, 280-284

Torbutton, 281

TorCheck, 281

Torpark, 281

Tournachon, Gaspar, Felix, 193

tracking, cross-site tracking, 206-209

TrackMeNot, 275

traffic

- driving with searches, 122

- traffic analysis, 162, 303

translation disclosures, 163

TriWest, loss of data, 18

TrueCrypt, 277

trust

- boot process chain of trust, 34-36

- in hardware/software, 33-34

- information disclosure and, 17

Tsai, Janice, 292

Turnbull, Alex, 188

Turnbull, James, 188

Tygar, J. D., 263

typographic errors, 146

Tzu, Sun, 259

## U

uniform resource locator (URL), 59-60

uniqueness, 89-91

United States, proposals for Internet filtering, 304

URL (uniform resource locator), 59-60

usable security, 263-265

Usable Security Blog, 265

user accounts, registered user accounts, 78-80

user space, 35

User-Agent rule (Googlebot), 245

users

- educating about information disclosure

  - overview, 260

  - raising awareness, 260

  - what users need to know, 262

- linking user classes via geographic

  - relationships, 184-186

## V

value of online information, xiii-xiv

Varian, Hal, 3

Vericept, 270

VHDL (Very High Speed Integrated Circuit  
Hardware Description Language), 35

Vidalia, 281

video messaging, 149-151

videos on YouTube, 226

voice messaging, 149-151

Voice over IP (VoIP), 149

Vonage, 149

## W

Web 2.0, 309

web-analytic services

- cross-site tracking, 206-209

- Google Analytics, 222-225

- overview, 205-206

web applications, future of

- office services, 308-309

- sensors and RFID, 309

- Web 2.0, 309

web-based applications, trend toward, 15-16

web-based information disclosure. *See*  
information disclosure



- web browsers
    - load web pages, 60
    - requests, 60
  - web crawlers. *See* webbots
  - web pages
    - browser requests, 60
    - DNS (Domain Name System), 63-64
    - HTML (Hypertext Markup Language), 60
    - HTTP (HyperText Transfer Protocol), 60
    - loading, 60
    - URL (uniform resource locator), 59-60
  - web servers
    - cookies, 72-76
      - definition, 73
      - elements of, 74
      - example, 74
      - lifespan, 76
      - minimum browser support, 76
      - name/value pairs, 75
      - persistent cookies, 73
      - session cookies, 73
    - definition, 62
    - logs
      - access to, 62-63
      - contents of, 64-65
      - Googlebot entries, 242-244
  - web sites
    - BrowserSpy, 71-72
    - navigation, 81
      - inter-web site navigation, 81-85
      - intra-web site navigation, 85-87
  - web spiders. *See* webbots
  - webbots
    - definition, 239
    - Googlebot. *See* Googlebot
    - popular search webbots, 247
    - risks
      - Google hacking, 251-252
      - overview, 247-248
      - placing sensitive information on web, 250-251
      - potential attacks with, 252-253
      - spoofing, 248-249
  - Webbots, Spiders, and Screen Scrapers: A Guide to Developing Internet Agents with PHP/CURL* (Schrenk), 242
  - Webmon, 247
  - Websense, 269-270
  - Whitten, Alma, 263
  - "Why Johnny Can't Encrypt" (Whitten and Tygar), 263
  - Wikimapia, 190
  - Wink, 305
  - wisdom, definition of, 32
  - Wright, Wilbur, 193
- ## X-Y-Z
- Yahoo! webbot, 247
  - Yee, Ka-Ping, 265
  - Yeti, 247
  - Young, John, 188
  - YouTube videos, 226
  - Zalewski, Michal, 253
  - zip code-based geolocation, 69