

1
2
3
4
5
6
7
8
9
10
11
12

Academic Data Collection in Electronic Environments: Defining Acceptable Use of Internet Resources¹

Gove N. Allen

Freeman School of Business
Tulane University
New Orleans, LA 70118
gallen@tulane.edu

Dan L. Burk

Law School
University of Minnesota
Minneapolis, MN 55455
burkx006@umn.edu

Gordon B. Davis

Carlson School of Management
University of Minnesota
Minneapolis, MN 55455
gdavis@csom.umn.edu

¹An earlier version of this paper was presented at the 4th annual meeting of the Association of Internet

Academic Data Collection in Electronic Environments: Defining Acceptable Use of Internet Resources

Gove N. Allen received his Ph.D. from the University of Minnesota in 2001 and currently serves as an assistant professor of e-business and information systems at Tulane University's A. B. Freeman School of Business. He holds bachelors and masters degrees in Accountancy from Brigham Young University. Dr. Allen has consulted major corporations on the implementation of information technology including Sony, AT&T, Sprint, Hewlett Packard, Micron, Intel, 3M, American Express, and NASA. More recently, he has made presentations at academic gatherings regarding automated data collection in the North America, Europe, and Australia. His continuing research utilizes various configurations of automated Internet data collection agents.

Gordon B. Davis is the Honeywell Professor of Management Information Systems, Emeritus at the Carlson School of Management at the University of Minnesota.. He is a pioneer in management information systems research and education. Professor Davis is the author of 23 texts and over 200 articles in management information systems, data processing, programming, and EDP auditing. His book, *Management Information Systems: Conceptual Foundations, Structure and Development*, 1974, Second Edition (with M. Olson), McGraw-Hill Book Company, 1985, has been ranked as a classic in the field

Dan L. Burk is the Oppenheimer, Wolf and Donnelly Professor of Law at the University of Minnesota, where he teaches courses in Patent Law, Copyright, and Biotechnology Law. An authority on issues related to high technology, he is the author of numerous papers on the legal and societal impact of new technologies, including articles on scientific misconduct, regulation of biotechnology, and the intellectual property implications of global computer networks. He is perhaps best known for his work in the area of "cyberlaw," where he has been a leading figure in the debates surrounding Internet jurisdiction, trespass to computers, and the deployment of digital rights management systems.

1 **Academic Data Collection in Electronic Environments:**
2 **Defining Acceptable Use of Internet Resources**
3

4 **Abstract**

5 Academic researchers access commercial websites to collect research data. This research
6 practice is likely to increase. Is this appropriate? Is this legal? Such commercial websites are
7 maintained to achieve business objectives; research access uses site resources for other purposes.
8 Website administrators may therefore deem academic data collection inappropriate. Is there a
9 process to make research access more open and acceptable to website owners and
10 administrators? These are significant issues. This article clarifies the problems and suggests
11 possible approaches to handle the issues with sensitivity and openness.

12 Research access to commercial websites may be manual (using a standard web browser)
13 or automated (using automated data collection agents). These approaches have different effects
14 on websites. Researchers using manual access tend to make a limited number of page requests
15 because manual access is costly to perform. Researchers using automated access methods can
16 request large numbers of pages at a low cost. Therefore, website administrators tend to view
17 manual access and automated access very differently.

18 Because of the number of accesses and non-business purpose, automated research
19 requests for data are sometimes blocked by site administration using a variety of means (both
20 technological and legal). This paper details the pertinent legal issues including trespass,
21 copyright violation, and breach of contract. It also explains the nature of express and implied
22 consent by site administration for research access.

23 Based on the issues presented, guidelines for researchers are proposed to reduce
24 objections to research activities, to facilitate communication with website administration, and to

1 achieve express or implied consent. These include **notification** to website administration of
2 intended automated research activity, **description** of the research project posted as a web page,
3 and clear **identification** of automated requests for web pages. In order to encourage good
4 research practices with respect to automated data collection, suggestions are made with respect to
5 disclosing methods used in research papers and for self regulation by academic associations.
6 **Keywords:** Internet, research, automated data collection

1. INTRODUCTION

The Internet is used for many academic research purposes. The focus of this paper is academic access to commercial website resources. Two examples of recent publications in this journal demonstrate some of the varied approaches to such data collection.

Ba and Pavlou (2002) conducted a study of trust building in electronic markets. Using a web browser, they collected data about 682 Ebay auctions and feedback information for their sellers over a 45 day period. They used this information to demonstrate that positive characteristics of seller feedback profile lead to price premiums for sellers. Bapna et al. (2004) wrote "an automated agent to capture information directly from auction websites" (p. 27). The agent collected data every five minutes over twelve months. The data was used to show five distinct classes of bidders and how their proportions changed from 1999 to 2000.

These articles are only two of many that rely on information from commercial websites. A review of three top outlets for information systems research (MIS Quarterly, Information Systems Research, and Management Science) for the years 2002-2004 indicates eight articles that make use automated data collection techniques from commercial websites (Clemons et al. 2002; Palmer 2002; Bapna et al. 2003a; Bapna et al. 2003b; Brynjolfsson et al. 2003; Easley & Tenorio 2004; Pavlou & Gefen 2004; Bapna et al. 2004). At least another twelve make some use of commercial website resources without automated interaction (Ba & Pavlou 2002; Zhu & Kraemer 2002; Chen & Hitt 2002; Kim et al. 2002; Koufaris 2002; Dewan et al. 2003; Dellarocas 2003; Pinker et al. 2003; Bolton et al. 2004; Bhargava & Choudhary 2004; Snir & Hitt 2003).

As the role of the Internet continues to increase, studies of commercial Internet sites are crucial to academic research. The conditions for collecting data required for such inquiry are

1 important for both practical and legal reasons. The benefit of using data from commercial
2 websites for academic research is clear. Since the sites are publicly accessible, our field has
3 apparently given little thought to issues of research access, especially automated access.
4 However, commercial websites almost universally post terms of service (TOS) documents that
5 expressly restrict the manner in which individuals may use the site. The strict adherence to these
6 terms would often exclude academic inquiry of any kind. In fact, with one exception (Koufaris,
7 2002) each article cited in the prior paragraph, would violate the terms of at least one website
8 TOS document if conducted today².

9 This paper seeks to illuminate the legal issues surrounding the use of commercial
10 websites for academic research. It suggests actionable guidelines for researchers who make use
11 of commercial websites. The rest of this paper is organized as follows: Section 2 summarizes the
12 typical ways that academic researchers make use of commercial websites. Section 3 describes the
13 means that commercial web sites employ to control how their data are accessed. Section 4
14 describes the current legal issues pertaining to unauthorized access of websites. Section 5 details
15 our recommendation for the field for the appropriate use of commercial websites in academic
16 research and Section 6 presents our conclusions.

17 **2. DESCRIPTION OF ACADEMIC USE OF COMMERCIAL INTERNET RESOURCES**

18 The academic use of commercial Internet resources varies widely. Some studies use
19 commercial sites only to gain an understanding of a context they model analytically; some use
20 them for background or general information; some study the site *per se*, sending subjects or
21 researchers to evaluate specific functions within the site; and, others collect data publicly

² Because Snir and Hitt (2003) do not disclose site identity, their study excluded

1 available on the site for analysis. Table 1 categorizes in to three groups the 20 articles found in
 2 *MIS Quarterly*, *Information Systems Research*, and *Management Science* from 2002 to 2004 that
 3 made use of commercial Internet resources.

Use	Study
General information, background, context	Bhargava & Choudhary 2004; Chen & Hitt 2002; Bolton et al. 2004; Pinker et al. 2003; Dellarocas 2003; Dewan et al. 2003; Pavlou & Gefen 2004
Study website <i>per se</i> .	Koufaris 2002; Kim et al. 2002; Zhu & Kraemer 2002; Palmer 2002; Agarwal & Venkatesh 2002
Analyze data collected from site	Ba & Pavlou 2002; Brynjolfsson et al. 2003; Bapna et al. 2003a; Bapna et al. 2003b; Easley & Tenorio 2004; Clemons et al. 2002; Bapna et al. 2004; Sinr & Hitt 2003

Table 1. Use of commercial Internet resources in recent IS literature

4 More than a third of the research reported in these articles made use of automated
 5 Internet data collection agents to accomplish their research objectives. Such automated practices
 6 merit special discussion because their use raises additional issues. Automated Internet data
 7 collection agents are simply computer programs designed to interact with web servers via the
 8 Internet to collect various kinds of data (Kauffman et al. 2000). They retrieve web pages and
 9 parse them to find data to be stored for analysis, references to images to be downloaded, and
 10 links to other web pages that might contain useful information. Because the agents have the
 11 capability to simulate data being entered into a web form (such as a customer order form) and
 12 posted by a user, they can dynamically interact with electronic commerce web servers and
 13 collect detailed data about various practices and behaviors in the online environment. Many
 14 programming environments provide tools for the development of such agents, so they are
 15 increasingly simple to develop and deploy in meaningful ways (Allen and March 2000).

16 Automated data collection agents vary widely in their architecture and sophistication.
 17 They can retrieve documents in either a serial or parallel fashion. When an automated agent

1 requests a page from a web server, it typically takes much longer to retrieve the page than it does
2 to parse it into data elements and process those data. An agent that processes its requests serially
3 (waiting for the first request to be retrieved and parsed before the second request is issued),
4 places a much lower burden on the web server with which it is interacting than an agent that
5 issues concurrent requests for a large number of pages and parses them as they are received.
6 Data collection agents can also be built to operate in a distributed manner in which several copies
7 of an agent can be running simultaneously on different computers, all working together to
8 complete the same data collection task. Because Internet data collection agents are designed to
9 interact with web servers over the Internet, they are easily adapted to cooperate with each other
10 over the same channel, making it relatively simple to deploy massively parallel, geographically
11 dispersed data collection networks with hundreds or thousands of nodes. Typically, details about
12 the particular data collection agent used are not disclosed in academic articles, although the
13 personal interactions of the authors with several researchers actively using this technology
14 indicate that current research employs the full range from the simplest to the most complex.

15 **3. MEANS USED BY WEBSITES TO CONTROL ACCESS TO DATA**

16 Although our experience indicates that few commercial websites have specifically
17 considered the implications of either allowing or preventing access for academic research,
18 virtually all have employed a variety of means to grant access to some while restricting access by
19 others. For competitive reasons, commercial websites employ various technological and legal
20 means to prevent unauthorized access to their data.

21 **Technological Means to Control Access to Website Content**

22 Technological mechanisms for such control include password protection to portions of
23 websites with restricted data, encryption of data transmissions to prevent unauthorized access to

1 their data while it is sent across the public network, copy protection, and traffic monitoring
2 coupled with page request refusal for access that may be undesirable because of its volume or
3 because the site objects to a particular use such as a competitor accessing information about
4 product offerings.

5 The technological mechanisms work directly to prevent unwanted access, that is, they
6 actively deny access to content unless the requests present the proper credentials. The use of
7 technological controls alone is unsatisfactory for four reasons. First, there is typically some way
8 to circumvent the technology: encryption can be cracked, copy protection is regularly defeated,
9 and passwords can be obtained under false pretenses. Second, even if access is correctly granted,
10 once content has been delivered, there may be inadequate means to control its subsequent use.
11 For example, a user may purchase a copy of a movie intended for in-home use and later use it in
12 a public setting. Third, the implementation of technical means may place extra burden on users,
13 as for example, requiring a user to read and enter a code rendered as a graphic with irregular type
14 in order to access specific content. Finally, technical means generally place increased demands
15 on the infrastructure. For example, encrypted transmissions are larger (sometimes by a factor of
16 two) than their unencrypted counterparts.

17 **Legal Means to Control Access to Website Content**

18 To complement the performance of technological means of controlling the access and use
19 of data, websites have several legal means at their disposal. Legal means currently in use include
20 seeking a court order to enjoin individuals or organizations against access or seeking the redress
21 of damages on the basis of trespass, copyright infringement, or breach of contract. Unlike their
22 technological counterparts, legal means work indirectly, relying on court injunctions, the threat

1 of court action, and the costs of potential litigation to deter individuals from inappropriate use of
2 system data and resources. There are three legal claims site owners typically make to control the
3 use of their resources: copyright infringement, trespass to chattel, and breach of contract.

4 Although issues surrounding copyright infringement are significant, they are typically of small
5 concern to academic researchers as will be discussed in the next section. Trespass to chattel and
6 breach of contract both hold significant implications for academic research and both depend
7 heavily on the terms specified in two documents maintained by sites concerned about access.
8 These are a site's terms of service (TOS) document and a robots.txt file.

9 **Terms of Service Document (TOS)**

10 The home page of a commercial website generally has a link to the site's terms of service
11 (TOS) document (also often called "terms of use" document). Along with specifying the details
12 of various site policies, this document describes acceptable use of the site.

13 Because of the central role that TOS documents play in asserting the legal rights of a
14 website, and because they are written with the intent of forming a binding contract with users of
15 the site, they are typically lengthy and written in precise language. Moreover, they tend to be
16 written as adhesion contracts (drafted entirely by one party and heavily restrict the other party)
17 and allow only a few very specific uses of the site. Because of the inherent inequality in such
18 contracts, courts tend to interpret any ambiguity in favor of the non-drafting party. Accordingly,
19 their terms often appear impractical. Consider the following excerpt from the walmart.com TOS
20 document:

21 Your use of the Site following any such change constitutes your unconditional agreement to follow and be
22 bound by these Terms of Use as changed. For this reason, we encourage you to review these Terms of Use
23 whenever you use this Site. (<http://www.walmart.com/catalog/catalog.gsp?cat=119985>, 3/28/2005).

1 Considering that when printed single-spaced in 12 point type, the walmart.com TOS
2 document is 12 pages, the statement that users should read it each time they visit the site is
3 unreasonable. In fact, companies must know that almost no users have even seen the site's TOS
4 document. An analysis of the web logs of one major corporation's web site indicated that out of
5 every one million visits to their home page, fewer than five followed the link to the TOS
6 document. This demonstrates one of the very interesting characteristics of the TOS document: It
7 is intended to support litigation rather than to enlighten users with respect to the proper use of the
8 site. When a website wants to teach users the appropriate way to make use of site resources, it
9 does so with techniques that are much more user accessible, such as frequently-asked-question
10 pages (FAQs) and Help pages.

11 The role of the TOS document as a legal tool rather than the definitive document that
12 details acceptable use of site resources is demonstrated in the following excerpt also taken from
13 the TOS document at walmart.com.

14 You may download or copy the Contents and other downloadable materials displayed on the
15 Site for your personal use only. No right, title or interest in any downloaded materials or
16 software is transferred to you as a result of any such downloading or copying. You may not
17 reproduce (except as noted above), publish, transmit, distribute, display, modify, create
18 derivative works from, sell or participate in any sale of or exploit in any way, in whole or in
19 part, any of the Contents, the Site or any related software
20 (<http://www.walmart.com/cservice/terms.gsp, 3/28/05>)

21
22 This policy clearly permits copies of the site's data for some personal use, and indicates that
23 commercial use of the data is prohibited. However, this statement directly conflicts with
24 another document that describes allowable use of the site: robots.txt.

25 **Robots.txt File**

26 The robots.txt file plays a role that is similar to TOS documents; however, its use is
27 limited to automated retrieval tools. Written in compliance with Standard for Robot Exclusion

1 (Koster 1994), the robots.txt file provides specific instructions for the kinds of automated
2 retrieval tools (called robots or bots) used by Internet search sites like google.com to index
3 Internet resources. By reading the file, which is located in the root directory of a web server, a
4 bot can determine what parts of the site it is allowed to access. Like a TOS document, robots.txt
5 does not exert any technological control over the use of site resources; individuals who write
6 bots may or may not implement the functionality to read and respect the restrictions it contains.
7 The robots.txt file for walmart.com follows:

```
8         # go away  
9         User-agent: *  
10        #Disallow: /  
11        Disallow: /solutions  
12        Disallow: /cservice  
13        Disallow: /reflect.gsp
```

14
15 (<http://www.walmart.com/robots.txt>, 3/28/05)

16
17 This policy indicates that all bots are specifically restricted from requesting documents that begin
18 with “walmart.com/solutions,” “walmart.com/cservice,” or “walmart.com/reflect.gsp,” but that
19 all other areas of the site are available for automated retrieval. In this case, walmart.com has
20 specifically allowed bots access to all product and pricing information located under
21 “walmart.com/catalog.” At the same time, it restricts bots from accessing parts of the site
22 dealing with order tracking, the store locator, and product recall information, all of which are
23 located under “walmart.com/cservice.” Clearly, the intent of the robots.txt file is to allow bots to
24 index the content of the site to make it available to be searched by users of Internet search
25 engines such as google.com, yahoo.com, and altavista.com—all commercial endeavors, even
26 though a strict reading of the TOS document clearly proscribes this kind of access.

1 Thus, walmart.com allows itself to be indexed by search engines while preserving the
2 strength of its TOS document. The contradiction is likely intentional and demonstrates that the
3 wording of TOS documents should not be taken as the absolute position of a site with regard to
4 the use of its resources.

5 **4. LEGAL ISSUES PERTAINING TO INTERNET ACADEMIC DATA COLLECITON**

6 The legal tools that websites employ to control access to their resources largely rely on
7 three separate legal principles: trespass, copyright infringement, and breech of contract.

8 **Trespass as a Basis for Legal Challenge to Academic Use of Commercial Websites**

9 The unauthorized use of website resources by commercial entities has been the subject of
10 considerable case law and courts have begun to formulate legal standards for such activity (Burk
11 2000; O'Rourke 2000; Elkin-Koren 2001). By far the most successful legal claim asserted
12 against on-line data gathering has been the theory of trespass to chattels, an ancient legal doctrine
13 intended to protect the owners of moveable property (chattels) against interference with the use
14 of their property. At common law, a claim of trespass to chattels requires contact be made with
15 the chattel, notification that contact was unwelcome, and interference with or dispossession of
16 the chattel, resulting in some harm or damage to the chattel, or pecuniary loss to the owner. This
17 definition of trespass has been reformulated for networked computers to hold that electrical
18 impulses satisfy the common law requirement of physical contact, and the increased load on the
19 networked system qualifies as interference or dispossession, and terms posted in the TOS
20 document or robots.txt fill the requirement of notice. Harm or pecuniary loss is often presumed
21 from the loss of processing cycles, the use of network transmission capacity, or the diversion of
22 data storage capacity because the owners of web sites bear the cost of responding to requests
23 made to their servers. While any single incremental request does not bear a direct variable cost,

1 clearly the infrastructure needed to respond to 10 million requests per day is more costly than the
2 infrastructure needed to respond to a thousand requests per day.

3 Several courts have now embraced this renovated legal theory. For example, in *eBay v.*
4 *Bidder's Edge*, a United States District Court relied on a theory of trespass to enjoin Bidder's
5 Edge, an aggregator of on-line auction data, from collection of data from the eBay site (eBay
6 2002). Similarly, in *Register.com v. Verio*, a trespass theory was used to penalize the collection
7 of ownership data from a publicly accessible domain name database (Register 2000). Most
8 recently, in *American Airlines v. Farechase*, a Texas state court enjoined the producer of
9 software used to search in real time for air fares on airline websites (American 2002).

10 Several trespass cases have as an alternate theory relied upon federal statutes prohibiting
11 unauthorized access to computers. Section 1030 of the U.S. criminal code penalizes
12 unauthorized access to a networked computer when damage results. This statute was originally
13 enacted in response to the Robert Morris Cornell "worm" incident. Unfortunately, Congress
14 failed to define "authorization," allowing this statute to be applied to a broad range of Internet
15 activity where permission to engage in the activity has been implicitly or explicitly denied. For
16 example, in the case of *EF Cultural Travel v. Explorica*, automated retrieval of travel pricing
17 data by a bot was held to constitute unauthorized access to a travel agency's website (EF Cultural
18 Travel 2001). Similarly, in the recent *Farechase* decision mentioned above, access to the airline
19 website was also enjoined on an alternate theory of unauthorized access under a state computer
20 crime statute. Many states have such statutes that may be applied in addition to the federal
21 unauthorized access statute.

1 These cases are significant for academic researchers because they establish an almost
2 absolute right for website owners to exclude others from their publicly available equipment and
3 information. Although these cases do not deal directly with the academic use of commercial
4 Internet resources, the legal theories advanced do not provide for any research or public benefit
5 exceptions. Any load placed on a networked server without prior consent may be considered
6 trespass or unauthorized access. Although liability for trespass or other legal claims depends to a
7 great extent on the terms of the researcher's employment contract, both the researcher and the
8 sponsoring institution could be targets for legal action.

9 **Copyright as a Basis for Legal Challenge to Academic Use of Commercial Websites**

10 In the commercial context, copyright claims have also been asserted against some
11 undesired use of system resources. However, this claim cannot extend to the mere collection of
12 facts. Under the United States Supreme Court holding in *Feist v. Rural Telephone*, facts are
13 generally not copyrightable because they lack originality (Feist 1991). Compilations of facts
14 may be protected by copyright, if the selection and arrangement of the facts is original.
15 Extracting the facts from their matrix of selection and arrangement takes nothing copyrightable,
16 so long as the selection or arrangement is not also taken. Consequently, copyright claims
17 asserted to stop unwanted access of system resources have not been successful to date.

18 A more credible copyright claim might be made for copies of web pages generated in the
19 course of access by automated agents. Computer networks function by reproducing digitized
20 files and transmitting them piecemeal to other computers where they are reassembled. Along the
21 way, all or part of the files are reproduced by routers and other intermediate computers. When
22 automated data agents request web pages from networked servers, the pages served typically are

1 reproduced in several places, including the recipient computer cache file and in the RAM of the
2 recipient computer. In operational terms, the Internet comprises a massive, distributed copy
3 machine and any access to files on the network produces copies.

4 Additionally, automated search agents may be designed to request a page and store that
5 page as originally downloaded, in addition to parsing it for the information the researchers wants.
6 This practice has significant academic value. When longitudinal data is analyzed, some
7 observations may be puzzling; by keeping an original record of the web pages, a researcher can
8 refer back to the source data in its original context to examine it more closely. The anomalous
9 observations may be accurate, they may appear because of an error in the program that generated
10 the file at the website, or they may appear because of an error in the data collection agent.
11 Without the original HTML file, determining the cause of the anomaly is difficult at best, and
12 most often impossible. Perhaps more importantly, when the original HTML file is stored, it can
13 be reparsed by a researcher without placing any load on the originating website. Accordingly,
14 this practice benefits the research process as well as the website owners. However, storage of
15 original HTML pages constitutes a conscious form of archiving or copying beyond the copies
16 made in the normal course of the Internet's technological function.

17 Under United States law, copies that are fixed long enough to be perceived, either
18 directly or with the aid of a machine, are subject to copyright law. The Copyright Act indicates
19 that copies existing only temporarily in computer memory fall outside the Act, but copies that
20 exist in more durable media, such as those written to a hard drive, are likely subject to copyright
21 law. Some courts have held that even RAM copies may be relevant for copyright purposes.

1 Thus, copies made when web pages are served up could constitute copyright
2 infringement if they are not authorized by the owners of the pages. Typically, such authorization
3 is implied from the public availability of the pages. Due to the nature of the technology, they can
4 only be viewed if copies are made, so we infer that the owners knew or intended that such copies
5 should be made in order for the pages to be viewed.

6 However, implied consent can usually be revoked by some sort of notice to users.
7 Continued use of the copyrighted material in the face of explicit notice that access is
8 unauthorized could then constitute infringement. In such a case (Kelly 2002), at least one court
9 has held that display of images served to a search engine constituted copyright infringement;
10 however, academic research agents are unlikely to publicly display the pages they search.

11 Alternatively, the creation of such copies in the course of Internet access may constitute
12 fair use. In the United States, unauthorized uses of copyrighted materials may be permissible
13 under the fair use provisions of the Act, which particularly favors uses for non-commercial
14 educational or research purposes. Several courts have held that creation of temporary or
15 intermediate digital copies in the course of a permissible end use may constitute fair use. Fair
16 use principles may also tend to favor academic archiving of web pages that are under study.

17 It must be remembered, however, that the Internet is international in its reach, and fair
18 use does not exist in the copyright law of any country other than the United States. Many other
19 countries have limited exceptions or “fair dealing” provisions that may permit certain otherwise
20 unauthorized uses. These provisions often include exceptions for educational or research
21 purposes. Not only do the exceptions vary by country, but copyright law in some regions is
22 currently undergoing substantial change. For example, the March 2001 European Union

1 Copyright Directive required member states to make changes to their national copyright laws.
2 Although the directive allowed member countries substantial latitude in their national adoption,
3 many member states have still not legislated the required changes but continue to debate over
4 various drafts. Thus, not only is the storage of web pages for analysis not permissible in some
5 countries from which the pages originate or to which they may be transmitted, but for many
6 countries, that status may be changing.

7 **Contract Breach as a Basis for Legal Challenge to Academic Use of Commercial Websites**

8 The terms posted in TOS documents purport to form a binding contract with those using
9 the site where the TOS document is posted. Under this theory, the TOS document is held to act
10 in the same manner as “shrinkwrap” or “clickwrap” licenses accompanying software. Under a
11 legal fiction of agreement, users are said to manifest assent to the terms of the contract by the act
12 of making use of the site. However, formation and enforceability of such contracts is
13 problematic, both because consumers may not have seen the terms, and because terms of the
14 agreement may be unconscionable or surprising. Courts have been split on the enforceability of
15 such purported contracts. Virginia and Maryland have recently amended their contract law by
16 adoption of the Uniform Computer Information Transaction Act, or UCITA, to make this type of
17 contractual assertion legally enforceable. There has been significant opposition to adoption of
18 this act in other states, including enactment by a few states of legislation explicitly prohibiting
19 the application of UCITA to their citizens. Thus, the contractual implications of posted TOS
20 document may vary not only from country to country, but from state to state in the United States

Achieving Express Consent to Avoid Legal Challenges

1
2 As the discussion of TOS documents demonstrates, each legal objection to automated
3 access hinges upon the degree of access consented to by the owner of a website, and the
4 inferences that can be drawn from website owner actions regarding consent. The challenges to
5 automated data collection agent access for research can be removed by obtaining express consent
6 to access the site. Although conceptually desirable, express consent is frequently impractical.

7 Specific express consent can be obtained by a researcher for a particular research project.
8 In any situation where a researcher interacts with a certain web server, the owner of the server
9 clearly has the right to expressly allow the researcher's data collection activity. It has been our
10 experience that when owners of small commercial sites are informed of an intention to request
11 small (and even moderate) amounts of data from their servers for purposes of academic research,
12 some are very willing to work with the researcher. Accordingly, one approach for a researcher to
13 take is to seek express consent. When express consent has been granted, many of the other legal
14 issues cease to be concerns.

15 Obtaining express consent may not be a satisfactory approach for several reasons. First,
16 in some situations, the resources consumed in the process of obtaining express consent (both on
17 behalf of the researcher and on behalf of the organization whose consent is sought) may far
18 outweigh the resources that would be used in obtaining the data. Consider, for example, what
19 the process might entail in gaining express consent from any substantial player in the online
20 community. Once an initial request for express consent is sent, it may several days (or weeks)
21 for the request to even reach the person or committee with the authority to grant consent. Once
22 the decision maker has the request there may be other issues that prohibit consent from being

1 granted. Decision makers may be concerned about setting a particular precedent within the
2 organization, or they may be unsure of the legal ramifications of granting express consent.
3 Moreover, since there is likely very little direct benefit to the organization, even the process of
4 considering such requests may outweigh the perceived benefits, leading to a policy of summarily
5 rejecting all such requests. One researcher sought express consent to use automated data-
6 collection agents to gather data from a major e-commerce site, even visiting the company
7 headquarters. Although upper management of the site was interested to learn of the results of the
8 studies the researcher had already conducted using data from their site, they would not grant
9 permission for continued automated access. However, after the meeting, one member of the
10 management team gave additional guidance to the researcher in a separate communication. The
11 executive told the researcher that the company's site had mechanisms in place to monitor
12 automated data access and that anytime automated data access became onerous, site
13 administrators take corrective action. He explained that because the researcher's prior data
14 collection activity had not been "noticed" by the site's administration, he could reasonably expect
15 to continue his activity at similar levels without objection by the site's management. Although
16 anecdotal, this incident demonstrates the tendency for organizations to be conservative in formal
17 decisions, even as informal channels may be less so.

18 A second concern with requesting express consent is that once the dialogue is opened, it
19 is entirely possible that, beyond not granting express consent, an officer of the company may
20 expressly forbid a researcher from conducting the study. In this case any implied consent that
21 might have reasonably been inferred by the connection of the server to a public network would
22 clearly be revoked.

1 A third concern with seeking express consent is that such contact with web site
2 administration may threaten the validity of some research studies. In much the same way that
3 people behave differently when they know they are being observed, websites may present
4 different information if they know they are being examined for purposes of an academic study.
5 In 2000, Amazon.com came under heavy fire for giving different prices to customers depending
6 on which browser was used to issue the request as well as the age of Amazon cookie in the
7 browser's cache on the customer computer (Rosencrance, 2000). Clearly the technology exists
8 to relay different information based on the characteristics of the request. If there is reason to
9 believe that a site, once aware of the study, may tailor the information sent to an academic
10 researcher, it may contaminate the data collection to notify the organization of the study by a
11 request for express consent.

12 There are significant differences between seeking express consent from a website to use
13 their data in academic research and negotiating research access to the resources of organizations
14 in a traditional research setting. The differences are based on the fact that in automated site
15 access, the data to be accessed for the research activity has already been placed before the public
16 by posting it on a website. In a traditional research setting, an organization being asked for
17 access to processes or data will typically have five concerns: access will disrupt normal
18 operations, direct cost associated with granting access, known or unknown legal issues with
19 releasing some information, publication or leaking of information may have adverse effects on
20 the company's competitive position, and allowing researchers access to particular locations may
21 expose the company to undue liability. In the case of automated access to websites,
22 organizations address each of these issues before making data and processes publicly accessible.

Implied Consent when Express Consent is Absent

In some cases, commercial websites post neither a TOS document nor a robots.txt file. This is the reason that one of the research articles listed in Table 1 (Koufaris, 2002) would not violate the terms of service if it were conducted today—it made use of booksamillion.com, which posts neither a TOS document nor a robots.txt file. When such is the case, one might reasonably infer that an organization grants “implied consent” for the public to access its web servers, even through the use of heavily automated means, because they are connected to a public network in which the use of automated retrieval mechanisms is pervasive.

5. RECOMMENDATIONS FOR RESEARCHERS COLLECTING DATA FROM COMMERCIAL WEBSITES

The following guidelines for researchers reflect our views regarding manual website access and researcher circumvention of technological restrictions. There are three recommendations for actions to communicate with website administration when using automated access to websites. The discussion of the recommendations reflects the perspectives of both website administration and academic researchers. The views and recommendations are intended to guide academic researchers conducting studies that make use of publicly available, commercial Internet resources or are contemplating such research.

Views on Manual Access to Websites and Circumvention of Technological Restrictions

Manual, non-automated access of information on publicly available web pages should be acceptable without special permissions or actions. Even though the website may not expressly permit such access for research, the load on the website is negligible. Such limited access outside the directly intended purpose of a website is expected by site owners. Therefore, limited access by those doing research fits within normal website expectations. Moreover, the strict

1 enforcement of TOS document terms in this situation would virtually close commercial websites
2 to any examination by academia. This would be bad public policy and may lead to corrective
3 legislation, weakening the efficacy of existing legal means for controlling access to site content.

4 Any action taken by an academic researcher to bypass technological means of controlling
5 access to website resources is unacceptable. Such actions show that the researcher is acting in
6 bad faith and will reduce the credibility that other academic researchers hold with commercial
7 websites. Perhaps more importantly, the circumvention of technological measures to gain access
8 to data clearly violates the Digital Millennium Copyright Act (DMCA) in the United States, as
9 well as provisions set forth in the European Union Copyright Directive (2001/29/EC).
10 Accordingly, such behavior may have criminal consequences. If access to data is restricted using
11 technological means, the only acceptable way for a researcher to access the data is to seek
12 express consent in cooperation with website management.

13 **Recommendations for Automated Access to Websites**

14 Automated access of information on publicly available web pages can place a
15 substantial load on site resources. For such access that complies with a site's robots.txt file, or in
16 its absence, specific terms about automated retrieval in the TOS document, we suggest two
17 principles to facilitate communication between the researcher and the website. These are
18 *description* (providing site administration a document that describes the research) and
19 *identification* (providing site administration with the ability to easily determine the effect of the
20 research on the site). In addition, when the robots.txt file (or in its absence, the TOS document)
21 does not expressly permit the planned access a recommended procedure is *notification* of
22 research activity.

1 **Provide Description of Research Activity as a Webpage**

2 Commercial websites are becoming increasingly aware of the need to limit competitor
3 access to their site resources (YiHua, et al. forthcoming). Unfortunately, the activity of a
4 competitor using an automated, data-collection agent often bears strong resemblance to a
5 researcher using such technology. For this reason, it is important to clearly describe that the data
6 collection activity is strictly for academic research purposes. A detailed description of the
7 research study (including potential findings) can be posted as a web page. This may make a
8 system administrator more likely to allow access than if no information is provided about the
9 research. The study's description document should provide the email address of the appropriate
10 member of the research team and perhaps a telephone number. It may also include a statement
11 of the researcher's intention to not disrupt the site's operations and to comply with the wishes of
12 site management, as well as a statement about whether the identity of the site will be disclosed in
13 any research report.

14 In a particular longitudinal study in which no description of the academic nature of the
15 study was provided to website administrators, one website blocked access to the researcher by
16 restricting the range of IP addresses allowed to access the site. Later when the researcher
17 apologized for causing trouble at the site, the site administrator indicated that if he had known
18 that the requests were part of an academic study, he would not have terminated access. There are
19 different approaches to making a website aware of the academic nature of a particular study; one
20 is to include a reference to the descriptive webpage with each automated request.

21 **Provide Identification of Source with each Automated Request**

1 A website administrator can usually determine the source of accesses on the site through
2 the analysis of log files stored by web servers. Commercial web servers have the functionality to
3 record information about each request made of the server. They record specified portions of
4 each HTTP (hypertext transfer protocol) request. The information in the request header is used
5 to specify the requested resources on the server and also to indicate information about the
6 program (e.g., a bot or a browser) that issued the request.

7 The User-agent field of the HTTP request header is used to uniquely identify the program
8 (and version) that is requesting the resources of the server. The field may also include a
9 comment enclosed in parentheses that can be used for research access identification. For
10 example, Netscape Navigator version 6.2 sends the following in the User-agent field of the
11 request header:

12 Mozilla/5.0 (Windows; U; Windows NT 5.1; en-US; rv:0.9.4.1) Gecko/20020508 Netscape6/6.2.3

13 This indicates that the user agent is based in the open source Mozilla code base and also supplies
14 information about the current operating environment in the comment portion of the field
15 (enclosed in parentheses). Finally, it has information regarding the version of Navigator.

16 Researchers should therefore use the User-agent field of the HTTP request to identify
17 their data-collection agent. Moreover, a researcher should use the comment portion of the field
18 to direct a site administrator to the web page that describes the current research endeavor. For
19 example, a researcher might formulate the user-agent field for all HTTP requests of a particular
20 study as follows:

21 AcademicAgent/1.0 (<http://research.tulane.edu/~gallen/Study7.html>)

22 This helps the site administrator identify how much traffic is being generated by the
23 particular agent and allows easy access to the document that describes the research project. This

1 approach works well because the User-agent field is almost universally logged by all web
2 servers, so when an administrator examines the logs, the required information is readily
3 available. By making requests identifiable in this manner, a researcher facilitates
4 communication with a concerned administrator.

5 Following these suggestions to describe the academic intent of the research access along
6 with the practice of identifying each automated request places some burden on the researcher but
7 it also provides them the benefit of initiating an open, cooperative relationship with the website
8 rather than a closed, adversarial one. Clearly, in this setting, website administration will favor
9 these recommendations because it provides more information as they decide which automated
10 requests to allow and which to deny using various technological tools for blocking access.

11 **In Some Cases, Send Notification of Research Activity**

12 If a researcher wishes to pursue a data collection protocol that is not expressly permitted
13 by the robots.txt file (or in its absence, TOS documents), he or she should send a message to the
14 webmaster of the site indicating his or her intentions. This notification should indicate the
15 expected duration of the data collection, the intended frequency of requests, a description of the
16 kinds of data to be collected, a means to contact the appropriate member of the research team,
17 and a statement of the researcher's intention to cooperate with the site's administration. The
18 notification need not be phrased as a request for permission, and a researcher should not expect
19 to receive a response. This way, the administrator can easily express any concerns to the
20 researcher; however, if the amount of data to be requested is of no concern, the administrator
21 may simply ignore it. From the perspective of the academic researcher, the requirement to send
22 notification is a burden, and may result in express notification that the research activity is

1 unwanted; however, it demonstrates that the researcher is acting in good faith which may be
2 important to the resolution of a website's concerns without legal action. From the perspective of
3 the website, this is an agreeable arrangement for several reasons. First, it removes any burden
4 the site would otherwise have in detecting unwanted automated access. Second, it allows
5 research activity with limited infrastructure impact to be allowed without the time required for a
6 formal approval process. Third, it allows the site to permit some research data collection activity
7 without expressly setting a precedent. Fourth, it reduces uncertainty by allowing site to observe
8 the effect of the data collection activity before rendering an official decision. Fifth, it preserves
9 the site's absolute ability to revoke access at any time because formal permission is never
10 granted.

11 We should mention that suits have been filed alleging violation of the DMCA because
12 provisions in the robots.txt file were ignored. However, none have been successful to date and
13 precedent is emerging (e.g. Lexmark 2005) that more clearly excludes specifications of a
14 robots.txt file from being considered "effective technical protections" as specified in the DMCA.

15 **Recommendations for a Generally Accepted Policy for Disclosing**
16 **Automated Data Collection Procedures in Research Papers**

17 Rather than waiting for legislation to define research access, it may be helpful for
18 academic researchers to agree on acceptable behavior. For example, a formal policy of accepted
19 behavior for automated data collection could be adopted by an appropriate organization such as
20 the Association for Information Systems (AIS) or the Association of Internet Researchers
21 (AoIR). This would provide the basis for wording in research papers indicating researcher
22 compliance with the policy or if the policy was violated, an explanation of why. Until such a
23 policy is in place, reviewers may reasonably require authors to provide a brief description of the

1 automated data collection system (perhaps in a footnote) and list any communication between
2 the researchers and web site administrators—such as requests for consent made by the researcher
3 or inquiries made by website administration into the data collection activity.

4 Establishing such a policy which could be broadly adopted on a voluntary basis may hold
5 significant future implications. Just as industry self-regulated adoption of accepted privacy
6 principles in the context of Fair Information Practices has led to proposed legislation that
7 complements self-regulation programs rather than supplants them, our field has the opportunity
8 to establish the standard that could influence potential future legislation in a manner that may
9 favor academic research as a preferred activity—especially if we as academic researchers can
10 demonstrate effective self-regulation.

11 **6. CONCLUSION**

12 Researchers are currently making substantial academic use of commercial Internet
13 resources. Such research activity is important in developing our understanding of the many
14 organizational aspects that are so deeply affected by the Internet. Because these commercial
15 Internet resources are publicly available, researchers have given little thought to the processing
16 load that their activity places on website infrastructure. Commercial sites have begun to exercise
17 various legal means to limit individuals from making automated access of their resources.
18 Academic research is not exempt from the legal arguments that have been successfully advanced.

19 Commercial websites define acceptable access of their resources in their terms of service
20 document often in conjunction with the robot.txt file located at the root directory of a web server.
21 Although the terms of service document may be useful for defining research access, it is often
22 written to support legal action against any possible site misuse. Consequently, its terms should
23 not be viewed as absolute with respect to academic research. The robot.txt file, when present,

1 defines specific areas of a site that are available to or restricted from automated interaction.
2 Accordingly, the robots.txt file should be viewed together with the terms of service document in
3 assessing the appropriateness of automated academic access.

4 Given some ambiguity in a website's position on allowing academic research access, we
5 recommend three polices to assist researchers in dealing with the practical problems of
6 automated data collection for research purposes: description (provide a webpage that describes
7 the research activity), identification (identify source of each page request along with a reference
8 to description), and notification (if research access is not expressly permitted) by sending a
9 message indicating purpose, scope, and frequency of requests.

10 If researchers can arrive at reasonable policies for automated research access of Internet
11 sites, then self regulation may be useful in dealing with site administration, as well as important
12 in influencing emerging legislation pertaining to academic use of commercial internet resources.
13 A professional organization is the appropriate place for such a policy to be formulated.

14 While the practices that we outline here will not grant academic researchers legal carte
15 blanche to crawl Internet websites or immunize researchers from legal liability, we expect that
16 adherence to these guidelines may allow researchers to avoid or defuse most legal conflicts.
17 Formal adjudication is costly, and is unlikely to be pursued if reasonable informal solutions are
18 available. In a situation where informal channels are not explored before legal action is taken by
19 a website owner, adherence to these practices demonstrates that a researcher has acted in good
20 faith, which may be critical in obtaining a favorable outcome.

21

22

REFERENCES

- 1
- 2 Agarwal, R. and V. Vankatesh. 2002. "Assessing a Firm's Web Presence: a Heuristic Evaluation
3 Procedure for the Measurement of Usability." *Information Systems Research*. 13 (2), 168-
4 186.
- 5 Allen, G. and S. March. 2000. "Developing Internet agents: a Tutorial using Visual Basic 6.0."
6 *The Proceedings of the International Conference on Information Systems 2000*. Brisbane,
7 Australia.
- 8 American Airlines v. Farechase, No. 067-194022-02 (Tx Dist. Ct. Tarrant Cty., March 8th 2002)
- 9 Ba, S. and P. A. Pavlou. 2002. "Evidence of the Effect of Trust Building Technology in
10 Electronic Markets: Price Premiums and Buyer Behavior." *MIS Quarterly*, 26(3), 243-
11 268
- 12 Bapna, R., P. Goes, and A. Gupta. 2003a. "Replicating Online Uankee Auctions to Analyze
13 Auctioneers' and Bidders' Strategies." *Information Systems Research*. 14 (3), 244-268.
- 14 Bapna, R., P. Goes, and A. Gupta. 2003b. "Analysis and Design of Business-to-Consumer
15 Online Auctions." *Management Science*. 49 (1), 85-101.
- 16 Bapna, R., P. Goes, and A. Gupta. 2004. "User Hetrogeneity and its Impact on Electronic
17 Auction Market Design: An Empirical Exploration." *MIS Quarterly*, 28(1), 21-43
- 18 Bhargava, H. K. and V. Choudhary. 2004. "Economics of an Information Intermediary with
19 Aggregation Benefits." *Information Systems Research*. 15 (1), 26-36.
- 20 Bolton, G. E., E. Katok, and A. Ockenfels. 2004. "How Effective are Electronic Reputation
21 Mechanisms? An Experimental Investigation." *Management Science*. 50 (11), 1587-
22 1603.
- 23 Brynjolfsson, E., Y. Hu, and M. D. Smith. 2003. "Consumer Surplus in the Digital Economy:
24 Estimating the Value of Increased Product Variety at Online Booksellers." *Management
25 Science*. 49 (11), 1580-1596.
- 26 Burk, D. L. (2000) The trouble with trespass. *Journal of Small & Emerging Business Law*, 4,
27 27-56.
- 28 Chen, P. Y. and L. M. Hitt. 2002. "Measuring Switching Costs and the determinants of Customer
29 Retention in Internet-Enabled Businesses: A Study of the Online Brokerage Industry."
30 *Information Systems Research*. 13 (3), 255-274.
- 31 Clemons, E. K., I. Hann, and L. M. Hitt. 2002. "Price Dispersion and Differentiation in Online
32 Travel: An Empirical Investigation." *Management Science*. 48 (4), 534-549.

1 Dellarocas, C. 2003. "The Digitization of Word of Mouth: Promise and Challenges of Online
2 Feedback Mechanisms." *Management Science*. 49 (10), 1407-1424.

3 Dewan, R., B. Jing, and A. Seidmann. 2003. "Product Customization and price Competition on
4 the Internet." *Management Science*. 49 (8), 1055-1070.

5 Easley, R. F. and R. Tenorio. 2004. "Jump Bidding Strategies in Internet Auctions."
6 *Management Science*. 50 (10), 1407-1420.

7 eBay Inc. v. Bidder's Edge, Inc., 100 F.Supp.2d 1058 (N.D. Cal. 2000).

8 EF Cultural Travel BV v. Explorica Inc, 274 F.3d 577 (1st Cir. 2001)

9 Elkin-Koren, N. 2001. Let the Crawlers Crawl: On Virtual Gatekeepers and
10 the Right to Exclude Indexing. *Univ. Dayton Law Rev.* 26, 179-209.

11 Feist Publications, inc. v. Rural Telephone Service Co., 499 U.S. 340 (1991)

12 Kauffman, R., S. March, and C. Wood. "Mapping Out Design Aspects for Data-Collecting
13 Agents." *International Journal of Intelligent Systems in Accounting, Finance, and*
14 *Management* 9 (4), 217-236.

15 Kelly v. Arriba, 366 F.3d 811, 822 (9th Cir. 2002)

16 Kim, J. , L. Jungwon, K. Han, and M. Lee. 2002. "Businesses as Buildings: Metrics of the
17 Architectural Quality of Internet Businesses." *Information Systems Research*. 13 (3), 239-
18 254.

19 Koster, M. 1994. "A Standard for Robot Exclusion." <http://www.robotstxt.org/wc/norobots.html>.
20 (on June 1, 2005).

21 Kougaris, M. 2002. "Applying the Technology Acceptance Model and Flow Theory to Online
22 Consumer Behavior." *Information Systems Research*. 13 (2), 205-223.

23 Lexmark v. Static Control 387 F.3d 522, 6th Circuit (2005)

24 O'Rourke, M.A. 2000. Shaping Competition on the Internet: Who Owns Product
25 and Pricing Information?. *Vanderbilt Law Rev.* 53, 1965-2006.

26 Palmer, J. W. 2002. "Web Site Usability, Design, and Performance Metrics." *Information*
27 *Systems Research*. 13 (2), 151-167.

28 Pavlou, P. A. and D. Gefen. 2004. "Building Effective Online Marketplaces with Institution-
29 Based Trust." *Information Systems Research*. 15 (1), 37-59.

30 Pinker, E. J., A. Seidmann, and Y. Vakrat. 2003. "Managing Online Auctions: Current Business
31 and Research Issues." *Management Science*. 49 (11), 1457-1484.

1 Register.com, Inc. v. Verio, Inc., 126 F. Supp.2d 238 (S.D.N.Y. 2000).

2 Rosecrance, L. 2000. "Amazon charging different prices on some DVDs." *Computer World*,
3 (Sept. 5).

4 Sheng, Y., P. Mykytyn, and C. Litecky, Forthcoming. "The Internet, Intelligent Agents,
5 Competitor Analysis and its Defense in the E-marketplace." *Communications of the*
6 *ACM*.

7 Snir, E. M. and L. M. Hitt. 2003. "Costly Bidding in Online Markets for IT Services."
8 *Management Science*. 49 (11), 1504-1520.

9 Zhu, K. and K. Kraemer. 2002. "e-Commerce Metrics for Net-Enhanced Organizations:
10 Assessing the Value of e-Commerce to Firm Performance in the Manufacturing Sector."
11 *Information Systems Research*. 13 (3), 275-295.

12

13