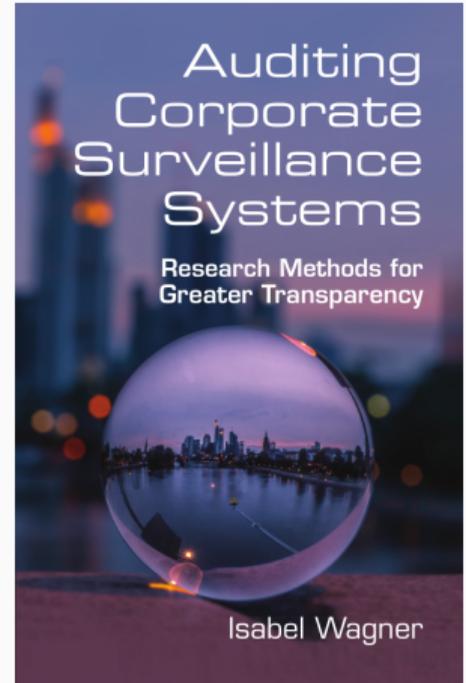


# DATA COLLECTION

---

Isabel Wagner

De Montfort University



Book design ©2022  
by Cambridge University Press

- Aim of data collection: execute the designed experiment and record response variables
- Ethical considerations: how to identify and minimize harmful impacts of data collection?
- Techniques for collecting data from different sources:
  - Archival sources
  - Passive data collection
  - Active data collection
    - Data collection from mobile apps
    - Data collection via crowdsourcing
- Methods to influence input variables
- Data storage

# ETHICAL CONSIDERATIONS

---

## RESEARCH ACTIVITY INVOLVING HUMANS

- Institutional Review Board (IRB): university process to consider/approve research designs
- When is IRB approval needed?
  - Gathering information from or about individuals and organizations.
  - Using archived data in which individuals are identifiable
  - Research into criminal activities
- Guidelines for considering ethics in different situations
  - Codes of ethics, e.g., ACM Code of Ethics, IEEE Code of Ethics
  - Privacy by Design<sup>1</sup>
  - Ethical research guidelines<sup>2</sup>

<sup>1</sup>A. Cavoukian, "Privacy by Design," Information and Privacy Commissioner of Ontario, Tech. Rep., 2013.

<sup>2</sup>E. Kenneally and D. Dittrich, "The Menlo Report: Ethical Principles Guiding Information and Communication Technology Research," U.S. Department of Homeland Security, Tech. Rep., Aug. 2012.

## ADNAUSEAM: ETHICAL JUSTIFICATION FOR A CONTROVERSIAL STUDY

- AdNauseam: browser extension hides ads, clicks on ads in the background<sup>3</sup>
- Three-part ethical and moral justification
- Is the aim laudable?
  - Ubiquitous surveillance violates tenets of liberal democracy, is against user's wishes
  - Privacy as societal value should be infused in systems
  - Protest against tracking by disrupting corporate surveillance business model (random clicking on ads pollutes user profiles, making them less valuable)
- Do alternatives exist?
  - Best alternative is regulation, which can seem unlikely to be realized
  - Cost induced by AdNauseam: server resources and bandwidth small compared to resources spent on ad serving
- Is it different from click fraud?
  - Fraud normally benefits the perpetrator (financially), not the case here

<sup>3</sup>D. C. Howe and H. Nissenbaum, "Engineering Privacy and Protest: A Case Study of AdNauseam," in *Proceedings of the 3rd International Workshop on Privacy Engineering*, vol. 1873, San Jose, CA, USA: IEEE, 2017, pp. 57–64.

## HARMFUL EFFECTS ON INDIVIDUALS

- Which research methods may cause harmful effects on individuals?
- Data collection from researcher-created app or browser extension
  - For example, to study cloud usage of mobile apps; price discrimination; mobile tracking; or ad pricing
- Studies with human participants
  - For example, studies that use surveys; crowdsourcing
- Data collection from websites/apps that display user data
  - For example, scraping from freelance marketplaces or résumé search engines
- Data collection through residential Internet connections
  - For example, crawling via Luminati to study vantage points in many countries
- Researcher-created synthetic data is visible to users
  - For example, researcher-run ad campaigns; measurement clients for Uber

## MINIMIZING HARMFUL EFFECTS ON INDIVIDUALS

- Informed consent
  - Inform participants about extent and purpose of data collection
  - Participants free to decide whether to participate
  - May not be possible if participants' awareness that they are participating biases the results
  - Additional duty to minimize harm; inform participants as soon as possible
- Data minimization
  - Only collect data that is necessary to achieve study objectives
  - For example, avoid collecting personally identifiable information
- Allow users to turn off data collection
- Process data on user device, only transmit summary data to researchers
- Limit sharing of datasets with user data
- Ensure no negative effects, e.g., no booking of tasks on freelance marketplaces, no requesting of rides on Uber

## HARMFUL EFFECTS ON ORGANIZATIONS

- User accounts
  - Limit account creation, minimize number of virtual personas
- Server resources and bandwidth
  - Rate limit requests to web servers
  - Respect robots.txt
- Responsible disclosure
  - Notify services before making vulnerabilities public, allow time to fix
- Ad impressions/clicks
  - Limit cost incurred by automated scraping of ads/clicking on ads
  - Estimate cost to advertisers, minimize

- Compliance with laws, regulations, policies
  - Different rules depending on location and research!
  - Computer Fraud and Abuse Act (CFAA)
  - Digital Millennium Copyright Act (DMCA)
  - General Data Protection Regulation (GDPR)
  - Terms of service<sup>4</sup> (may forbid automated queries, reverse engineering, scraping)
  - “non-commercial research for the public good that deals with issues of societal importance must be able to access public Web resources for research purposes as long as automated processes do not produce an unreasonable load.”<sup>5</sup>
- Compliance with ethical codes, e.g. ACM code of ethics
  - Need to comply with terms of service, “unless there is a compelling ethical justification to do otherwise”

<sup>4</sup>K. Vaccaro, K. Karahalios, C. Sandvig, et al., “Agree or cancel? research and terms of service compliance,” in *ACM CSCW Ethics Workshop: Ethics for Studying Sociotechnical Systems in a Big Data World*, Vancouver, Canada: ACM, 2015.

<sup>5</sup>G. Soeller, K. Karahalios, C. Sandvig, et al., “MapWatch: Detecting and Monitoring International Border Personalization on Online Maps,” in *Proceedings of the 25th International Conference on World Wide Web*, ser. WWW '16, Montréal, Québec, Canada: International World Wide Web Conferences Steering Committee, 2016, pp. 867–878. doi: [10.1145/2872427.2883016](https://doi.org/10.1145/2872427.2883016).



About Issues Our work News Take action **Donate**



- Lawsuit filed 2016
  - Does First Amendment grant researchers the right to provide false information to websites in the course of testing for discrimination on the basis of race, gender, and other characteristics protected under civil rights laws?
- 2020: federal court rules that research aimed at uncovering whether online algorithms result in racial, gender, or other discrimination does not violate the Computer Fraud and Abuse Act (CFAA)<sup>6</sup>

<sup>6</sup><https://www.aclu.org/press-releases/federal-court-rules-big-data-discrimination-studies-do-not-violate-federal-anti>

<sup>7</sup><https://www.aclu.org/cases/sandvig-v-barr-challenge-cfaa-prohibition-uncovering-racial-discrimination-online>

## DATA SOURCES

---

## ARCHIVAL DATA SOURCES: WEBSITES

- Internet archive's Wayback Machine
  - <https://web.archive.org>
  - Archive of website snapshots since 1996
  - APIs available to find available snapshots
  - Some issues: not all sites archived; some embedded resources not archived; resources may be archived at different times<sup>8</sup>
- OpenWPM: Monthly crawls of top 1million websites, 2015-2018<sup>9</sup>
- Open Observatory of Network Interference (OONI): Blocking and censorship measurements from many vantage points<sup>10</sup>
- CommonCrawl: regular snapshots of 3 billion websites since 2013, HTML pages without embedded resources<sup>11</sup>

<sup>8</sup>A. Lerner, A. K. Simpson, T. Kohno, *et al.*, "Internet Jones and the Raiders of the Lost Trackers: An Archaeological Study of Web Tracking from 1996 to 2016," in *25th USENIX Security Symposium (USENIX Security 16)*, Austin, TX, USA: USENIX, 2016.

<sup>9</sup><https://webtransparency.cs.princeton.edu/webcensus/data-release/>

<sup>10</sup><https://ooni.org/>

<sup>11</sup><https://commoncrawl.org/>

- WhoTracks.me
  - Dataset of tracker domains and who owns them
  - Monthly data starting 2017<sup>12</sup>
- Ad block filter lists: EasyList<sup>13</sup>
- Tracker filter lists: EasyPrivacy
- Anti-ad block filter lists: Anti-Adblock Killer list<sup>14</sup>, ad block warning removal list<sup>15</sup>

<sup>12</sup><https://github.com/cliqz-oss/whotracks.me>

<sup>13</sup><https://easylist.to/>

<sup>14</sup><https://github.com/reek/anti-adblock-killer>

<sup>15</sup><https://easylist-downloads.adblockplus.org/antiadblockfilters.txt>

## PASSIVE TRAFFIC CAPTURE

- Record real live traffic that was not influenced by researchers
  - netflows or full traffic at ISP
  - application-layer traffic (search, email) at search/email provider
- Advantages:
  - Service use by real users, no artificially induced effects
- Challenges:
  - Cannot vary experimental factor => analysis is more limited
  - Data can be difficult to obtain
  - Typically not possible to decrypt HTTPS traffic

## ACTIVE TRAFFIC CAPTURE

- Recording traffic that was initiated by researchers' measurement
- Desktop
  - Instrument browser to record cleartext HTTP, HTML, JavaScript, screenshots
- Mobile
  - Depending on platform, can be hard to instrument apps or OS
  - Can record outgoing network traffic, use VPN permissions/mitmproxy to decrypt HTTPS traffic
- Advantages:
  - Easy to influence input variables, define control/experimental groups
  - Record plaintext traffic at application layer
  - Easier control of noise factors
- Challenges:
  - Traffic may not resemble real user traffic (e.g., crawlers see more third-party requests than humans<sup>16</sup>)

<sup>16</sup>D. Zeber, S. Bird, C. Oliveira, et al., "The Representativeness of Automated Web Crawls as a Surrogate for Human Browsing," in *Proceedings of The Web Conference 2020*, ser. WWW '20, Taipei, Taiwan: ACM, Apr. 2020, pp. 167–178. doi: 10.1145/3366423.3380104.

# TECHNIQUES FOR DATA COLLECTION

---

- Most widely used: Selenium
  - Allows to script web browsing sessions
  - Advantage: full browser including JavaScript execution closely imitates real user experience
  - Many academic frameworks built on top, e.g. openWPM, AdFisher
- openWPM<sup>17</sup>
  - Flexible instrumentation, including browser extensions
  - Data storage
- AdFisher<sup>18</sup>
  - Instrumentation for advertising (ad settings, scraping of ads)
  - Permutation tests

<sup>17</sup>S. Englehardt and A. Narayanan, "Online Tracking: A 1-million-site Measurement and Analysis," in *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, ser. CCS '16, Vienna, Austria: ACM, 2016, pp. 1388–1401. doi: [10.1145/2976749.2978313](https://doi.org/10.1145/2976749.2978313).

<sup>18</sup>A. Datta, M. C. Tschantz, and A. Datta, "Automated Experiments on Ad Privacy Settings," *Proceedings on Privacy Enhancing Technologies*, vol. 2015, no. 1, pp. 92–112, Apr. 2015. doi: [10.1515/popets-2015-0007](https://doi.org/10.1515/popets-2015-0007).

- “Monkey” tools automate user interaction with apps
  - Android UI/Application Exerciser Monkey: creates random input events
  - DroidBot<sup>19</sup>: creates app-specific state-transition model to trigger sensitive app behaviors at higher rate
  - Both can run in emulator or on real device
- Disadvantages:
  - Additional work needed for apps that require login
  - Random events may not trigger all app behaviors

<sup>19</sup>H. Jin, M. Liu, K. Dodhia, et al., “Why Are They Collecting My Data?: Inferring the Purposes of Network Traffic in Mobile Apps,” *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.*, vol. 2, no. 4, 173:1–173:27, Dec. 2018. doi: [10.1145/3287051](https://doi.org/10.1145/3287051).

- Browser extensions, apps
  - Idea: create browser extension/app that offers real benefit to users, hope that users consent to participating in research in exchange (e.g. by contributing some of their data)
  - Examples:
    - FDVT displays value users create for Facebook by viewing/clicking on ads<sup>20</sup>
    - \$heriff allows users to compare prices for products from vantage points in other countries<sup>21</sup>
- Crowdsourcing platforms
  - Idea: pay crowd workers small amounts for short tasks
  - Amazon Mechanical Turk, Crowd Flower, Prolific Academic

<sup>20</sup>A. A. Galán, J. G. Cabañas, Á Cuevas, *et al.*, "Large-Scale Analysis of User Exposure to Online Advertising on Facebook," *IEEE Access*, vol. 7, pp. 11959–11971, 2019. doi: [10.1109/ACCESS.2019.2892237](https://doi.org/10.1109/ACCESS.2019.2892237).

<sup>21</sup>C. Iordanou, C. Soriente, M. Sirivianos, *et al.*, "Who is Fiddling with Prices?: Building and Deploying a Watchdog Service for E-commerce," in *Proceedings of the Conference of the ACM Special Interest Group on Data Communication*, ser. SIGCOMM '17, Los Angeles, CA, USA: ACM, 2017, pp. 376–389. doi: [10.1145/3098822.3098850](https://doi.org/10.1145/3098822.3098850).

# CROWDSOURCING

- Fair compensation for crowd workers
  - 96% of workers on Amazon Mechanical Turk earn less than US federal minimum wage of \$7.25<sup>22</sup>
  - For fair compensation, ensure at least minimum wage for task + required unpaid work (task selection, screening tasks)
- Crowd worker tasks
  - Surveys, labeling images or other data, performing searches or similar to measure personalization based on real user profiles
- Quality control
  - Screening of workers: approval rating, location
  - Task replication: each task is performed by more than one crowd worker, majority vote to filter bad results
  - Controls: subtasks with known results are included in tasks, exclude results from workers who do controls incorrectly

<sup>22</sup>K. Hara, A. Adams, K. Milland, et al., "A Data-Driven Analysis of Workers' Earnings on Amazon Mechanical Turk," in *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, Montréal, Québec, Canada: Association for Computing Machinery, Apr. 2018, pp. 1–14.

# DATA STORAGE

---

- Collected data may be large
- Goals for data storage:
  - Permanent record
  - Fast to access for data analysis
- US Library of Congress recommends platform-independent, character-based formats based on well-known schemas, such as CSV and SQLite
- Combination of storage systems can help to achieve performance and permanence

# SUMMARY

---

## SUMMARY: METHODS FOR DATA COLLECTION

- Ethical considerations for data collection
- Techniques to minimize harm to individuals and organizations
- Data sources
  - Archival sources
  - Passive traffic capture
  - Active traffic capture
- Automation of browsers and apps
- Crowdsourcing
- Data storage

## ABOUT THIS SLIDE DECK

- These slides are designed to accompany a lecture based on the textbook “Auditing Corporate Surveillance Systems: Research Methods for Greater Transparency” by Isabel Wagner, published in 2022 by Cambridge University Press.
- Except where otherwise noted (e.g., logos and cited works) this slide deck is Copyright © 2017-2022 Isabel Wagner
- The slides are free to use for non-commercial purposes, provided that the source of the slides, i.e. the textbook and its companion website, are cited appropriately
- Please leave this slide intact, but indicate modifications below.
  - Version 2022-04
  - Improved version for release on book website (Isabel Wagner)
- Updated versions of the original slide deck are available online: [corporatesurveillance.org](https://corporatesurveillance.org)