# Forensics, Data and AI

1 - Background

2 - Recent research projects

3 - Overview of teaching experience

4 - Q&A

February 11, 2025 • Simson Garfinkel

# Background

# Simson Garfinkel: Background and Bio

## Career #1: Science writer (1985-)
- Newspapers, Magazines, Books
- Most recently: History of computing — Technology Review & CACM

## Career #2: Entrepreneur (1992-)
- SGAI — 1992-1993 — Commercialized AI approach from MIT Media Lab
- Vineyard.NET — 1995-2002 — ISP on Martha's Vineyard
- Sandstorm Enterprises — 1998-2001, 1998-2006 (board) — Security tools
- Broadband2Wireless — 2000-2001 — Wireless ISP

## Career #3: CS Researcher (1985-87, 90-91, 2002-)
- MIT Media Lab 1985-1987, 90-91
- MIT PhD  2003-2005
- Harvard SEAS CRCS — 2005-2006
- Naval Postgraduate School — 2006-2014
- NIST — 2015-2016

## Career #4: Government Innovator
- US Census Bureau — 2017-2021 — Differential privacy for the 2020 Census
- US DHS — 2021-2022 — DHS Chief Data Officer – Data Inventory

**I'm currently Chief Scientist at BasisTech LLC, a startup accelerator in Somerville.**



3

# Major research projects

## Bulk_extractor (2006-2020) — processing bulk data for digital forensics

- Developed novel approaches to extracting formatted data from bulk data.
- $15 million in funding from DOD, DHS, NSF, FBI and others.
- Produced a widely used digital forensics tool, multiple peer-reviewed publications, hundreds of training videos on YouTube.

## Digital Corpora (2006-) — realistic data for digital forensics research and education

- Funded by NSF, NIST and Amazon
- Major contribution to DARPA SafeDocs project.
- 26TB of open use data — hard drive images, worked "scenarios" (with solutions), packets, files

## Disclosure Avoidance System for the 2020 Census (2017-2021)

- Most complex deployment of differential privacy to date.
- Lead computer scientist for the project. Developed experimental and production framework
- Instituted modern software development practices.

## PlantTracer (2023-) — Using computer vision to watch plant movement.

- An online environment for high school and college students.

# What's next:
# Building a regional tech abuse center

## Tech abuse is a major threat facing many users.

- Like "physical abuse" and "metal abuse," but with technology.

## Examples:

- One partner configures the other's phone with "family controls."
- Using router's "family controls" to monitor (or disable) internet access.
- Eavesdropping on emails sent to friends (or lawyers)
- Turning off cell phone service as "punishment."
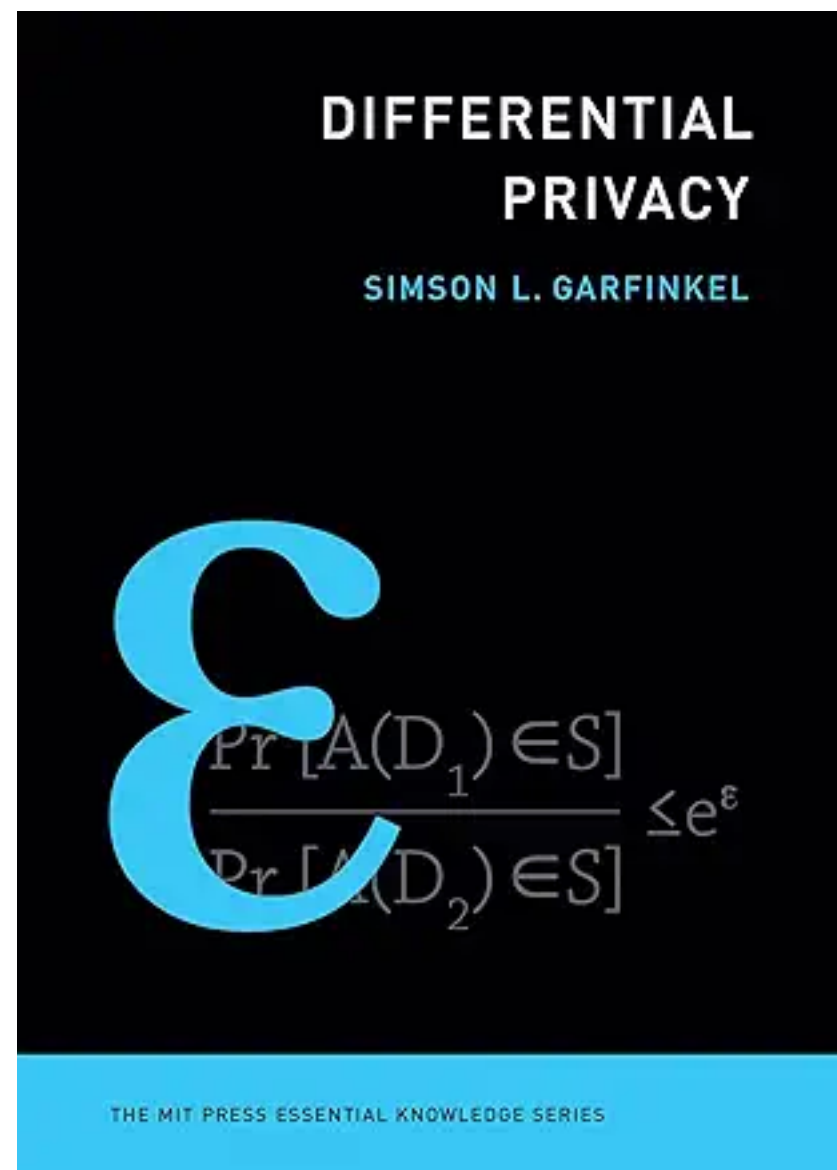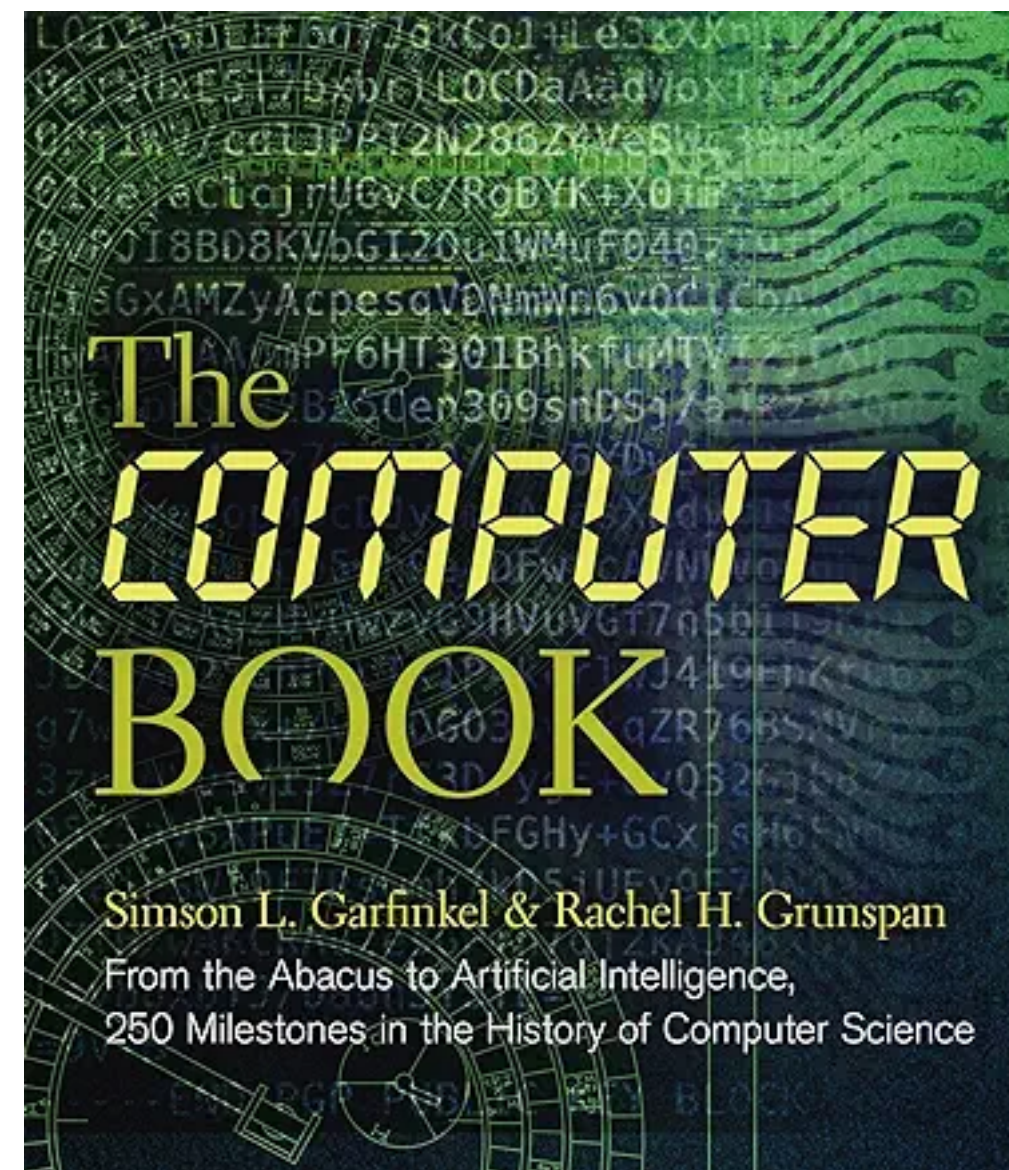- Hidden cameras

## What can we do?

- **CS research:** Identify patterns for defense and detection
- **Social research:** Identify scale of the problem
- **Policy research:** Inform policymakers & propose solutions
- **Tech clinic:** Help people clear their devices
- **Legal clinic:** Help with protective orders
- **Community building:** Organize academics & activists in the Boston area.



6 Oct. 2023
OPINION

IoT & Personal Devices    Personal Privacy

# Privacy professionals need to be aware of tech abuse

Related stories

Opinion: Limiting end-to-end encryption to only verified users

Notes from the IAPP Canada: Balanced thinking needed on biometrics

Online consent: How can it be made valid in practice?

Simson Garfinkel
Contributor
CIPP/US
7 Minute Read

Features designed to improve privacy and protect children in online services, apps and networked devices also make it easier for abusers to maintain control in abusive relationships.

"Ever since caller ID and GPS became part of our lives, we've known that digital technologies can be used by abusers to harm or track their victims, and that's only become more complicated and more prevalent
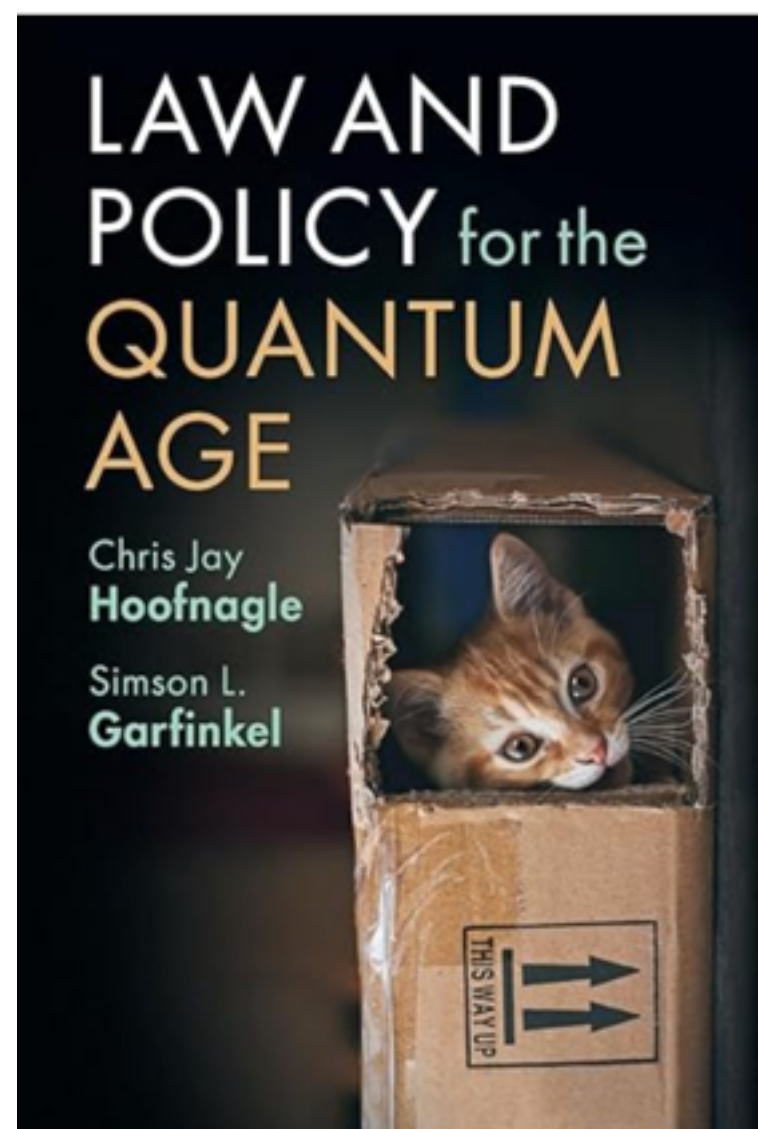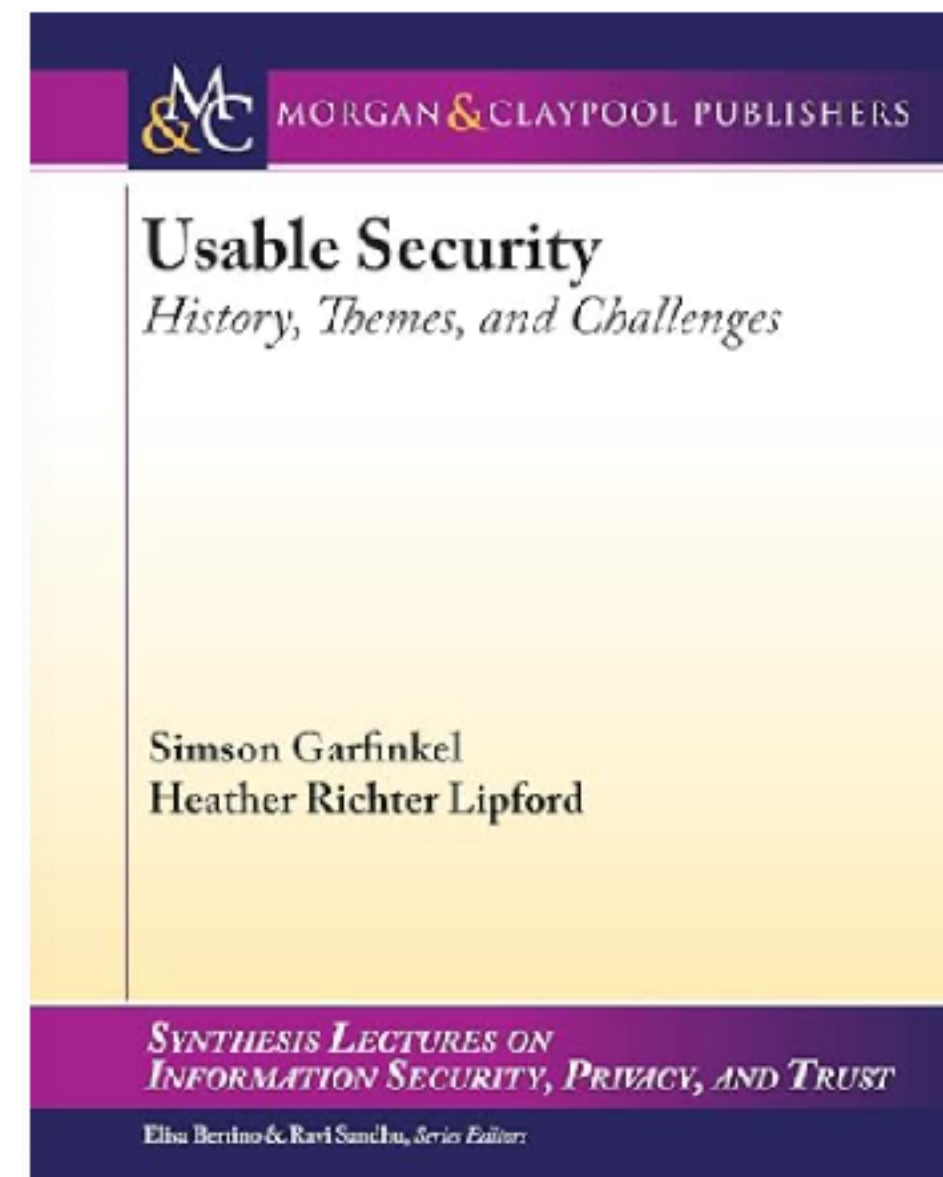
# Recent books

Differential Privacy
MIT Press, 2025

The Computer Book: From the Abacus to Artificial Intelligence, 250 Milestones in the History of Computer Science (Sterling Milestones), by Simson L. Garfinkel and Rachel H. Grunspan. 2018 (Sterling)

Law and Policy for the Quantum Age,
Chris Jay Hoofnagle and Simson L. Garfinkel,
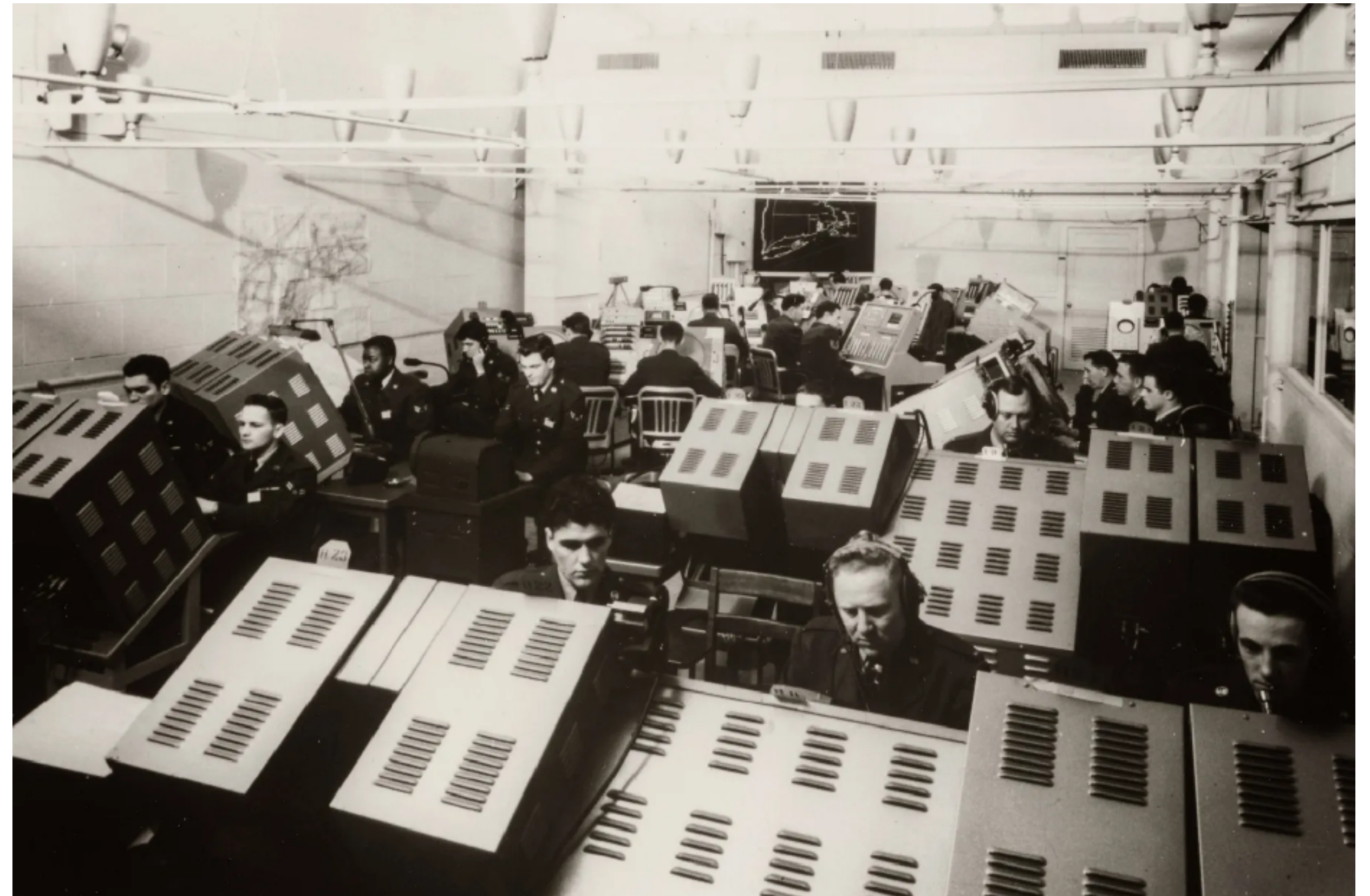2021 (Cambridge)

Usable Security: History, Themes, and Challenges,
Simson Garfinkel and Heather Lipford, 2014.
(Morgan & Claypool, part of the Synthesis Lectures on Information Security, Privacy and Trust series.)

# Recent articles in the history of computing and MIT.

**History of Computing, Technology and MIT**
1. Garfinkel, S. MIT's (mostly) secret society, December 23, 2024
2. Garfinkel, S. This is MIT and yes, we have bananas, August 27, 2024
3. Garfinkel, S. Editor of The Tech becomes president of MIT, Technology Review, June 25, 2024
4. Garfinkel, S. How Technology Review got its start, Technology Review, January 4, 2024
5. Garfinkel, S. MIT's First Divorce (how MITRE was created and got its name), Technology Review, June 27, 2023
6. Garfinkel, S. Cold Trick Indeed (dorm room set up on the Charles, 1985), Technology Review, December 19, 2022
7. Garfinkel, S. How an MIT Marxist weathered the Red Scare (Dirk Struik), Technology Review, June 29, 2022
8. Garfinkel, S. In praise of the Feistel network (Horst Feistel '37), Technology Review, April 27, 2022
9. Garfinkel, S. The man no one knows who changed Boston (Charles Hayden), Technology Review, February 23, 2022
10. Garfinkel, S. 5 MIT patents that changed computing, Technology Review, February 23, 2022
11. Garfinkel, S. Walker and the "Indian Question:" Before arriving at MIT, Francis Amasa Walker had twice led the US Census—and helped justify the troubling US policy of containing Native Americans on reservations. Technology Review, August 24, 2021
12. Garfinkel, S. Tomorrow's computer, yesterday. Four decades ago at Endicott House, an MIT professor convened a conference that launched quantum computing. Technology Review, April 27, 2021
13. Garfinkel, S. Punching In: Bored teaching at MIT, Herman Hollerith left to launch the information age for the US Census. Technology Review, August 18, 2020
14. Garfinkel, S. Everything is a Punch Card. ;login:, Fall 2020
15. Garfinkel, S. The Tricky Cryptographic Hash Function. ;login:, Winter 2020
16. Garfinkel, S. Shafi Goldwasser, Technology Review, August 21, 2019
17. Garfinkel, S. Radia Perlman '73, SM '76, PhD '88, Technology Review, August 21, 2019
18. Garfinkel, S. The Geek (Chris Schmandt), Technology Review, April 24, 2019



Garfinkel, S. MIT's First Divorce (how MITRE was created and got its name), Technology Review, June 27, 2023

# Outline for the rest of this talk

## Research projects:

- Bulk_extractor (2006-2020) — processing bulk data for digital forensics
- Digital Corpora (2006-) — realistic data for digital forensics research and education
- Disclosure Avoidance System for the 2020 Census (2017-2021)
- PlantTracer (2023-) — Using computer vision to watch plant movement.

## Two fun classes

- "Defense Against the Dark Arts (cybersecurity edition)"
- Artificial Intelligence, Internet of Things, and Cybersecurity,

## What's Next: Two big projects for 2025-2026

- Tech Abuse Research Center
- AI for Data Management

# Research Projects

Bulk_extractor (2006-2020) — processing bulk data for digital forensics

Digital Corpora (2006-) — realistic data for digital forensics research and education

Disclosure Avoidance System for the 2020 Census (2017-2021)

PlantTracer (2023-) — Using computer vision to watch plant movement.

# In 2003, I bought 200 used hard drives

The goal was to find drives that had not been properly sanitized.

## First strategy:

- DD all of the disks to image files
- run **strings** to extract printable strings.
- **grep** to scan for email, CCN, etc.
  - *VERY SLOW!!!!*
  - *HARD TO MODIFY!*

## Second strategy:

- Use SBook approach!
- Read disk 1MB at a time
- Pass the *raw disk sectors* to flex-based scanner.
- Big surprise: scanner didn't crash!

# Simple flex-based scanners required substantial post-processing to be useful

Techniques include:

- Additional validation beyond regular expressions (CCN Luhn algorithm, etc).
- Examination of feature "neighborhood" to eliminate common false positives.



The technique worked well to find drives with sensitive information.

Could it be of use in digital forensics?

# Between 2005 and 2008, I interviewed law enforcement officers regarding their use of forensic tools.

Law enforcement officers wanted a *highly automated* tool for finding:

- Email addresses
- Credit card numbers (including track 2 information)
- Search terms (extracted from URLs)
- Phone numbers
- GPS coordinates
- EXIF information from JPEGs
- All words that were present on the disk (for password cracking)



https://www.americanscientist.org/article/digital-forensics

# I also learned about their requirements for the user experience.

The tool had to:

- Run on Windows, Linux, and Mac-based systems
- Run with *no* user interaction
- Operate on raw disk images, split-raw volumes, E01 files, and AFF files
- Allow user to provide additional regular expressions for searches
- Automatically extract features from compressed data such as gzip-compressed HTTP
- Run at maximum I/O speed of physical drive
- Never crash

# Starting in 2008, I made a series of limited releases

- January 2008 — Created Subversion Repository
- April 2010 — Initial public release - 0.1.0
- May 2010 — Initial multi-threading release - 0.3.0
  - *Each thread runs in its own process*
- Sept. 2010 — Stop lists - 0.4.0
- Oct. 2010 — Context-based stop-lists - 0.5.0
- Dec. 2010 — Switch to POSIX-based threads — 0.6.0
- Dec. 2010 — Support for Windows HIBERFIL.SYS decompression — 0.7.0
- Jun. 2010 — First 1.0.0 Release

Tool capabilities result from substantial testing and user feedback.

Moving technology from the lab to the field was challenging:

- Must work with evidence files of *any size* and on *limited hardware.*
- Users can't provide their data when the program crashes.
- Users are *analysts* and *examiners*, not engineers.

San Luis Obispo is
"the happiest
place in America"

Watch the video to find out why National
Geographic named San Luis Obispo the top spot.

http://www.sanluisobispovacations.com/

15

## District Attorney filed charges against two individuals:

- Credit Card Fraud
- Possession of materials to commit credit card fraud.



## Defendants:

- Arrested with a computer.
- Expected to argue that defends were unsophisticated and lacked knowledge.

## Examiner given 250GiB drive *the day before preliminary hearing.*

- Typically, it would take several days to conduct a proper forensic investigation.

# bulk_extractor found actionable evidence in 2.5 hours!

Examiner given 250GiB drive *the day before preliminary hearing.*



Bulk_extractor found:

- Over 10,000 credit card numbers on the HD (1000 unique)
- Most common email address belonged to the primary defendant (possession)
- The most commonly occurring Internet search engine queries concerned credit card fraud and bank identification numbers (intent)
- Most commonly visited websites were in a foreign country whose primary language is spoken fluently by the primary defendant.

Armed with this data, the DA was able to have the defendants held.

# So why was bulk_extractor a success?

Open source  ⬅ *Not the whole story!*

- Government users could download it from the Internet and use it immediately.
  - *Existing authorities allowed for open source digital forensics tools to be used on specific systems.*

Plug-in architecture ⬅ *Not the story at all*

- Allowed students to create modules for student projects.
- Successful projects could be adopted into the main branch.

**Delivered results that no other program could deliver**
- Recursive analysis of coded and compressed data.
- Recovery of data from file fragments.

**Did not compete with existing software — and other software did not compete with it!**
- Because it was free, the only cost to using bulk_extractor was time and computational resources.
- Eliminates the need to implement a complete forensic stack — BE does not compete with existing tools.
  - *In fact, at least one existing tool incorporated BE into its analysis pipeline.*

**Super easy-to-use!**

# Stream-Based Disk Forensics with bulk_extractor.

Scan the disk from beginning to end; do your best.



0                         **1TB**

**3 hours, 20 min
to *read* the data**

1. Read all of the blocks in order.
2. Look for information that might be useful.
3. Identify & extract what's possible in a single pass.

# bulk_extractor splits the disk into 16M "pages" (blocks) and processes each page independently.

**THREAD1**

| 0-15M | → scan_email → | XYZ@COMPANY.COM |

**THREAD2** 16-31M → scan_email → ABC@company.com

**THREAD3** 32-47M → scan_email → DEF@company.com

This finds obvious email addresses in bulk data:

```
a097 83a1 ed96 26a6 3c69 3d0f 750a 2399    ......&.<i=.u.#.
a2b5 bea7 692f 5847 a38a dd53 082c add5    ....i/XG...S.,..
5061 b64c 721d 864b 90b6 b55f bb04 735c    Pa.Lr..K..._..s\
9448 6730 5453 df64 813e b603 5795 2242    .Hg0TS.d.>..W."B
e9c8 7454 7322 7cdc b60e 97af 2f64 2728    ..tTs"|...../d'(
3cfb 84bd 2a84 2dfe 50ea 5935 c349 1513    <XYZ@COMPANY.COM
a9e9 e92c a3f8 6e46 0530 8a88 c7a2 5d2b    ...,..nF.0....]+
d89d 77cc fe1e f637 f3f3 d0af 1b47 c09b    ..w....7.....G..
```

# bulk_extractor examines every byte to see if it is the beginning of an "encoded" region.

Once the region is found, it's decoded, then processed.

THREAD1

0-15M

zlib? → scan_email → XYZ@COMPANY.COM

RAR? → scan_email

HIBER? → scan_email

BASE64? → scan_email

Compressed blocks are discovered, decompressed, and recursively processed.

This "optimistic" approach also recovers data from fragments of files.

# Students saw that open source made innovation easier!
# (About half of these videos were created by students)

# Students saw that open source made innovation easier!
(About half of these videos were created by students)

# I used bulk_extractor as a platform for research and education

## 2012 — Statistical sampling breakthrough

- US Patent 8,433,959 granted April 30, 2013

## By 2016 bulk_extractor was

- FBI approved tool
- Incorporated into two products — one commercial, one GOTS (government-off-the-shelf).
- Widely used in digital forensics education.
- Incorporated into multiple digital forensics boot DVDs.

## Bulk_extractor — helped teach students to innovate

- Showed students how to introduce advanced technology into US Government agencies that were resistant to change.
- Provided a testbed for students to develop their own modules.
- Showed how to pitch sponsored research and transition it to the field.



October 1, 2012
https://crcs.seas.harvard.edu/event/simson-l-garfinkel-digital-forensics-innovation-searching-terabyte-data-10-minutes

# Digital forensics tools require constant maintenance.
## OS Creep • Language creep • Forensic science progress • O&M (operations & maintenance) "tail"

## From 2018-2021 I upgraded BE:

- Moved from C++ STL to C++17

- Added CI/CD testing on GitHub

- Improved multi-threading for modern CPUs.

- Wrote an article for *ACM Queue* and *Communications of the ACM*.

# Research Projects

Bulk_extractor (2006-2020) — processing bulk data for digital forensics

Digital Corpora (2006-) — realistic data for digital forensics research and education

Disclosure Avoidance System for the 2020 Census (2017-2021)

PlantTracer (2023-) — Using computer vision to watch plant movement.

# Digital forensics research and education had a data problem in 2009.

Digital forensics practitioners must be able to

- analyze *any* digital data,
- from any computer,
- that has ever been used,
- anywhere.

The data problem: getting data that are **ecologically valid**

- **representative** of the diversity of systems found on computers collected by law enforcement and defense practitioners.
- **complex enough** to present students and researchers with more than toy problems.
- **simple enough** that the problems can be solved in *hours or days*, rather than *weeks or months*.

The solution

- Get students to create complex scenarios *as a learning exercise.*
- Allow free downloads of the dataset.
- Track usage through the "teacher's solutions."

# I created the Digital Corpora —
# a collection of complex digital artifacts for forensics education and tool testing.

https://digitalcorpora.org/

## Initial funding:

- NIST/NPS Inter Agency Agreement
- NSF Grant No. 0919593

## Today:

- Scenarios and data contributed by cybersecurity programs and practitioners all over the world
- Corpus hosting by Amazon's Open Data Sponsorship Program

# The corpus has many scenario-based digital artifacts.

## Complex, deep datasets

- Scripted scenarios.
- Multiple characters with clearly defined motivations
- Specific challenges for the investigator to uncover
- Multiple problems that require different levels of skill and analysis to solve
- Created in "real-time" over weeks or months
- "Teachers guides" and "solutions" are available for many of the datasets.

## Multi-modality

- Disk images
- Cell phone images
- Memory dumps
- Log files from servers
- Packet dumps (wiretaps)

# A few scenarios in the corpus available for download

A "Lone Wolf" who becomes self-radicalized on YouTube and plans a school shooting.

- He was turned in by his brother.
- You have the laptop
- https://downloads.digitalcorpora.org/corpora/scenarios/2018-lonewolf/

A macOS/iOS terrorist recruitment scenario with multiple personas and international travel

- Picked up by FBI
- You have the Mac and iPod Touch backup
- https://downloads.digitalcorpora.org/corpora/scenarios/2019-tuck/

A planned defacement of art at the DC National Gallery by a direct action group, combined with a nasty divorce proceeding.

- You have disk images, phone images, captured packets, and a bungled wiretap
- https://downloads.digitalcorpora.org/corpora/scenarios/2012-ngdc/

+ many others contributed by educators around the world.

# Constructed, scenario-based artifacts are better for research and education.

## No privacy-sensitive data! No PII!

- Computer users are not real people, they are personas

## No pornography! No illegal content!

- We know that there's no pornography in the data
- Especially an issue with students under 18 years old

## No child exploitation scenarios!

- CSAM scenarios are a big turn-off!

## There are solutions!

- Solutions are distributed on the website as encrypted PDFs
- Decrypt keys are available on a case-by-case basis to faculty at accredited institutions, law enforcement, and partners

# GOVDOCS1M — The first ecologically valid "files" corpus.

Developed in 2008, a corpus of 1 million files downloaded from US Government web servers.

- US Government websites to avoid copyright issue.

## Includes:

- Image formats (JPEG, TIFF, PNG, etc)
- Document formats (PDF, MSOffice)
- Text files
- Log files
- SQL dumps

## At the time, this let me teach…

- Approaches for working within the copyright law
- How to handle legal missteps
- Scientific principles of reproducibility

… by sharing the issues with students



(one of many research articles have used the corpus.)

# GOVDOCS was the seed for the DARPA SafeDocs program

Goal of SafeDocs: build an exploit-proof PDF reader using formal methods.



2018-2023

When SafeDocs shut down, DARPA donated 8M PDFs to the Digital Corpora

- SafeDocs became open data!

We now have 24TB of data…

- We had to be entrepreneurial in dealing with storage requirements!
- Today we are hosted by Amazon's Open Data program.
- With minimal copyright and privacy issues, this Internet snapshot can power the creation of tools for the digital humanities.



May 16, 2023

# Digital Corpora: Educational Impact

Solutions to the scenarios are distributed as an encrypted PDF.

Faculty can request the decryption key; so far over 325 have.

We surveyed those requesting the key; 92 completed our survey.

How did you learn about the digitalcorpora.org website?

Which materials did you use?

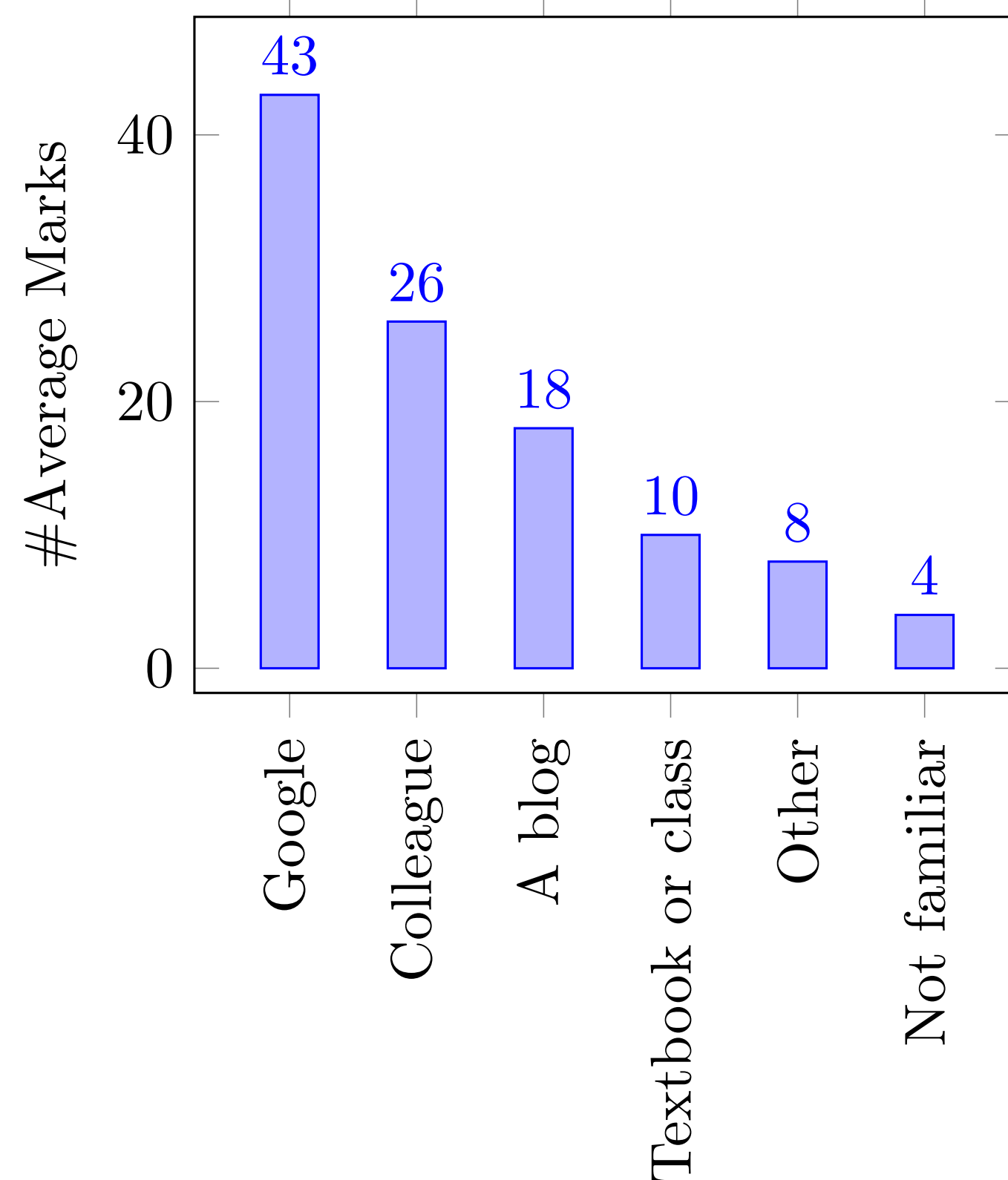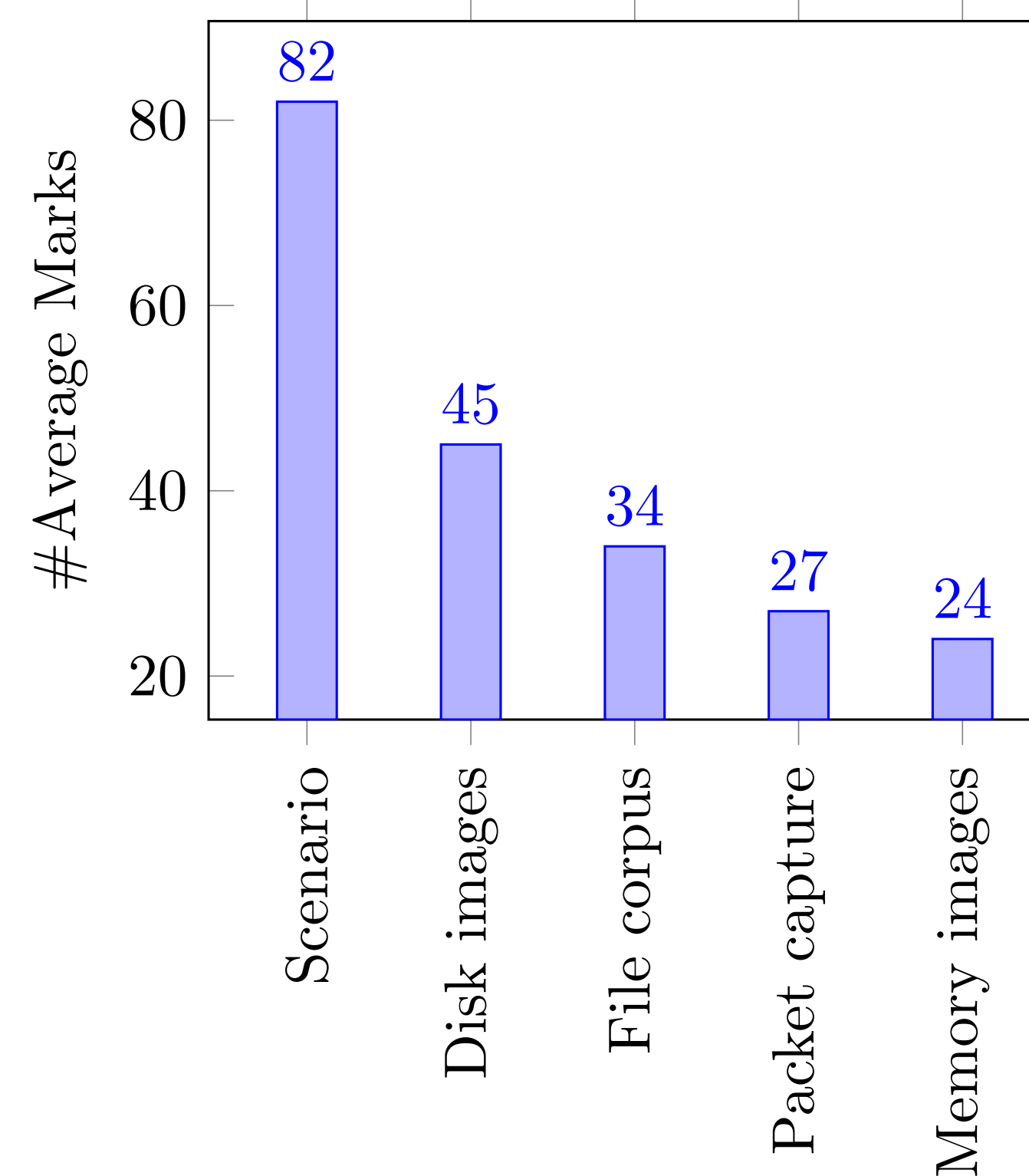|  | Guide | Count |
|---|---|---|
|  | Lone Wolf Scenario | 55 |
|  | M57 patents | 51 |
| For which did you download a teacher's guide? | Nitroba university | 34 |
|  | DC art gallery | 12 |
|  | Narcos | 3 |
|  | M57 Jean | 1 |

|  | | |
|---|---|---|
|  | Education and Training | 81 |
|  | Tool testing | 22 |
|  | R&D for new tools and features | 13 |
| We used the datasets for: | Practice to prepare for casework | 11 |
|  | Proficiency testing | 9 |
|  | Research on DF investigative practices | 1 |
|  | Analysis and exploratory research | 1 |

# Research Projects

Bulk_extractor (2006-2020) — processing bulk data for digital forensics

Digital Corpora (2006-) — realistic data for digital forensics research and education

Disclosure Avoidance System for the 2020 Census (2017-2021)

PlantTracer (2023-) — Using computer vision to watch plant movement.

## In 2017 I started at the US Census Bureau with the mission of modernizing the bureau's "Disclosure Avoidance."

"Disclosure Avoidance" — aka "Statistical Disclosure Limitation"
- aka "Privacy Preserving Data Publishing" and "Privacy Preserving Data Analysis."

Chief of the Center for Disclosure Avoidance (2017) — GS15 management

Senior Computer Scientist for Confidentiality and Data Access (2017-2021) - ST
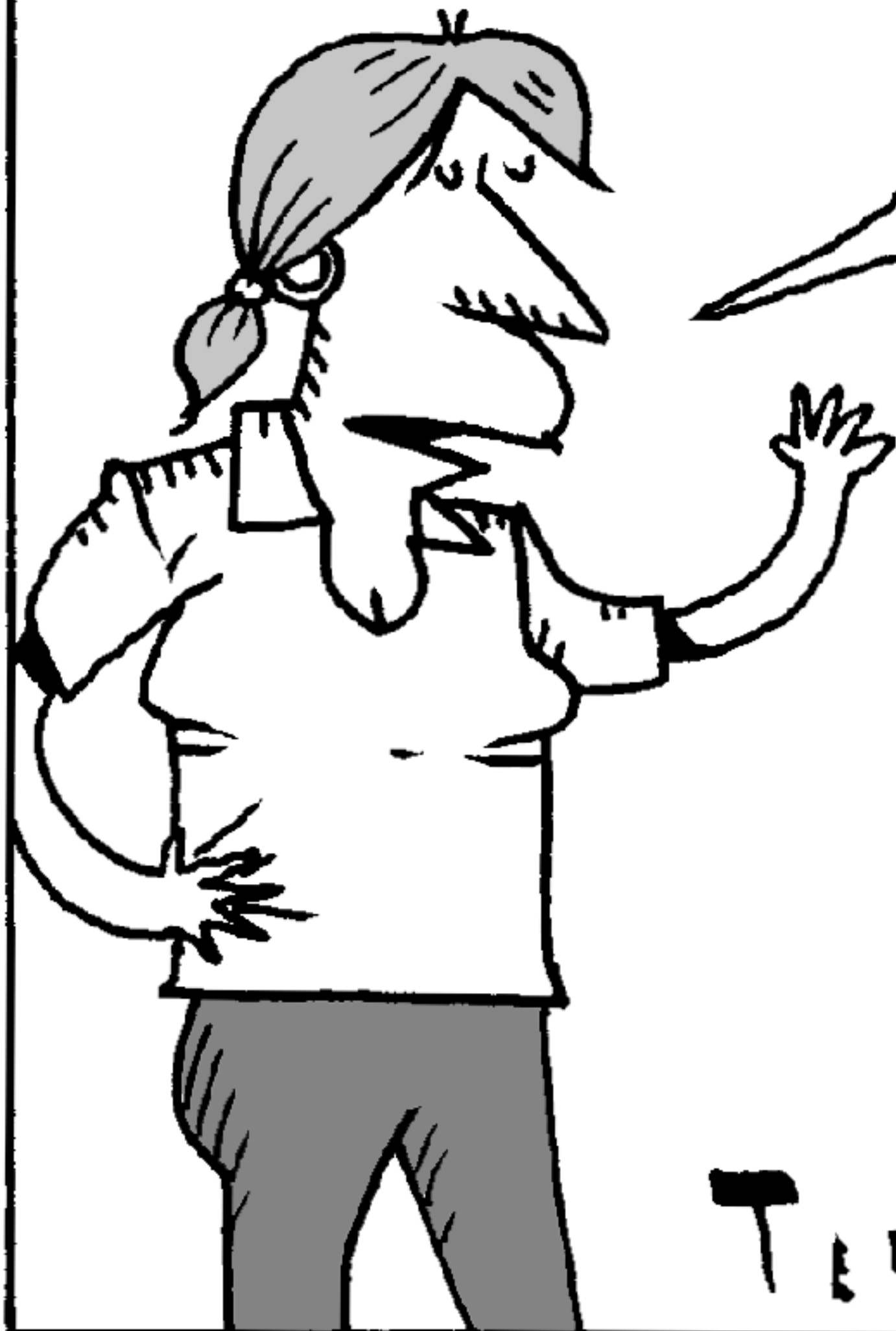
Chair, Disclosure Review Board (2017-2019)

Accomplishments:
- Brought differential privacy to the 2020 Census

- Modernized DRB

- Laid ground work for modernizing privacy protections in American Community Survey, Federal Housing Survey, Economic Census, and many other data products.

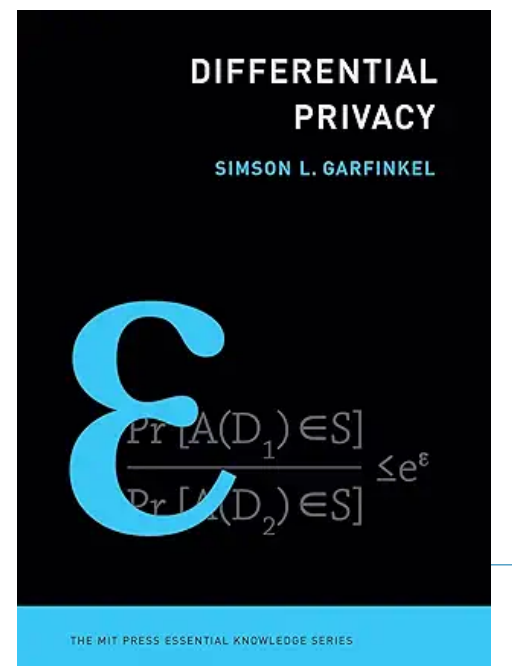- Educated Census senior leadership on differential privacy.

# Data flow in the 2020 Census (Original vision)

**DRF** — Decennial Response File

**CUF** — Census Unedited File

**CEF** — Census Edited File

**DAS** — Global Confidentiality Protection Process / Disclosure Avoidance System

**MDF** — Microdata Detail File

Pre-specified tabular summaries: **PL94-171, SF1**

**Special tabulations** and post-census research

Privacy-loss Budget, Accuracy Decisions

Confidential data

Public data

The top-down mechanism:
Each histogram provides statistical accuracy to those underneath.

# Research Projects

Bulk_extractor (2006-2020) — processing bulk data for digital forensics

Digital Corpora (2006-) — realistic data for digital forensics research and education

Disclosure Avoidance System for the 2020 Census (2017-2021)

PlantTracer (2023-) — Using computer vision to watch plant movement.

# Plant Tracer  (w/ Eric Brenner, Pace Univ.)

- Converted iOS app to web-based app
- Designed learning management system for sharing and analyzing video
- Developing open source web-based video capture system & ESP32-based camera.



## Movie #6286 is traced!

| Marker | Name | Location (pixels) | Location (cm) |
|---|---|---|---|
| ● | Apex | (137,81) | n/a |
| ● | Ruler 0mm | (238,97) | n/a |
| ● | Ruler 20mm | (243,138) | n/a |

**Download trackpoints**

# Plant Tracer  (w/ Eric Brenner, Pace Univ.)

- Converted iOS app to web-based app
- Designed learning management system for sharing and analyzing video
- Developing open source web-based video capture system & ESP32-based camera.





Movie #6286 is traced!

| Marker | Name | Location (pixels) | Location (cm) |
|--------|------|-------------------|---------------|
| ● | Apex | (137,81) | n/a |
| ● | Ruler 0mm | (238,97) | n/a |
| ● | Ruler 20mm | (243,138) | n/a |

Download trackpoints

# Two fun classes

# Defense Against the Dark Arts

Course goal

Provide students with a basic understanding of how they can get hacked and how to defend themselves

## First-year seminar course.

- 12 first semester students, ~ 6 with background in computing
- 2.5 hours each week
- 90 minutes discussion of reading
- 90 minute "lab"

## Books:

- *Fancy Bear Goes Phishing: The Dark History of the Information Age in Five Extraordinary Hacks* by Scott J. Shapiro (2023) (entire book)
- *This Is How They Tell Me the World Ends: The Cyberweapons Arms Race* by Nicole Perlroth (2021) (selected chapters)
- *Crypto: How the Code Rebels Beat the Government Saving Privacy* in the Digital Age by Steven Levy (2002) (selected chapters)

## Technical Papers

- Light Commands: Laser-Based Audio Injection on Voice-Controllable Systems
- An Empirical Analysis of the Commercial VPN Ecosystem
- Security at the End of the Tunnel: The Anatomy of VPN Mental Models Among Experts and Non-Experts in a Corporate Context
- Extending a Hand to Attackers: Browser Privilege Escalation Attacks via Extensions

# Course Outline

Week 1 — A Hacker's Bestiary 1: The Great Worm and Computer Viruses

Week 2 — A Hacker's Bestiary 2: The Wizards

Week 3 — Magic Words (passwords)

Week 4 — The Oldest Dark Art (Secret Key Cryptography)

Week 5 — New Directions in Cryptography (Public Key Cryptography)

Week 6 — Keeping Your Data Out of Azkaban (Ransomware and Malware)

Week 7 — Finding Spells in the Library: Read the Computer Security Literature

Week 8 — Your Cloak of Invisibility (VPNs)

Week 9 — Sneaking through the forest (Browser Privacy, and Private Browsing Mode)

Week 11 — Securing your broom (Browser Extensions)

Week 12 — Transfiguration (phishing)

Week 13 — There's magic in SoK (Tech abuse SoK)

# Labs

week 1 - lab - command line

week 2 - lab - hashing

week 3 - lab - password cracking

week 4 - lab - encryption

week 5 - lab - digital signatures

week 6 - lab - Finding Hidden Data

week 7 - lab - voice commands

week 8 - lab - HTTP Toolkit

week 9 - traveling around the world with ping and traceroute

week 10 - Browser Extensions

**Each lab designed to be accomplished in 45 minutes by first-year students with varying levels of background.**

**Everybody wins!  (Everybody gets something out of the lab.)**

# Other significant courses I've taught

## This semester:

- Critical Thinking in Data Science (Data Science Ethics) — Harvard Data Science Master's Program
- Artificial Intelligence, Internet of Things, and Cybersecurity — Harvard Division of Continuing Education

## In DC:

- Mac Forensics, Cloud Forensics, Document Forensics — George Mason University Digital Forensics Master's Program
- Data Science Ethics — George Washington University Data Science Program
- Massive Data Analytics, Data Privacy — Georgetown University

## Naval Postgraduate School:

- Information Crime, Law and Ethics
- Data Fusion with Online Information Systems
- Network Security
- Java as a Second Language
- Special topics in computer security
- Advanced Computer Architecture
- Automated Document and Media Exploitation

# What's next

**Tech Abuse Center** — A multidisciplinary research center and clinic to address the issue of technology-enabled domestic abuse.

**Embedded Forensics** — Deep-dive into the embedded microelectronics ecosystem.

# Q&A