# Driving Innovation from the University to Government and the Enterprise: Tricks, Traps and Techniques

Monday, June 3, 2024

Simson Garfinkel

These slides can be downloaded from https://simson.net/ref/2024

# Abstract

Innovation and entrepreneurship in science and technology requires seizing new ideas, identifying their distinctive core, and finding enthusiastic users and customers (**tricks**).

It also requires avoiding mistakes that can sideline amazing technology and promising careers (**traps**).

In this talk, I'll present some tricks I've learned, traps that I've been fortunate enough to avoid or escape, and which approaches I think can be taught and generalized (**techniques**) based on my career in academia, government, and several startups.

# Tricks

Embrace the leading edge.

Embrace open source & open data.

# 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 Innovation
- US Census Bureau — 2017-2021 — Differential privacy
- US DHS — 2021-2022 — DHS Data Inventory

**I'm currently Chief Scientist at BasisTech LLC, a startup accelerator in Somerville.**
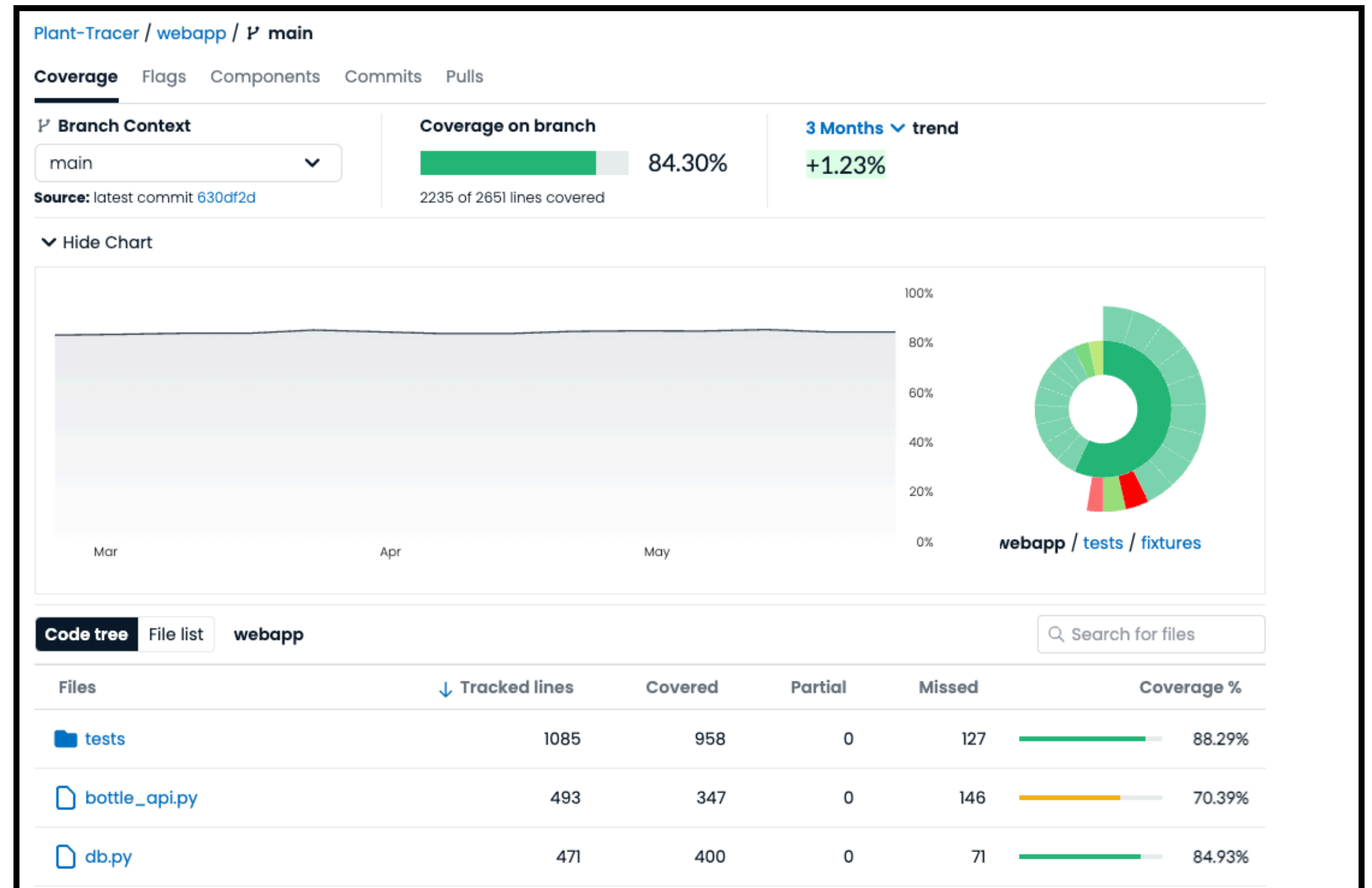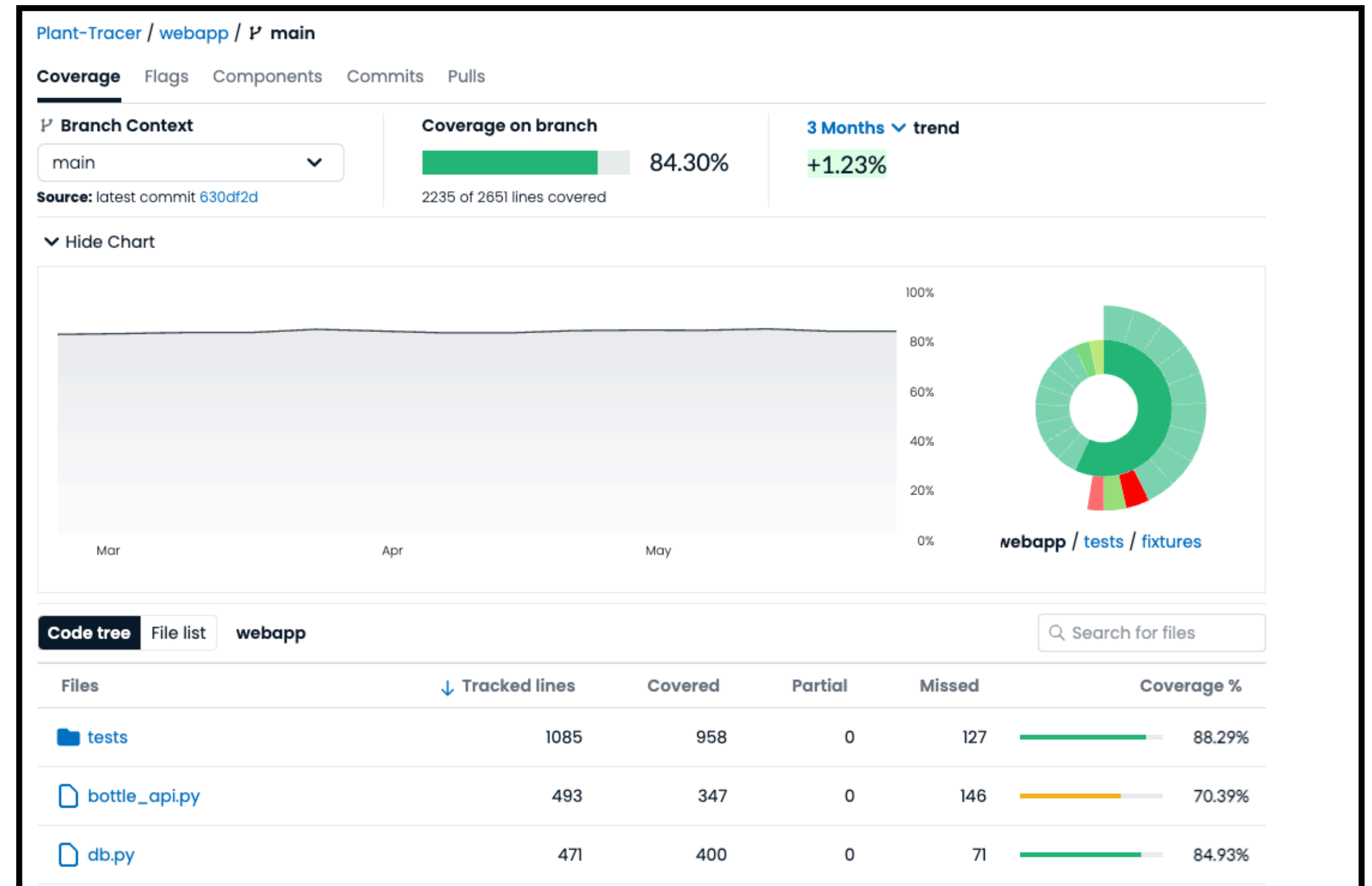
# Trick #1 — Embrace the leading edge

Teach and embrace the modern programming languages & development tools.

Example: AWS in Georgetown "massive data" course (2016-2017)
- Student got a job with AWS because he learned the stack in my class.

Example: Plant Tracer Project (w/ Pace Univ.)
- Computer vision for tracking plant movement
- Unit tests & code coverage
- Function-as-a-service (AWS Lambda)



Plant-Tracer / webapp / ⅄ main

**Coverage**  Flags   Components   Commits   Pulls

⅄ **Branch Context**                    **Coverage on branch**           3 Months ∨ **trend**
main ∨                                   ▇▇▇▇▇▇ 84.30%                   +1.23%
**Source:** latest commit 630df2d        2235 of 2651 lines covered

∨ Hide Chart

| Files | ↓ Tracked lines | Covered | Partial | Missed | Coverage % |
|---|---|---|---|---|---|
| 📁 tests | 1085 | 958 | 0 | 127 | 88.29% |
| 📄 bottle_api.py | 493 | 347 | 0 | 146 | 70.39% |
| 📄 db.py | 471 | 400 | 0 | 71 | 84.93% |

# Trick #1 — Embrace the leading edge

Teach and embrace the modern programming languages & development tools.

Example: AWS in Georgetown "massive data" course (2016-2017)

- Student got a job with AWS because he learned the stack in my class.

Example: Plant Tracer Project (w/ Pace Univ.)

- Computer vision for tracking plant movement
- Unit tests & code coverage
- Function-as-a-service (AWS Lambda)

# Trick #2 — Embrace open source and open data

## Bulk_extractor - Open source digital forensics tool

Initial development: 2006-2010

First deployment: 2010

Version 2.0 rewrite: 2018-2022 (C++20)

"CV impact:" 2 articles



156 citations      0 citations (14K DL downloads)

Real world impact: education, law enforcement & defense

## Digital Corpora - Open data set for digital forensics

Initial development: 2006-2010

First deployment: 2007-

Migration to Amazon Open Data program: 2020

"CV impact:" 1 article



482 citations

Real world impact: digital forensics education & research

#1
Innovating with open source

# Stream-Based Disk Forensics:
## 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.

NPS Presentation from 2011-06-14

**THREAD1**

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

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

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

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

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

In 2011,
I didn't stress that bulk_extractor was open source.

Open source was critical to bulk_extractor's success.

# Because bulk_extractor was open source, it was widely adopted

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



13

# 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 innovation and entrepreneurship.

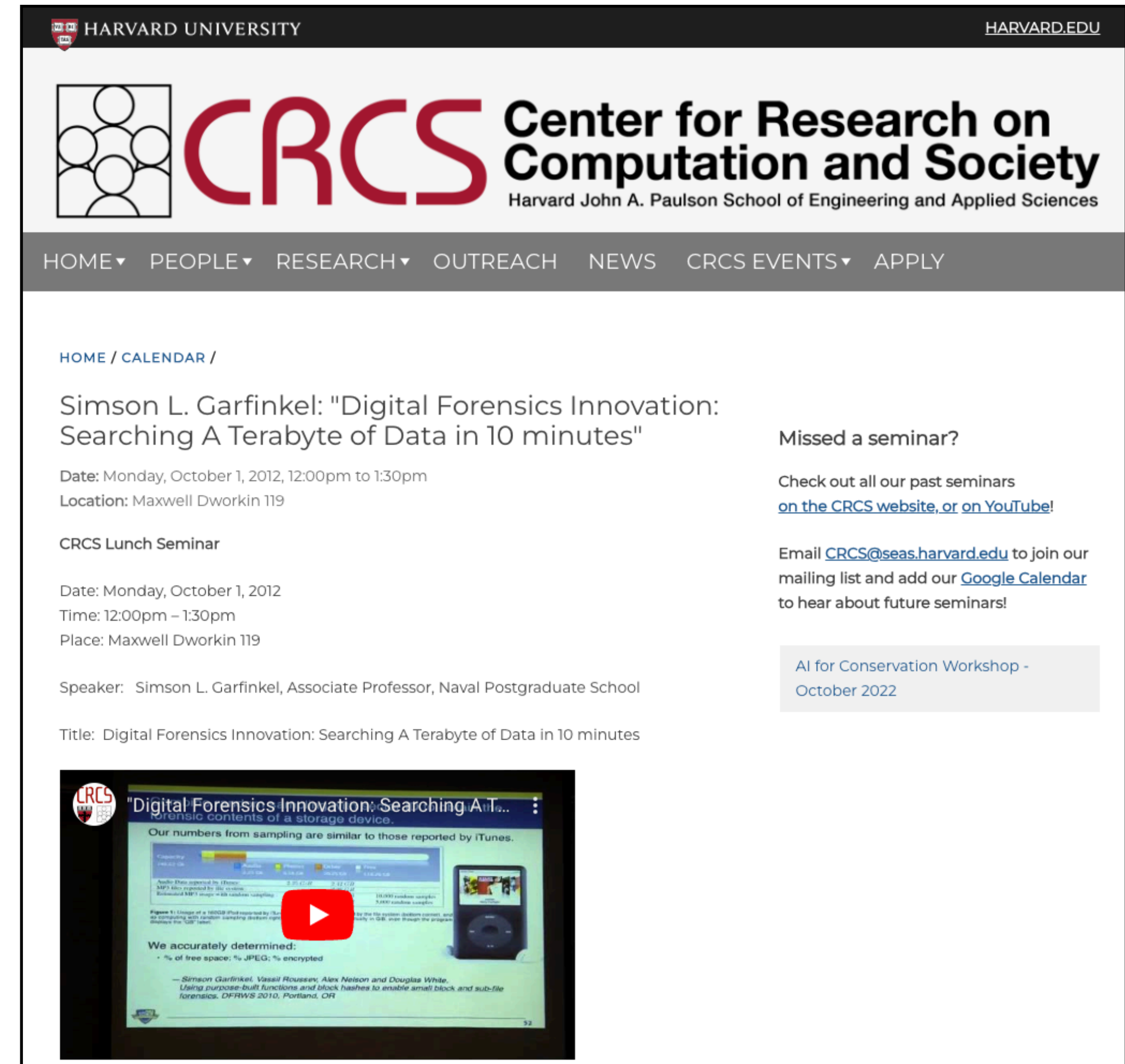## 2012 — Statistical sampling breakthrough

- 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

# bulk_extractor is an important tool in digital humanities.

# bulk_extractor is an important tool in digital humanities.

#2
Innovating with open data

# Creating, maintaining and distributing open data became another aspect of teaching entrepreneurship.



**18**

# 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



**Open Data on AWS**
Share any volume of data with as many people as you want

Contact us

# 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.



DARPA | DEFENSE ADVANCED RESEARCH PROJECTS AGENCY | ABOUT US / OUR RESEARCH

> Defense Advanced Research Projects Agency  >  Our Research  >  Safe Documents
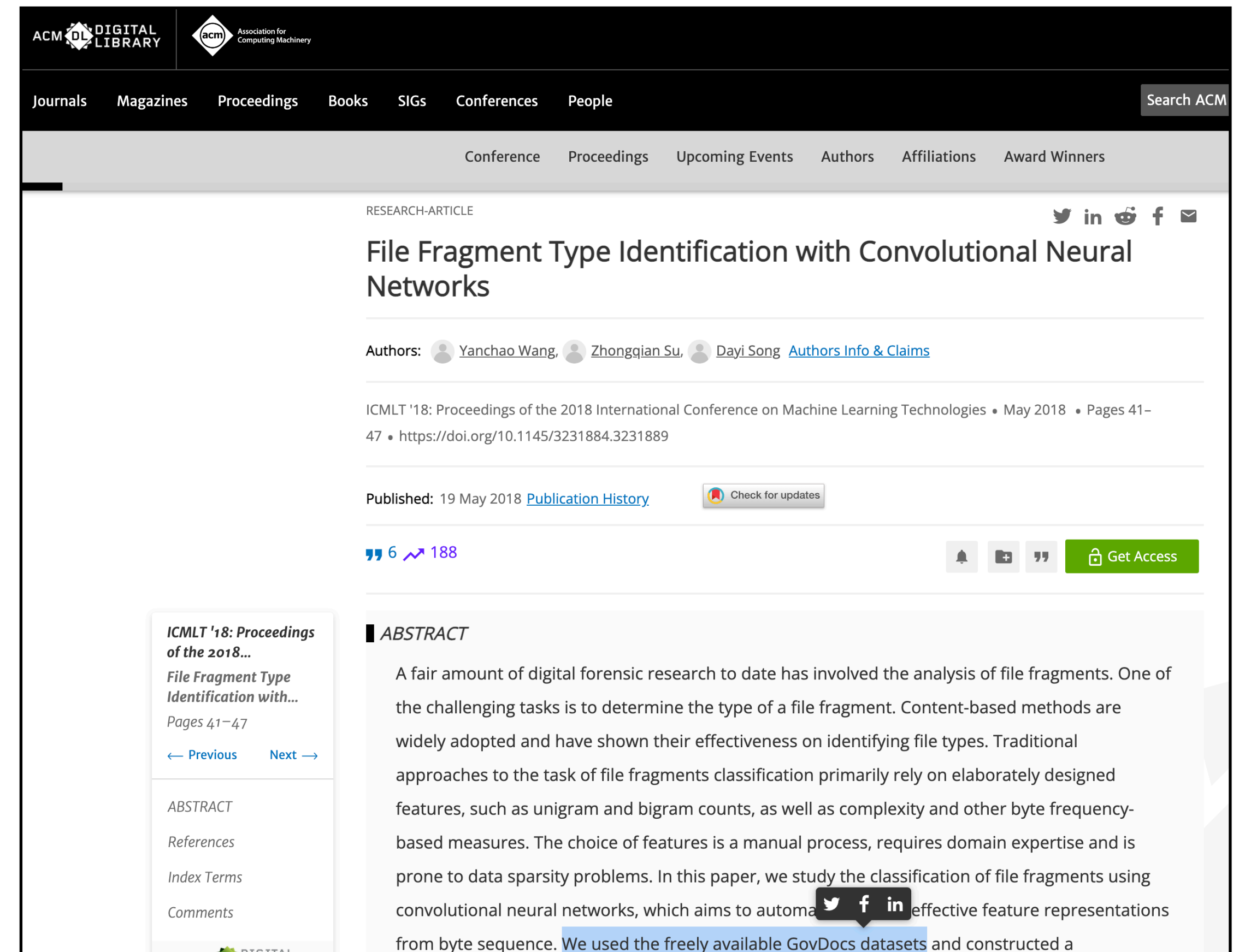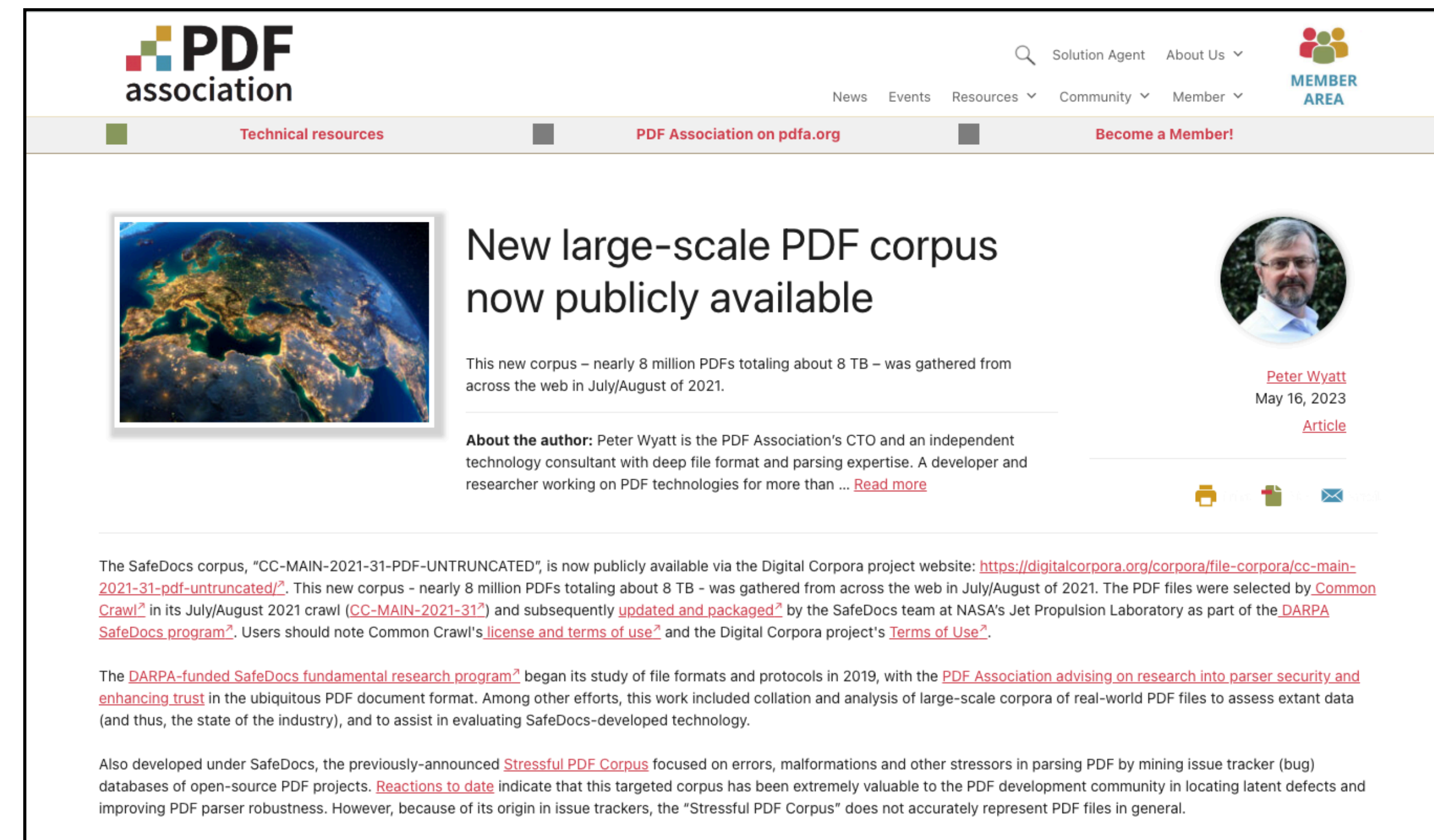
## Safe Documents (SafeDocs)

Dr. Sergey Bratus

2018-2023



PDF association

News    Events    Resources ∨    Community ∨    Member ∨    MEMBER AREA

Technical resources        PDF Association on pdfa.org        Become a Member!

### New large-scale PDF corpus now publicly available

This new corpus – nearly 8 million PDFs totaling about 8 TB – was gathered from across the web in July/August of 2021.

Peter Wyatt
May 16, 2023
Article

**About the author:** Peter Wyatt is the PDF Association's CTO and an independent technology consultant with deep file format and parsing expertise. A developer and researcher working on PDF technologies for more than … Read more

The SafeDocs corpus, "CC-MAIN-2021-31-PDF-UNTRUNCATED", is now publicly available via the Digital Corpora project website: https://digitalcorpora.org/corpora/file-corpora/cc-main-2021-31-pdf-untruncated/↗. This new corpus - nearly 8 million PDFs totaling about 8 TB - was gathered from across the web in July/August of 2021. The PDF files were selected by Common Crawl↗ in its July/August 2021 crawl (CC-MAIN-2021-31↗) and subsequently updated and packaged↗ by the SafeDocs team at NASA's Jet Propulsion Laboratory as part of the DARPA SafeDocs program↗. Users should note Common Crawl's license and terms of use↗ and the Digital Corpora project's Terms of Use↗.

The DARPA-funded SafeDocs fundamental research program↗ began its study of file formats and protocols in 2019, with the PDF Association advising on research into parser security and enhancing trust in the ubiquitous PDF document format. Among other efforts, this work included collation and analysis of large-scale corpora of real-world PDF files to assess extant data (and thus, the state of the industry), and to assist in evaluating SafeDocs-developed technology.

Also developed under SafeDocs, the previously-announced Stressful PDF Corpus focused on errors, malformations and other stressors in parsing PDF by mining issue tracker (bug) databases of open-source PDF projects. Reactions to date indicate that this targeted corpus has been extremely valuable to the PDF development community in locating latent defects and improving PDF parser robustness. However, because of its origin in issue trackers, the "Stressful PDF Corpus" does not accurately represent PDF files in general.

May 16, 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.
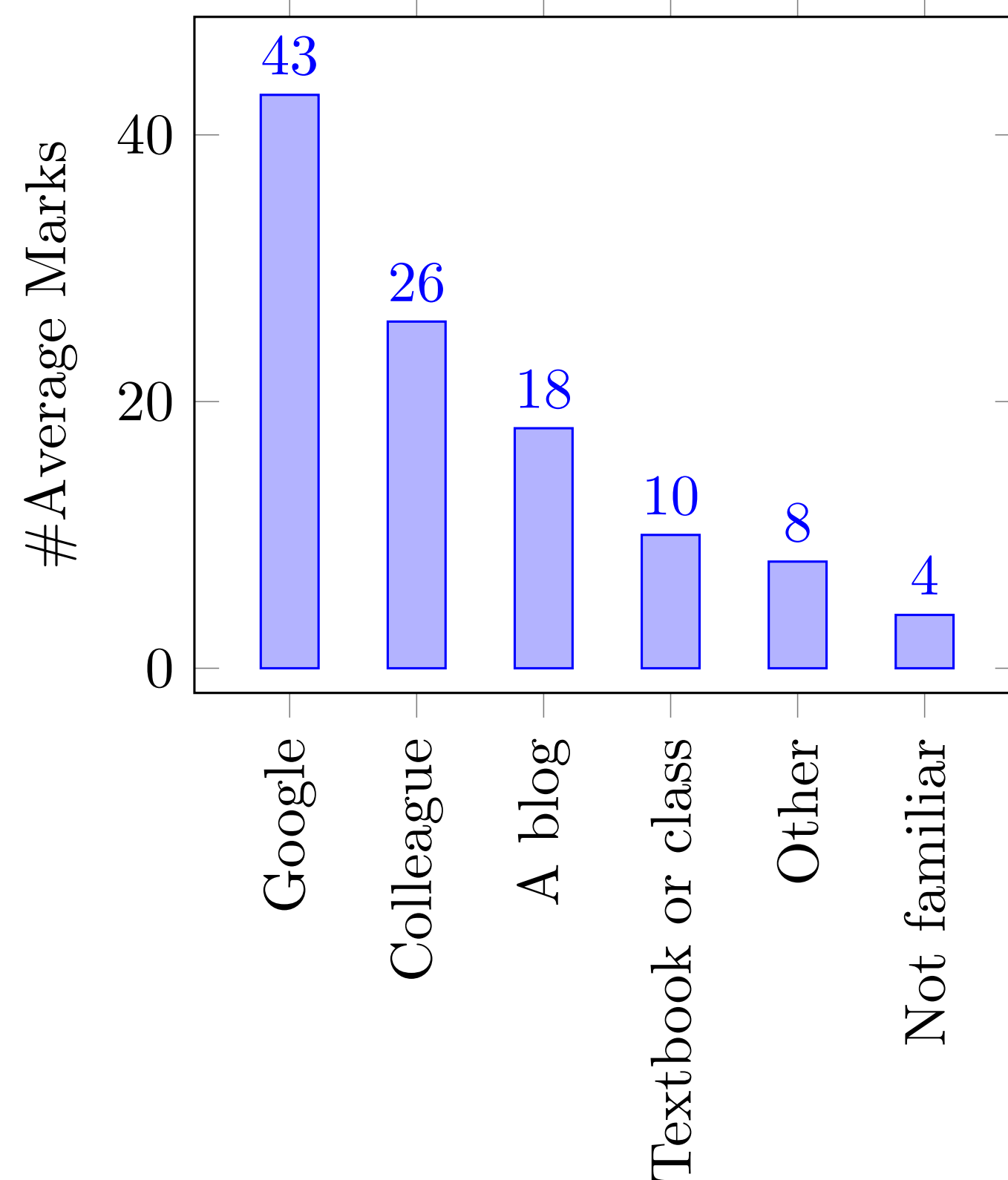
# Digital Corpora: Educational Impact

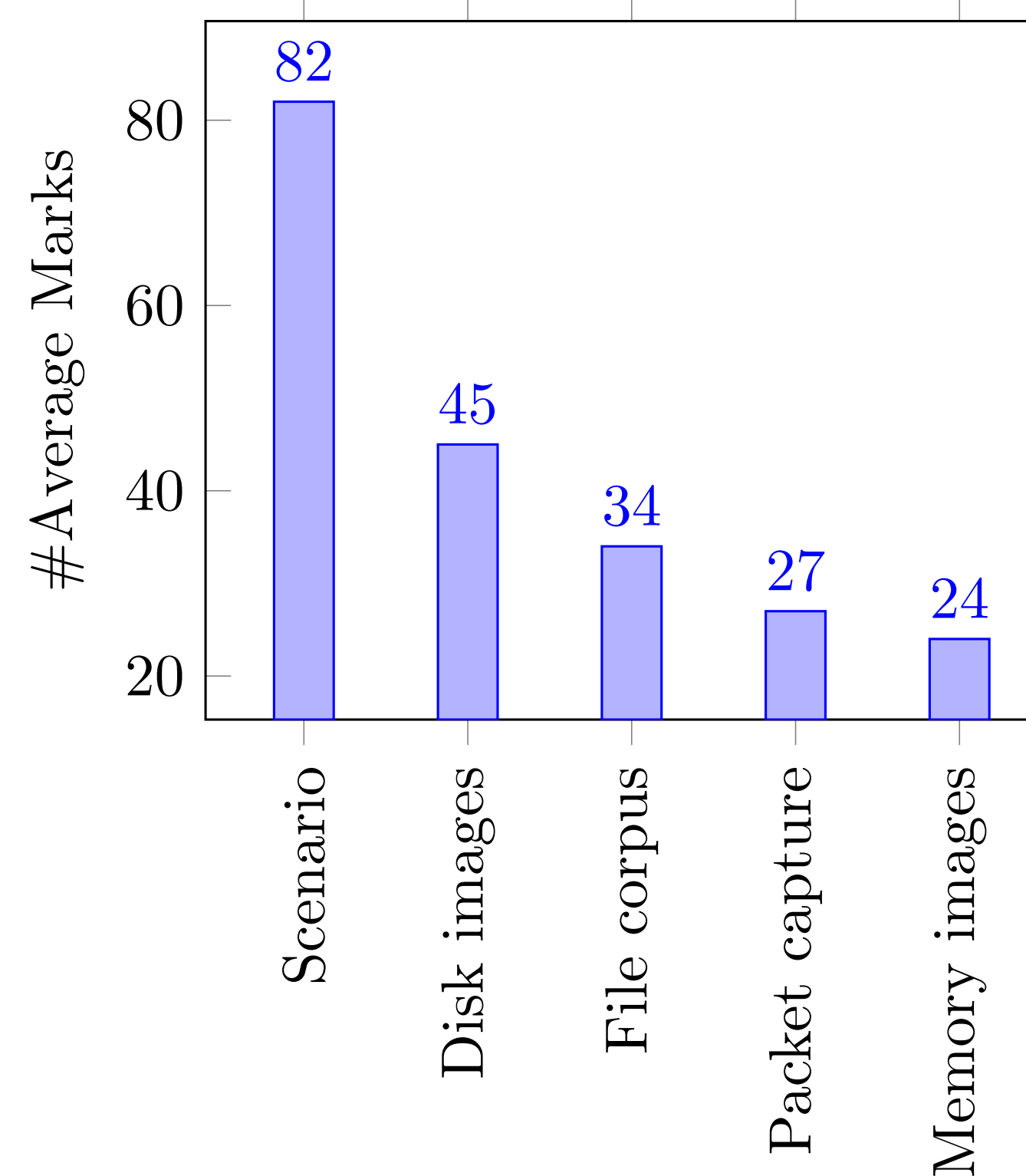Solutions to the scenarios are distributed as an encrypted PDF.

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?

26

# Datasets were used for education, tool testing, and a little research…

| | 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 |

# The digital corpora project also teaches innovation and entrepreneurship.

Developing scenarios that will be useful to others.

Planning and executing a complex project.

Quantifying the impact

- Woods, Kam, Christoper Lee, Simson Garfinkel, Extending Digital Repository Architectures to Support Disk Image Preservation and Access, JCDL 2011, June 13-17, 2011, Ottawa, Canada.
- Woods, K., Christopher Lee, Simson Garfinkel, David Dittrich, Adam Russel, Kris Kearton, Creating Realistic Corpora for Forensic and Security Education, 2011 ADFSL Conference on Digital Forensics, Security and Law

# My career is a testament to the power of <u>open source</u> and <u>open data</u>.

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

> Open data and free access to historical archives powers all of my journalism & historical work.

## 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 —
- Broadband2Wireless — 2000-2001

> Open source policy got my code out of MIT.

> Open source software was critical for the success of Vineyard.NET and Sandstorm.

## 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

> Open source policies within the US Government made it possible to rapidly transition software from my lab to DOD, FBI, Secret Service, and other government agencies.

## Career #4: Government Innovation
- US Census Bureau — 2017-2021 — Differential privacy
- US DHS — 2021-2022 — DHS Data Inventory

> Open source and open data policies made it easy to share code and data with the American people.

I use open source and open data to teach innovation and how to be entrepreneurial.

# Traps

Beware the support tail

Beware the sunk cost fallacy

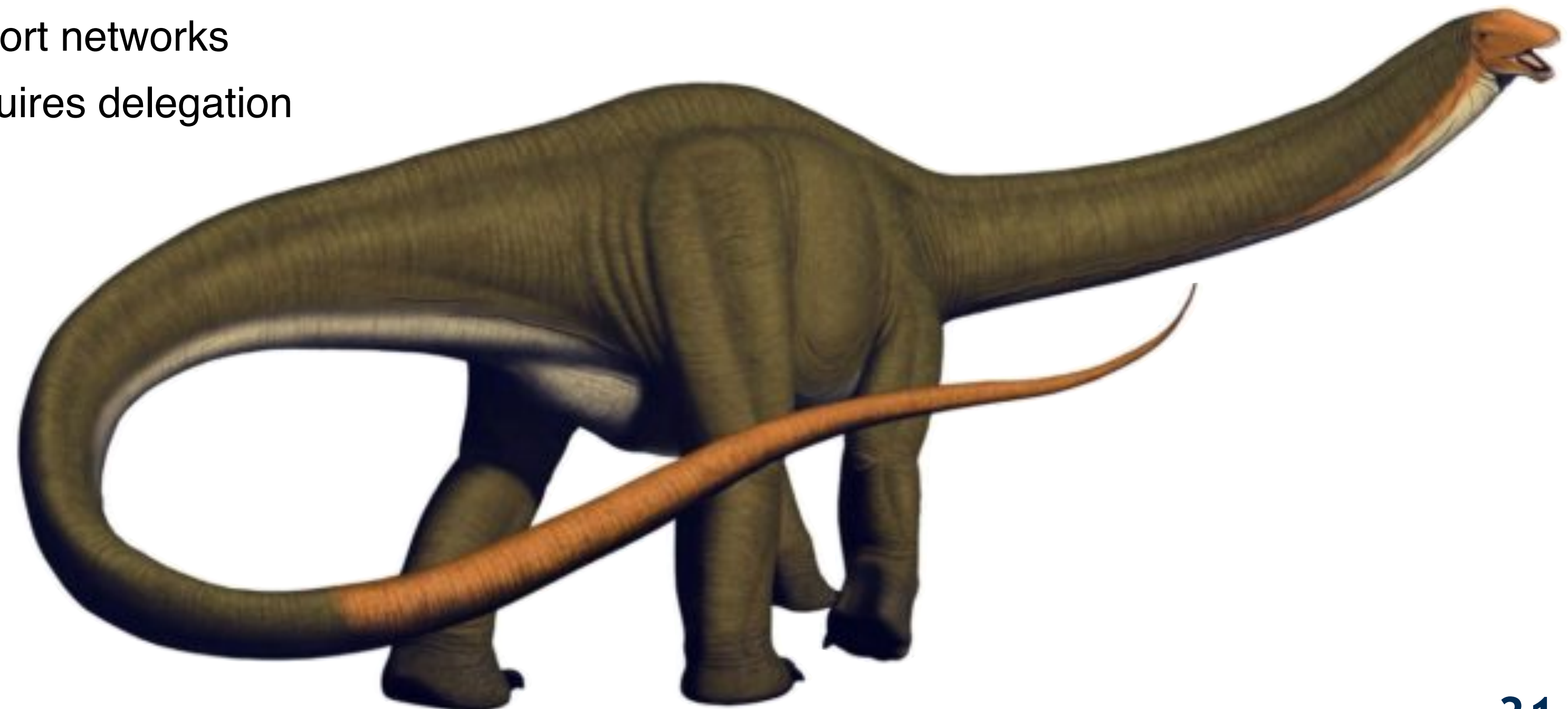Don't over-estimate your customers

Don't be greedy

# Trap — Beware the support tail

Early advice (bad): "You don't want to have users."

- Actually, you do.
- You don't want to have users who demand support

Avoid this trap: You don't want to support your users — you want your users' support.

- Get more people involved
- Clearly distinguish research & deployment
- Plant the seeds for self-sustaining mutual support networks
- Delegate, delegate, delegate — innovation requires delegation

# Trap — Beware the sunk cost fallacy

"The reason that God was able to create the world in seven days is that he didn't have to worry about the installed base."

## Things that get in the way of innovation:

- Existing code-base
- Deep expertise in tech stacks no longer in vogue*

## Example — Apple iPhoto & Aperture

- 2002 - 2015 — iPhoto was Apple's primary digital photo application
- 2005 - 2014 — Aperture was Apple's "professional" ($$$) photo editor
- 2015 — Apple killed both; replaced with Photos
- Photos was — less complex, integrated w/ iPhone and iCloud
- Cleaner Photos made possible more innovation

## Avoid this trap: Never be afraid to start over

—*a.k.a. "kill your darlings."*

Gordon Bell. 2003. Sink or Swim: Know When It's Time to Bail: A diagnostic to help you measure organizational dysfunction and take action. Queue 1, 9 (December/January 2003-2004), 60–67. https://doi.org/10.1145/966789.966806

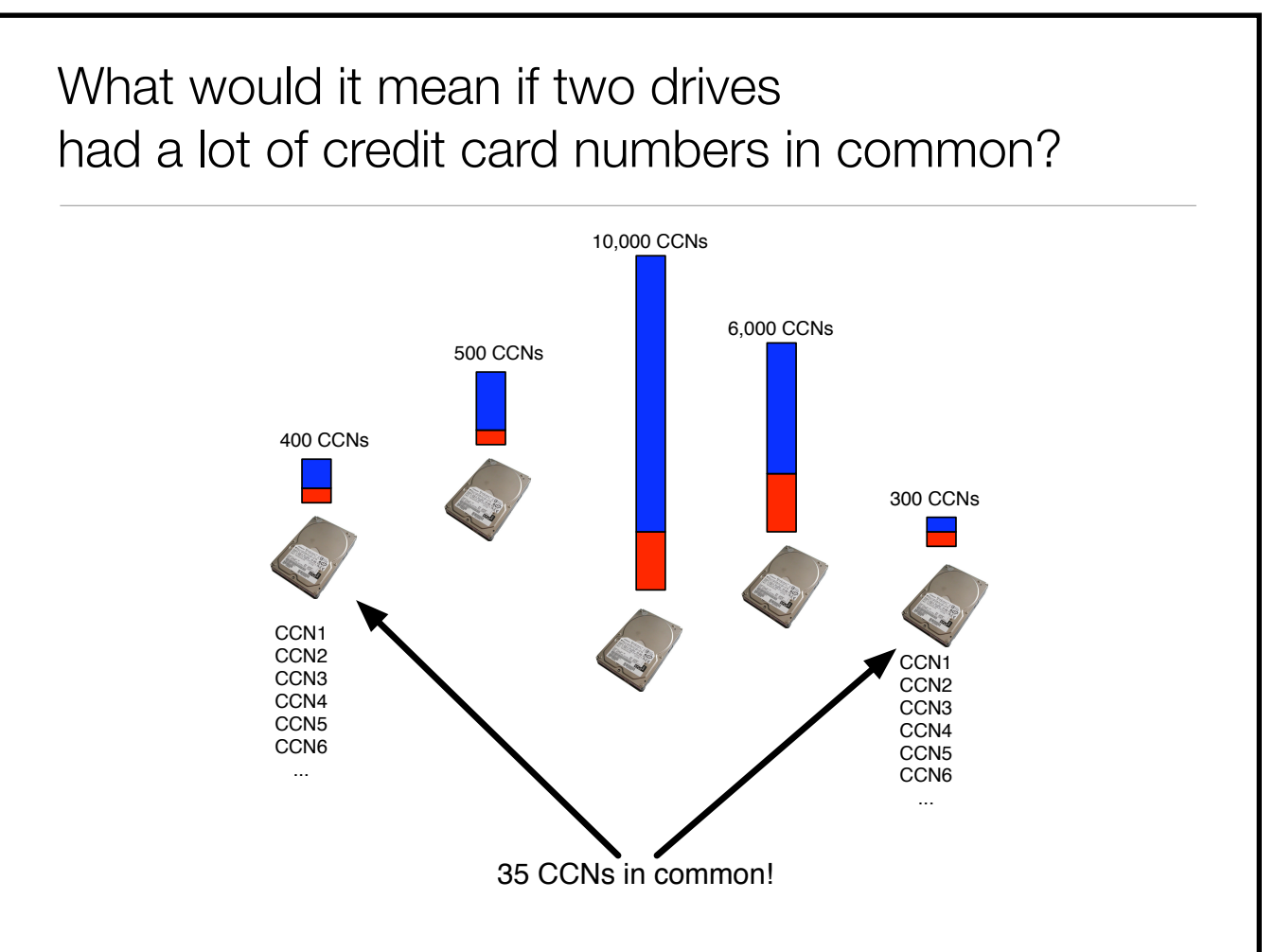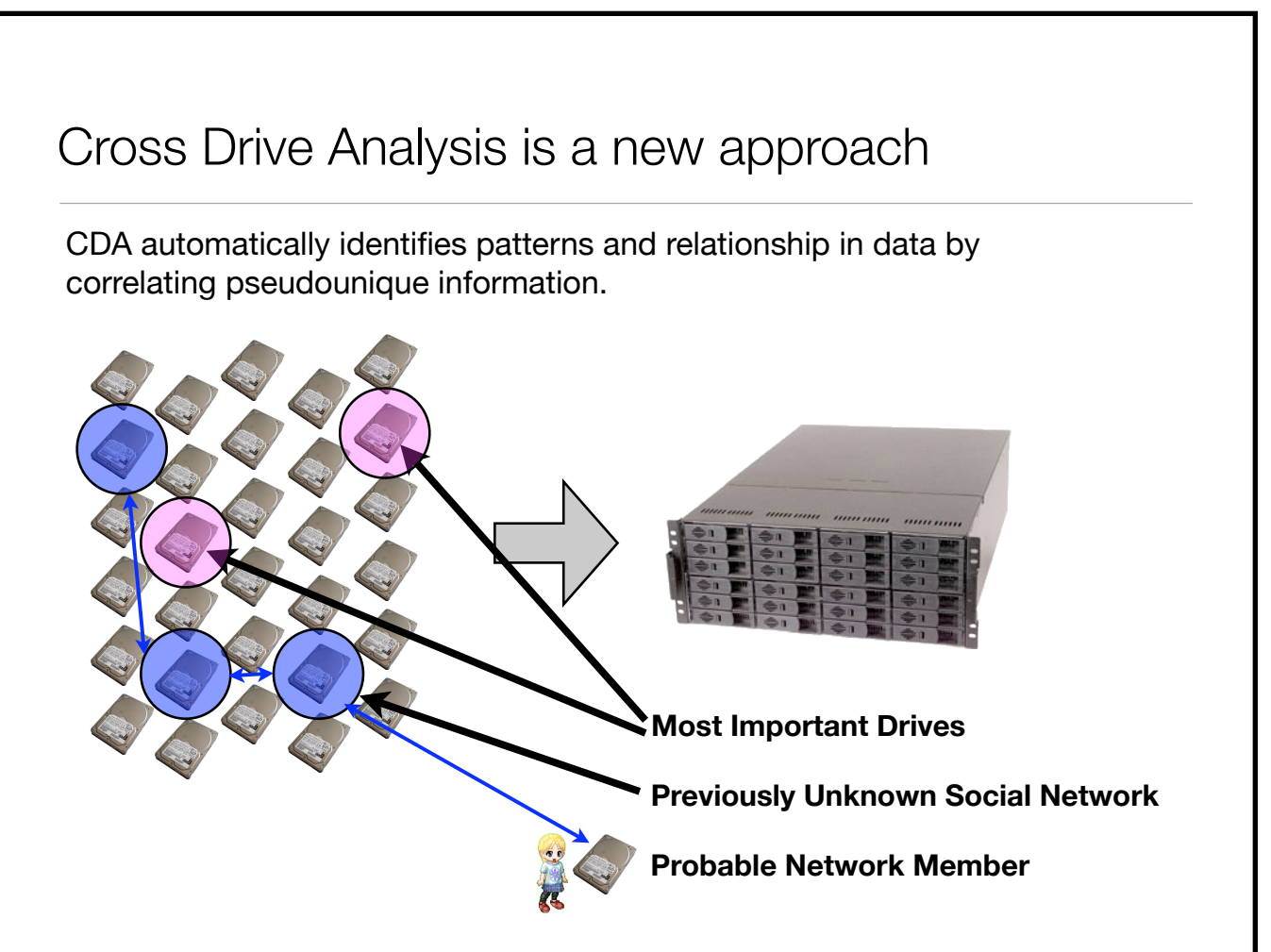# Trap — Don't over-estimate your customers (users)

## 2005 — I developed "cross-drive analysis" (CDA)

- Technique for finding hard drives & cell phones used by criminals in an investigation.
- Based on correlating **identifiers** (email, credit card #, etc.) between devices.
- Uses TF-IDF

## 2008 — I'm describing the power of CDA to a potential user.

- Me:
  - *We extract all of the phone numbers and email addresses from the drives.*
  - *We put them in a massive database.*
  - *We run this O(n2) algorithm*
  - *We use TF-IDF*

- Customer:
  - *You can extract phone numbers and email addresses from a drive automatically?'*
- Deploying automatic extraction at scale was transformative.

Avoid this trap: match your customer's technology readiness level.



Cross Drive Analysis is a new approach

CDA automatically identifies patterns and relationship in data by correlating pseudounique information.

Most Important Drives
Previously Unknown Social Network
Probable Network Member



What would it mean if two drives
had a lot of credit card numbers in common?

10,000 CCNs
500 CCNs
6,000 CCNs
400 CCNs
300 CCNs

CCN1
CCN2
CCN3
CCN4
CCN5
CCN6
...

CCN1
CCN2
CCN3
CCN4
CCN5
CCN6
...

35 CCNs in common!

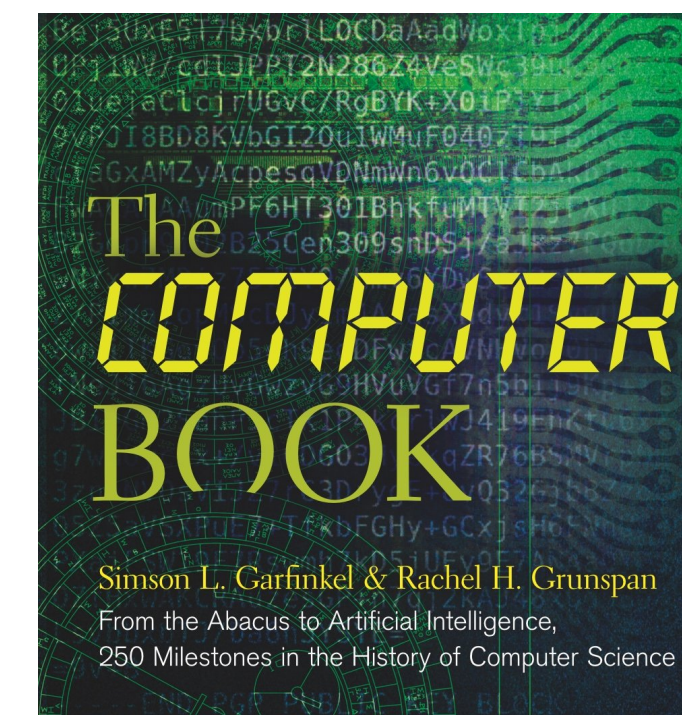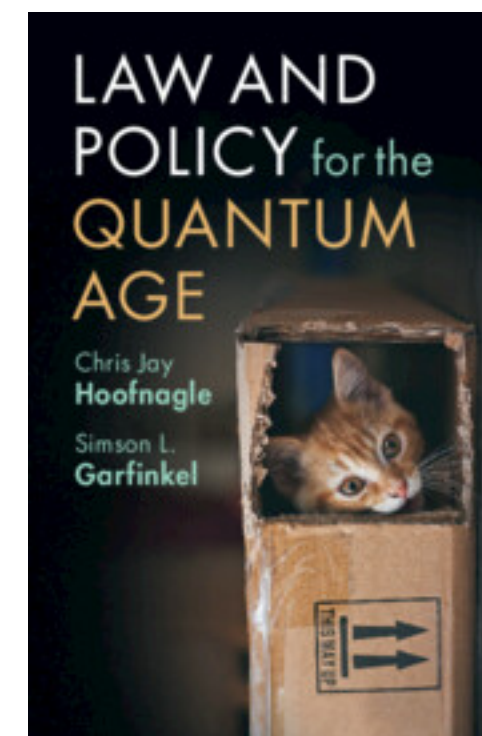Slides from 2005 presentation

# Trap — Don't be greedy

## Academic credit — be generous with authorship.

- Instead of listing students in acknowledgements, work to make them co-authors!
- I have multiple masters' students, undergrads, and 1 high school student as co-authors.
- Engaging with these students about the work is a great way to come up with new ideas.
  - *The best new ideas are conceptual, rather than incremental.*

## Academic work — seek out co-authors

- Of my 18 books, 15 have co-authors
- My "favorite" books all have co-authors

## Innovation requires multiple POVs

## Avoid this trap: Give credit, lower prices, exert less control.

# Techniques

Know your customer

Pursue extreme usability

Design and build for maintainability

Collaborate for scale

Teaching innovation by innovating

#1
Know your customer

# Technique — Know your customer

You've got a great idea — now what?

Product Market Fit (PMF) — What stands out in the lab is rarely what the market needs.

- "Market" — Customers with money / Users with a need / Government agencies with a need

Market research is critical.

Surprisingly, many would-be "innovators" people don't do this:

- Inventor / Entrepreneur has a great idea and goes out to find a market.
- "Build it and they will come" does not always work.
- Losing money to build market share only works if you can pivot and start making money.
  - *The success of Google, Facebook, Tesla, and others makes it hard to see all of the similar companies that failed.*

How to teach PMF? Eyeball interviews (something I learned in journalism school)

- Go out and talk to people. In person.
- This is terrifying for many students!

# The bulk_extractor market research study

**Bulk_extractor was 20 years in the making. It was not the tool that I planned to create!**

- In 1991 I developed SBook, a free-format address book for NeXT computers (pre-cursor to MacOS X).



- SBook used "Named Entity Recognition" to find addresses, phone numbers, email addresses *while you typed.*

# 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.

# 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
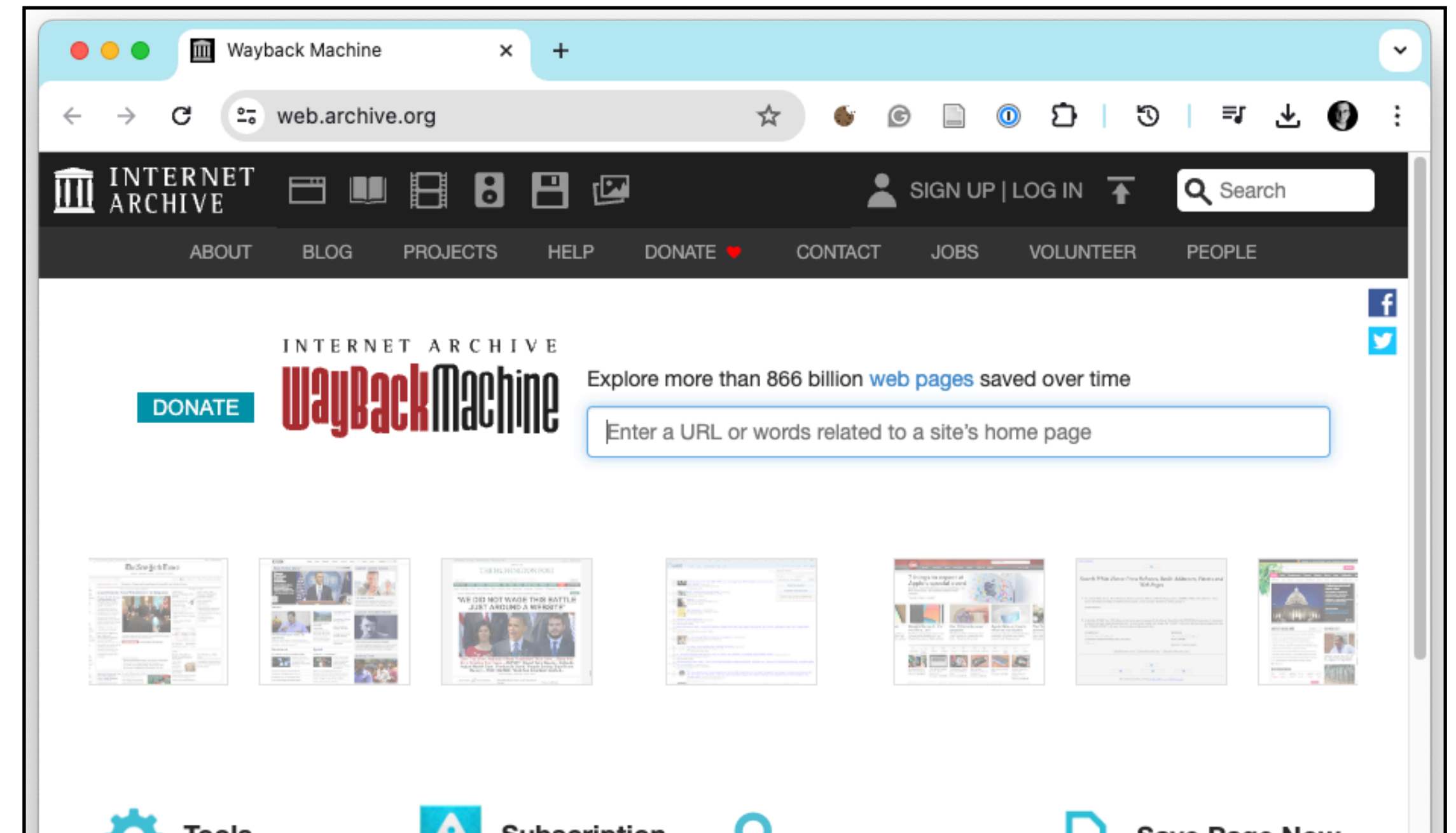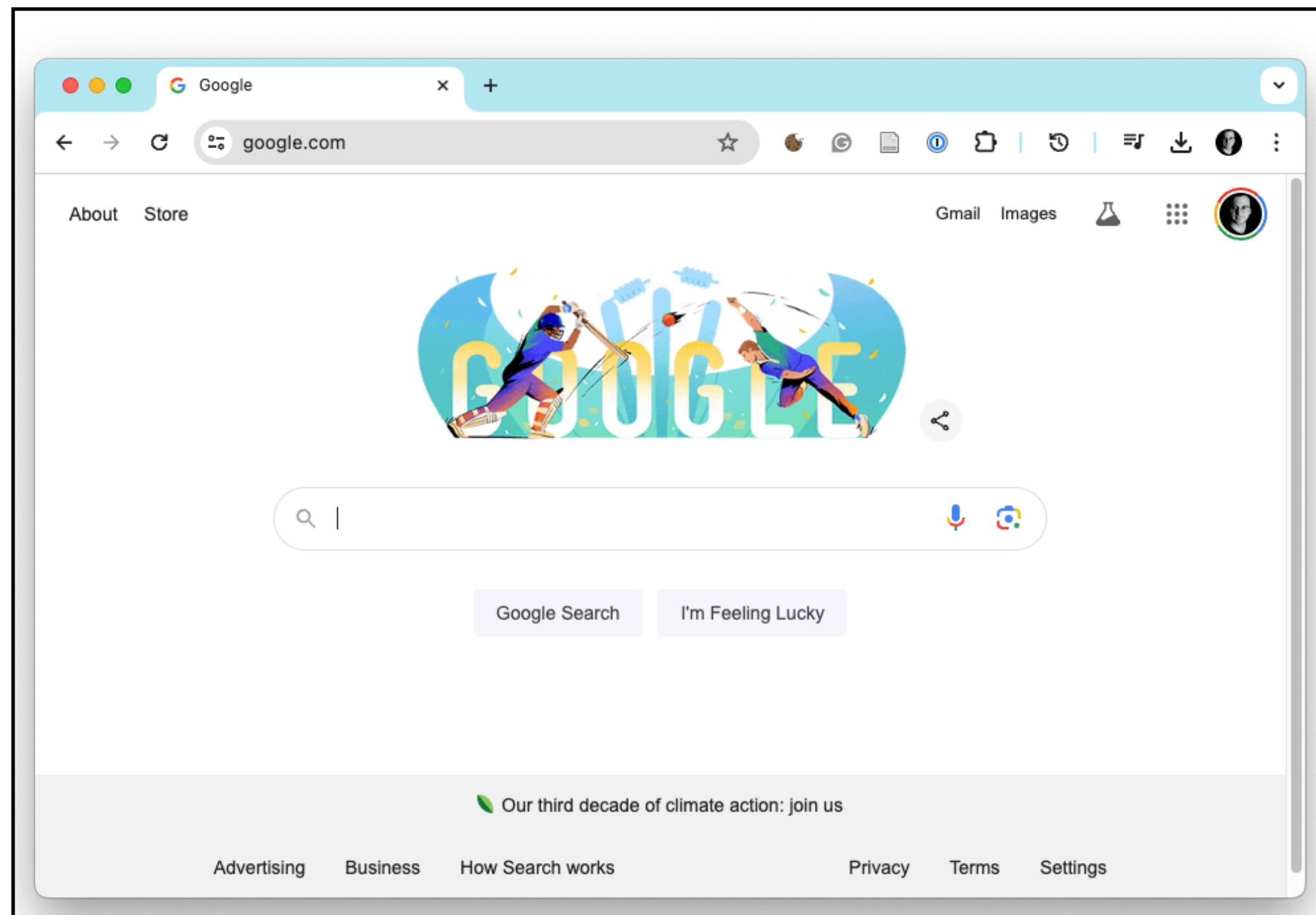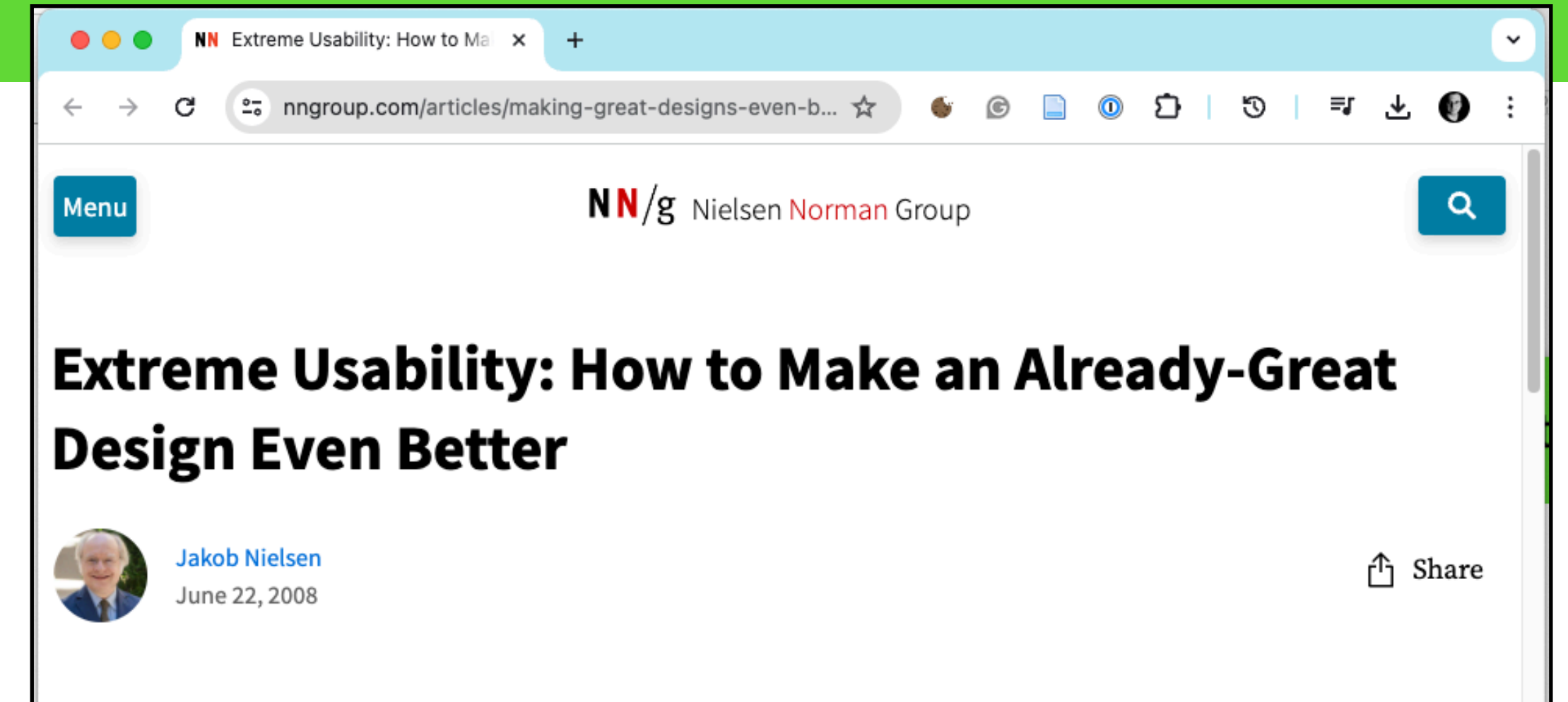- Run at maximum I/O speed of physical drive
- Never crash

# 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.

http://www.sanluisobispovacations.com/

A bulk_extractor
Success Story

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.

# Faster than conventional tools.
# Finds data that other tools miss.

Runs 2-10 times faster than EnCase or FTK *on the same hardware.*

bulk_extractor is multi-threaded; EnCase 6.x and FTK 3.x have little threading.

## Finds stuff others miss.

- "Optimistically" decompresses and re-analyzes all data.
- Finds data in browser caches (downloaded with zip/gzip), and in many file formats.

## Presents the data in an easy-to-understand report.

- Produces "histogram" of email addresses, credit card numbers, etc.
- Distinguishes primary user from incidental users.

## Faster than conventional tools.
## Finds data that other tools miss.

Runs 2-10 times faster than EnCase or FTK *on the same hardware.*

bulk_extractor is multi-threaded; EnCase 6.x and FTK 3.x have little threading.

### Finds stuff others miss.

- "Optimistically" decompresses and re-analyzes all data.
- Finds data in browser caches (downloaded with zip/gzip), and in many file formats.

Presents the data in an easy-to-understand report.

## This is the primary thing that mattered.

- Produces histogram of email addresses, credit card numbers, etc.
- Distinguishes primary user from incidental users.

# 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!**

#2
Extreme Usability

# Technique — Pursue extreme usability

Most bulk_extractor requirements from my study were *usability requirements*:

> ## 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

## Other usability requirements:

- Run "out of the box" with no training and no configuration.
- Run with no user interaction — A "get evidence" button

# Technique — Pursue extreme usability

## What is "extreme usability"

- Look for what's great, and replicate it.
- Look for what can go wrong, and ensure it never happens.
- Go beyond user experience and "consider enterprise usability"
- "Discover unmet needs"

# Usability goes everywhere — so does accessibility.

All materials should be *available* and *accessible.*

- Open source data helps!
- My book on quantum computing is Open Access (next book is also open access)

Design for accessibility:

- Materials should work with screen readers.
- Do not make assumptions about the ability of to students to see, hear, walk, etc.
- Understand how everyone on a team makes their contribution.
- *Model this for students — it's an important skill for entrepreneurs.*

LAW AND POLICY for the QUANTUM AGE

Chris Jay **Hoofnagle**

Simson L. **Garfinkel**

**Driving Innovation from the University to Government and the Enterprise: Tricks, Traps and Techniques**

Monday, June 3, 2024
Simson Garfinkel

These slides can be downloaded from https://simson.net/ref/2024

# #3
# Design and build for maintainability

# Technique — Design and Build for Maintainability

## Approaches for maintainability

- Correctness

- Clean, modular design

- Documentation

- Identify code that is <u>untested</u> and <u>dead</u>.

- Attention to non-technical issues — copyright, patents, privacy practices, etc.

## Innovation requires maintainability

- Technical debt — without maintainability, the cost of adding new functionality steadily rises.

- Clean documentation makes due-diligence easier and faster

# Design and build for maintainability
# Case Study — the 2020 Disclosure Avoidance System

#4
Collaborate for Scale

# Technique — Collaborate whenever you can

Students —

Faculty —

Businesses —

People in Government—

Collaboration is about the flow of *ideas,* not the flow of money or other resources.

*—And don't be greedy.*

#5
Teach innovation by innovating

# Teaching innovation by innovating

## "Advanced Computing for the Digital Humanities"

- Deploying private set intersection (PSI) to find connections between closed collections.

- Radically improving OCR for older texts.

- Large language models for ancient texts.

  —*The most innovative and transformative classes are research classes!*

  —*Digital humanities is an intelligence problem!*

## "Moving ideas from the lab to the marketplace" — case studies and student projects.

- Projects with which I've been personally involved — digital forensics, and differential privacy and the 2020 Census.

- Companies started and run by people in my network.

- Case studies that are both historically important and relevant to SEAS (e.g. MITRE, Digital Equipment Corp., etc)

## "Open Source Intelligence"

  —*aka "advanced web-scraping"*

  —*aka "data fusion with online information"*

- How do to it. Systematically. Bringing order to the world of online information

- Natural language processing, AI, scale, cloud processing, authentication, user interface — this course has it all!

# Questions (for you)

Q1 — The installed base.

Q2 — The cost of innovation.

# Question: When is it okay for an innovator to break their installed base?

```
From: Stuart Feldman <...@google.com>
Date: Mon, 20 Apr 2015 at 15:51
Subject: Re: make versus tabs
To: Michael Stillwell <...@google.com>

Story is only partly true.

I used tabs because I was trying to use Lex (still in first
version) and had trouble with some other patterns.

(Make was written over a weekend, rewritten the next weekend ...)

So I gave up on being smart and just used a fixed pattern (^\t)
to indicate rules.

Within a few weeks of writing Make, I already had a dozen friends
who were using it.

So even though I knew that "tab in column 1" was a bad idea, I
didn't want to disrupt my user base.

So instead I wrought havoc on tens of millions.

I have used that example in software engineering lectures.

Side note: I was awarded the ACM Software Systems Award for Make
a decade ago.  In my one minute talk on stage, I began "I would
like to apologize". The audience then split in two — half started
laughing, the other half looked at the laughers.

A perfect bipartite graph of programmers and non-programmers.
```

https://beebo.org/haycorn/2015-04-20_tabs-and-makefiles.html

Stu Feldman, author of 'make'

# Question — Is it "innovation" when you spend 4 years updating a code base?



Modern code base (C++20) • extensive unit tests (from 0% to 50% code coverage) improved threading model

This creates an infrastructure for future innovation.