# OBSERVING ARCHIVES
## Web Archival Labour as Socio-technical Practice

### JESSICA OGDEN, SUSAN HALFORD, & LES CARR

UNIVERSITY OF Southampton | Web Science DTC
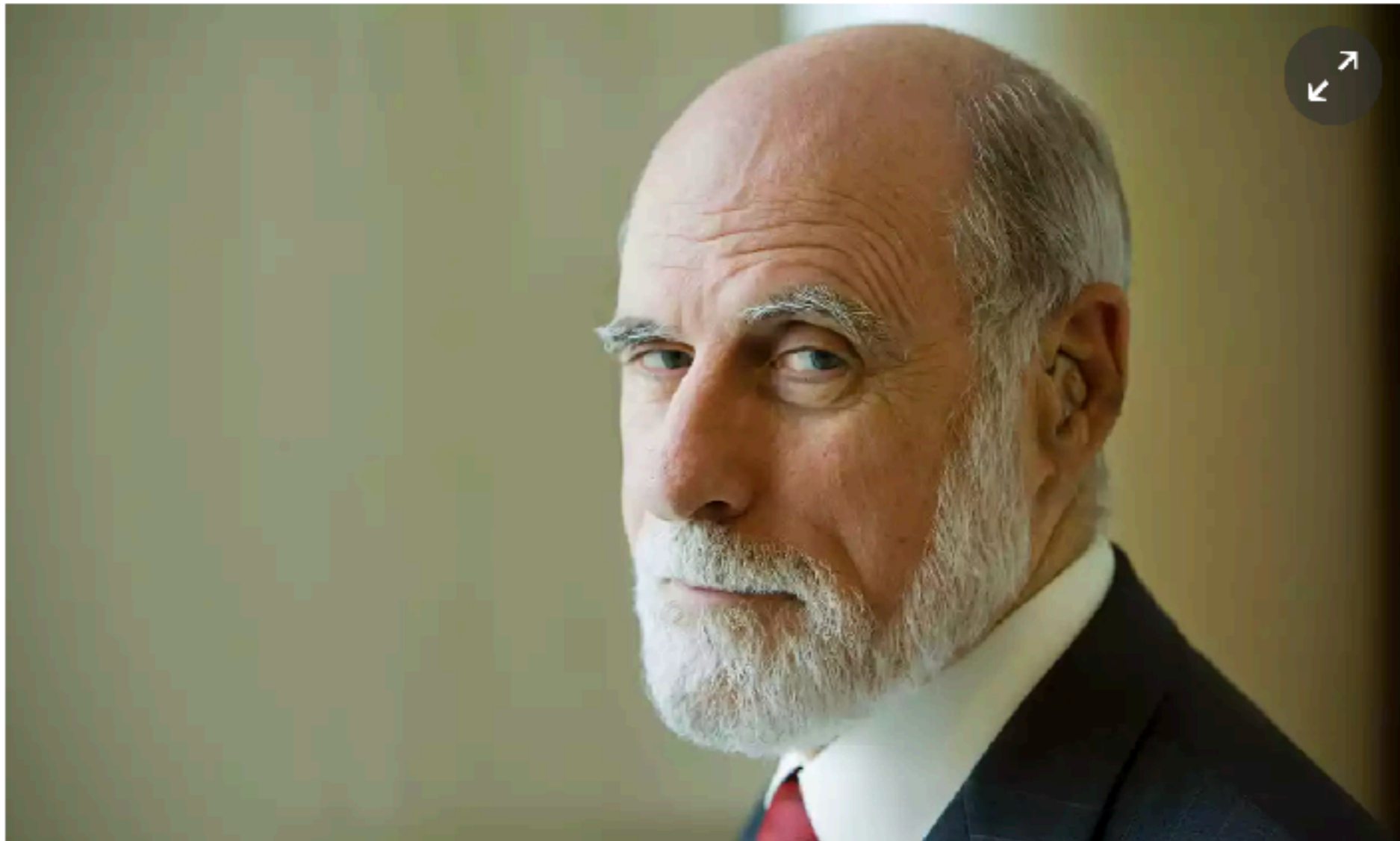
EPSRC

Research Councils UK
Digital Economy
Transforming Business and Society

# Google boss warns of 'forgotten century' with email and photos at risk

Digital material including key historical documents could be lost forever because programs to view them will become defunct, says Vint Cerf



Vint Cerf: 'We are nonchalantly throwing all of our data into what could become an information black hole.'
Photograph: Murdo Macleod

Guardian: https://www.theguardian.com/technology/2015/feb/13/google-boss-warns-forgotten-century-email-photos-vint-cerf
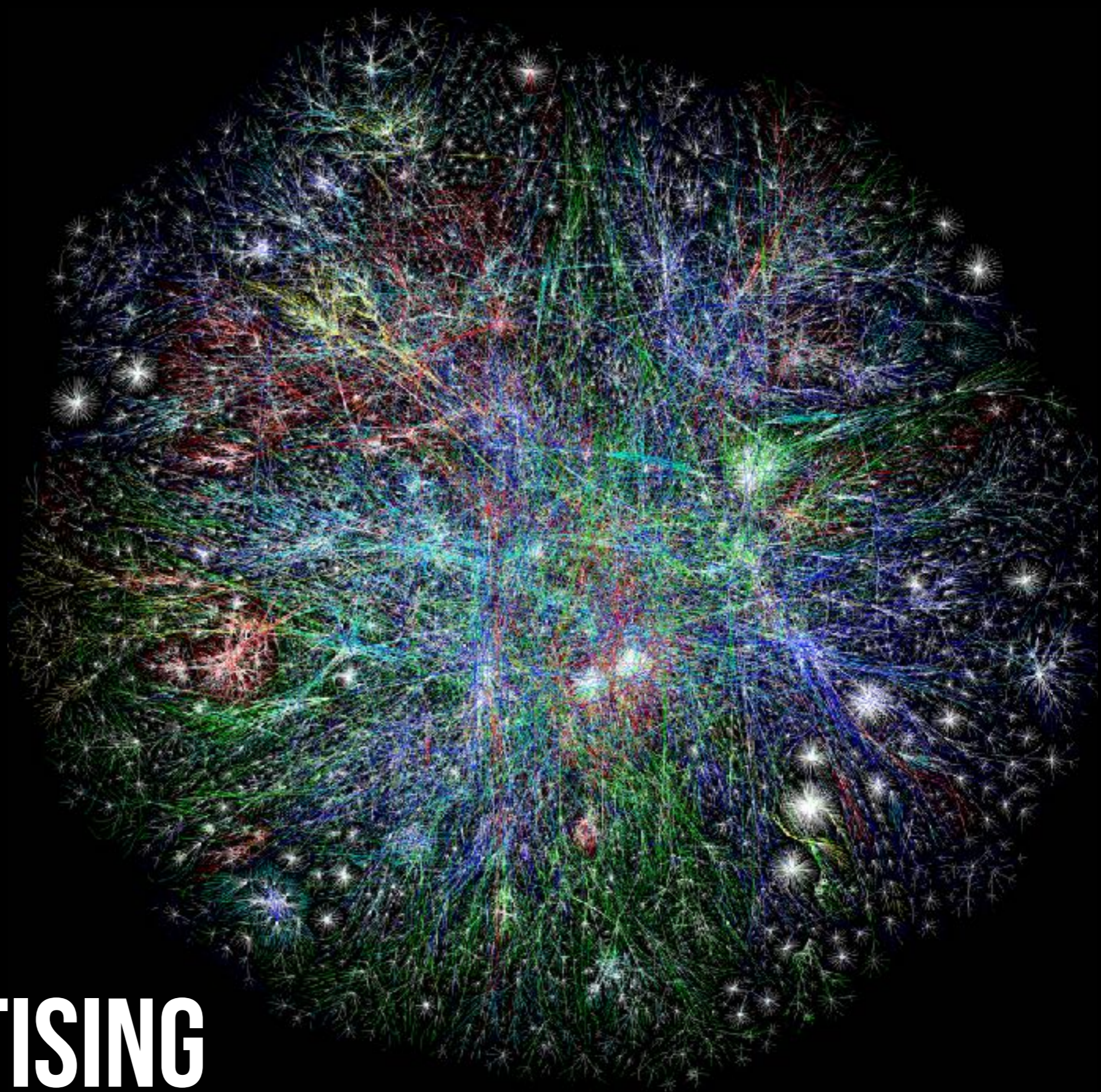
Library of Congress. Image: https://flic.kr/p/dUt1ir

WEB ARCHIVING

The Internet Archive

Data Rescue NYC

Image: Sam Hodgson, NYTimes https://nyti.ms/2mxKAGM

PROBLEMATISING

# WEB ARCHIVAL LABOUR

REPAIR

KNOWLEDGE WORK

WEB ARCHIVAL LABOUR

BREAKDOWN

TRANSLATION

Gregory J. Downey. 2014. Making Media Work: Time, Space, Identity, and Labor in the Analysis of Information and Communication Infrastructures. In Media Technologies: Essays on Communication, Materiality, and Society, Tarleton Gillespie, Pablo J. Boczkowski, and Kirsten A. Foot (Eds.). MIT Press, Cambridge, Massachusetts; London, England, 141–165.

# METHODOLOGY

- 16 un/semi-structured (ethnographic) interviews

- Observation records - what was done, made and used, what was said

- Documentary sources - wiki, policies, reports

- Two-tiered consent, 4 weeks

Explore more than 304 billion web pages saved over time

Saved **474 times** between April 12, 2004 and August 30, 2017.

## Summary of 4sonline.org

997 | 1998 | 1999 | 2000 | 2001 | 2002 | 2003 | 2004 | 2005 | 2006 | 2007 | 2008 | 2009 | 2010 | 2011 | 2012 | 2013 | 2014 | 2015 | 201

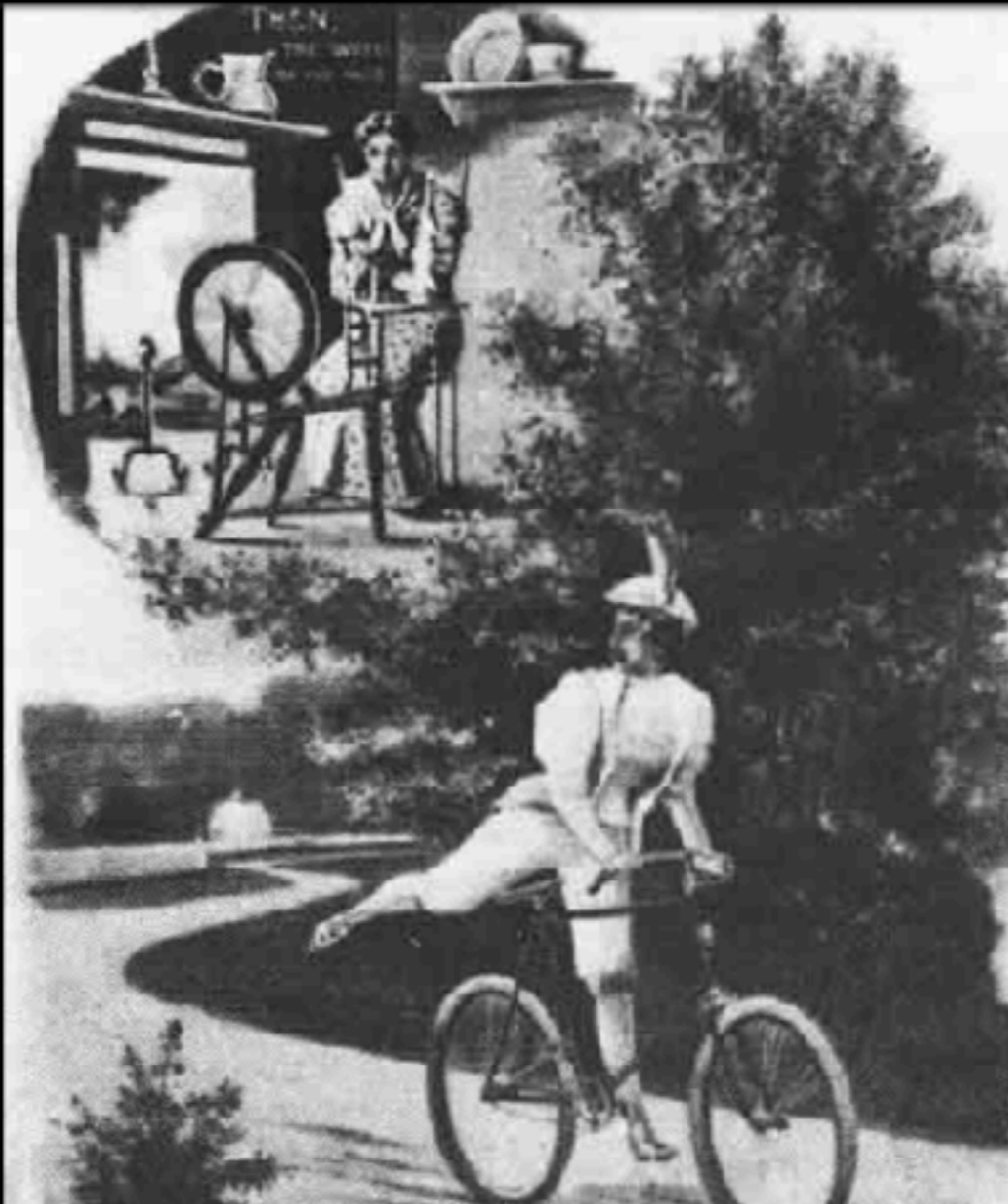| | JAN | | | | | | | | FEB | | | | | | | | MAR | | | | | | | APR |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | | | 1 | 2 | **3** | 4 | | | 1 | **2** | 3 | 4 | | | | | | 1 |
| 8 | 9 | 10 | 11 | 12 | 13 | 14 | | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| 15 | 16 | 17 | 18 | 19 | 20 | 21 | | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |

the profession

the society

scholarly
resources

annual meeting

for students

**4S News**
Announcing **4sonline.org**
4S now has a domain name all its own! Click to **bookmark this new location**.

Annual 4S & EASST
**Meeting in PARIS**, 25-28 August, 2004

**Travel grants** for graduate students

Graduate Student **Essay competition**

Science, Technology & Society
**Course Syllabi** available online

The Society for Social Studies of Science (4S) is a nonprofit, professional association. It was founded in 1975 and now has an international membership of about 1000.

The main purpose of 4S is to bring together those interested in understanding science, technology, and medicine, including the way they develop and interact with their social contexts.

2006

# MAPPING ROLES & ACTIVITIES

**External crawling**

Alexa Toolbar

**Directed crawling**

Wayback Machine

Heritrix

Archive-It

**Self-directed crawling**

Brozzler

Python Wayback

- User-directed

- Domain Crawls

- Thematic Crawls

- Event-based Crawls

- Wide Crawls

- Survey Crawls

- Platform/service specific

'Librarians'
Engineers
Web Archivists

Engineers
Web Archivists
'The Crowd'

*Schematic, not to scale and probably not even chronological!

# BREAKDOWN AND REPAIR

## Examples

- Breakdown as moments when contingencies of assemblage are revealed (Star & Ruhleder 1996)

- Repair and maintenance reveals 'ethics of care' afforded to technologies over time (Jackson 2014)

‣ crawler traps

‣ missing capture elements

‣ debugging capture vs replay

‣ patch crawling

STAR, SUSAN LEIGH, AND KAREN RUHLEDER. 1996. 'STEPS TOWARD AN ECOLOGY OF INFRASTRUCTURE: DESIGN AND ACCESS FOR LARGE INFORMATION SPACES'. INFORMATION SYSTEMS RESEARCH 7 (1).

STEPHEN J. JACKSON. 2014. RETHINKING REPAIR. IN MEDIA TECHNOLOGIES: ESSAYS ON COMMUNICATION, MATERIALITY, AND SOCIETY, TARLETON GILLESPIE, PABLO J. BOCZKOWSKI, AND KIRSTEN A. FOOT (EDS.). MIT PRESS, CAMBRIDGE, MASSACHUSETTS; LONDON, ENGLAND, 221–239.

# VALUES & 'ETHICS OF CARE'

- Popularity, novelty and levels of 'precariousness' all inform collection tasks

- Organisation/archivists value a Web with temporal depth - a **'Google for the historical Web'**

Reflect (expanding?) public insecurities around web-based truth claims

Highlight controversies over power dynamics, sustainability questions

# ECONOMIES OF WEB ARCHIVING

## Examples

- Emergent economy of external (repair) services around the Wayback Machine

  ‣ download and revive ageing websites

- Sharing economy and reliance on volunteer labour; knowledge production

  ‣ capturing links shared on social media; Wikipedia

- Politics of knowledge production beyond selection to consider repair

  ‣ broken link repair; browser extensions

# SUMMARISING

- web archiving as a form of maintenance

- ability to repair and maintain WAs intimately tied to innovation on the Web

- changing the landscape of the Web; creating stakeholders

- increasing reliance on parallel project of WAs - power dynamics



If you're interested in more:

OGDEN, JESSICA, SUSAN HALFORD, AND LESLIE CARR. 2017. 'OBSERVING WEB ARCHIVES: THE CASE FOR AN ETHNOGRAPHIC STUDY OF WEB ARCHIVING'. IN PROCEEDINGS OF WEBSCI'17, TROY, NY, USA., JUNE 25–28, 2017. ACM. DOI: 10.1145/3091478.3091506.

# ACKNOWLEDGEMENTS

EPSRC

Research Councils UK
Digital Economy
Transforming Business and Society

ARCHIVE
INTERNET