# OBSERVING WEB ARCHIVES

## The Case for an Ethnographic Study of Web Archiving

JESSICA OGDEN | SUSAN HALFORD | LES CARR

UNIVERSITY OF Southampton | Web Science DTC

Corresponding Author: jessica.ogden@soton.ac.uk

# OUTLINE

- Web archiving as a field of practice

- Problematising practice and the field; a theoretical framework

- Methodology and empirical work

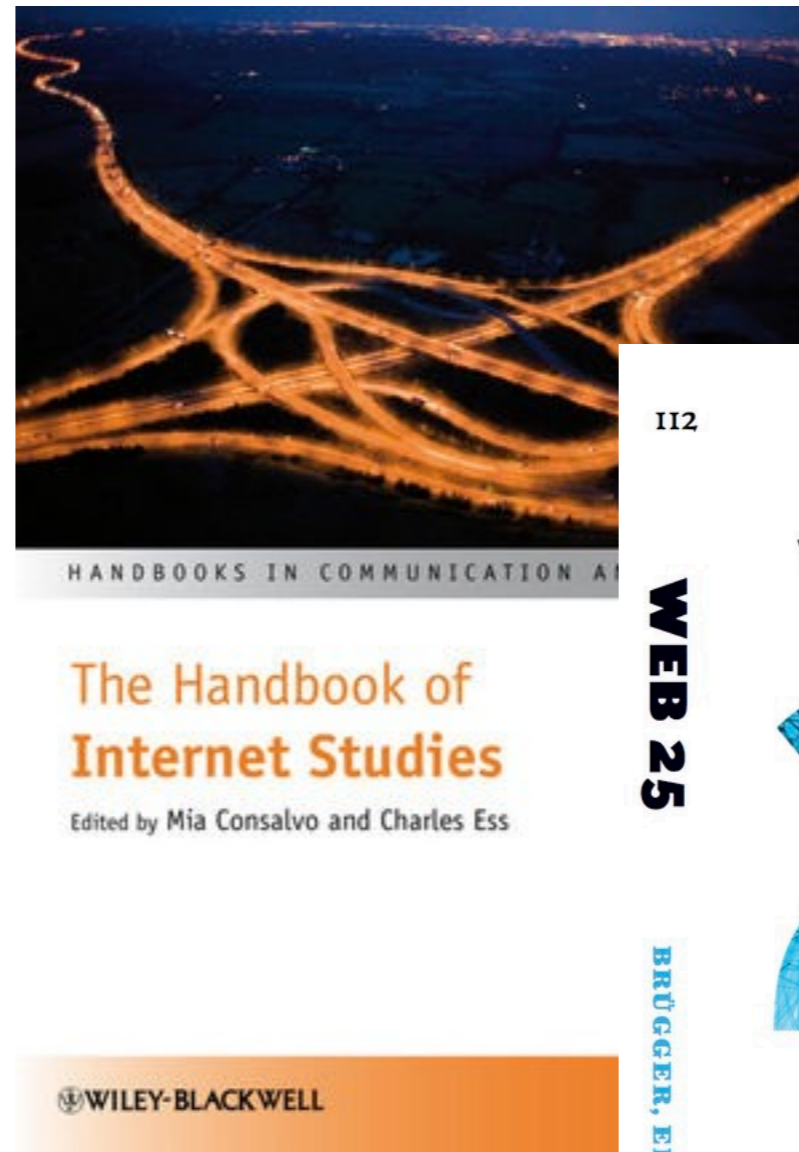- Preliminary results of a study at the Internet Archive

# WEB ARCHIVING

Web Archiving

Julien Masanès (Ed.)

archiving websites

a practical guide for information management professionals

adrian brown

The Handbook of Internet Studies

Edited by Mia Consalvo and Charles Ess

HANDBOOKS IN COMMUNICATION A

WILEY-BLACKWELL

112

WEB 25

WEB 25

Histories from the First 25 Years of the World Wide Web

NIELS BRÜGGER, EDITOR
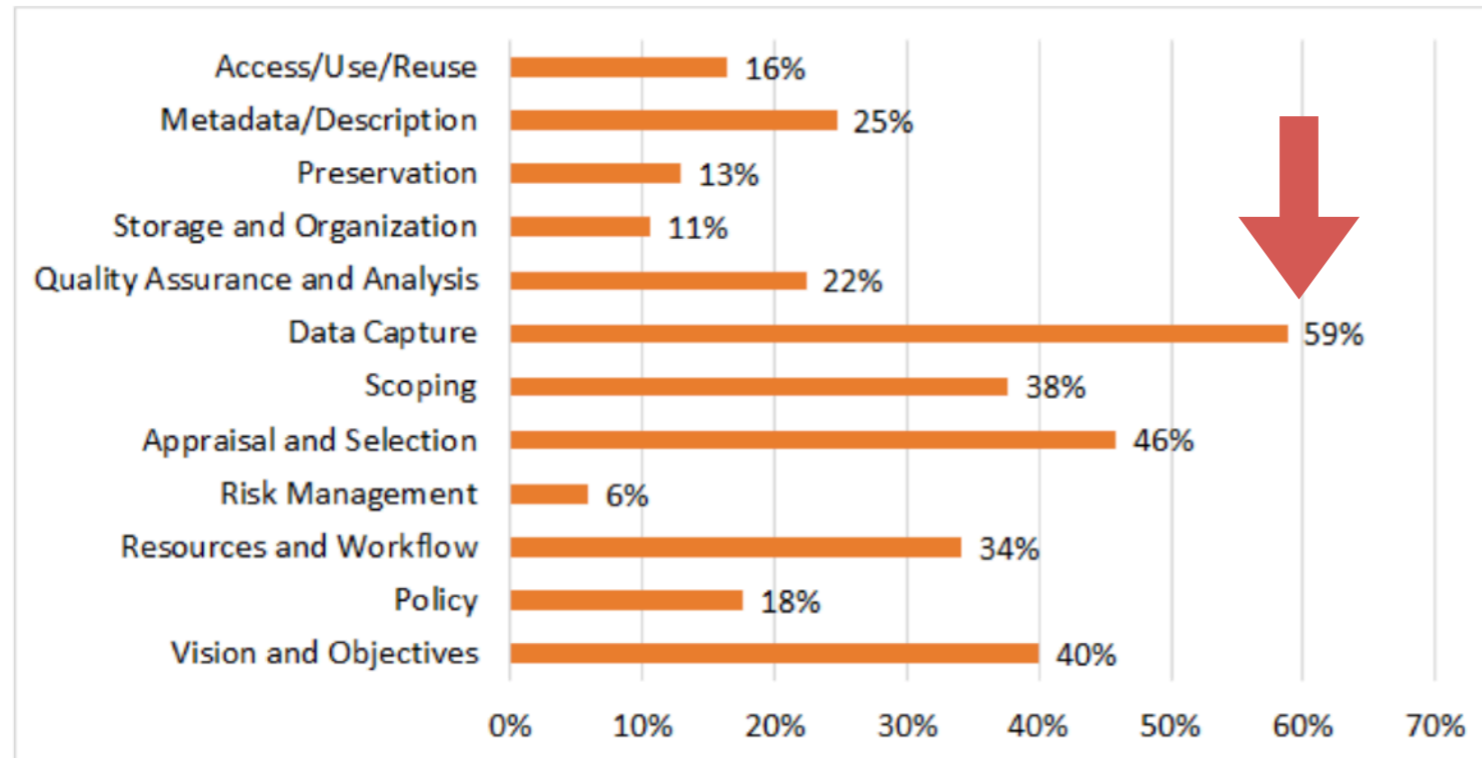
BRÜGGER, ED.

PETER LANG

# PRACTICE



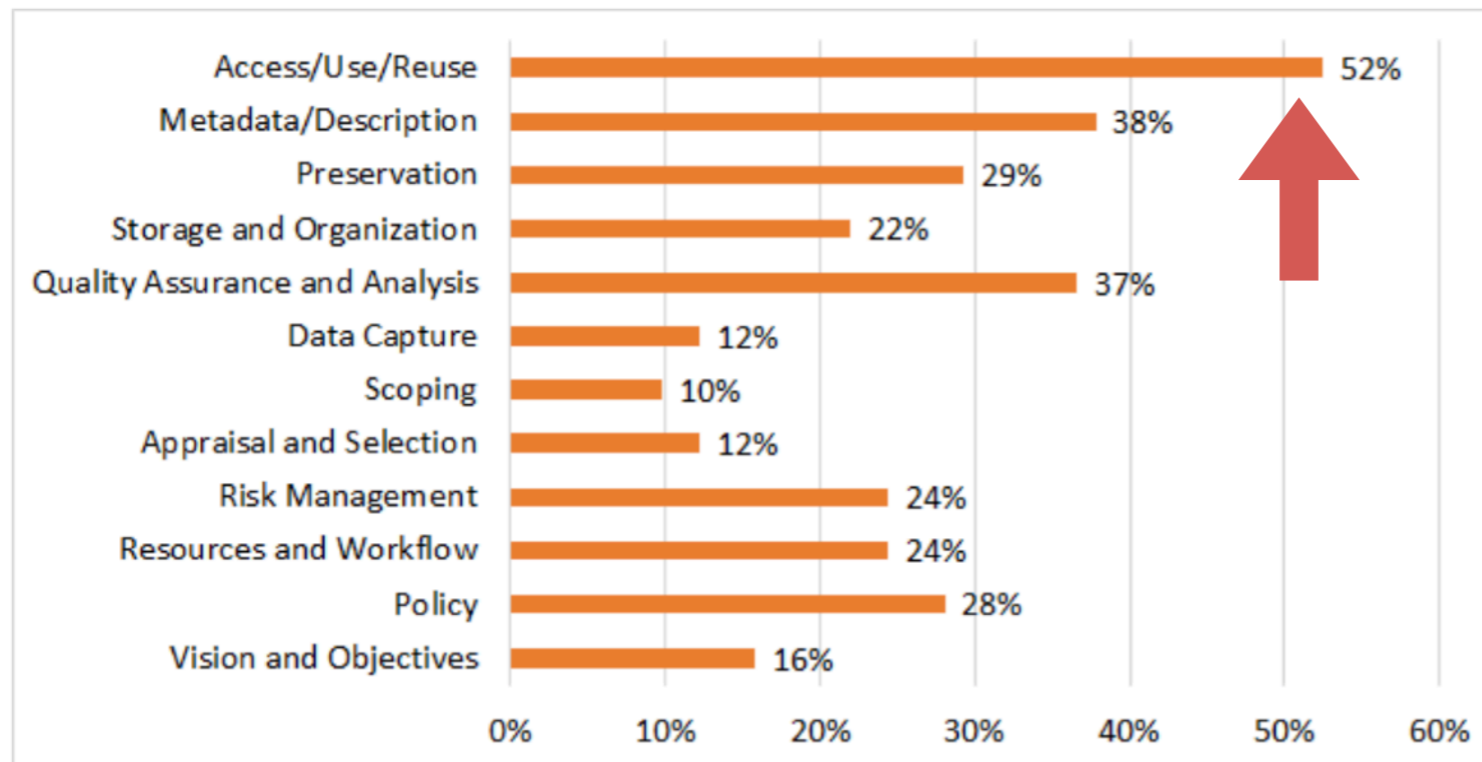FIGURE 3: PERCEPTIONS OF MOST PROGRESS IN LAST TWO YEARS



FIGURE 4: PERCEPTIONS OF LEAST PROGRESS IN LAST TWO YEARS

Bailey et al. 2017

Daniel Gomes et al. (2011) A Survey on Web Archiving Initiatives.

National Digital Stewardship Alliance (2012) Web Archiving Survey Report.

Bailey et al. (2014) Web Archiving in the United States: A 2013 Survey

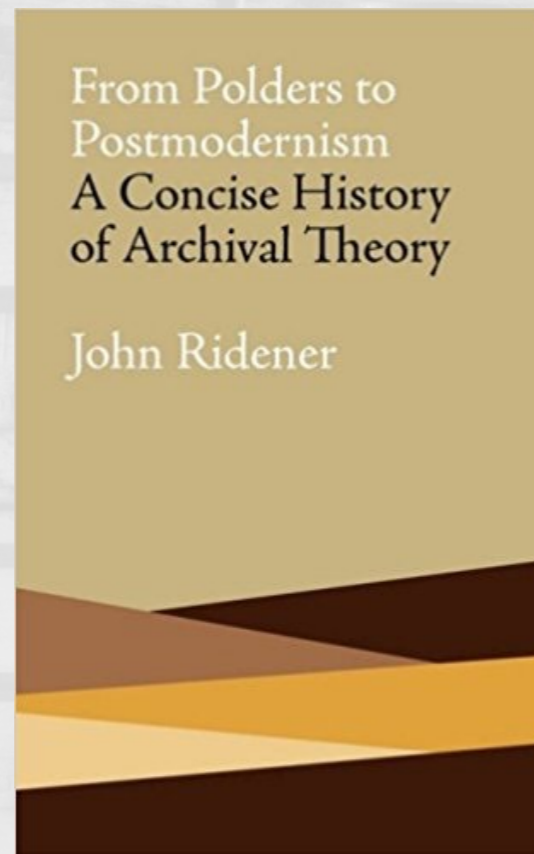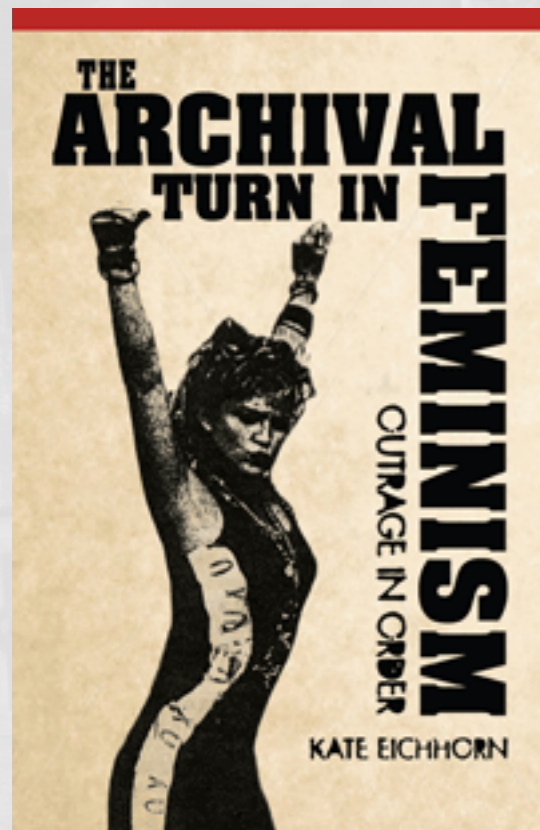Bailey et al. (2017) Web Archiving in the United States: A 2016 Survey

# PROBLEMATISING

- Existing overviews, best practice documents, surveys are partial

- Web archiving itself is inherently selective, contingent and ≠ a copy of the Web(s)

- Need to understand the role and agency of non/human actors in process

- Understanding the why, when + how is important for claims
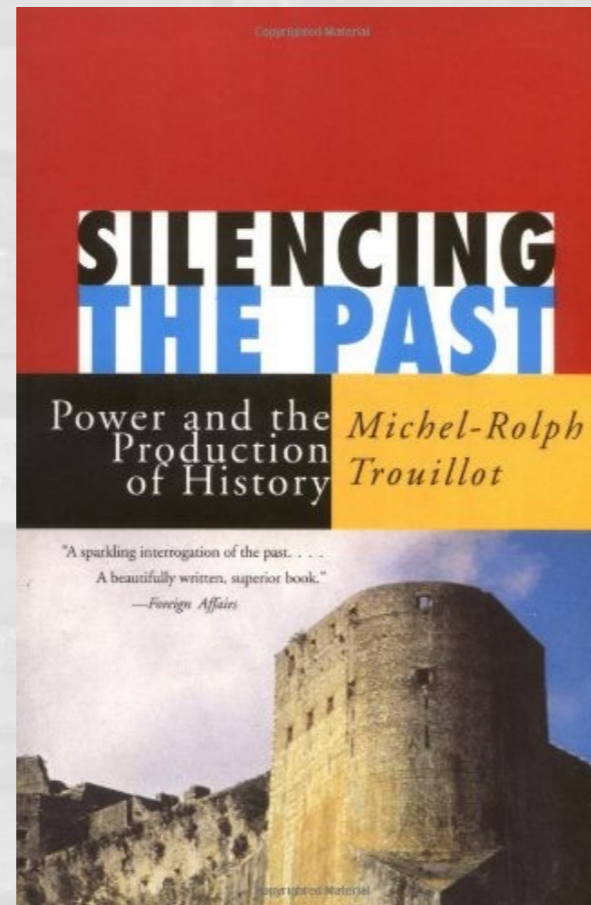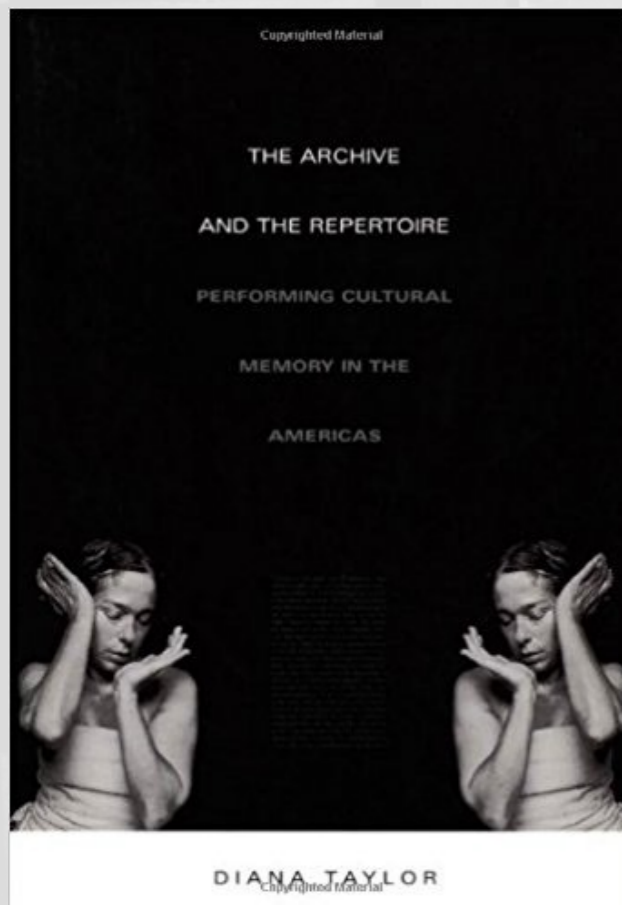
# FRAMING ENGAGEMENT

" *We will move from databases to knowledge bases. We will move, in the language of the post-modernists, to re-contextualize our activities: we will reorient ourselves from the content to the context, and from the end result to the original empowering intent, that is, from the artifact (the actual record) to the creating processes behind it, and thus to the actions, programmes, and functions behind those processes.*
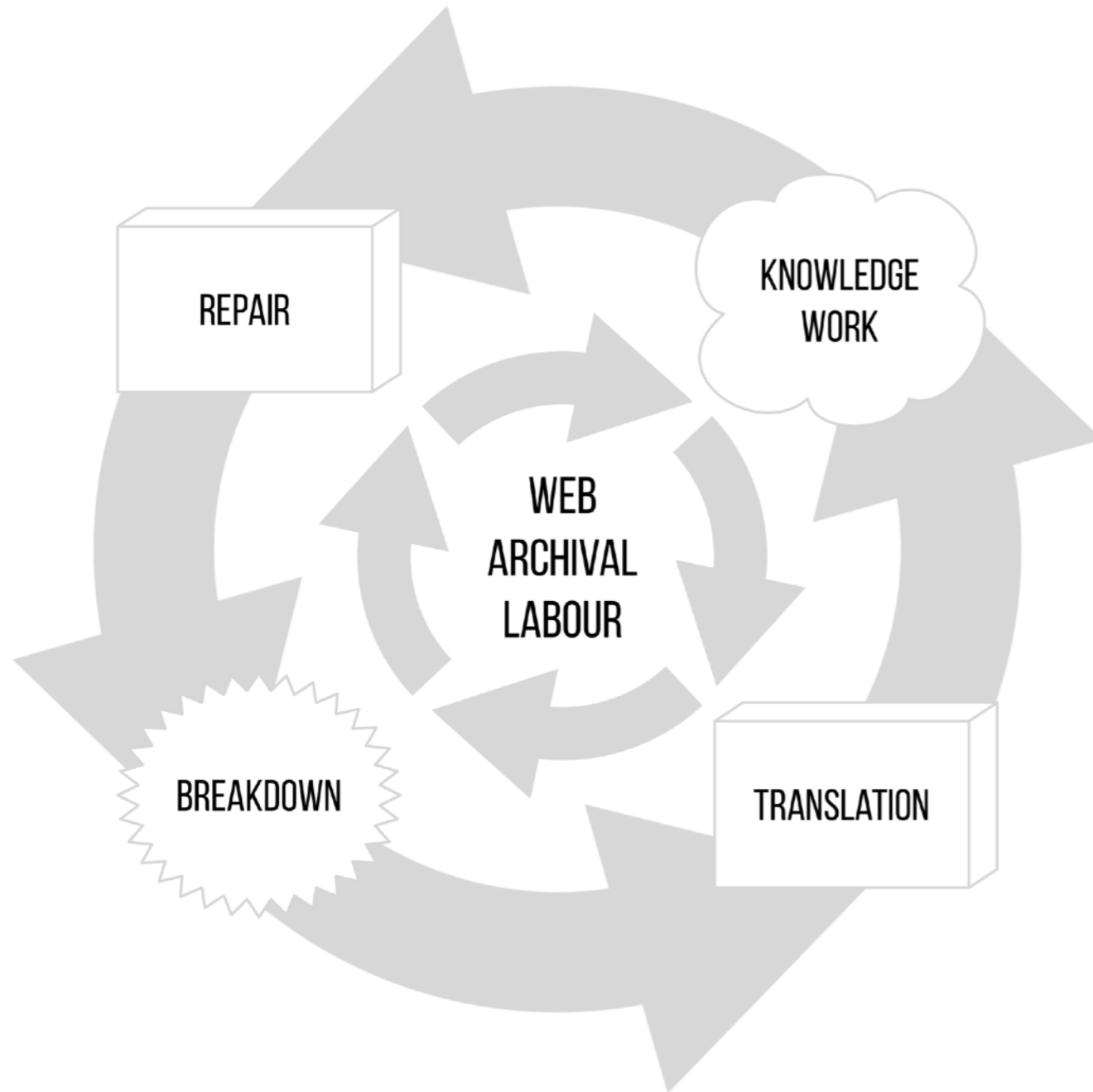
*-Terry Cook*

# ROLE OF THEORY

- From archival science to the archival turn

- The archive as generative, subjective not a 'view from nowhere'

- Understanding non/professional practice, role of politics

- Materiality of knowledge - practice as labour, tangible implications of intervention

# LABOUR



- Drawing on 'information labor' as outlined by Downey (2014)

- Labour = the work it takes to produce and transform web archives into information sources

- Process of making information 'useful'; putting it into 'circulation'

- Human and algorithmic; often obscure, hidden

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

- Ethnographic methods were chosen to document patterns of work

- 'Archival ethnography' (Gracy 2004)

- Understanding socio-cultural meaning, not a behavioural study

- Draws on extensive body of STS literature around record creation, scientific labs

# 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

- Thematic analysis (Spradley 1979) of 'things informants know'

# THE SITE

- Non-profit digital library founded in 1996 by Brewster Kahle

- Headquarters in San Francisco

- Began as mechanism to capture web pages as by-product of indexing

- Expanded remit to include books, audio, film/video, images, docs, video games and software

INTERNET ARCHIVE

# INTERNET ARCHIVE
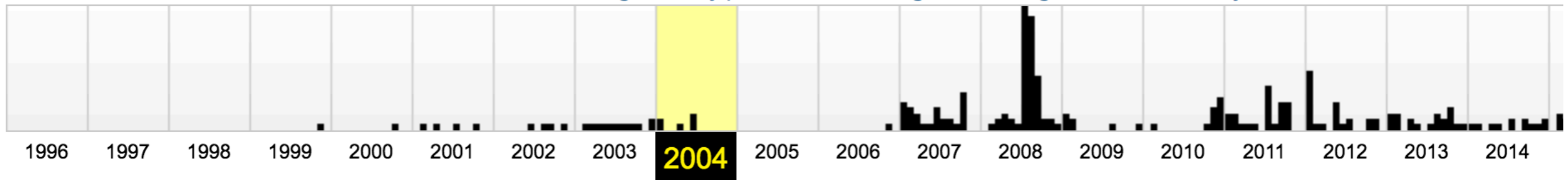## WayBackMachine

`http://webscience.org/`  ✕

Explore more than 298 billion web pages saved over time

Saved **319 times** between November 27, 1999 and June 16, 2017.

## Summary of webscience.org

**PLEASE DONATE TODAY.** Your generosity preserves knowledge for future generations. Thank you.

1996  1997  1998  1999  2000  2001  2002  2003  **2004**  2005  2006  2007  2008  2009  2010  2011  2012  2013  2014

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

# WSRI
web science research initiative

▣ Home   ▣ Publications   ▣ Events   ▣ Contact

# Web Science

The Web Science Research Initiative (WSRI) is a joint endeavour between the Computer Science and Artificial Intelligence Laboratory (CSAIL) at MIT and the School of Electronics and Computer Science (ECS) at the University of Southampton. The goal of WSRI is to facilitate and produce the fundamental scientific advances necessary to inform the future design and use of the World Wide Web.

The initiative will have four founding directors: Tim Berners-Lee, director of the World Wide Web Consortium, senior research scientist at MIT and professor at the University of Southampton; Wendy Hall, professor of computer science and head of the School of Electronics and Computer Science at the University of Southampton; Nigel Shadbolt, professor of artificial intelligence at the University of Southampton and director of the Advanced Knowledge Technologies Interdisciplinary Research Collaboration; and Daniel J. Weitzner, Technology and Society Domain leader of the World Wide Web Consortium and principal research scientist at MIT. Jim Hendler, Professor of computer science department at Rensselaer Polytechnic Institute, will serve as Associate Director.

"Since its inception, the World Wide Web has changed the ways scientists communicate, collaborate, and educate. There is, however, a growing realization among many researchers that a clear research agenda aimed at understanding the current, evolving, and potential Web is needed. If we want to model the Web; if we want to understand the architectural principles that have provided for its growth; and if we want to be sure that it supports the basic social values of trustworthiness, privacy, and respect for social boundaries, then we must chart out a research agenda that targets the Web as a primary focus of attention.

When we discuss an agenda for a science of the Web, we use the term "science" in two ways. Physical and biological science analyzes the natural world, and tries to find microscopic laws that, extrapolated to the macroscopic realm, would generate the behavior observed. Computer science, by contrast, though partly analytic, is principally synthetic: It is concerned with the construction of new languages and algorithms in order to produce novel desired computer behaviors

## Upcoming Events

[Memories for Life Colloquium: The Future of our Pasts](#)
12th December, 2006

## Latest News

[Southampton and MIT launch Web Science collaboration](#)
November 2, 2006

## News about WRSI on other sites

[Via Google News](#)

# DOWNHOMER
## Magazine

### Visit the Beaten Path

*A little part of Newfoundland and Labrador for people everywhere*

**Shoppe & Gallery**

**Newfoundland Kitchen**
*Our Guestbook*

**Life's Funny Experiences**

**Scramble Game**

**Trivia Game**

**Events Listing**

**Photo Gallery**

You are visitor
**799926**

CURRENT ISSUE

Shop in our On-Line
**CATALOGUE**
*Books • Music • Videos*

**LEGACY STONES**

*About our Magazine*     *Free Copy/Subscription*     *Letters to the Editor*     *Classifieds*

*Contact Us*     *Clubs and Associations*     *Search Engines*     *Porthole Links*

[Employment Opportunities with the Downhomer](#)

**vocm**

*Site maintained by* [*Vince Marsh*](#) *and* [*Grant Young*](#)
*Last updated November 26, 1999*

**2006**

# MAPPING PRACTICE 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!

# WEB ARCHIVAL LABOUR

| Knowledge Work | Examples |
|---|---|
| • Defining selection priorities | ‣ popularity, novelty, 'precarity' |
| • Allocating resources | ‣ staff, domain/host storage limits |
| • Analysing corpus; curating | ‣ link-based analysis, tagging |
| • Maintaining crawls | ‣ monitoring activities |

# WEB ARCHIVAL LABOUR

## Breakdown & Repair

- 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)

## Examples

‣ crawler traps

‣ missing capture elements

‣ patch crawling

‣ other quality assurance tasks

# SUMMARY & FUTURE WORK

- Complex system of knowledge/
  maintenance work for prioritising
  web archiving

- Archive is leveraging corpus and
  crowd for identifying domains

- Developed multiple tagging and
  analysis tools for curation

Further examining algorithmic (non-
human) labour

Expand to include other
communities, case studies

# ACKNOWLEDGEMENTS

EPSRC

Research Councils UK
Digital Economy
Transforming Business and Society

INTERNET ARCHIVE