

MAINTAINING THE WEB

Web Archiving, Labour and the Internet Archive

UNIVERSITY OF
Southampton

Web Science
DTC

Jessica Ogden (@jessogden)

Susan Halford

Les Carr

ASSOCIATION OF INTERNET RESEARCHERS (AOIR) CONFERENCE 2017 | TARTU, ESTONIA | 18 - 22 OCTOBER 2017

Library of Congress. Image: <https://flic.kr/p/dUt1ir>



The Internet Archive



Data Rescue NYC

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

WEB ARCHIVING

- a lot of rhetoric around web archives
- don't know a lot about *how* web archiving operates at scale, for whom and for what purpose
- inherently selective, contingent and \neq a copy
- increasingly used to make claims about the Web



PROBLEMATISING

<https://www.newscientist.com/article/dn4434-internet-mapping-project-weaves-colourful-web/>

CONCEPTUALISING MAINTENANCE & REPAIR

Web archiving as a form of maintenance
work for the archived Web.

Web archiving as a form of
repair for the Web.

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





association(of).internet.researchers

- academic association
- resource and support network
- critical and scholarly Internet research
- independent from traditional disciplines
- existing across academic borders
- international in scope

goals

faq

join

virtual graduate seminar (Spring 2000)

member database

downloadable flier (Acrobat PDF format)

listserv archive

action items - your help needed!

Heritrix: Admin Console - Mozilla Firefox

File Edit View Go Bookmarks Tools Help

http://foo.edu:8080/index.jsp

HERITRIX Status as of **Jun. 23, 2006 19:20:37 GMT** Alerts: no alerts

Admin Console CRAWLING JOBS RUNNING job: *150-sites-deciding-10MB*
0 jobs pending, 9 completed 9108 URIs in 4m43s (25.5/sec)

Console Jobs Profiles Logs Reports About Help

Crawler Status: **CRAWLING JOBS** | Hold

Jobs Memory
Running: *150-sites-deciding-10MB* 102217 KB used
0 pending, 9 completed 117056 KB current heap
Alerts: 0 (0 new) 236224 KB max heap

Job Status: **RUNNING** | Pause | Checkpoint | Terminate

Rates
25.5 URIs/sec (32.27 avg)
595 KB/sec (545 avg)

Time
4m43s elapsed
3m34s remaining (estimated)

Totals
downloaded 9108 **56%**
15997 total downloads
151 MB uncompressed

[Refresh](#)

[Shut down Heritrix software](#) | [Logout](#)

Done



ARCHIVE-IT Home Collections Crawls Archives ARS Help Center Welcome, AIT Demo

Home

Archive-It Demo Account

3.6 TB archived since Sep 22, 2010

Current Subscription

628.9 GB Archived

Data Budget Usage

▶ Current Subscription Details

▶ Past Subscription Totals

Active Collection List (3 Active Collections)

Type to Filter Active Collections

[Create a Collection](#)

[Download Collection List](#)

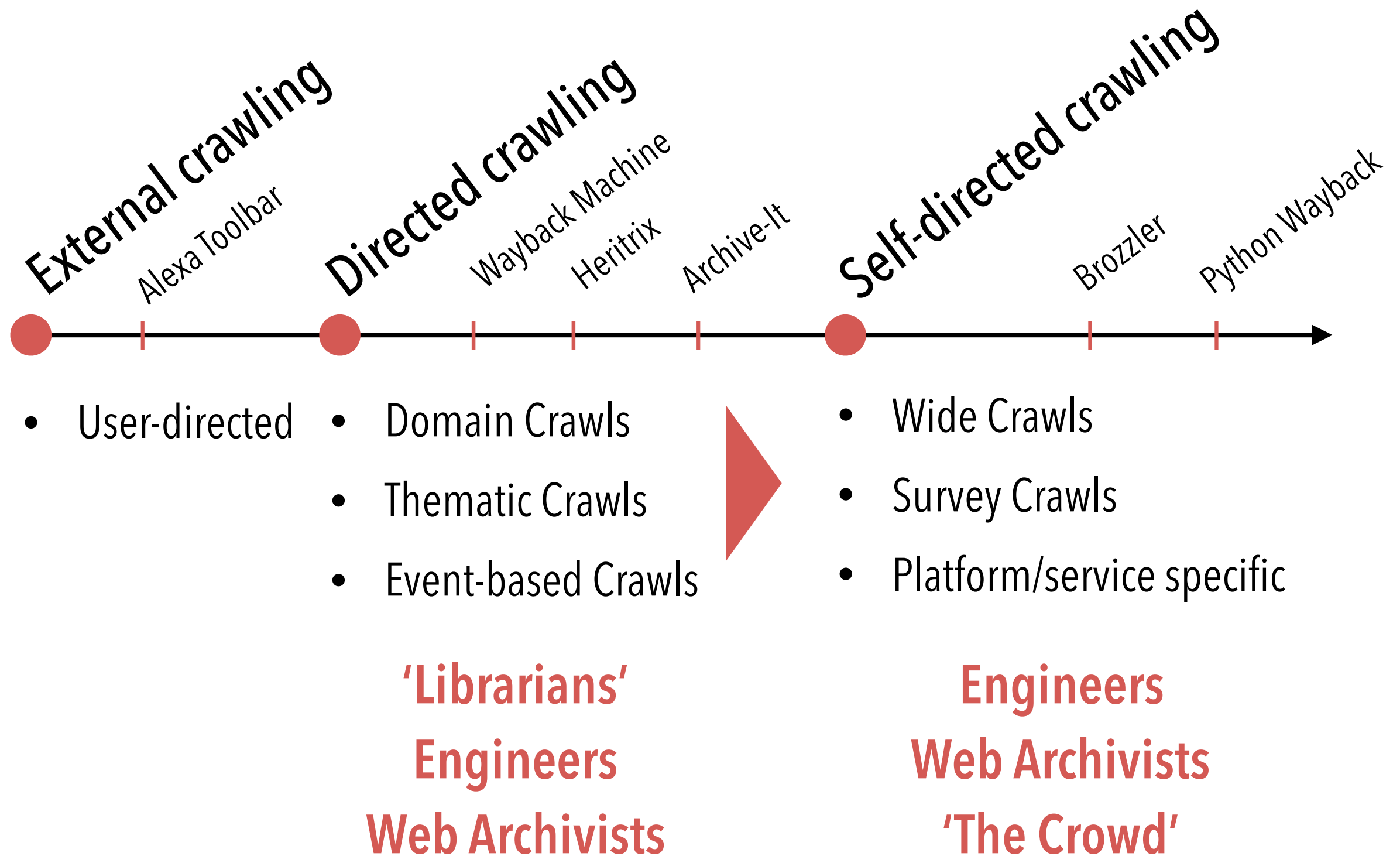
Collection Name	Data (this period)	Docs (this period)	Active Seeds	Last Crawl
Archive-It Websites and Press	72.4 GB	1,187,978	57	Jan 6, 2016
2015 United Nations Climate Ch	42.6 GB	726,589	5	Dec 10, 2015
Climate Change	1.9 GB	92,932	13	Dec 27, 2015

FMCCOWN AT ENGLISH WIKIPEDIA, [HTTPS://COMMONS INDEX.PHP?CURID=10700055](https://commons.index.php?curid=10700055)

2006

MAPPING ROLES & ACTIVITIES

"strategically mishmash"



*SCHEMATIC, NOT TO SCALE AND PROBABLY NOT EVEN CHRONOLOGICAL!

**WEB ARCHIVAL
LABOUR**



COLLECTION WORK

SELECTION & MODERATION WORK

REASSEMBLY WORK

REPAIR WORK

***WORK IN PROGRESS!**

COLLECTION WORK

- e.g. Wide Crawls

```
2011-06-23T17:12:08.802Z 200 1299 http://content-5.powells.com/robots.txt LREP http://content-5.powells.com/cgi-bin/imageDB.cgi?isbn=9780385518635 text/plain #014 20110623171208574+225 sha1:YI
UOKDGOLGI5JYHDTXRFFQ5FF4N2EJRV - -
2011-06-23T17:12:09.591Z 200 15829 http://www.identitytheory.com/etexts/poetics.html L http://www.identitytheory.com/ text/html #025 20110623171208546+922 sha1:7AJUMSDTOMT4FN7MBFGGNJU3Z56MLCMW
- -
```



SELECTION & MODERATION WORK

- e.g. Domain Browser Tool

“The domain browser manual tool is for identifying undesirable domains. It’s used to establish and prioritise ‘shades of gray,’ for example only crawl this site if there are no other sites to crawl. It’s used as a ranking mechanism for prioritising domains based on time, resources and place in the queue, as certain important URLs can get blocked by many instances of unimportant URLs. [...]Each domain/host is assigned a budget and the crawler is paused if it reaches its budget.”

SELECTION & MODERATION WORK

- e.g. Domain Browser Tool

“What we did was hired half a dozen people - they would just go through it and get the top 30,000 hosts [...] and they go through 4-5,000, that’s what one person can do in a month or so. And then we get actual human interaction to say yes, this is a good website. And then we would delete or modify or prioritise based on that input. So having humans actually spend a little bit of time at the top really helped. We’d love to do it further of course.”

SELECTION & MODERATION WORK

- e.g. Studying link graphs

“I do a lot of link analysis where I study the hyperlink structure of our crawls and try to figure out in certain pockets, use some rank methodologies to figure out ‘oh these are important resources,’ for instance they have a lot of links to them or traffic is really high - let’s seed the crawl with those. The most recent wide crawl I took the most linked to pages from every single website, so 230 million websites [...] and instead of crawling the Alexa top million, let’s crawl this bit.”

SELECTION & MODERATION WORK

- e.g. Studying link graphs

“The way I think of it [...] there’s three branches, there’s popularity, there’s novelty and there’s the risk of going away. How do you achieve that balance? You want to get stuff that people are using - not just junk that is on the Web that you’re just filling up the servers with that won’t ever be found useful, like calendar pages, things like that, crap [...] - there’s no novelty. It’s new? We want to make sure it’s preserved because it just came out, it’s a new article, it’s a new website. And then there’s the risk of going away [...] - if you’re going to shut down this service - Vine is going away - we jump in and crawl. So as we’re crawling the Web can we do a good job of sort of achieving that balance? We don’t quite know what the solution is to achieving that balance.”

'GOOGLE FOR THE ARCHIVED WEB'

Internet Archive Blogs



Save Page Now

SAVE PAGE

Capture a web page as it appears now for use as a trusted citation in the future.

Only available for sites that allow crawlers.

← In the news: Trump Archive, end-of-term preservation, & link rot

If You See Something, Save Son In the Wayback Machine



Posted on [January 25, 2017](#) by [Alexis Rossi](#)

In recent days many people have shown interest in copies of the web pages they care about most linked to – and they will continue to exist even if removed from the web.

There are several ways to save pages and the Wayback Machine. Here are 6 of them.

Too bad archive.org waxes political; and far left at that.

I really appreciate the fact as a conservative Christian, I can find content here.

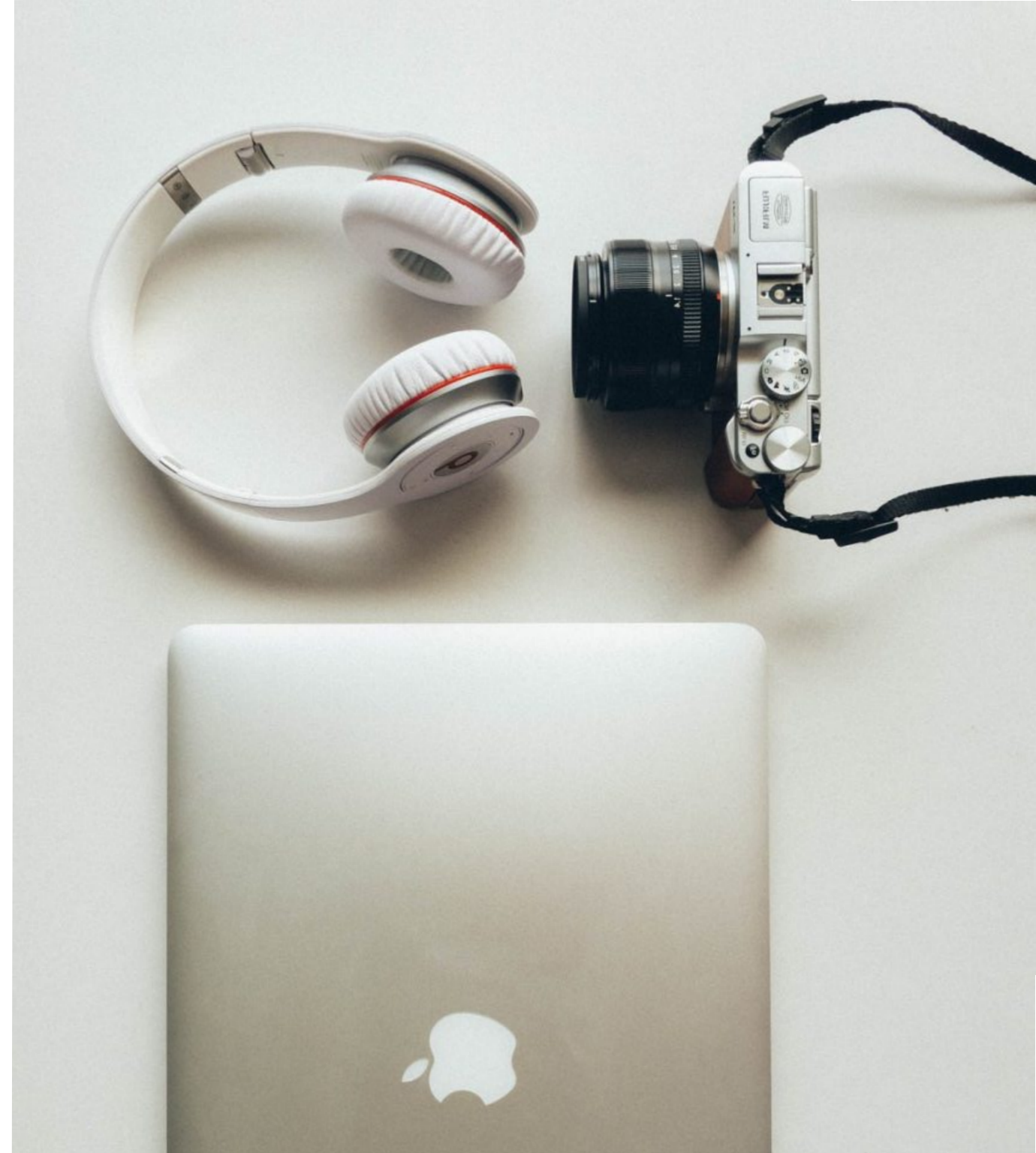
I wonder when that material will be removed because it is both Christian and Conservative?

One of the last bastions of public cyber space has gone off in HATE TRUMP land. How disappointing! Can't we keep politics out?

<https://blog.archive.org/2017/01/25/see-something-save-something/>

SUMMARISING

- tensions between automation and manual labour reflect complexities of work
- values embedded at inflection points in selection, collection
- underlying role of power in shaping web imaginaries - for whom?
- future work



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

This work was supported by the **UK Engineering and Physical Sciences Research Council** and the **Web Science Centre for Doctoral Training**, Grant No. EP/G036926/1.

The authors would also like to thank the **Internet Archive and staff** for opening their doors and being so generous with their time and feedback.

The logo for the Engineering and Physical Sciences Research Council (EPSRC), consisting of the letters 'EPSRC' in a bold, sans-serif font, underlined by two horizontal lines.