

Participatory Web Archiving

Opening the Black Box of Save Page Now

JESSICA OGDEN

University of Southampton

jessica.ogden@soton.ac.uk

@jessogden

ED SUMMERS

University of Maryland

edsu@umd.edu

@edsu

SHAWN WALKER

Arizona State University

shawn.w@asu.edu

@walkeoh

THE WEB THAT WAS - RESAW 2019 | AMSTERDAM | JUNE 22, 2019

Outline

- Introduce SPN Project
 - Centring participation - *participatory web archiving*
 - What is SPN
 - Related Work - IA/Wayback Machine, WA research tech, WA use
- Pilot RQs
- Methodology - collaborative ethos, data collection, analysis tools
- Foregrounding Limitations
- Results
 - What gets saved/when → diversity
 - Liveness/archival distribution → ephemerality
 - Role of automation → novelty (as one characteristic of contribution)
- Discussion
- Future Work



**We met at
Archives Unleashed**



How is SPN web archival
labour

How SPN responds to external
events such as disasters and social
movements?

Web Archive
built by Robots
and 1,000 librarians
Save Page Now



Image: David Rinehart (2016)

https://archive.org/details/ia20thanniversaryevent_images/page/n29



Brewster Kahle

@brewster_kahle

Following



651,621,510,000 web URL's now in the
Wayback Machine by [@internetarchive](#) .
Billions and Billions of web pages! users
hitting "save page now" at 100 per second:
web.archive.org

7:56 PM - 9 May 2018

68 Retweets 159 Likes



3



68



159



If You See Something, Save Something – 6 Ways to Save Pages In the Wayback Machine

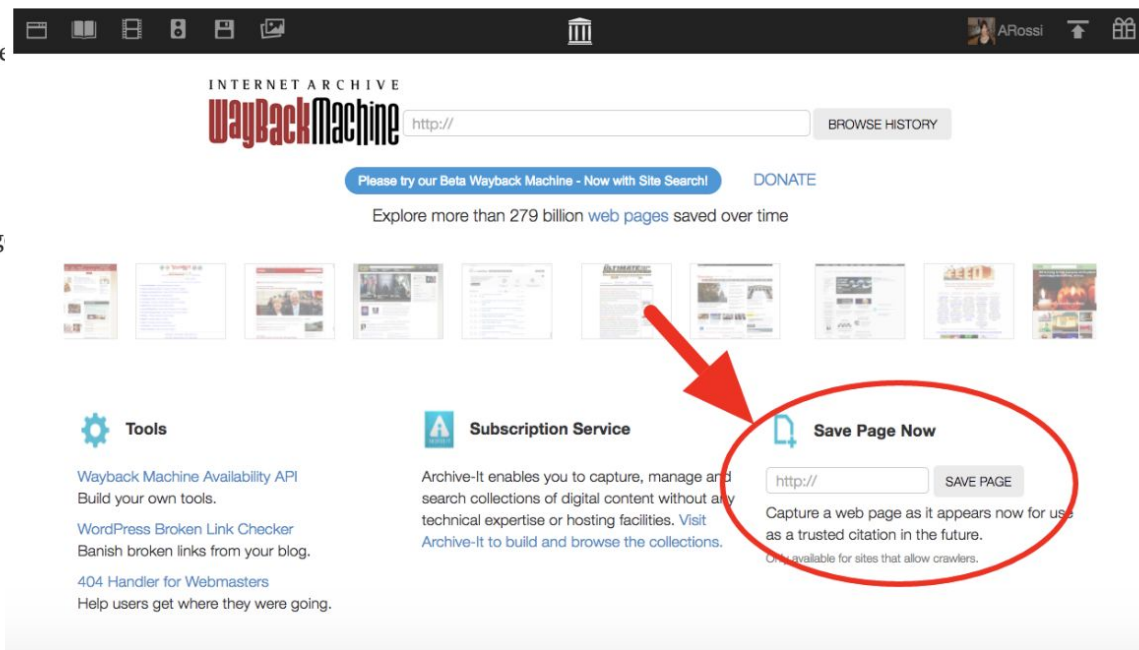
Posted on [January 25, 2017](#) by [Alexis Rossi](#)

In recent days many people have shown interest in making sure the [Wayback Machine](#) has copies of the web pages they care about most. These saved pages can be cited, shared, linked to – and they will continue to exist even after the original page changes or is removed from the web.

There are several ways to save pages and whole sites so that the Machine. Here are 6 of them.

1. Save Page Now

Put a URL into [the form](#), press the button, and we save the page permanent URL for your page.



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



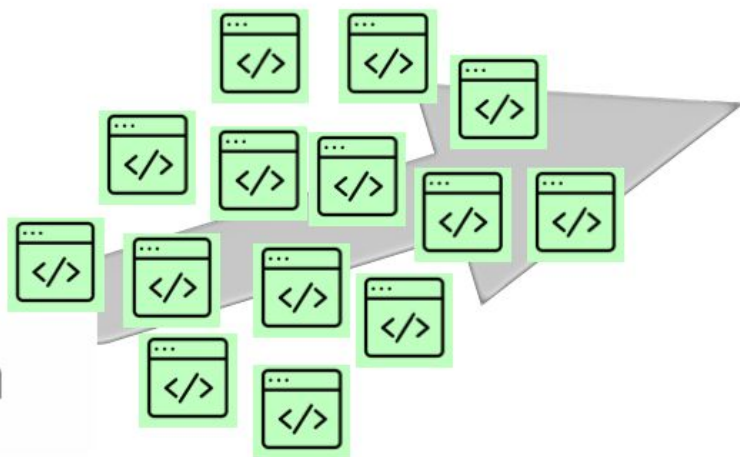
Browser



INTERNET ARCHIVE



Bot



SPN has changed over time

SPN v1

- Heritrix
- *No difference between web & API submissions*

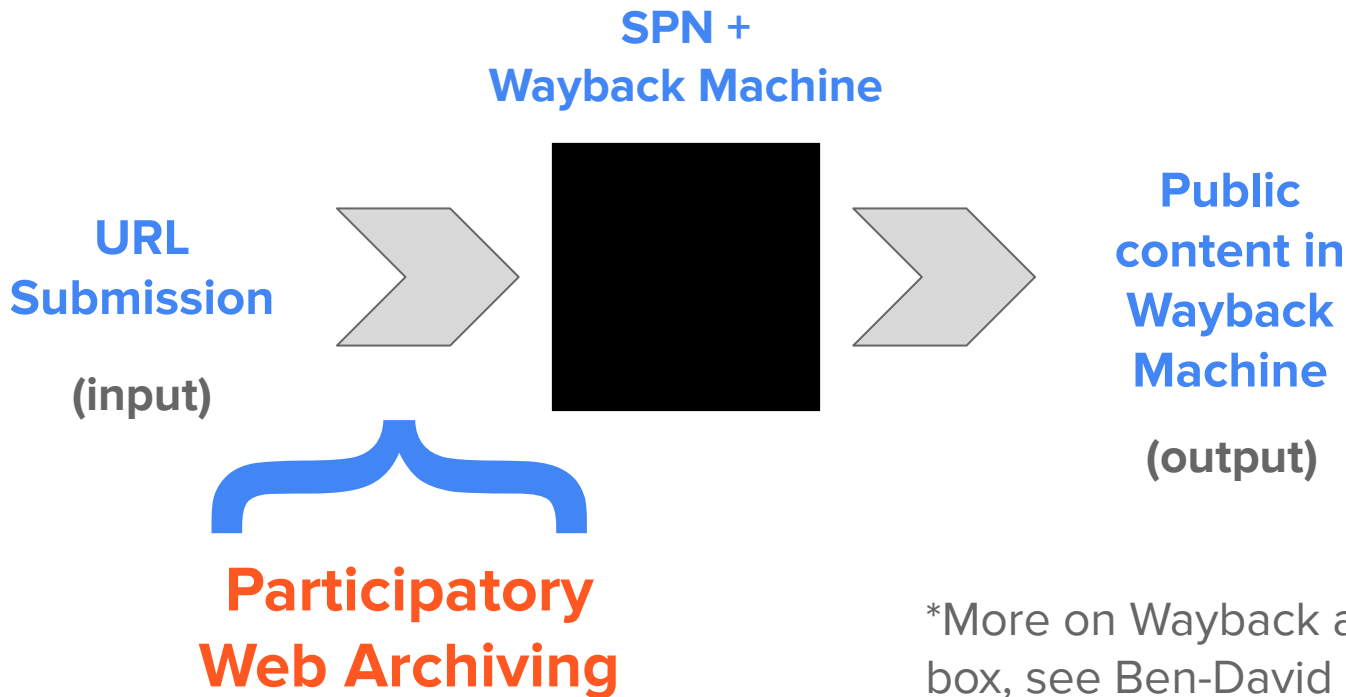
SPN v2

- Heritrix for API submissions
- Browser-based archiving for web submissions

SPN v3

- Server-based headless browser for API submissions
- Browser-based archiving for web submissions

SPN as black box



Motivating Research Questions

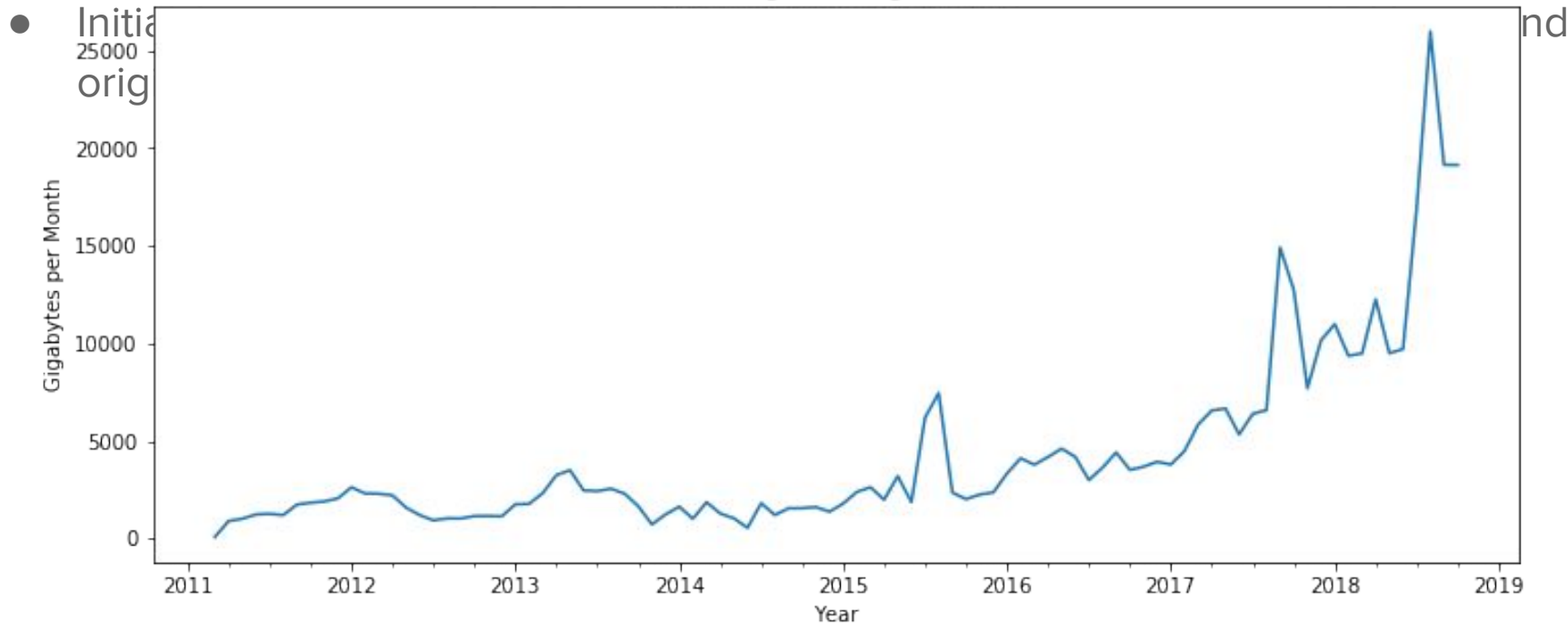
Aim: to understand SPN as form of participatory WA infrastructure

- **RQ1:** What is saved via SPN and how has the ‘collection’ changed over time?
- **RQ2:** To what extent are SPN resources available on the live Web and in other web archives?
- **RQ3:** In what ways is automation a factor in archival production and what purposes does it serve?

Methodology

Methodological Frame (High Level)

Save-Page-Now Ingest Rate



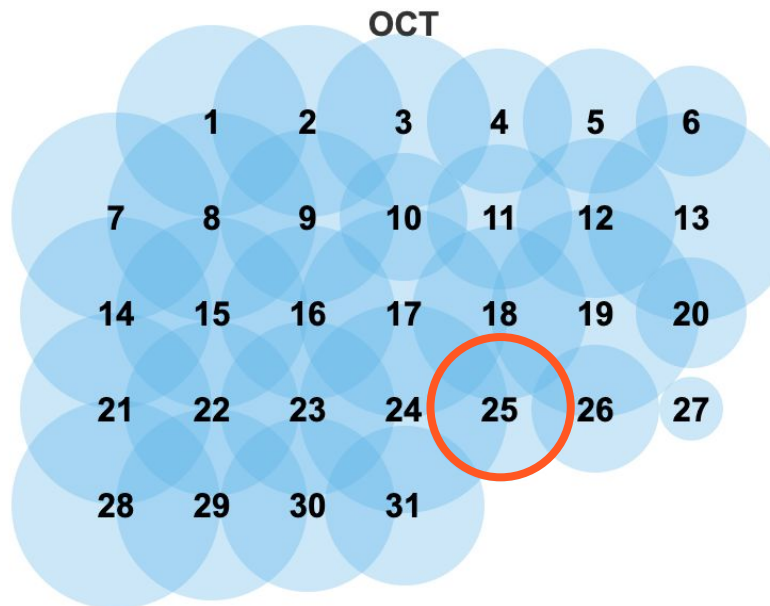


WayBackMachine

Sample:

- One day per year
- 5 years of data
- Oct 25, 2013 - 2018

SPN first made public on
homepage on October 25, 2013





WARCs
downloaded via
API



+



+



Related Work

- Studying Web Archival Practices
 - ?Provenance (Maemura et al.)
 - Seeding practices in WAs (Summers and Punzalan 2017)
 - Ethnographic research at the Internet Archive - characterising ‘backstage’ of Wayback Machine and projected role of SPN in diversifying seeding (Ogden et al. 2017)
 - *(That new report...)*
- Characterising the Internet Archive Wayback Machine (IAWM)
 - language and geographic distribution of IAWM content (AlSum et al., 2014; Thelwall and Vaughan, 2004)
 - IAWM use - examination of IA access logs (AlNoamany et al., 2013); circulation of IAWM links on live Web (Zannettou et al., 2018)
 - IAWM as sociotechnical assemblage and web archival research as forensics (Ben-David and Amram 2018)
- (But nothing studying SPN specifically - SPN as platform?)

Pilot Research Question - v3

Aim: to understand SPN as form of participatory WA infrastructure

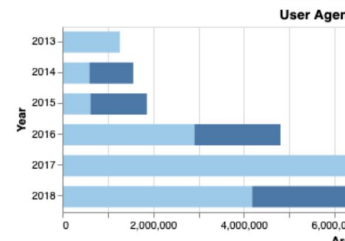
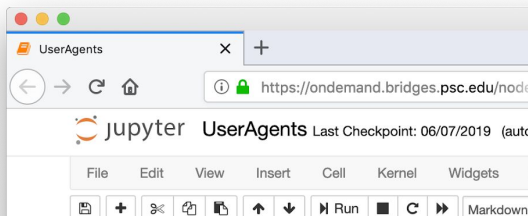
- RQ: What is the role of automation in the use of SPN?

We will look at this through 3 dimensions:

- Diversity - what is saved
- Liveliness/Ephemerality - how at risk are target resources
- Novelty - how new/novel are target resources

Methodological Frame (High Level)

- Importance of collaborative working



This is a helpful visualization because it shows that account for about 1/2 of the total requests. So there

But 2017 is unusual in that by far the top-10 User-A of WARC request records that year, as compared to same trend.

```
In [9]: from os.path import getsize  
  
sizes = {'year': [], 'gb': []}  
for year in range(2013, 2019):
```

The screenshot shows a Zoom meeting window with three participants: a woman, Shawn Walker, and a man. Below the video feeds is a terminal window displaying a list of user agents. The terminal output shows a series of user agent strings, including 'mediawords bot', 'Mozilla/5.0 (Windows NT 6.1; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/69.0.3497.100 Safari/537.36', and 'python-requests/2.18.1'.

```
archive.org/details/archive.org_bot)  
User-Agent: mediawords bot (http://cyber.law.harvard.edu)  
User-Agent: Mozilla/5.0 (Windows NT 6.1; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko)  
Chrome/69.0.3497.100 Safari/537.36  
User-Agent: mediawords bot (http://cyber.law.harvard.edu)  
User-Agent: Mozilla/5.0 (Windows NT 6.1; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko)  
Chrome/69.0.3497.100 Safari/537.36  
User-Agent: Mozilla/5.0 (Windows NT 6.1; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko)  
Chrome/69.0.3497.100 Safari/537.36  
User-Agent: Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko)  
Chrome/69.0.3497.100 Safari/537.36  
User-Agent: Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/537.36 (KHTML, like Gecko) Chrom  
e/60.0.0.1508 Safari/537.36  
User-Agent: Mozilla/5.0 (X11; Linux x86_64) AppleWebKit/537.36 (KHTML, like Gecko) SamsungBr  
owser/7.0 Chrome/59.0.3071.125 Safari/537.36  
User-Agent: Mozilla/5.0 (Macintosh; Intel Mac OS X 10_12_6) AppleWebKit/537.36 (KHTML, like  
Gecko) Chrome/69.0.3497.100 Safari/537.36  
User-Agent: python-requests/2.18.1  
User-Agent: Mozilla/5.0 (compatible; archive.org bot; Wayback Machine Live Record; +http://a
```

URLs

What URLs in the WARC data represent user intention to archive something?

HTTP Content-Type

HTTP Status Code

Eliminating .js, .css, .jpg, etc

Need for WARC request records.

Top Hostnames / URLs lists (slippery).

- Problems with isolating the targeted intention (wanted to use referrers, probs of using mimetype?)

RQ1: SPN Collection Development



→ Hey guys!



Check out my imageboard

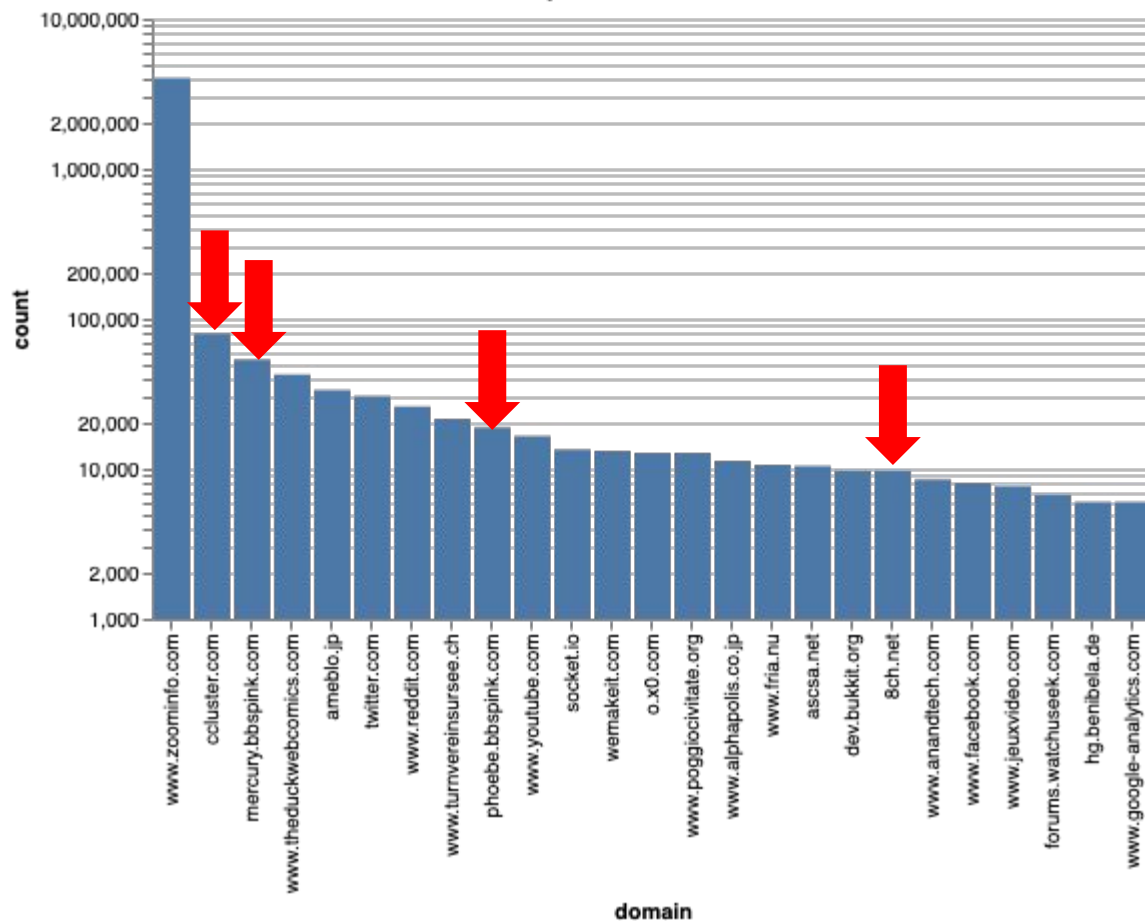
It's called 600chan

It has a /b/ board!!!!

It's also better than 4chan and 8chan combined

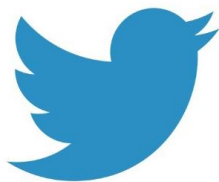
Plz visit my website, no one goes there...

Popular Domains 2017

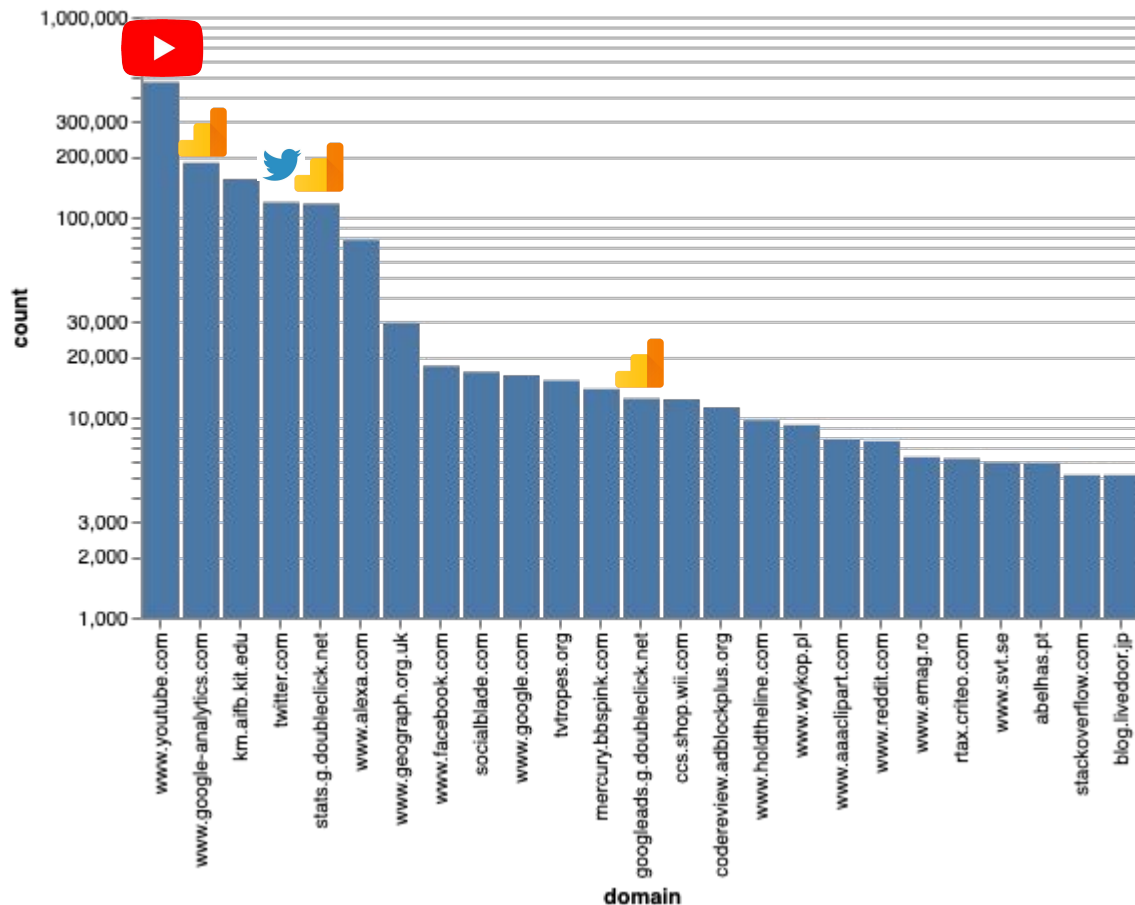




Google Analytics



Popular Domains 2018



User-agents and Automation

	date	user_agent
0	2018-10-25T11:18:39Z	python-requests/2.13.0
1	2018-10-25T11:37:07Z	Mozilla/5.0 (compatible; archive.org_bot; Wayback Machine Live Record; +http://archive.org/details/archive.org_bot)
2	2018-10-25T12:00:22Z	python-requests/2.18.1
3	2018-10-25T12:01:45Z	Mozilla/5.0 (compatible; archive.org_bot; Wayback Machine Live Record; +http://archive.org/details/archive.org_bot)
4	2018-10-25T12:11:19Z	python-requests/2.18.1
5	2018-10-25T12:21:36Z	python-requests/2.18.1
6	2018-10-25T12:31:56Z	python-requests/2.18.1
7	2018-10-25T12:43:07Z	python-requests/2.18.1
8	2018-10-25T13:19:18Z	Chrome 41.0.2227.0
9	2018-10-25T13:54:52Z	Firefox 40.1
10	2018-10-25T13:58:42Z	Mozilla/5.0 (compatible; archive.org_bot; Wayback Machine Live Record; +http://archive.org/details/archive.org_bot)
11	2018-10-25T13:58:43Z	Mozilla/5.0 (compatible; archive.org_bot; Wayback Machine Live Record; +http://archive.org/details/archive.org_bot)
12	2018-10-25T14:32:02Z	Safari 5.1.7
13	2018-10-25T15:13:25Z	Chrome 41.0.2227.0
14	2018-10-25T15:37:12Z	Mozilla/5.0 (compatible; archive.org_bot; Wayback Machine Live Record; +http://archive.org/details/archive.org_bot)
15	2018-10-25T15:37:12Z	Mozilla/5.0 (compatible; archive.org_bot; Wayback Machine Live Record; +http://archive.org/details/archive.org_bot)
16	2018-10-25T15:53:29Z	Firefox 33.0
17	2018-10-25T16:32:46Z	Chrome 41.0.2228.0
18	2018-10-25T17:08:11Z	Chrome 41.0.2228.0
19	2018-10-25T17:44:50Z	Firefox 36.0
20	2018-10-25T18:21:49Z	Chrome 41.0.2228.0
21	2018-10-25T18:58:06Z	Safari 6.0
22	2018-10-25T19:34:05Z	Chrome 41.0.2227.1

Automated API capture 20 seconds after tweet

A very big part of the Anger we see today in our society is caused by the purposely false and

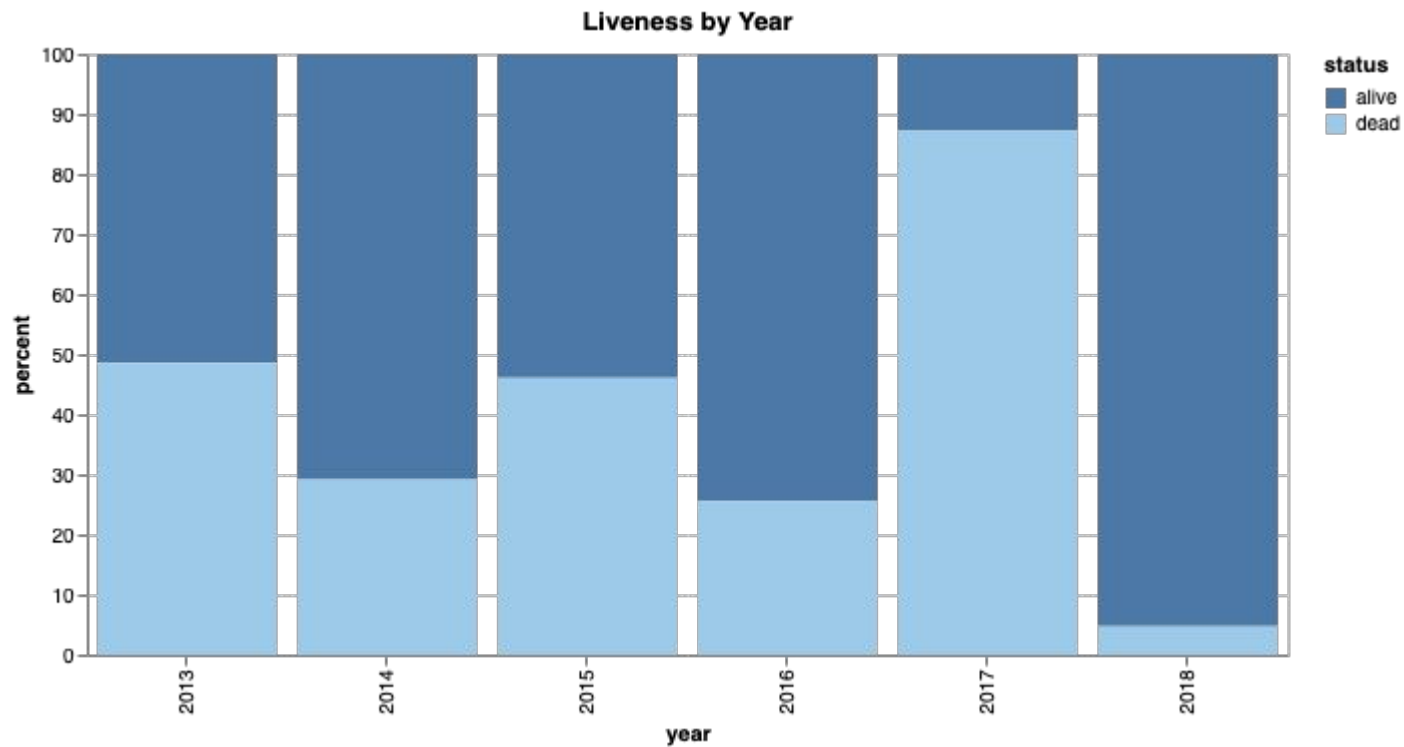
User-Agent spoofed automated API capture

No user-agent was supplied

All of these are automated captures

RQ2: Measuring Live(li)ness

- Sample of the URLs in dataset
- 2017 - ZoomInfo URLs die



RQ3: Assessing the Role of Automation

Assessing the Role of Automation

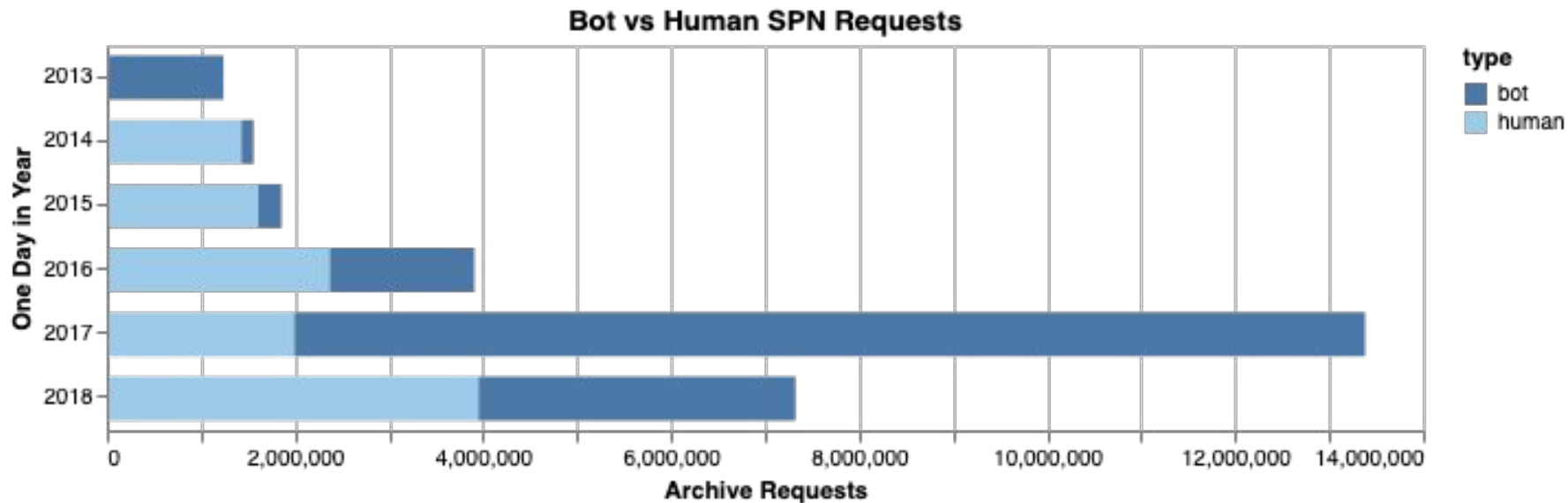
1) Hypothesised two types of SPN interactions -

- Software as **intermediary** (home page tool, browser extension) - SPN *mediates* intent of user and transports to archive
- Software as **actant** (API, cron) - SPN *transforms* selection and transports to archive

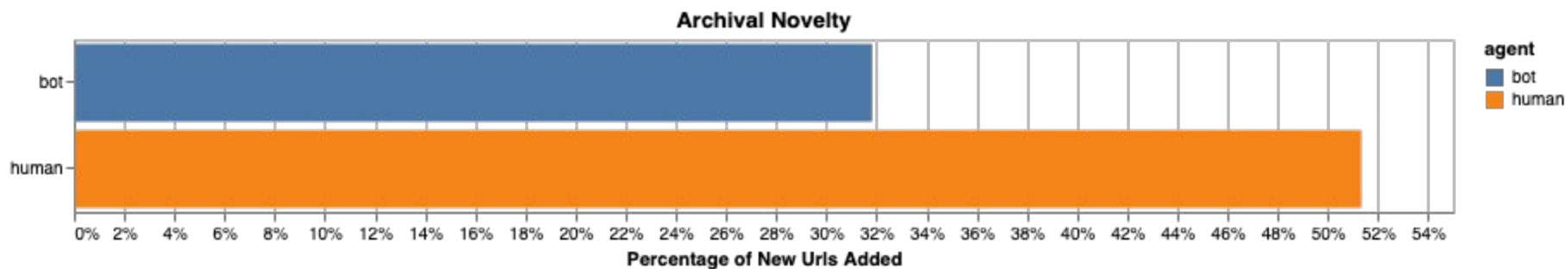
2) User-Agents as proxy for detecting automation

- User-Agents and grouping User-Agent Families
- Percentage of traffic from different types of Users

Bots vs Browsers?



Novelty



Discussion - Reading the Tea Leaves (methods)

- Assessing archival intention - Answering questions like ‘why’ require qualitative research methods to identify why things are being archived.
 - How to combine with quantitative?
- Zooming in / Zooming out - close/distant reading (Nicolini 2009)
 - Tools - usefulness of AUT and ArchiveSpark; python + jupyter notebooks made results more legible, accessible for collaboration
 - Detecting anomalies - e.g. Zoominfo - ?impact of sample
- Shifting infrastructure - added difficulty of studying a system like SPN over time, in a constant ‘state of becoming’ (Barad 2003)
 - e.g. the way that SPN records user agents changed over time - highlighted dangers of looking at data as disconnected from software that produces it
 - Trace experiments helped (some). SPN3 is coming...

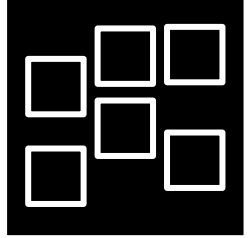
Discussion - Reading the Tea Leaves

- Diversity
- Ephemerality
- Novelty

Limitations

- Sample strategy - it's **one day a year**.
- Detecting automation
- Conceptualizing web archival attributes (e.g. 'novelty')
- URLs - delineating targets → inferring intentionality

Reflexivity as Strategy



- Complex layering of data and findings (boxes within boxes), ways that proxies/data abstraction often creates more questions + uncertainty
- Problematizing *situatedness* (Haraway 1988) of so-called big data (+ extra complexities of SPN) - recognizing that data views are always partial, ‘*cooking data with care*’ (Bowker 2005; Geiger and Halfaker 2018)
- Value of triangulation, epistemological flexibility - WA research as *boundary work* (Star and Griesemer 2015; Gieryn 1983) - requires a lot of translation

Future work...

Future Work

Wider statistical analysis?

warcio toolkit?

Interviews (Brügger, 2012)

Access to WARC data in web archives.

Research raised more questions than it answered (as typical of WA research)

Acknowledgements



XSEDE

Extreme Science and Engineering
Discovery Environment

This work used the Extreme Science and Engineering Discovery Environment (XSEDE) Bridges and Bridges Storage at the Pittsburgh Supercomputing Center through allocation TG-ECS180012. Extreme Science and Engineering Discovery Environment (XSEDE), which is supported by National Science Foundation grant number ACI-1548562.