# Digital Humanities: Birds of a Feather session

eResearch NZ 2019

Auckland 18-20 February 2019

Steve Knight, Programme Director, Preservation Research & Consultancy, National Library of New Zealand
steve.knight@dia.govt.nz

Today – 3.30 – 4.30

3.30    Introduction and welcome

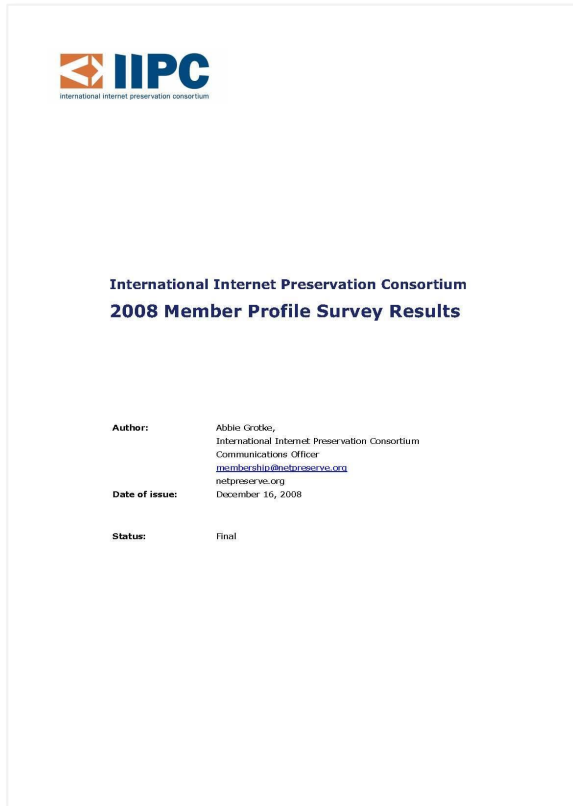3.35    What w ewant to have a look at today

3.40    A whip round the room

3.45    Three brief commentaries

4.00    Group exercise

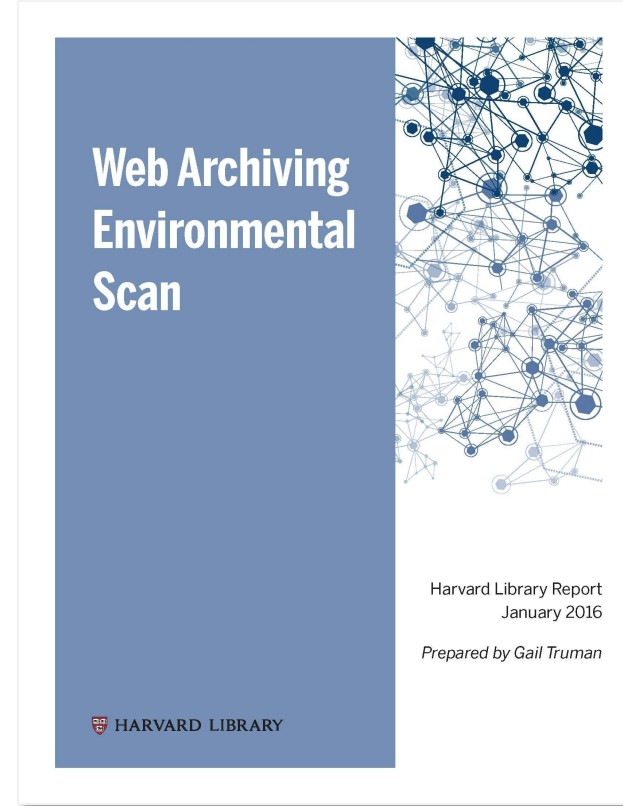4.20    Wrap-up and next steps

# Landscape of web archiving activities



IIPC Survey, 2008
IIPC Survey, 2018

Internet Memory Foundation, 2010

Harvard Library, 2016

National Digital Stewardship Alliance Web Archiving Surveys
2011, 2013, 2016, 2018

Born Digital Legal Deposit Policies and Practices, 2017

Preserving Online Multiple Information: Towards a Belgian Strategy, 2018
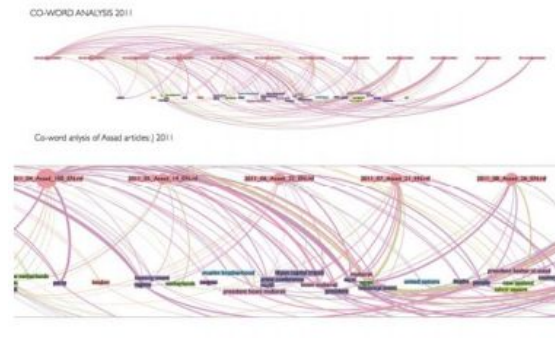
# Case studies

## LINK ANALYSIS



Hannes Mühleisen
Visualization of linking between websites of different languages,
Babel 2012 Web Language Connections

https://github.com/norvigaward/2012-naward25/wiki/Babel-2012---Web-Language-Connections
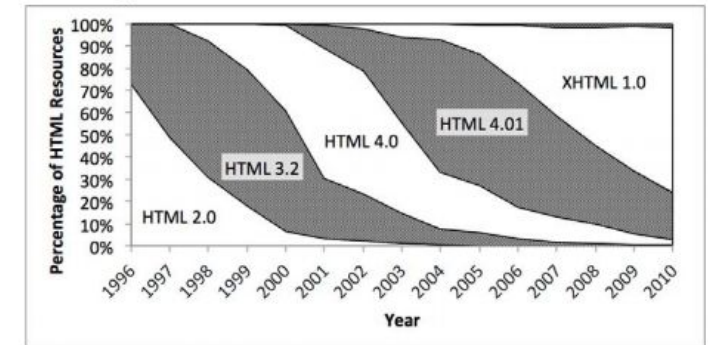
## LINK ANALYSIS



N-gram Search, UK Web Archive (British Library)

Searching the (News) Archives, Web Archive Retrieval Tools (University of Amsterdam)

Sentiment Analysis and the Reception of the Liverpool Poets, Helen Taylor (Royal Holloway)

Footprint of the world top companies, Alexandre Marah and Ferenc Szabó (University of Tewnte)

## ANALYSIS OF TECHNOLOGY TRENDS



HTML Version Usage Over Time, UK Web Archive

Format Profile, UK Web Archive (British Library)

Web Data Commons, Christian Christian Bizer et al. (University of Mannheim & Karlsruhe Institute of Technology) http://webdatacommons.org/

An analysis of the use of JavaScript libraries on the web, Dennis Pallett et al. (University of Twente)

Discovery & Access, NLNZ, 16 Nov 2018

# Researching web archives

**SURVEYS, REPORTS, PRESENTATIONS** (selected)**:**

- Maria-Dorina Costea. **Report on the scholarly use of web archives**. NetLab. 2018.

- Emily Maemura. **Explorations of Netarkivet: Preliminary Findings**. 2018.

- Janne Nielsen, **Using web archives in research**, NetLab 2017.

- Web Archives as Scholarly Sources: Issues, Practices and Perspectives.

- Harriet Riley & Mark Crookston. **Awareness and use of the New Zealand Web Archive**. 2015.

- Meghan Dougherty, Eric T. Meyer, Christine McCarthy Madsen, Charles van den Heuvel, Arthur Thomas, and Sally Wyatt. **Researcher Engagement with Web Archives: State of the Art**. JISC Report, 2010.

- Miguel Costa & Mario J Silva**. Understanding the Information Needs of Web Archive Users**. 2010.

- Jane Winters. **Big UK Domain Data for the Arts and Humanities**. IIPC WAC 2015.

- Helen Hockx-Yu. **Up close and personal - Researchers and the UK Web Archive Project**. IIPC WAC 2011.

- Jefferson Bailey. **Program Models for Research Services**. IIPC WAC 2016.

**PROJECTS, NETWORKS, SERVICES, CASE STUDIES** (examples):

- **ARS: Archive-It Research Services**

- **AU**: **Archives Unleashed Project**

- **RESAW.eu**: **Research Infrastructure for the Study of Archived Web Materials**, 2012-

- **BnF**: **Néonaute** 2018

- **ATI & UKWA**: Detecting semantic shift in large corpora, 2018

- **Arquivo.pt**: **Investiga XXI**, 2017 & **Awards**, 2018

- Niels Brügger and Ralph Schroeder, eds., **The Web as History: Using Web Archives to Understand the Past and the Present**. London: UCL Press. 2017.

- **BUDDAH**: **Big UK Domain Data for the Arts and Humanities**, 2013-4 □ **SHINE**

- **WebART**: **Web Archive Retrieval Tools, 2012-2016**

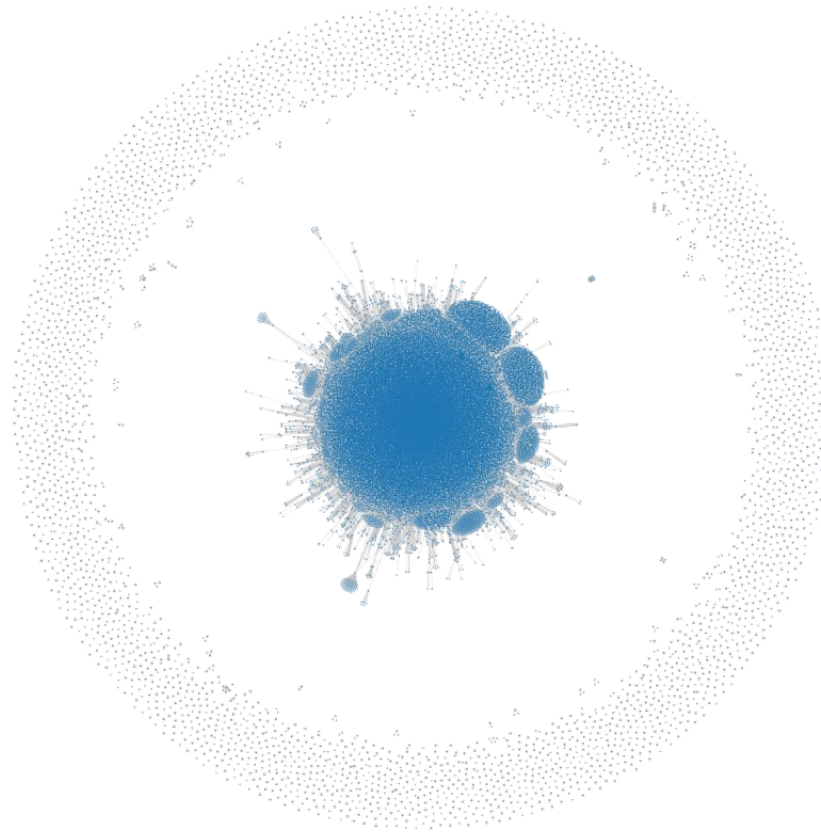# Program Models for Research Services, 2016 – lessons learned

- Researchers don't always know what they want

- Researchers default to wanting access to raw/all data

- Researchers will have varying levels of technical resources or support

- Address upfront issues of technical proficiency, nonarchive technical support and/or methodological stuff

- Will require reference/resources to explain and contextualize web archive tools and processes

- More data doesn't equal better analysis

- Focus on derivation, portability, and access

- Focus on scalable partnerships & decentralization

- Research support expectations often != with available resources or services

- Research methodologies (conceptual, practical, technical) often != with data, collecting, tools

- Service models or death (though yet to emerge for most data-driven LAM-ish research)

Jefferson Bailey. **Program Models for Research Services**. IIPC WAC 2016.

# Use cases: Buddah

**Big UK Domain Data for the Arts and Humanities**

Institute of Historical Research, UoL
Oxford Internet Institute
Aarhus University
British Library

**Topics**
disability
Euroscepticism
public archaeology
Ministry of Defence
companies, corporations
reception of Beat literature
French expatriates in London
poetry community & networks

**Data:** the UK web domain crawl 1996-2013

**Project team**

Jane Winters
Helen Hockx-Yu
Eric Meyer
Ralph Schroeder
Niels Brügger
Jonathan Blaney
Andrew Jackson
Peter Webster

Rainer Simon: Visualisation of links between websites from the UK crawled during 1996, Discovery & Access, NLNZ, 16 Nov 2018

# Arquivo.pt

- **DIOGO DUARTE**, The Study of punk culture through the Portuguese Web Archive

- **DIOGO SILVA DA CUNHA**, Today's news to be forgotten tomorrow?

- **RICARDO BASÍLIO**, Memory of the online presence of a Faculty: an exhibition

netpreserveblog.wordpress.com/tag/investiga-xxi

ARQUIVO.PT

# Questions to dhn researchers

**DHN — DIGITAL HUMANIORA I NORDEN — DIGITAL HUMANITIES IN THE NORDIC COUNTRIES**

## DATA

- **What** to archive (seed selection)?

- **Thematic crawls** – mainstream, "alternative" or both?

- Types of **datasets**: raw data / WARC files?

- Crawling: frequency, depth, width, completeness (formatting, images, etc.) – is **documentation** concerning this of interest to researchers? If so, how detailed? In the form of an overview or more closely attached to the harvested material?

## CONTENT

- How should we focus on **validity and reliability** of the archives?

- How should we **document** the way our archives (and single crawls) have been constructed?

## TOOLS

- Our tools or your tools?

- Are there any tools you want us to make available which would be useful without stepping into legal issues (e.g. **N-gram services**)?
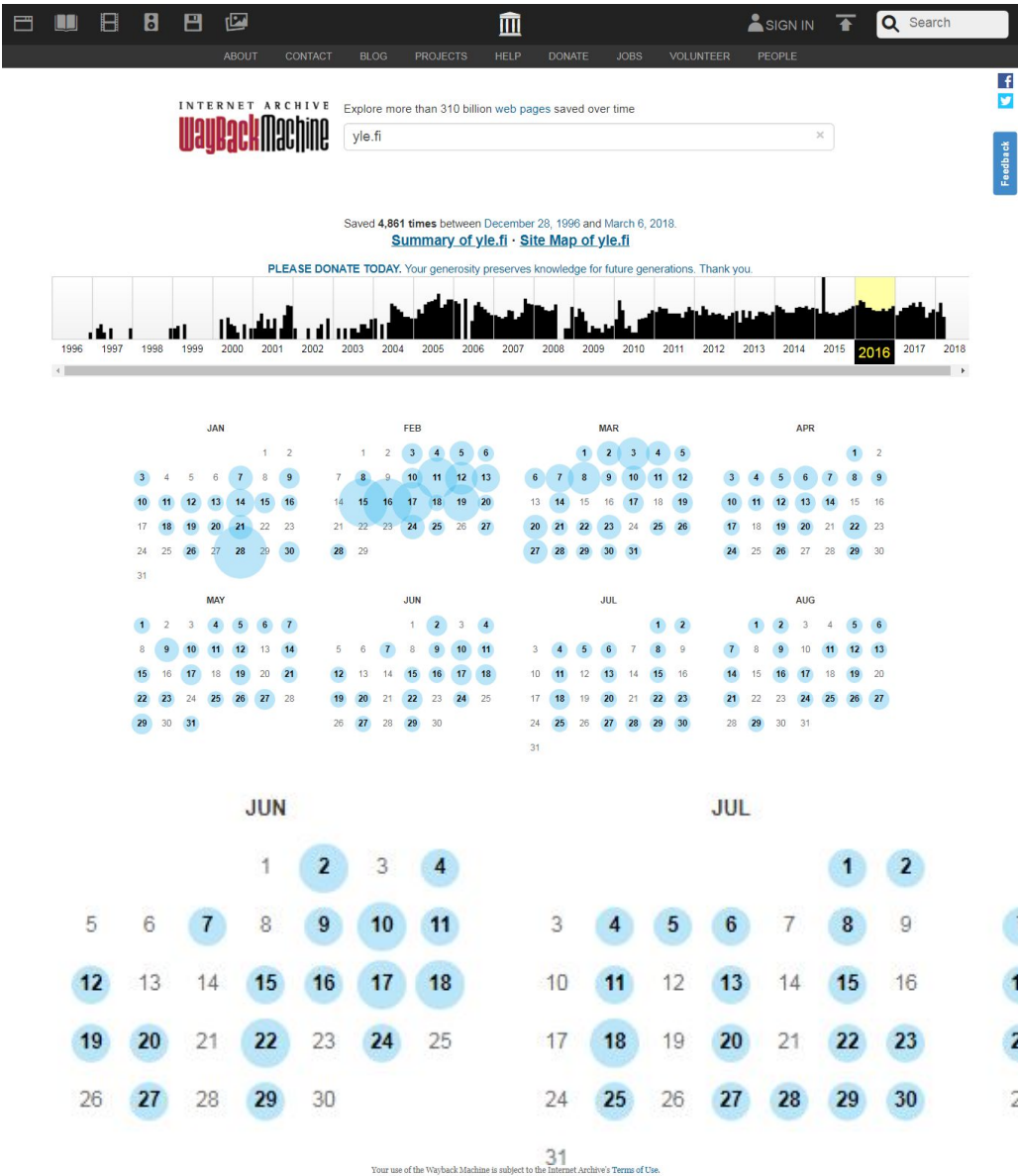
Vefsafn.is

Suomalainen verkkoarkisto

Norsk Nettarkiv

Kulturarw3

netarkivet.dk

**WEB ARCHIVES:**
**WHAT'S IN THEM FOR DIGITAL HUMANISTS?**
PANEL ON WEB ARCHIVING IN THE NORDIC COUNTRIES

National Library of Finland, **Lassi Lager**
National Library of Norway, **John Erik Halse**
National Library of Sweden, **Pär Nilsson**
The Royal Danish Library, **Caroline Nyvang**

# Looking for Brexit in the Finnish web archive
# DHN2018: Hacking the News



Front page of yle.fi: captures in IA vs FWA (Brexit)

http://verkkoarkisto.kansalliskirjasto.fi/wayback/*/yle.fi

132 captures (FI) vs. 68 (IA)

WARC files: 377MB extracted for this timelapse via AUT for yle.fi (+ 3 newspapers: satakunnankansa.fi, aamulehti.fi, pohjolansanomat.fi)

# JISC Report, 2010: Key points

 gap between the ***potential* community of researchers** who have good reason to engage with creating, using, analysing and sharing web archives, and **the *actual* (generally still small) community of researchers currently doing so**

- **Encourage the creation of communities** that increase the accessibility and usability of web archiving tools

- Efforts should be made to diversify tool and interface development beyond preservation and into use

- Tools should be developed that are able to execute query searches over multiple web archives

- Standards, protocols and methods of quality control are need for interoperability, but not at the cost of flexibility

- The possibilities of web archives should be communicated to a much broader research community

Meghan Dougherty, Eric T. Meyer, Christine McCarthy Madsen, Charles van den Heuvel, Arthur Thomas, and Sally Wyatt. **Researcher Engagement with Web Archives: State of the Art**. JISC Report, 2010.

# RESEARCH WORKING GROUP

Scope
- The Research Working Group seeks to promote the use of web archives and IIPC collections among researchers, share information about web archiving research projects at IIPC member organisations, including workflows and lessons learnt, and facilitate ways for dissemination and discussion of use cases.

- The group will collaborate with relevant research communities to help make the IIPC collaborative collections and IIPC member collections available to researchers. The Research Working Group will include not only IIPC members  but it will be open to researchers.

# The toolkit

| Platform/UI | Extraction | Textual analysis | Graphing |
|---|---|---|---|
| Shine https://github.com/ukwa/shine | Web-archive discovery https://github.com/ukwa/webarchive-discovery | Spacy https://spacy.io/ | Gephi https://gephi.org/ |
| SolrWayback https://github.com/netarchivesuite/solrwayback | Archives Unleashed Toolkit https://github.com/archivesunleashed/aut | Sem https://github.com/YoannDupont/SEM | SigmaJS http://sigmajs.org/ |
| Warclight https://github.com/archivesunleashed/warclight | Web-archive indexing https://github.com/ikreymer/webarchive-indexing | Open NER http://www.opener-project.eu/ | |
| UKWA UI https://github.com/ukwa/ukwa-ui | | Stanford CoreNLP https://stanfordnlp.github.io/CoreNLP/ner.html | |
| AUT Cloud https://github.com/archivesunleashed/aut | | Voyant Tools https://voyant-tools.org/ | |
| | | http://tapor.ca/home | |
| | | https://cloud.archivesunleashed.org/derivatives/text-sentiment | |
| | | https://cloud.archivesunleashed.org/derivatives/text-antconc | |