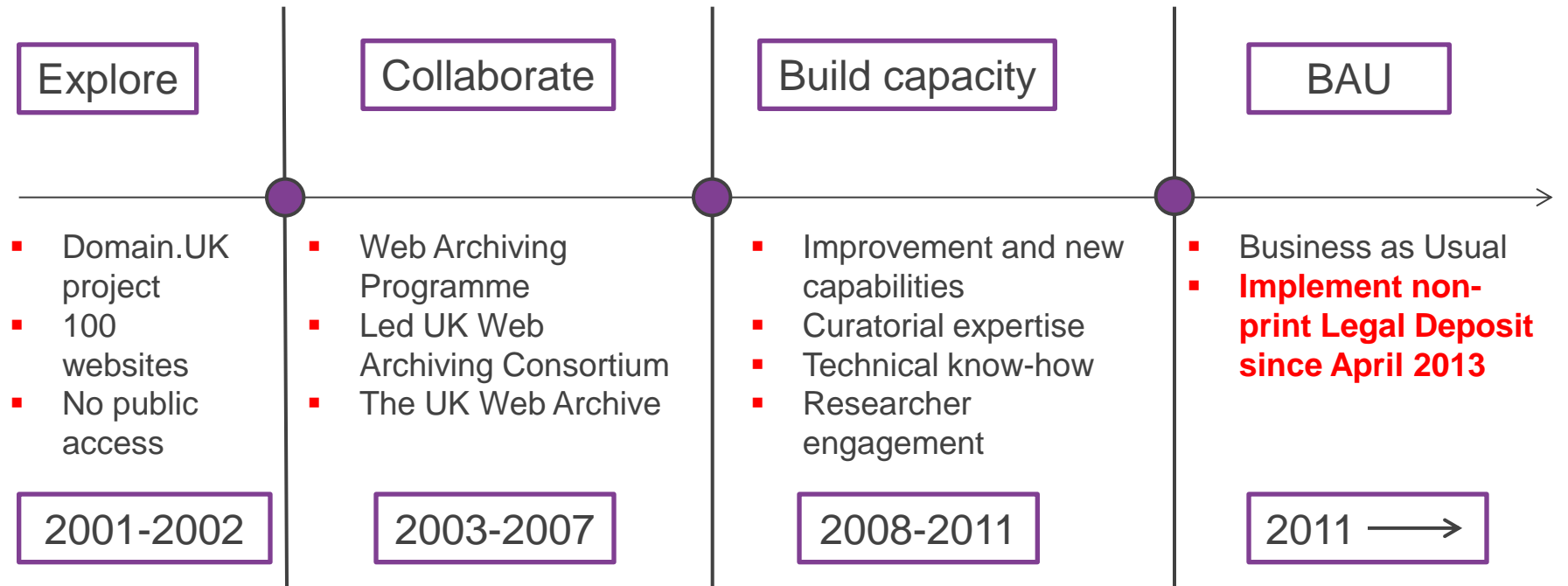# 10 Years of archiving the UK Web

Helen Hockx-Yu

Head of Web Archiving, British Library

# Web Archiving @ the British Library

| Explore | Collaborate | Build capacity | BAU |
|---|---|---|---|
| ▪ Domain.UK project<br>▪ 100 websites<br>▪ No public access | ▪ Web Archiving Programme<br>▪ Led UK Web Archiving Consortium<br>▪ The UK Web Archive | ▪ Improvement and new capabilities<br>▪ Curatorial expertise<br>▪ Technical know-how<br>▪ Researcher engagement | ▪ Business as Usual<br>▪ **Implement non-print Legal Deposit since April 2013** |
| 2001-2002 | 2003-2007 | 2008-2011 | 2011 ⟶ |

Collect UK digital heritage and provide continued access to archived web resources
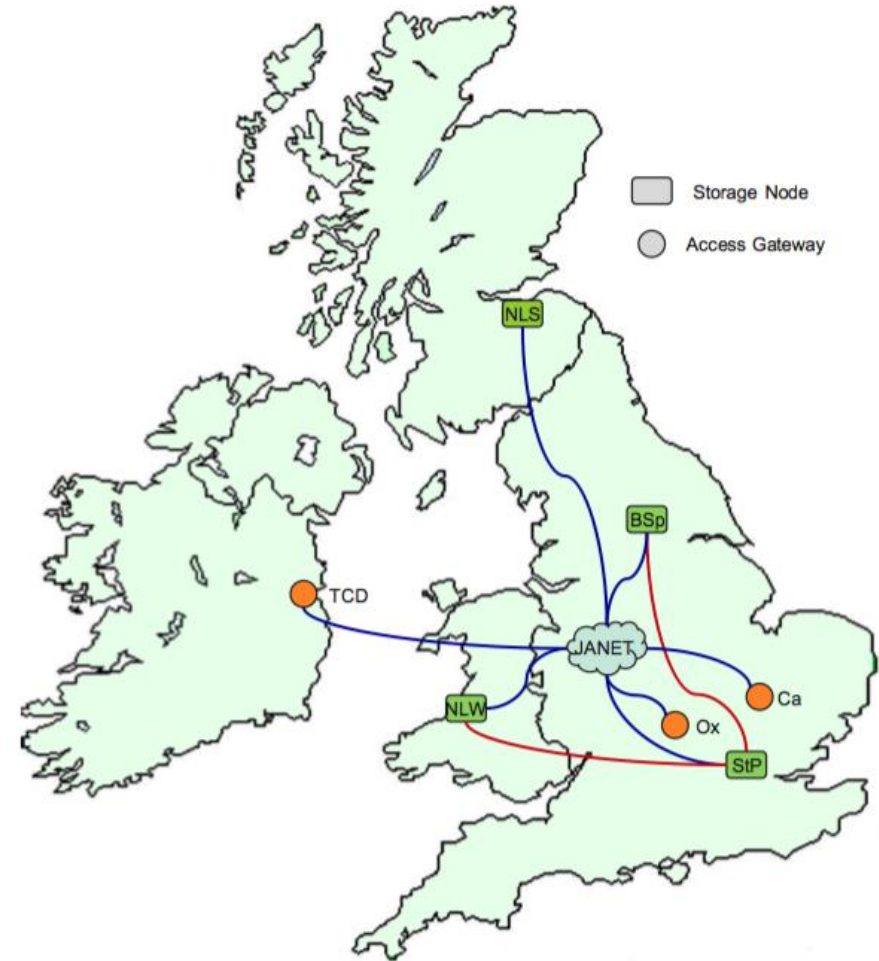
# The UK Web Archiving Consortium

# 6th April 2013…

- Legal Deposit Libraries (Non-Print Works) Regulations 2013

- Extension of existing legal framework

- Systematic collection of UK's published output for heritage & preservation

- By 6 UK Legal Deposit Libraries

# The Digital Library System

- 4 Nodes (Complete Copies)
  - British Library, St. Pancras
  - British Library, Boston Spa
  - National Library of Wales
  - National Library of Scotland

- Additional Access Points
  - Bodleian Library, Oxford
  - Cambridge University Library
  - Trinity College Library, Dublin

# The UK Web Domain

- Over 10 million .uk registered domain

- UK organisations also use non .uk domain names (eg .com or .org)
  - Exact scale unknown

- Non-print Legal Deposit applies to
  - the open (freely available) web: .uk
  - other UK-published (non .uk) websites, such as .com, .org…
    - made available to the public by a person or an organisation and the activities relating to the creation or the publication of the work take place within the United Kingdom

# Collecting strategy

**Domain Crawl**

**Events**

**Key sites**

**News**

Special collection

Special collection

Special collection

Special collection

Domain crawl:
- Broad crawl of UK domain
- Once or twice a year

Events & key sites and news:
- Events of UK interest
- High value, high impact sites
- National & regional news

Special Collection:
- Focused, thematic collections
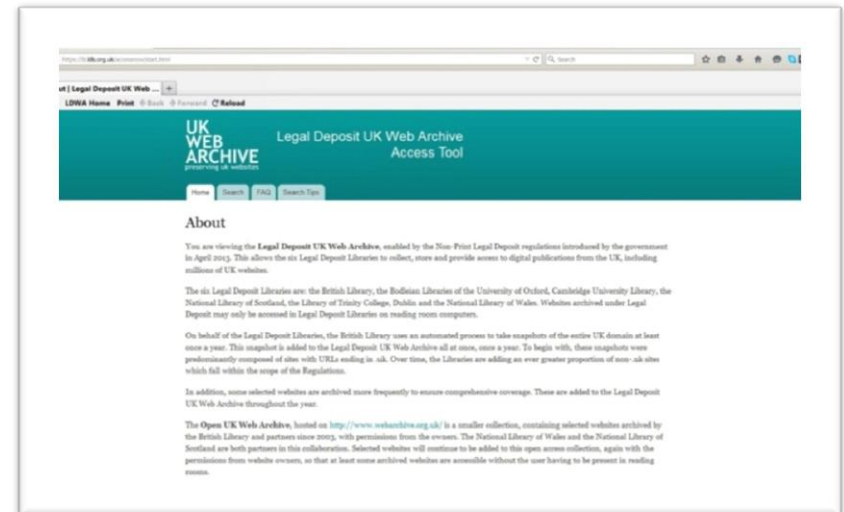- Support priority subjects

# Open UK Web Archive

- Based on websites owners' permissions

- 24/7 access

- ~68, 000 point-in-time snapshots of over 15,000 selected UK websites

- 21TB

- Continues to grow

- Going forward: showcase the Legal Deposit  content

# Legal Deposit UK Web Archive

- Annual UK Domain Crawl
    - 2013: 31TB, 1.6 billion URLs, 4 million hosts
    - 2014: 57TB including 2.5 non .uk domains

- Events, news sites, key sites
    - ~15TB

- Discoverable through British Library's online catalogue

- Full-text searchable in the reading rooms

# Same job, just bigger and more complex

"Hockx-Yu is now on the front line of the most ambitious expansion of the British Library's archiving capability in more than 300 years. At the stroke of midnight on April 5, 2013, legislation known as the Legal Deposit regulations came into force, charging the library with capturing the contents of the entire U.K. web domain—every site carrying the .uk suffix—preserving the material and making it publicly accessible."

Jacobson (2014)



TECH & SCIENCE

**Inside the Struggle to Preserve the World's Data**

BY **PHILIP JACOBSON** / JULY 2, 2014 12:21 PM EDT

## Newsweek

" the British Library's web archiving operation, [is] rated by experts in the field as among the best in the business. "

# Web archiving: the global picture

- Survey of web archiving initiatives (Daniel Gomes et al 2010)
    - 42 web archiving initiatives across 26 countries since 1996
    - 6.6PB archived web resources

- International Internet Preservation Consortium (IIPC)  -  from 12 founding members in 2003 to nearly 50 in 2015
    - Most members are national memory institutions

- How much  of the Web is Archived (Scott Ainsworth et al, 2012)

| Percentage archived | # of copies in public archive |
|---|---|
| 35% -90% | At least one |
| 17-49% | 2-5 |
| 1%-8% | 6-10 |
| 8%-63% | >10 |

- BL and CDNL (association of chief executives of national libraries) survey: 30 nations have legislations by 2012 enabling web harvesting

# Types of web archives

- Global: Internet Archive's Wayback Machine

- National: UK, Denmark, France, Finland and many other countries

- Sub-set of a national domain: eg UK Government Web Archive, Canadian Government Web Archive

- Topical collections: eg Human rights

- Collaborative: Archive–It service by Internet Archive

- Corporate: a single organisation's web estate, eg UK Parliament Web Archive

# Reflections / observations

- High concentration of use of technology and practice

- Modern browsers can render virtually any webpage. But the same pages often defeat our purpose-built crawlers

- National approach to archiving the global web

- Does size matter?

- Role of Legal Deposit

- Library/Archive discourse

- What is driving web archiving?

# Access and use of web archives

- Completeness versus openness
  - Many countries don't have Legal Deposit
  - Legal Deposit national collections have restricted access

- Not enough versus too much
  - Pre-selected or defined collections not relevant to all researchers; difficulty in finding relevant content in large scale web archive

- Arbitrary (national) boundaries often irrelevant to research question

- (Single) envisaged use case based on browsing of individual websites page by page

# An archived website

# Collections

# Subjects

# Legal Deposit UK Web Archive

- Full-text search, results can be refined by facets
    - Content type (html, pdf…)
    - Crawl year
    - (individual) Domain (bbc.co.uk; newsforscotland.com…)
    - Domain suffix (.co.uk; .org.uk; .ac.uk; .com)
    - Author (based on metadata extracted from web pages)
    - …

- Still has "collections" – selected and curated by curators, usually crawled within a fixed period of time, relating to an event

- Subjects information – added to selected sites only – included in the index

# Issues of "document-centric" approach

- Rise of "digital scholarship", taking advantage of the possibilities offered by technology

- Primacy of "text" as object of study no longer exists

- "Paratexts" play a crucial role in textual coherence of a website: header & footer drop-down menu, site map, breadcrumb, etc. (Brügger, 2010)

- "Distant reading" focus on units that are much smaller or much larger than the text: devices, themes, tropes—or genres and systems (Moretti, 2000)

# "Resource Not in Archive"

- Common error message but appears for different reasons

- Intended boundary
  - No permission for linked content
  - Not allow by robots.txt
  - Edge of a archive
  - Data limitation

- Technical limitations, e.g. dynamic content the crawler could not collect

- In UK Web Archive: avoid dead end in navigation
  - Search
  - Link to live web
  - Find archived copies elsewhere

# There is so much more…

- Statistical overview, scale and distribution of a web national domain

- Evolution, pattern of change over time, eg domain names

- Space: geo location, postcodes

- Type of content, eg file format, language

- Correlation, eg between certain term and historical event

- Structure, linked entities and networks

- Other web archives

# JISC UK Web Domain dataset (1996-2013)

- Copies of UK websites extracted from the Internet Archives collection
    - 57TB, 4 billion URLs

- Research agreement between JISC and IA, upholding IA's Terms of Use
    - Access via IA's Wayback Machine
    - Allows replication / extraction of secondary datasets: Format profile, Geoindex, Host Link Graph

- BL hosts the dataset on behalf of JISC

- Serves as research test bed and used by various research projects
    - JISC funded: Analytical Access to the Domain Dark Archive (AADDA); Big data for political science
    - AHRC funded: Big UK Domain Data for Arts and Humanities
    - ERC funded: Alexandria

# The UK Web in 2013



Largest UK Hosts (in GB)



Total Data Volume



Nmber of URLs

# Evolution of the UK web (2004 -2013)

# Postcode-based access

# HTML version analysis

# Image Format Analysis

# Using N-gram for scholarly research





- ▪ *Courtesy of Dr Peter Webster, Institute of Historical Research, University of London*

# Exploring links



Figure 5: Network diagram of hyperlinks between universities. Different colors indicate different university affiliations.

No affiliation
Million+ Group
Russell Group
University Alliance
Cathedrals Group
1994 Group

Government and governmental institutions
Blogs, social networks, entertainment
Corporations
Academic
NGOs, associations, groups

Courtesy of Peter Webster, Rainer Simon and Jules Mataly

# Visualising links (to and from bl.uk 1996)



Wednesday, 01 January, 1997 00:00:00

www.lib.ed.ac.uk
uk/ac/ed/lib/www/

**Interactive version**
**How it is done**

# Visualising links (to and from bl.uk 2010)



Friday, 01 January, 2010 00:00:00

**Interactive version**
**How it is done**

# How was the UK web linked in 1996?



- By Rainer Simon using UK Host-Level Link Graph (1996-2010) dataset.

-  Based on the 1996 portion: 58,842 hosts (nodes); 184,433 host-to-host links (edges)

-  UK web as part of the global web

- Scalability issues with large dataset over time

# Mementos service (http://www.dcc.ac.uk)

# Collaboration with researchers

- Building collections
  - Researchers' involvement in scoping collections, selecting and describing websites
  - Creation of specific, (narrow) topical collections

- Formulating research question
  - Brain-storm sessions, workshops, discussion, surveys etc.
  - Lack of awareness & baseline knowledge
  - Challenging: you don't know what you don't know
  - Analytical access: web archive as dataset

# Co-development

- Institute of Historical Research project: Analytical Access to the Domain Dark Archive (AADDA)
    - Experimental user interface

- Big UK Domain Data for Arts and Humanities (2014-2015)
    - Institution of Historical Research, Oxford Internet Institute, British Library and Aarhus University
    - Develop theoretical and methodological framework for the study of web archives
    - Build on ADDAA: researchers and the BL co-produce access tools
    - Study of the history of UK web space from 1996 to 2013
    - Research use cases

- Changing the way we archive the web

# Web archiving researcher bursaries

## Welcome to our 11 bursary holders

Posted on **May 27, 2014** by **Jane Winters**

One of the main aims of the project is to involve arts and humanities researchers in the development of tools for analysing web archives, thereby ensuring that those tools meet real rather than perceived researcher needs. We recently ran an open competition inviting researchers to submit proposals across a range of disciplines which focus on the archived web, and have selected 11 from a tremendously strong and varied set of applications. The topics that will be studied over the next eight months are:

- Rowan Aust – Tracing notions of heritage
- Rona Cran – Beat literature in the contemporary imagination
- Richard Deswarte – Revealing British Euroscepticism in the UK web domain and archive
- Chris Fryer – The UK Parliament Web Archive
- Saskia Huc-Hepher – An ethnosemiotic study of London French habitus as displayed in blogs
- Alison Kay – Capture, commemoration and the citizen-historian: Digital Shoebox archives relating to P.O.W.s in the Second World War
- Gareth Millward – Digital barriers and the accessible web: disabled people, information and the internet
- Marta Musso – A history of the online presence of UK companies
- Harry Raffal – The Ministry of Defence's online development and strategy for recruitment between 1996 and 2013
- Lorna Richardson – Public archaeology: a digital perspective
- Helen Taylor – Do online networks exist for the poetry community?

INSTITUTE OF HISTORICAL RESEARCH — University of London School of Advanced Study

BRITISH LIBRARY

oiioiioii oxford internet institute university of oxford

AARHUS UNIVERSITY

Arts & Humanities Research Council

# Shine



- Query building
- Corpus formation and handling
- Annotation and curation
- In-corpus analysis
- Whole-dataset analysis

# Scholarly use of web archives



http://netpreserve.org/web-archiving/videos

# Requirements for web archives

| Scholarly requirements | Requirements for web archives |
|---|---|
| Availability | No access restriction, available online |
| Text, paratext and context | - Access to text and any contextual information at various granularity.<br>- Access to secondary datasets |
| Persistence and citability | - Persistent identifiers<br>- Standards of citing archived websites<br>- Integration with bibliographical management tools (eg Zotero) |
| Formation of research corpus | - Archiving of research corpora on demand<br>- Means to mix, match and reassemble corpora |
| Quality | - Represents as much as possible the live website in completeness, intellectual content, behaviour and look and feel<br>- Information on what is "missing" to help interpret incomplete sources |
| Close and distant reading | Multiple access methods including tools for data analytics and visualisations |
| Boundary & format-independent | - Interlinked web archives<br>- integration with other digital and printed sources |

# IDCC 2025

- The UK Web Archive
    - Single search for all collections
    - Base-line knowledge self-explanatory, text and much contextual information
    - Link to many more web archives
    - Inactive websites highlighted
    - Open datasets / content

- "Macroscope" of the UK web history
    - "a single data point, .. both visualised at scale in the context of a billion other data points, and drilled down to its smallest compass"

- Flexibility to reconstruct context, big or small

- Seamlessly linked web archives

- Crawlers as capable as browsers

- Independent use of web archives

**CAREERS**

most precious items

Careers home    People and projects    Benefits and development    Working at the Library    **Vacancies and how to apply**

bl.uk > Careers home > Vacancies and how to apply > Job Details

## Job Details

Login

Search and apply

Job alerts

**Software Developer - Web Archiving**

**Ref**            OPS00109
**Location**       Boston Spa, Yorkshire
**Position Type**  Permanent
**Specialism**

£28,450 - £33,041 plus benefits

Permanent

The role

This is a great opportunity to join the Library's IT department working as part of the Remote Services Development team you will develop, test, and help deploy key web archiving applications, meeting current operational needs and contributing to the Library's ongoing capability of archiving the evolving web.

Essential Criteria

- Strong web application development experience in Java, including use of Eclipse, JUnit, Maven and Git.

- Ability to plan and execute unit, system, integration and regression testing.

- Good working knowledge of SQL and database design or legacy equivalent.

- Strong understanding of the Web stack (HTTP, HTML, JavaScript, CSS).

- Experience using server software including Apache Tomcat and Httpd.

- Experience developing on Linux.

- Knowledge and experience of shell scripts and regular expressions.

- Experience using the Play Framework (or similar frameworks) to develop Java web applications.

- Experience with Apache Solr or similar systems (client side usage).

- Be able to communicate clearly and effectively in various ways, including

**People and projects**

Our amazing projects and the people who made them happen.

▶ Find out more

**Interactive fact-finder**

wing  collections  silence
cient exhibitions news scan
serials  successes  pod
nic  science  maps  liter

What have trains, podcasts and thrash metal got to do with the British Library

▶ Find out more

## Java Developer for Web Archiving

**Closing date Midnight 15th February**

**Please help spread the word or Apply!**

**Contact: Helen.hockx-yu@bl.uk**
**Twitter: @ukwebarchive**

# References

- Ainsworth, Scott G., AlSum, Ahmed, SalahEldeen, Hany, Weigle, Michele C., Nelson, Michael L. (2012) 'How much of the web is archived?', http://arxiv.org/abs/1212.6177 .

- Brügger, Niels (2010) *Website Analysis: Elements of a conceptual architecture*. Aarhus: The Centre for Internet Research.

- Gomes, D., Miranda, J., Costa, M., A survey on web archiving initiatives. In *Proc.* TPDL'11, Springer-Verlag Berlin, Heidelberg (2011), 408-420, http://dl.acm.org/citation.cfm?id=2042590 .

- Jacobson, P (2014), 'Inside the Struggle to Preserve the World's Data'. http://www.newsweek.com/2014/07/11/inside-struggle-preserve-worlds-data-257020.html.

- Moretti, Franco (2000) 'Conjectures on world literature'. *New Left Review*, 1, http://newleftreview.org/II/1/franco-moretti-conjectures-on-world-literature.