

# Co-developing access to the UK Web Archive

Helen Hockx-Yu  
Head of Web Archiving, British  
Library

# Ten years of archiving the UK Web Archive

- Started web archiving in 2004, non-print Legal Deposit since April 2013
- Three collections: over six billion resources and over 100TB compressed data
- Focus not just on content collection
- Proactive development of access and use, through close engagement with researchers
  - User survey
  - Content selection and curation
  - Brain-storming sessions and workshops to formulate research questions
  - Research projects

# JISC UK Web Domain Dataset 1996-2013

- Funded by JISC to create a research collection of historical UK websites
- Collaboration between the Internet Archive, JISC and the British Library
- Copy of subset of the Internet Archive's web collection that relates to the UK
- c.300 million resources, 60TB in total
- No local access – possible through the Internet Archive
- Can be used to generate secondary datasets

# Co-design at every stage

- Research use case articulated
- Generic user requirements abstracted
- Requirements refined following feedback
- Iterative development cycles: Develop -> user testing -> feedback -> develop ...

# Use cases (generalised)

- Full-text/facet search -> individual resource
- Full-text/facet search -> analysis/visualisation
- Search -> corpus creation -> annotation/curation
- Corpus creation -> full-text search -> individual resource
- Corpus -> search -> analysis/visualisation
- [Derived datasets -> take-away]
- [Direct access to WARC/CDX -> take-away]

# High-level requirements

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

# Prototype: Shine

- Full-text search, with proximity options, and to exclude specified text strings
- Apply and remove multiple facet filters to result sets
  - Content type, public suffix, domain, crawl year
  - Also available: postcode, links to public suffix, language, links domains
- Exclude single resources, or whole hosts from result sets
- Save a query
- Export basic query results, as CSV or similar
- Available at: <http://webarchive.org.uk/shine>

# Advanced Search

← → ↻ [www.webarchive.org.uk/shine/search/advanced](http://www.webarchive.org.uk/shine/search/advanced)

**Search Terms:**

**Proximity:**    If you wish to search for word that appear together, try a proximity search. For example, the following finds resources where the words 'coffee' and 'java' occur within 25 words of each other.

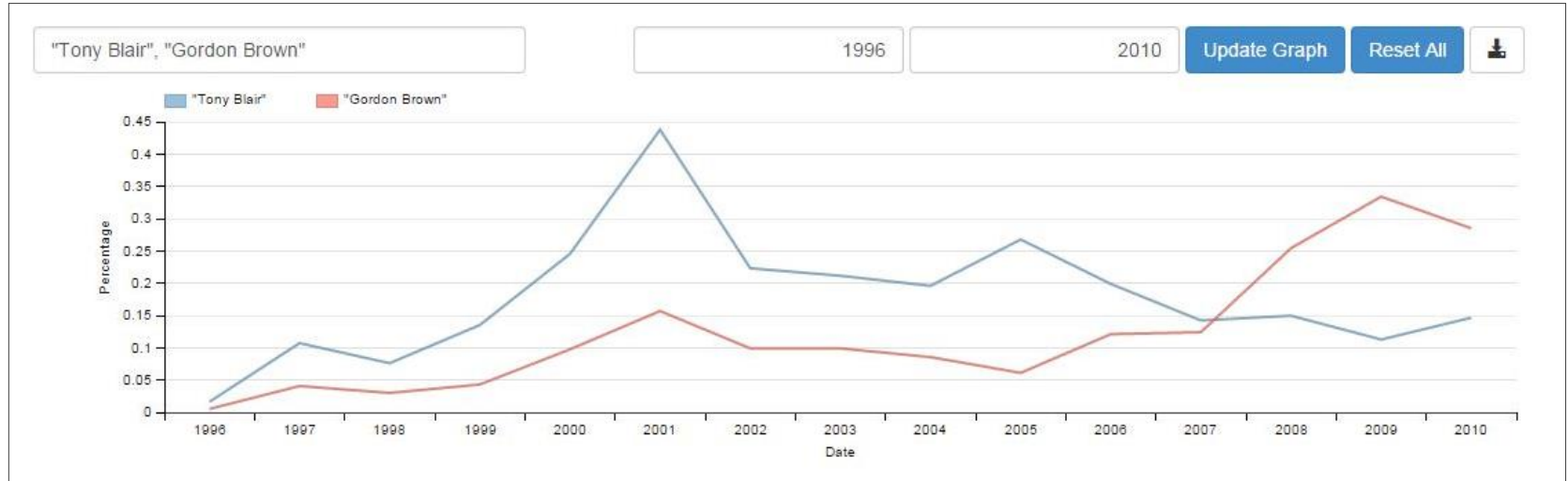
**None of these words:**

**Within Resources**

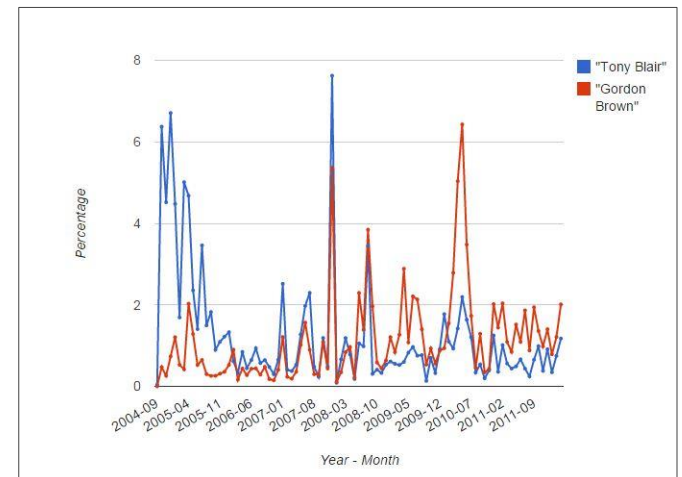
<b>Date Range:</b>	<input type="text" value="dd/mm/yyyy"/> <input type="text" value="dd/mm/yyyy"/>	Restrict by date (Format: dd/mm/yyyy)
<b>URL:</b>	<input type="text" value="URL"/>	URL
<b>Host, Domain or Public Suffix:</b>	<input type="text" value="Host, Domain or Public Suffix"/>	Match the values in the 'host', 'domain' or 'public_suffix' fields
<b>File Format:</b>	<input type="text" value="File Format"/>	File format
<b>Website Title:</b>	<input type="text" value="Website Title"/>	Website title
<b>Page Title:</b>	<input type="text" value="Page Title"/>	Page title
<b>Author:</b>	<input type="text" value="Author"/>	Author



# Ngram



- Same search terms, different datasets
- Broadly similar trends
- Interesting to examine turning point
- Not useful without understanding of scope
- Visualisation not the end point



# Pages mentioning “Gordon Brown” (2007)

Found 100 samples matching ' "Gordon Brown"' from 2007.

	Matching Text	Link
levels before the end of 2008. So earlier this year, UK Chancellor	Gordon Brown changed the tax regime for	bbc.co.uk
...celebrities, money, dreams, celebrity, satire, drawing, tax, bud ,	gordon brown , lampoonEurovision Song Contest Disaster!Á Á Á Á Á Á Á Á Á Á ...	humorousarts.co.uk
...Press: 1997-55HelpContact usAccess keysSite mapA-ZSearch	Gordon Brown announced today.Mr Brown said:"The Informal ECOFIN meetin...	hm-treasury.gov.uk
...calls to Downing Street going unanswered.advertisementHe sai	Gordon Brown , in July at the British Irish Council. Picture: PAGordon Brown w... Gordon Brown seems to be stuck in the	theherald.co.uk
all.Should those figures be taken as good news or bad news for	Gordon Brown is enough of a change to meet that public demand. Clearly the... Gordon Brown ? Only somewhere between 30-40% of	ukpollingreport.co.uk
calling on	Gordon Brown to build an extra 1,170 council or housing association homes in...	melfonmowbraytoday.co.uk
...BennettThursday April 19, 2007The GuardianWhoever wrote C	Gordon Brown , he, she or they are to be	guardian.co.uk
...doing his job as Deputy PM and Minister for DETRSatisfied 32	Gordon Brown	icmresearch.co.uk
...moreDate:27/06/2007 Author:NewsBrown to be held to acc .....	Gordon Brown prepares to take	theecologist.co.uk
...putative successor	Gordon Brown , on moderation. He deplored the naive language of counter-te...	guardian.co.uk
...according to key staff in the immigration system.Members of t ,	Gordon Brown is interested in grabbing headlines not solving problems.' Davis... Gordon Brown this summer as a major new counter-terrorism initiative, claim t...	guardian.co.uk
...topDEPUTY ROLEHilary Benn: Wants to be deputy party leade	Gordon Brown who is deputy	bbc.co.uk
Harman. Argues her good relationship with	Gordon Brown - she was once shadow Chief Secretary to the	
responsible public servant should do thisAt the same time	Gordon Brown gave this money with one hand, he was	worcesterstandard.co.uk
challenge	Gordon Brown andDavid Cameron to back up their rhetoric by endorsing our p...	wellsibdems.org.uk
...huge levels of fraudand error in the Tax Credits system is disgr	Gordon Brown haswalked into	
...hints at change of heart over supercasinos - Foster Commenti	Gordon Brown &rsquo;s hint during Prime	
...amount of money wasted in the first three years is on course to	Gordon Brown walked into	
- Huhne Responding to	Gordon Brown &rsquo;s announcement of a Â£14 million supportpackage for ...	
...Ministry of Defence   Defence News   Brown: UK will continue h	Gordon Brown has reiterated the UK's position on Iraq following an announce...	www.mod.uk
"Prime Minister	Gordon Brown "I intend to make a more detailed statement when Parliament r...	
CentreBasra handover is imminentPrime Minister	Gordon Brown , making a surprise visit to troops in Iraq, has	
...to persuade the chancellor	Gordon Brown , and the prime minister to make childcare a theme for a	guardian.co.uk
...doorstep. It's time for a change of Government, as that pasty-'	Gordon Brown is just as bad as his	thesun.co.uk
comes from no less a source than	Gordon Brown . While the Penguinás were away in Germany over Christmas	politicalpenguin.org.uk
...speak fluently. However, a brief 15 minutes of catching CNN br	Gordon Brown in the	
...burmaÁ show   burmaÁ websiteGordon Brown's Lord Mayor...	Gordon Brown will address	britishblogs.co.uk
...sparked severe flood warnings and evacuations on England's e	Gordon Brown has held	
guttled when	Gordon Brown invites Maggie Thatcher to No 10. How many mosques would gi...	guardian.co.uk
speakers. Just like a lot of Labour people feel gutted when	Gordon Brown invites Maggie Thatcher to No 10	
planning gain supplement2005 issue 48Chancellor	Gordon Brown set to announce further consultation in pre	building.co.uk
responsible public servant should do thisAt the same time	Gordon Brown gave this money with one hand, he was	worcesterstandard.co.uk
...Poll ResultsNewsÁ  Á Site MapÁ  Á About UsÁ - Select Country	Gordon Brown 4.67%Lifestyle Change 1.17%Total voters: 6/info/poll_results.a...	apartmentsjavea.co.uk
...ANALYSISBrown's agendaWhat sort of British prime minister ,	Gordon Brown make?Departures loungeOutgoing British	bbc.co.uk
know the score here or that, come the 27th June, it will be	Gordon Brown weeks after the Conservative leader passed on the idea in priv... Gordon Brown they'll be dealing with to	ministryoftruth.org.uk

# Trends analysis

www.webarchive.org.uk/shine/graph?query="google"&year\_start=1996&year\_end=2010&action=update

UK Web Archive Search Trends Login

**Warning!** This is a research prototype for a web archive search service, and may be taken down at any time.

"google" 1996 2010 Update Graph Reset All

Date	Percentage
1996	0.0
1997	0.0
1998	0.0
1999	0.0
2000	0.0
2001	0.1
2002	0.3
2003	0.6
2004	1.2
2005	1.5
2006	1.3
2007	1.4
2008	2.5
2009	3.0
2010	3.8

Found 76 samples matching "google" from 1997.

Matching Text	Link
google { print " Inherit::mumble();}sub:: google {my \$self = shift;\$self->Foo::Inherit::	soton.ac.uk
here's the goo\n" }package Bar; @ISA = qw( Buz );sub" google { print " shift;\$self->SUPER::mumble();}sub = google {my \$self = shift;\$self->SUPER::	rdg.ac.uk

# Access to data supporting trends

Internet Archive Wayback Machine  
4 captures  
9 May 97 - 18 Jan 98

http://www.hawksand.co.uk/gaba/membroo2.html

MAY JUL OCT  
1996 24 1997 1998

Name	Gary Jerram		
GABA Mail-Lister	No		To come...
Membership Number	00457		
Comments			

Name	Shaun Ewing		
GABA Mail-Lister	No		To come...
Membership Number	00481		
Comments			
<i>I am 13 years old and no I am not one of those teenagers who just like to google at Chliian Anderson pictures all day.</i>			

Name	Rebecca Appleton		
GABA Mail-Lister	No		No Picture
Membership Number	00482		
Comments			

Name	Ev		
GABA Mail-Lister	No		
Membership Number			
Comments			

# Next steps

- Inclusion of the full JISC dataset – seamless interface to all 3 components of UK Web Archive
- Better support for corpus creation (eg combination of existing corpus)
- Annotation and sharing of corpus
- (standard) analysis and visualisation of corpus
- Faceted search within user-define corpus
- (semantic) clustering of search results

# Lessons learnt

- A learning process for both
- Not a choice between “big data” or “small data”
- “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”
- Context and paratext just as important
- User expectation / assumption
- Maximum transparency