

Big Data and History: Seeing the Past through a Macroscope

Danish Society for Research on Contemporary History,
Copenhagen

Ian Milligan
Assistant Professor
[@ianmilligan1](https://twitter.com/ianmilligan1)



UNIVERSITY OF WATERLOO
FACULTY OF ARTS
Department of History

Hi everybody and thanks so much for having me here today - to Professor Sorensen for the gracious invitation, and to the Danish Society for Research on Contemporary History for making my attendance here possible.

What I want to do here today is talk about why I think Big Data matters, and walk you through several major points:

- (1) That Big Data is already reshaping our profession - and it is something that will continue to accelerate thanks to the advent of web archives, the area I work on with many others including Professor Niels Brugger;
- (2) That the shift from historical scarcity to abundance is a fundamental shift;
- (3) And that we need new tools, the Macroscope, to access this information. I'll give examples from my work, with links and notes towards different tools and walkthroughs that I've used to access it.

And then we will have time for questions!

Historians are largely unprepared to engage with the quantity of digital sources that will fundamentally transform their trade.

So let's begin with my thesis:

Because, as I've explained in other venues as well as here, I believe that historians are unprepared to engage with the quantity of digital sources that will fundamentally transform their trade. The advent of the World Wide Web, both with born-digital sources that can shed so much light on human culture and activity after the advent of web archived material in approximately 1996, as well as digitized traces of the past, presents a challenge to us as professionals.

**... we need to think
about data ...**

It means that we need to start thinking about data. How to preserve it, how to access it, how to interpret it, and how to facilitate its reuse.

Today's Talk

- 1. **Prologue:** Big Data is everywhere
- 2. **The Web Age:** Will accelerate this process
- 3. **What can we do with big data?**

So what I want to do today, in the time we have together, is to explore a few different dimensions of historical research and how data is transforming it. Let's start with the big point about why thinking about data matters for us, go into the particulars of the Big Data of the Web Age - a plug for my own research - and then conclude around some of the problems that I think we're going to have to tackle.

A Prologue: Big Data is Everywhere

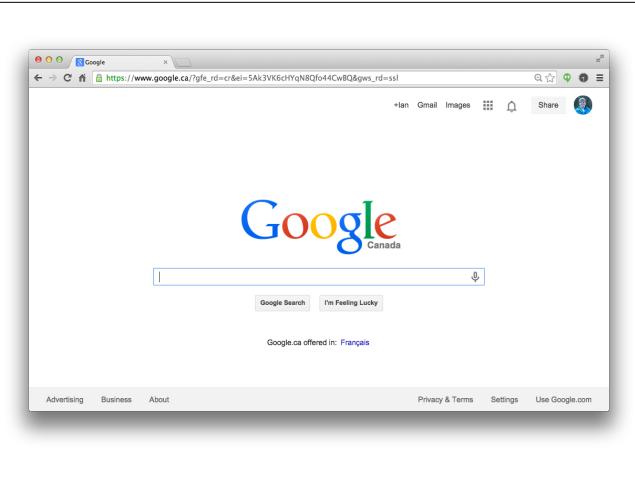
Let's start with a prequel to the Web Age.

What do we mean by Big Data?

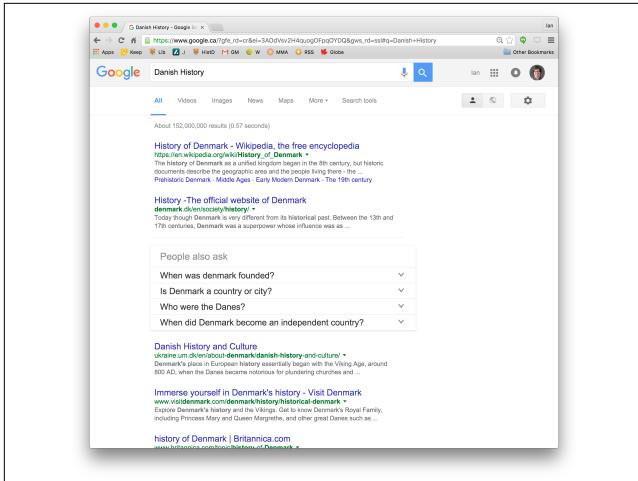
- Computational definition: the 5 Vs (Volume, Velocity, Variety, Veracity, and Value)
- “For us, as humanists, big is in the eye of the beholder. If it’s more data that you could conceivably read yourself in a reasonable amount of time, or that requires computational intervention to make new sense of it, it’s big enough!” (Shawn Graham, Ian Milligan, Scott Weingart, *Exploring Big Historical Data*)

two kinds of definitions

**Why is it
everywhere?**

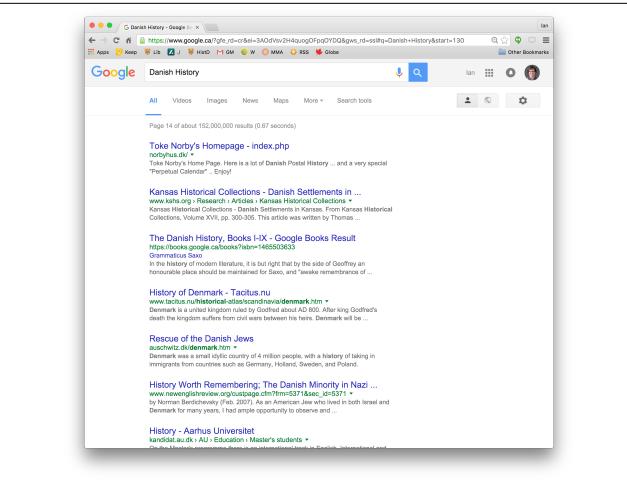


The first point I want to make is that no matter the kind of research historians are doing, chances are they're using data - and at the mercy of those who interpret it for them. Every time you run a search result...



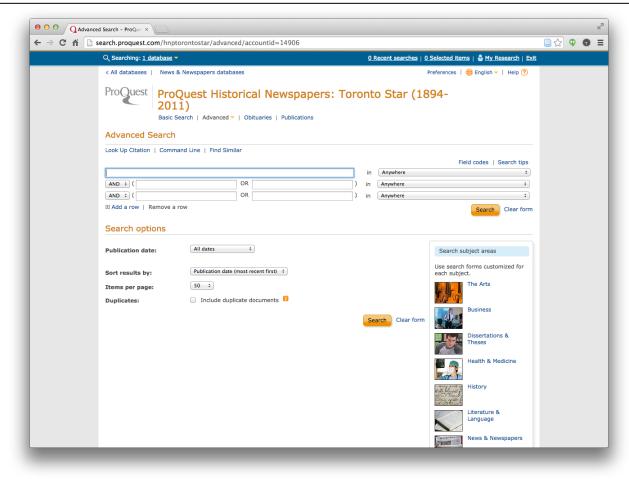
There's a reason that the wikipedia article is at the top of this search on Google Canada for Danish History, that the 100th result is a home page of a Danish researcher, and that beyond that we have resources that we never know about - who here actually visits the 1000th most popular site on a Google search?

Most humanists probably don't think about it, but Google's skewing their research. ITHAKA S+R did a survey, and at least in North America, this is our new **research assistant**.



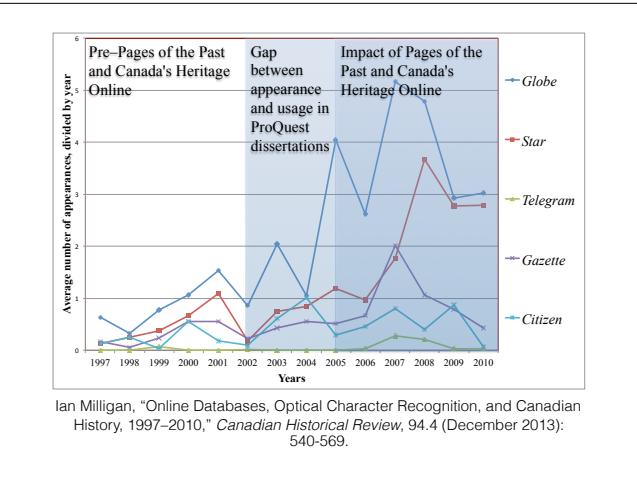


But in the Web age, historians are accessing most of their sources through mediated search portals, which necessitates computational thinking. For decades, historians have used microfilms, and now in literally the last decade, we've been using Web-based primary sources.



On a whim, I was wondering what Web-based primary sources might mean - mainly because I was reading somebody's dissertation, they were talking about events in Oshawa, Ontario, and they were using the Toronto Star - not the Oshawa Daily Times. Or, to be honest, the first time I'd started researching a strike in Hamilton, I started using the Globe and Mail until my supervisor pointed out that the Hamilton Spectator would be better.

Well, luckily I wasn't alone.



As things are digitized, people use them more. And the things that aren't digitized, are used less. That's fine, except:

Digitized sources don't reflect the geographic diversity of Canada or the GLOBAL NORTH - i.e.

- * Hamilton not digitized, but Toronto is. That makes decisions made in 2001-era Canada impacting our histories of 1890s Canada.
- * The OCR was generated in 2001 (!) based on digitized microfilms
- * The original OCR is inaccessible to the creators
- * The search engine is dreadful
- * Our citation methods don't differentiate generally between whether we got it in a collection file, a microfilm, an original copy, or an online database - and all these things are different

In short, **digital sources are mediated in different ways than the original sources**, and arguably for a decade Canadian historians were not recognizing the role that algorithms and data play in the mediation of our research.

Digital resources in the Web age are useful, but we're a profession not terribly equipped to deal with this.

... this is our long-term *track record* w/ digital resources ...

I mention this by way of introduction because it sort of sets out the track record that historians have when we're dealing with this material. Digital resources came, and we uncritically used them for years, not just in Canada but across the historical profession. I've made this point in a Canadian Historical Review article and it's not immodest of me to say that it's appearing on syllabi and being discussed in France, Australia, and the United States because it's a relatively novel point for historians.

Our history with digital
sources is the **unreflective**
use of technology.

... we've become, in some ways, a discipline defined by the keyword ...

We've become a keyword-based discipline almost overnight, from finding aids to digitized newspapers.

A process that is only
now **beginning to**
accelerate.

And this becomes incredibly pressing, because soon historians are going to have a TON of data to draw on. This brings me into the main topic of my presentation: the web age. If we can't handle newspapers critically once they're digitized and presented as DATA, how can we handle the 1990s?

**First - more data than
ever before being
preserved;**

**Second - it'll be
saved/delivered to us
in **very different ways****



Since the 1890s, this has been the bread-and-butter of the historical profession. An archival box. Good, old fashioned, 'analog' to use an imprecise word, technology. Historians go to archives, we flip these open, we see a bunch of documents, and from these traces of the past we attempt to reconstruct narratives that we call history.

While historians have challenged some of the preconceptions of this since the 1960s, the old gold standard was government documents. Even since the 1960s, we largely - not wholly, as oral historians and others would testify - use archives deposited with organizations: Library and Archives Canada, the University of Waterloo's Special Collections, etc.

This all means that historians are working with scant traces of the past. 99.9% of things that happen are never recorded (something which philosophically really isn't changing with new sources). Events are transitory, they happen, and unless it ends up in an archive box, or a newspaper, they are generally lost.



The problem is, with these, in a word, that of scarcity, as the late great American historian Roy Rosenzweig noted. So if that has been the landscape for well over a hundred years of professional development -- I want to hold to you today that the arrival of something new portends to change and alter the work that we do. And it is this:



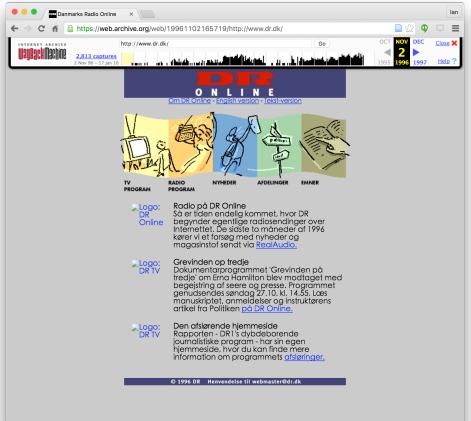
WebARChive (WARC) File

The WARC file, or WebARChive file.

This, I argue, is what historians need to wrap their heads around. Because it represents the biggest challenge to our profession since that old man in the last slide first tried to professionalize our discipline over a hundred years ago.

So what is a WARC file? It's an International Standard, ISO 28500:2009 for those keeping track, and has become the international standard to preserve content and information "from mainstream Internet application layer protocols, such as the Hypertext Transfer Protocol (HTTP), Domain Name System (DNS), and File Transfer Protocol (FTP)" with attached metadata.

What does this really mean? In a nutshell, I can take the entirety of the University of Copenhagen's website today, and pop it all into a WARC file. It preserves the website as it is today, and I can either access it through the WaybackMachine to experience it, or I can use a suite of text-based tools to extract meaningful information quickly.



The screenshot shows the Occupy Wall Street NYC website. At the top, there's a navigation bar with links for News, Livestream, #OccupyWallSt, Forum, Chat, User Map, NYCQA, and About. Below the navigation is the Occupy Wall Street logo with the tagline "The revolution continues worldwide!". A main headline reads "Farmers Join Occupy Wall Street, Calling for Food Justice". It includes a sub-headline: "As Wall Street's corporate influence on the economy has grown, the corporate ownership of our food system has hurt the health and livelihoods of some of our most vulnerable communities. This Sunday, December 4th food justice activists and occupiers will be traveling from far and wide to call for change." Below this is a section titled "MARCH!" with a list of demands. Another headline below is "NATIONAL DAY OF ACTION DEC. 6, 2011". At the bottom, there's a note about a solidarity action in Brooklyn and a link to the NYCQA committee meeting times.

So what does this mean?

Imagine a future historian, taking on a central question of social and cultural life in the middle of the first decade of the twenty-first century, such as how Canadians understood Idle No More or the Occupy movement through social media? What would her archive look like?



I like to imagine it as boxes stretching out into the distance, tapering off, without any immediate reference points to bring it into perspective.

Many of these digital archives will be without human-generated finding aids, would have perhaps no dividers or headings within “files,” and certainly no archivist with a comprehensive grasp of the collection contents.



Now this all means that historians need to get involved a bit earlier, of course. Of all those Occupy Wall Street sites - created in the heat of a movement, etc., they didn't have data management plans. And so two years later, only 41% of those sites are still active today.

I like to throw this in because it shows that we need to actively maintain our digital sources, and historians need to begin thinking about retaining the data that we generate today.

If I wrote a book in 1991 and put it on a shelf in my basement, and came back in 2014, chances are I can read it. I'm damn sure that if I wrote a digital object, it would be unusable. We need to be active.



But I digress. Some of it will be retained, and it's going to be retained on a massive scale thanks to the efforts of institutions like the Internet Archive and some legal deposit institutions in Europe. Remember scarcity? **We're in the era of historical abundance.**



Take Twitter, for example. During the #IdleNoMore protests, there were an stoning 55,334 tweets on 11 January 2013. Each tweet can be up to 140 characters. Through a complicated bit of math that I whipped together, I argue that that's over 1,800 pages if we take 300 words per page. That's a MASSIVE book, and you've got one for every day of the big protest. You can't read that yourselves, you're going to have to learn how to program. And that, my friends today, I hold to you is why historians need to be leading the charge up the digital humanities hill. But we're not.



Early theories of information saw the Library of Congress as the pinnacle, the largest quantity of information conceivable by information theorist Claude Shannon in 1949. It is, at least by some measures - the British Library disagrees - the largest repository of traditional print information in the world. Simply walking along its 838 miles, or 1,349 kilometres of shelving would take weeks – without even stopping to open up a single book.

The Library of Congress is no longer the pinnacle of information storage.



For that, we have the Internet Archive.

Now, comparing miles of shelves to the Internet Archive's 11,000 hard drives is an “apples versus oranges” issue, but it can be done as a rough thought experiment. Trying to put a firm data figure on an analog collection is difficult: a widely distributed figure is that the Library of Congress’s print collection amounts to 10 terabytes. That is too low, and a more accurate figure is somewhere in the ballpark of 200 to 250 terabytes (if one digitized each book at 300 DPI, leading to a rough figure of eighty megabytes per book). As a petabyte is a thousand terabytes, if we take the latter figure we arrive at a 1:40 ratio.

1 (Library of Congress):
40 (Internet Archive)

The Internet Archive is massively bigger. The Wayback Machine continues to grow, and grow, and grow, and historians are now confronted with historical sources on an entirely new order of magnitude.

“.... [n]ow expectations have inverted. Everything may be recorded and preserved, at least potentially.”

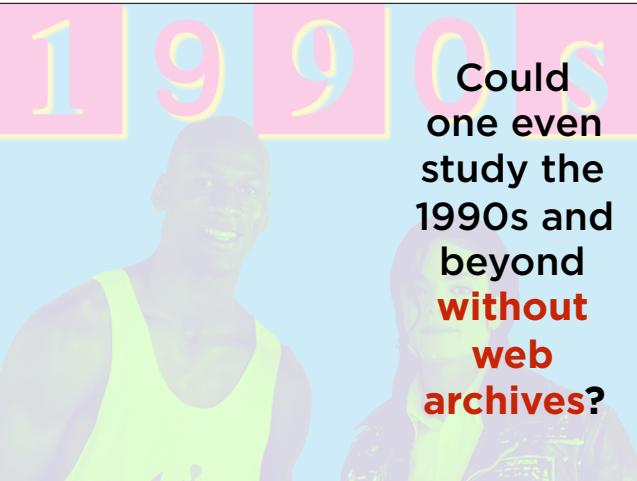
- James Gleick

These sources present both a boon and a challenge to historians. If the norm until the digital era was to have human information vanish, “[n]ow expectations have inverted. Everything may be recorded and preserved, at least potentially.” Useful historical information is being preserved at a rate that is accelerating with every passing day.

I think James is perhaps overstating the point here. Remember the previous process I outlined to you before - event happens, some traces of the past are left behind, and historians take those little traces to write their histories.

On a philosophical level, this is still true.

But, more traces than ever before are being left behind.



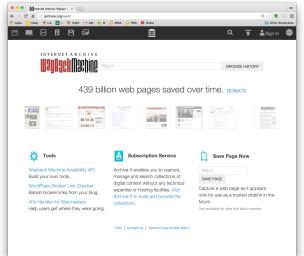
Could one even study the 1990s and beyond without web archives?

Let's not fool ourselves, this stuff matters. Things that would never have been preserved before are, and there are unparalleled opportunities for social historians to capture the broad context of the time that they are studied.

And, most importantly, I want to underscore that I don't think you can do justice to the 1990s and beyond if you do not consider the World Wide Web. This holds true for all branches of our historical profession. Political historians cannot do justice to elections without understanding the tweets, the blogs, the websites that surround not only elections, but the everyday process of making policy, understanding public sentiment, and reaching out to the electorate on a new level. Military historians will have the voices of rank-and-file soldiers, playing out on discussion boards and other parts of the Web. Cultural and social historians, the source base is even more evident: the impact of online and offline culture playing out, and the voices of everyday people that would not be lost.

Yes, the Web is not a perfect democracy: there are many lower-income people who still do not have access to the Web, and there are age and racial cleavages as well. We cannot forget those. But we are still expanding our percentage of 'traces of the past' that are preserved, and we cannot forget that either.

Nightmare Scenario



This won't be enough!

Because if we don't, I want to outline a nightmare scenario.

[make sure to note the URL only restriction of the Wayback Machine]



... but what will our
search engines look like?

Just like digitized newspapers.

My fear is that we'll query large databases, of literally hundreds of millions of documents, look at the first 100 results or even 1000, and we'll just be reaffirming what we already think.

Nightmare Scenario

- Historians rely uncritically on **date-ordered keyword search results**, putting them at mercy of search algorithms they do not understand (similar to digitized newspapers);
- Historians are completely left out of post-1996 research, letting everybody else do the work (a la Culturomics project/*Nature* magazine article);
- Our profession gets left behind...

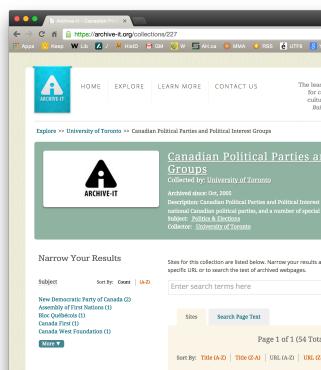
[and emphasize the comparison with historians and newspapers! i.e. we have these black boxes]

We already see some disturbing trends in this area, as with the Culturomics project.

**What can we do to
access this
information?**

Building Portals

- Democratizing access so that historians can use them.
- Building **transparent indexes**.
- But they have to be useful and tested...



Pivotal Changes in Canadian Politics, 2005-2015

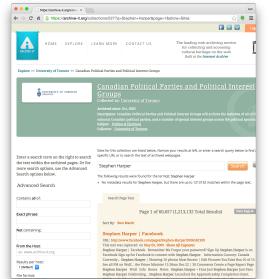
- Militarization of Canadian society?
- Change from 'natural governing party' of Liberals to Conservatives
- Major policy changes on foreign policy, environment, etc.
- How to measure?

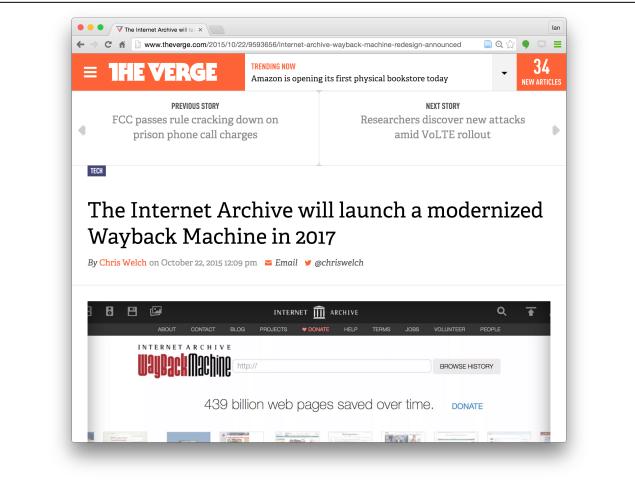


collections SHOULD have been used a lot more

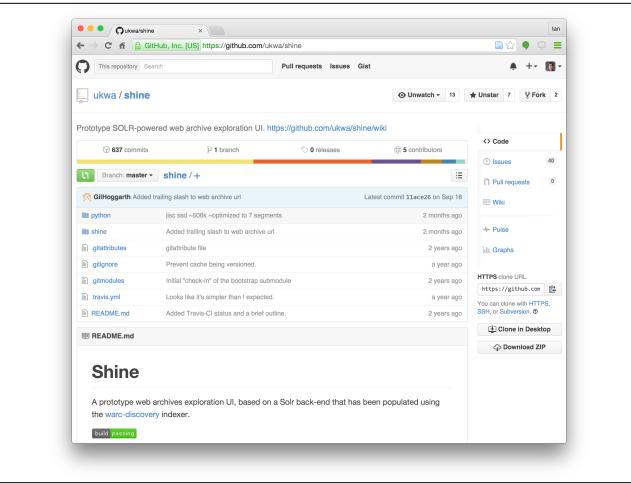
Current Interface

- Very limited - simple search engine, some advanced options; no facets
- Great collections.. but nobody uses them!





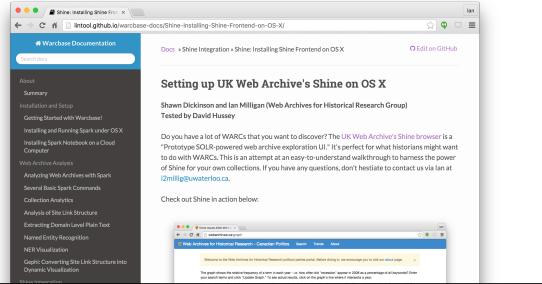
Some of this is in the area of providing sophisticated researcher access to web archives, and let people do what they need to do without having to be programmers. The upcoming rebuild of the Internet Archive's Wayback Machine is extremely exciting, and I've already set my countdown timer to 2017. I'm a bit afraid that it will deprive me of the straw man argument that I like to use around its URL-centric shortcomings, but I think the field will benefit.



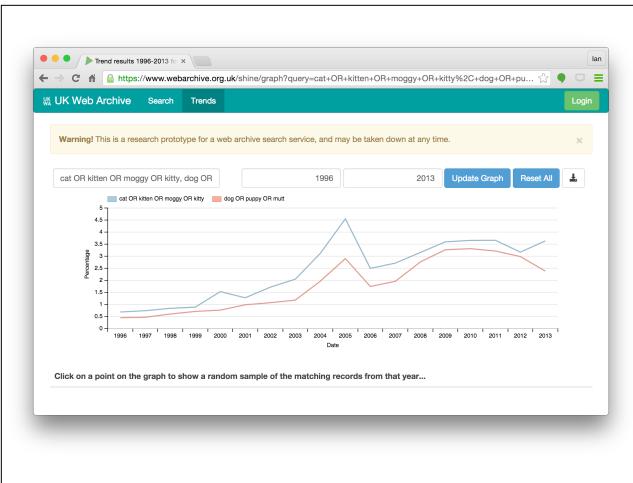
It's been things like the United Kingdom Web Archive's SHINE interface, though, that I think herald a lot of promise in the realm of DEMOCRATIZING ACCESS.

<describe shine>

Walkthroughs at:
**[http://lintool.github.io/
warcbase-docs/Shine-Installing-
Shine-Frontend-on-OS-X/](http://lintool.github.io/warcbase-docs/Shine-Installing-Shine-Frontend-on-OS-X/)**



Jason Webber's cat example

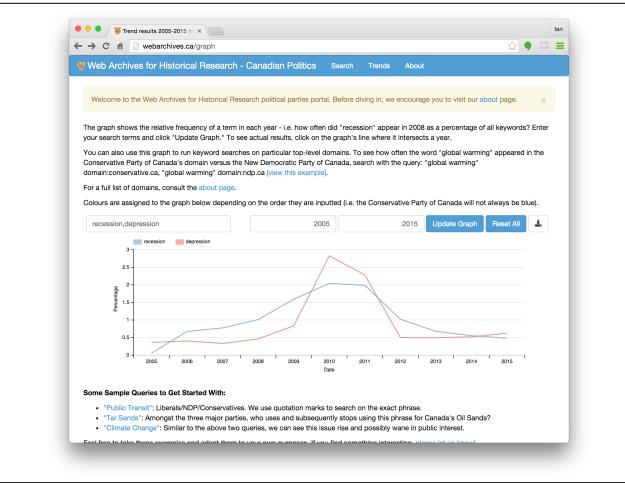


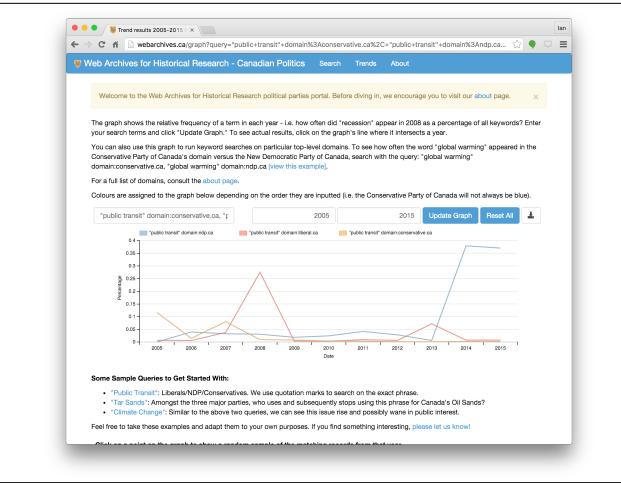


With Nick Ruest (York University), Nich Worby (University of Toronto), Jimmy Lin (University of Waterloo)

We were really inspired by this as a way to build user engagement, and so on the eve of the recent Canadian federal election, we launched WebArchives.ca using the Shine build.

We used the University of Toronto's CPP collection, which was some 50 websites: all major political parties, minor ones, and a nebulous group of political interest groups, from David Suzuki's environmental organization to the Assembly of First Nations, to groups fighting for equal marriage.



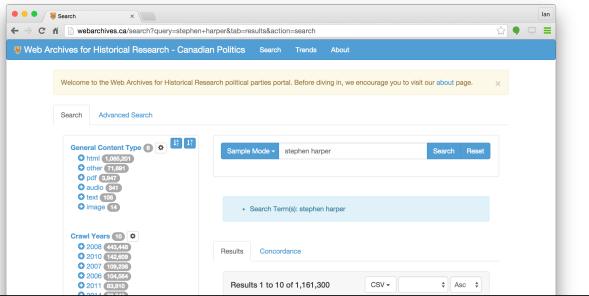




Five Things We Learned

- Political parties delete content
- User-generated comments were more common in political parties
- Absences can be more informative than presences
- We can see the rise/fall of prominent people
- Enabling user access is truly transformative

Good for public engagement - but limited for scholarship....



**Getting over my bias
towards content **and**
embracing metadata**

Gephi 0.9
[\(http://gephi.github.io/\)](http://gephi.github.io/)

**Walkthrough at
ianmilligan.ca: “From
Dataverse to Gephi” -
try it on this data!**

Step-by-Step Walkthrough

Once you've downloaded the file, open up Gephi. On the opening screen, you want to select "Open a Graph File..." and select the all-links-to-gephi-link.gexf file that you downloaded from our Dataverse page. You then want to click "OK" on the next page. Create a tree graph.

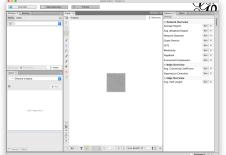
Do you want to make this link graph yourself from our data? Read on...

You should now see what I (metaphorically) call a bong cube. That's good, because it means that the data is in there. We need to make it usable, however.

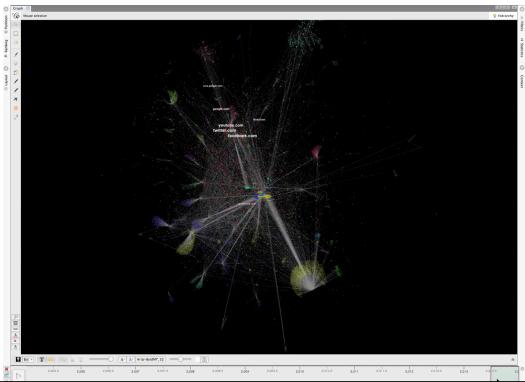
Click on the "Data Laboratory" tab at the top.

Click on "Nodes" above. When it is shaded behind it, that means that it is selected.

Click on "Copy Data to another Column" > Select ID, and then select "Node" on the drop down menu.



Metadata Extraction



December 2006
Stephane Dion Elected Leader of Party









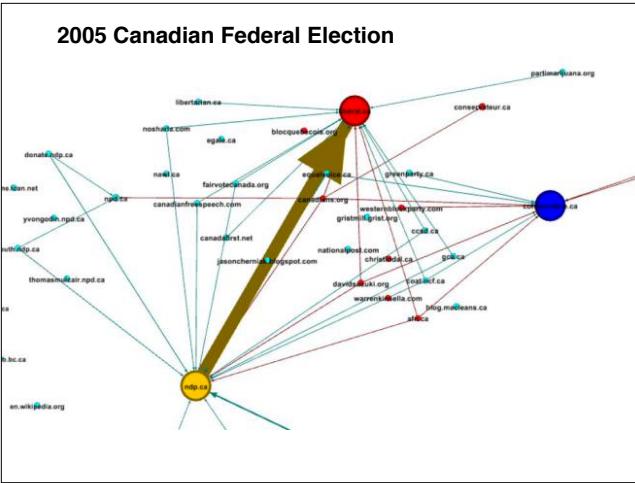
October 2008
Election Campaign - Advertisement Sites



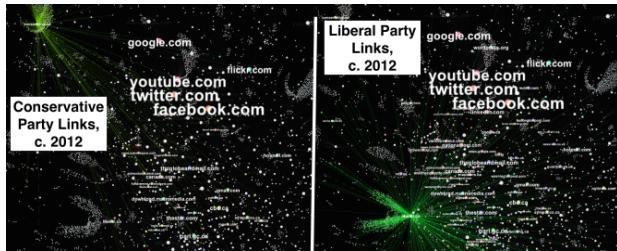
December 2008
Election campaign Ends; Attacking Harper
on Anti-American Grounds (bushharper)



2005 Canadian Federal Election



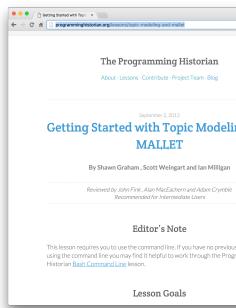
Metadata Extraction



**Topic Modelling using
MALLET ([http://
mallet.cs.umass.edu/](http://mallet.cs.umass.edu/))**

Walkthrough:

**[http://
programminghistorian.org/
lessons/topic-modeling-
and-mallet](http://programminghistorian.org/lessons/topic-modeling-and-mallet)**



Metadata Extraction

liberal.ca	27
liberal.ola.org	27
liberal.us1.list-manage.com	27
liberal.us1.list-manage1.com	27
liberal.us1.list-manage2.com	27
liberaluniversity.liberal.ca	27
license.icopyright.net	27
live.cbc.ca	27
lpc.ca	27
macleans.ca	27
masses.tao.ca	27
mcss.gov.on.ca	27
mediagnite.com	27
mediasales.cbc.ca	27
membercentre.cbc.ca	27
mentalhealthcommission.ca	27
metrics.mmailhost.com	27
mondesdefemmes.ca	27
music.cbc.ca	27
navl.ca	27
newswire.ca	27
nowtoronto.com	27
nrd.ca	27

colincarriemp.ca	12
colincarriemp.ca&lang=fr	12
colinmayers.ca	12
colinmayers.ca&lang=fr	12
congrespec.ca	12
conservateur.ca	12
conservateur.us5.list-manage.com	12
conservative.ca	12
conservative.us5.list-manage.com	12
consumersfirst.ca	12
cornellichius.ca	12
cornellichius.ca&lang=fr	12
costasmeneagakis.ca	12
costasmeneagakis.ca&lang=fr	12
cpcconvention.ca	12

Metadata Extraction

- Results @ <http://ianmilligan.ca/2015/02/05/topic-modeling-web-archive-modularity-classes/>

Metadata Extraction

- Conservative themes (2014): economic development, family, immigration, legislation, women's issues, senior issues, Ukrainians, constituency offices, some prominent (and not-so-prominent) MPs, and of course, our economic action plan.
- Liberal themes (2014): Justin Trudeau (the new leader), cuts to social programs, child poverty, mental health, municipal issues, labour, workers, Stop the Cuts, and housing.

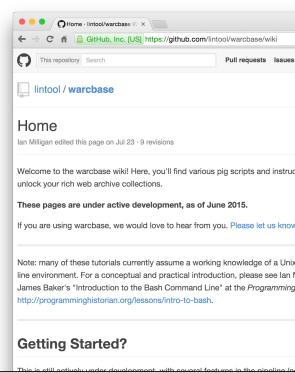
Metadata Extraction

- Conservative themes (2006): education, university, but **tons of information on Aboriginal issues**;
- Liberal themes (2006): community questions, electoral topics, universities, human rights, child care support.

Interdisciplinary

Warcbase

- Jimmy Lin (main developer), Ian Milligan, Jeremy Wiebe, Alice Zhou, all University of Waterloo
- Developing code, walkthroughs, and instructions for historians to be able to take a directory of WARCs and...

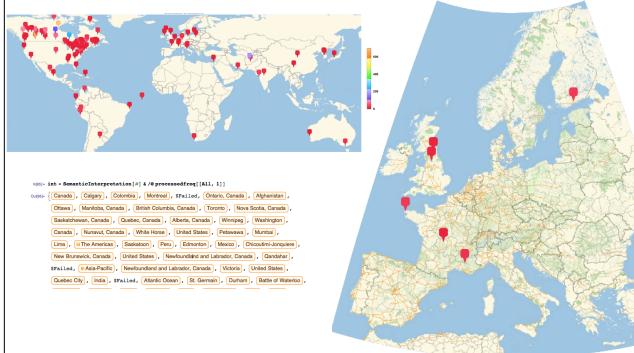


Extract all Plain Text

Extract Entities



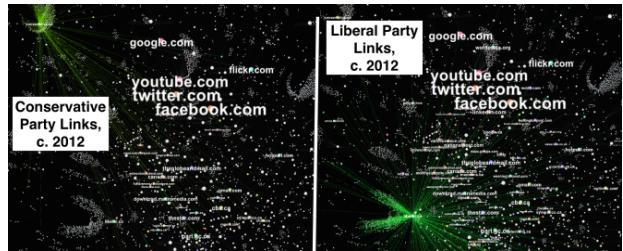
Extract Entities



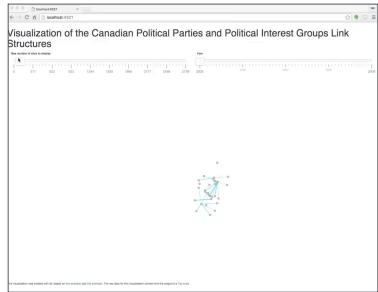
Extract Entities



Extract Links/Gephi Connector



Or D3.js link networks in browser



All walkthroughs at:
docs.warcbase.org

Bringing it all together in
a notebook environment



The screenshot shows a Jupyter Notebook interface with the title "Government Information Day - Demo". The notebook contains the following Scala code:

```
In [1]: !cp /Users/iamwilligan1/dropbox/varcbase/target/varcbase-0.1.0-SNAPSHOT-fatjar.jar ...
...
In [2]: import org.varcbase.spark.matchbox_
import org.varcbase.spark.rdd.RecordRDD_...
import org.varcbase.spark.matchbox_
import org.varcbase.spark.rdd.RecordRDD_...
Out[2]: 161 milliseconds
In [3]: var arc = "/Users/iamwilligan1/dropbox/warc-workshop/227-20051004191331-00000-crawling015.archive"
var warc = "/Users/iamwilligan1/dropbox/warc-sample-data/src-warc/ARCHIVE017-227-QUARTERLY-XU-GECCV-4.warc"
var arcdir = "/Users/iamwilligan1/dropbox/warc-workshop";
arc: String = /Users/iamwilligan1/dropbox/warc-workshop/227-20051004191331-00000-crawling015.archive
warc: org.warc.WarcRecordReader = org.varcbase.VarcRecordReader@332112831172-00000-crawling015.archive
arcdir: String = /Users/iamwilligan1/dropbox/warc-workshop
981 milliseconds
In [4]: val r =
  RecordLoader.loadarc(arc,
    .keepValidPages()
    .map(r => ExtractTopLevelDomain(r.getUrl))
...

```

Where to learn?



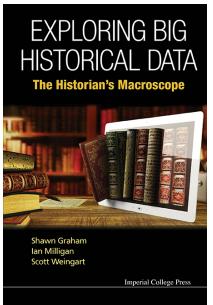
Programming Historian

- Network Analysis Lessons
- Topic Modeling Lessons
- Command Line Lessons
- etc.

Exploring Big Historical Data

- Check out our draft at macroscope.org

- Conceptual introduction to topic modelling
- Network analysis
- Visualizations
- Field of digital humanities

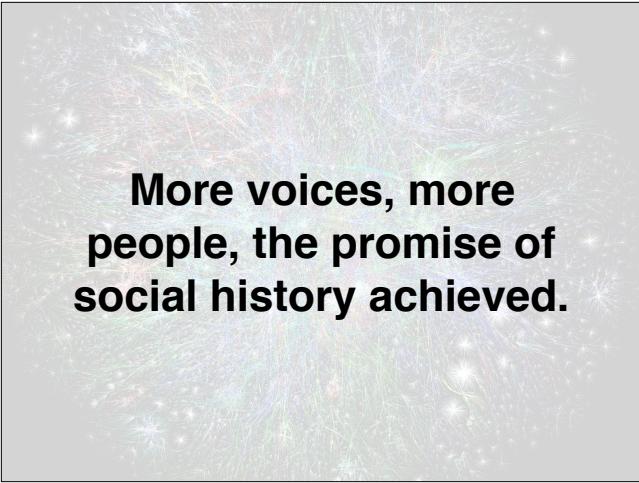


Events

- **Software Carpentry** – in-person events, looking into building connections with *Programming Historian*
- **Interdisciplinary hackathons** - *Archives Unleashed* (Toronto, March 2016; Washington, June 2016 - TBA)
- **Conferences** - Like this one, or others

... but most of all, a
willingness to learn and
fail.

Because, as I hope I
have shown today..
it's worth it.



**More voices, more
people, the promise of
social history achieved.**

Thanks very much!

Questions?

Ian Milligan
Assistant Professor
@ianmilligan1



UNIVERSITY OF WATERLOO
FACULTY OF ARTS
Department of History

And would love your questions today.