

Web Histories & Web Archives

Government Information Day 2016

Ian Milligan
Assistant Professor
@ianmilligan1



UNIVERSITY OF WATERLOO
FACULTY OF ARTS
Department of History

Why?

The sheer amount of social, cultural, and political information generated every day presents new opportunities for historians.



MATCH YOUR INTEREST TO A NEIGHBOR

FREE HOME PAGES AND E-MAIL	ARTS AUTOS BUSINESS COMPUTERS CULTURE	EDUCATION ENTERTAINMENT ENVIRONMENT FAMILY FASHION	FOOD GAMES GAY & LESBIAN GOVERNMENT HEALTH	KIDS MUSIC PEOPLE RECREATION SCIENCE F...
--	---	--	--	---



199

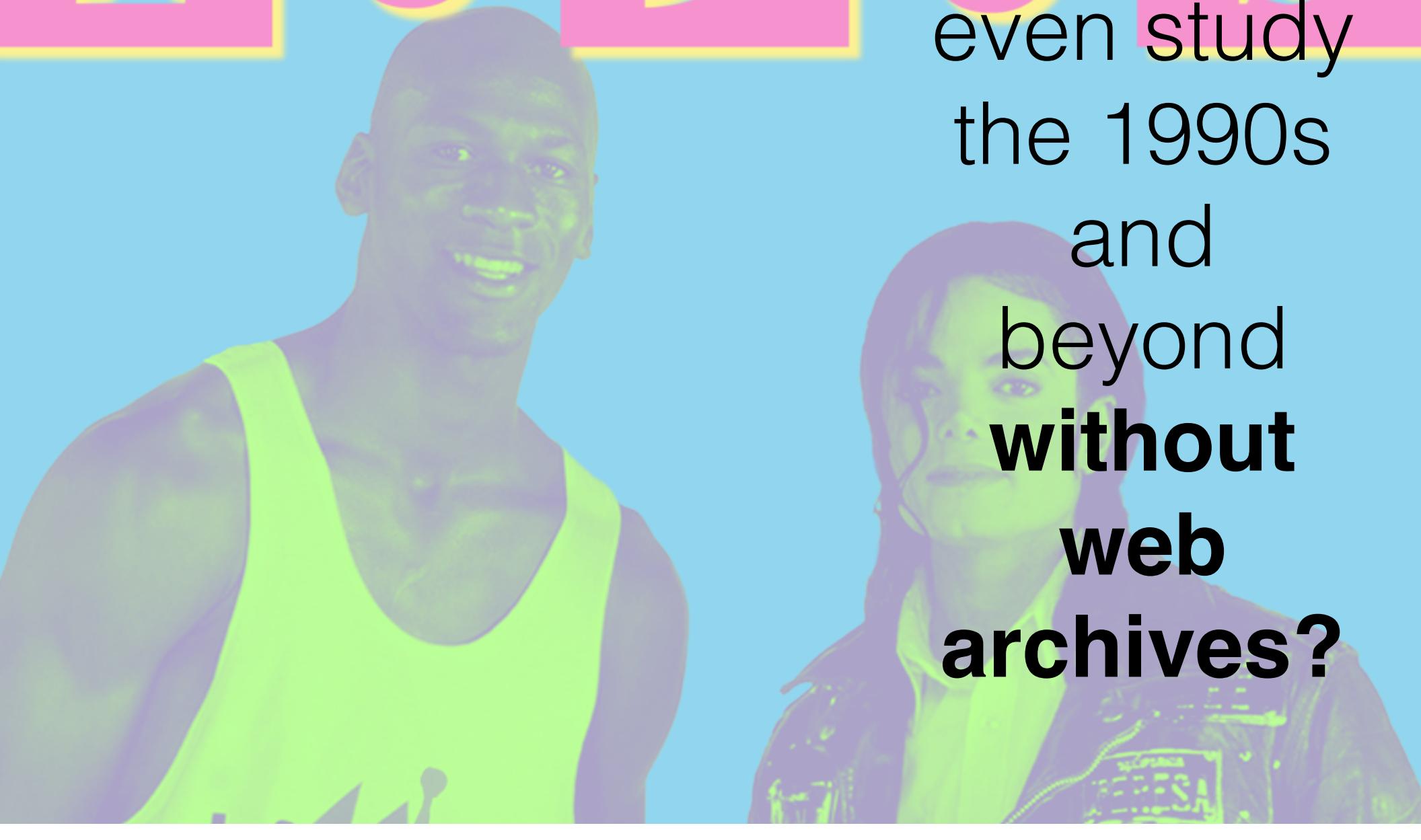
99

99

0

S

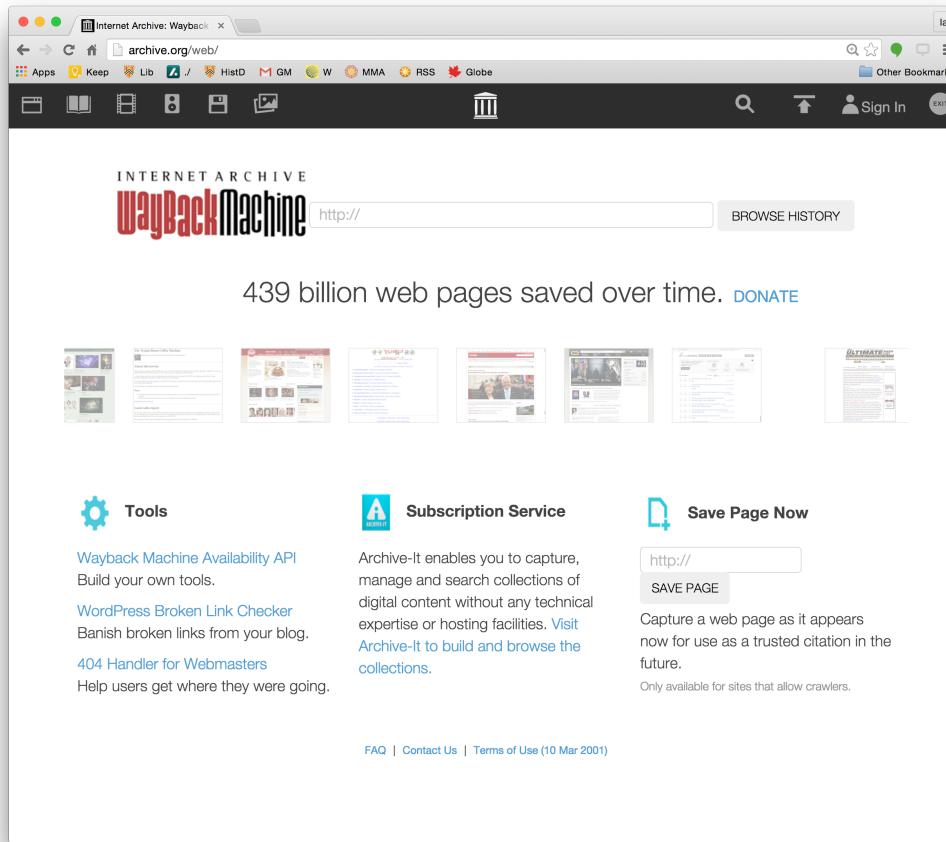
Could one
even study
the 1990s
and
beyond
without
web
archives?



No.

Historians need to do this now, or
we're going to be left behind.

Nightmare Scenario



This won't be enough!



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

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...

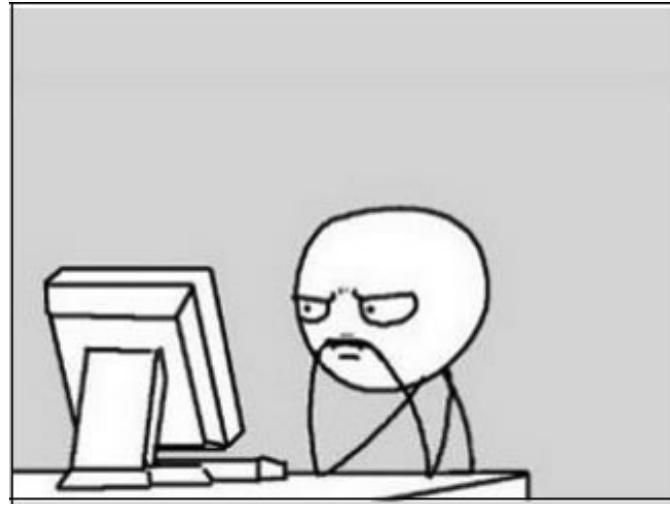
But what will web archives look like?

- Two Case Studies
 - **Wide Web Scrape**, March - December 2011 (Internet Archive) (sample of 80TB WARC collection);
 - **Archive-It Longitudinal Collections, Canadian Political Parties & Interest Groups**, 2005-present (Archive-It/University of Toronto)

Similarities -

Windows into the lives of
everyday people.





Differences -
Incredible range of technical
skills/no common platform!

A prologue to the
political/government
data..

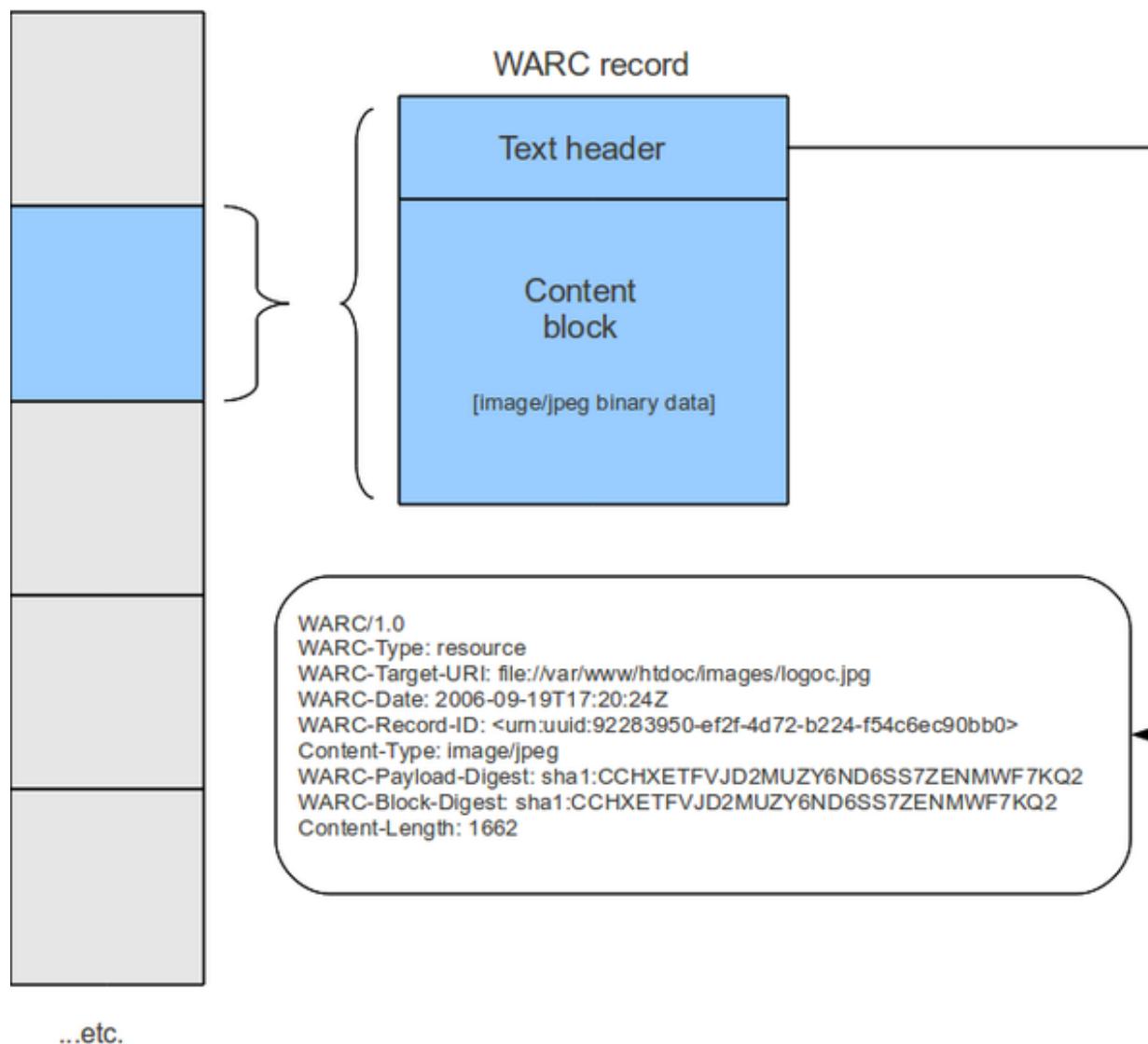
Case Study One

- The **Wide Web Scrape** (~ 80TB) - Snapshot of the Web
- **85,570** WARC files, CDX metadata
- Using it here as **an introduction to what the scale of this data can look like**

The screenshot shows a web browser window with a dark theme. The address bar displays the URL <https://archive.org/details/wide00002&tab=about>. The main content area is titled "Wide Crawl started March 2011". Below the title, there is a brief description: "Web wide crawl with initial seedlist and crawler configuration from March 2011. This uses the new HQ software for distributed crawling by Kenji Nagahashi." A "MORE" link is visible. At the bottom of this section, there are three navigation links: "About" (which is underlined), "Collection", and "Forum". To the right of this main content, there is a sidebar with a yellow header "Created on October 5 2010" and a profile picture of a woman. Below this, there are sections for "ADDITIONAL CONTRIBUTOR" (with a profile picture of a man and the name "brewste Archivist") and "VIEWS" (with a profile picture of a person and the name "kngenie Archivist"). At the very bottom of the sidebar, there is a note: "However this was a somewhat experimental crawl for us," followed by a small set of navigation icons.

WARC File Format

WARC file



ca,yorku,justlabour)/ 20110714073726
<http://www.justlabour.yorku.ca/> text/html
302 3I42H3S6NNFQ2MSVX7XZKYAYSCX5QBYJ
[http://www.justlabour.yorku.ca/index.php?
page=toc&volume=16](http://www.justlabour.yorku.ca/index.php?page=toc&volume=16) - 462 880654831
WIDE-20110714062831-crawl416/
WIDE-20110714070859-02373.warc.gz

Top-Level Domain	Number of Distinct URLs Downloaded in Sample	Number of Overall URLs in Wide Web Scrape (selected domains)	Percentage of URLs Captured
.com	29,219,706	1,260,409,874	2.32%
.org	2,489,050	96,681,268	2.57%
.net	2,438,903	140,726,805	1.73%
.edu	350,482	6,620,283	5.29%
.gov	97,484	2,205,332	4.42%
.mil	10,268	103,507	9.92%
.ca	622,365	8,512,275	7.31%
.uk	464,991	21,870,821	2.13%
.fr	239,160	13,654,404	1.75%
.in	105,287	3,736,316	2.82%
.cn	5,499,593	133,105,864	4.13%
.ke	4883	37,871	12.89%
TOTAL	41,542,172	1,687,664,620	2.46%

CDX Files (finding aids)



WARC File

Plain Text
Conversion

Indexing

Carrot2 Workbench

Search

Source: Solr

Algorithm: Lingo

Basic

Query (Required):

Read Solr clusters if present

Results: 1000

[Aduna Cluster Map Visualization](#) [Circles Visualization](#) [FoamTree Visualization](#)

children (1000 documents from Solr, 47 clusters from Lingo)

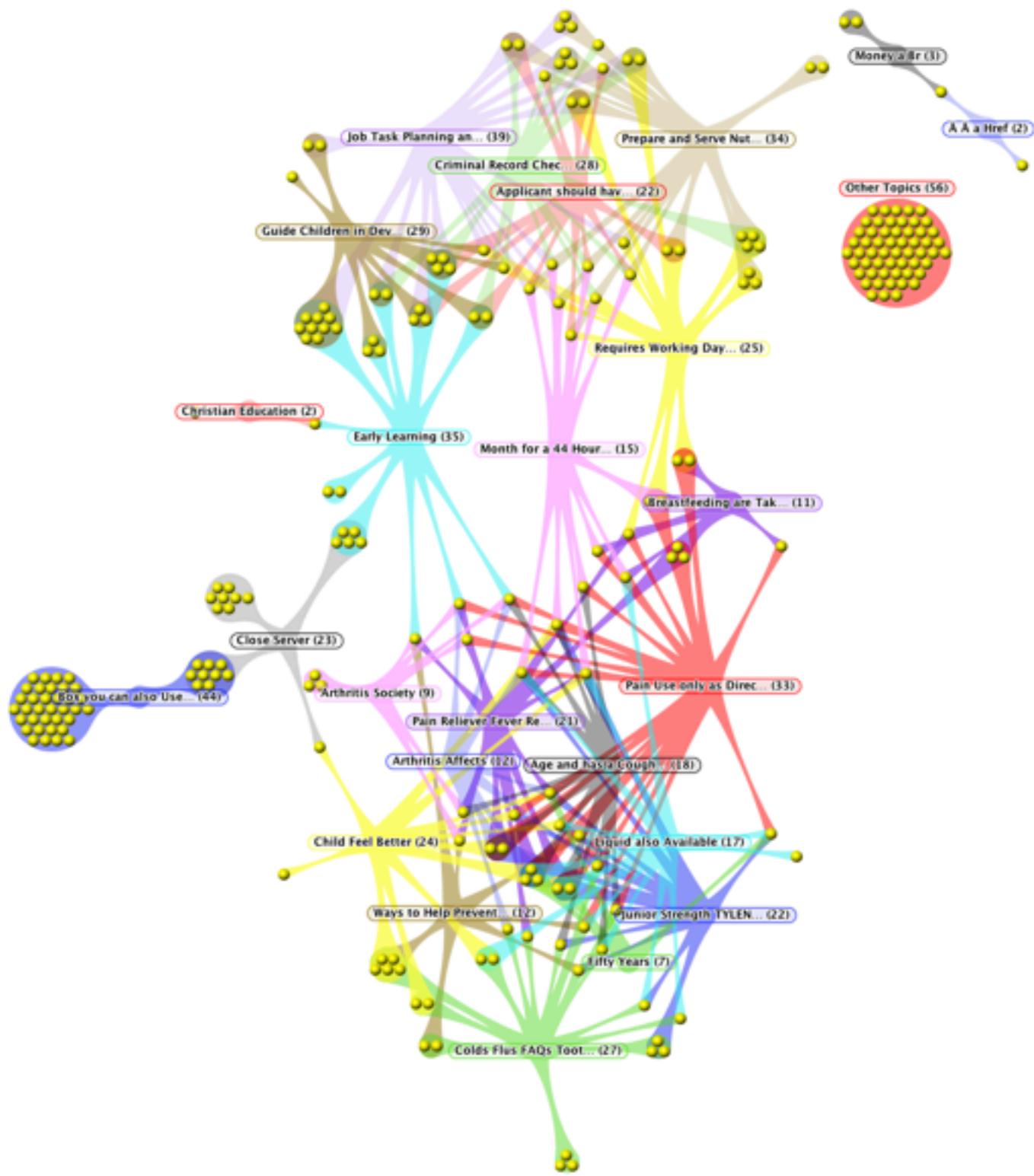
Clusters

- Child Health (192)
- Canada Service (168)
- Left side of the Page (161)
- Document Input (158)
- Research Research (147)
- Health Centre (138)
- Children Value (127)
- Services Community (123)
- Consumer Product (122)
- Providing Services (120)
- Health Community (113)
- School Services (111)
- Health and Wellness (105)
- Health Services (103)
- New Image (101)
- Returns List (98)
- Support Services (98)
- Public Health (97)
- Health and Safety (95)
- Family Services (93)
- Education Document (92)
- Service Days (91)
- Research Programs (88)
- Health Promotion (84)
- Development Research (83)
- Research will Help (82)
- Youth Services (82)
- Services Community Education (74)
- Health Professionals (74)
- Research Resources (69)
- Areas of Health (63)
- University of Ottawa (58)
- Community Health Centre (54)
- Research and Events (56)
- Mental Health (53)
- Health Issues (54)
- Research Interests (50)
- Invitation Templates (48)
- University University of Ottawa f (46)
- Flu Is Available (38)
- Natural Health Products (38)
- Products and Services (35)
- Birthday Party Invitations (27)
- Centre for Research on Commun (24)
- Birthday Age (5)
- Youth Services Bureau of Ottawa (5)
- Other Topics (365)

Documents

- [1] <http://www.tylenol.ca/children/children-6-11-years/cough-cold-products>: text/html; charset=utf-8 Infants 0-24 Months Children 2-5 Years Children 6-11 Years Tylenol logo Home | Contact us | Français Search _____ Search * Adult * Children * Products... /Users/kanniligan1/Desktop/output/76-Canadian-456.html
- [1] <http://www.tylenol.ca/children/children-6-11-years/children-6-11-years>: text/html; charset=utf-8 Infants 0-24 Months Children 2-5 Years Children 6-11 Years Tylenol logo Home | Contact us | Français Search _____ Search * Adult * Children * Products... /Users/kanniligan1/Desktop/output/76-Canadian-1721.html
- [1] <http://www.tylenol.ca/children/children-3-5-years/children-3-5-years>: text/html; charset=utf-8 Infants 0-24 Months Children 2-5 Years Children 6-11 Years Tylenol logo Home | Contact us | Français Search _____ Search * Adult * Children * Products... /Users/kanniligan1/Desktop/output/72-Canadian-1170.html
- [1] <http://www.tylenol.ca/children/products>: text/html; charset=utf-8 For Adults For Children Tylenol logo Home | Contact us | Français Search Search * Adult * Children * Products... About Tylenol * News & Information All Children's Products... /Users/kanniligan1/Desktop/output/25-Canadian-3512.html
- [1] <http://blogs.afortunecookie.ca/tag/children/feed/>: text/html; charset=utf-8 Infants 0-24 Months Children 2-5 Years Children 6-11 Years Tylenol logo Home | Contact us | Français Search _____ Search * Adult * Children... /Users/kanniligan1/Desktop/output/23-Canadian-2494.html
- [1] <http://www.tylenol.ca/children/children-6-11-years/aches-pains/about-aches-pain>: text/html; charset=utf-8 Infants 0-24 Months Children 2-5 Years Children 6-11 Years Tylenol logo Home | Contact us | Français Search _____ Search * Adult * Children... /Users/kanniligan1/Desktop/output/70-Canadian-886.html
- [1] <http://www.tylenol.ca/children/children-3-5-years/aches-pains/reducing-your-child-s-aches-pains>: text/html; charset=utf-8 Infants 0-24 Months Children 2-5 Years Children 6-11 Years Tylenol logo Home | Contact us | Français Search _____ Search... /Users/kanniligan1/Desktop/output/29-Canadian-2278.html
- [1] <http://www.tylenol.ca/children/children-6-11-years/cough-cold-flu/symptoms>: text/html; charset=utf-8 Infants 0-24 Months Children 2-5 Years Children 6-11 Years Tylenol logo Home | Contact us | Français Search _____ Search... /Users/kanniligan1/Desktop/output/40-Canadian-1.html
- [1] <http://www.tylenol.ca/children/children-6-11-years/aches-pains/reducing-your-child-s-aches-pains>: text/html; charset=utf-8 Infants 0-24 Months Children 2-5 Years Children 6-11 Years Tylenol logo Home | Contact us | Français Search _____ Search... /Users/kanniligan1/Desktop/output/37-Canadian-2224.html

Other Topics



children (250 documents from Solr, 26 clusters from Lingo)

Clusters

- Box you can also Use it Program
- Job Task Planning and Organizizi
- Early Learning (35)
- Prepare and Serve Nutritious Me
- Pain Use only as Directed (33)

Documents

[190] <http://www.lutheranchurch.ca/missions.php?s=nicaragua&p=6&print=yes> : text/html; charset=latin1_swedish_ci

CLWR funds Nicaraguan medical and dental clinic, scholarships

2010 [Nicaraguan_medic... /Users/ianmilligan1/Desktop]

Services

- Open Link
- Open Link in New Window
- Download Linked File
- Copy Link

Search With Google

WaybackMachine

New TextWrangler Document with Selection

EasyFind: Find Selection...

Add to iTunes as a Spoken Track

Open URL

Add to Reading List

- Age and has a Cough or Cold (1)
- Liquid also Available (17)
- Month for a 44 Hour Week (15)
- Arthritis Affects (12)
- Ways to Help Prevent Earaches (
- Breastfeeding are Taking (11)
- Arthritis Society (9)
- Fifty Years (7)
- Money a Br (3)
- Christian Education (2)
- Ã¢â€ša Href (2)
- Other Topics (56)

Lutheran Church-Canada

http://www.lutheranchurch.ca/news.php?id=158&print=yes

INTERNET ARCHIVE WaybackMachine 3 captures 5 Dec 10 - 14 Jul 11 DEC JUL 14 2010 2011

LUTHERAN CHURCH-CANADA ÉGLISE LUTHÉRIENNE du CANADA

CLWR funds Nicaraguan medical and dental clinic, scholarships

Friday, January 22, 2010

WINNIPEG – Canadian Lutheran World Relief (CLWR) has announced \$36,500 in funding for two Lutheran Church-Canada (LCC) programs in Nicaragua this year.

The announcement was made as Iglesia Luterana Sinodo de Nicaragua (ILSN) prepares for its first biennial convention and includes new money for a medical and dental clinic and increased school scholarships.

The medical clinic, which began operations in May 2009, is open every Thursday beginning at 8 a.m. and remains open until all patients have been seen.

The clinic is staffed by a doctor and a dentist, who see an average of 40-45 patients each week, and provides common medications because many patients are too poor to purchase them.

CLWR will continue supporting the Christian Children Education Program. The program, conducted in all 23 congregations of ILSN, provides an average of 25 scholarships in each community to the neediest children. The scholarships include the required school uniforms, shoes, backpacks and school supplies.

Each child is also enrolled in the tutoring and Christian-education class held five days a week when children are not in school (Children attend school in the morning or in the afternoon.)

These classes, held in the churches and led by teachers and deaconesses, provide tutoring and homework support for the children in math, Spanish and other subjects. A portion of the time is also set aside for Christian education and cultural activities.

More than 750 children are enrolled in the program. CLWR has provided support for about 250 children.

Since 1999, CLWR has partnered with LCC to support community-development projects.

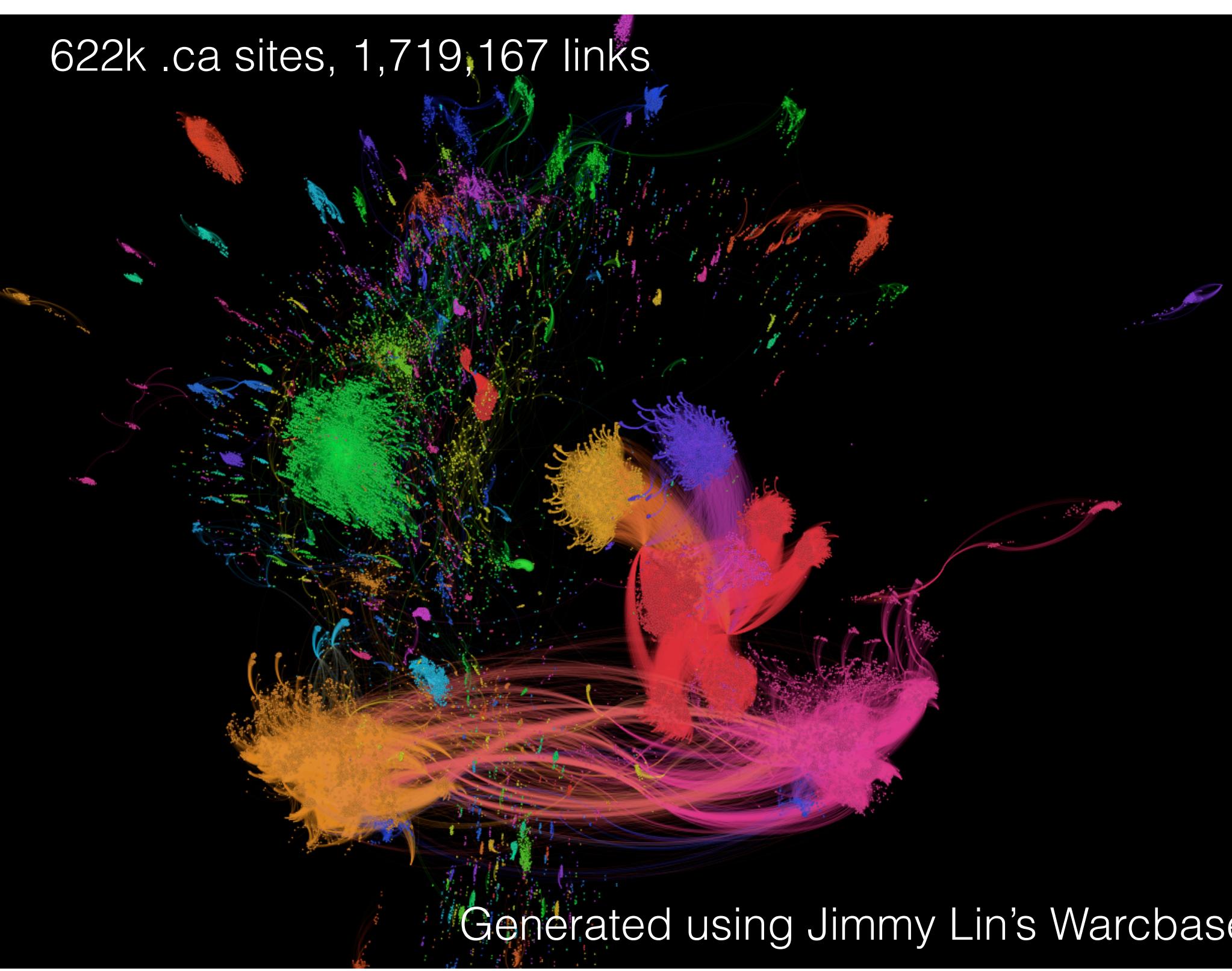
Robert Granke, executive director of CLWR, visited congregations of the ILSN in November. You can read more about his visit at www.lccontheroad.ca, The Canadian Lutheran or in the forthcoming issue of CLWR's Partnership newsletter due out in early February.



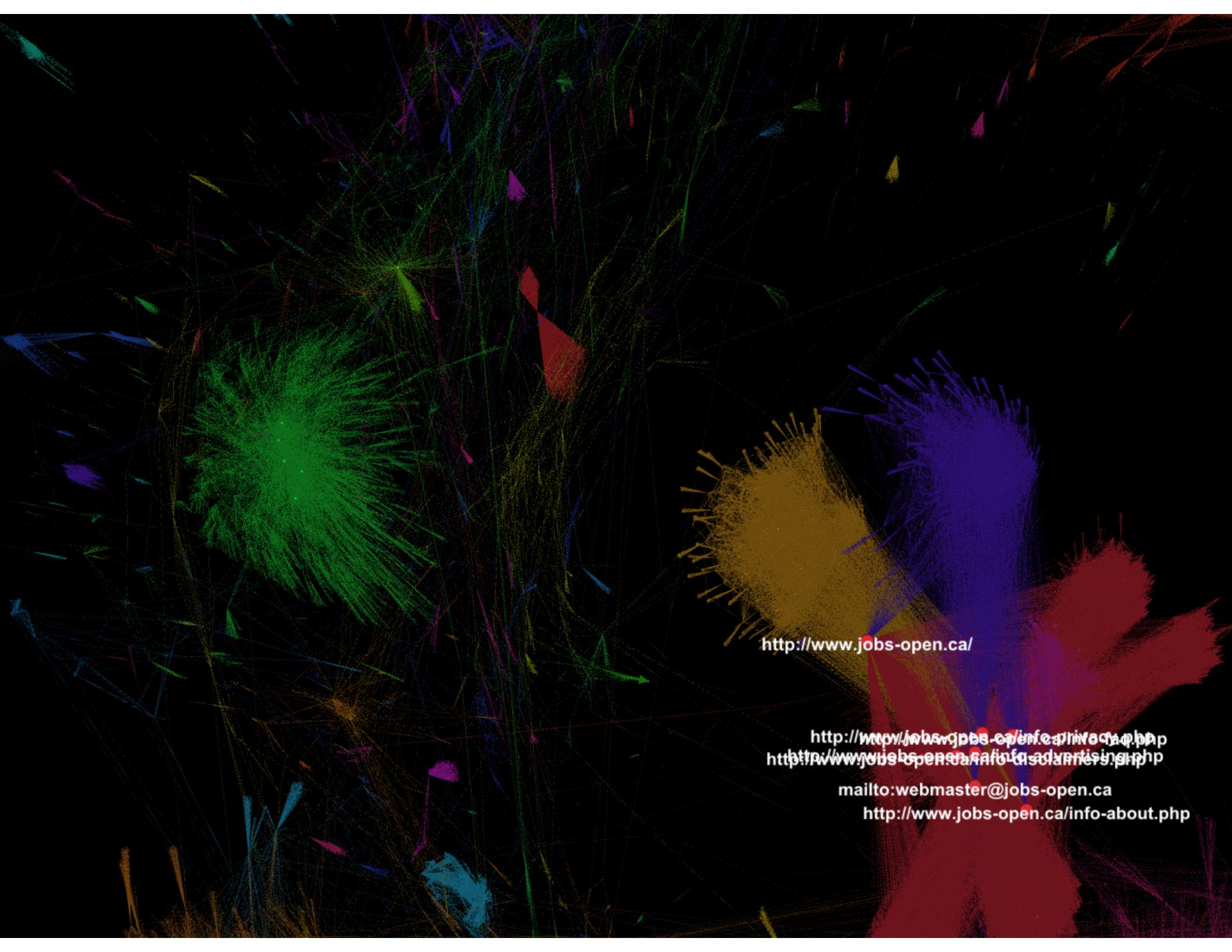
A medical clinic in Nicaragua.

Problem is.. you need to
know what you're looking
for!

622k .ca sites, 1,719,167 links



Generated using Jimmy Lin's Warchbase

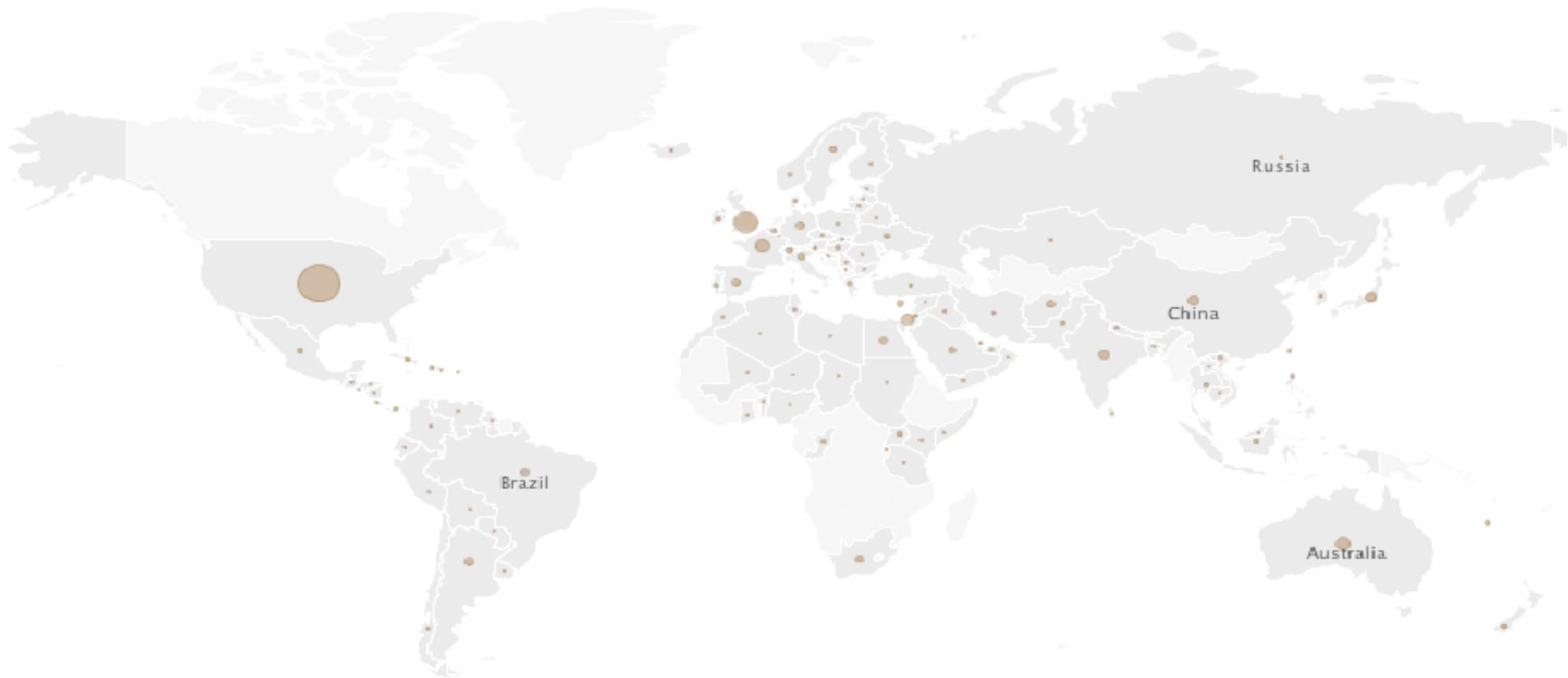


<http://www.jobs-open.ca/>

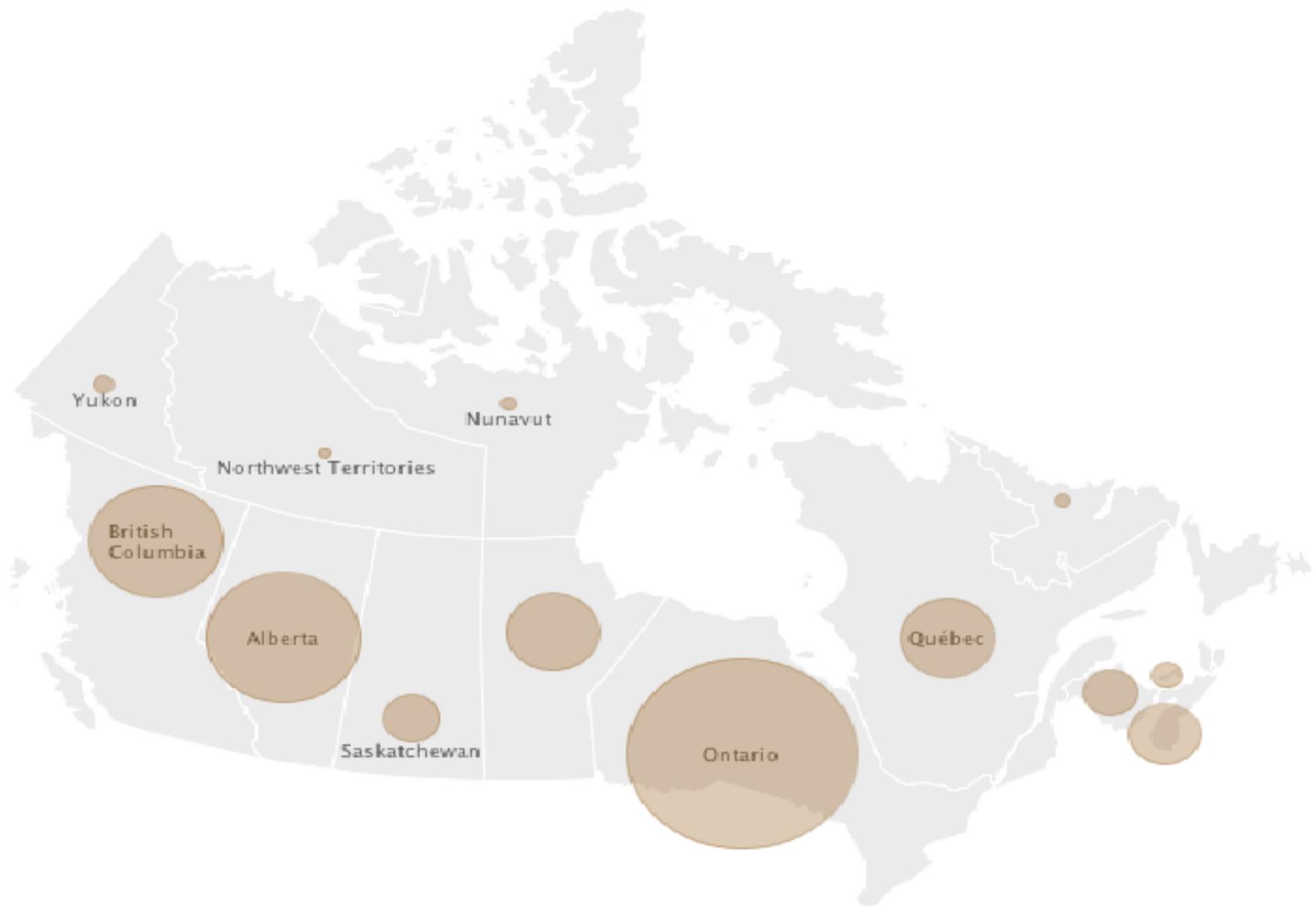
<http://www.jobs-open.ca/info-about.php>
<http://www.jobs-open.ca/advertising.php>

mailto:webmaster@jobs-open.ca
<http://www.jobs-open.ca/info-about.php>

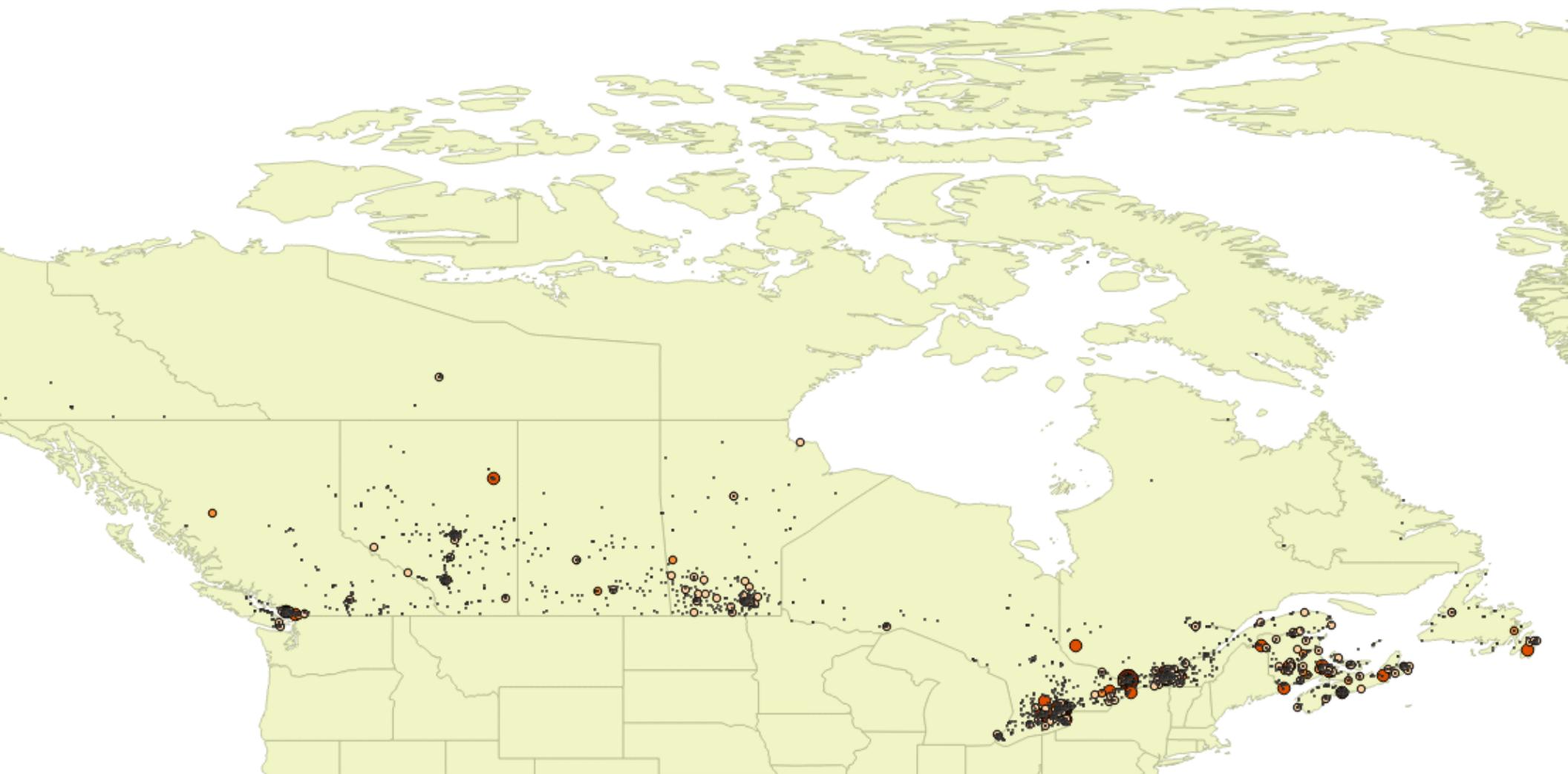
Countries Mentioned in .ca TLD (excluding Canada)



Provinces Mentioned in .ca TLD



Canadian Postal Codes visualized



With longitudinal analysis, we could do fantastic things - contextualizing subsequent studies.

Wide Web Scraps and the **Dream of Social History.**

Case Study Two

- **Archive-It Research Services:** “Canadian Political Parties and Political Interest Groups” Collection
(Collaborative effort w/
Archive-It & University of
Toronto’s Nicholas Worby)
- 2005 - 2015
- WAT & WARC files

The screenshot shows a web browser window with the URL <https://archive-it.org/collections/227>. The page title is "Archive-It - Canadian Political Parties and Political Interest Groups". The header includes links for HOME, EXPLORE, LEARN MORE, and CONTACT US. Below the header, it says "Explore > University of Toronto > Canadian Political Parties and Political Interest Groups". A large green box features the Archive-It logo and the text "Canadian Political Parties and Political Interest Groups" along with "Collected by: University of Toronto". It also notes "Archived since: Oct, 2005" and "Description: Canadian Political Parties and Political Interest Groups". The main content area is titled "Narrow Your Results" and lists categories like "New Democratic Party of Canada (2)", "Assembly of First Nations (1)", etc. At the bottom, there are buttons for "Sites" and "Search Page Text", and a footer note "Page 1 of 1 (54 Total)".

WAT Files

**Hyperlink Information
(URLs/Anchor Text)**

WARC file



WARC record



...etc.

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?



Current Interface

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

The screenshot shows a web browser window displaying the Archive-It.org collection page for "Canadian Political Parties and Political Interest Groups". The URL in the address bar is <https://archive-it.org/collections/227?q=Stephen+Harper&page=1&show=Sites>. The page features a header with the Archive-It logo, navigation links for HOME, EXPLORE, LEARN MORE, and CONTACT US, and a tagline about collecting and accessing cultural heritage on the web. Below the header, it says "Explore >> University of Toronto >> Canadian Political Parties and Political Interest Groups". The main content area displays the "Canadian Political Parties and Political Interest Groups" collection, which was collected by the University of Toronto and archived since Oct, 2005. It includes a brief description, subject terms (Politics & Elections), and a collector note. A search bar at the bottom allows users to search within the collection results.

The need to **democratize**
access, so that ~~huddites~~
historians can use them.

The Internet Archive will la Ian

www.theverge.com/2015/10/22/9593656/internet-archive-wayback-machine-redesign-announced

THE VERGE

TRENDING NOW
Amazon is opening its first physical bookstore today

34 NEW ARTICLES

PREVIOUS STORY
FCC passes rule cracking down on prison phone call charges

NEXT STORY
Researchers discover new attacks amid VoLTE rollout

TECH

The Internet Archive will launch a modernized Wayback Machine in 2017

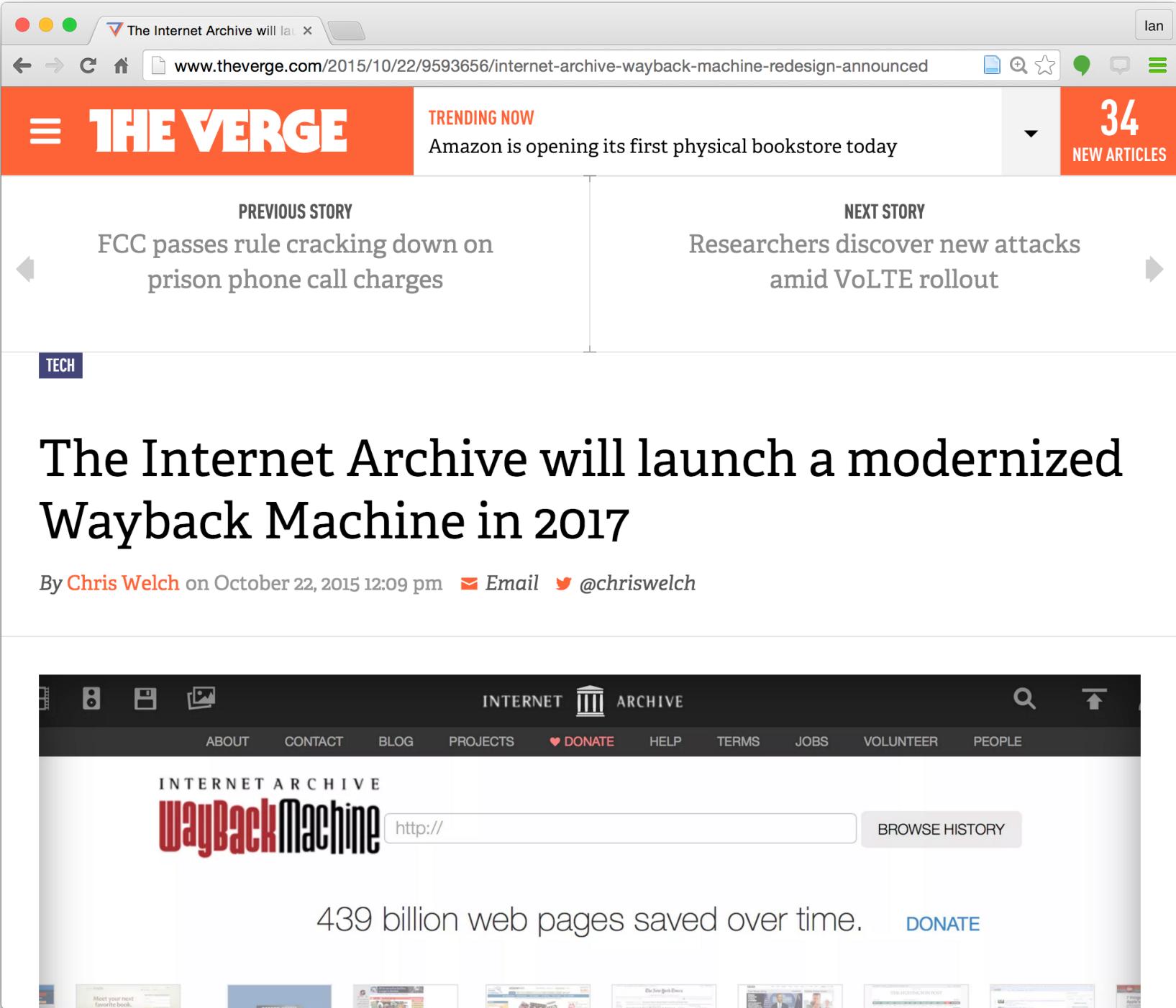
By Chris Welch on October 22, 2015 12:09 pm Email @chriswelch

INTERNET ARCHIVE

ABOUT CONTACT BLOG PROJECTS DONATE HELP TERMS JOBS VOLUNTEER PEOPLE

INTERNET ARCHIVE WayBack Machine http:// BROWSE HISTORY

439 billion web pages saved over time. DONATE



ukwa/shine lan

GitHub, Inc. [US] <https://github.com/ukwa/shine>

This repository Search Pull requests Issues Gist

ukwa / shine Unwatch 13 Unstar 7 Fork 2

Prototype SOLR-powered web archive exploration UI. <https://github.com/ukwa/shine/wiki>

637 commits 1 branch 0 releases 5 contributors

Branch: master → shine / +

GilHoggarth Added trailing slash to web archive url Latest commit 11ace26 on Sep 18

File	Description	Time
python	jisc ssd ~506k ~optimized to 7 segments	2 months ago
shine	Added trailing slash to web archive url	2 months ago
.gitattributes	gitattribute file	2 years ago
.gitignore	Prevent cache being versioned.	a year ago
.gitmodules	Initial "check-in" of the bootstrap submodule	2 years ago
.travis.yml	Looks like it's simpler than I expected.	a year ago
README.md	Added Travis-CI status and a brief outline.	2 years ago

README.md

Shine

A prototype web archives exploration UI, based on a Solr back-end that has been populated using the [warc-discovery](#) indexer.

build passing

Code

- Issues 40
- Pull requests 0
- Wiki
- Pulse
- Graphs

HTTPS clone URL
<https://github.com>

You can clone with [HTTPS](#), [SSH](#), or [Subversion](#).

[Clone in Desktop](#)

[Download ZIP](#)



The Canadian Political Party webarchives.ca

Welcome to the Web Archives for Historical Research political parties portal. Before diving in, we encourage you to visit our [about](#) page. x

The Canadian Political Parties and Political Interest Groups Portal

On this website, you can search web archived content from 50 political parties and political interest groups, from October 2005 to March 2015.

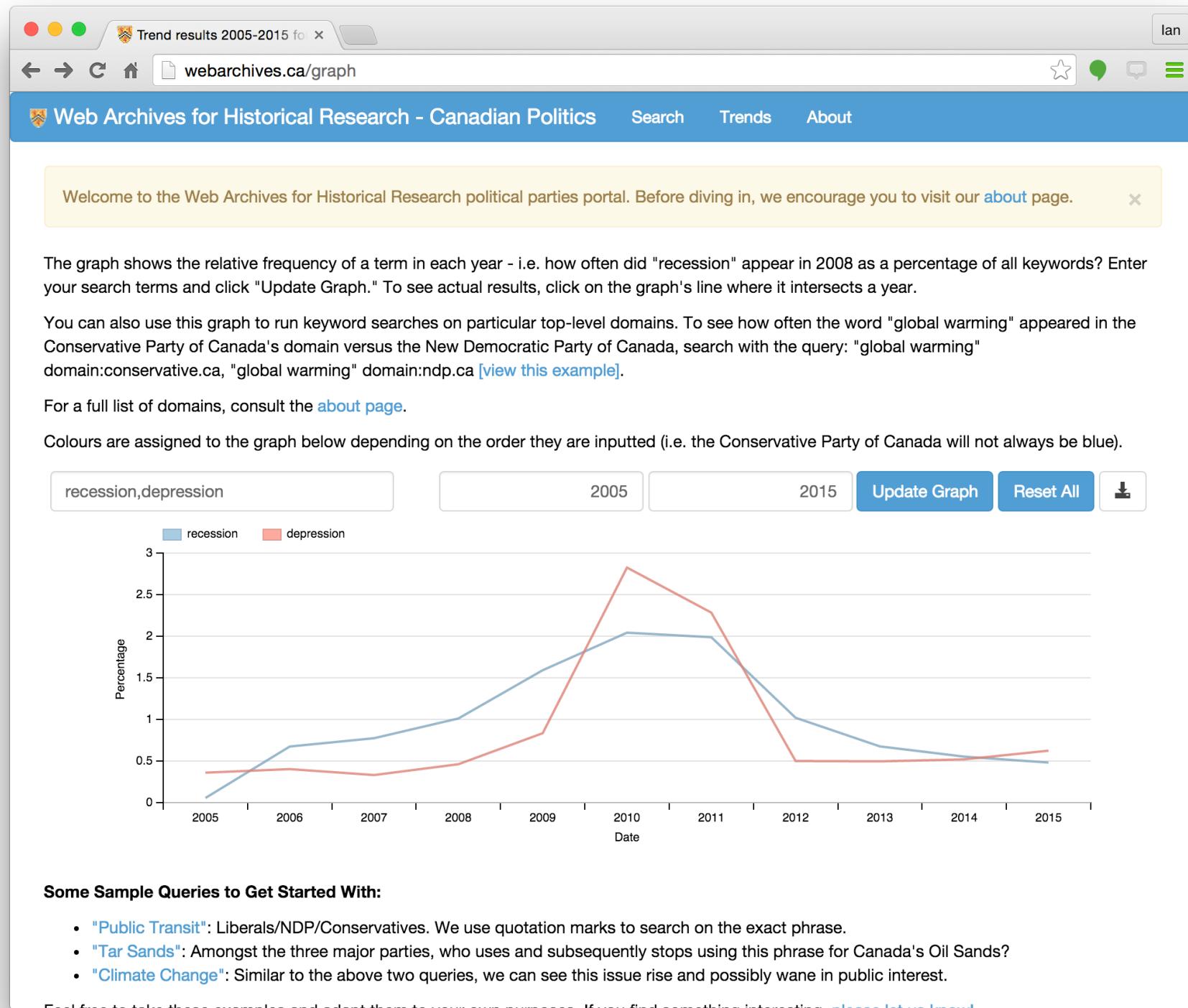
Curious how the Liberal Party of Canada responded to the 2008 financial crisis ([a search for "recession" in 2008, liberal.ca](#))? How the Canadian Centre for Policy Alternatives [reacted to Michael Ignatieff](#)? Now you can check it all out.

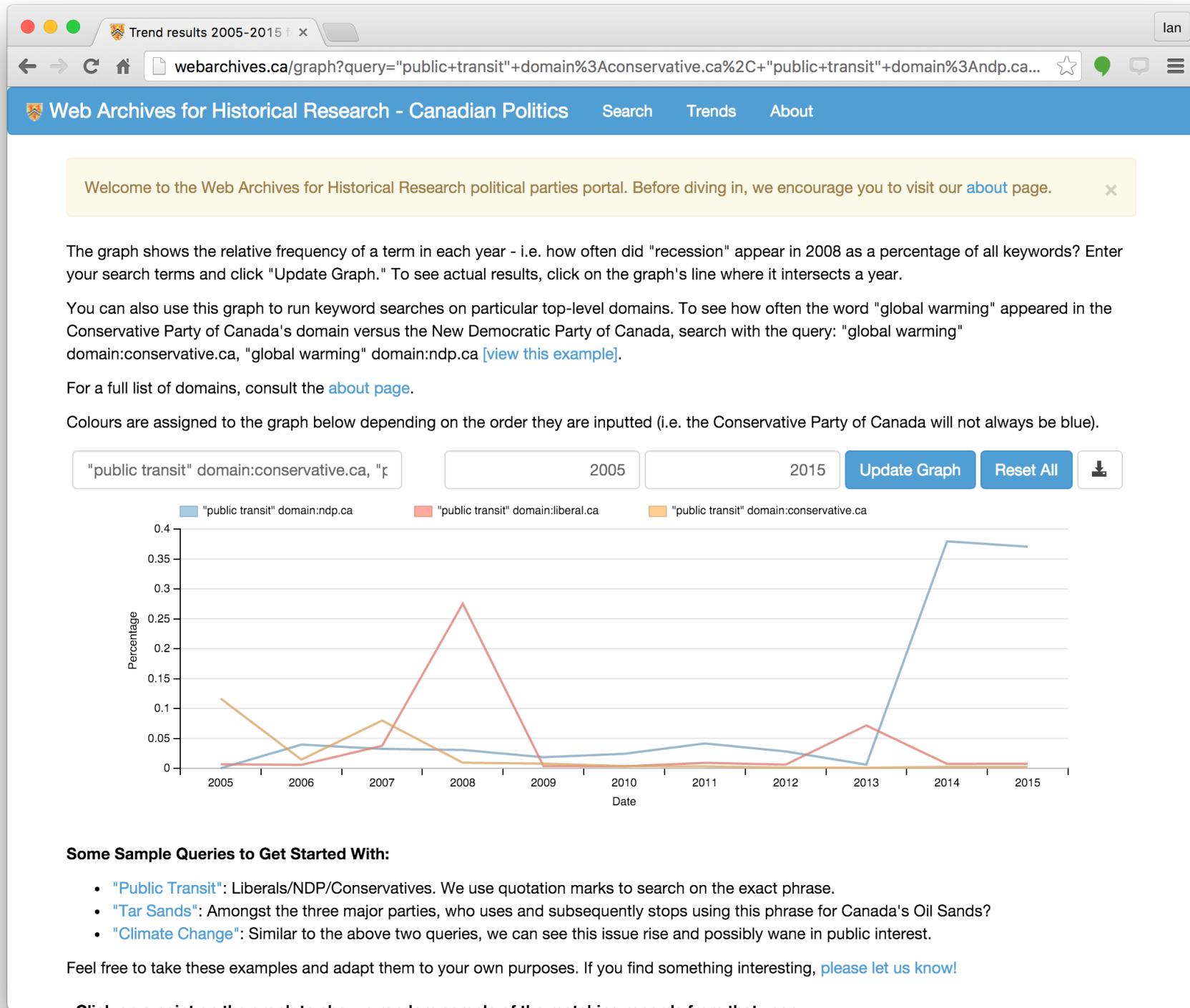
Options include:

- [Basic keyword searching](#) [Example: "Rob Ford", only Liberal.ca]
- [Graphing trends over time](#) [Example: Liberal Opposition Leaders, 2005-2015]
- [Advanced search, including words in proximity to each other](#) [Example: environmental and tax within 25 words of each other]

Below, here are all of the links for the entire time period, visualized below.

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





Welcome to the Web Archives for Historical Research political parties portal. Before diving in, we encourage you to visit our [about](#) page. X

The Canadian Political Parties and Political Interest Groups Portal

On this website, you can search web archived content from 50 political parties and political interest groups, from October 2005 to March 2015.

Curious how the Liberal Party of Canada responded to the 2008 financial crisis ([a search for "recession" in 2008, liberal.ca](#))? How the Canadian Centre for Policy Alternatives reacted to [Michael Ignatieff](#)? Now you can check it all out.

Options include:

- [Basic keyword searching](#) [Example: "Rob Ford", only Liberal.ca]
- [Graphing trends over time](#) [Example: Liberal Opposition Leaders, 2005-2015]
- [Advanced search, including words in proximity to each other](#) [Example: environmental and tax within 25 words of each other]

Below, here are all of the links for the entire time period, visualized below.



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....

The screenshot shows a web browser window with the following details:

- Address Bar:** webarchives.ca/search?query=stephen+harper&tab=results&action=search
- Page Title:** Web Archives for Historical Research - Canadian Politics
- Header:** Search, Trends, About
- Welcome Message:** Welcome to the Web Archives for Historical Research political parties portal. Before diving in, we encourage you to visit our [about](#) page.
- Search Options:** Search, Advanced Search
- General Content Type:** html (1,085,201), other (71,691), pdf (3,947), audio (341), text (106), image (14)
- Sample Mode:** stephen harper (Search, Reset)
- Search Term(s):** stephen harper
- Crawl Years:** 2008 (443,448), 2010 (142,609), 2007 (109,236), 2006 (104,564), 2011 (83,910), 2014 (70,746)
- Navigation:** Results, Concordance
- Results Summary:** Results 1 to 10 of 1,161,300
- Download Options:** CSV ▾, Asc ▾

Shine

- **Advantages:** accessible to the general public, easy to use, interactive trend diagram allows digging down for context, can move down to level of document itself.
- **Disadvantage:** keyword searching requires you know what to look for; random sampling misleading when tens of thousands of records; etc.
- Doesn't take advantage of what makes web sources so powerful: hyperlinks

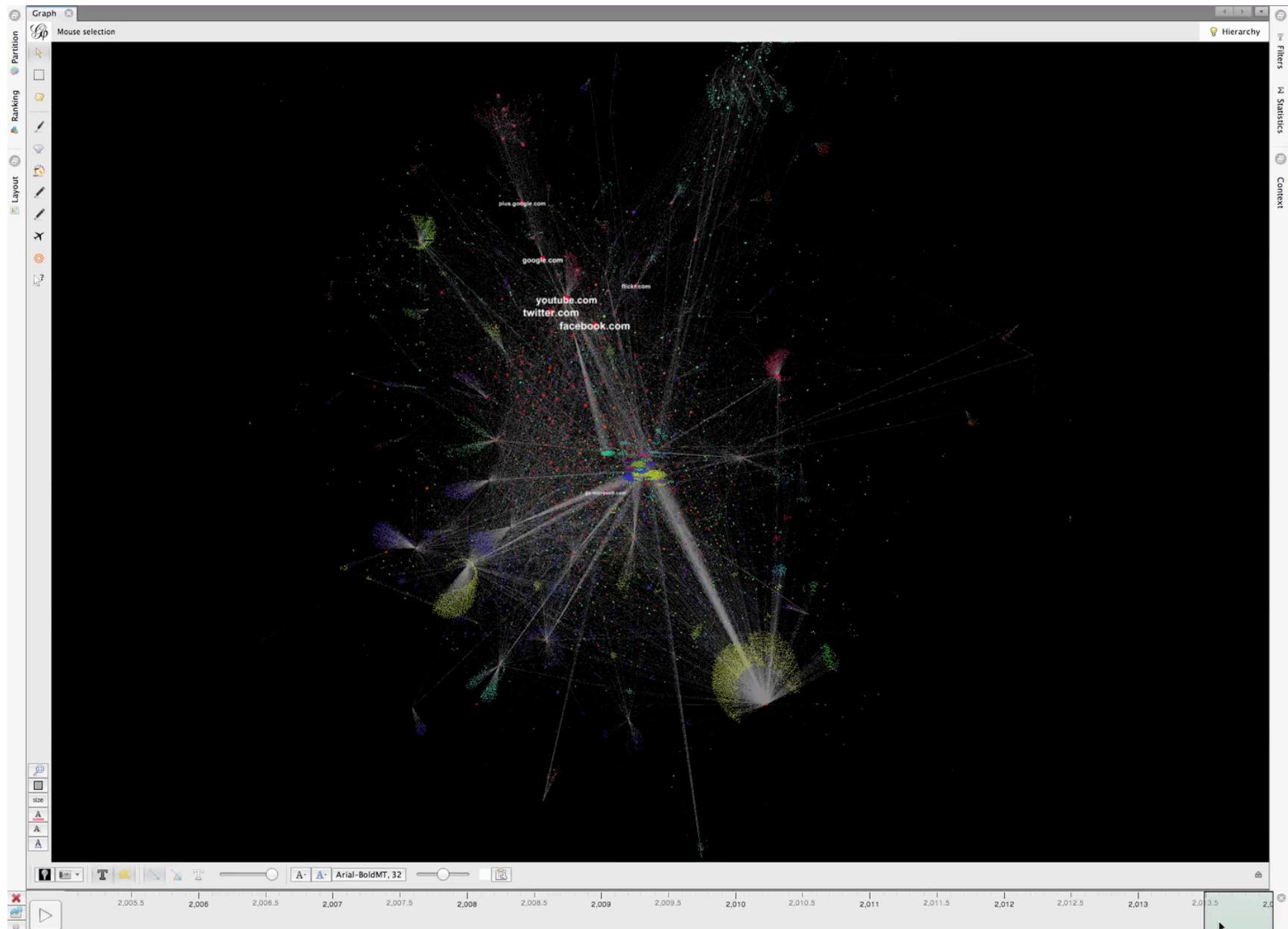
Getting over my bias
towards content **and**
embracing metadata

WATs vs. WARCs?



WATs help us find the files
we need to use - and to
contextualize them

Metadata Extraction



December 2006

Stephane Dion Elected Leader of Party



December 2007
Rise of Social Media



April 2008

Fundraising with the Victory Fund/ Fonds de la Victoire



July 2008

The Green Shift Announced!



October 2008

Election Campaign - Advertisement Sites

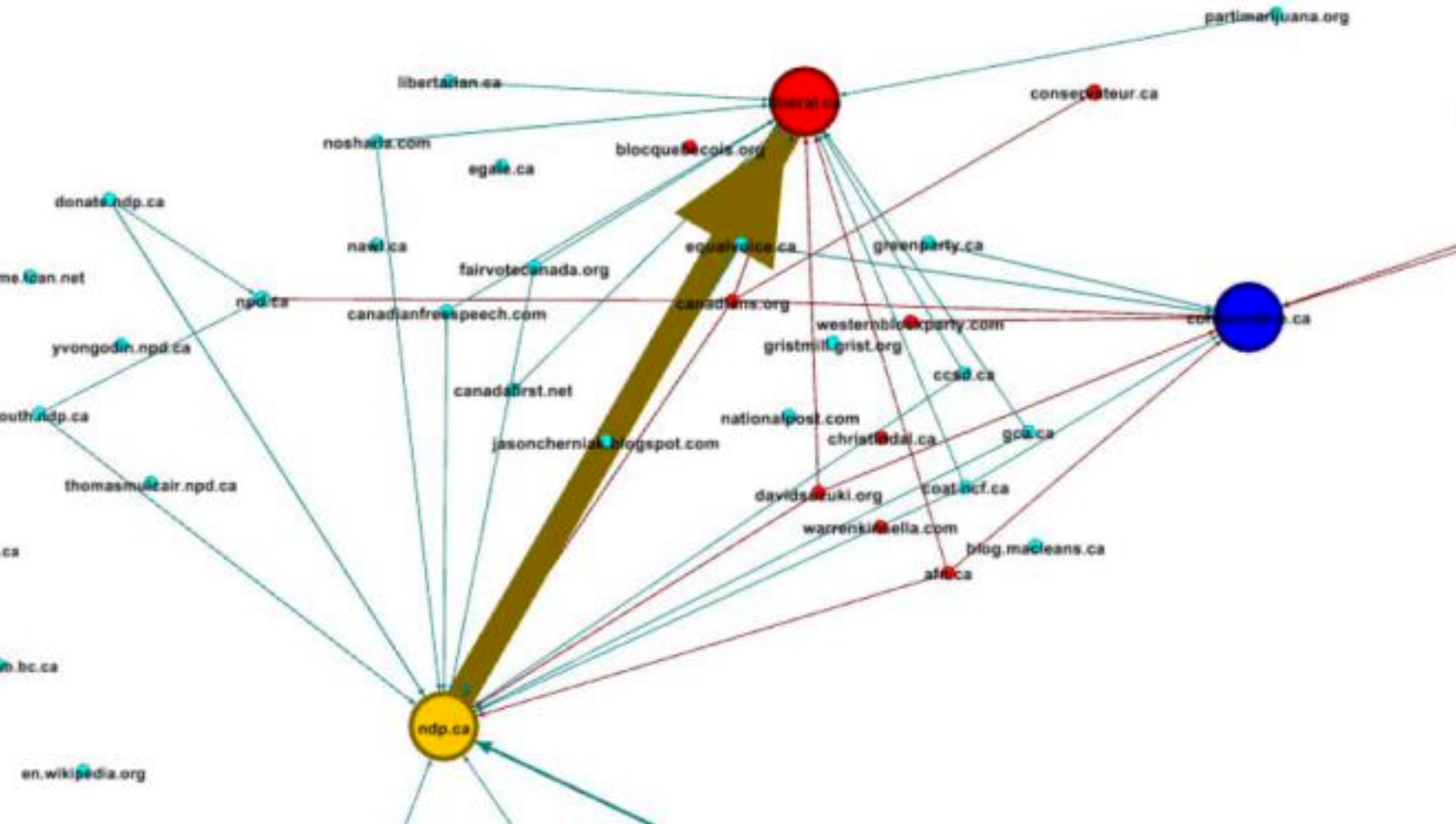


December 2008

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



2005 Canadian Federal Election

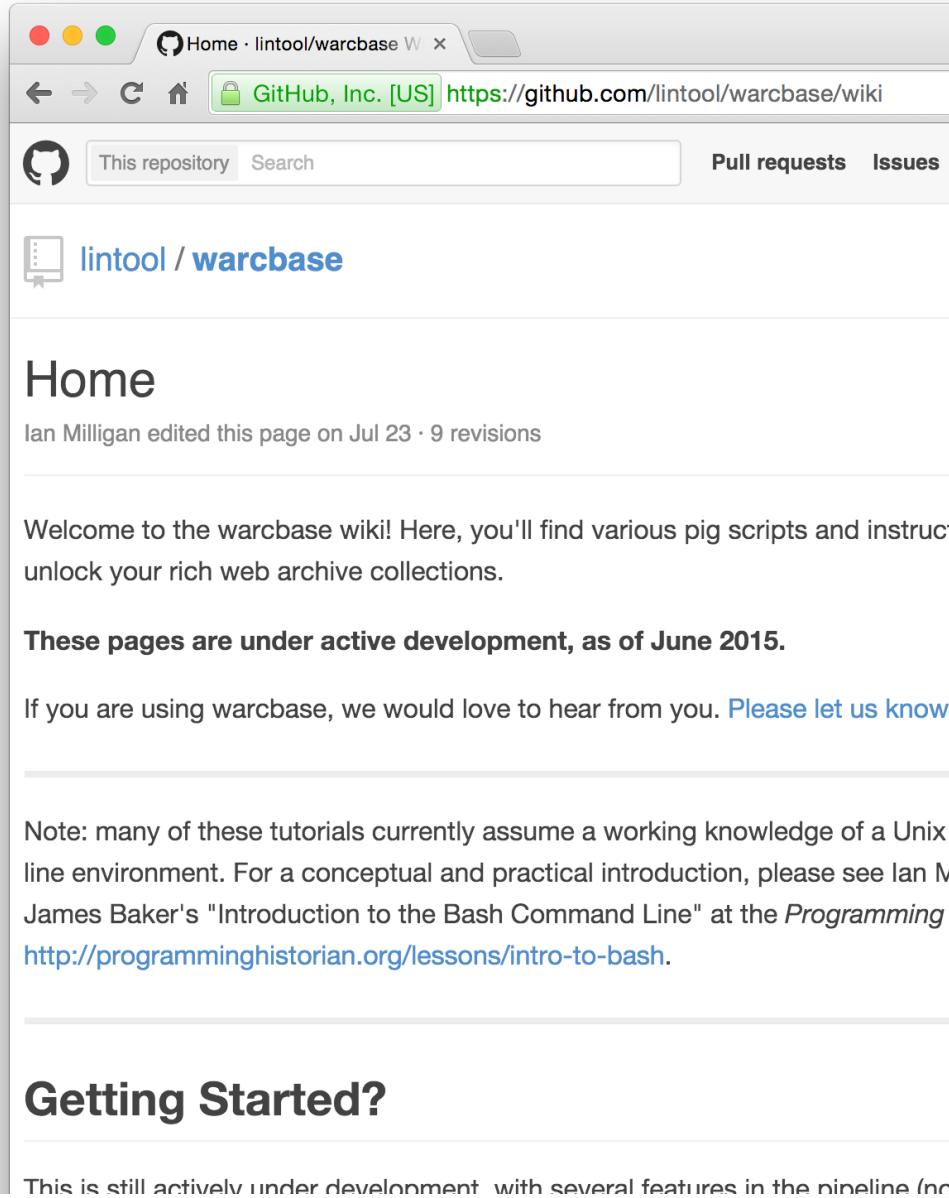


We need to get over our **content fixation** as historians - and
embrace metadata!
(make the two talk to each other)

Interdisciplinarity

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...



The screenshot shows a GitHub wiki page for the 'warcbase' repository. The title bar indicates the page is at <https://github.com/lintool/warcbase/wiki>. The main content area is titled 'Home' and features a welcome message: 'Welcome to the warcbase wiki! Here, you'll find various pig scripts and instructions to unlock your rich web archive collections.' It also states, 'These pages are under active development, as of June 2015.' Below this, there's a note about Unix knowledge and a link to James Baker's 'Introduction to the Bash Command Line'. At the bottom, there's a 'Getting Started?' section and a footer note about the page being actively developed.

Home · lintool/warcbase

GitHub, Inc. [US] <https://github.com/lintool/warcbase/wiki>

This repository Search Pull requests Issues

lintool / warcbase

Home

Ian Milligan edited this page on Jul 23 · 9 revisions

Welcome to the warcbase wiki! Here, you'll find various pig scripts and instructions to unlock your rich web archive collections.

These pages are under active development, as of June 2015.

If you are using warcbase, we would love to hear from you. [Please let us know](#)

Note: many of these tutorials currently assume a working knowledge of a Unix line environment. For a conceptual and practical introduction, please see Ian M. James Baker's "Introduction to the Bash Command Line" at the [Programming Historian](http://programminghistorian.org/lessons/intro-to-bash).

Getting Started?

This is still actively under development with several features in the pipeline (no

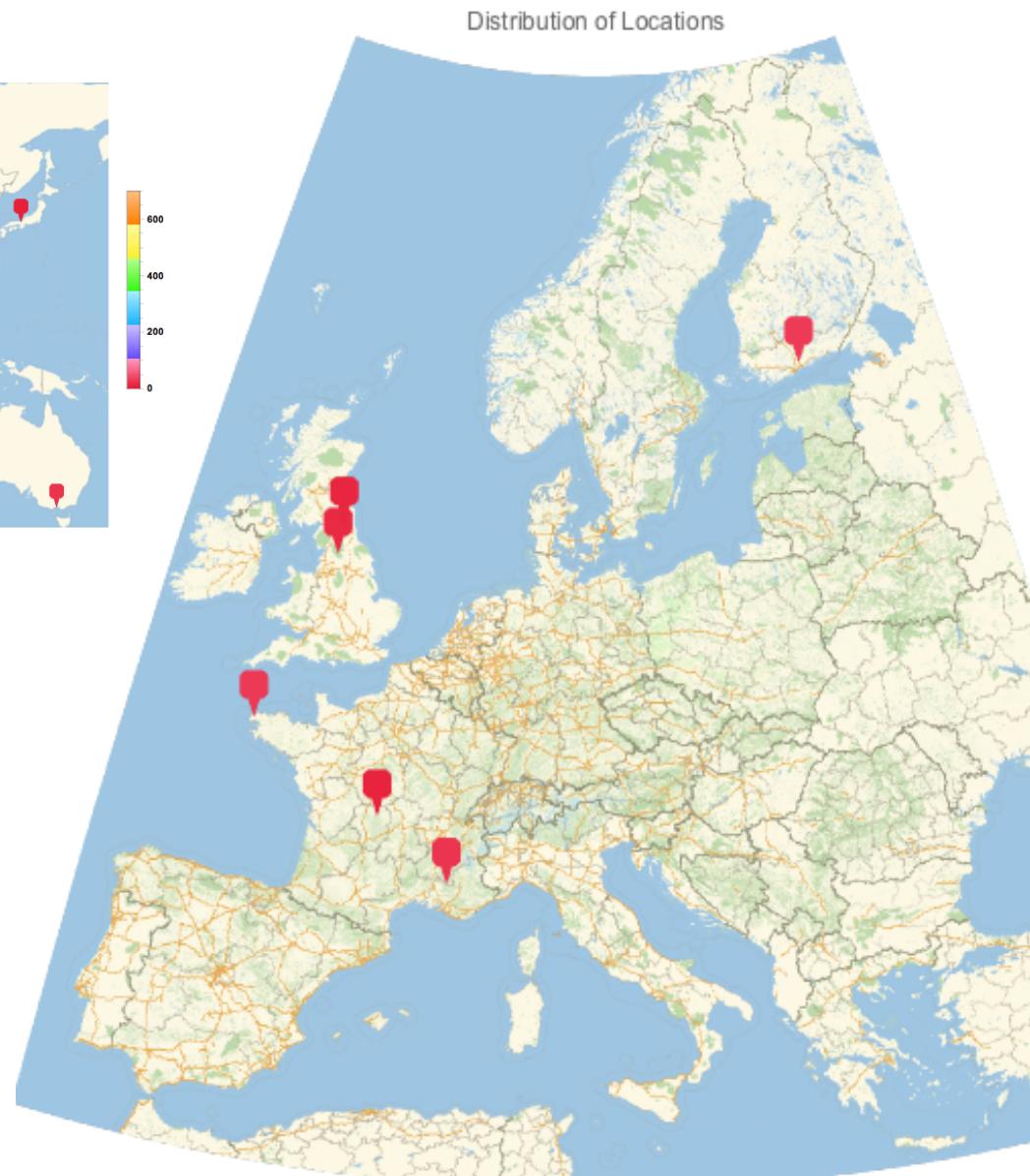
Extract all Plain Text

```
1. i2millig@rho: ~/derivatives/cpp.all.plaintext (ssh)
Python      bash      bash      bash      i2millig@rho:... i2millig@rho:... bash
(20060222,liberal.ca,http://liberal.ca/bio_e.aspx?id=35049,Liberal Party of Canada HOME THE TEAM THE P
ARTY MEDIA CENTRE COMMISSIONS YOUR RIDING      Omar Alghabra www.omaralghabra.com Home > Mississauga--E
rindale Riding Map (PDF) Omar Alghabra came to Canada at a very young age, and immediately knew Canada
was his home. He was first elected 2006 as the Member of Parliament for Mississauga-Erindale. Mr. Al
ghabra is an experienced entrepreneur. For the past six years, he has worked for a large multinational
corporation, carrying out different responsibilities including quality assurance, project management, s
ales, contract management and management of a complete department handling a global mandate. Mr. Alghab
ra is an active member of his community. He is the former National President of the Canadian Arab Feder
ation (2004-2005) and a former member of the Community Editorial Board for the Toronto Star. (2003-2004
). Mr. Alghabra is currently a member of the Diversity Council for General Electric Canada and is activ
e in Junior Achievement for the Toronto Region. He was a member of the Multicultural Inter-Agency of Pe
ople from 2001 to 2002. Mr. Alghabra has a degree in Mechanical Engineering from Ryerson University and a
Masters in Business Administration (MBA) from York University.    Omar Alghabra 790 Burnamthorpe West,
Unit 10 905-276-2806   info@omaralghabra.ca Riding President Elias Hazineh Send an email
Home | News | Your Riding | Contact Us | français This website is the property of the
Liberal Party of Canada and may not be reproduced in whole or in part without express written permission.
© Liberal Party of Canada 2006. All rights reserved. Authorized by the registered agent for the Libe
ral Party of Canada. Privacy Policy)
(20060222,liberal.ca,https://liberal.ca/news_e.aspx?id=11470,Liberal.ca HOME THE TEAM THE PARTY MEDIA C
ENTRE COMMISSIONS YOUR RIDING  Celebrating our National Flag  February 15, 2006 February 15 is Nationa
l Flag Day in Canada, which marks the 41st anniversary of the first raising of the maple leaf flag on P
arliament Hill. Today is a celebration of our shared values, common citizenship and sense of pride in t
his great country we call home. The Canadian flag is one of the most recognizable symbols in the world
and flies proudly to remind us all of who we are and where we come from. The maple leaf's symbolic orig
ins date back to the beginning of our nation's history, while the red and white bars on the flag repres
ent strength and unity. Canada's flag was adopted in 1964 under the courageous leadership of Liberal Pr
ime Minister Lester B. Pearson. The idea of changing the Red Ensign which featured the Britain's Union
Jack, was very controversial at the time, with the Conservative Party strongly opposed to changing the
status quo. Facing strong Conservative resistance in the House of Commons, Pearson's minority governme
nt fought hard in the name of national unity and Canada's multicultural future to make the new flag a r
eality. In an impassioned speech to the House of Commons, Pearson said: "Mr. Speaker, it is for this g
eneration, for this Parliament, to give them and to give us all a common flag; a Canadian flag which, w
hile bringing together but rising above the landmarks and milestones of the past, will say proudly to t
he world and to the future: I stand for Canada." Thanks to Pearson's courageous leadership, Canadians a
cross this great nation celebrate our flag and what it stands for – a country and a citizenship that ar
e the envy of the world.
Home | News | Your Riding | Contact Us | fran
cais This website is the property of the Liberal Party of Canada and may not be reproduced in whole or
```

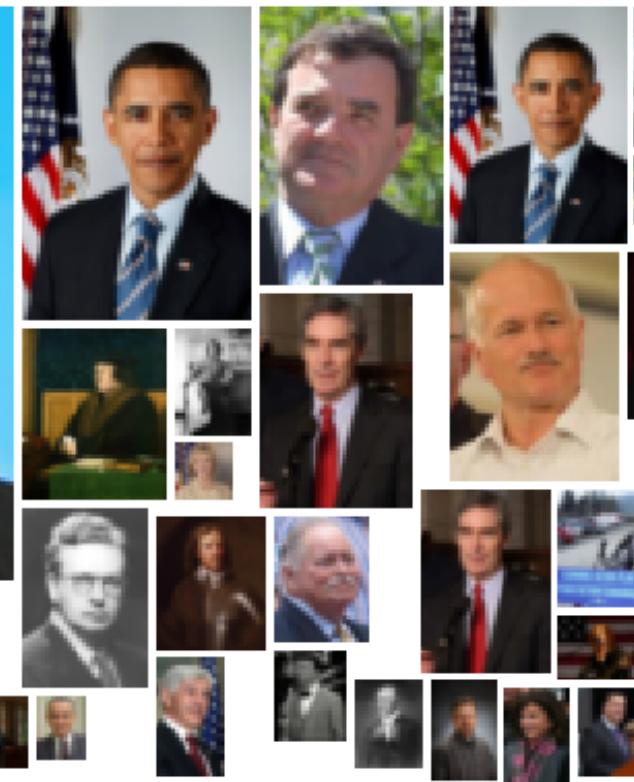

Extract Entities



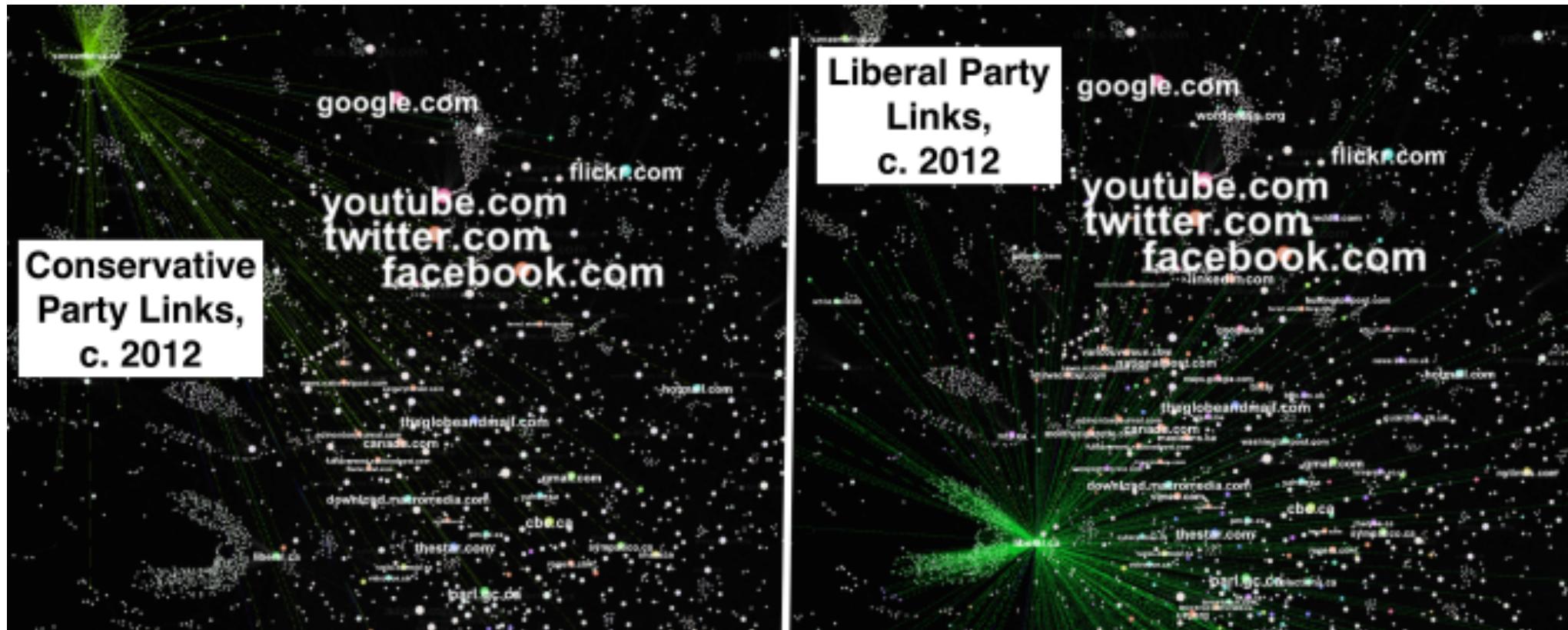
```
In[95]:= int = SemanticInterpretation[#] & /@ processedfreq[[All, 1]]  
Out[95]= {Canada, Calgary, Colombia, Montreal, $Failed, Ontario, Canada, Afghanistan,  
Ottawa, Manitoba, Canada, British Columbia, Canada, Toronto, Nova Scotia, Canada,  
Saskatchewan, Canada, Quebec, Canada, Alberta, Canada, Winnipeg, Washington,  
Canada, Nunavut, Canada, White Horse, United States, Petawawa, Mumbai,  
Lima, The Americas, Saskatoon, Peru, Edmonton, Mexico, Chicoutimi-Jonquiere,  
New Brunswick, Canada, United States, Newfoundland and Labrador, Canada, Qandahar,  
$Failed, Asia-Pacific, Newfoundland and Labrador, Canada, Victoria, United States,  
Quebec City, India, $Failed, Atlantic Ocean, St. Germain, Durham, Battle of Waterloo,
```



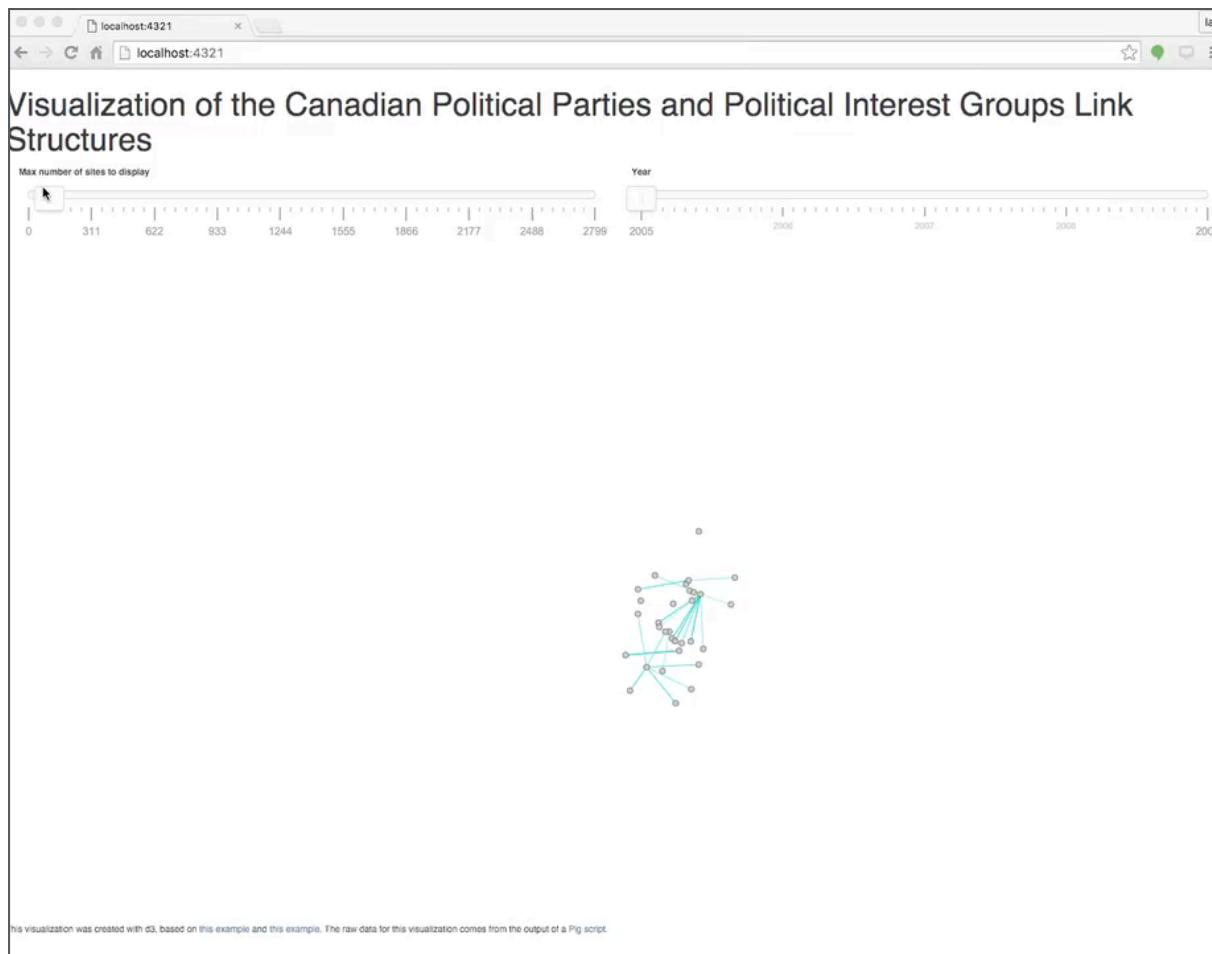
Extract Entities



Extract Links/Gephi Connector



Or D3.js link networks in browser



Bringing it all together in
a notebook environment



Spark Notebook Government Information Day - Demo (unsaved changes) Ian

localhost:9000/notebooks/Government%20Information%20Day%20-%20Demo.snb#

SPARK NOTEBOOK Government Information Day - Demo (unsaved changes)

File Edit View Insert Cell Kernel Help Scala [2.10.4] Spark [1.3.0] Hadoop [2.6.0]

In [1]: :cp /Users/ianmilligan1/dropbox/warcbase/target/warcbase-0.1.0-SNAPSHOT-fatjar.jar

In [2]: import org.warcbase.spark.matchbox._
import org.warcbase.spark.rdd.RecordRDD._

import org.warcbase.spark.matchbox._
import org.warcbase.spark.rdd.RecordRDD._

Out[2]: 161 milliseconds

In [3]: var arc="/Users/ianmilligan1/Dropbox/warcs-workshop/227-20051004191331-00000-crawling015.archive.org.arc.gz"
var warc="/Users/ianmilligan1/dropbox/wahr/sample-data/arc-warc/ARCHIVEIT-227-QUARTERLY-XUGECV-20091218231727-00039-crawling06.us.archive.org-8091.warc.gz"
var armdir="/Users/ianmilligan1/dropbox/warcs-workshop";

arc: String = /Users/ianmilligan1/Dropbox/warcs-workshop/227-20051004191331-00000-crawling015.archive.org.arc.gz
warc: String = /Users/ianmilligan1/dropbox/wahr/sample-data/arc-warc/ARCHIVEIT-227-QUARTERLY-XUGECV-20091218231727-00039-crawling06.us.archive.org-8091.warc.gz
armdir: String = /Users/ianmilligan1/dropbox/warcs-workshop

Out[3]: /Users/ianmilligan1/dropbox/warcs-workshop 961 milliseconds

In [4]: val r =
RecordLoader.loadArc(arc,
sc)
.keepValidPages()
.map(r => ExtractTopLevelDomain(r.getUrl))
....

Walkthroughs at
[https://github.com/lintool/
warcbase/wiki](https://github.com/lintool/warcbase/wiki)

And if you have
suggestions, let us know!

Building connections
between Warcbase and
Shine

Shared Problems

- Never have enough processing power or memory;
- Web archive tools often designed for clusters - less than ten historians in North America probably can use one...
- **Tools**
 - Some work on **WARCs**;
 - Some work on **ARCs**;
 - Some work on **WATs**;
 - And some work on **live-web material**;

Should we effect large-scale non-print legal deposit?

(we've got the power in *Library and Archives Canada Act* (2005), right?)

Overlaps on #elxn42?

Could WAT files be
shared to researchers as
derivative data?

Let's figure it out together!

“Archives Unleashed”
Hackathon

March 3 - 5 2016, University
of Toronto Library

archivesunleashed.ca

Travel funds for grad
students/contingent faculty &
researchers



End-user tools and co-operation with **CS, librarian, and archivist** colleagues is key.

The screenshot shows a GitHub repository page for 'lintool/warcbase'. The repository has 449 commits, 4 branches, and 0 releases. The branch is set to 'master'. A list of recent commits includes:

- .settings: Tweaked settings.
- src: Added option to change MAX_CONTENT_SIZE in IngestFiles, Issues #112
- .gitignore: Added .iml files
- README.md: Error in README
- pom.xml: Updated versions of some artifacts.

Warcbase

Warcbase is an open-source platform for managing web archives built on Hadoop and HBase. The platform provides a flexible data model for storing and managing raw content as well as extracted knowledge. Tight integration with Hadoop provides powerful tools for analysis and data processing.

Getting Started

Clone the repo:

What happens if we don't?

- Do we need to become programmers?
- No - but historians need to know how to speak the lingo - to know the possible, the impossible, and what to strive for
- Not enough just to be friendly - need to work in teams, learn numbers, etc.

Because, as I hope I can
show 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