# Researcher Access for Web Archives

## Our Experiences

**MARAC 2017**
**Buffalo, NY**
**27 October 2017**

**Ian Milligan**
**Associate Professor**
**@ianmilligan1**

**UNIVERSITY OF WATERLOO**
**FACULTY OF ARTS**
Department of History

# Why does a historian care so much about web archives?

Scarcity

Scarcity

Abundance

# 1996

**You can't study the 1990s and beyond without web archives!**

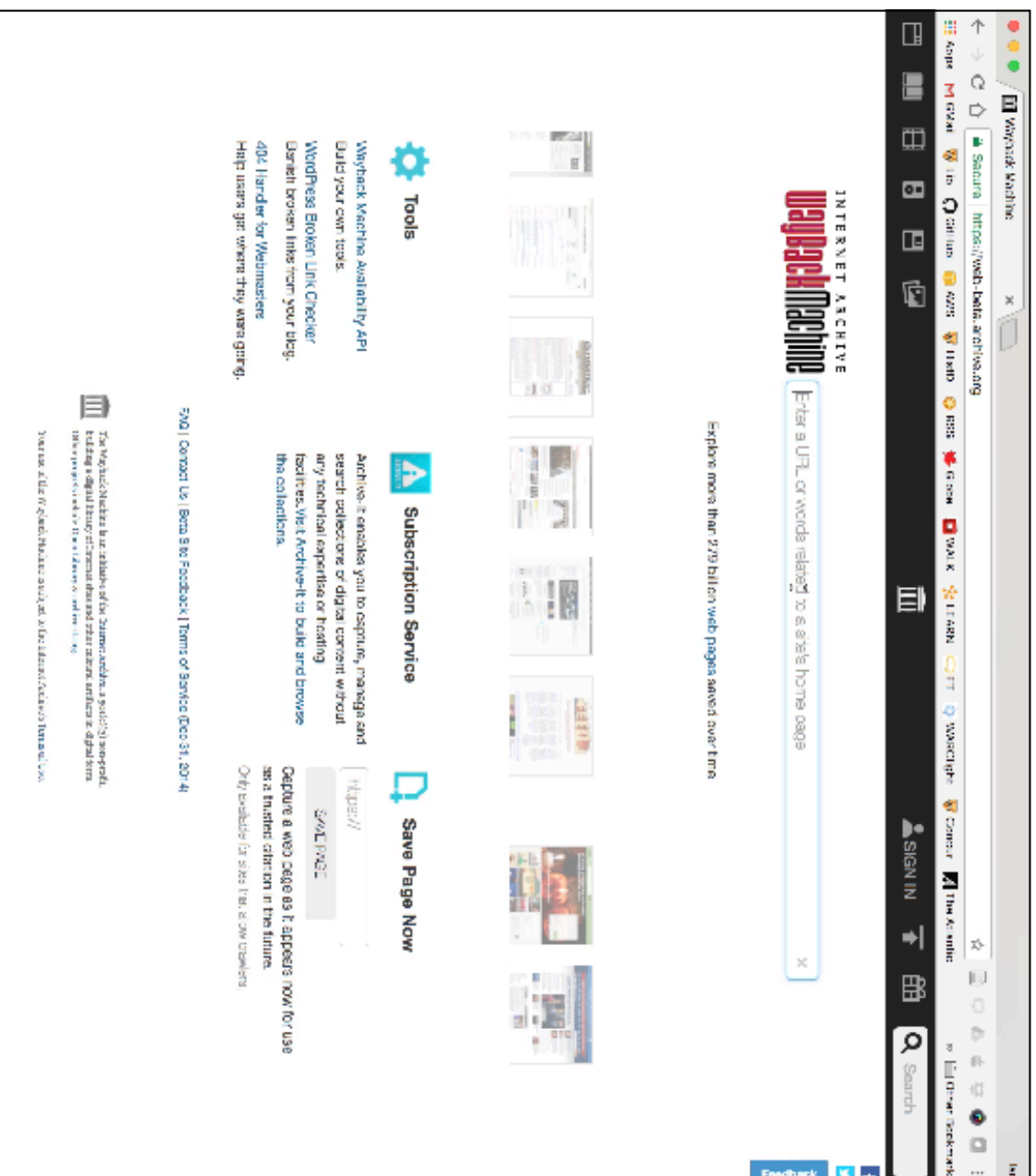**And the 1990s are history. The Web is 25 years old.**

# Current Access Method

This won't be enough!

So what can we do?

# Ask Nich Worby for data

# Canadian Political Parties & Political Interest Group Collection

- 50 Websites

- All major political parties

- Minor political parties

- Political interest groups

- Collected quarterly between 2005 & present.

# How could we provide better access?

**ukwa / shine**

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

⟳ 637 commits    ⌥ 1 branch    ◇ 0 releases    👥 5 contributors

Branch: **master** ▾    **shine** / +

GitHoggarth Added trailing slash to web archive url    Latest commit 11acc26 on Sep 18

| 📁 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

https://www.webarchive.org.uk/shine/graph?query=cat+OR+kitten+OR+moggy+OR+kitty%2C+dog+OR+pu...

Ian

## 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.

cat OR kitten OR moggy OR kitty, dog OR

1996 | 2013 | **Update Graph** | **Reset All**

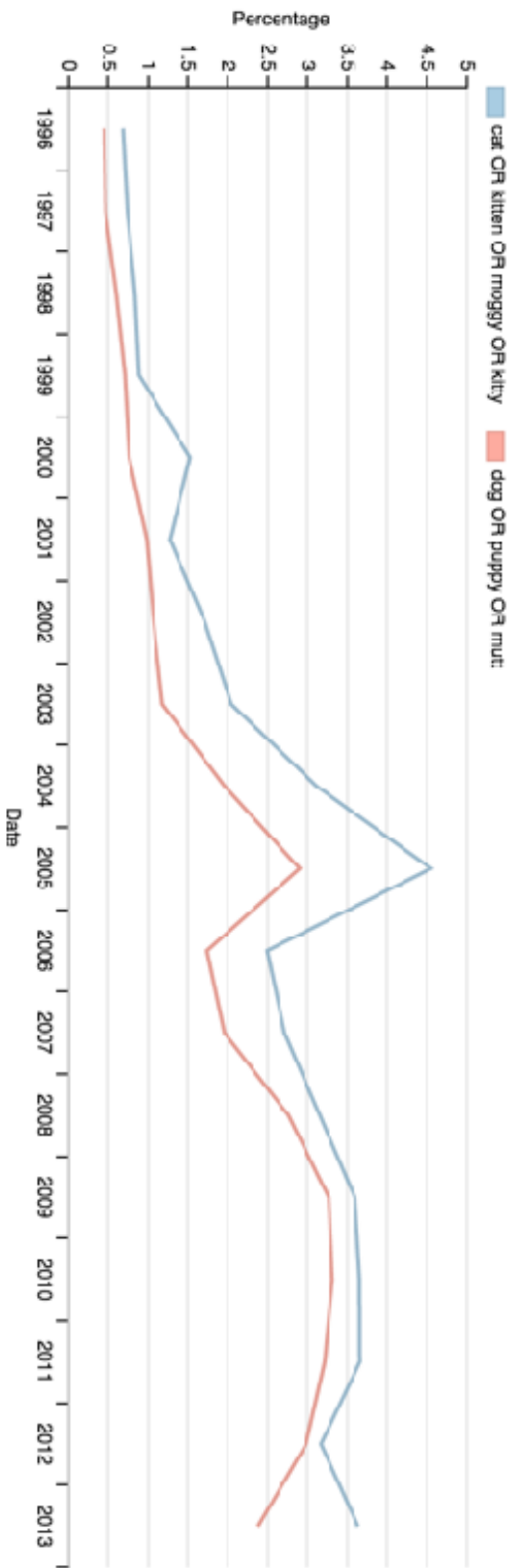- cat OR kitten OR moggy OR kitty
- dog OR puppy OR mut:

Click on a point on the graph to show a random sample of the matching records from that year...



Percentage — Date (1996–2013)

Search    Trends    About    Datasets

Welcome to the Web Archives for Historical Research political parties portal. Before diving in, we encourage you to visit our about page.

# 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.

We also wanted to explore derivative datasets for researcher access with our new project.

# Extract all Text

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 PARTY MEDIA CENTRE COMMISSIONS YOUR RIDING    Omar Alghabra www.omaralghabra.com Home > Mississauga--Erindale Riding Map (PDF) Omar Alghabra came to Canada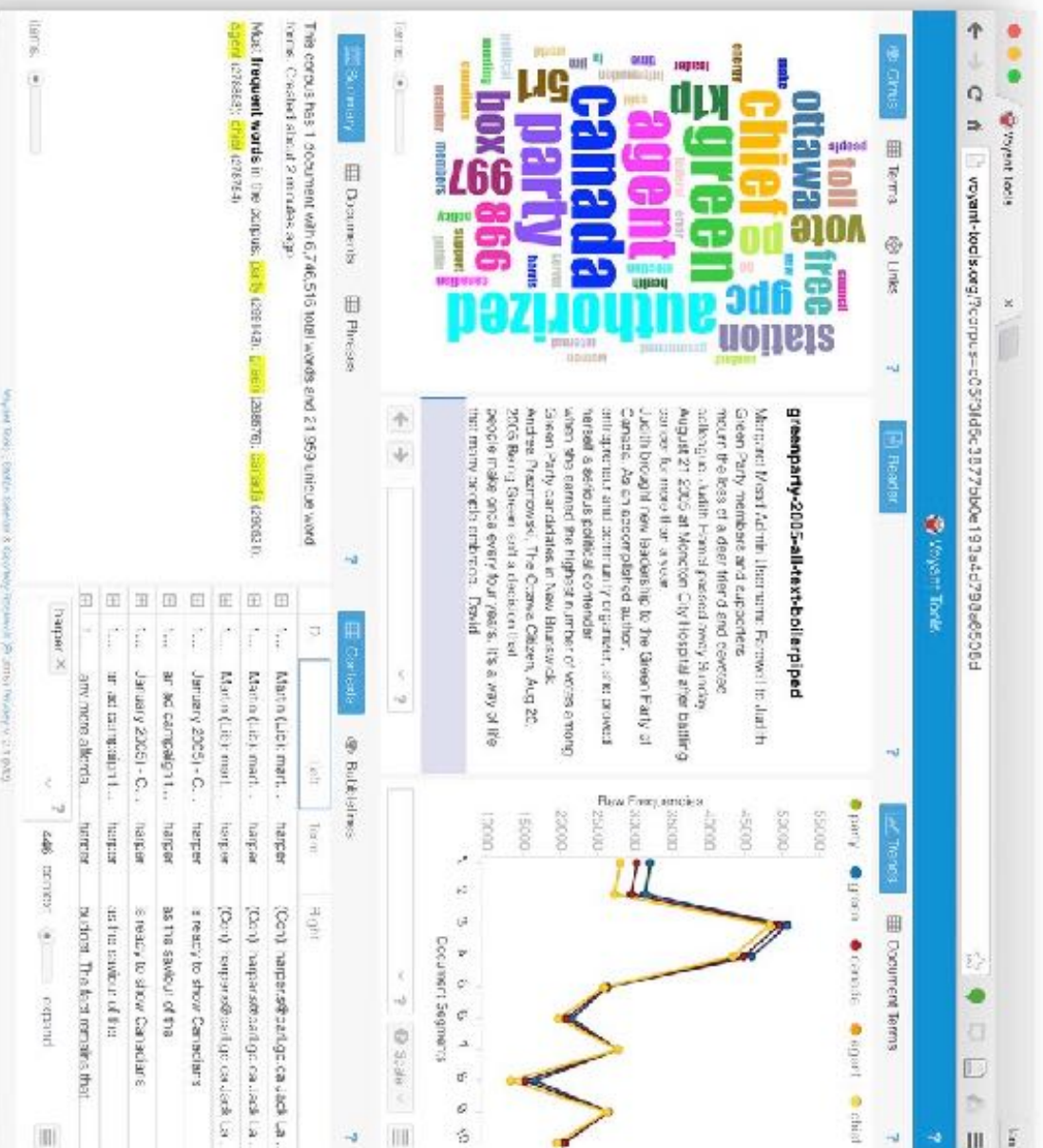 at a very young age, and immediately knew Canada was his home. He is was first elected 2006 as the Member of Parliament for Mississauga-Erindale. Mr. Alghabra is an experienced entrepreneur. For the past six years, he has worked for a large multinational corporation, carrying out different responsibilities including quality assurance, project management, sales, contract management and management of a complete department handling a global mandate. Mr. Alghabra is an active member of his community. He is the former National President of the Canadian Arab Federation (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 active in Junior Achievement for the Toronto Region. He was a member of the Multicultural Inter-Agency of Peel 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 Burnhamthorpe West, Unit 10 905-276-2806

Home | News | Your Riding | Contact Us | français This website is the property of the
info@omaralghabra.ca Riding President Elias Hazineh Send an email
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 Liberal Party of Canada. Privacy Policy

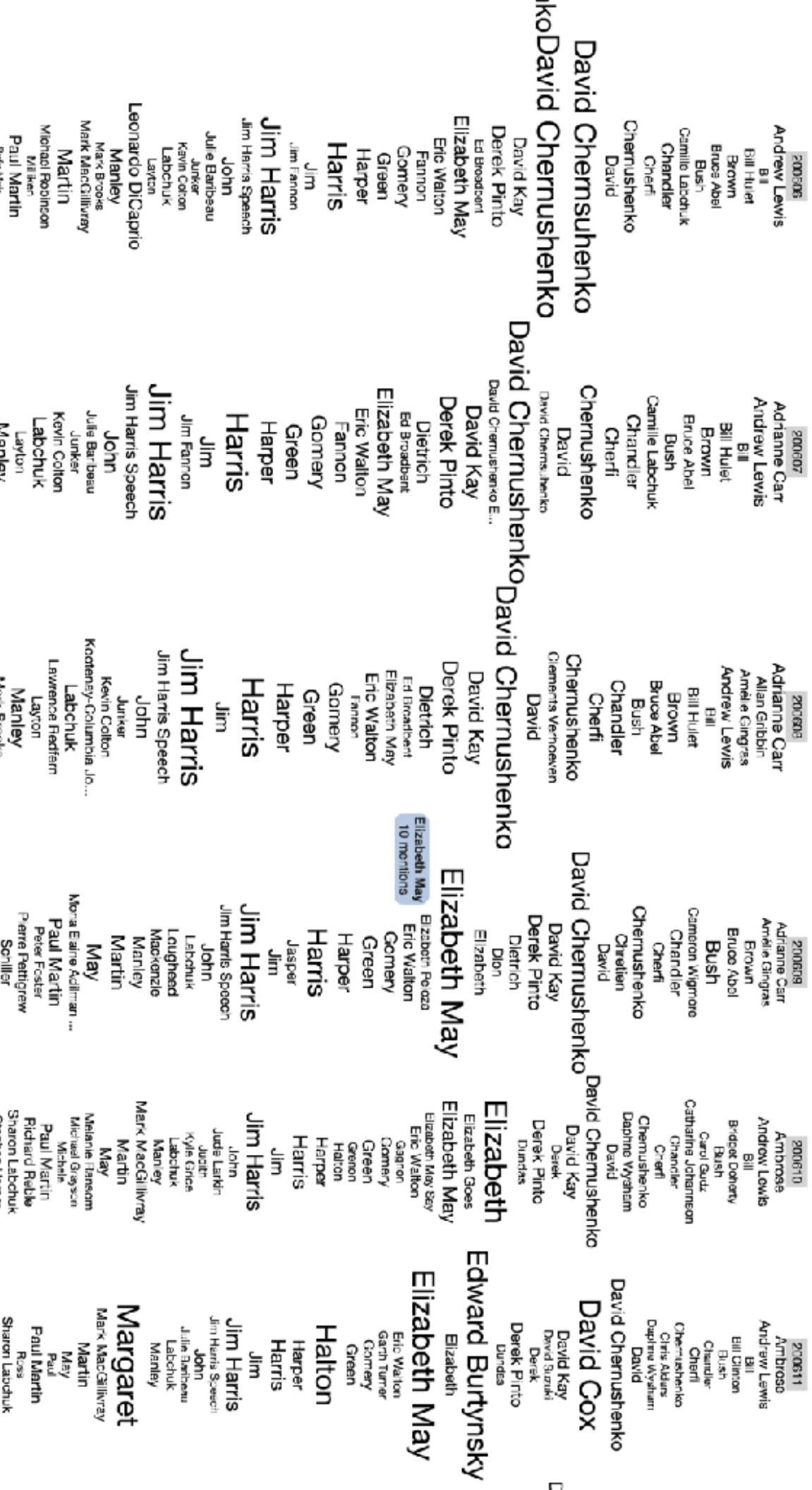(20060222,liberal.ca,https://liberal.ca/news_e.aspx?id=11470,Liberal Party of Canada HOME THE TEAM THE PARTY MEDIA CENTRE COMMISSIONS YOUR RIDING    Celebrating our National Flag February 15, 2006 February 15 is National Flag Day in Canada, which marks the 41st anniversary of the first raising of the maple leaf flag on Parliament Hill. Today is a celebration of our shared values, common citizenship and sense of pride in the 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 origins date back to the beginning of our nation's history, while the red and white bars on the flag represent strength and unity. Canada's flag was adopted in 1964 under the courageous leadership of Liberal Prime 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 Government fought hard in the name of national unity and Canada's multicultural future to make the new flag a reality. In an impassioned speech to the House of Commons, Pearson said: "Mr. Speaker, it is for this generation, for this Parliament, to give them and to give us all a common flag; a Canadian flag which, while bringing together but rising above the landmarks and milestones of the past, will say proudly to the world and to the future: I stand for Canada." Thanks to Pearson's courageous leadership, Canadians across this great nation celebrate our flag and what it stands for — a country and a citizenship that are the envy of the world.
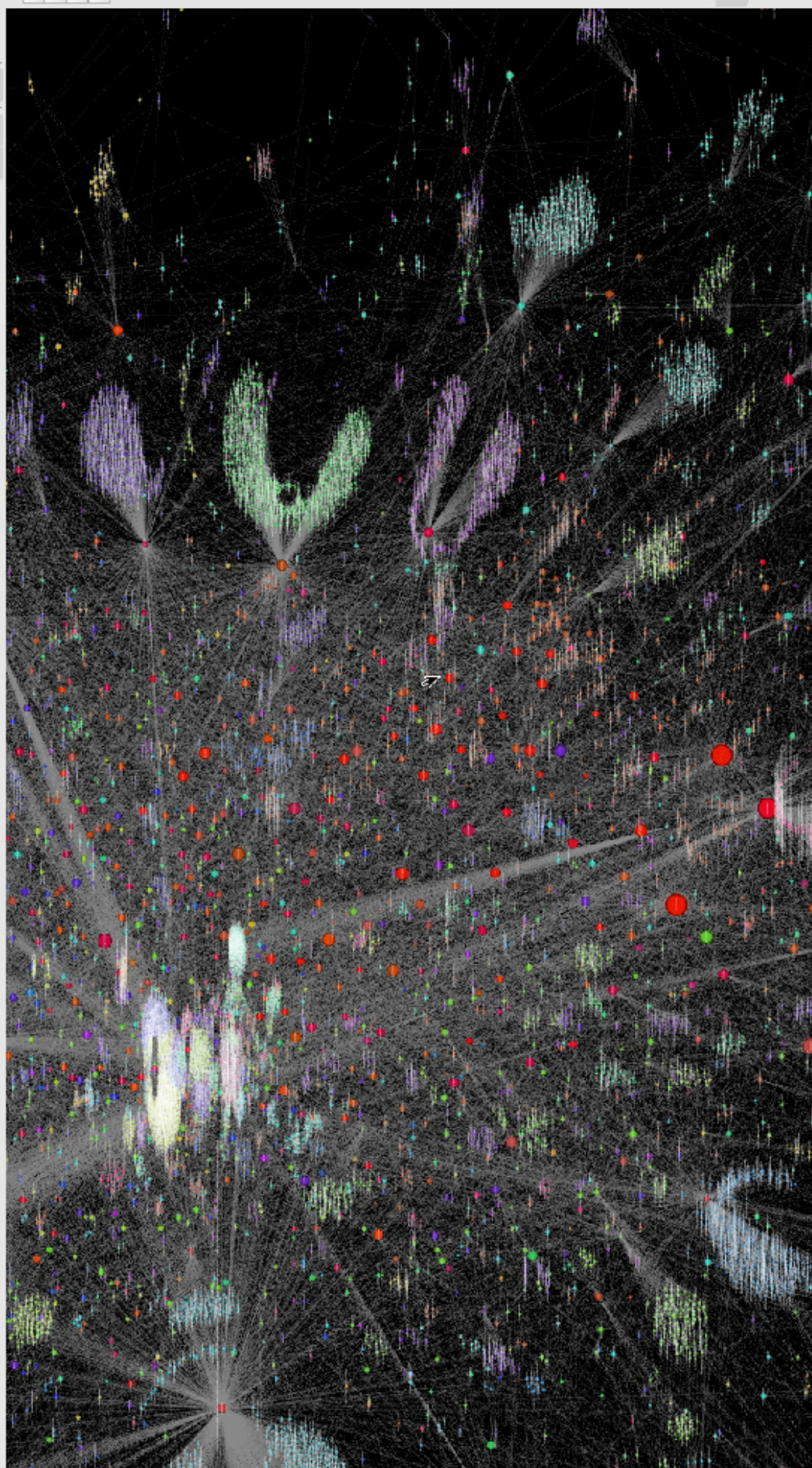
Home | News | Your Riding | Contact Us | fran...
als.This website is the property of the Liberal Party of Canada and may not be reproduced in whole or

Extract all Text

# Extract Entities

Andrew Lewis
Bill
Bill Hulet
Brown
Bruce Abel
Bush
Camille Labchuk
Chandler
Cherfi
Chernushenko
David

David Chernsuhenko
ikoDavid Chernushenko

David
David Kay
Derek Pinto
Elizabeth May
Ed Broadcert
Eric Walton
Fannon
Gomery
Green
Harper
Harris

Jim
Jim Fannon

Jim Harris
Jim Harris Speech
John
Julie Baribeau
Junker
Kevin Colton
Labchuk
Layton
Leonardo DiCaprio
Manley
Mark Brooks
Mark MacGillivray
Martin
Michael Robinson
Miller
Paul Martin

Adrianne Carr
Andrew Lewis
Bill
Bill Hulet
Brown
Bruce Abel
Bush
Camille Labchuk
Chandler
Cherfi

Chernushenko
David

David Chernushenkc
Clements Verhoeven
David

David Kay
Derek Pinto
Dietrich
Ed Broadbent
Elizabeth May
Eric Walton
Fannon
Gomery
Green
Harper
Harris

Jim
Jim Fannon

Jim Harris
Jim Harris Speech
John
Julie Baribeau
Junker
Kevin Colton
Labchuk
Layton
Manley

David Chernushenko E...
David Kay
Derek Pinto

Adrianne Carr
Allan Gribbin
Amélie Gingras
Andrew Lewis
Bill
Bill Hulet
Brown
Bruce Abel
Bush
Chandler
Cherfi

Chernushenko
David

David Chernushenko

David Kay
Derek Pinto
Dietrich
Ed Broadbent
Elizabeth May
Eric Walton
Fannon
Gomery
Green
Harper
Harris

Jim
Jim Harris
Jim Harris Speech
John
Junker
Kootenay-Columbia Jo...
Kevin Colton
Labchuk
Lawrence Redfern
Layton
Manley

Elizabeth May
10 mentions

Elizabeth Peroza
Eric Walton
Gomery
Green
Harper
Harris
Jasper
Jim

Elizabeth May

David Chernushenko
David Kay
Derek Pinto
Dietrich
Dillon
Elizabeth

Adrianne Carr
Amélie Gingras
Brown
Bruce Abel
Bush
Cameron Wigmore
Chandler
Cherfi
Chernushenko
Chretien
David

Jim Harris
Jim Harris Speech
John
Labchuk
Lougheed
Mackenzie
Manley
Martin

Nora Elaine Adkiman ...
Paul Martin
May
Peter Foster
Pierre Pettigrew
Schiller

Edward Burtynsky
Elizabeth

Elizabeth
Elizabeth Goes
Elizabeth May

Elizabeth May Say
Eric Walton
Gagnon
Gomery
Green
Grenon
Halton
Harper
Harris
Jim

Jim Harris
Jim Harris Speech
John
Jude Larkin
Justin
Kyle Grice
Labchuk
Manley

Halton
Harper
Harris

David Chernushenko
David Kay
Derek
Derek Pinto
Dundas

Adrianne Carr
Amélie Gingras
Bush
Bridget Doherty
Carol Gudz
Catharine Johanneon
Chandler
Cherfi
Chernushenko
Chris Alders
Daphne Wysham
David

Margaret
Mark MacGillivray
Martin

Mark MacGillivray
Martin
May
Melanie Ransom
Michael Grayson
Michelle
Paul Martin
Richard Reble
Sharon Labchuk

David Cox

David Chernushenko
David Kay
David Suzuki
Derek
Derek Pinto
Dundas

Andrew Lewis
Ambrose
Bill
Bill Clinton
Bush
Chandler
Cherfi
Chernushenko
David Suzuki
Daphne Wysham
David

Eric Walton
Garth Turner
Gomery
Green

Jim Harris
Jim Harris Speech
John
Julie Baribeau
Labchuk
Manley

Martin
May
Paul
Paul Martin
Ross
Sharon Labchuk

# Next Step

- The **Archives Unleashed Cloud**

- Under development - moving Shine to a Blacklight-based system for community involvement;

- A GUI-based workbench for uploading WARCs or using Internet Archive APIs to extract derivatives from your collections to do text analysis, networks, other DH work.

- **Set of regional events and workshops - with travel funding - stay tuned!**

Scholars can directly use the
Archives Unleashed Toolkit
using local resources of via
the Archives Unleased Cloud

Scholars can download
results, derived data,
etc. for additional local
analyses if desired

Data ingestion from Archive-It,
scholar's local collections, etc.

**Altiscale**

**Compute Canada**

Local
Spark Cluster
(can be
single-node)

Archives
Unleashed
Toolkit

Archives
Unleashed
Cloud
Portal

Archives
Unleashed
Toolkit

Spark Cluster

Search portal
(Project Blacklight)

Visit us at archivesunleashed.org

ARCHIVES UNLEASHED PROJECT

AUT PROJECT

Our goal is to make petabytes of historical internet content accessible to scholars and others interested in researching the recent past.

ABOUT    SOFTWARE    GET INVOLVED    PROJECT TEAM    ADVISORY BOARD    CONTACT US

# Thanks very much!

**Ian Milligan**
Associate Professor
@ianmilligan1

UNIVERSITY OF WATERLOO
FACULTY OF ARTS
Department of History