

# Between Metadata and Content

**Exploring Canadian Political History with Archive-It's  
Research Services**

---

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



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

**Jimmy Lin**  
Professor and David R. Cheriton Chair  
[@lintool](https://twitter.com/lintool)



**UNIVERSITY OF WATERLOO**  
FACULTY OF MATHEMATICS  
David R. Cheriton School  
of Computer Science

**Jeremy Wiebe**  
PhD Candidate  
[@jeremyw](https://twitter.com/jeremyw)



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

199

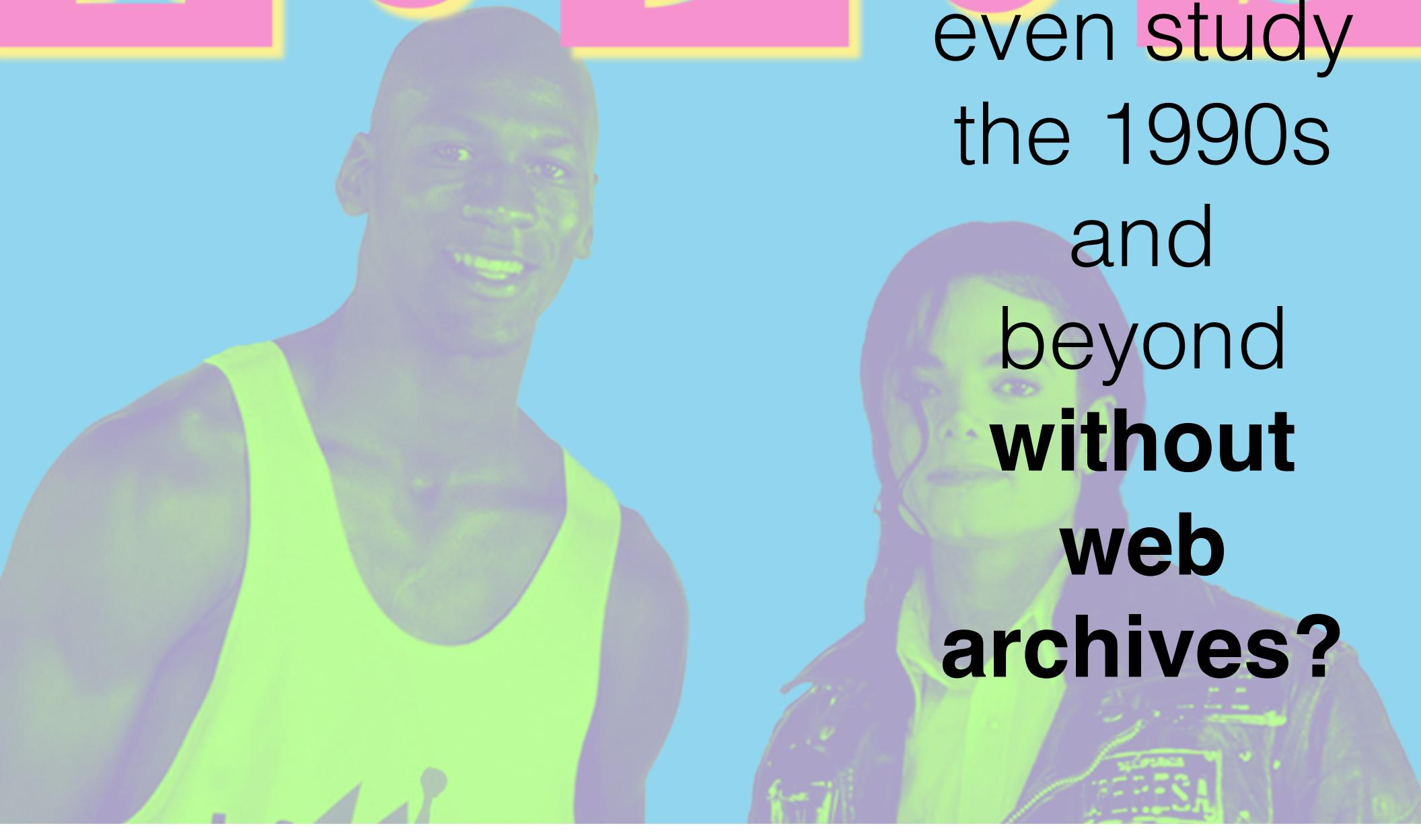
99

99

0

S

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



# The Case Study

- **Archive-It Research Services:** “Canadian Political Parties and Political Interest Groups”
- 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 banner features the Archive-It logo and the text "Canadian Political Parties and Political Interest Groups" and "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 "Subject", "Sort By: Count (A-Z)", and "New Democratic Party of Canada (2)", "Assembly of First Nations (1)", etc. There is a search bar at the bottom right with the placeholder "Enter search terms here". The footer indicates "Page 1 of 1 (54 Total)" and "Sort By: Title (A-Z) | Title (Z-A) | URL (A-Z) | URL (Z-A)".

# 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!
- <https://archive-it.org/collections/227>

The screenshot shows a web browser window displaying the Archive-It collection 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 header includes 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.

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 note about the collector.

Below the collection summary, there is a search bar and search results for "Stephen Harper". The search results page shows a total of 60,657 results. The results list includes a link to "Stephen Harper | Facebook" with a URL of <http://www.facebook.com/pages/Stephen-Harper/9106562109>. The page also features filters for search terms, host, and file format, and a "Sort By" dropdown set to "Best Match".

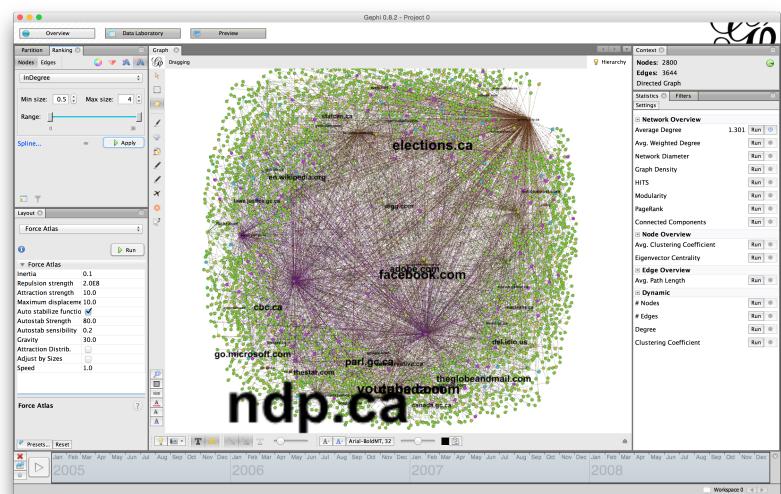
How to provide  
access?

# Warcbase

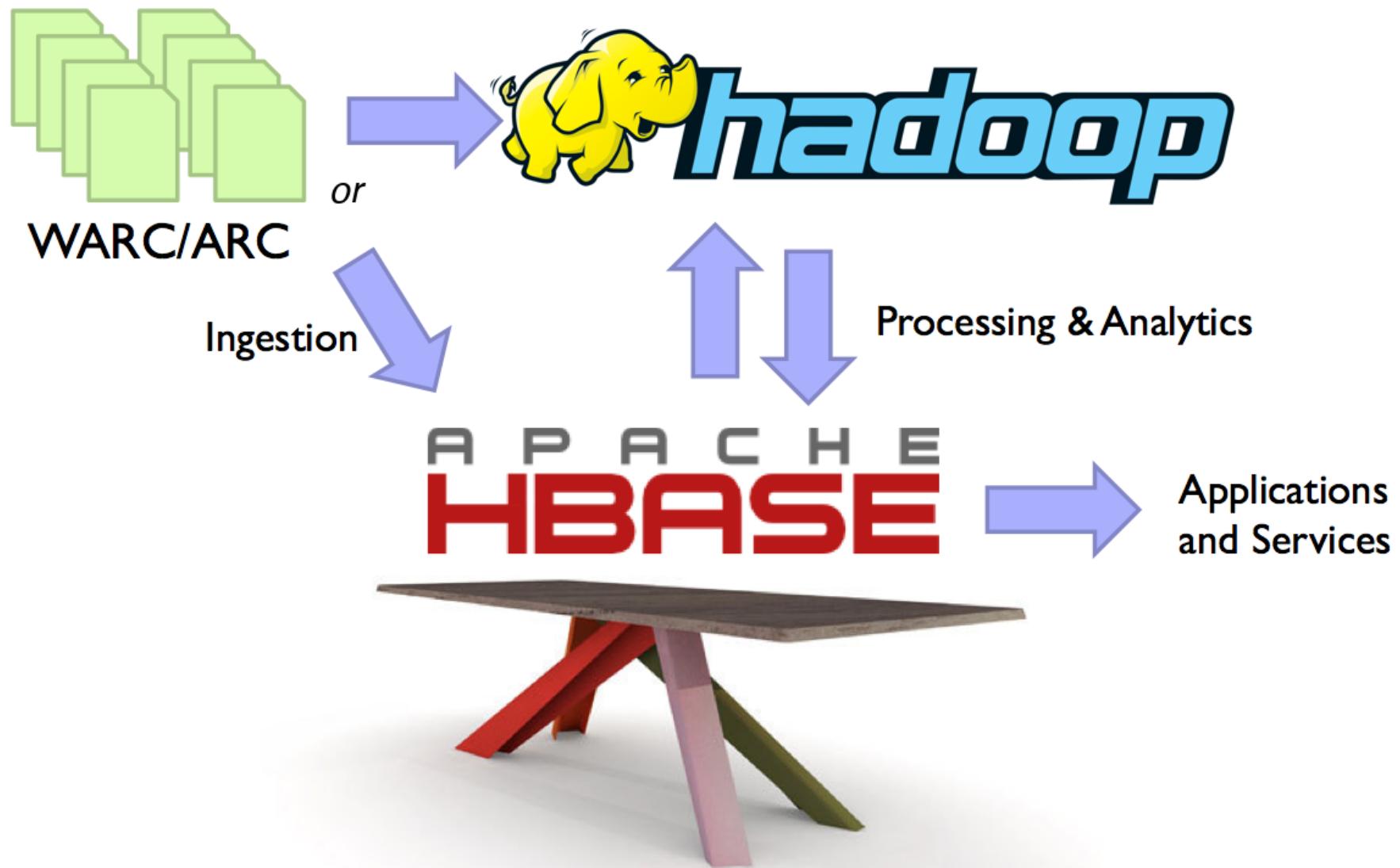
An open-source platform for managing web archives  
<http://warcbase.org>

Two main facets

- A flexible data store: your own Wayback Machine
- Scriptable analytics and data processing



# Warcbase



# Warcbase

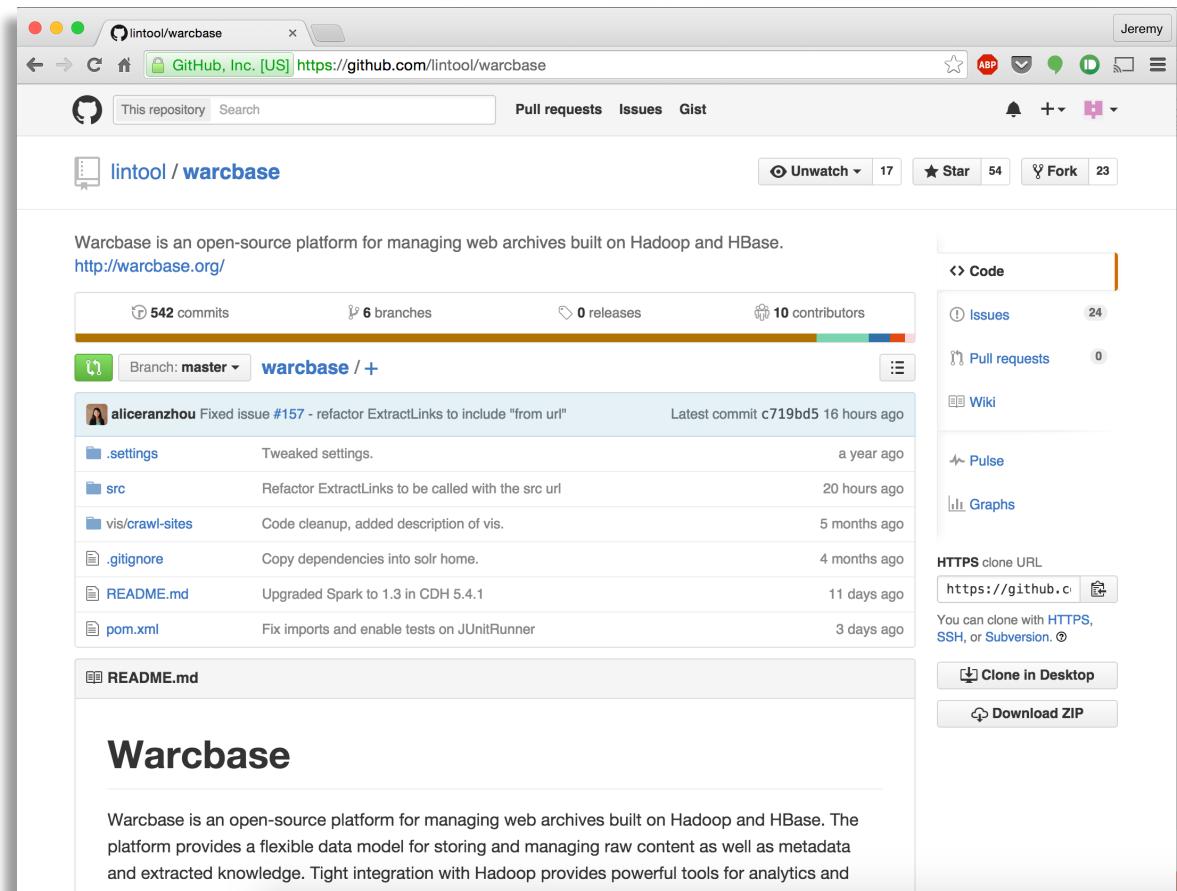
- Framework for distributed storage and distributed processing of very big data
- Scalable
- Potentially very powerful
  - *Trantor*: 1.2PB of disk, 25 compute nodes (each w/ 128GB memory, 2×6-core Intel Xeon E5 v3 = 3.2TB memory and 300 current-generation Intel cores)
- In active development, led by Jimmy Lin, collaborator with Web Archives Historical Research Group



# You Can Warcbase Too!

<http://warcbase.org>

- Flexible requirements:  
designed for powerful  
distributed systems,  
but will also run on  
your laptop
  - Linux or OS X
- Comprehensive  
documentation, with  
script examples



# Analysis: Scripting

- Previously, Pig scripts ran Hadoop MapReduce jobs
  - Deprecation imminent
- Currently transitioning to Spark, scripts written in Scala
  - Better performance
  - API and libraries for Python (*PySpark*) – hope to facilitate adoption



# A Simple Script

## Domain level plain text extraction (Spark)

```
import org.warcbase.spark.matchbox.ArcRecords
import org.warcbase.spark.matchbox.ArcRecords._

val r = ArcRecords.load("src/test/resources/arc/example.arc.gz", sc)
  .keepMimeTypes(Set("text/html"))
  .discardDate(null)
  .keepDomains(Set("greenparty.ca"))
  .extractDomainUrlBody()

r.saveAsTextFile("out/")
```

# Named Entity Recognition

## Some interesting results

- Stanford Named Entity Recognition library
- Call by script
- A useful discovery tool

```
%default I_PARSED_DATA_DIR '/user/jrwiebe/cpp.text-greenparty/*.txt';
%default O_ENTITIES_DIR '/user/jrwiebe/cpp.text-greenparty/entities.gz/';
%default I_NER_CLASSIFIER_FILE 'english.all.3class.distsim.crf.ser.gz';

SET mapred.max.map.failures.percent 10;
SET mapred.reduce.slowstart.completed.maps 0.9

REGISTER 'target/warcbase-0.1.0-SNAPSHOT-fatjar.jar';

DEFINE NER3CLASS org.warcbase.pig.piggybank.NER3ClassUDF('$I_NER_CLASSIFIER_FILE');

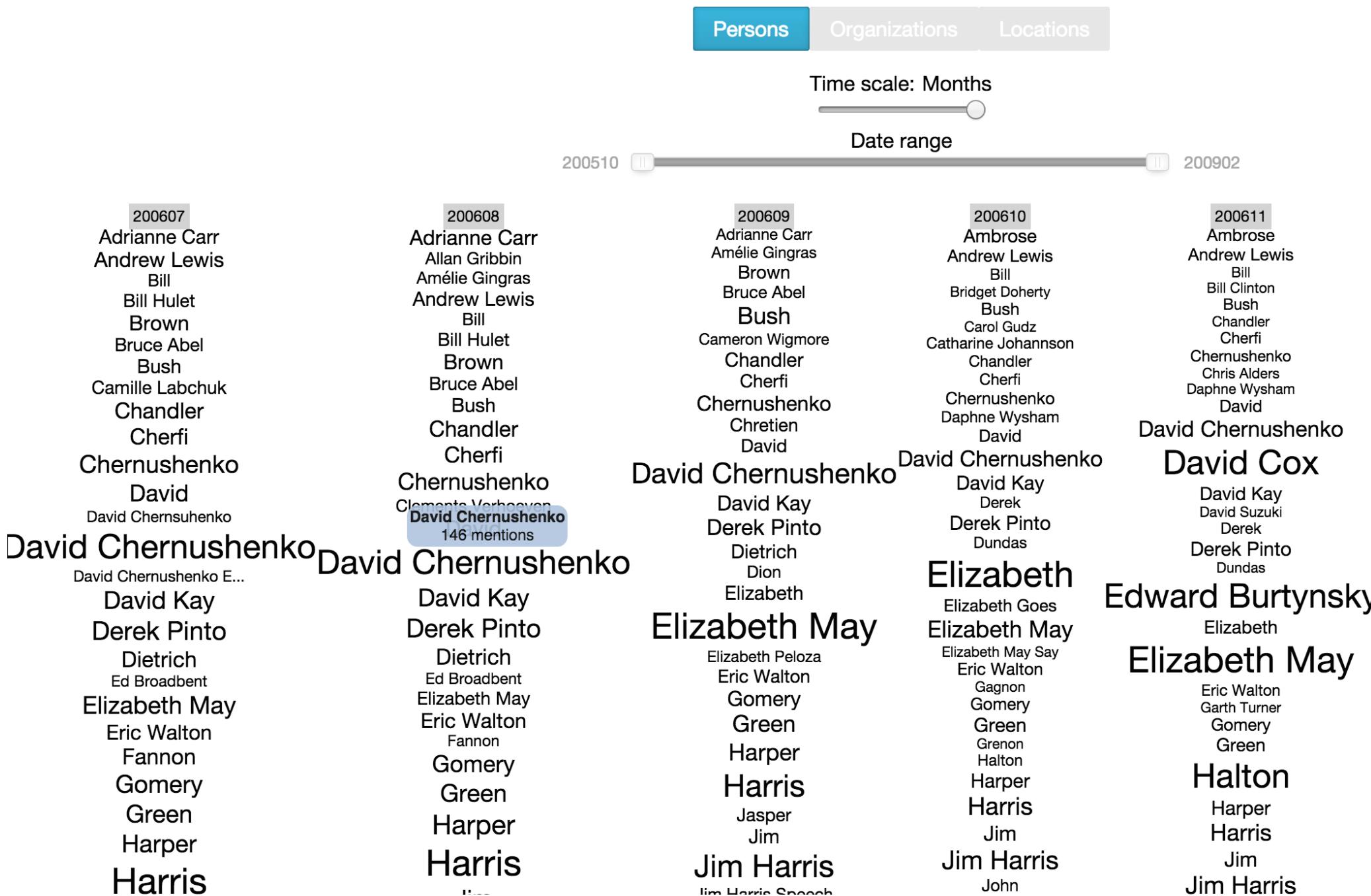
Scrapes = LOAD '$I_PARSED_DATA_DIR' AS (date:chararray, url:chararray, content:chararray);

Entities = FOREACH Scrapes GENERATE date, url, NER3CLASS(content) AS entityString;

STORE Entities into '$O_ENTITIES_DIR';
```

# Named Entity Visualization

Data source: [greenparty.csv](#)





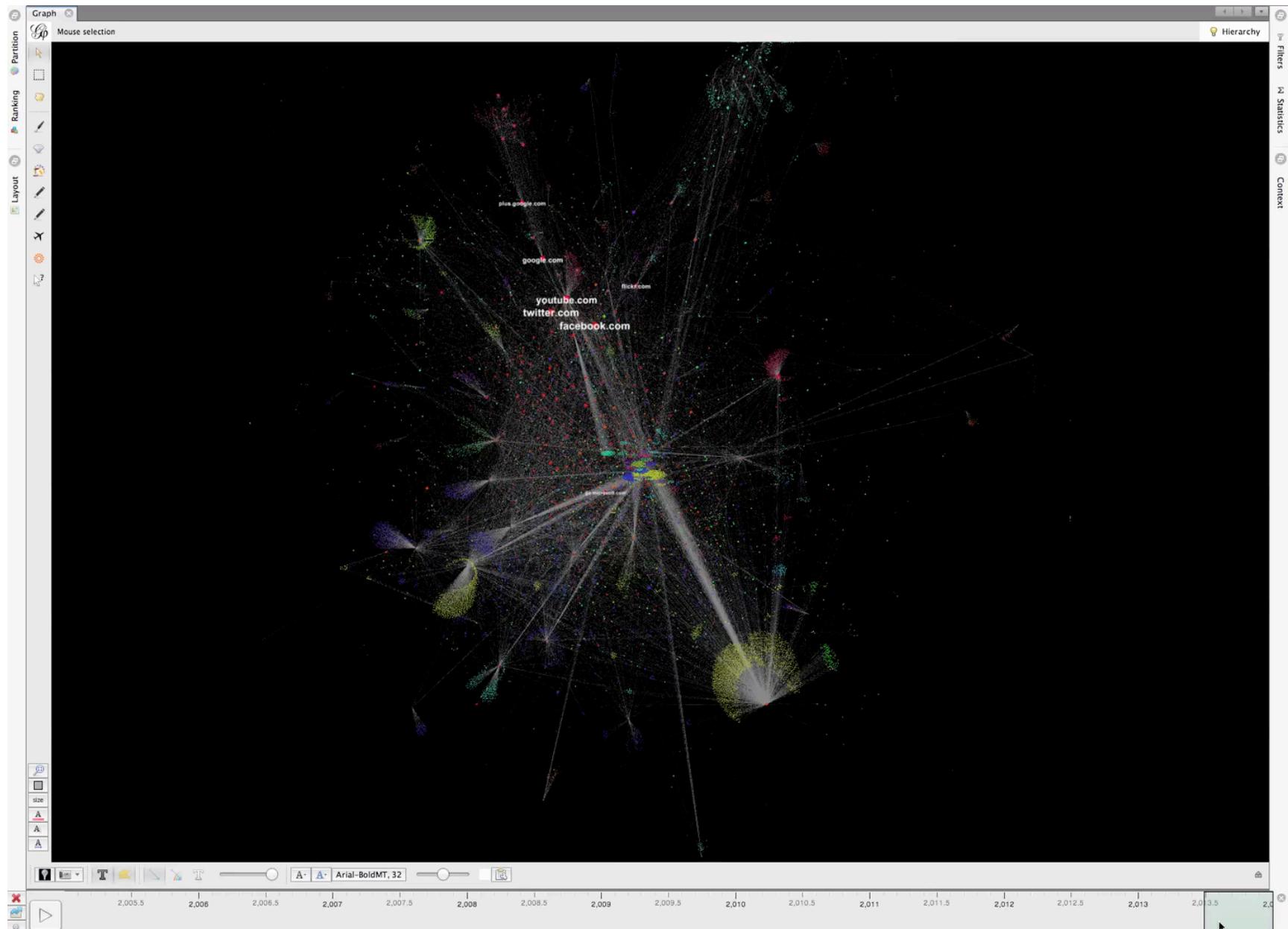


# **Historians want content**

- but metadata might  
actually be more useful!

So here's some cool stuff  
we can do with  
warcbase.

# Metadata Extraction



Short stories from  
metadata as well

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

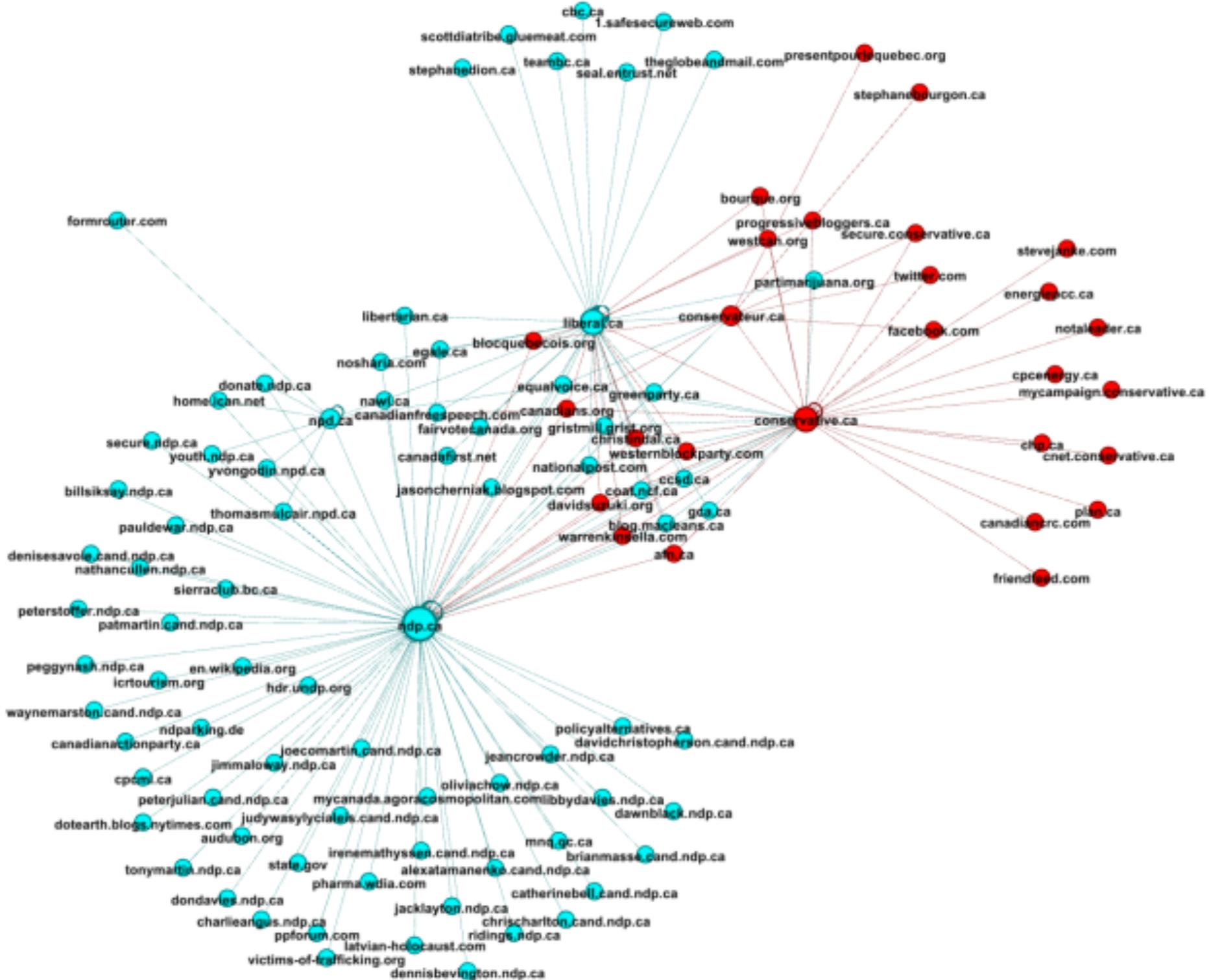
## Election Campaign - Advertisement Sites



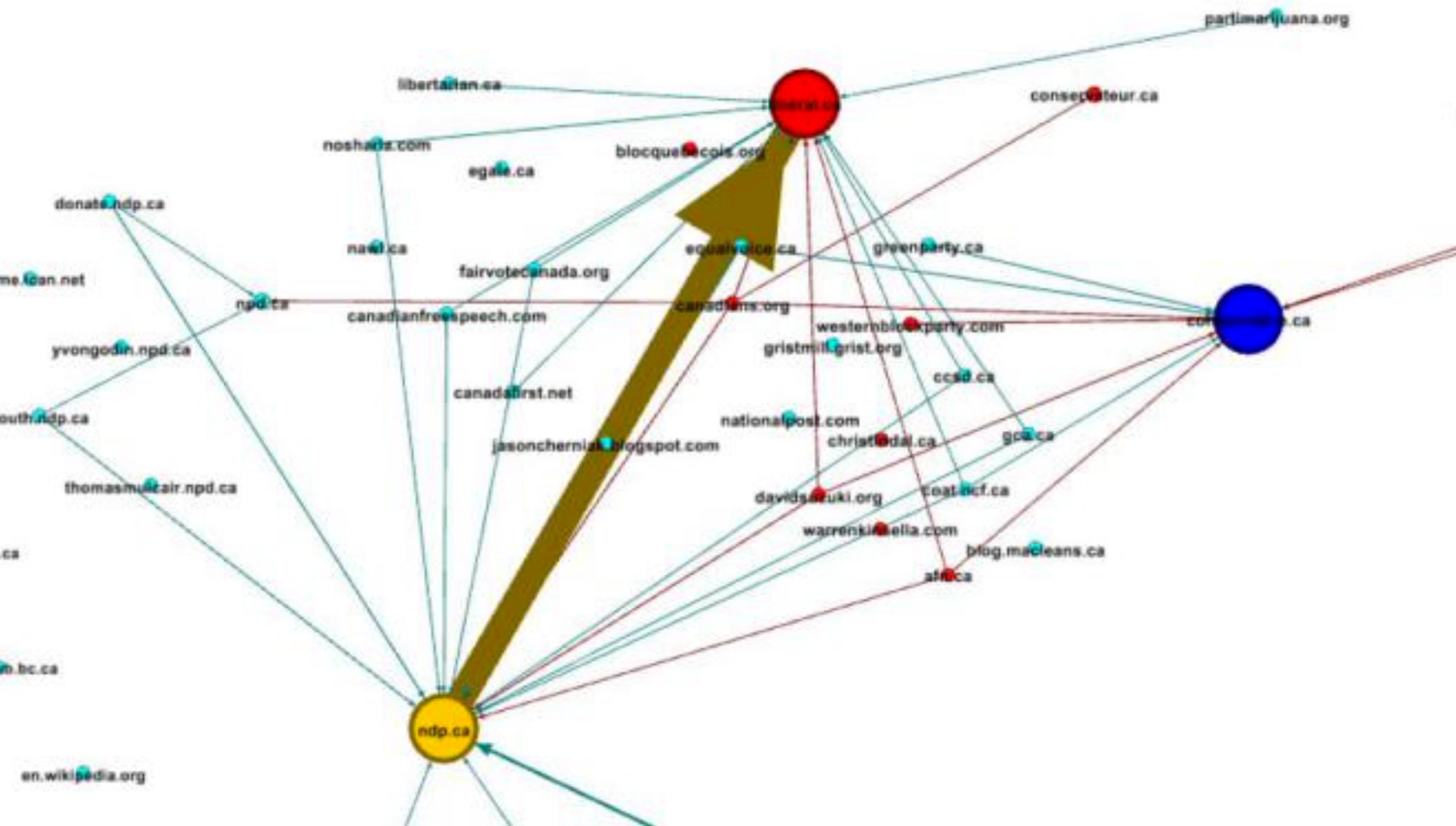
December 2008

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





# 2005 Canadian Federal Election



So metadata, of course,  
helps historians find the  
content they need!

But content is always  
easier to sell to the  
public..

# Shine/WebArchives.ca

- UK Web Archive's Shine (<https://github.com/ukwa/shine>)
- Indexing as bottleneck; warcbase includes a walkthrough for generating Lucene indexes - ~ 5 days locally, few hours with a cluster



# Demo

# Five Things I've 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

Bringing it all together?  
(stay tuned for tomorrow)

# Thank You!

Any Questions?

---

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



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

**Jimmy Lin**  
Professor and David R. Cheriton Chair  
[@lintool](https://twitter.com/lintool)



**UNIVERSITY OF WATERLOO**  
FACULTY OF MATHEMATICS  
David R. Cheriton School  
of Computer Science

**Jeremy Wiebe**  
PhD Candidate  
[@jeremyw](https://twitter.com/jeremyw)



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