

Walking the WALK

Facilitating Interdisciplinary Web Archive
Collaboration

UNIVERSITY OF
WATERLOO



Nick Ruest (@ruebot)
Ian Milligan (@ianmilligan1)

YORK
UNIVERSITÉ
UNIVERSITY





**Why should we
even care about
web archives?**

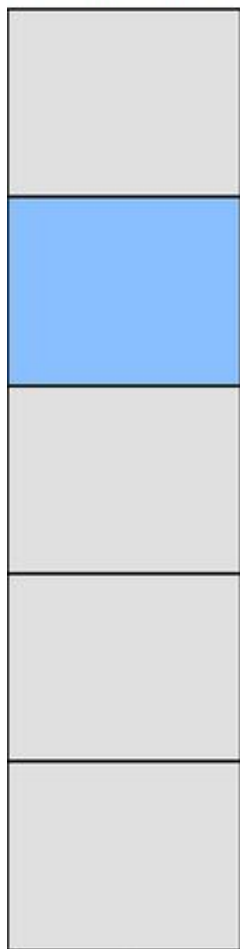
**First, more data than
ever before is being
preserved...**

**Second, it'll be saved
and delivered to us in
very different ways**



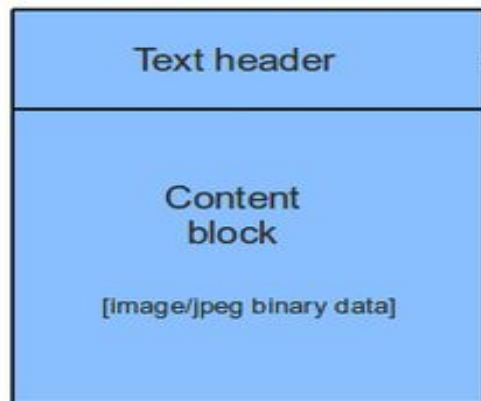
WARC (ISO 28500:2009)

WARC file



...etc.

WARC record



```
WARC/1.0
WARC-Type: resource
WARC-Target-URI: file://var/www/htdocs/images/logoc.jpg
WARC-Date: 2006-09-19T17:20:24Z
WARC-Record-ID: <urn:uuid:92283950-ef2f-4d72-b224-f54c6ec90bb0>
Content-Type: image/jpeg
WARC-Payload-Digest: sha1:CCHXETFVJD2MUZY6ND6SS7ZENMWF7KQ2
WARC-Block-Digest: sha1:CCHXETFVJD2MUZY6ND6SS7ZENMWF7KQ2
Content-Length: 1662
```

You are viewing an archived web page, collected at the request of [Internet Archive Global Events](#) using [Archive-It](#). This page was captured on 5:07:27 Dec 03, 2011, and is part of the [Occupay Movement 2011/2012](#) collection. The information on this web page may be out of date. See [All versions](#) of this archived page. [Videos](#) [Metadata](#)



Welcome [login](#) | [signup](#)
Language [en](#) [es](#) [fr](#)

- News
- LiveStream
- #HowToOccupy
- Forum
- Chat
- User Map
- NYCGA
- About
- Donate
-
-
-

Farmers Join Occupy Wall Street, Calling for Food Justice

Posted 5 hours ago on Dec. 2, 2011, 6:21 p.m. EST by [OccupayWallSt](#)



As Wall Street's corrupt influence on the economy has grown, the corporate ownership of our food system has hurt the health and livelihood's of some of our most vulnerable communities. This *Sunday, December 4th* food justice activists and occupiers will be traveling from as far as Colorado, Iowa, Maine and Upstate New York to join together for the **Occupy Wall Street FARMERS' MARCH**. Through a day of dialogue, musical performances, and a march, farmers and their urban allies working for food justice in their communities will form alliances to fight and expose corporate control of the food supply.

Events throughout the day will call and inspire participants to fight against the corporate manipulation of the agriculture system. An industry that is responsible for using chemical toxins tied to soaring obesity rates, heart disease and diabetes and limiting access to affordable, wholesome food to the country's poorest citizens.

[Read More...](#)

[30 Comments](#)

Occupy Wall Street Goes Home

Posted 1 day ago on Dec. 1, 2011, 3:04 p.m. EST by [OccupayWallSt](#)



On December 6th Occupy Wall Street will join in solidarity with a Brooklyn community to re-occupy a foreclosed home. The day of action marks a national kick-off for a new frontier for the occupy movement: the liberation of vacant bank-owned homes for those in need. The banks not

General Inquiries:
general@occupywallst.org
Press Inquiries:
press@occupywallst.org
Press Phone: +1 (347) 292-1444
Help & Directions: +1 (516) 708-4777
Watch: The world we're building
Read: This call to action
Liberty Square Eviction Defense:
Text "@occupyalert" to 23559 to receive alerts in the event of imminent emergency.

Occupy Wall Street is leaderless resistance movement with people of many colors, genders and political persuasions. The one thing we all have in common is that **We Are The 99%** that will no longer tolerate the greed and corruption of the 1%. We are using the revolutionary **Arab Spring** tactic to achieve our ends and encourage the use of nonviolence to maximize the safety of all participants.

This #ows movement empowers real people to create real change from the bottom up. We want to see a **general assembly** in every backyard, on every street corner because we don't need Wall Street and we don't need politicians to build a better society.

the only solution is WorldRevolution

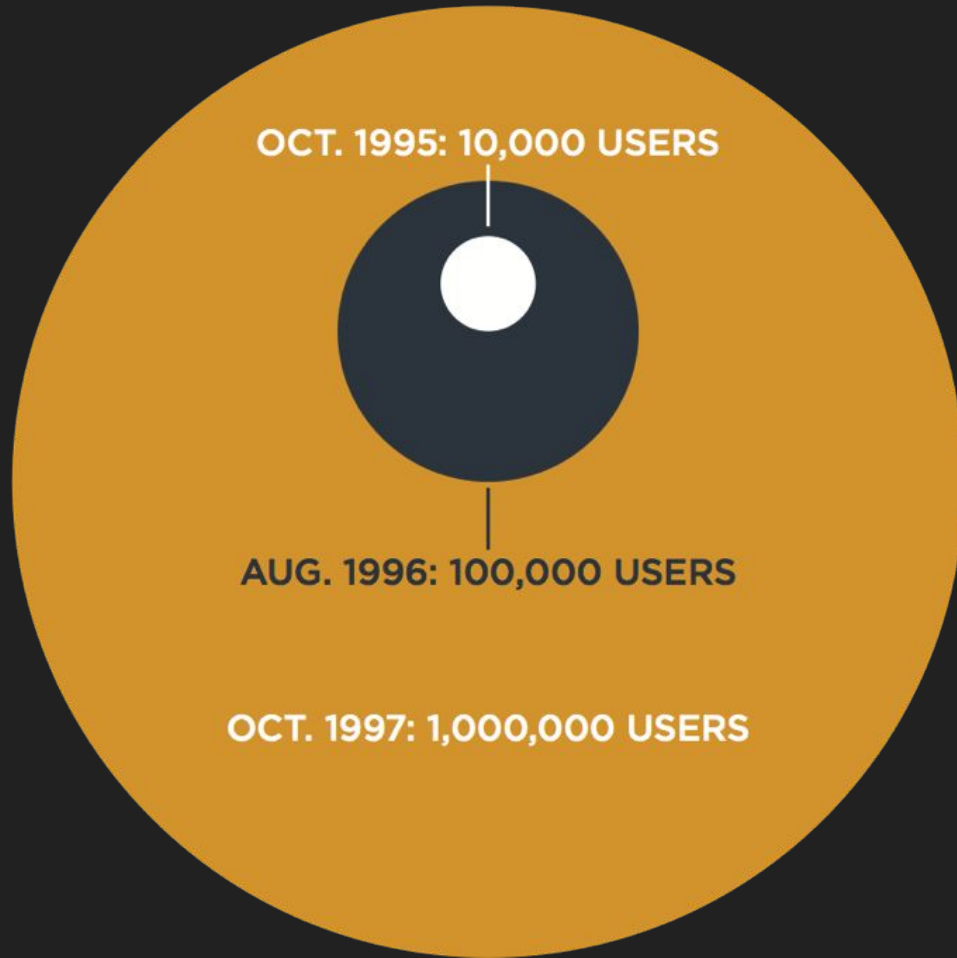
[Click here](#) for NYCGA committee meeting times.





~~Scarcity~~
Abundance

GEOCITIES USERS:



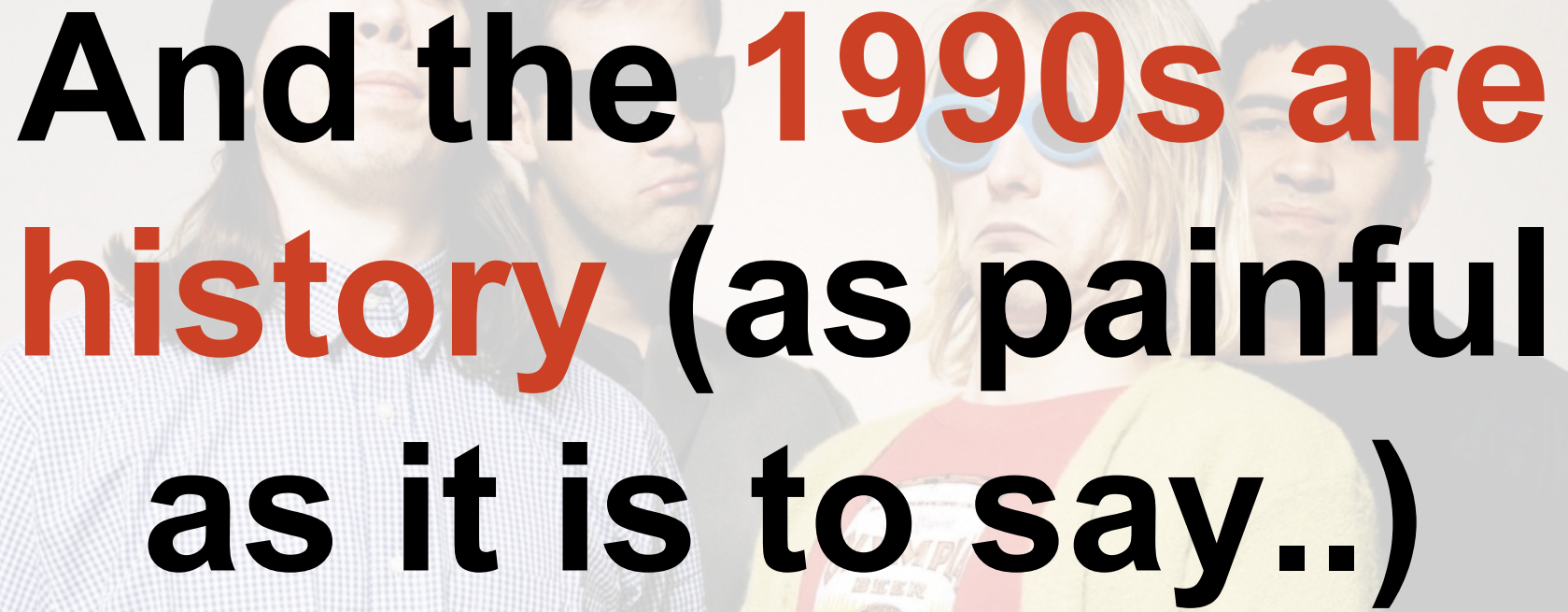
OCT. 1995: 10,000 USERS

AUG. 1996: 100,000 USERS

OCT. 1997: 1,000,000 USERS

The background features the text '1990s' in a large, stylized font at the top. Below it, the faces of two people are visible, one on the left and one on the right, both appearing to be smiling or speaking. The overall color palette is light and pastel.

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



**And the 1990s are
history (as painful
as it is to say..)**

**But right now you
have to use the
Wayback Machine -
requiring you know
the URL!**

University of Waterloo | Home

web.archive.org/web/20160401000145/http://www.uwaterloo.ca

Home | About | News | Contact | Search

University of Waterloo
Ontario, Canada



[On the students for this year](#)

[Campus Waterloo](#)

- prospective students
- students

UW-led team measuring ozone from space

The Atmospheric Chemistry Experiment is a small satellite, designed at Waterloo, that orbits the Earth testing the concentration of gases and particles. "ACE measures the constituents of our atmosphere by collecting spectra as the sun sets through the atmosphere," explains UW chemistry professor Peter Samson. [Details](#)

Today is Thursday, January 21, 2016

[Daily Bulletin \(Thursday\): Distinguished teachers are named](#)

[Critical review tips for new Assistant Centre](#)

[UW launches distance degree in business and technology ... **More news**](#)

[UW welcomes 'five report' on post-secondary education](#)

[Summer programs for children on UW's campus](#)

SEARCH

loading for web.archive.org... | page 1/1

Microsoft Visual Basic 5.0 CONTROL CREATION EDITION

Visual Basic 5.0 Control Creation Edition Free!

GEOCITIES YOUR HOME ON THE WEB



- ENTER HERE
- INFORMATION
- NEIGHBORHOODS
- WHAT'S NEW
- WHAT'S COOL
- WHAT IS GEOCITIES?

G
Our communities are home to the most popular collection of **FREE HOME PAGES & E-MAIL** on the web. Please join or visit one of our 29 neighborhoods today.



[* Free Home Pages & Free Member Email](#)

[Advertiser Information](#)

[GeoCities Daily Audio Update](#) -- Sponsored by [IBM VoiceType Simply Speaking](#)

[Happy Holidays from all of us at Geocities!](#)

[A message from our CEO](#)

Today's Cool Homestead
[WallStreet1456](#)

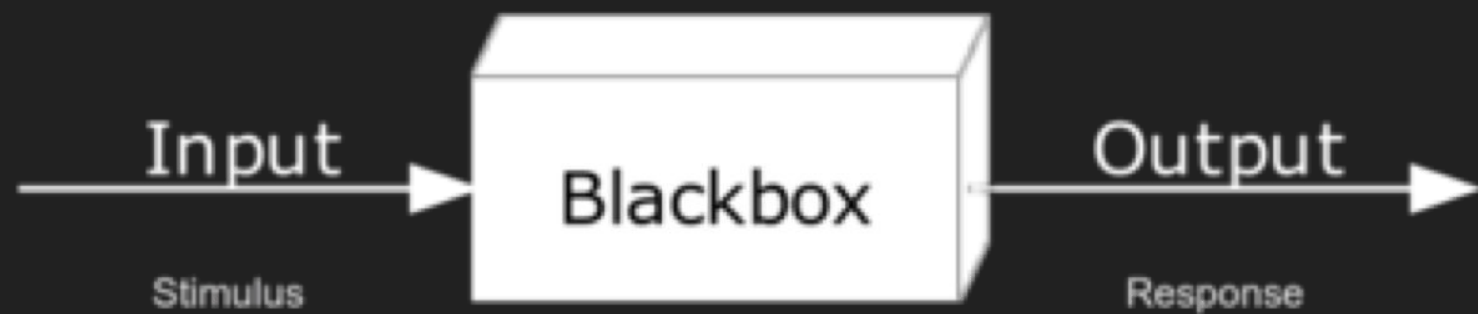
Beginning investors and speculators won't want to miss the Working Class Investor Newsletter.

GeoCities News of the Day - 12/25/96



Building a home page for the holidays?

Submit your letters to Santa, favorite holiday recipes and other holiday cheer to our special [North Pole neighborhood](#). And share your holiday spirit with GeoCitizens around the world [in our virtual holiday tree!](#)



**We need
interdisciplinary
collaboration to
tackle this problem!**



Social Sciences and
Humanities Research
Council of Canada

Conseil de recherches
en sciences humaines
du Canada

Canada



compute
canada

calcul
canada

Team(s)

We form like Voltron

Ian Milligan

History Faculty Member

Jimmy Lin

Computer Science Faculty Member

Jeremy Wiebe

History PhD Candidate

Alice Zhou

Computer Science Undergraduate

Youngbin Kim

Computer Science Undergraduate

Ryan Deschamps

PhD Candidate, Public Policy

Nick Ruest

Digital Assets Librarian

Collaboration

#Slack & GitHub -> Mentoring

Platforms

CLI tools

awk, sed, grep, parallel, sort, uniq, wc, jq

Shine

<https://github.com/ukwa/shine/>

Canadian Political Parties & Political Interest Group Collection (ARCHIVE-IT/Toronto)

- 50 Websites
 - All major political parties
 - Many minor political parties
 - Political interest groups
- Collected quarterly between 2005 and present

HOME THE TEAM THE PARTY ISSUES MEDIA CENTRE YOUR RIDING DONATE



ADDRESS BY PRIME MINISTER PAUL MARTIN

TAKE ACTION TODAY!



Your Excellencies, Honourable Members, Ladies and Gentlemen:

Let me begin by expressing, on behalf of all Canadians, our appreciation to the Right Honourable Adrienne Clarkson and John Ralston Saul. With warmth, intelligence, and wit, they have honoured this high office and made an indelible contribution to our nation.

Over the course of six years, Madame Clarkson recognized achievement, decorated bravery, bore witness to tragedy and grief, and encouraged the disadvantaged. She welcomed foreign visitors and eloquently explained before audiences abroad what it is that makes Canada special. She took great interest in our cities and towns, and especially the north. She traveled to more than 200 communities across Canada; in some of them, it was the first-ever visit by a representative of the Crown.

[Full Story](#)

Stay Informed
GO

Top Stories

- September 29, 2005**
Statement by the Prime Minister on the retirement of John Hamm, Premier of Nova Scotia
- September 28, 2005**
Charity Barbecue Raises \$125,000 for Hurricane Katrina Victims
- September 27, 2005**
Address by Prime Minister Paul Martin at the installation of the new Governor General

[Complete List of Stories](#)

Commissions

-  Young Liberals of Canada
-  National Women's Liberal Commission
-  Aboriginal Peoples' Commission
-  Senior Liberals Commission

[Home](#) | [News](#) | [Your Riding](#) | [Issues](#) | [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 2005. All rights reserved. Authorized by the registered agent for the Liberal Party of Canada.
[Privacy Policy](#)

The Current Interface..

- Very limited - simple search engine, some advanced options; no facets
- Great collection.. But nobody uses them.

The screenshot shows the Internet Archive website interface. At the top, there are navigation links: HOME, EXPLORE, LEARN MORE, and CONTACT US. On the right, there is a logo for the Internet Archive and a tagline: "The leading web archiving service for collecting and accessing cultural heritage on the web. Built at the Internet Archive." Below this, there is a breadcrumb trail: "Explore >> University of Toronto >> Canadian Political Parties and Political Interest Groups".

The main content area features a header for the collection: "Canadian Political Parties and Political Interest Groups", collected by the University of Toronto. It includes a description: "Canadian Political Parties and Political Interest Groups will archive the websites of all of the national Canadian political parties, and a number of special interest groups across the political spectrum." and a collector: "University of Toronto".

Below the header, there is a search bar with the text "stephen harper" entered. To the right of the search bar are buttons for "Search" and "Clear". Below the search bar, there is a message: "The following results were found for the term(s): stephen harper". A bullet point indicates: "No metadata results for stephen harper, but there are up to 1233056 matches within the page text."

There is an "Advanced Search" section with several input fields: "Contains all of:", "Exact phrase:", "Not containing:", and "From the Host:". Below these fields, there is a "Results per host:" dropdown set to "1 (default)", a "File format:" dropdown set to "All formats", and a "Capture date range:" section with "From:" and "To:" dropdowns.

At the bottom of the search section, there are buttons for "Advanced Search" and "Help with Search".

On the right side of the search results, there is a "Search Page Text" section. It shows "Page 1 of 61,653 (1,233,056 Total Results)" and a "Next Page" button. Below this, there is a "Sort By:" dropdown set to "Best Match".

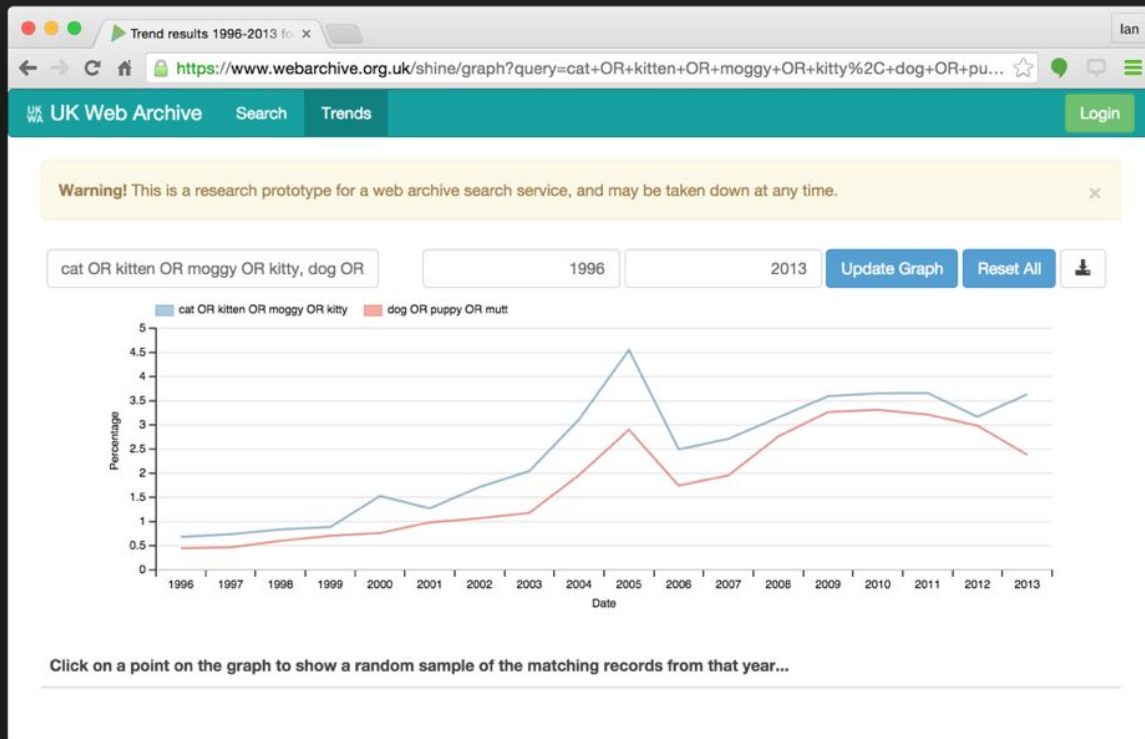
The search results are displayed in a list. The first result is "Stephen Harper | Facebook". The URL is "http://www.facebook.com/pages/Stephen-Harper/9106562109". The text is captured on May 02, 2009. The result includes a snippet of text: "Facebook Sign up for Facebook to connect with Stephen Harper. Information Country: Canada Currently... Stephen Harper | Showing 10 photos Most Recent | Edit Pictures YouTube Box 10 of 13 See all PM on Wolf... the Prime Minister 11:28am Dec 22 | 30 Comments Create a Page Report Page Stephen Harper Wall Info Boxes Notes Stephen Harper + Fans Just Stephen Harper Just Fans Stephen Harper Celebrating... Stephen Harper Launched the Apprenticeship Completion Grant. \$2000 to eligible apprentices. http://tinyurl.com/cqzvyv April 9 at 11:47am Stephen Harper 'Lest we forget.' Statement on the 92nd anniversary of the battle of Vimy Ridge. http://bit.ly/ERb1 April 9 at 11:25am Stephen Harper Announced new... Content: text/html Size: 108 KB More Results from facebook.com".

The second result is "Stephen Harper (pmharper) on Twitter". The URL is "http://twitter.com/PMHarper". The text is captured on Aug 03, 2010. The result includes a snippet of text: "Stephen Harper (pmharper) on Twitter Skip past navigation On a mobile phone? Check out m.twitter.com ! Skip to navigation Skip to sign in form Have an account? Sign in Username or email Password Remember me Forgot password? Forgot username? Already using Twitter on your phone? Get short, timely messages from Stephen Harper. Twitter is a rich source of instantly updated information. It's easy to stay updated on an incredibly wide variety of topics. Join today and follow @pmharper - Get updates via SMS by texting follow pmharper to 40604 in the United States Codes for other countries Two-way (sending and receiving) short codes: Country Code For customers of Australia 0198089488 Telstra Canada... Account Name Stephen Harper Location Ottawa, Ontario Web http://www.conser... Bio Prime Minister of... Content: text/html Size: 46 KB More Results from twitter.com".

The third result is "The Walrus » The Man Behind Stephen Harper » Tom Flanagan » politics". The URL is "http://www.walrusmagazine.com/articles/the-man-behind-stephen-harper-tom-flanagan/". The text is captured on Aug 03, 2010. The result includes a snippet of text: "The Man Behind Stephen Harper... politics".

14 Million Solr docs!

Shine



News - Canada's NDP

wayback.archive-it.org (27/201411072244) http://www.ndp.ca/news/archive/2014-03/special/modules/breaking_news

You are viewing an archived web page, collected at the request of [University of Toronto](#) using [Archive-It](#). This page was captured on 22:44:57 Nov 07, 2014, and is part of the [Canadian Political Parties and Political Interest Groups](#) collection. The information on this web page may be out of date. See [All versions](#) of this archived page.

Home Tom Mulcair Party News Photos and Videos Take Action **DONATE**

News



All News Types

- Press Releases
- Statements
- Speeches
- Reality Checks

Recent News

- October 2014
- September 2014
- August 2014

2014 03 31
The NDP still fighting to save Canada Post
In order to bring the drastic restructuring at Canada Post back to the spotlight, the NDP held a press conference this morning to condemn the Conservatives' new price hikes for stamps, which will bring the cost to a record \$1 and affect consumers across the country.

2014 03 31
Conservatives leaving municipalities behind
The New Building Canada plan will in no way reverse the infrastructure and public transit deficiencies that communities are facing throughout the country. The NDP is urging the Conservative government to take immediate action to ensure that the plan meets the needs of municipalities.

2014 03 28
Official Opposition statement on Earth Hour
Official Opposition Environment Critic Megan Leslie (Halifax) made the following statement on Earth Hour:

"On Saturday night (8:30-9:30 PM local time), New Democrats will join millions of people and businesses around the world in turning off our lights to raise awareness about the need to tackle

Webpage by wayback.archive-it.org

Twitter

twarc

twarc-report

#elxn42

Search**Current Issue**

Issue 32, 2016-04-25

Previous Issues

Issue 31, 2016-01-28

Issue 30, 2015-10-15

Issue 29, 2015-07-15

Issue 28, 2015-04-15

[Older Issues](#)**For Authors**[Call for Submissions](#)[Article Guidelines](#)

An Open-Source Strategy for Documenting Events: The Case Study of the 42nd Canadian Federal Election on Twitter

This article examines the tools, approaches, collaboration, and findings of the Web Archives for Historical Research Group around the capture and analysis of about 4 million tweets during the 2015 Canadian Federal Election. We hope that national libraries and other heritage institutions will find our model useful as they consider how to capture, preserve, and analyze ongoing events using Twitter.

While Twitter is not a representative sample of broader society – Pew research shows in their study of US users that it skews young, college-educated, and affluent (above \$50,000 household income) – Twitter still represents an exponential increase in the amount of information generated, retained, and preserved from 'everyday' people. Therefore, when historians study the 2015 federal election, Twitter will be a prime source.

*On August 3, 2015, the team initiated both a Search API and Stream API collection with *twarc*, a tool developed by Ed Summers, using the hashtag #elxn42. The hashtag referred to the election being Canada's 42nd general federal election (hence 'election 42' or elxn42). Data collection ceased on November 5, 2015, the day after Justin Trudeau was sworn in as the 42nd Prime Minister of Canada. We collected for a total of 102 days, 13 hours and 50 minutes.*

*To analyze the data set, we took advantage of a number of command line tools, utilities that are available within *twarc*, *twarc-report*, and *jq*. In accordance with the [Twitter Developer Agreement & Policy](#), and after ethical deliberations discussed below, we made the tweet IDs and other derivative data available in a data repository. This allows other people to use our dataset, cite our dataset, and enhance their own research projects by drawing on #elxn42 tweets.*

Our analytics included:

- *breaking tweet text down by day to track change over time;*
- *client analysis, allowing us to see how the scale of mobile devices affected medium interactions;*
- *URL analysis, comparing both to Archive-It collections and the [Wayback Availability API](#) to add to our understanding of crawl completeness;*
- *and image analysis, using an archive of extracted images.*

Our article introduces our collecting work, ethical considerations, the analysis we have done, and provides a framework for other collecting institutions to do similar work with our off-the-shelf open-source tools. We conclude by ruminating about connecting Twitter archiving with a broader web archiving strategy.

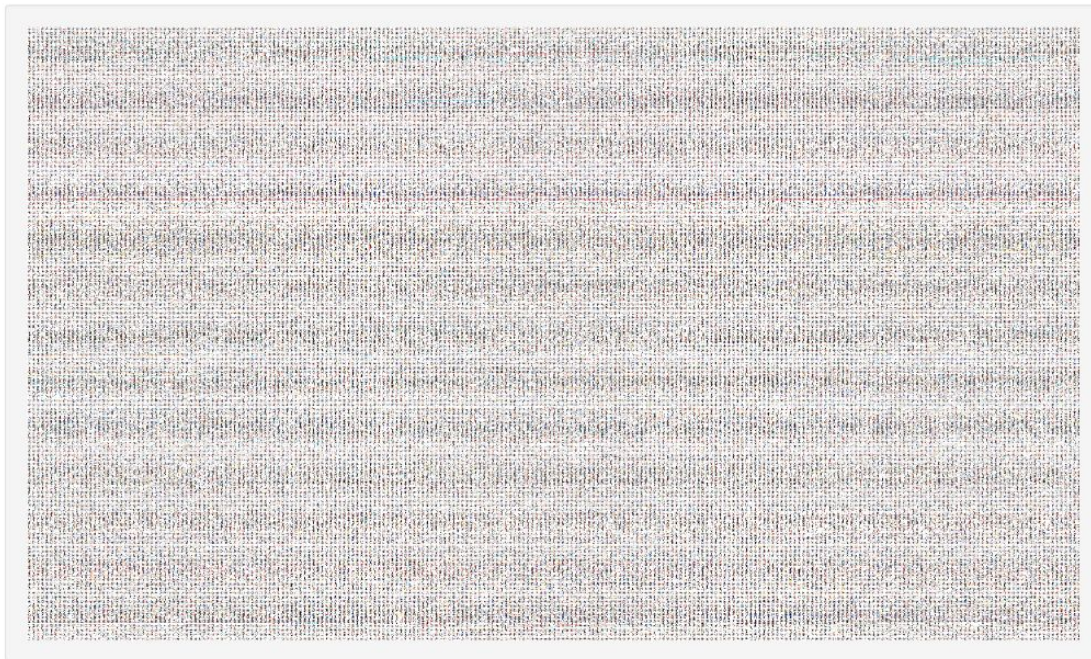
by Nick Ruest and Ian Milligan

Introduction

During the 2015 Canadian federal elections, we captured 3,918,932 tweets written using the #elxn42 hashtag: thoughts on the nature and stature of political candidates or parties, live running commentary during leader debates, exhortations to vote, and witty ripostes or jokes to liven up the long campaign. Political scientists, journalists, and other researchers can use these tweets as evidence of sentiment amongst a

1,203,867 #elxn42 images

Dataset is available [here](#). | Original 32G png available [here](#). | Tiled with [deepzoom.py](#).



This work is licensed under a [Creative Commons Attribution 4.0 International License](#).

[Nick Ruest, Web Archives for Historical Research](#)



SP Dataverse Network

#ELXN42 TWEETS (42ND CANADIAN FEDERAL ELECTION)

hdl:10864/11311

Version: 2 -- Released: Tue Jan 26 17:34:45 EST 2016

[< View Previous Study Listing](#)[Cataloging information](#) **DATA & ANALYSIS** [Comments \(0\)](#) [Versions](#)

i Use the check boxes next to the file name to download multiple files. Data files will be downloaded in their default format. You can also download all the files in a category by checking the box next to the category name. You will be prompted to save a single archive file. Study files that have restricted access will not be downloaded.

 Select all files [Download Selected Files](#)Total Number of Files: **5** Total Downloads: **13**

<input type="checkbox"/>			
<input type="checkbox"/>	elxn42-tweet-ids.txt Plain Text - 71 MB - 8 downloads MD5 Checksum: 98b204a8fc0dbae70e5480c5d4a40a50	Download	line-oriented #elxn42 tweet ids
<input type="checkbox"/>	elxn42-tweets-images.txt Plain Text - 54 MB - 1 download MD5 Checksum: 79376146858835c268cec312c4f7d945	Download	line-oriented #elxn42 image urls tweeted
<input type="checkbox"/>	elxn42-tweets-urls.txt Plain Text - 154 MB - 1 download MD5 Checksum: 78c6372996c69f24ba7ea78a0c2a65a8	Download	line-oriented #elxn42 urls tweeted
<input type="checkbox"/>	elxn42-tweets-users.txt Plain Text - 43 MB - 2 downloads MD5 Checksum: 4821eeb35d9b6b78a4d1f3723c5645ef	Download	line-oriented #elxn42 twitter users
<input type="checkbox"/>	elxn42-tweet-tags.txt Plain Text - 1016 KB - 1 download MD5 Checksum: 59b8161616ddd122f844dd82a8eb4d9	Download	#elxn42 hashtags with counts



[Advanced search](#)

Username *

Password *

Log in

#elxn42 web crawl

Description

Consists of a web crawl of unique URLs tweeted with the #elxn42 hashtag. #elxn42 collection took place from August 3, 2015 - November 5, 2015. Unique URLs were extracted from the dataset, and harvested with Heritrix on January 29, 2016 - February 8, 2016.

Download

- [war: #elxn42 web crawl.gz](#)
- [cdx: #elxn42 web crawl.cdx](#)
- [wat: #elxn42 web crawl.wat.gz](#)
- [Seed list: #elxn42 web crawl.txt](#)
- [Heritrix configuration: crawler-beans.xml](#)

In collections

- [#elxn42](#)

Details

Title:	#elxn42 web crawl
Creator(s):	Nick Ruest
Note:	Tweet ids: http://hdl.handle.net/10864/11311
Identifier (local):	WEB-20160208134917869-00013-3991~rho.library.yorku.ca~9191-ELXN42
Identifier (md5):	63d707352a6fb45a62889c448154610f
Type:	Website
Subject(s):	#elxn42 Canadian federal election, 2015 Canadian politics Federal politics Canada
Date captured:	2016-01-29
Size:	13GB
File size:	12972947583
PUID:	x-fmt/266
Funding:	This research was supported by a research grant -- 435-2015-0011 -- issued by Social Sciences and Humanities Research Council.
Rights:	Use of this resource is governed by the terms and conditions of the Creative Commons "Attribution" License (http://creativecommons.org/licenses/by/2.0/)

#panamapapers

#NDP2016

#MakeDonaldDrumpAgain

#YMMfire

Is this *a* future of collection development?

Institutional vs Research Data

An At-Scale Case Study: Geocities (4.1TB)



Join  **GeoCities**

Neighborhoods Members' Area Shopping Center Search

Search the At Hand® Network Yellow Pages

CATEGORY CITY STATE [AL] [GO]

NEIGHBORHOODS

Visit These Neighborhoods

GeoCities members, or Homesteaders, create their home pages within themed communities called Neighborhoods. Find a Neighborhood that interests you, and see how our Homesteaders use their pages to showcase their interests and creative content for millions of people to see.

- [Arcade](#): Science fiction and fantasy
- [Athens](#): Education, literature, poetry, philosophy
- [Avalon](#): Golf and the finer side of the fairways
- [Baja](#): Four-wheeling, SUVs, off-roading, adventure travel
- [BourbonStreet](#): Jazz, Cajun food, Southern culture
- [Broadway](#): Theater, musicals, show business
- [CanoCanaveral](#): Science, mathematics, aviation
- [Capitol Hill](#): Government, politics, and lots of strong opinions
- [CollegePark](#): University life, from academics to extracurriculars
- [Colosseum](#): Sports and recreation
- [EnchantedForest](#): A neighborhood for and by kids
- [Europe](#): Small businesses, home offices
- [FashionAvenue](#): Top designers, beauty and fashion

Search GeoCities [GO]

Yellow Pages | Domain Stocks | Other Pages | Maps

Home
Help
Info

Doom Level Design with DEU 5 - Netscape

File Edit View Go Communicator Help

Back Forward Reload Home Search Netscape Print Security Stop

Bookmarks Location http://www.geocities.com/Hollywood/29734

DOOM LEVEL DESIGN WITH DEUS

Table of Contents

- [Getting Started](#)
 - Introduction
 - Software/Hardware Requirements
- [What Editor to Choose?](#)
 - Understanding the Editor
 - Understanding the Construction Features
- [Building a New Level](#)
 - Inserting Vertices

Document: Done

Start Doom Level Design w...

00:36

GEOCITIES

Our communities are home to the most popular collection of **FREE HOME PAGES & E-MAIL** on the web. Please join or visit one of our 29 neighborhoods today.

FREE PARIS HEARTLAND TITAN QUARE

- NEIGHBORHOODS
- WHAT'S NEW
- WHAT'S COOL
- WHAT IS GEOCITIES?

* [Free Home Pages & Free Member Email](#) [Advertiser Information](#)



Today's Cool Homestead  [GeoCities Daily Audio Update](#)

[HotSprings 1837](#)

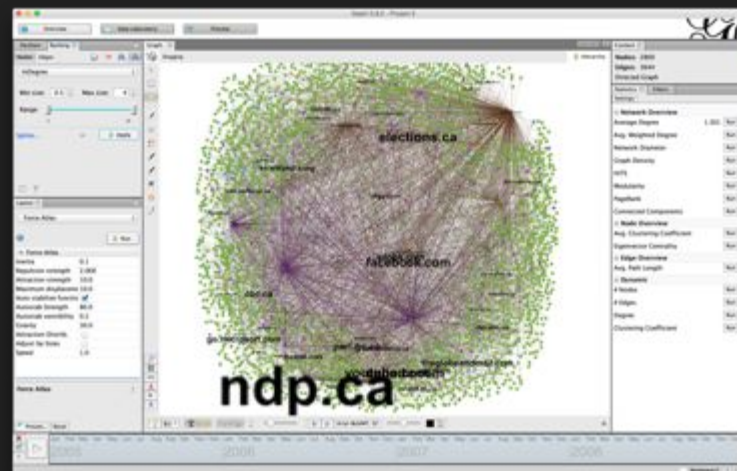
So you hit the snooze bar ten times every morning. You might be lazy. But then again, you might have a sleep disorder. Find out here.

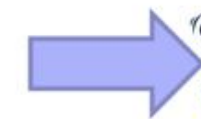
GeoCities News of the Day - 10/22/96

Warcbase

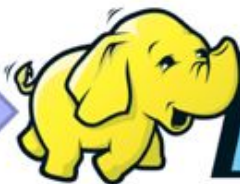
Warchbase

- An open-source platform for managing web archives
- Two main components
 - A flexible data store: your own Wayback Machine
 - Scriptable analytics and data processing





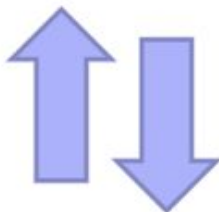
or



hadoop

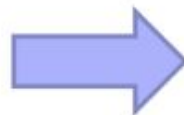
WARC/ARC

Ingestion



Processing & Analytics

A P A C H E
HBASE



Applications
and Services

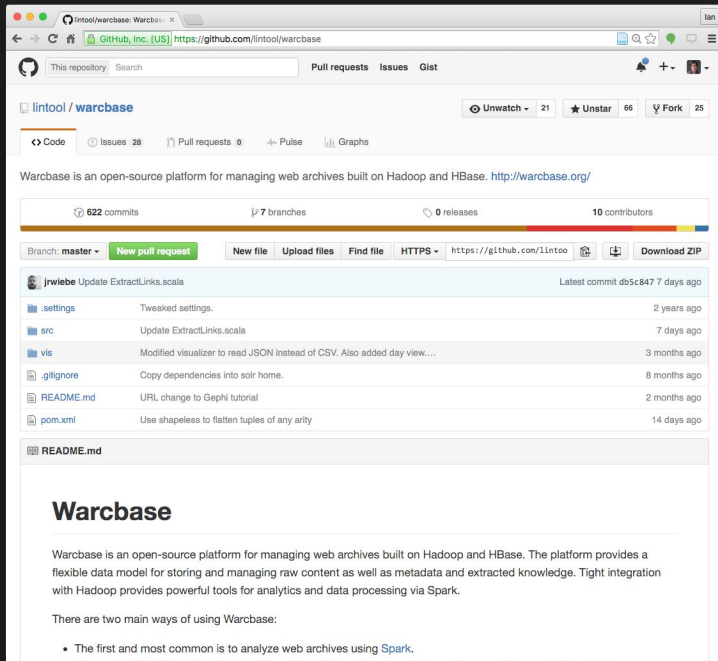


Warcbase

- Scalable
 - From Raspberry Pi to Desktop Computer to Server to Cluster, **all with same scripts and commands**
- 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 at the University of Waterloo

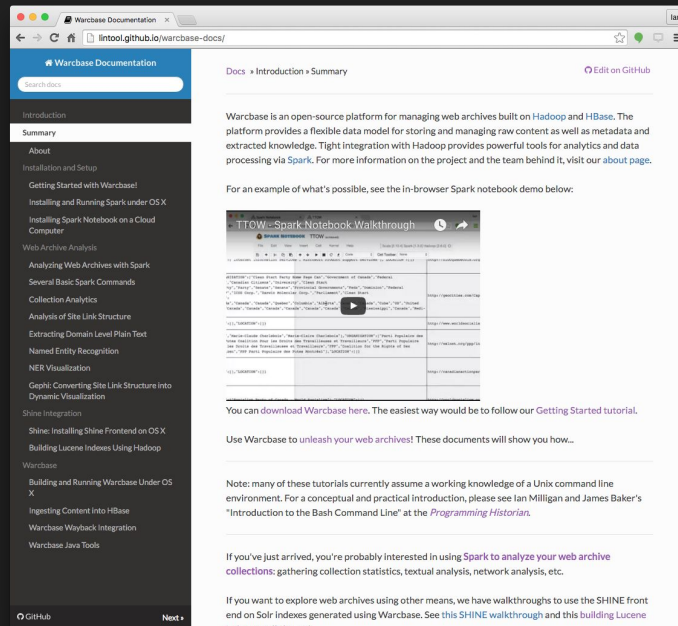


You can Warcbase Too! (...and Twarcbase soon!)



The screenshot shows the GitHub repository page for `lintool/warcbase`. At the top, it displays repository statistics: 622 commits, 7 branches, 0 releases, and 10 contributors. Below this, there's a list of files and folders including `.settings`, `src`, `vis`, `.gignore`, `README.md`, and `pom.xml`. The `README.md` file is selected and its content is visible below. The README text reads: "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 metadata and extracted knowledge. Tight integration with Hadoop provides powerful tools for analytics and data processing via Spark. There are two main ways of using Warcbase: • The first and most common is to analyze web archives using Spark."

warcbase.org



The screenshot shows the Warcbase Documentation page. The page title is "Warcbase Documentation" and it includes a search bar. The main content area is titled "Introduction" and "Summary". It describes Warcbase as an open-source platform for managing web archives built on Hadoop and HBase. Below the text, there's a video player for "TTOW - Spark Notebook Walkthrough". The page also includes a list of documents and a note about the tutorials assuming a working knowledge of a Unix command line environment.

docs.warcbase.org

**Let's do a quick
walkthrough of how
we've used it on one
collection: GeoCities**



```
1. i2millig@rho: /mnt/vol1/data_sets/geocities/warcs (ssh)
bash bash i2millig@rho: /mnt/vol1/data...
GEOCITIES-20091029114236-00191-ia400110.us.archive.org.warc.gz
GEOCITIES-20091029115416-00171-crawling08.us.archive.org.warc.gz
GEOCITIES-20091029123034-00172-crawling08.us.archive.org.warc.gz
GEOCITIES-20091029130439-00173-crawling08.us.archive.org.warc.gz
GEOCITIES-20091029134536-00174-crawling08.us.archive.org.warc.gz
GEOCITIES-20091029140344-00192-ia400110.us.archive.org.warc.gz
GEOCITIES-20091029141553-00193-ia400110.us.archive.org.warc.gz
GEOCITIES-20091029141726-00175-crawling08.us.archive.org.warc.gz
GEOCITIES-20091029144445-00176-crawling08.us.archive.org.warc.gz
GEOCITIES-20091029152151-00177-crawling08.us.archive.org.warc.gz
GEOCITIES-20091029160824-00178-crawling08.us.archive.org.warc.gz
GEOCITIES-20091029164941-00179-crawling08.us.archive.org.warc.gz
GEOCITIES-20091029165037-00194-ia400110.us.archive.org.warc.gz
GEOCITIES-20091029170431-00195-ia400110.us.archive.org.warc.gz
GEOCITIES-20091029171605-00180-crawling08.us.archive.org.warc.gz
GEOCITIES-20091029174154-00181-crawling08.us.archive.org.warc.gz
GEOCITIES-20091029180818-00182-crawling08.us.archive.org.warc.gz
GEOCITIES-20091029182725-00183-crawling08.us.archive.org.warc.gz
GEOCITIES-20091029185858-00184-crawling08.us.archive.org.warc.gz
GEOCITIES-20091029193728-00185-crawling08.us.archive.org.warc.gz
GEOCITIES-20091029194541-00196-ia400110.us.archive.org.warc.gz
GEOCITIES-20091029195911-00197-ia400110.us.archive.org.warc.gz
GEOCITIES-20091029202041-00186-crawling08.us.archive.org.warc.gz
GEOCITIES-20091029221340-00198-ia400110.us.archive.org.warc.gz
GEOCITIES-20091029222459-00199-ia400110.us.archive.org.warc.gz
GEOCITIES-20091030021147-00197-ia400103.us.archive.org.warc.gz
GEOCITIES-20091030021444-00198-ia400103.us.archive.org.warc.gz
GEOCITIES-20091030022413-00171-ia400104.us.archive.org.warc.gz
i2millig@rho:/mnt/vol1/data_sets/geocities/warcs$ du -h
4.1T
i2millig@rho:/mnt/vol1/data_sets/geocities/warcs$
```

```
ianmilligan@Ians-MacBook-Pro:~$ rho
i2millig@rho.library.yorku.ca's password:
Welcome to Ubuntu 14.04.2 LTS (GNU/Linux 3.13.0-32-generic x86_64)

 * Documentation:  https://help.ubuntu.com/

System information as of Mon Mar  7 13:43:20 EST 2016

System load:  0.99           Users logged in:      1
Usage of /:   34.7% of 744.67GB   IP address for em1:   130.63.180.18
Memory usage: 16%             IP address for em2:   10.0.0.18
Swap usage:   6%              IP address for docker0: 172.17.0.1
Processes:   359

Graph this data and manage this system at:
https://landscape.canonical.com/

242 packages can be updated.
130 updates are security updates.

Last login: Mon Mar  7 13:43:21 2016 from 38.123.136.254
i2milli@rho:~$ ./spark-1.5.1/bin/spark-shell --jars ~/warcbase/target/warcbase-0.1.0-SNAPSHOT-fatjar.jar
WARN NativeCodeLoader - Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Welcome to

  ____
 /  __ \
 \  / __/
  \  /_ 
   \____/
          version 1.5.1

Using Scala version 2.10.4 (Java HotSpot(TM) 64-Bit Server VM, Java 1.8.0_45)
Type in expressions to have them evaluated.
Type :help for more information.
WARN MetricsSystem - Using default name DAGScheduler for source because spark.app.id is not set.
Spark context available as sc.
SQL context available as sqlContext.

scala> :paste
// Entering paste mode (ctrl-D to finish)

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

val r =
RecordLoader.loadWarc("/mnt/voll/data_sets/geocities/warcs/GEOCITIES-20090808133634-04399-crawling08.us.archive.org.warc.gz", sc)
.keepValidPages()
.map(r => ExtractTopLevelDomain(r.getUr1))
.countItems()
.take(10)

// Exiting paste mode, now interpreting.

INFO WacWarcInputFormat - Loading file:/mnt/voll/data_sets/geocities/warcs/GEOCITIES-20090808133634-04399-crawling08.us.archive.org.warc.g
z
import org.warcbase.spark.matchbox._
import org.warcbase.spark.rdd.RecordRDD._
r: Array[(String, Int)] = Array((geocities.com,3748), (www.geocities.com,240), (www.myfilehut.com,12), (asiarooms.com,7), (us.geocities.com
,6), (www.theginge.com,3), (www.angelfire.com,3), (images.quizilla.com,3), (pub28.bravenet.com,3), (ss.webring.yahoo.com,2))

scala>
```

LIVE DEMO

Extracting all URLs

```
1 import org.warcbase.spark.matchbox._
2 import org.warcbase.spark.rdd.RecordRDD._
3
4 val r = RecordLoader.loadWarc("/mnt/vol1/data_sets/geocities/
   warcs", sc)
5 .keepValidPages()
6 .map(r => r.getUrl)
7 .saveAsTextFile("/mnt/vol1/derivative_data/geocities/url-list")
```

Results = 186,761,346 URLs, 9.9GB text file

Extracting a Link Graph

```
1 import org.warcbase.spark.matchbox.{ExtractTopLevelDomain,
   ExtractLinks, RecordLoader}
2 import org.warcbase.spark.rdd.RecordRDD._
3
4 RecordLoader.loadArc("/mnt/vol1/data_sets/geocities/warcs/*", sc)
5 .keepValidPages()
6 .map(r => (r.getCrawldate, ExtractLinks(r.getUrl, r.
   getContentString)))
7 .flatMap(r => r._2.map(f => (r._1, ExtractTopLevelDomain(f._1).
   replaceAll("^\\s*www\\.",""), ExtractTopLevelDomain(f._2).
   replaceAll("^\\s*www\\.",""))))
8 .filter(r => r._2 != "" && r._3 != "")
9 .countItems()
10 .filter(r => r._2 > 5)
11 .saveAsTextFile("/mnt/vol1/data_sets/geocities/geocities.
   sitelinks")
```

Results

- 1 ((20090903, <http://geocities.com/saganaki2000/ADSLGR/adslgr.htm>, <http://www.adslgr.com>), 15337)
- 2 ((20091026, <http://geocities.com/saganaki2000/ADSLGR/adslgr.htm>, <http://www.adslgr.com>), 15337)
- 3 ((20091027, <http://geocities.com/spankbank69hard/>, http://pg.photos.yahoo.com/ph/spankbank69hard/my_photos/), 9807)
- 4 ((20090903, <http://geocities.com/spankbank69hard/index.html>, http://pg.photos.yahoo.com/ph/spankbank69hard/my_photos/), 9807)
- 5 ((20091027, <http://geocities.com/CollegePark/Locker/8187/>, <http://www.comercialuruapan.com>), 8056)
- 6 ((20090903, <http://geocities.com/CollegePark/Locker/8187/>, <http://www.comercialuruapan.com>), 8056)

Creating Entities

403GB of link graph data.

- <http://www.geocities.com/EnchantedForest/Grove/1234/index.html>
- <http://www.geocities.com/EnchantedForest/Grove/1234/pets/cats.html>
- <http://www.geocities.com/EnchantedForest/Grove/1234/pets/dogs.html>
- <http://www.geocities.com/EnchantedForest/Grove/1234/pets/rabbits.html>

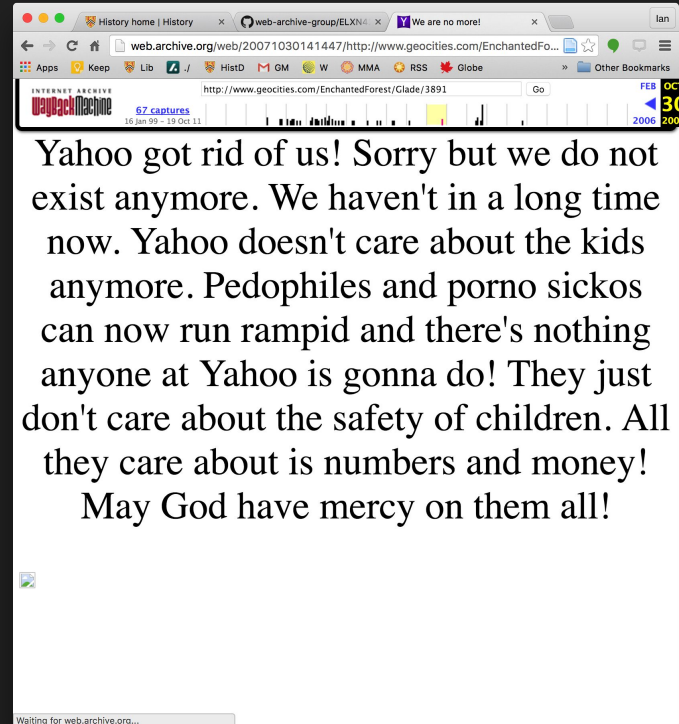
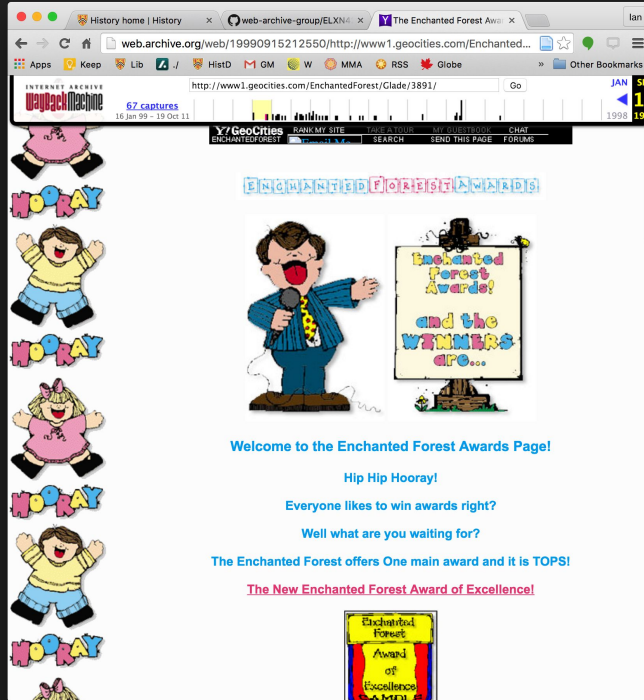
Link Structure

- 1 Source, Target, Weight
- 2 <http://www.geocities.com/EnchantedForest/Meadow/1134>, <http://www.geocities.com/EnchantedForest/1004>, 83
- 3 <http://www.geocities.com/EnchantedForest/Meadow/1134>, <http://www.geocities.com/EnchantedForest/1004>, 83
- 4 <http://www.geocities.com/Area51/Stargate/1357>, <http://www.geocities.com/Area51/EnchantedForest/4213>, 33
- 5 <http://www.geocities.com/Area51/Stargate/1357>, <http://www.geocities.com/Area51/EnchantedForest/4213>, 33
- 6 <http://www.geocities.com/Eureka/1309>, <http://www.geocities.com/EnchantedForest/Tower/7555>, 27
- 7 <http://www.geocities.com/Eureka/1309>, <http://www.geocities.com/EnchantedForest/Tower/7555>, 27

<http://geocities.com/WestHollywood/Castro/6902>



EnchantedForest/Glade/3891



WALK Project

Web Archives for Longitudinal Knowledge

WALK

ARCHIVE-IT (Internet Archive
subscription service)

25 Canadian partners

130 Publicly-accessible
collections

~ 35TB of ARCs/WARCs files



Could we do for Canadian
Archive-It Partners what we
did for GeoCities and
webarchives.ca?

Preparing our dataset!

Project

Compute

Overview

Instances

Volumes

Images

Access & Security

Network

Orchestration

Identity

Overview

Limit Summary



Instances
Used 2 of 10



VCPUs
Used 18 of 20



RAM
Used 75GB of 75GB



Floating IPs
Allocated 2 of 2



Security Groups
Used 1 of 10



Volumes
Used 1 of 10



Volume Storage
Used 15.8TB of 15.8TB

Usage Summary

Select a period of time to query its usage:

From: 2016-06-01

To: 2016-06-01

Submit

The date should be in YYYY-mm-dd format.

Active Instances: 2 Active RAM: 75GB This Period's VCPU-Hours: 437.23 This Period's GB-Hours: 11246.43 This Period's RAM-Hours: 1865498.75

Usage

[Download CSV Summary](#)

Instance Name	VCPUs	Disk	RAM	Time since created
compute-canada-will-accidentally-delete-this	16	412GB	60GB	3 weeks, 1 day
...

Public Access

Shine

Voyant-Tools

```
[nruest@rho:nruest]$ curl "http://localhost:8080/solr/select?q=rob+ford&start=0&rows=2&fl=url%2Ccrawl_date%2Ccontent%2Ctitle%2Ccrawl_year%2Cdomain&wt=json&indent=true"
{
  "responseHeader":{
    "status":0,
    "QTime":0,
    "params":{
      "q":"rob ford",
      "indent":"true",
      "fl":"url,crawl_date,content,title,crawl_year,domain",
      "start":"0",
      "rows":"2",
      "wt":"json"}},
  "response":{"numFound":4616,"start":0,"docs":[
    {
      "url":"http://canadians.org/sites/default/files/styles/large/public/wp-content/uploads/steven-and-rob-ford-240x180_0.jpg?itok=ePOBxrbr",
      "domain":"canadians.org",
      "crawl_date":"2013-08-08T20:45:38Z",
      "crawl_year":"2013",
      "content":[""]},
    {
      "url":"http://www.youtube.com/watch?v=c0sEV985ScI",
      "domain":"youtube.com",
      "crawl_date":"2014-11-05T22:50:05Z",
      "crawl_year":"2014",
      "content":["Ford bedbug - YouTubeUpload Sign inSearchLoading...This video is unavailable.Watch QueueTV QueueWatch QueueTV QueueRemove allDisconnectLoading...Watch QueueTV Queue_count_/to
tal_Find out why CloseFord bedbugAnimation: Mike ConstableSubscribeSubscribedUnsubscribe 15Subscription preferencesLoading...Loading...Working...Add toWant to watch this again later?Sign in to add
this video to a playlist.Sign inShareMoreReportNeed to report the video?Sign in to report inappropriate content.Sign inStatistics161 views1Like this video?Sign in to make your opinion count.Sign in0
Don't like this video?Sign in to make your opinion count.Sign inLoading...Loading...Loading...Sign inRatings have been disabled for this video.Rating is available when the video has been rented.This
feature is not available right now. Please try again later.Uploaded on Feb 1, 2011A look at a capitalist as mayor of torontoCategoryComedyLicenseStandard YouTube LicenseShow moreShow lessLoading...
40:51Play nextPlay nowThe Rob Ford Story - the fifth estateby CBC News5,330 views24:11Play nextPlay nowThe Funniest MAYOR ROB FORD Compilation!by Aiden McLean949,863 views30:44Play nextPlay nowCoffe
Run - Rob Fordby deadmau5815,474 views9:19Play nextPlay nowJon Stewart DESTROYs Obama and Mayor Rob Ford 11/14/2013 |by richardnyavor35,226 views6:16Play nextPlay nowJon Stewart's Crack Smoking May
or Of Toronto Rob Ford segment 11 14 2013by Peter Tsirlis70,575 views5:36Play nextPlay nowDel Grande Speaks on Appointmentby Matt Elliott178 views4:45Play nextPlay nowJon Stewart on Rob Ford Smoking
Crack Again (May 2014) New Videoby LittleDragon7789,240 views5:01Play nextPlay nowRob Ford's Brother Doug: Top 10 WTFord Momentsby Truth Mashup6,791 views12:58Play nextPlay nowctv news kitchener ro
b ford mayor of toronto crack cocaine nov 4 2013by derek otto260 views5:42Play nextPlay nowToronto Mayor Rob Ford on Jimmy Kimmel Live PART 3by Jimmy Kimmel Live953,911 views6:01Play nextPlay nowCan
adian Lifestyle is Shit, Bedbugs, Mold in Toronto; Major Nourhighighl 's Proofs 0^U 0±U†0^U by CanadianLifeStyle469 views4:36Play nextPlay nowToronto Mayor Rob Ford on Jimmy Kimmel Live PART 1by Jim
my Kimmel Live2,247,246 views2:51Play nextPlay nowI Have Bed Bugs...by Raya55,292 views0:08Play nextPlay nowCranked cruiserby pippyconstable78 views4:45Play nextPlay now500 Dawes road-the worst buil
ding in Torontoby kegfiwuegfu hsegrfcyweg10,900 views2:34Play nextPlay now200 Wellesley Apartment Fire Victims Struggle To Resettle TCHCby 200Wellesley2,242 views3:45Play nextPlay nowCouncillor Fletc
her and the Toronto Bed Bug Projectby PaulaFletcher30158 views2:30Play nextPlay nowToronto Community Housing Corporation Public Scams And Corruption TCHC Part 26 of 78by 200Wellesley1,070 views2:40P
lay nextPlay nowFast Ford Nation! Tasha Hearts ROB FORDby Tasha James33,791 views1:25:41Play nextPlay nowToronto 2014 mayoral debate with Rob Ford, John Tory, Olivia Chow and David Soknackiby The GL
obe and Mail31,295 viewsLoading more suggestions...Show moreLanguage:EnglishCountry:WorldwideSafety:OffHelpLoading...Loading...Loading...AboutPress & BlogsCopyrightCreators & PartnersAdvertisingDeve
lopersYouTubeTermsPrivacyPolicy & SafetySend feedbackTry something new!Loading...Working...Sign in to add this to Watch LaterAdd to"],
    "title":"Ford bedbug - YouTube"}]
  }}
[nruest@rho:nruest]$ █
```

[http://voyant-tools.org/?input=http://rho.library.yorku.ca:8080/solr/select%3Fq%3Ddayton%26start%3D0%26rows%3D250%26fl%3Durl,crawl_date,content,title,crawl_year, domain%26wt%3Dxml%26indent%3Dtrue&xmlDocumentsXpath=//doc&xmlContentXpath=//arr\[@name%3D%27content%27\]&xmlTitleXpath=//str\[@name%3D%27title%27\]&xmlAuthorXpath=//str\[@name%3D%27domain%27\]&xmlPubDateXpath=//str\[@name%3D%27crawl_years%27\]](http://voyant-tools.org/?input=http://rho.library.yorku.ca:8080/solr/select%3Fq%3Ddayton%26start%3D0%26rows%3D250%26fl%3Durl,crawl_date,content,title,crawl_year, domain%26wt%3Dxml%26indent%3Dtrue&xmlDocumentsXpath=//doc&xmlContentXpath=//arr[@name%3D%27content%27]&xmlTitleXpath=//str[@name%3D%27title%27]&xmlAuthorXpath=//str[@name%3D%27domain%27]&xmlPubDateXpath=//str[@name%3D%27crawl_years%27])

Cirrus Terms Links

Reader

Trends

Document Terms



Jack Layton | Leadership(ndp.ca)

Jack Layton | Leadership

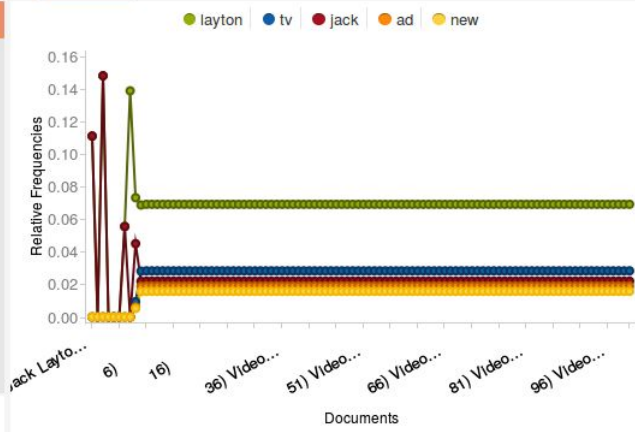
Jack Layton, Leaders Tour - Tournée du Chef - Jack Layton Jack Layton, Leaders Tour - Tournée du Chef - Jack Layton(davidsuzuki.org)

Jack Layton, Leaders Tour - Tournée du Chef - Jack Layton Jack Layton, Leaders Tour - Tournée du Chef - Jack Layton

Jack Layton, Leaders Tour - Tournée du Chef - Jack Layton Jack Layton, Leaders Tour - Tournée du Chef - Jack Layton(davidsuzuki.org)

Jack Layton, Leaders Tour - Tournée du Chef - Jack Layton Jack Layton, Leaders Tour - Tournée du Chef - Jack Layton

INSPIRE awards honour Jack Layton - YouTube(youtube.com)



Scale Terms

Navigation controls

Scale Frequencies

Summary Documents Phrases

Contexts

This corpus has 250 documents with 72,685 total words and 255 unique word forms. Created 2 seconds ago.

Document Length:
 • Longest: Video Gallery | NDP (320); Video Gallery | NDP (320); Video Gallery | NDP (320); Video Gallery | NDP (318); Video Gallery | NDP (318)
 • Shortest: (0); (0); (0); (0); (0)

Vocabulary Density:
 • Highest: INSPIRE awards honour... (1.000); Layton at Central Station... (1.000); Jack Layton | Leadership (1.000); Comments on: People's... (0.667); Video Gallery | NDP (0.603)
 • Lowest: 0; 0; 0; 0; 0

Items

Document	Left	Term	Right
1) Jack ...	Jack	lay...	Leadership
7) Jack ...	Jack	lay...	, Leaders Tour - Tournée du Chef
7) Jack ...	Tour - Tournée du Chef - Jack	lay...	Jack Layton, Leaders Tour - Tournée
7) Jack ...	du Chef - Jack Layton Jack	lay...	, Leaders Tour - Tournée du Chef
7) Jack ...	Tour - Tournée du Chef - Jack	lay...	
8) Jack ...	Jack	lay...	, Leaders Tour - Tournée du Chef
8) Jack ...	Tour - Tournée du Chef - Jack	lay...	Jack Layton, Leaders Tour - Tournée
8) Jack ...	du Chef - Jack Layton Jack	lay...	Leaders Tour - Tournée du Chef

5,037 context expand Scale

warcbase

Again :-)

Future

- The future is now!
- Ways of categorizing websites?
- Include a third set category (eg. by archiving institution)?
- Expand the number of collections for analysis.
- Refine the tool with better UX so it could be turned into a usable website.

Links!

- <https://uwaterloo.ca/web-archive-group/>
- <https://github.com/web-archive-group/>
- <https://github.com/ianmilligan1/>
- <https://github.com/ruebot>
- <http://dataverse.scholarsportal.info/dvn/dv/wahr>

Contact

Nick Ruest: @ruebot

ruestn@yorku.ca

Ian Milligan: @ianmilligan1

i2milligan@uwaterloo.ca