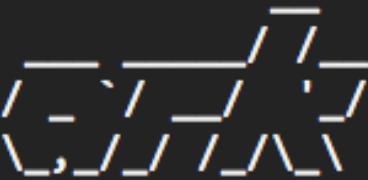


YOU'VE GOT WEB ARCHIVES.. WHAT TO DO WITH THEM?

Hacking the News Workshop



version 2.3.0

ersion 2.11.8 (Java HotSpot(TM) 64-Bit Server VM v1.8.0_111-b14) on 2017-09-12 11:30:21
essions to have them evaluated.
r more information.

aste mode (ctrl-D to finish)

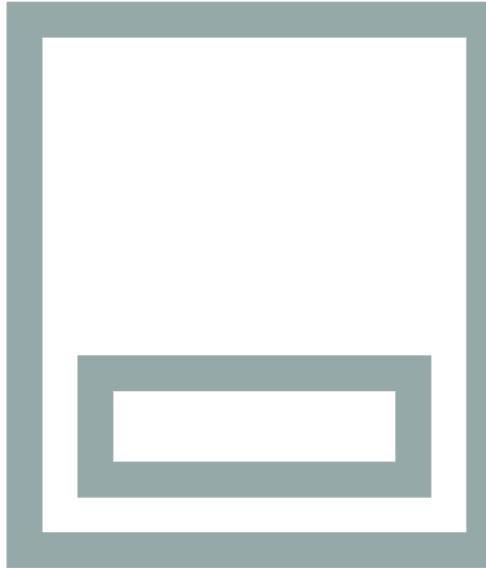
hivesunleashed._
hivesunleashed.matchbox._

dLoader.loadArchives("example.arc.gz")
es()
ractDomain(r.getUrl())

2. ssh

WEB ARCHIVE ANALYSIS

- The Archives Unleashed Toolkit was designed to tackle this problem.
- Allows a user to take WARCs and:
 - Determine elemental statistics about a collection;
 - Extract particular images, domains, URLs, pages with keywords, etc.
 - Do sophisticated Apache Spark-powered network analysis; and
 - Write custom Scala scripts to do almost anything you imagine with our set of custom web archiving User Defined Functions



LET'S SEE HOW THIS WORKS.

HACKING THE NEWS

Safari File Edit View History Bookmarks Develop Window Help

archivesunleashed.org

Welcome to the Archives Unleashed Project



The Archives Unleashed Project

Welcome

Archives Unleashed aims to make petabytes of historical internet content accessible to scholars and others interested in researching the recent past. Supported by a grant from the [Andrew W. Mellon Foundation](#), we are developing web archive search and data analysis tools to enable scholars, librarians and archivists to access, share, and investigate recent history since the early days of the World Wide Web.

Interested in the project? Subscribe to our [newsletter!](#) Or you can follow the links at left for information about the project, the [Archives Unleashed Cloud](#), [Archives Unleashed Toolkit](#), [Archives Unleashed Jupyter Notebooks](#), or our [events](#).

We're always looking for [ways to engage](#) archivists, librarians, researchers, developers, or any others interested in born-digital heritage!

Contact Us

Questions? Comments? Please contact us, either by leaving an issue on one of our [GitHub projects](#) or by sending us an e-mail. Are you a Slack user? Join our [Slack team!](#)

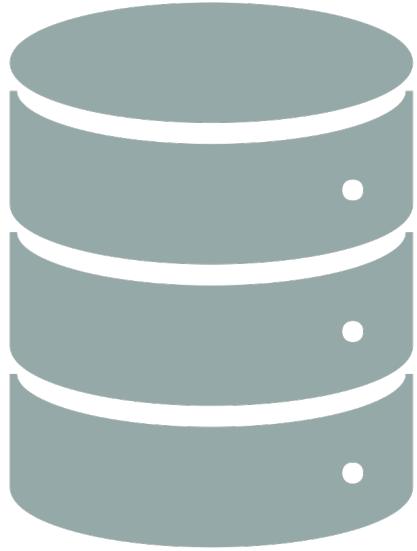
2. ssh

i2milligan@tuna:~\$



THIS IS DIFFICULT (REQUIRES COMMAND LINE)

HACKING THE NEWS



THE DATASETS IN “HACKING-THE-NEWS.ZIP”
WERE GENERATED BY THIS

The screenshot shows the Archives Unleashed Cloud interface. At the top, there's a navigation bar with icons for back, forward, search, and user profile (Ian Milligan). The main title is "Nova Scotia Municipal Governments". Below the title is a blue button labeled "Analyze Collection".

On the left, there's a sidebar with a user profile picture of Ian Milligan, his email (i2milligan@uwaterloo.ca), and his AU Cloud Account information. It also shows he has 13 jobs run and 67 GB of disk usage.

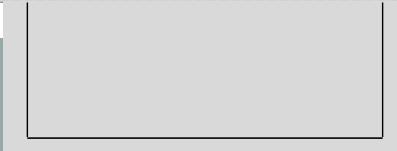
The main content area has two sections:

- Hyperlink Diagram:** A circular network graph where nodes represent web pages and edges represent hyperlinks between them. Nodes are color-coded by domain.
- Domains:** A bar chart showing the number of pages per domain. The domains listed are: www.cbpm.ns.ca, www.digbyisland.ca, www.chester.ca, www.inno.ca, www.mod.ca, www.town.benwick.ns.ca, www.town.yarmouth.ca, www.westhants.ca, www.town.windes.ns.ca, and www.wolfville.ca. The y-axis ranges from 0 to 10,000.

At the bottom, there are logos for the Mellon Foundation, University of Waterloo, and York University. A note says: "For more information on our project and sponsors, visit [archivesunleashed.org](#)". There are also links for "About", "Privacy Policy", "Documentation", and "FAQ".

AN EASIER WAY

- The **Archives Unleashed Cloud** was designed to make this easier..
- Allows a user to take WARCs and:
 - Use a modern UI to sync their collections from the provider;
 - Run basic analyses in the browser to find major sites of interest;
 - Download derivative file formats that can integrate with standard workflows.
- In other words: let's get the WARC out of the equation and *translate* it into a standard file format.

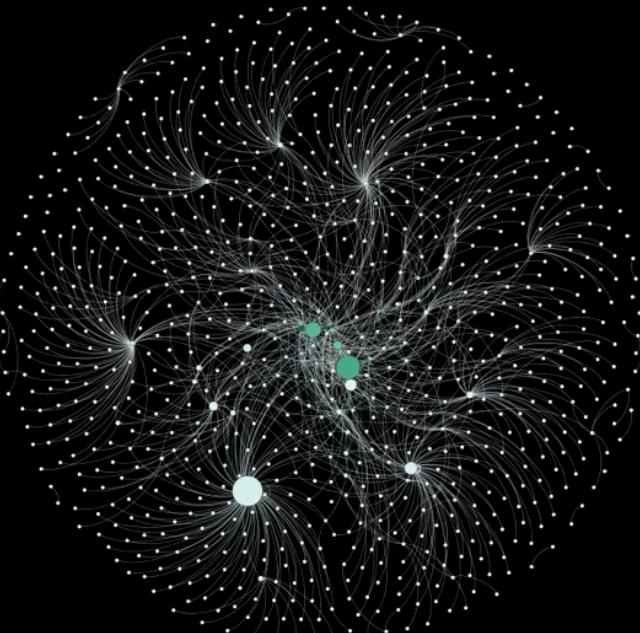


IF YOU HAVE AN ARCHIVE-IT
ACCOUNT, TRY IT OUT YOURSELF AT
CLOUD.ARCHIVESUNLEASHED.ORG

cloud.archivesunleashed.org

Archives Unleashed Cloud

New user? You currently need [Archive-It](#) credentials to use our service. Please read our [FAQ](#).



The Andrew W. Mellon Foundation

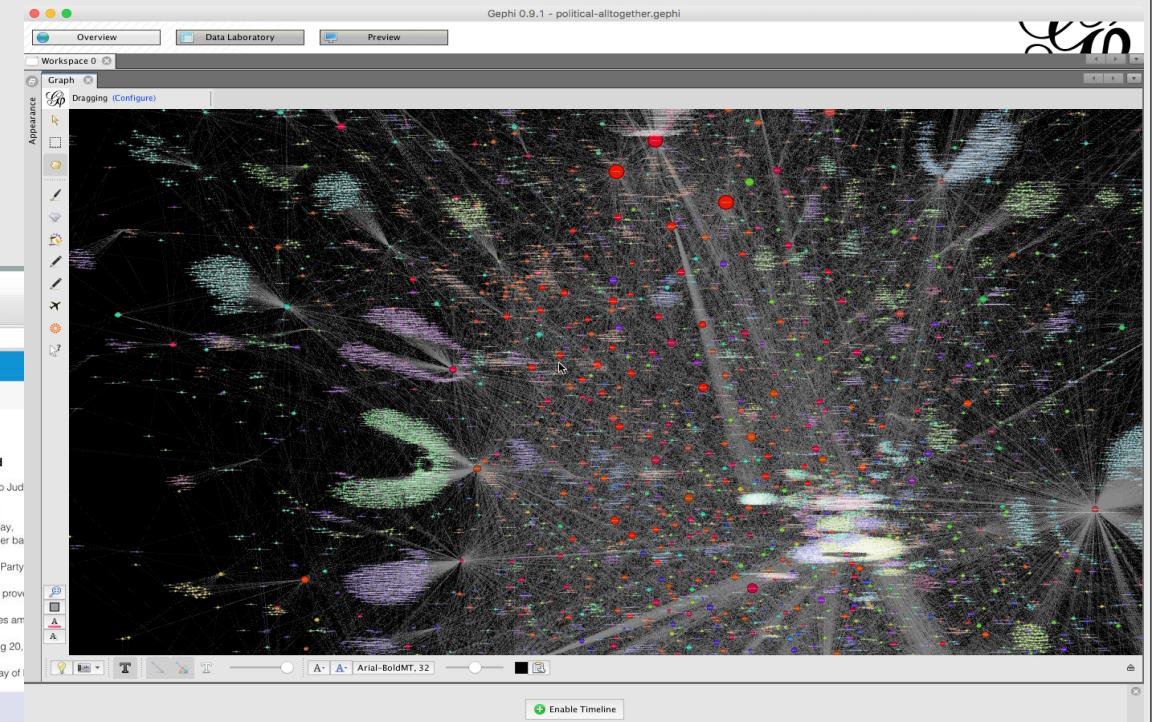
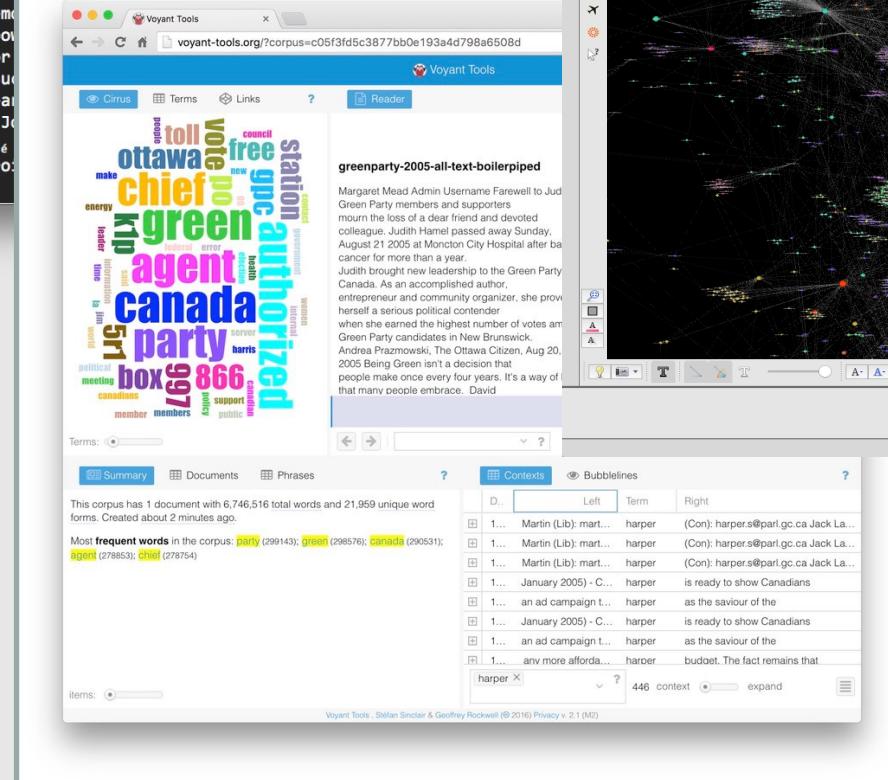
UNIVERSITY OF
WATERLOO

YORK UNIVERSITY

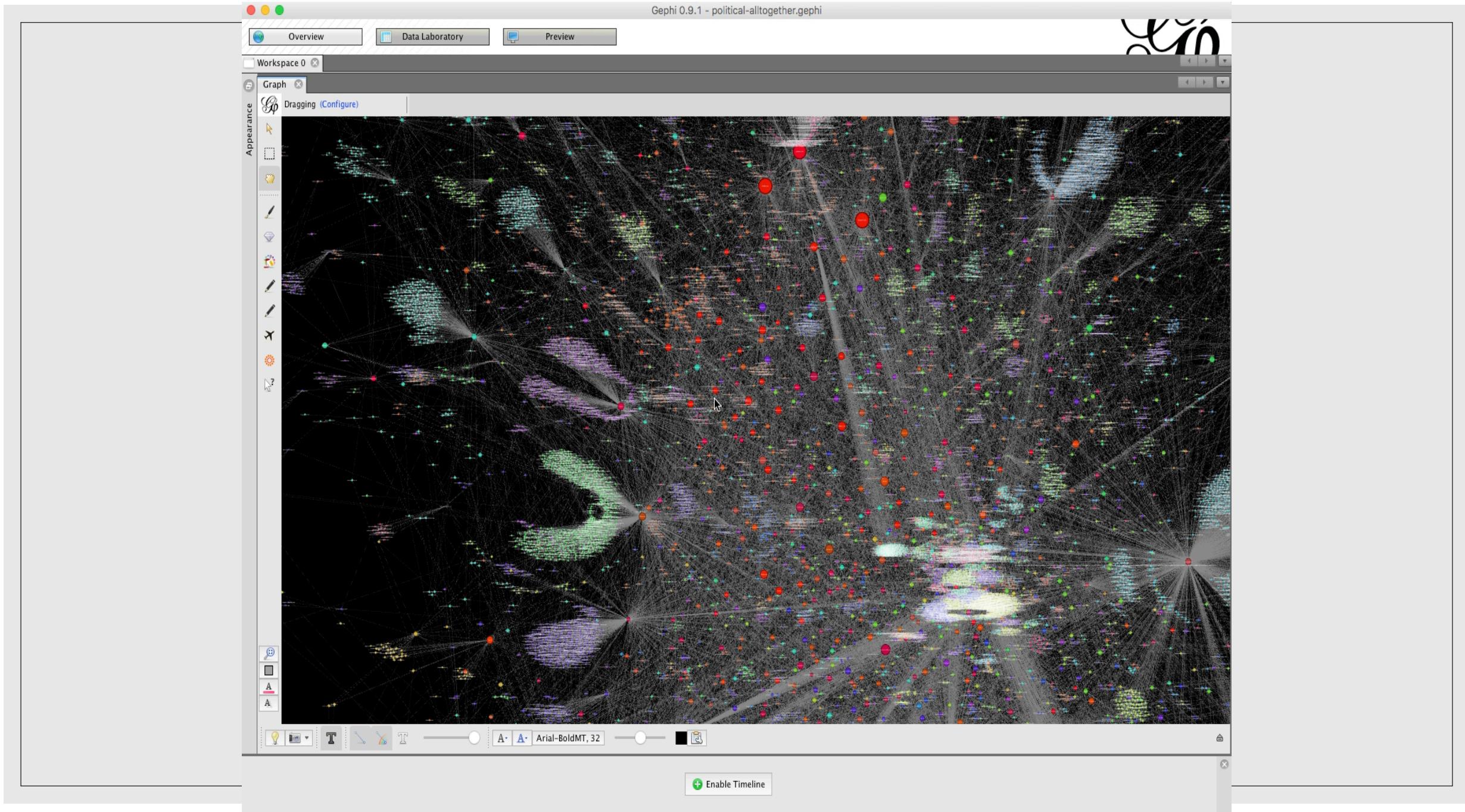
For more information on our project and sponsors, visit [archivesunleashed.org/](#).

[About](#) | [Privacy Policy](#) | [Documentation](#) | [FAQ](#)

1 ssh
 {20150805,communist-party.ca,http://communist-party.ca/,HTTP/1.1 200 OK Date: Wed, 05 Aug 2015 22:09 :10 GMT Server: Apache X-Pingback: http://communist-party.ca/xmlrpc.php Link: ; rel=shortlink x-frame-options: SAMEORIGIN Vary: Accept-Encoding,User-Agent Connection: close Content-Type: text/html; charset=UTF-8 Communist Party of Canada - Parti Communiste du Canada Parti Communiste du Québec People's Voice Newspaper The Spark! Theory Journal Young Communist League International Communist Movement Communist Party of Canada - Parti Communiste du Canada About Short history The Figueroa case Our Constitution Party Program Our aim is socialism Capitalism in Canada Canada in a changing world The Canadian state, nations and peoples of Canada, and the crisis of democracy The working class and people's struggle For a People's Government Building Socialism The Communist Party Campaigns Repeal Bill C-51! For a People's Recovery! Save Canada Post Hands off Syria! Contact us How to donate Join the Communist Party We need to move Canada to the left! Nous devons recentrer le Canada vers la gauche! Vote Communist! The Harper Conservative government's austerity agenda has been a disaster for the people. Meanwhile, the profits of the big banks and the largest corporations have ballooned. Most people are now convinced that Harper and his gang must go. What is really required in Parliament are truly progressive voices who will challenge right-wing, pro-corporate policies. We need to move Canada in a new, progressive direction, with a People's Alternative Platform based on the needs of working people and nature! Everywhere, the gap is widening between the majority of the people and a handful of the super-rich. Today, the top 1% own and control over 50% of the entire wealth of our planet. Racism and intolerance are also spreading. The crisis we face isn't just about the economy; it is about capitalism. It is time capitalism was replaced with a democratic system - with socialism. A new society, based on working class power, towards true democracy - the rule of the people, by the people and for the people. Such economic wealth are owned and controlled by the working people. Such a collective struggle. Voting Communist sends the strongest and clearest message: another world is possible, urgent and worth fighting for. Join us for a better future! Read more. Votons communiste! Les politiques d'austérité, flâne pour la population canadienne qui constate du même coup la croissance "part-00000" 851L, 4611113C



SO WHAT CAN YOU DO
WITH LINKS?



LET'S DO A QUICK
GEPHI
WALKTHROUGH.

HACKING THE NEWS



WANT TO KNOW MORE?

[HTTPS://CLOUDARCHIVESUNLEASHED.ORG/DERIVATIVES](https://cloudarchivesunleashed.org/derivatives)

HACKING THE NEWS