

Archives Unleashed!



Collections as Data | September 27, 2016 | Library of Congress

Problem: Web archives are a tremendous untapped resource for scholars, but many are unaware of how archives can be utilized

Opportunity: Develop a forum to engage and train scholars to work with archived Web data

Web archives are an amazing resource
but access and usage are often a problem.



**KEEP
CALM
AND
HACK**





It takes an army...



compute | calcul
canada | canada

LIBRARY OF
CONGRESS



UNIVERSITY OF
TORONTO



Library and Archives
Canada

Bibliothèque et Archives
Canada



Social Sciences and Humanities
Research Council of Canada

Conseil de recherches en
sciences humaines du Canada

Canada



LIBRARY OF CONGRESS

Format:

- **2- to 3-day event**
 - **Day 1:** Introductions, networking, team formation
 - **Day 2 – 3:** Work sessions
 - **Day 3:** Final presentations & awards



Archives Unleashed 1.0

March 3-5, 2016

University of Toronto

Toronto, CA

The image shows the grand interior of the Library of Congress. The architecture is highly ornate, featuring large white columns, arched doorways, and intricate ceiling decorations. In the foreground, a group of people is gathered around a digital display. A woman in a black dress is presenting to the group. The lighting is warm and focused on the display area.

Archives Unleashed 2.0

June 14- 15, 20 16

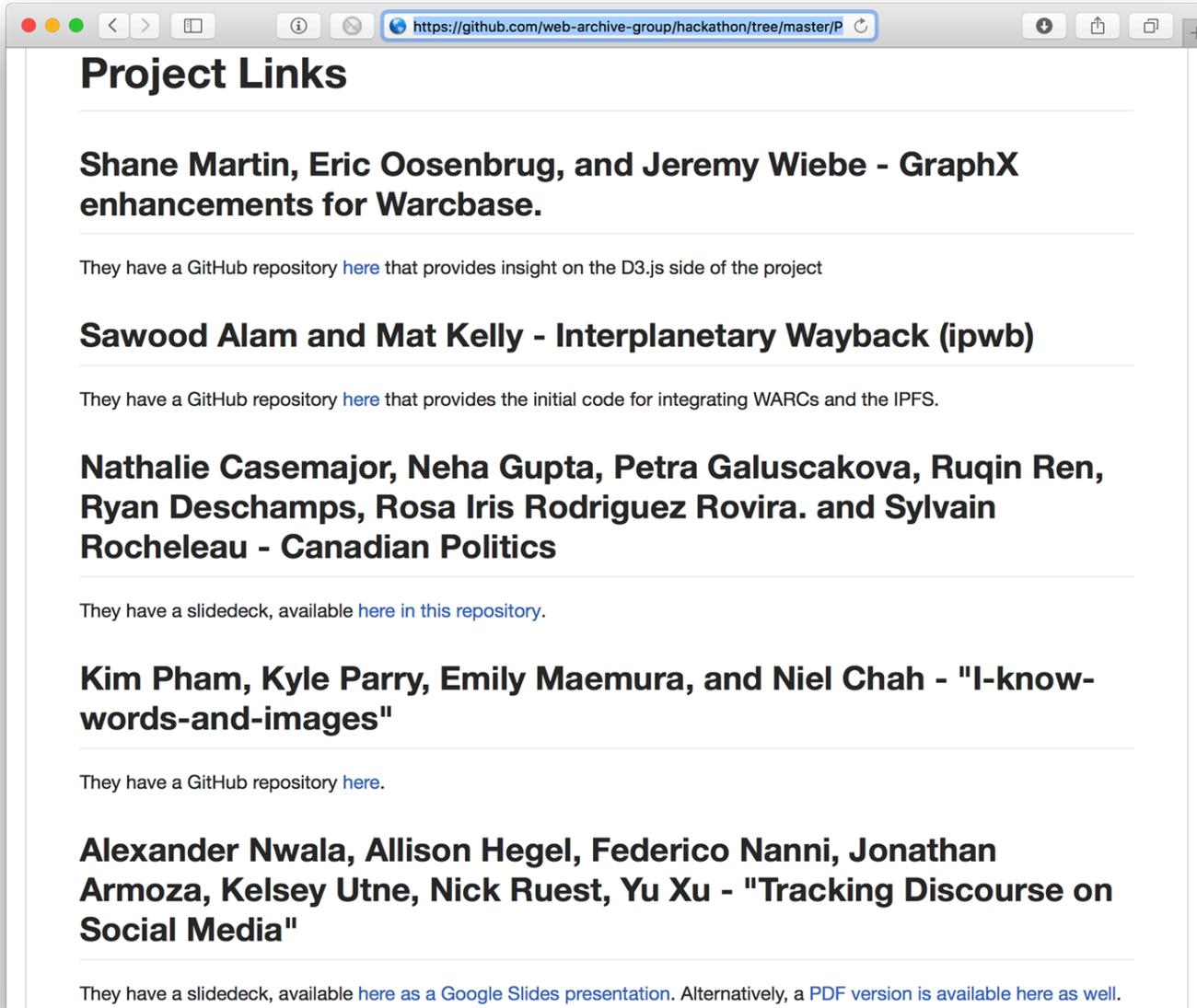
Library of Congress

Washington, D.C.

Putting the Plan into Action



Outcomes



The screenshot shows a web browser window with the address bar containing the URL <https://github.com/web-archive-group/hackathon/tree/master/P>. The page content is titled "Project Links" and lists five projects, each with a title and a brief description.

Project Links

Shane Martin, Eric Oosenbrug, and Jeremy Wiebe - GraphX enhancements for Warcbase.

They have a GitHub repository [here](#) that provides insight on the D3.js side of the project

Sawood Alam and Mat Kelly - Interplanetary Wayback (ipwb)

They have a GitHub repository [here](#) that provides the initial code for integrating WARCs and the IPFS.

Nathalie Casemajor, Neha Gupta, Petra Galuscakova, Ruqin Ren, Ryan Deschamps, Rosa Iris Rodriguez Rovira. and Sylvain Rocheleau - Canadian Politics

They have a slidedeck, available [here in this repository](#).

Kim Pham, Kyle Parry, Emily Maemura, and Niel Chah - "I-know-words-and-images"

They have a GitHub repository [here](#).

Alexander Nwala, Allison Hegel, Federico Nanni, Jonathan Armoza, Kelsey Utne, Nick Ruest, Yu Xu - "Tracking Discourse on Social Media"

They have a slidedeck, available [here as a Google Slides presentation](#). Alternatively, a [PDF version is available here as well](#).



Day One: Socialize



Day Two - Three: Work



Day Two - Three: Work



Day Two - Three: Working and Sharing

The Final Countdown:

Talking Points by Place and Party in the last days of the 2004 Election

Research Question:

Which places and topics did the two presidential candidates mentioned most in the days leading up to the 2004 election?

Method:

1. Parse and clean the text from WARCs of party/candidate websites
2. Extracted Named Entities (candidate names, state names)
3. Topic modeling (LDA, TF-IDF)
4. Semantic Web - Google Knowledge Graph API (for next time)

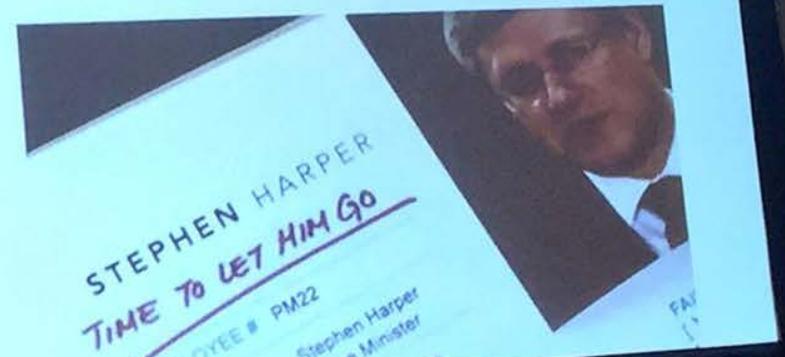
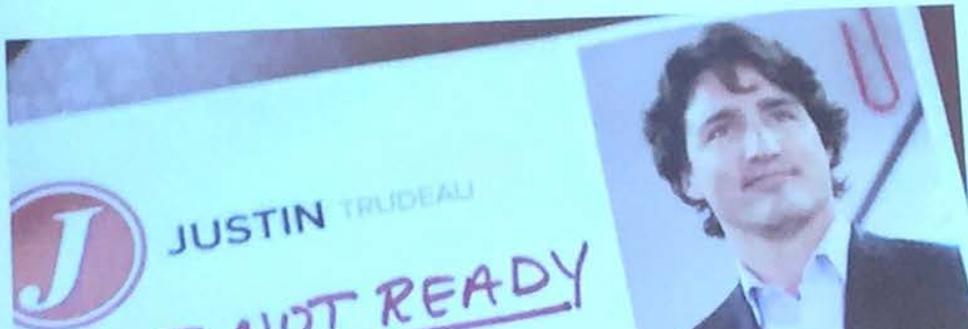
Day Two - Three: Working and Sharing



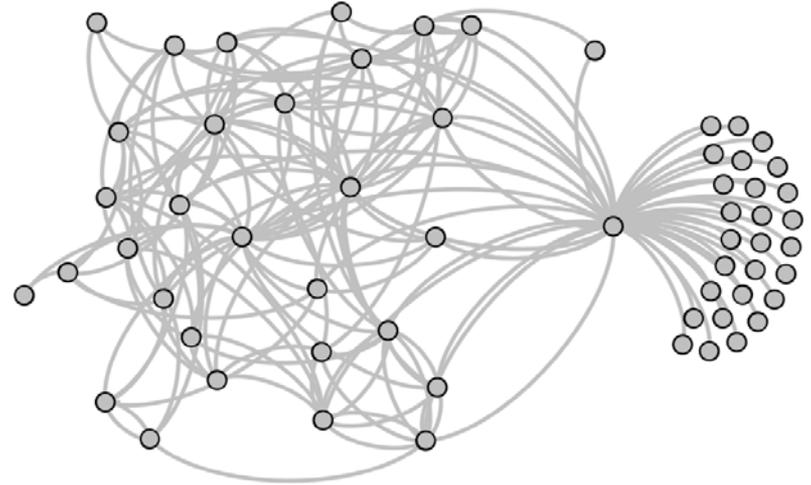
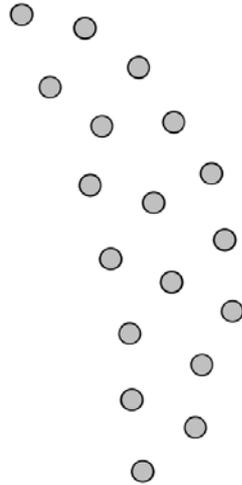
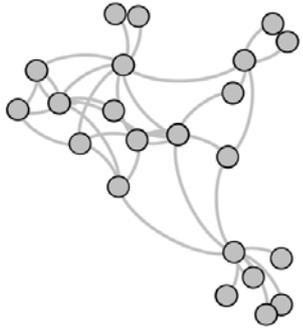
Day Three: Share!

IMAGES: NEGATIVE ADS

How often do political parties
use images
of their opponents?



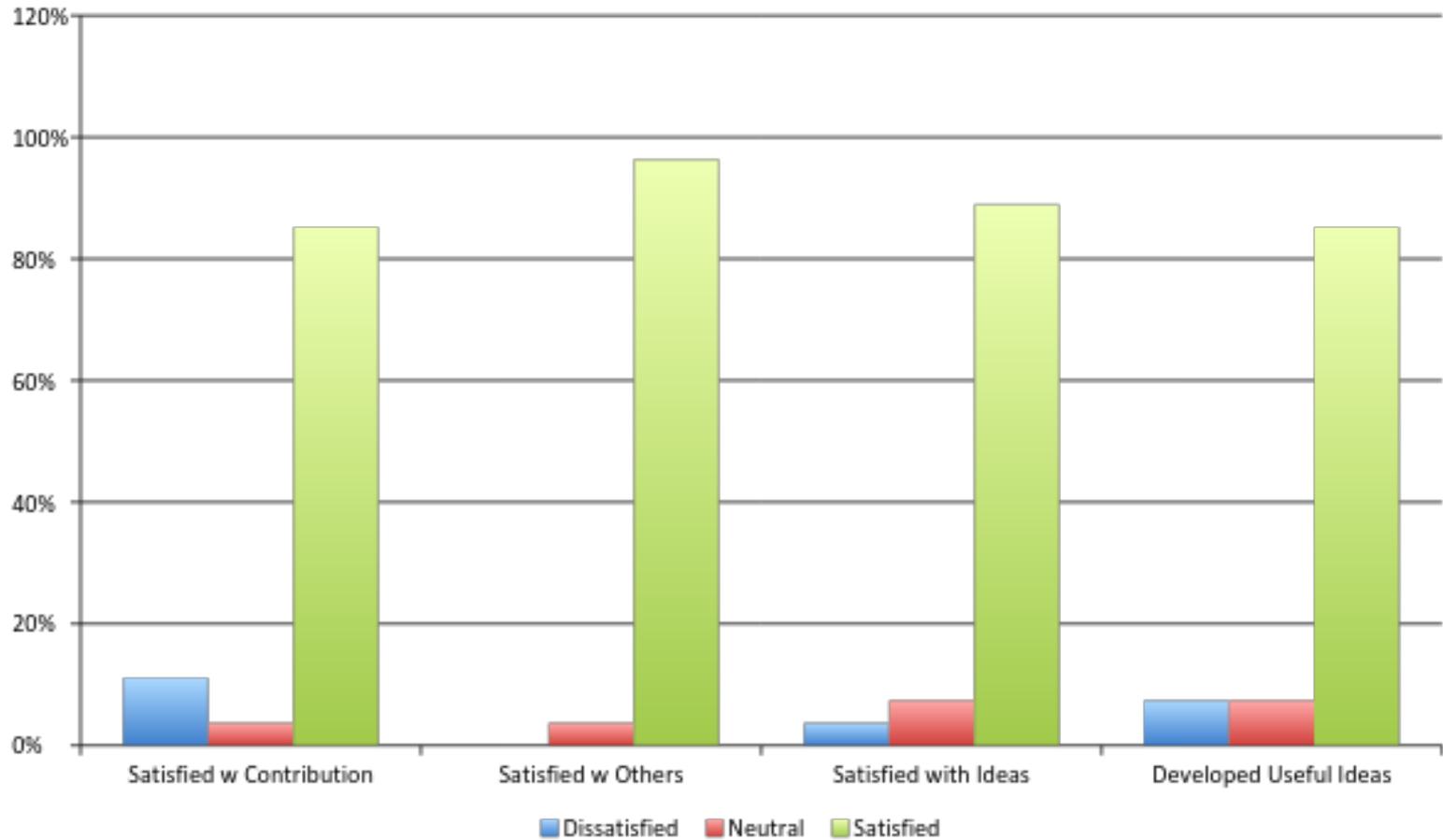
Forming Connections



Exchanging ideas *before...*

and after.

Satisfaction with Outcomes



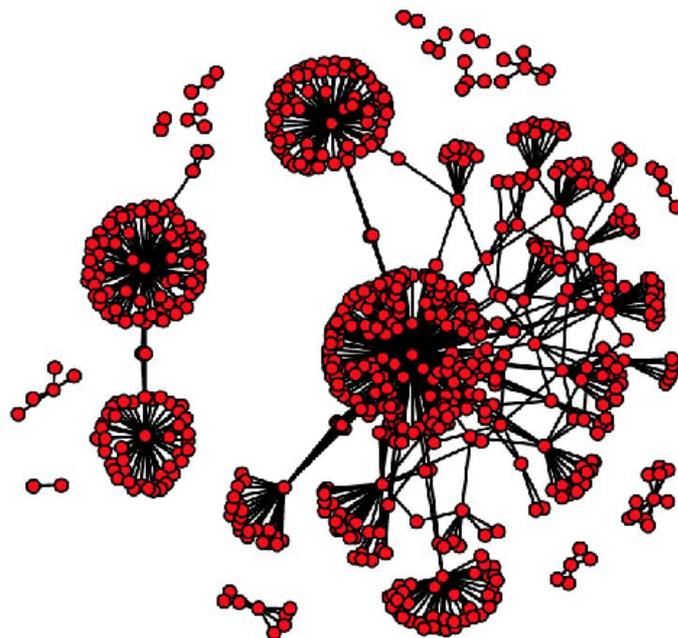
Results

- Patterns appear to validate Program on Extremism anecdotal analysis that there are "generators" and "amplifiers"
- Generated ideology scores on small sample

Next Steps

- Visualize relationship weights
- Distinguish between retweets vs. mentions
- Time-evolve relationship graphs
- Correlate with ideology scores
- Reconsider stack: Too computationally intensive for R on Mac → leverage cloud instance w/ map-reduce cluster (e.g. use `pig`)

January 2015
Mention-relationships



Archives Unleashed 3.0

February, 20 17

Internet Archive

San Francisco, CA



Archives Unleashed 4.0

June 20 17

British Library

London, UK



Creating a Sustainable Model



Contact:

Matthew S. Weber

matthew.weber@rutgers.edu

@docmattweber

Archives Unleashed 2.0: archivesunleashed.com

Archives Unleashed 1.0 Projects: <http://bit.ly/1ScAwaa>