

History in the Age of Abundance:

Skills, Tools, and Methods
for the 21st-Century
Historian

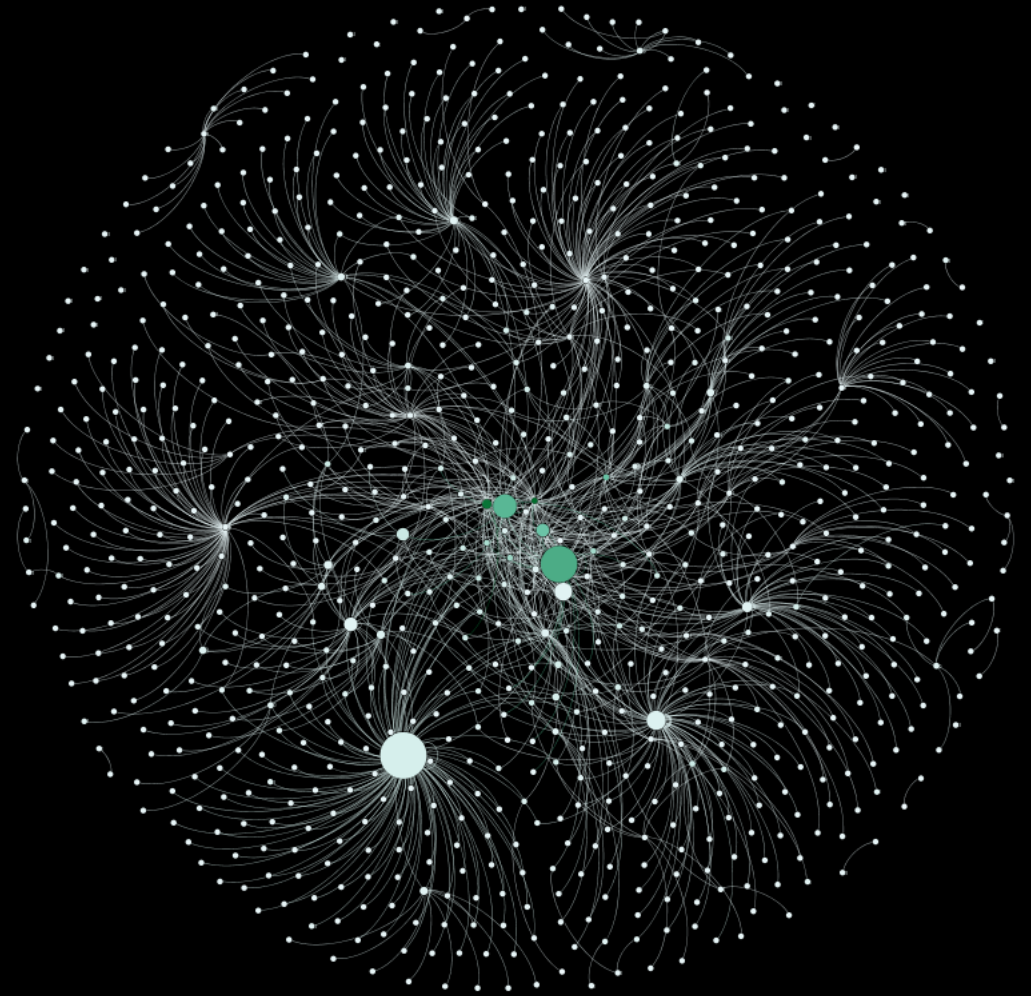
Ian Milligan

Associate Professor

Department of History



UNIVERSITY OF
WATERLOO

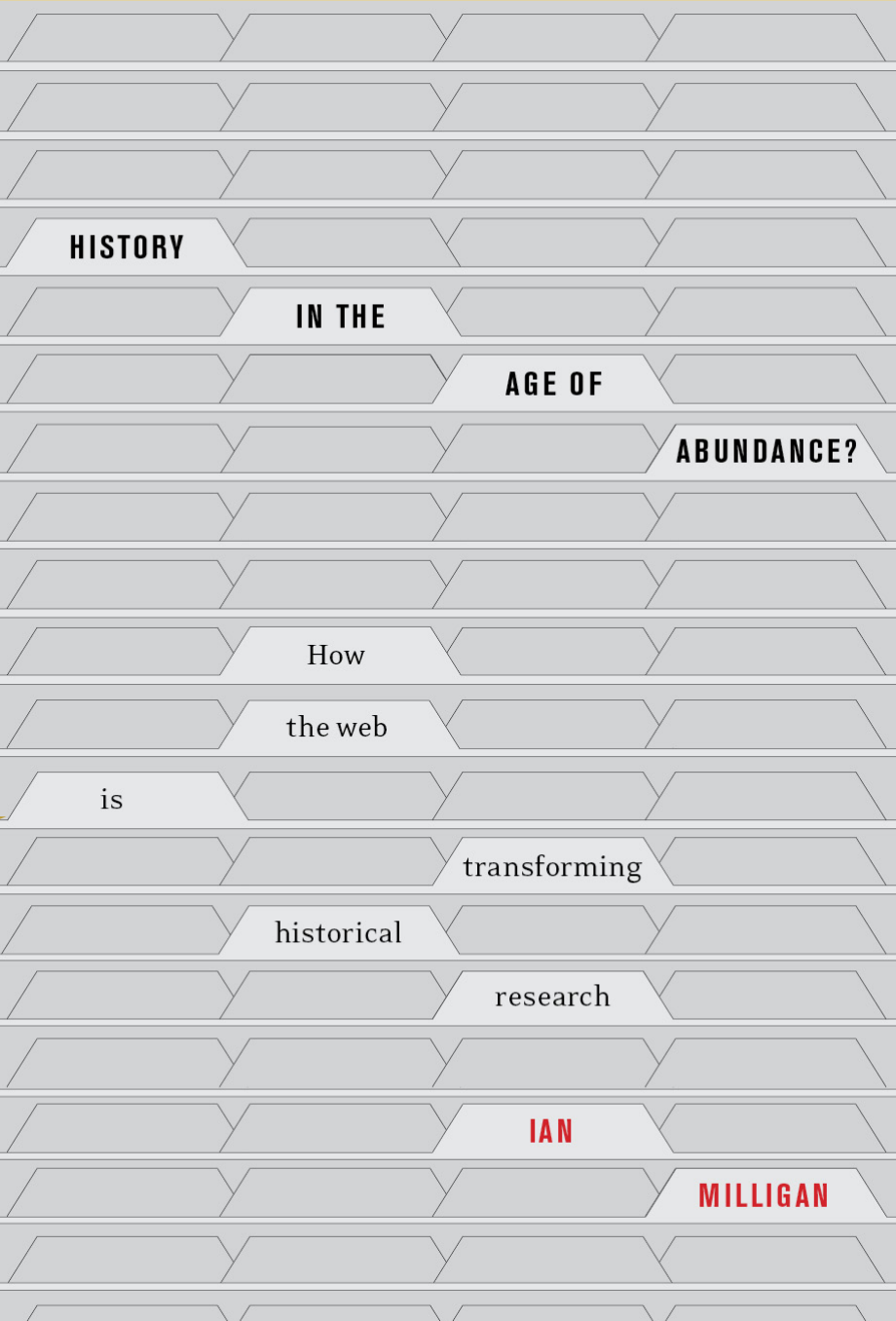


SLIDES PDF AVAILABLE AT

<https://www.ianmilligan.ca/talk/york-hist/>

Shameless Plug!

- *History in the Age of Abundance? How the Web is Transforming Historical Research.* McGill-Queen’s University Press, 2019.
- \$28 on Amazon, or there’s a copy at YUL or a few at Toronto Public Library



PLAN FOR THE TALK

- The Problem of TMI
- The Broad Solution to this Problem
- The Current Status of this Problem
- What to do through three personas
 - Computational Social Scientist/Humanist
 - Digital Social Scientist/Humanist
 - Social Scientist/Humanist
- Conclusions around history in the 21st century
- Questions



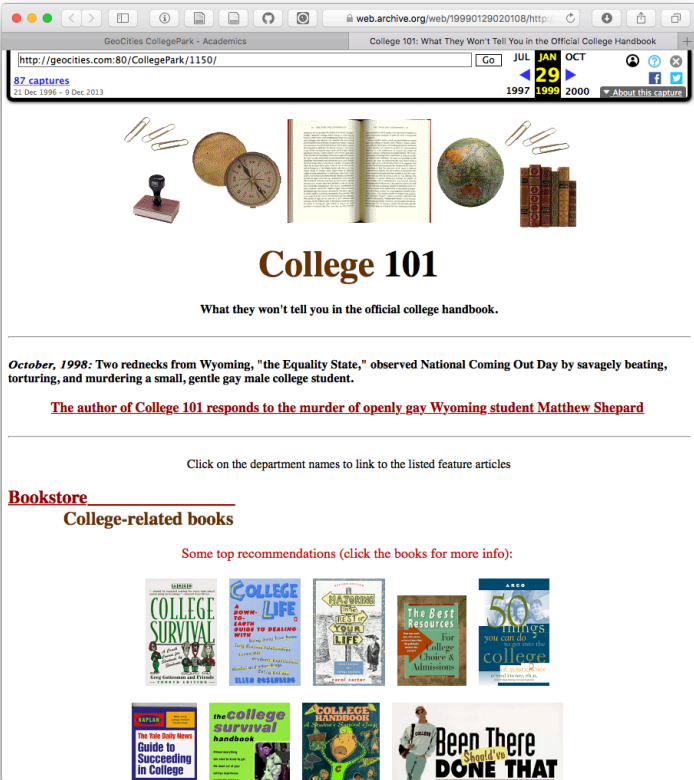
WHAT'S THE PROBLEM OF HISTORY IN THE “AGE OF ABUNDANCE”?

THE PROBLEM

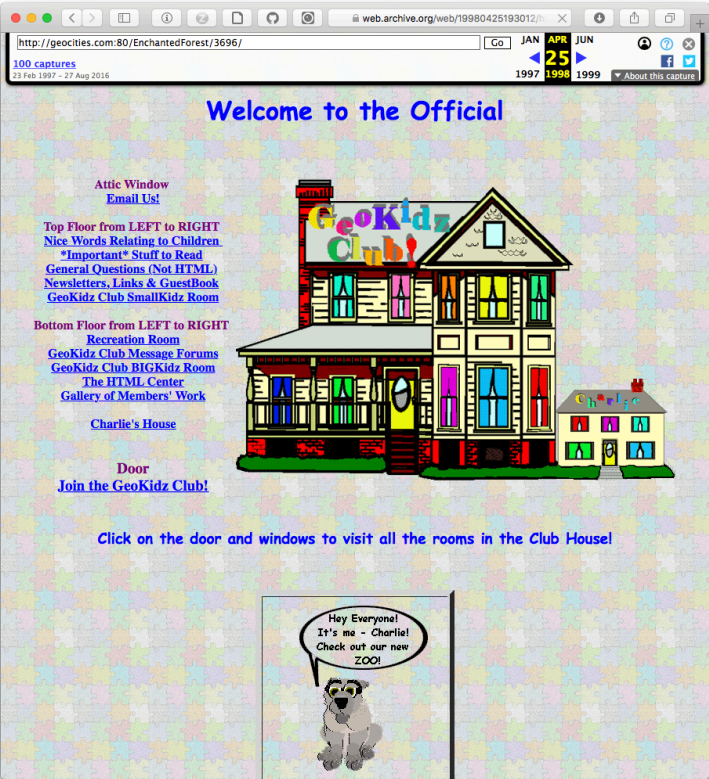


LAC Preservation Centre, Gatineau QC

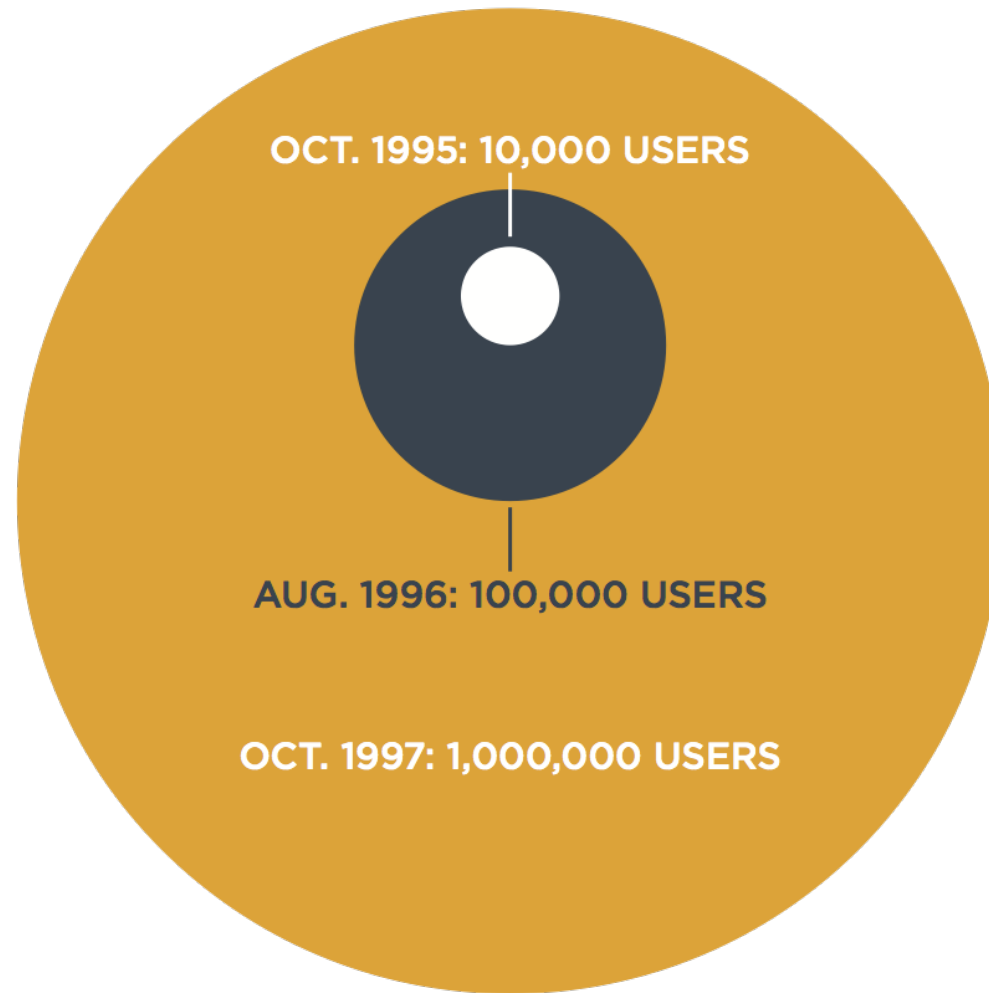
THE PROBLEM



GeoCities.com




THE PROBLEM



In other words...


**WEB ARCHIVES REPRESENT A NEW MEANS
OF KNOWLEDGE DISTRIBUTION AND
ACQUISITION**



A PERSONAL EXAMPLE OF OUR CHANGING HISTORICAL RECORD

[rec.games.miniatures.warhammer](#) ›
Civil War in Warhammer 40,000
7 posts by 7 authors 

★ **Nyarlathotep** Who can say, but it certainly sounds amusing. Later days, Nyar.

12/7/95

 **Peter Milligan**

12/7/95  

★ - hide quoted text -

I was recently playing a game of Warhammer 40,000 when the opponents had a disagreement and decided that they would have a civil war. The two armies in question were Space Wolf, Blood Angels and Imperial Guard. The Space Wolf player opened up with his Cyclone and using 12 missiles, wiped out the 1500 point Imperial army, and the survivors were picked off by the Blood Angel's thud gun.

Although I was delighted by this, it made for a short game, as the only model I got to use was my Leman Russ, who picked off the survivors after the Blood angels picked off the Space Wolves.

Do you think Civil War is an acceptable rule?

★ **Bub** Well If you're in another universe why not. But in terms of 40K no :) – Bub 7262...@compuserve.com

12/8/95

★ **Joel E Slovacek** once in a while, so long as the people fighting the war understand that they will lose (since they have

12/8/95

★ **Chad Lubrecht** What is this civil war rule? I seem to have missed the original posting, sorry!



12/8/95

A PERSONAL EXAMPLE OF OUR CHANGING HISTORICAL RECORD

- **First**, Don't Read the Post;
- **Secondly**, my discomfort is perhaps in and of itself significant;
- **Third**, the musings of an eleven-year old child are now available..
- **Fourth** – in many countries today, such musings would be subject to legal deposit.

[rec.games.miniatures.warhammer](#) >

Civil War in Warhammer 40,000

7 posts by 7 authors  

★ **Nyarlathotep** Who can say, but it certainly sounds amusing. Later days, Nyar. 12/7/95

 **Peter Milligan** 12/7/95  

★ - hide quoted text -

I was recently playing a game of Warhammer 40,000 when the opponents had a disagreement and decided that they would have a civil war. The two armies in question were Space Wolf, Blood Angels and Imperial Guard. The Space Wolf player opened up with his Cyclone and using 12 missiles, wiped out the 1500 point Imperial army, and the survivors were picked off by the Blood Angel's thud gun.

Although I was delighted by this, it made for a short game, as the only model I got to use was my Leman Russ, who picked off the survivors after the Blood angels picked off the Space Wolves.

Do you think Civil War is an acceptable rule?

★ **Bub** Well If you're in another universe why not. But in terms of 40K no :) – Bub 7262...@compuserve.com 12/8/95

★ **Joel E Slovacek** once in a while, so long as the people fighting the war understand that they will lose (since they have 12/8/95

★ **Chad Lubrecht** What is this civil war rule? I seem to have missed the original posting, sorry! 12/8/95



UNIVERSITY OF
WATERLOO

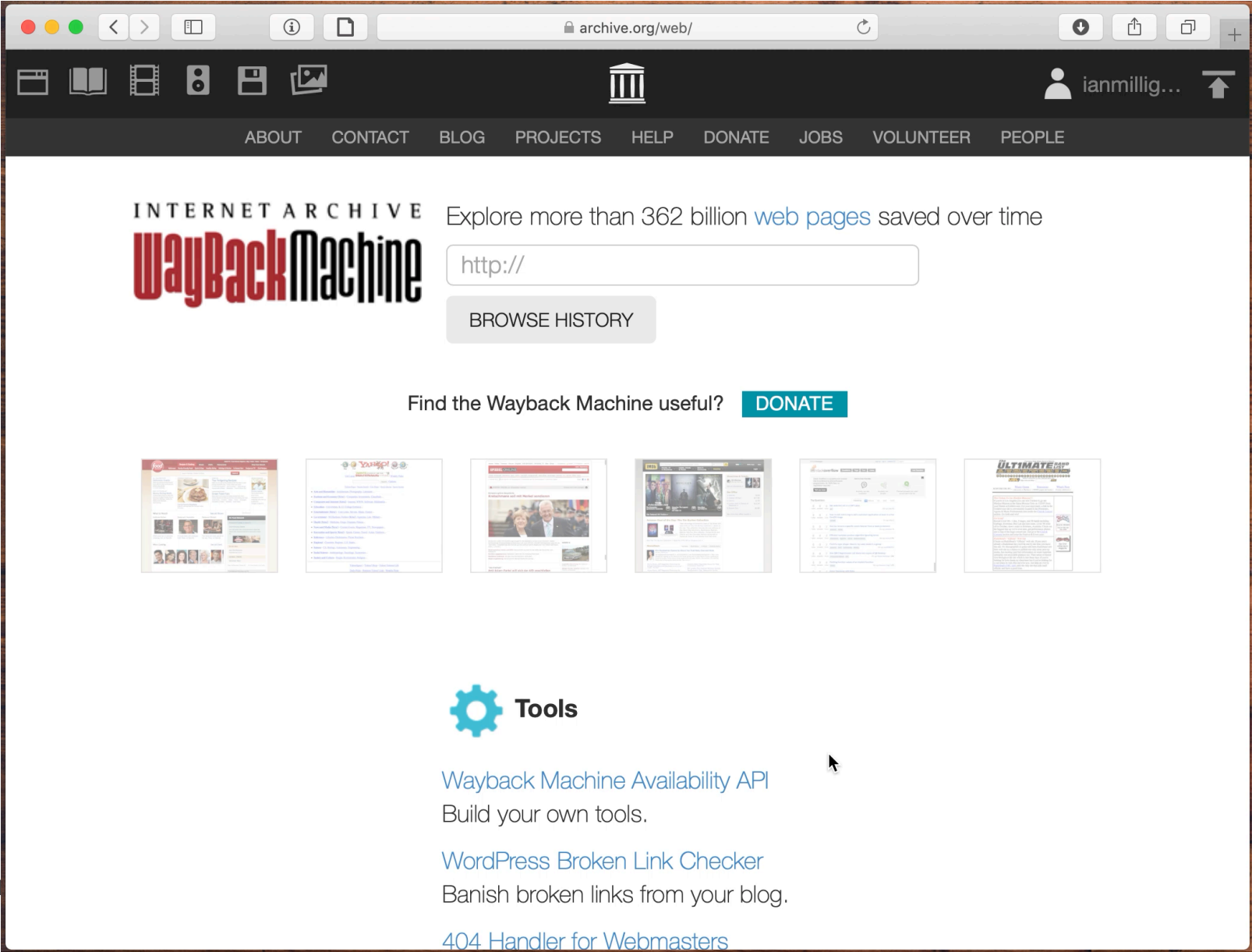
**THIS IS NOT A NICHE CONCERN FOR ANY
POST-1990 TOPIC**

WHAT OPPORTUNITIES DO THEY PRESENT?

- The way that we preserve our culture is changing;
 - **Scale:** Internet Archive has 635 billion+ URLs; 40PB of unique data (and non-Internet Archive collectors probably have about the same again).
 - **Scope:** Data that never before would have been collected is now being collected about people who aren't traditionally in the historical record.
- Any researcher tackling post-1996 topics will realistically need to understand the vast arrays of text, image, etc. that comprise our modern cultural record.
- The **Wayback Machine** isn't enough; will need to explore data at scale.

WHO HERE HAS USED THE INTERNET ARCHIVE?

THE WAYBACK MACHINE IN ACTION



HISTORY IN THE AGE OF ABUN



DOWNSIDERS OF THE WAYBACK MACHINE

- **Wayback Machine** is great if you know what you're looking for;
 - **Ever-improving keyword search functionality**
 - Represents a great stride in accessibility more generally
- But it isn't great for more detailed research queries:
 - You may want to do complicated queries (i.e. websites that say X and link to Y);
 - You may want to do exploratory text mining;
 - You may want to work with images en masse;
 - Etc.

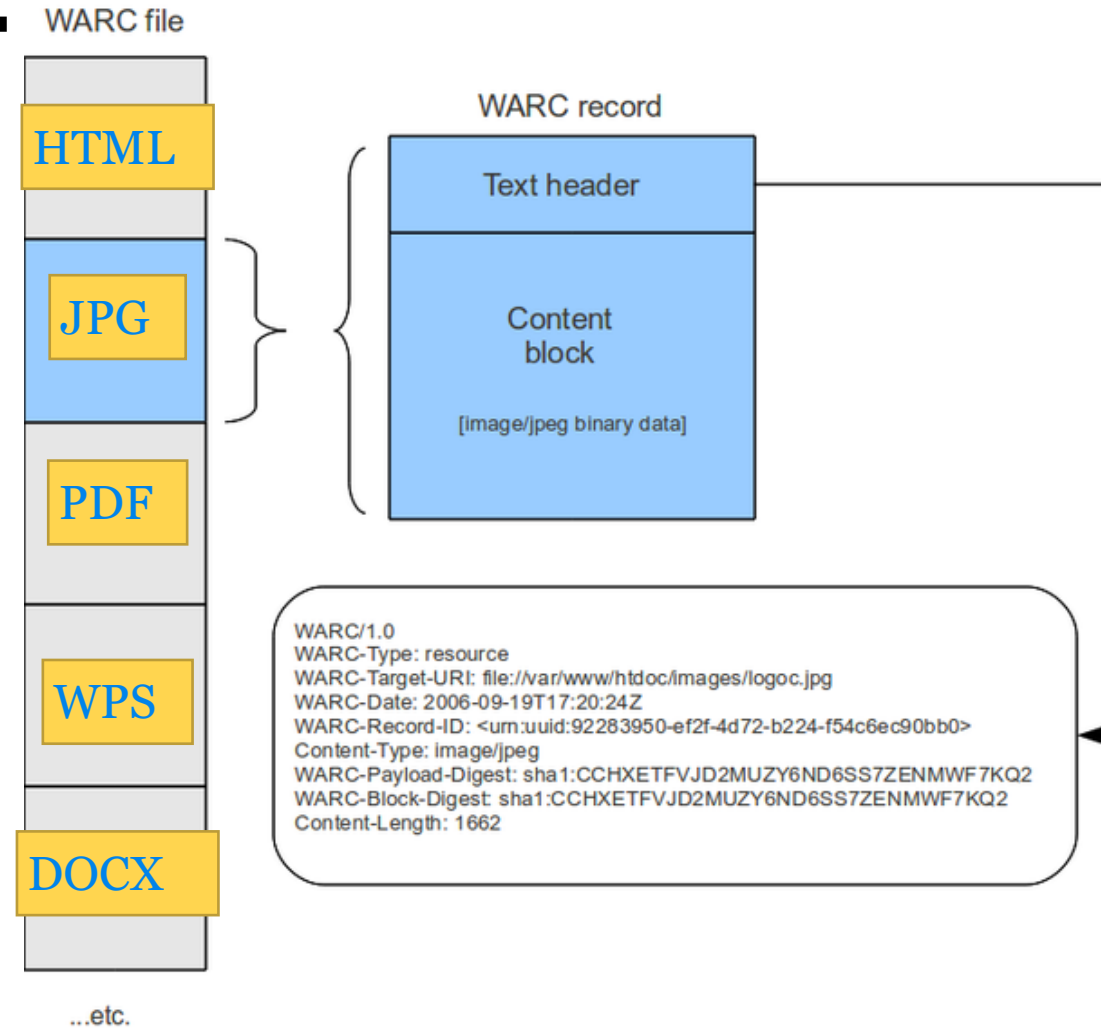


**SO INSTEAD, WE'LL NEED TO WORK WITH
DATA AT SCALE.**

What does this entail?

WORKING WITH DATA AT SCALE

- The **WARC** (ISO 28500:2009) file
- Pictured at right
- Hard to use and a bit idiosyncratic, with a smaller user base, so the first step is to usually transform the data into something that's a bit more common!

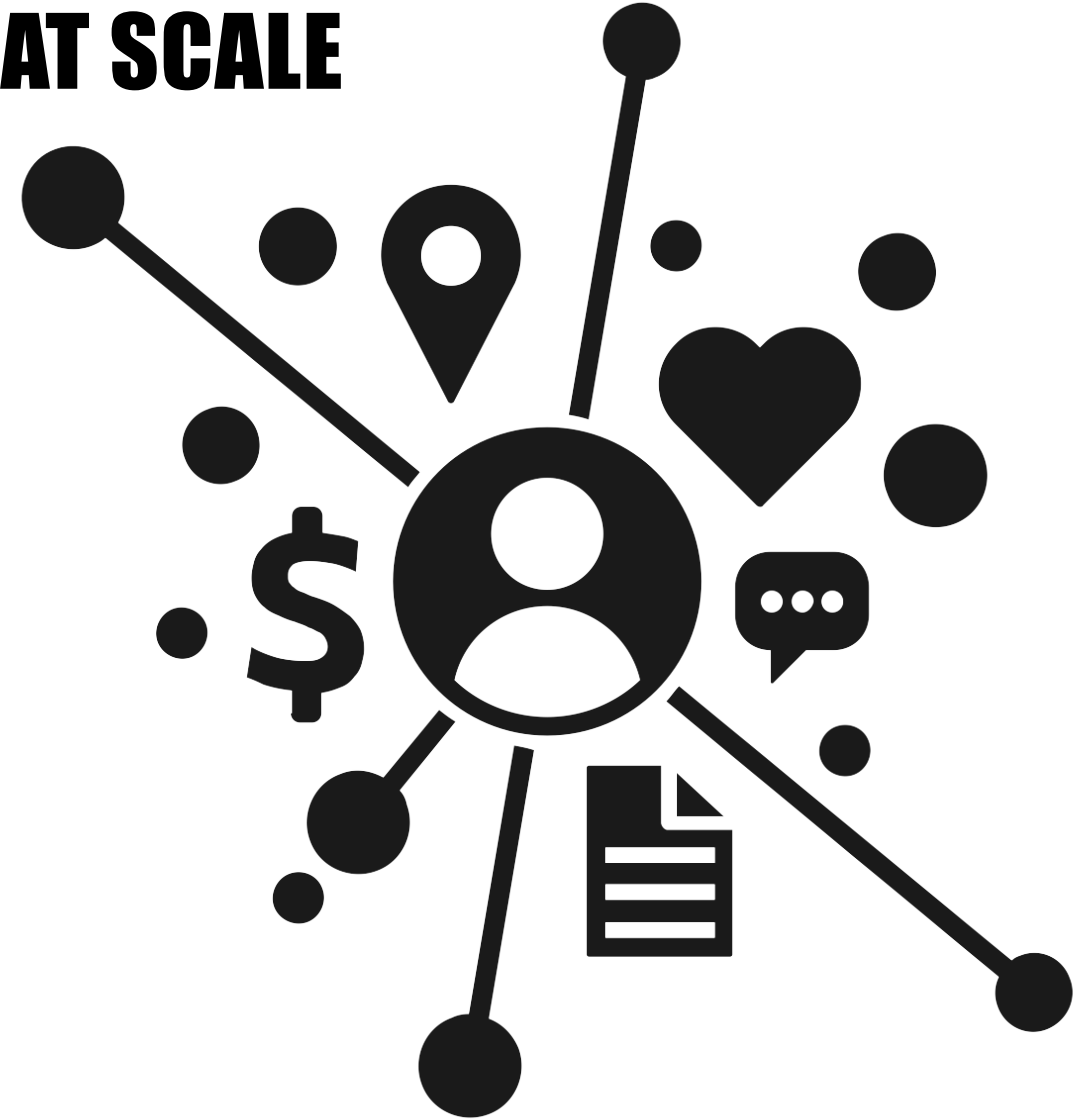


YOUR COMPREHENSIVE EXAMS PREPPED YOU FOR WORKING WITH DATA, RIGHT?

Uh oh... what kind of skills are you going to need?

SKILLSET ONE: WORKING WITH DATA AT SCALE

- An understanding of:
 - The Basics of **Natural Language Processing** (NLP)
 - **Basic Statistical Knowledge** to work with Quantitative Data (frequency of websites, terms, normalizing numbers, etc.)
 - Flexible **Data Science** skills (or StackOverflow skills..)
- In other words, being equipped with the skills and capacities to analyze text/data at scale.



SKILLSET TWO: UNDERSTANDING HOW DATA IS CONSTRUCTED

- They also need to have a solid understanding of:
 - **How and why the data was collected**, i.e. selection criteria;
 - What data **wasn't collected**;
 - How the **software used to create the dataset** has changed over time;
- How to **clean or normalize data** when necessary (i.e. as a URL changes from <http://www.ndp.ca> to <https://www.ndp.ca> to <http://ndp.ca> in the crawl, recognizing that those are probably the same website! You see that, the computer doesn't unless you teach it to).



OH YEAH...

**PLUS THE NORMAL SKILLS OF THE
HISTORIAN TOO..**

SO LET'S TAKE STOCK

- **Historians will need to understand and study the Web** in order to come to grips of history after the mid-1990s – not just for history of the Web, of course, but for the history of our society and culture as reflected on the Web
- **Existing tools like the Wayback Machine aren't enough** to tackle this problem
- **Historians will need new skills** for working with and understanding data, plus their traditional competencies



ARE HISTORIANS READY?

(Leading question alert!)

NO...
(Sad trombone)

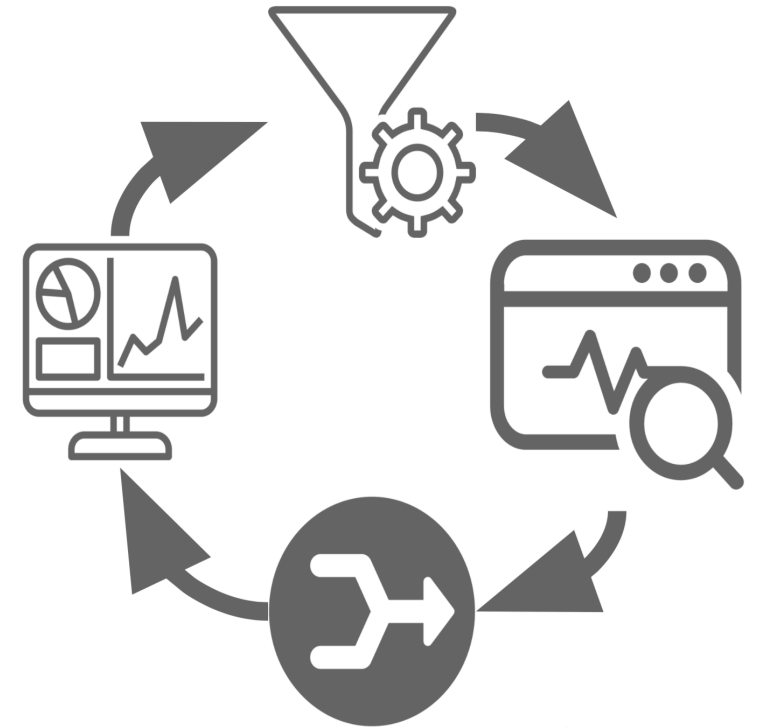
FIRST OF ALL, “WE” CAN’T DO IT ALONE

- I’ll return to this point, but when I use “we” it’s not the **royal we**; it’s the Archives Unleashed “we”
 - Nick Ruest (York University Libraries)
 - Jimmy Lin (Waterloo Computer Science)
 - + undergrads in CS; grads in CS and History
- I’ll return to some smiling mugshots in a few



THE CURRENT STATUS OF WORKING WITH WEB ARCHIVES AT SCALE?

- We have developed something called the **Archives Unleashed Toolkit**, based on working with scholars at seven datathons + numerous workshops
- It allows scholar to use the **FAAV** cycle to explore it:
 - **Filter**
 - **Analyze**
 - **Aggregate**
 - **Visualize**
- Allows for complicated research queries using Apache Spark



PROBLEM SOLVED, RIGHT?

Check out our cutting-edge interface....

THE ARCHIVES UNLEASHED TOOLKIT

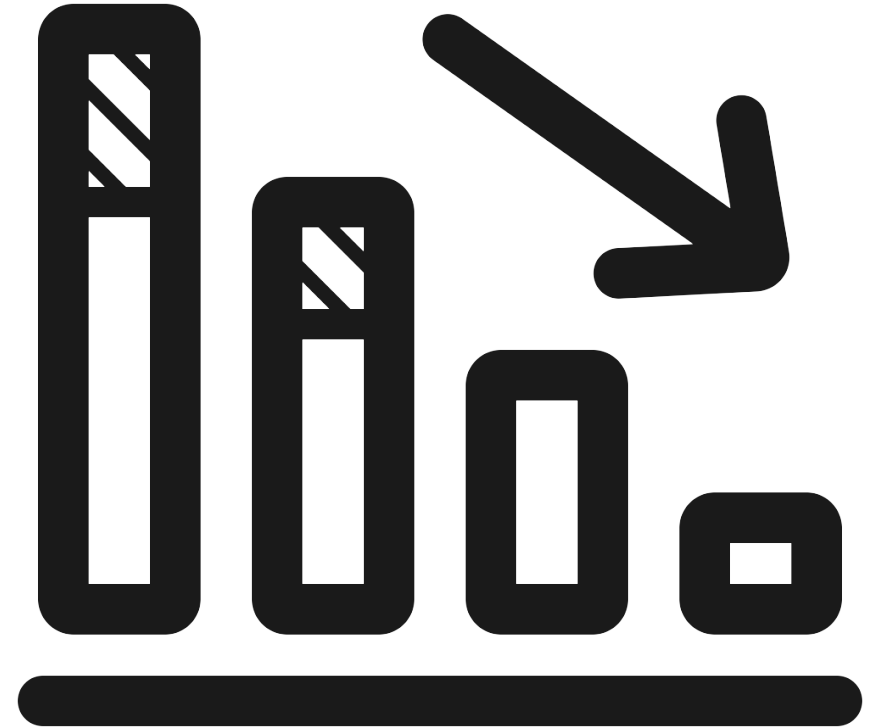
- It's easy to use, as long as you:
 - Know how to use the command line;
 - How to access a server;
 - How to use the Spark Shell;
 - How to code, at least somewhat, in Scala;
 - And, have a lot of patience for open-source documentation!

(OK, it's not easy to use..)

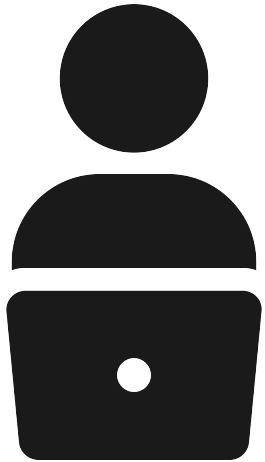
```
import io.archivesunleashed._ import
io.archivesunleashed.matchbox._
RecordLoader.loadArchives("example.a
rc.gz", sc) .keepValidPages()
.keepDate(List("200804"),
ExtractDate.DateComponent.YYYYMM)
.map(r => (r.getCrawlDate,
r.getDomain, r.getUrl,
RemoveHTML(r.getContentString)))
.saveAsTextFile("plain-text-date-
filtered-200804/")
```

... AND HISTORIANS HAVE BEEN TURNING AWAY FROM QUANT TO QUAL

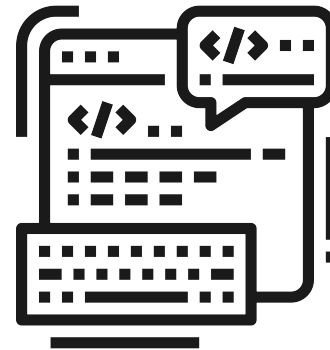
- Studies of our introductory historiography textbooks show this diminishing.
 - John Tosh, *Pursuit of History*
 - 1st, 2nd edition: “History by Numbers” (entire chapter)
 - By 5th edition, no quantitative history at all.
- **“Nevertheless, it is curious that at a time when both the use of and the breadth of humanities data is growing, quantitative skills ... seem to no longer form a core component of our undergraduate history programmes.”** (James Baker, <https://blogs.bl.uk/digital-scholarship/2014/04/digital-history-and-the-death-of-quant.html>)



WHERE TOOLS END AND USERS BEGIN

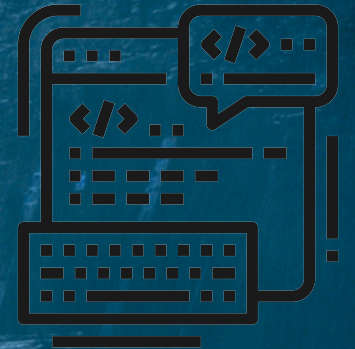


**Historian with
Research Question**



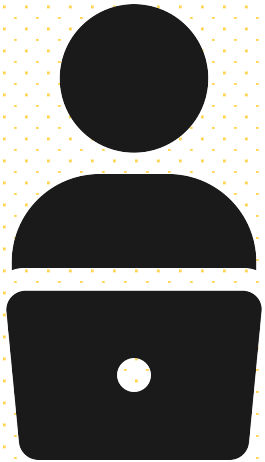
**Advanced
Web Archive Analysis**

WHERE TOOLS END AND USERS BEGIN



A now-dated cultural reference, thanks for killing the magic,
Game of Thrones Season 8.

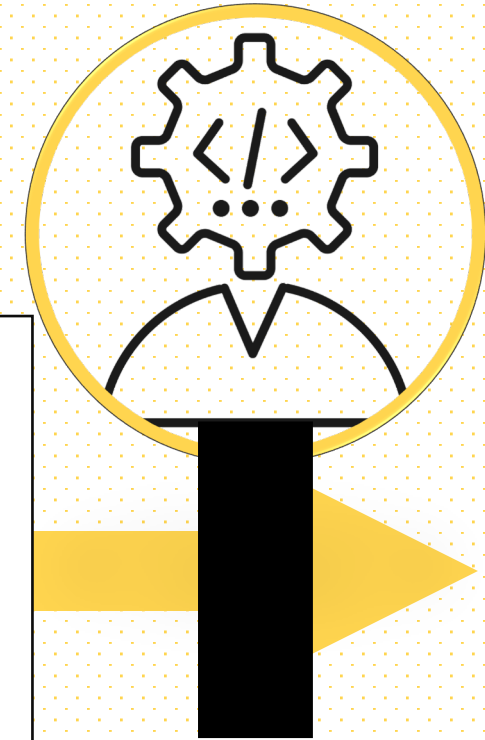
WHERE TOOLS END AND USERS BEGIN



**Historian with
Research Question**

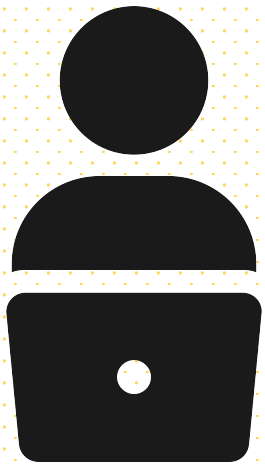
OPTION ONE: Historians become programmers, capable of advanced web archive analysis.

- Writing code;
- Contributing to open-source projects;
- Maybe need a bit of support but in general are self-sufficient in a computational environment.

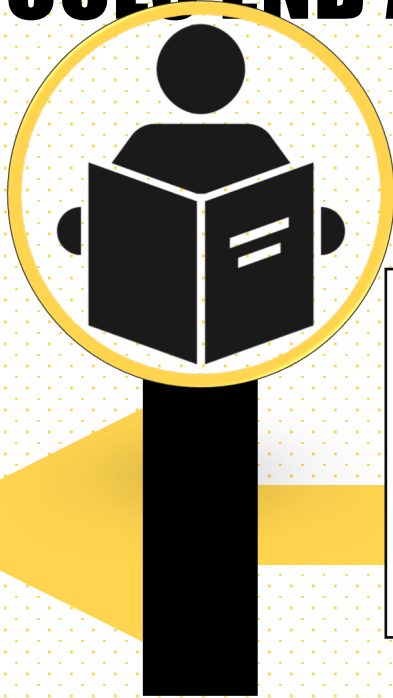


**Advanced
Web Archive Analysis**

WHERE TOOLS END AND USERS BEGIN

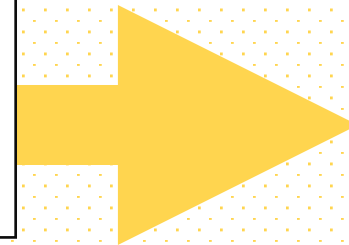


**Historian with
Research Question**



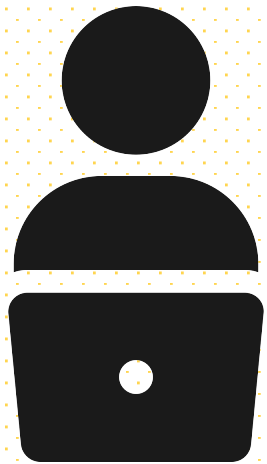
OPTION TWO: Historians do not develop technical skills, continuing solely as subject-matter experts.

- An approach that sees using a web archive as similar to using PROQUEST or JSTOR – a bit of work but nothing out of the ordinary.



**Advanced
Web Archive Analysis**

WHERE TOOLS END AND USERS BEGIN



**Historian with
Research Question**



OPTION THREE: The middle ground

- Some computational skills, but platforms designed to resemble conventional research processes as much as possible
- Positioned in a world where you can draw on standard library research support (i.e. not web archiving particularly, but “working with text”)
- No command lines!



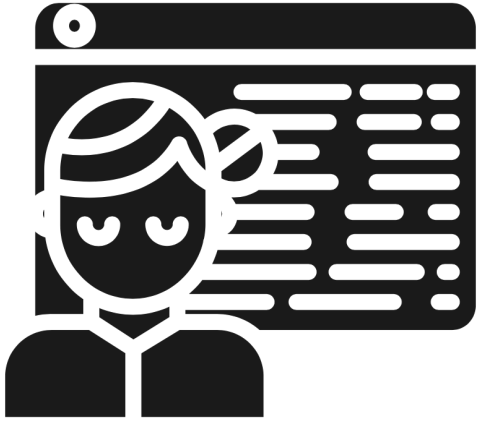
**Advanced
Web Archive Analysis**

WHERE TOOLS END AND USERS BEGIN

- We can draw on a model of successful **interdisciplinary cooperation** to see how this can be effected
- If a computer scientist and a historian work together, it isn't just "historian: now you become a computer scientist;" nor is it "computer scientist: now you become a historian." Compromise is worked out in terms of:
 - Workflow (Platforms? Collaboration?)
 - Publication venues
 - The shape of the work
- It really is often meeting half way, and being conscious

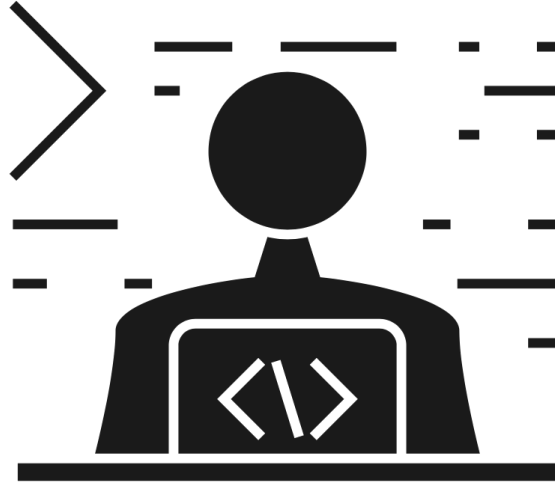


THREE PERSONAS



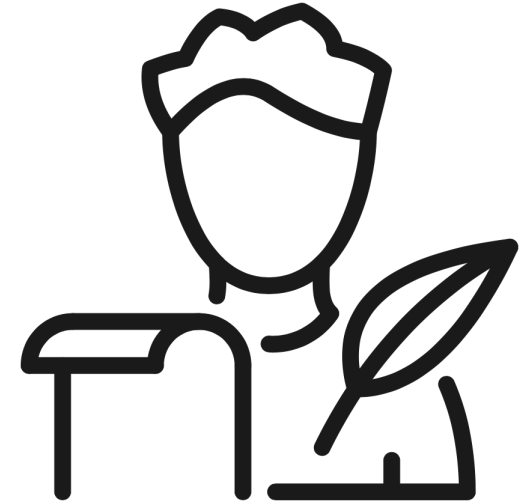
Person A

Computational
Humanist



Person B

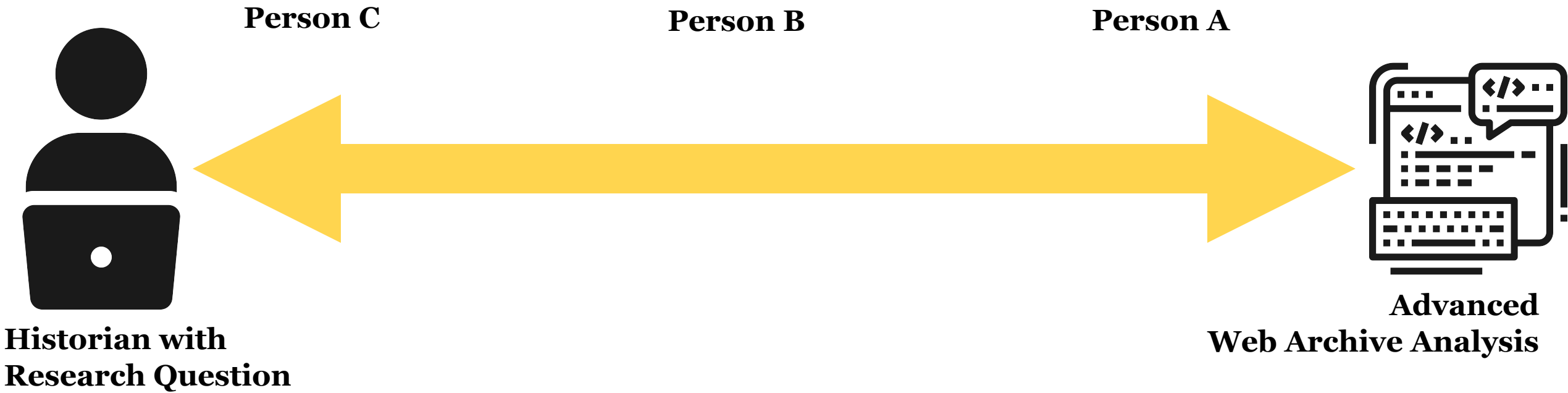
Digital Humanist



Person C

Conventional
Historian

IN OTHER WORDS...



**HOW CAN WE SUPPORT THIS
COMPUTATIONAL TURN FOR ALL THREE
PERSONAS?**

Enter our project...

Archives Unleashed

A white line-art icon on a black background. It features a rectangular briefcase with a handle. A test tube is positioned diagonally against the right side of the briefcase. Above the test tube, a burst of fireworks or sparks is depicted with several short lines radiating from a central point.

FIRST, WE NEED AN INTERDISCIPLINARY TEAM



Ian Milligan

Historian



Nick Ruest (YORK!)

Library/Archives



Jimmy Lin

Computer Science

WE CAN'T GO IT ALONE...!

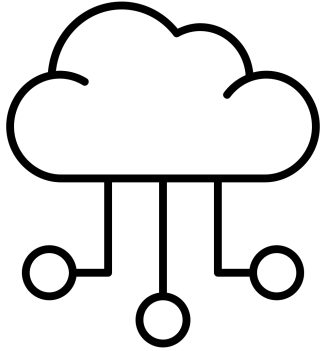
(So why would we expect our users to go it alone??)

AND A MISSION

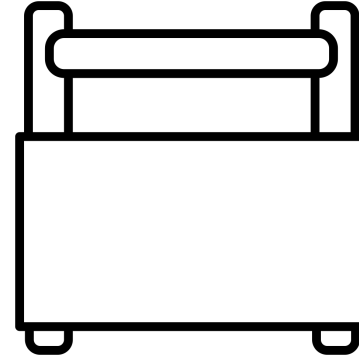
Archives Unleashed aims to make petabytes of historical internet content accessible to scholars and others interested in researching the recent past.

ARCHIVES UNLEASHED PROJECT

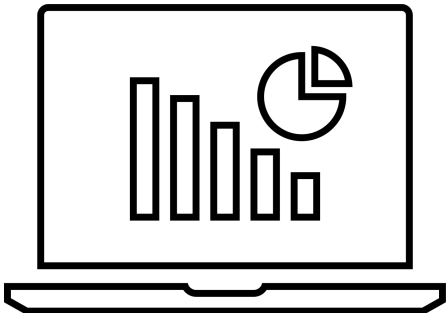
AND A SET OF TOOLS...



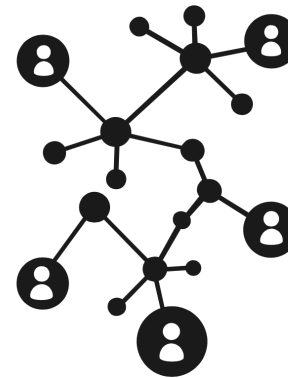
Cloud



Toolkit



Notebooks

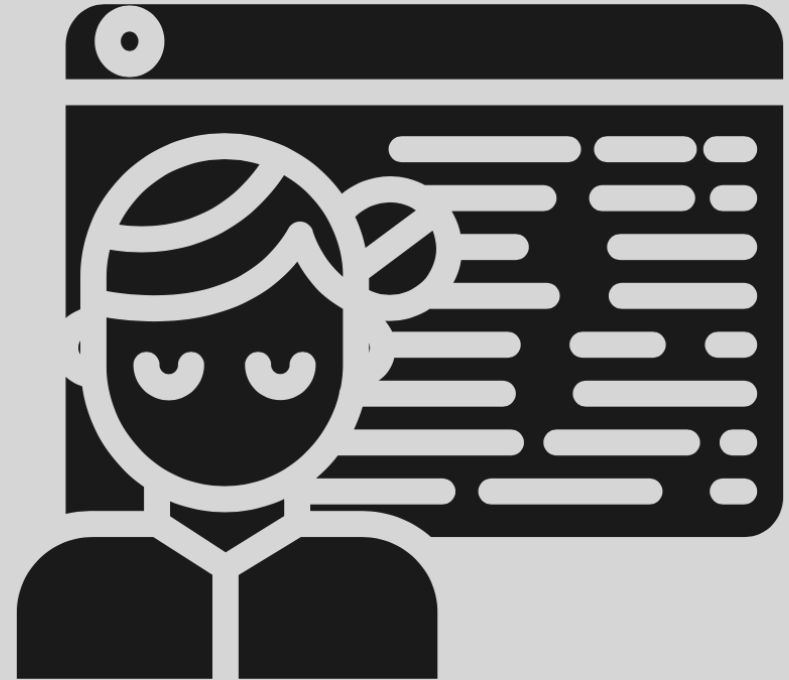


Datathons

PERSONA A

The sort of scholar we would have called a
“computational humanist” or
“computational social scientist.”
Comfortable installing packages,
understanding dependencies, fluent on
the command line, and can Stack
Overflow like a master.

THE COMPUTATIONAL HISTORIAN



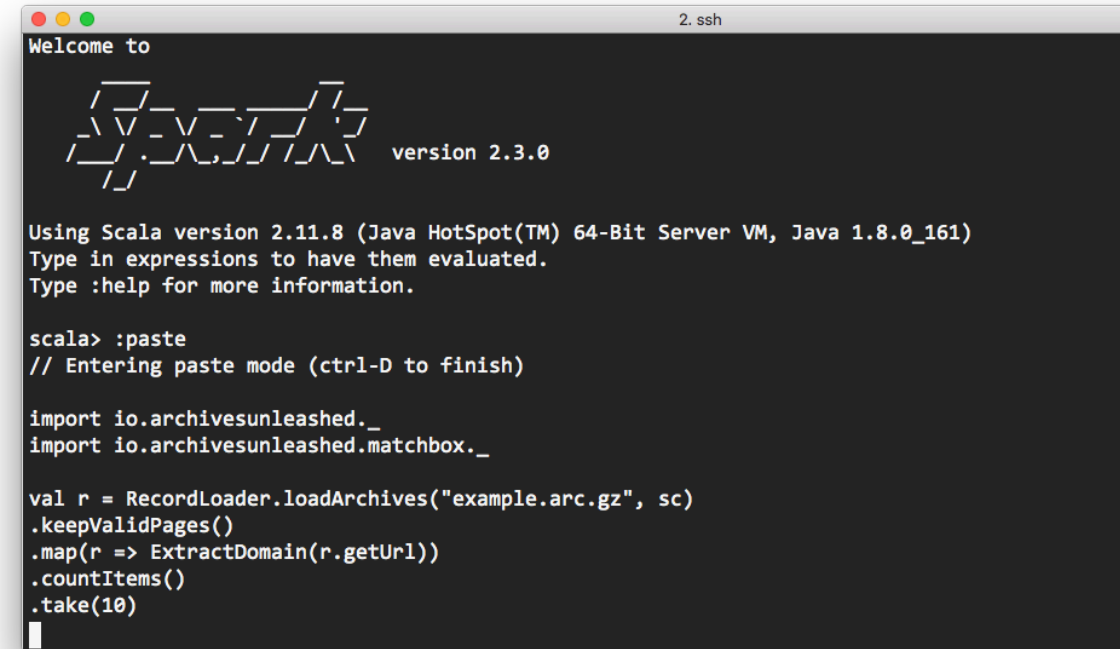
COMPUTATIONAL HUMANITIES

- Needs to have some sort of toolkit to translate WARC files into something they can work with
- Can interpret tested yet still dense documentation
- Can troubleshoot error messages (i.e. knows what to Google – this is a hard-earned skill)
- Can take sample scripts, execute them to see outcome, and then change them to run on their own data with own questions
- Can use the command line!

```
import io.archivesunleashed._ import
io.archivesunleashed.matchbox._
RecordLoader.loadArchives("example.arc.g
z", sc) .keepValidPages()
.keepDate(List("200804"),
ExtractDate.DateComponent.YYYYMM) .map(r
=> (r.getCrawlDate, r.getDomain,
r.getUrl,
RemoveHTML(r.getContentString)))
.saveAsTextFile("plain-text-date-
filtered-200804/")
```

COMPUTATIONAL HUMANITIES

- The **Archives Unleashed Toolkit** was designed around this persona.
- 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



```
2. ssh
Welcome to
Archives Unleashed Toolkit version 2.3.0

Using Scala version 2.11.8 (Java HotSpot(TM) 64-Bit Server VM, Java 1.8.0_161)
Type in expressions to have them evaluated.
Type :help for more information.

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

import io.archivesunleashed._
import io.archivesunleashed.matchbox._

val r = RecordLoader.loadArchives("example.arc.gz", sc)
  .keepValidPages()
  .map(r => ExtractDomain(r.getUrl))
  .countItems()
  .take(10)

```

LET'S SEE THIS IN ACTION...

🍏

Safari

File Edit View History Bookmarks Develop Window Help

🔍

archivesunleashed.org

🔄

📄

📁

+

☰

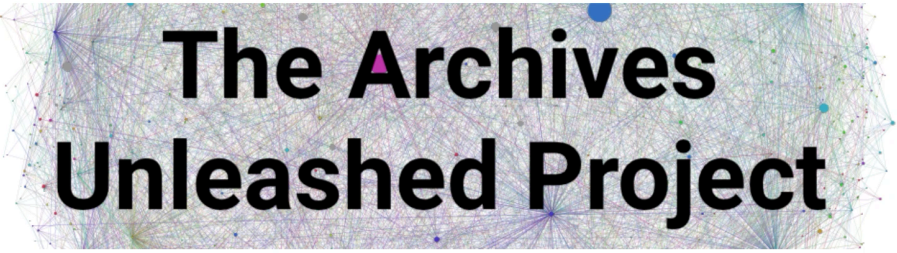
Welcome to the Archives Unleashed Project

🐦

🐙

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 email. Are you a Slack user? Join our Slack team!

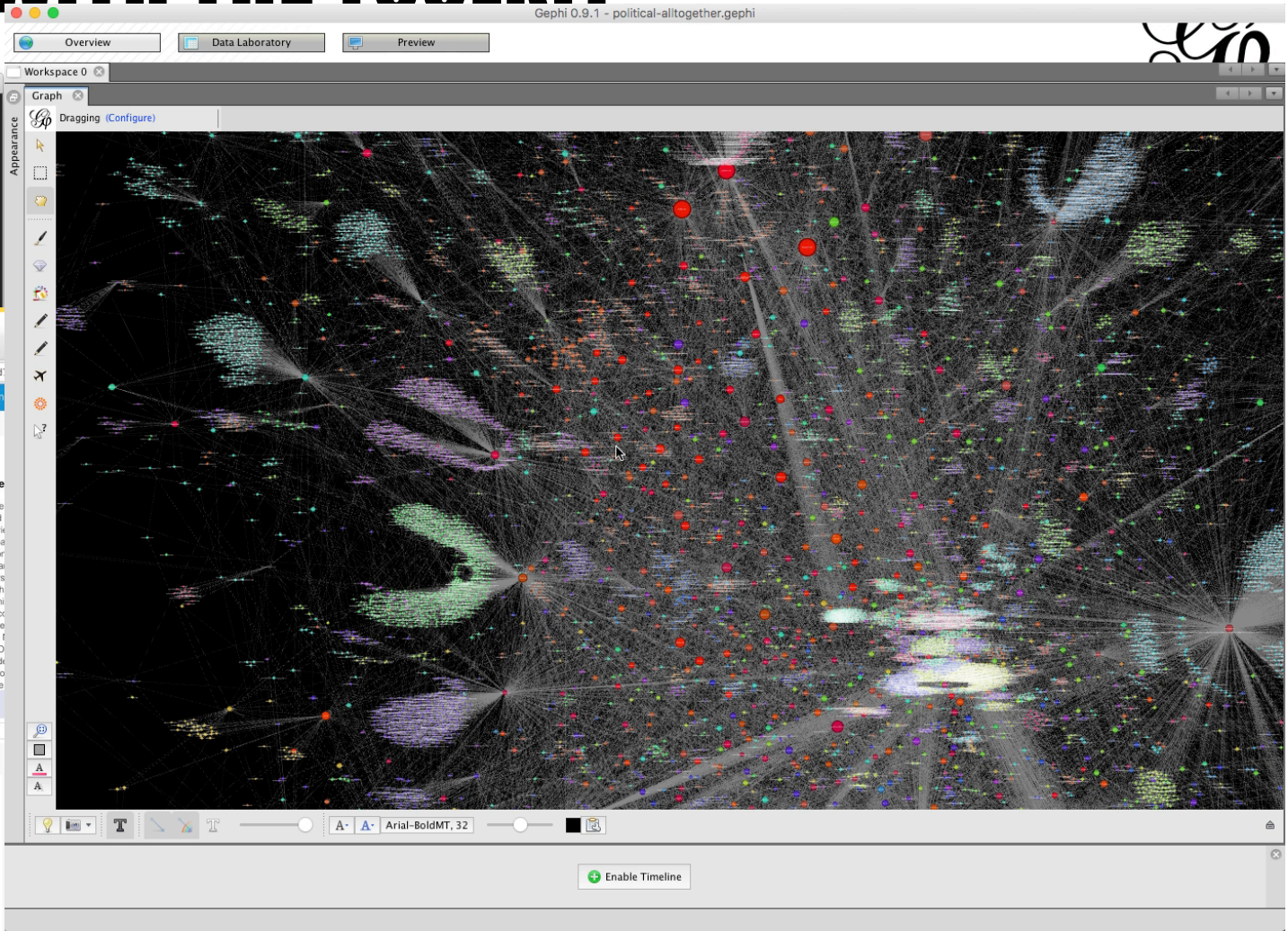
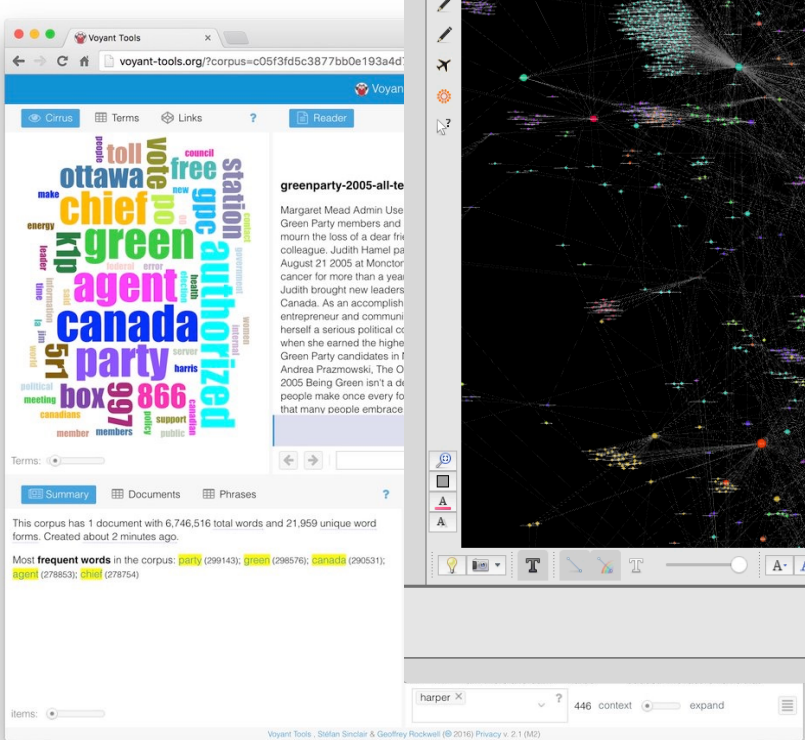
ssh 361

2. ssh

i2milligan@tuna:~\$

CAN DO SOME COOL STUFF WITH THE TOOLKIT

```
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
e-options: SAMEORIGIN Vary: Accept-Encoding,User-Agent Connection: close Content-Type: text/html; ch
arset=UTF-8 Communist Party of Canada - Parti Communiste du Canada Parti Communiste du Qu bec 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 Co
nstitution Party Program Our aim is socialism Capitalism in Canada Canada in a changing world The Ca
nadian state, nations and peoples of Canada, and the crisis of democracy The working class and peopl
e'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 C
ommunist Party We need to move Canada to the left! Nous devons recentrer le Canada vers la gauche!
Vote Communist! The Harper Conservative government's a
people. Meanwhile, the profits of the big banks and the
ple are now convinced that Harper and his gang must go.
uly progressive voices who will challenge right-wing, p
in a new, progressive direction, with a People's Alter
people and nature! Everywhere, the gap is widening betw
of the super-rich. Today, the top 1% own and control ov
acism and intolerance are also spreading. The crisis we
it is about capitalism. It is time capitalism was replac
ble system - with socialism. A new society, based on wor
owards true democracy - the rule of the people, by the p
nd economic wealth are owned and controlled by the work
gh collective struggle. Voting Communist sends the stro
a - and another world - is possible, urgent and worth fi
tter future! Read more. Votons communiste! Les politiques
fl au pour la population canadienne qui constate du m 
"part-00000" 851L, 4611113C
```





**WE MADE THE MISTAKE OF ASSUMING
THESE WERE OUR MAIN USERS...**

IN OTHER WORDS...



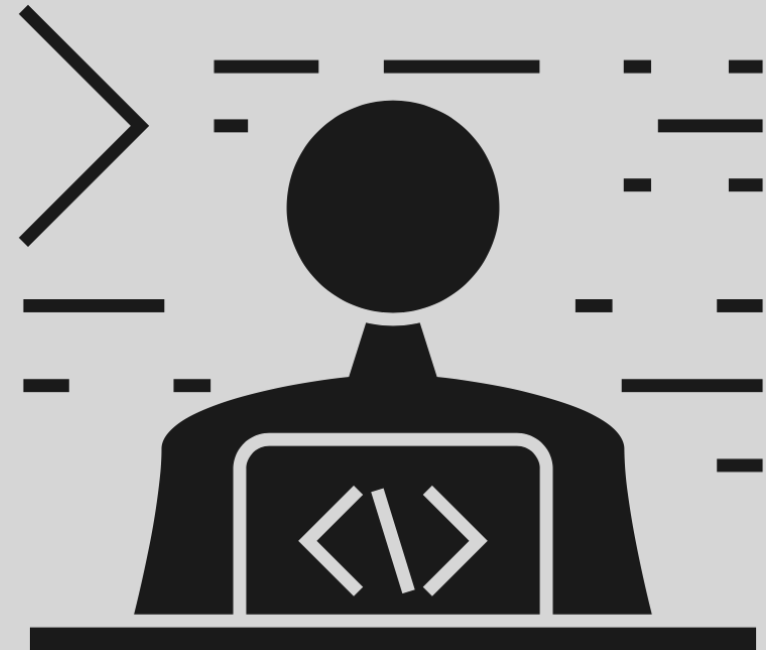
WE DIDN'T GO FAR ENOUGH... EXPECTED THEM TO COME TO US!

ENTER THE ANDREW W. MELLON FOUNDATION AND OUR GOAL OF MAKING AN ACCESSIBLE TOOL

PERSONA B

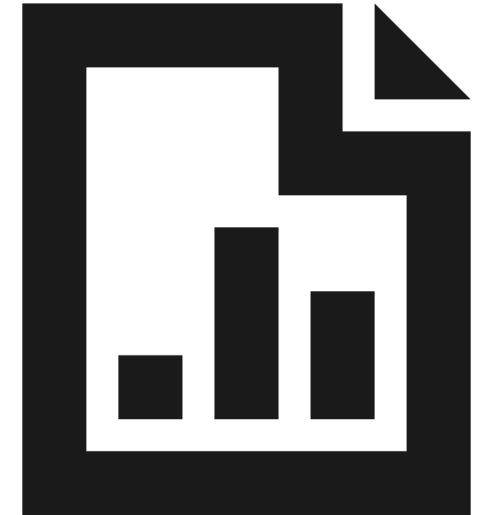
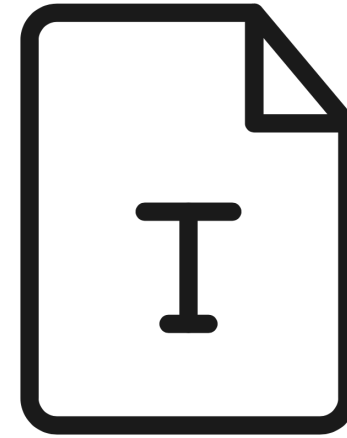
The sort of scholar we would have called a “digital humanist.” They’re comfortable with computers, can use some off-the-shelf tools like Voyant or Gephi to work with text/networks, and can use tutorials like those of the *Programming Historian* to learn some basic Python or R.

THE DIGITAL HISTORIAN



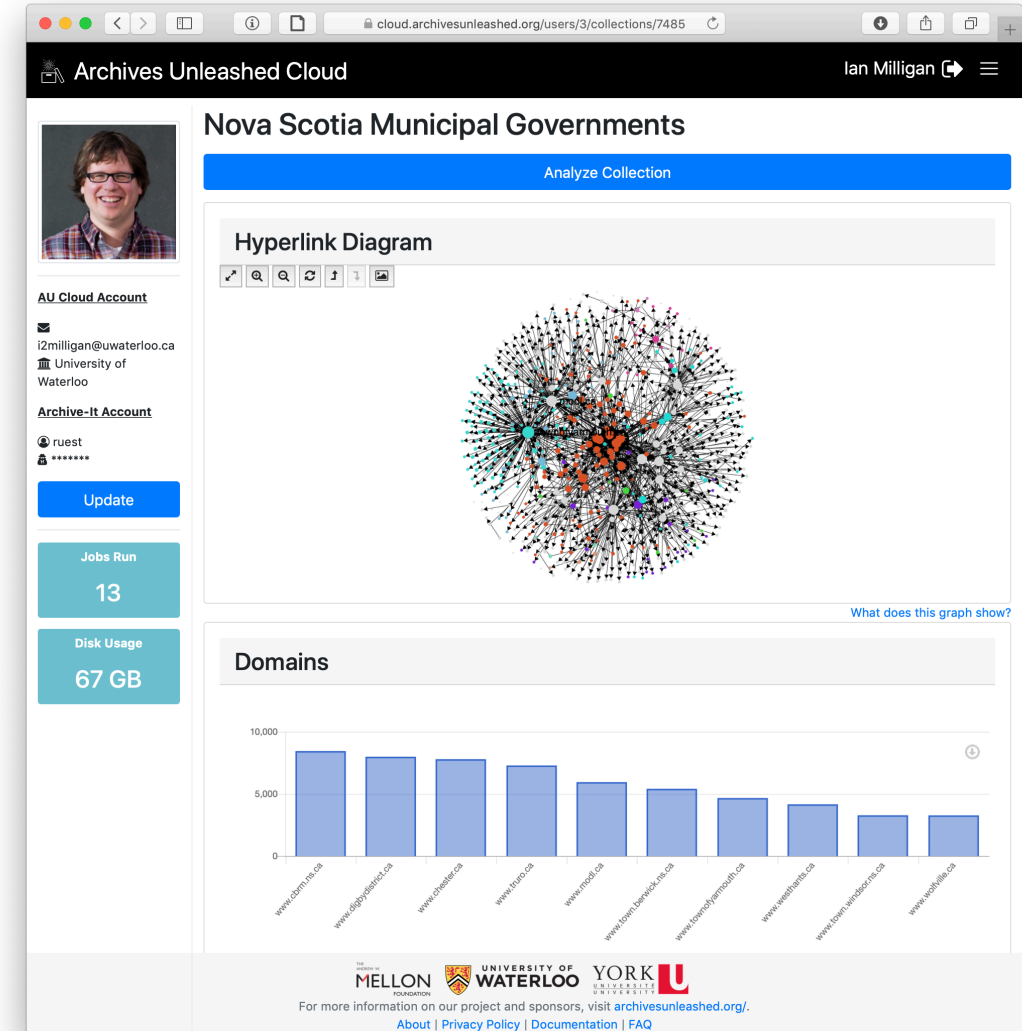
DIGITAL HUMANIST

- Needs to have some sort of toolkit to translate WARC files into something they can work with
- Can think critically about data
- Can use off-the-shelf tools to work with text/data, i.e. Voyant Tools or Gephi or other things based out of the *Programming Historian*
- In general doesn't want to use the command line or write custom programs



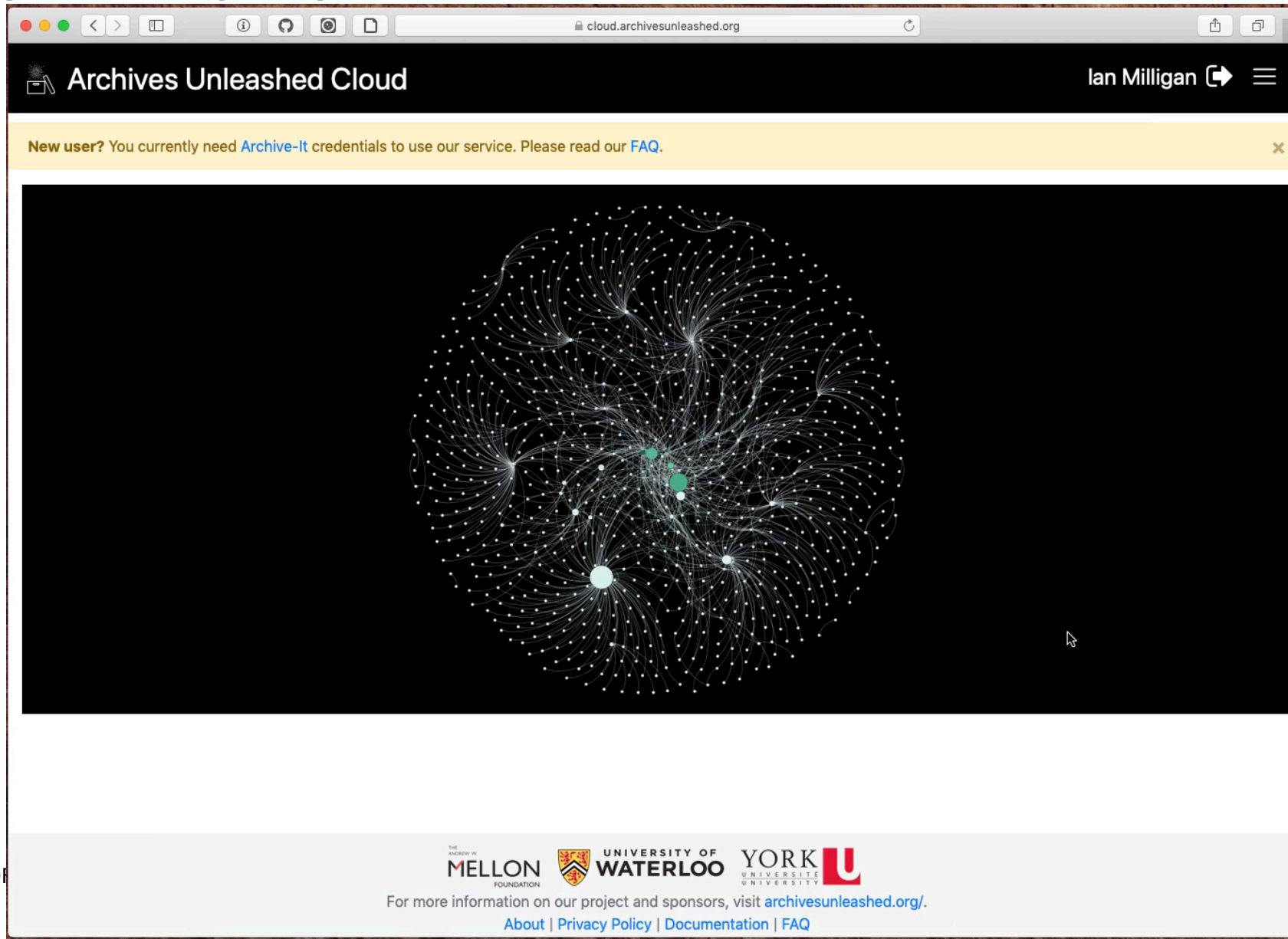
COMPUTATIONAL HUMANITIES

- The **Archives Unleashed Cloud** was designed around this persona.
- 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.



LET'S SEE THIS IN ACTION...

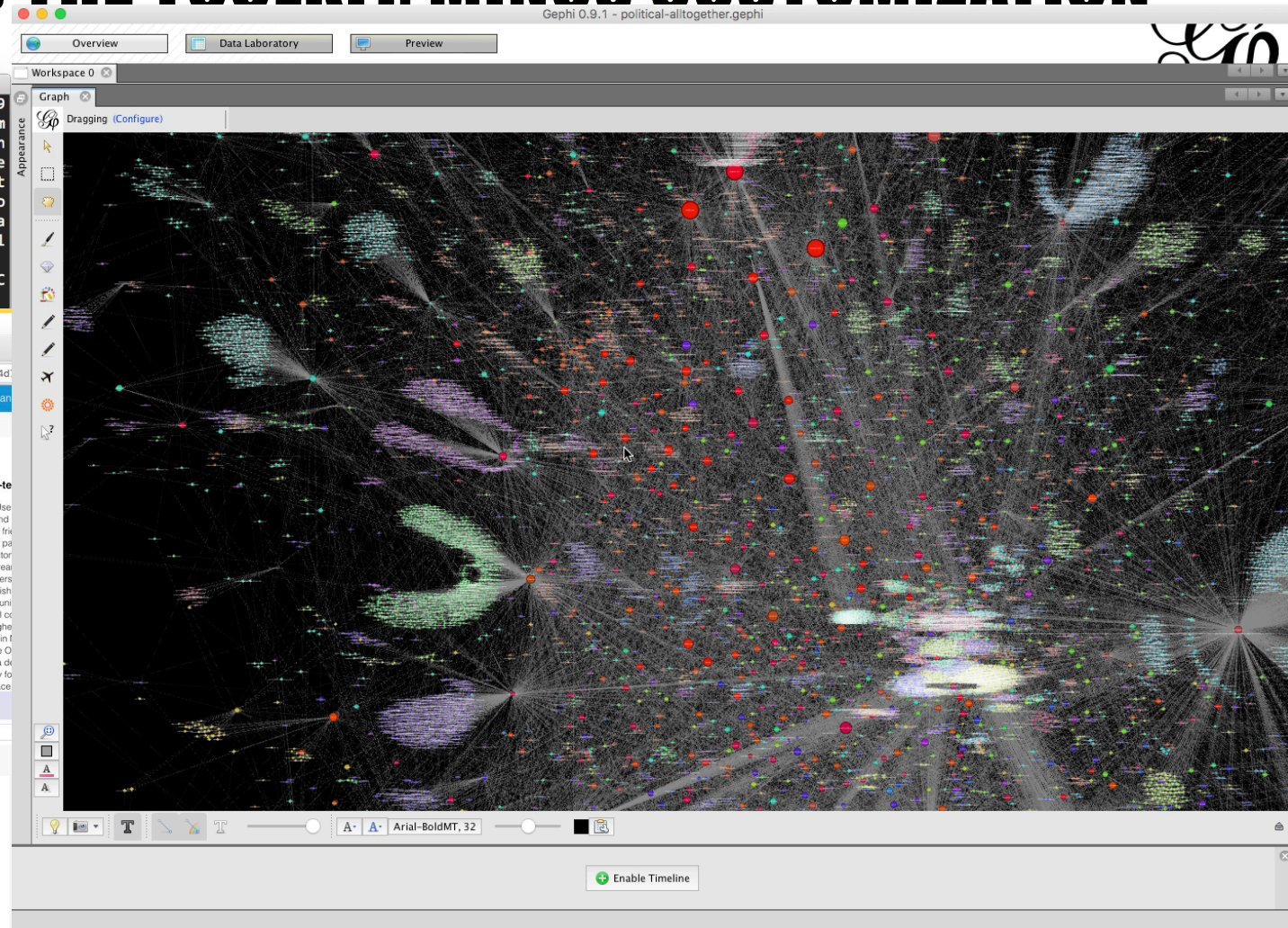
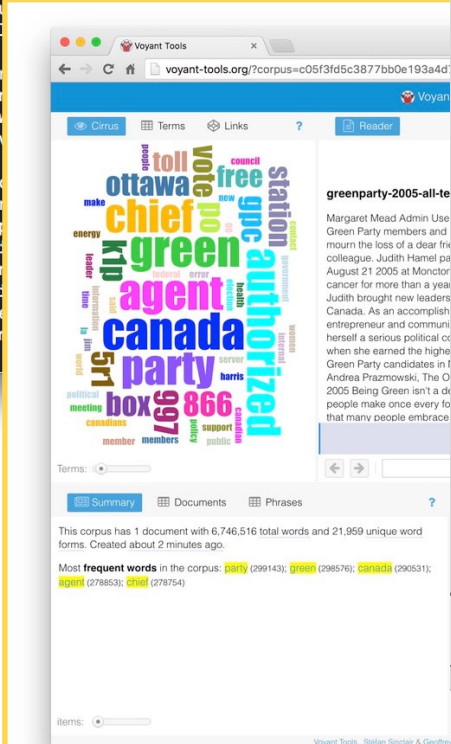
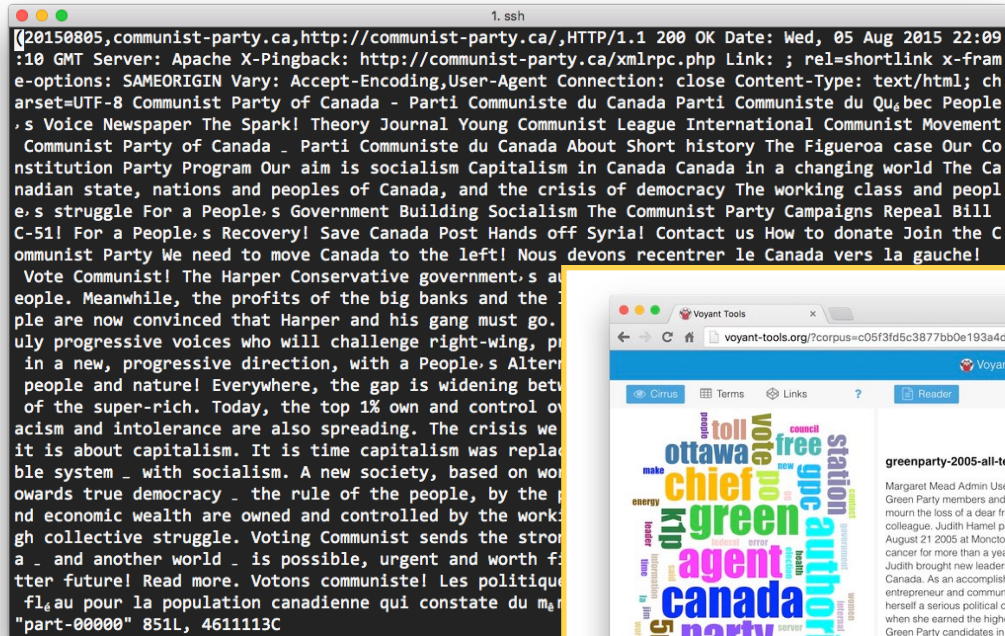
THE CLOUD IN ACTION



HISTORY IN THE AGE OF



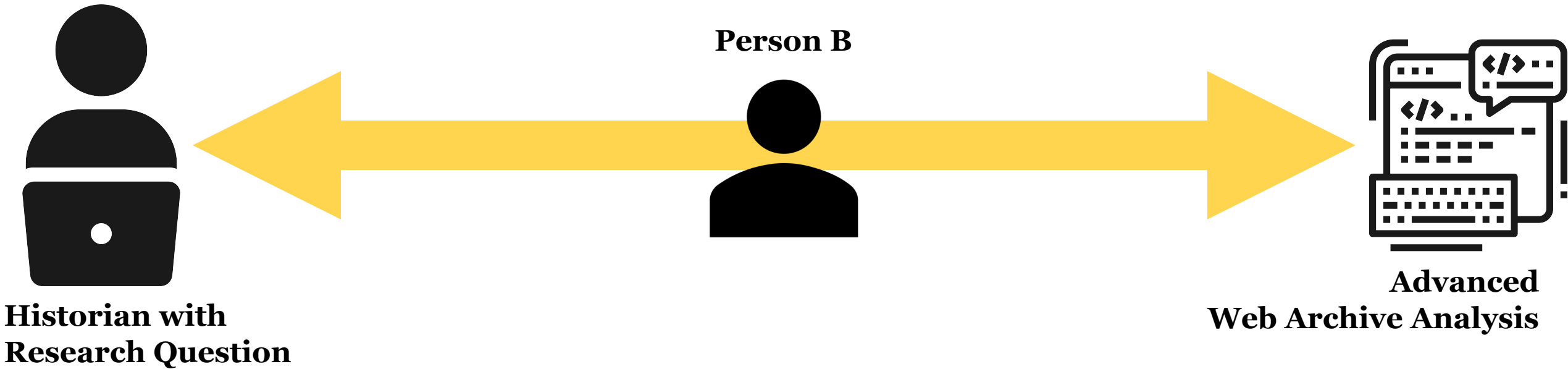
ENABLES THE SAME THINGS AS THE TOOLKIT. MINUS CUSTOMIZATION





**WE'RE MOVING TOWARDS OUR USERS...
BUT THEY STILL REQUIRE SOMEWHAT
SPECIALIZED SKILLS TO INTERPRET DATA.**

IN OTHER WORDS...

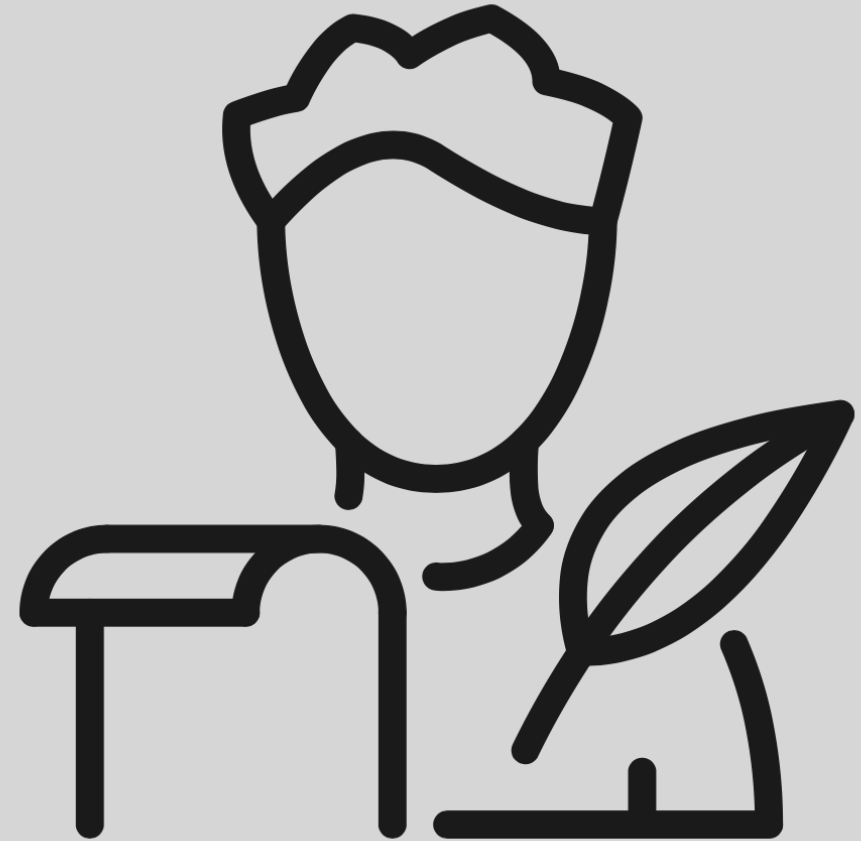


WE'RE GETTING THERE.. THEY COME PART WAY, WE COME PART WAY.

PERSONA C

The sort of scholar who might use computers – like for Word, some light Excel – but otherwise generally just wants to do historical research without earning new technical skills.

THE CONVENTIONAL HISTORIAN



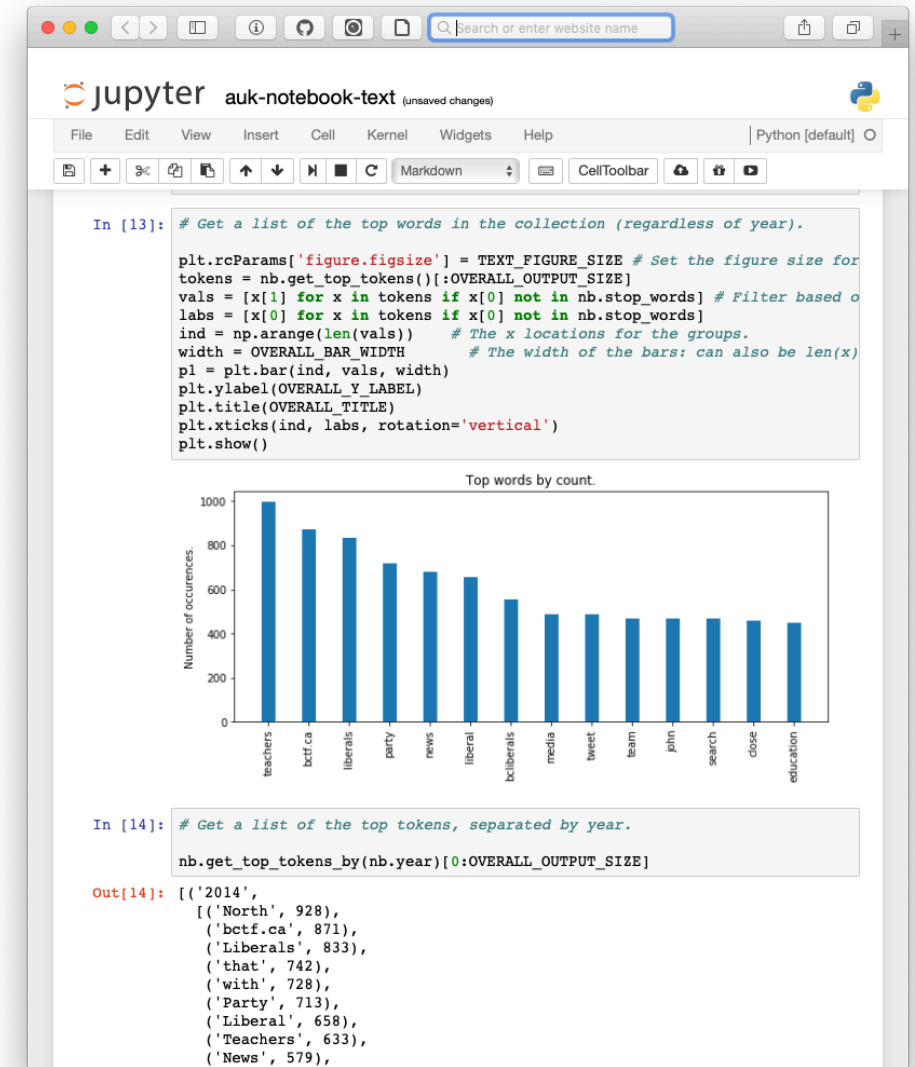
CONVENTIONAL HISTORIAN

- This is the toughest group to serve – if in the Toolkit model we weren't going out far enough to serve people in a reasonable area; this person might not be coming out far enough to meet us...
- But we will try.



CONVENTIONAL HISTORIAN

- The **Archives Unleashed Jupyter Notebooks** were designed with this in mind.
 - Uses web browser (which we can all use);
 - Slightly confusing layout of pressing play buttons and executing code, but thanks to Markdown integration we can mark it up.
 - Extensively playtested (and playtesting) with PhD students and MA students who have varying levels of computational comfort.
- **Exploring one-click hosting with Google Colab/Google Drive**

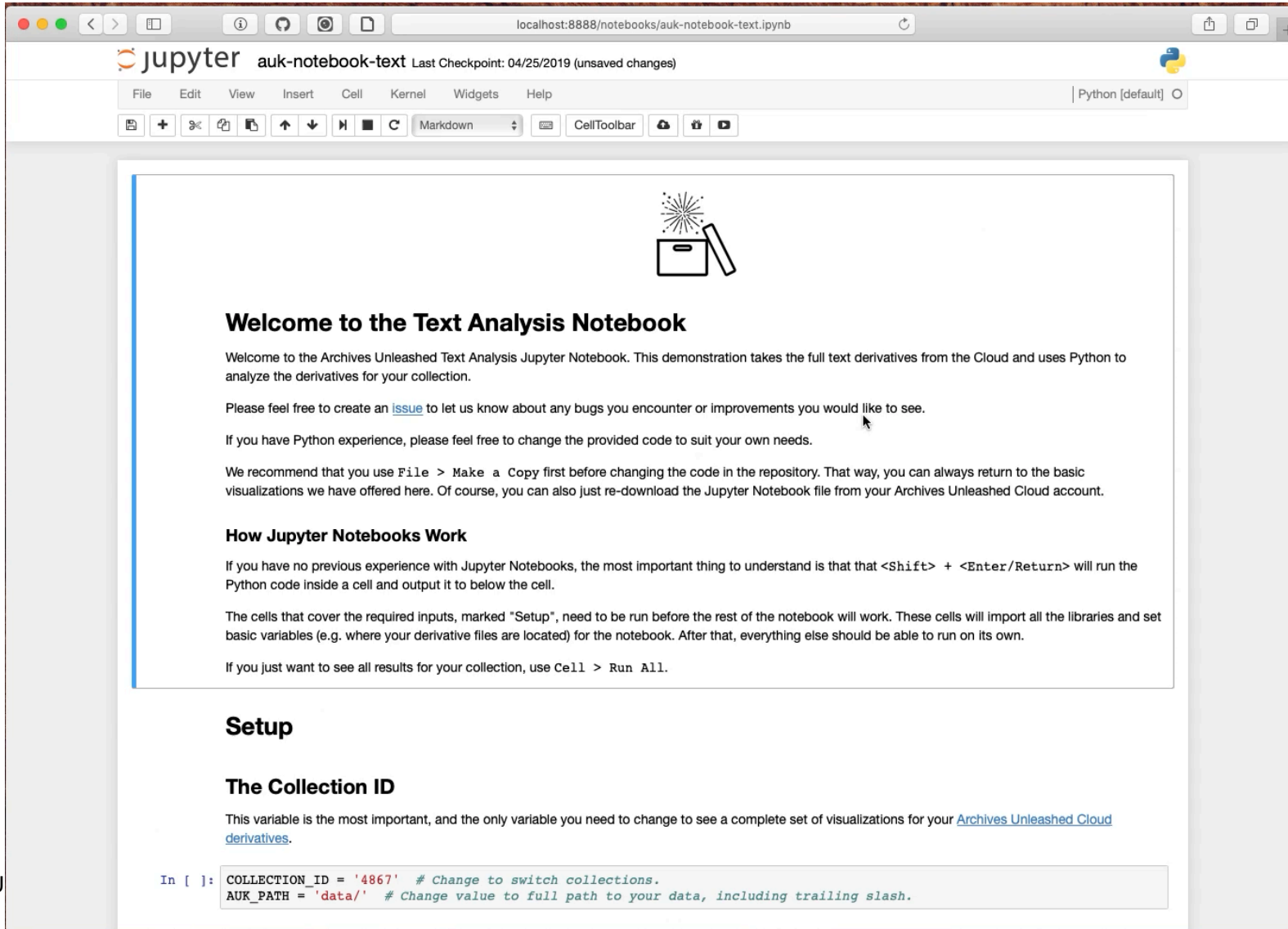


STILL NOT THE EASIEST TO USE..

But we're trying.

LET'S SEE THEM IN ACTION!

THE NOTEBOOKS IN ACTION




localhost:8888/notebooks/auk-notebook-text.ipynb

jupyter auk-notebook-text Last Checkpoint: 04/25/2019 (unsaved changes)

File Edit View Insert Cell Kernel Widgets Help Python [default]

Markdown CellToolbar



Welcome to the Text Analysis Notebook

Welcome to the Archives Unleashed Text Analysis Jupyter Notebook. This demonstration takes the full text derivatives from the Cloud and uses Python to analyze the derivatives for your collection.

Please feel free to create an [issue](#) to let us know about any bugs you encounter or improvements you would like to see.

If you have Python experience, please feel free to change the provided code to suit your own needs.

We recommend that you use File > Make a Copy first before changing the code in the repository. That way, you can always return to the basic visualizations we have offered here. Of course, you can also just re-download the Jupyter Notebook file from your Archives Unleashed Cloud account.

How Jupyter Notebooks Work

If you have no previous experience with Jupyter Notebooks, the most important thing to understand is that that <Shift> + <Enter/Return> will run the Python code inside a cell and output it to below the cell.

The cells that cover the required inputs, marked "Setup", need to be run before the rest of the notebook will work. These cells will import all the libraries and set basic variables (e.g. where your derivative files are located) for the notebook. After that, everything else should be able to run on its own.

If you just want to see all results for your collection, use Cell > Run All.

Setup

The Collection ID

This variable is the most important, and the only variable you need to change to see a complete set of visualizations for your [Archives Unleashed Cloud derivatives](#).

```
In [ ]: COLLECTION_ID = '4867' # Change to switch collections.
        AUK_PATH = 'data/' # Change value to full path to your data, including trailing slash.
```

IN OTHER WORDS...



WE'RE GETTING THERE... BUT THEY STILL NEED TO COME OUT A BIT..!

**FINALLY, WE BUILD CAPACITY TO HELP
PEOPLE COME OUT TO US THROUGH A
SERIES OF DATATHONS.**

DATATHONS



DATATHONS

- Helping to lower barriers
- Bringing people together
- Establishing a true community of practice to work with web archives



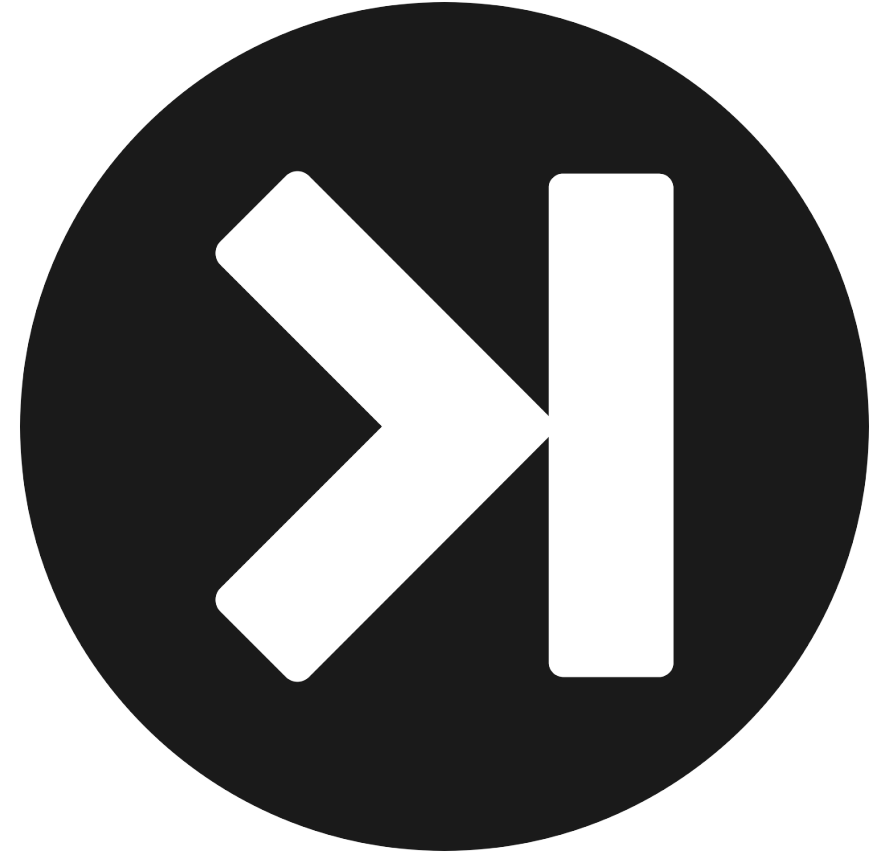
THE TL;DR

Historians in the future will need to understand the Web.

We need to make sure they're ready.

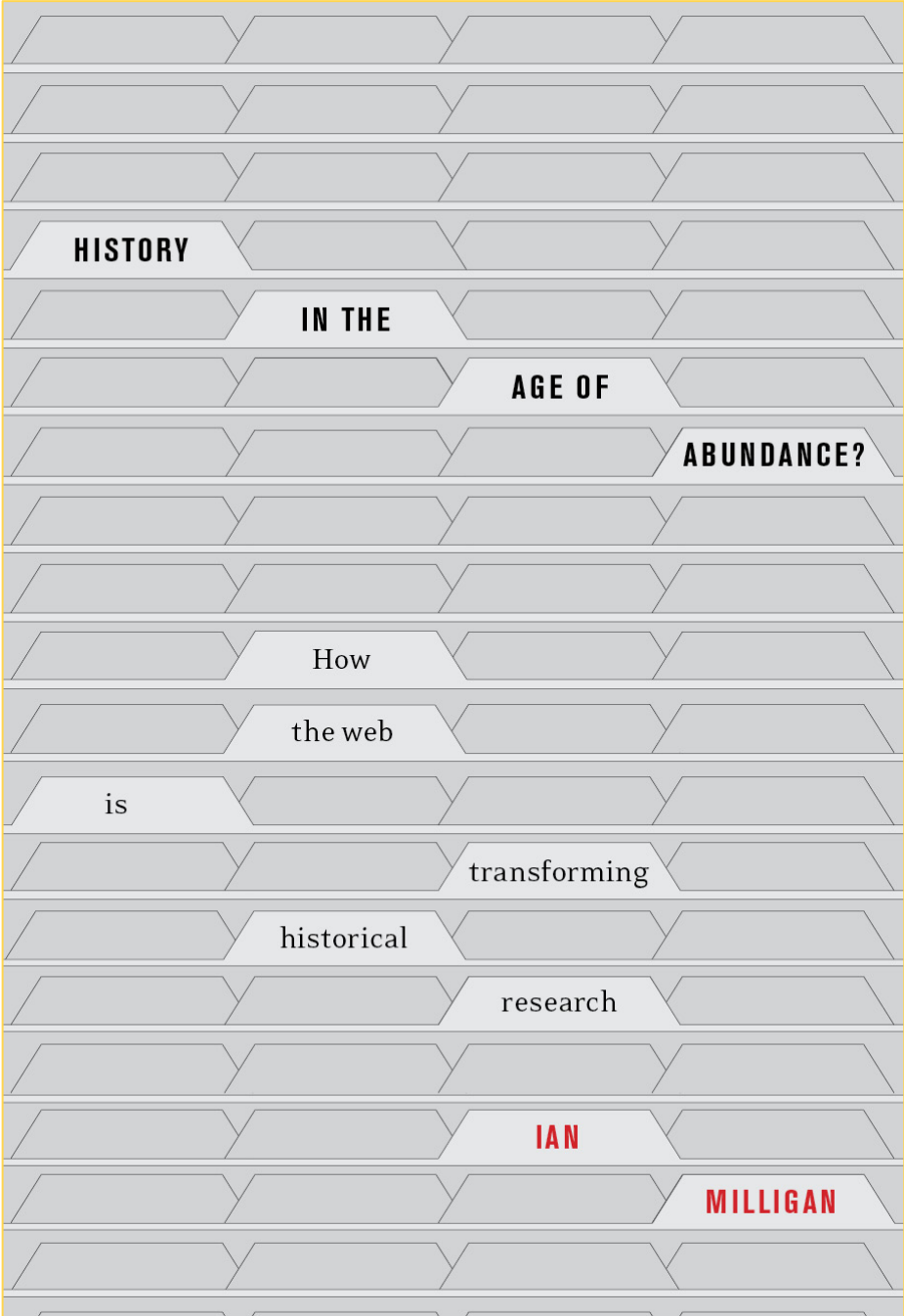
CONCLUDE

- Part of this is new, usable tools... (i.e. we come out towards them);
- Part of this is new cultures in the humanities and social sciences (i.e. they come out towards US)
- But overall, we all need to work together in light of the recognition that you can't study the 1990s without web archives – and we'll be studying the 1990s soon.



WANT TO LEARN MORE?

- You can ask me questions.
- Or you can funnel like a few royalty cents to me (or yeah, just pick up a copy at York University Libraries or Toronto Public Library, I can manage).
- *History in the Age of Abundance? How the Web is Transforming Historical Research.* McGill-Queen’s University Press, 2019.



THANKS TO OUR FUNDERS!



Social Sciences and Humanities
Research Council of Canada

Conseil de recherches en
sciences humaines du Canada

Canada