

# Supporting Computational Research on Large Digital Collections

Internet Archive & Archives Unleashed

Jefferson Bailey & Nick Ruest



CNI Fall 2022



# Presentation Outline

---

- General Challenges Supporting Computational Research on Large Digital Collections
- ARCH (Archives Research Compute Hub) Background, Goals, Status
- ARCH Technical Overview & Walkthrough
- Supporting Scholarly Use



# The Challenge

How to understand and address the technical, conceptual, and practical issues inherent in supporting data-driven uses of large, complex, heterogeneous (born) digital collections

# Technical Issues in Research Use

- Data format complexities (WARC, data, codecs)
- Data volume complexities (transfer, store)
- Data processing complexities (local v. cluster)
- Data “visibility” complexities (outputs)
- Data derivation complexities (pipelines)



# Conceptual Issues in Research Use

- Provenance complexities
- Acquisition complexities
- Border/boundary complexities
- Breadth complexities
- Ellipses (elisions)



# Practical Issues in Research Use

- Where is the front door? (requests)
- Scoping interviews (support)
- Research Agreements (paperwork)
- Time/cost modeling (budgeting)
- Responsible parties (staffing)
- Program Development (organization)



# Typology of Computational Research Services Models

- Bulk Data Model
- Cyberinfrastructure Model
- Roll Your Own Model
- Middleware Model
- Prepackaged Model
- Support & Community Model



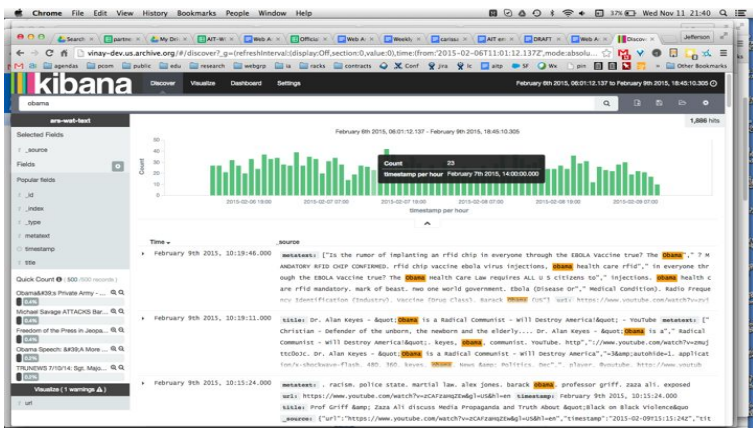
# Practical Lessons

- Large digital collections:
  - can be unkind to traditional methods of scholarly inquiry
    - Create guides or expertise to help users know the datasets, subsets, extractions to suite their research
  - impose many technical and conceptual issues
    - Be prepared to discuss/document these at the beginning of research request process
  - are alluring but can be deceitful
    - More data isn't better research; informed computational research services are critical
  - may be part of an even larger research corpus aggregation
    - Research support will require cross-functional skills





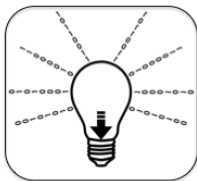
# Background: IA Computational Research Services



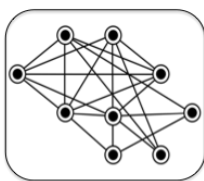
```

{
  "id": 5266,
  "account": 421,
  "created_by": "sbsreen",
  "created_date": "2015-02-03T00:51:21Z",
  "last_updated_by": "shallcross",
  "last_updated_date": "2015-03-02T16:38:43Z",
  "owner": "U.S. Presidential Election 2016",
  "tag": "",
  "state": "INACTIVE",
  "publicly_visible": true,
  "one_hop_off": false,
  "topics": "Government-USFederal;politicsAndElections;government-Nation",
  "tool_exported": false,
  "metadata": {
    "Contributor": [
      {
        "id": 926936,
        "value": "Shallcross, Michael"
      }
    ],
    "Description": [
      {
        "id": 926937,
        "value": "The 2016 U.S. Presidential Election web archive"
      }
    ],
    "Title": [
      {
        "id": 926939,
        "value": "U.S. Presidential Election 2016"
      }
    ],
    "Creator": [
      {
        "id": 926938,
        "value": "Green, Sarah; Nofziger, Cinda; and Thomas, Rob"
      }
    ],
    "Date": [
      {
        "id": 926935,
        "value": "2015-02-03"
      }
    ],
    "Type": [
      {
        "id": 926934,
        "value": "Web Archive"
      }
    ]
  }
}

```



**WAT Datasets**  
(Web Archive Transformation)  
Key Metadata from Every Resource



**LGA Datasets**  
(Longitudinal Graph Analysis)  
What Links to What over Time



**WANE Datasets**  
(Web Archive Named Entities)  
Names of People, Places, Organizations

```

In [6]: # Index by date
mime_types = ['text', 'image', 'audio', 'video', 'application', 'revisit', 'other']
days = pd.Series([0] * len(mime_types))

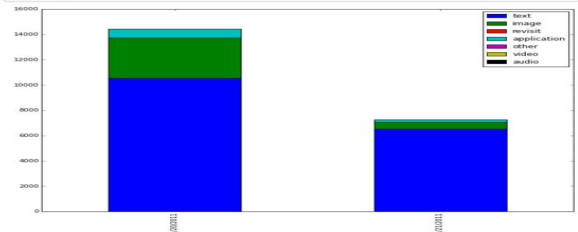
# Initialization
day_mime_dict = dict()
for day in days:
    for mime in mime_types:
        day_mime_dict[day][mime] = 0

for day, mime in day_mime_metadata.iterrows():
    day_mime_dict[day][mime] = day_mime_dict[day][mime] + 1
print day_mime_dict

[12/25/2011: {'text': 6542, 'image': 528, 'revisit': 0, 'application': 382, 'other': 0, 'video': 0, 'audio': 0}, 12/26/2011: {'text': 10955, 'image': 3168, 'revisit': 0, 'application': 706, 'other': 9, 'video': 1, 'audio': 0}]

In [7]: import pandas as pd
day_mime_df = pd.DataFrame.from_dict(day_mime_dict, orient='index')
columns_index = day_mime_df.columns
day_mime_df.plot(kind='bar', stacked=True)

```



# Background: The “Archives Unleashed” Project

- A long-running project (est. 2017) that “*aims to make petabytes of historical internet content accessible to scholars and others interested in researching the recent past,*” primarily supported by The Andrew W. Mellon Foundation.
- In 2020, Archives Unleashed partnered with the Internet Archive to develop “ARCH”.
- Three main partners: Internet Archive, University of Waterloo (hi Ian!), and York University.

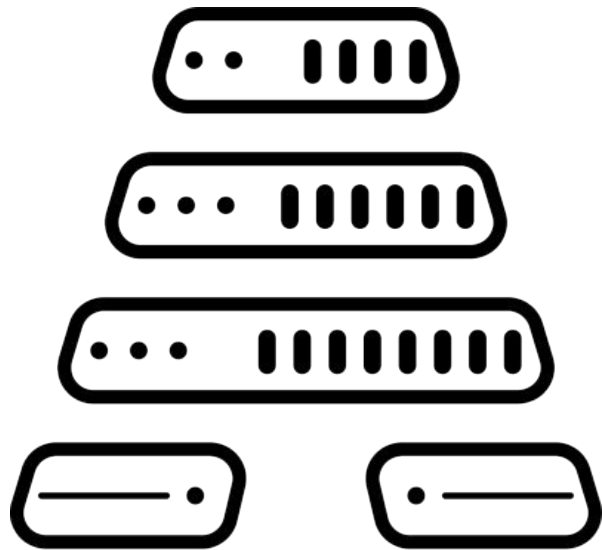
**Archives Unleashed** 



# ARCH Goals, Development, Roadmap

- Goal: merge and standardize our tools and services into one web/text/data mining platform in a OSS SaaS product model
- Goal: colocate data & compute in centralized IA-run infrastructure and data centers (eventually in US & Canada)
- Goal: embed researcher support into the platform dev process with funded cohorts teams and invited scholars
- Development: Engineering by AU & IA teams 2022-current
- Development: Supporting AU cohorts & IA pilot partners; 20+ LAMs & 50+ scholars, 16 countries; 200+TB processed
- Roadmap: Add text/image to collection types; query filtering; aggregated/uploaded collections; small/mid LAM outreach
- Roadmap: Production release Q1 2023



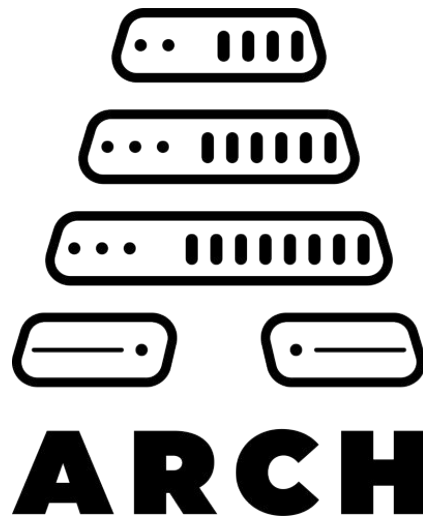


**ARCH**

Archives Research Compute Hub

# ARCH Platform Details

- Interactive web application interface both for collection curators and scholarly researchers
- Generate and download over 20 derivative datasets, and connect to Google Colab
- Three rounds of UI/UX testing
- In-browser visualizations and data previews that presents a glimpse into collection content
- Located in the Internet Archive data center, ARCH has quick access to the petabytes of content collected

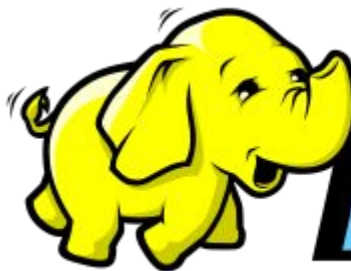
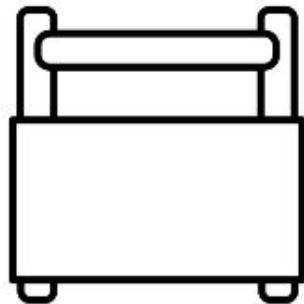




## ARCH Stack



Scalatra



**hadoop**

APACHE  
**Spark**<sup>TM</sup>

# Archives Research Compute Hub

Scala version 2.12.8 Scala version 2.5.4 license MIT

## About

Web application for distributed compute analysis of Archive-It web archive collections.

## Building

☰ README.md

- `sbt "prod/clean" "prod/assembly" "prod/assemblyPackageDependency"`

## Docker

1. Create a config ( `config/config.json` ) for your Docker setup, e.g., by copying the included template: `cp config/docker.json config/config.json`
2. `docker build --no-cache -t arch .`
3. `docker run -it --rm -p 12341:12341 -p 54040:54040 -v /home/nruest/Projects/au/sample-data/arch:/data -v /home/nruest/Projects/au/arch:/app arch`

Web application will be available at: <http://localhost:12341/ait>, and Apache Spark interface will be available at <http://localhost:54040>.

## Citing ARCH

How to cite ARCH in your research:

Helge Holzmann, Nick Ruest, Jefferson Bailey, Alex Dempsey, Samantha Fritz, Peggy Lee, and Ian Milligan. 2022. ABCDEF: the 6 key features behind scalable, multi-tenant web archive processing with ARCH: archive, big data, concurrent, distributed, efficient, flexible. In Proceedings of the 22nd ACM/IEEE Joint Conference on Digital Libraries (JCDL '22). Association for Computing Machinery, New York, NY, USA, Article 13, 1–11. <https://doi.org/10.1145/3529372.3530916>

Your citations help to further the recognition of using open-source tools for scientific inquiry, assists in growing the web archiving community, and acknowledges the efforts of contributors to this project.

## License

MIT

## Open-source, not open-contribution

Similar to [SQLite](#), ARCH is open source but closed to contributions.

The level of complexity of this project means that even simple changes can break a lot of other moving parts in our production environment. However, community involvement, bug reports and feature requests are warmly accepted.

## Acknowledgments

This work is primarily supported by the [Andrew W. Mellon Foundation](#). Other financial and in-kind support comes from the [Social Sciences and Humanities Research Council](#), [Compute Canada](#), [York University Libraries](#), [Start Smart Labs](#), and the [Faculty of Arts](#) at the [University of Waterloo](#).

Any opinions, findings, and conclusions or recommendations expressed are those of the researchers and do not necessarily reflect the views of the sponsors.

# ABCDE F - The 6 key features behind scalable, multi-tenant web archive processing with ARCH: Archive, Big Data, Concurrent, Distributed, Efficient, Flexible

Helge Holzmann<sup>1</sup>, Nick Ruest<sup>2</sup>, Jefferson Bailey<sup>1</sup>, Alex Dempsey<sup>1</sup>, Samantha Fritz<sup>2</sup>, Peggy Lee<sup>1</sup>, and Ian Milligan<sup>3</sup>

<sup>1</sup> Internet Archive

<sup>2</sup> Digital Scholarship Infrastructure Department, York University

<sup>3</sup> Department of History, University of Waterloo

## ABSTRACT

Over the past quarter-century, web archive collection has emerged as a user-friendly process thanks to cloud-hosted solutions such as the Internet Archive's Archive-It subscription service. Despite advancements in collecting web archive content, no equivalent has been found by way of a user-friendly cloud-hosted analysis system. Web archive processing and research require significant hardware resources and cumbersome tools that interdisciplinary researchers find difficult to work with. In this paper, we identify six principles - the ABCDEFs (Archive, Big data, Concurrent, Distributed, Efficient, and Flexible) - used to guide the development and design of a system. These make the transformation of, and working with, web archive data as enjoyable as the collection process. We make these objectives - largely common sense - explicit and transparent in this paper. They can be employed by every computing platform in the area of digital libraries and archives and adapted by teams seeking to implement similar infrastructures. Furthermore, we present ARCH (Archives Research Compute Hub)<sup>1</sup>, the first cloud-based system designed from scratch to meet all of these six key principles. ARCH is an interactive interface, closely connected with Archive-It, engineered to provide analytical actions, specifically generating datasets and in-browser visualizations. It efficiently streamlines research workflows while eliminating the burden of computing requirements. Building off past work by both the Internet Archive (Archive-It Research Services) and the Archives Unleashed Project (the Archives Unleashed Cloud), this merged platform achieves a scalable processing pipeline for web archive research. It is open-source and can be considered a reference implementation of the ABCDEF, which we have evaluated and discussed in terms of feasibility and compliance as a benchmark for similar platforms.

<sup>1</sup><https://github.com/internetarchive/arch>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).  
JCDL '22, June 20-24, 2022, Cologne, Germany  
© 2022 Copyright held by the owner/authors. Publication rights licensed to ACM.  
ACM ISBN 978-1-4503-9345-4/22.06...\$15.00  
<https://doi.org/10.1145/3529372.3530916>

## CCS CONCEPTS

• Information systems → Digital libraries and archives; Data extraction and integration.

## KEYWORDS

web archives, big data, data processing, distributed computing

## ACM Reference Format

Helge Holzmann<sup>1</sup>, Nick Ruest<sup>2</sup>, Jefferson Bailey<sup>1</sup>, Alex Dempsey<sup>1</sup>, Samantha Fritz<sup>2</sup>, Peggy Lee<sup>1</sup>, and Ian Milligan<sup>3</sup>. 2022. ABCDEF - The 6 key features behind scalable, multi-tenant web archive processing with ARCH: Archive, Big Data, Concurrent, Distributed, Efficient, Flexible. In *The ACM/IEEE Joint Conference on Digital Libraries in 2022 (JCDL '22)*, June 20-24, 2022, Cologne, Germany. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3529372.3530916>

## 1 INTRODUCTION

Web archiving is an important component of modern digital libraries. It is essential for enabling future research into contemporary history and ensuring the long-term preservation of our documentary heritage [11] [9]. Yet while collecting web archive content has matured into a user-friendly process, thanks in no small part to cloud-hosted solutions such as the Internet Archive's Archive-It service, this ease-of-use has not been matched on the analysis side. We accordingly need a user-friendly system that can enable the creation of research datasets from web archives so that researchers can work with material at scale.

In this paper, we present the Archives Research Compute Hub (ARCH), a production system tightly integrated with the Internet Archive infrastructure and services. ARCH grew out of the Archives Unleashed Cloud, a proof-of-concept platform that demonstrated the ability of a web-browser-based system to power backend Apache Spark-driven jobs on web archival datasets [12]. Powered by the Archives Unleashed Toolkit and the Internet Archive's Sparkling data processing library, the ARCH platform will become a complementary component of the Internet Archive's Archive-It system. ARCH is built around six key principles: archive, big data, concurrent, distributed, efficient, and flexible. We present these principles as considerations for projects and teams developing similar systems.

## 2 RELATED WORK AND PROJECT CONTEXT

Established in 2017, the Archives Unleashed project recognizes the collective need among researchers, librarians and archivists for analytical tools, community infrastructure, and accessible web archival

I am a BORING banner...

from adsenger.com





# Archive Research Compute Hub



Learn More: [ARCH Documentation](#)

## Collections

Collection Name	Public Collection	Recently Created	Last Created	Size
<a href="#">AU Team</a>	Yes	<a href="#">Extract plain text of webpages (Sample)</a>	2022-11-08 18:27:59	2.2 GB
<a href="#">Canada U15</a>	Yes	<a href="#">Extract image information</a>	2022-04-28 18:21:24	318.3 GB
<a href="#">Ontario COVID-19</a>	Yes	-	-	34.6 GB

## U.S. LGBTQ Web Collection Analysis

Learn More: [ARCH Documentation](#)

### Job Summary

[Generate Datasets](#)

### Collection Overview

257 seeds



Crawled Jan 23, 2019



1.1 TB



Public Collection



Public Collection Link: <https://archive-it.org/collections/02778>

There are currently no active jobs, [generate a new dataset](#).

### Completed Jobs

Job	Finished
<a href="#">Domain frequency</a>	2022-10-13 08:33:03
<a href="#">Extract domain graph</a>	2022-10-13 08:52:04
<a href="#">Extract domain graph (Example)</a>	2022-09-23 00:16:40
<a href="#">Extract image graph (Example)</a>	2022-09-23 00:04:17
<a href="#">Extract plain text of webpages</a>	2022-10-13 13:58:02
<a href="#">Extract plain text of webpages (Example)</a>	2022-09-21 14:03:54
<a href="#">Extract web graph (Example)</a>	2022-09-23 00:03:49



## U.S. LGBTQ Web Collection Jobs

Learn More: [ARCH Documentation](#)

Job Summary

**Generate Datasets**



### Collection #

These jobs provide a basic overview of the collection.

#### Domain frequency

Create a CSV with the following columns: domain and count.

Generate Example Dataset

View Dataset

#### Generate web archive transformation (WAT) files

Creates Web Archive Transformation (WAT) files that include a brief header which identifies its corresponding URL via "WARC-Target-URI," corresponding W/ARC file via "WARC-Refers-To," and additional mapping information.

Generate Example Dataset

Generate Dataset



## Network #

These jobs produce files that provide network graphs for analysis, and offer an opportunity to explore the way websites link to each other.

### Extract domain graph

Create a CSV with the following columns: crawl date, source domain, target domain, and count.

[View Example Dataset](#)[View Dataset](#)

### Extract image graph

Create a CSV with the following columns: crawl date, source of the image (where it was hosted), the URL of the image, and the alternative text of the image.

[View Example Dataset](#)[Generate Dataset](#)

### Extract longitudinal graph

Creates Longitudinal Graph Analysis (LGA) files which contain a complete list of what URLs link to what URLs, along with a timestamp.

[Generate Example Dataset](#)[Generate Dataset](#)

### Extract web graph

Create a CSV with the following columns: crawl date, source, target, and anchor text. Note that this contains all links and is not aggregated into domains.

[View Example Dataset](#)[Generate Dataset](#)



## Text #

These jobs produce files that allow the user to explore text components of a web archive, including extracted "plain text", named entities, HTML, css, and other web elements.

### Extract named entities

Creates Web Archive Named Entities (WANE) files which contain the named entities from each text resource, organized by originating URL and timestamp.

[Generate Example Dataset](#)[Generate Dataset](#)

### Extract plain text of webpages

Create a CSV with the following columns: crawl date, last modified date, web domain, URL, MIME type as provided by the web server, MIME type as detected by Apache TIKA, and content (HTTP headers and HTML removed).

[View Example Dataset](#)[View Dataset](#)

### Extract text files (html, text, css, js, json, xml) information

Create a CSV with the following columns: crawl date, last modified date, URL of the text file, filename, text extension, MIME type as provided by the web server, MIME type as detected by Apache TIKA, text file MD5 hash and text file SHA1 hash, and text file content.

[Generate Example Dataset](#)[Generate Dataset](#)



## File Formats #

These jobs produce files that contain information on certain types of binary files found within a web archive.

### Extract audio information

Create a CSV with the following columns: crawl date, last modified date, URL of the audio file, filename, audio extension, MIME type as provided by the web server, MIME type as detected by Apache Tika, audio MD5 hash and audio SHA1 hash.

Generate Example Dataset

Generate Dataset

### Extract image information

Create a CSV with the following columns: crawl date, last modified date, URL of the image, filename, image extension, MIME type as provided by the web server, MIME type as detected by Apache Tika, image width, image height, image MD5 hash and image SHA1 hash.

Generate Example Dataset

Generate Dataset

### Extract PDF information

Create a CSV with the following columns: crawl date, last modified date, URL of the PDF file, filename, PDF extension, MIME type as provided by the web server, MIME type as detected by Apache Tika, PDF MD5 hash and PDF SHA1 hash.

Generate Example Dataset

Generate Dataset

### Extract PowerPoint (e.g., .ppt, .odp, .key) information

Create a CSV with the following columns: crawl date, last modified date, URL of a PowerPoint or similar file, filename, PowerPoint or similar file extension, MIME type as provided by the web server, MIME type as detected by Apache Tika, PowerPoint or similar file MD5 hash and PowerPoint or similar file SHA1 hash.

Generate Example Dataset

Generate Dataset

### Extract spreadsheet information

Create a CSV with the following columns: crawl date, last modified date, URL of the spreadsheet file, filename, spreadsheet extension, MIME type as provided by the web server, MIME type as detected by Apache Tika, spreadsheet MD5 hash and spreadsheet SHA1 hash.

Generate Example Dataset

Generate Dataset

### Extract video information

Create a CSV with the following columns: crawl date, last modified date, URL of the video file, filename, video extension, MIME type as provided by the web server, MIME type as detected by Apache Tika, video MD5 hash and video SHA1 hash.

Generate Example Dataset

Generate Dataset

### Extract Word Documents (e.g., .doc, .odt, .rtf, .wpd) information

Create a CSV with the following columns: crawl date, last modified date, URL of the word document or similar file, filename, word document or similar file extension, MIME type as provided by the web server, MIME type as detected by Apache Tika, word document or similar file MD5 hash and word document or similar file SHA1 hash.

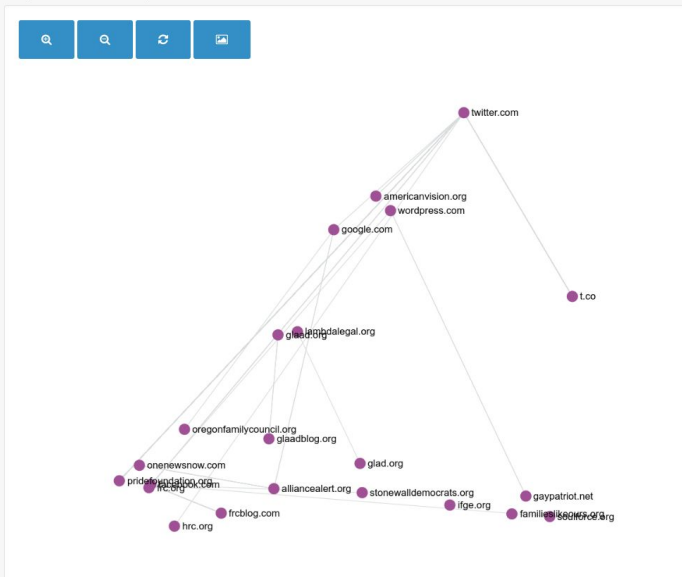
Generate Example Dataset

Generate Dataset

## U.S. LGBTQ Web Collection: Extract domain graph

[Learn More: ARCH Documentation](#)

### Top Hosts Sample



This graph shows up to the top 100 nodes with highest degree (number of inlinks + outlinks) along with their edges.

## Dataset(s)

A CSV with the following columns: crawl date, source domain, target domain, and count.

**File name:** domain-graph.csv.gz

**File size:** 42.6 MB

**Result count:** 6,327,594 lines

**Date completed:** 2022-10-13

**Checksum(s):** md5:d542459b517b1fc90c1ff01c69a132b4



Download

## Preview

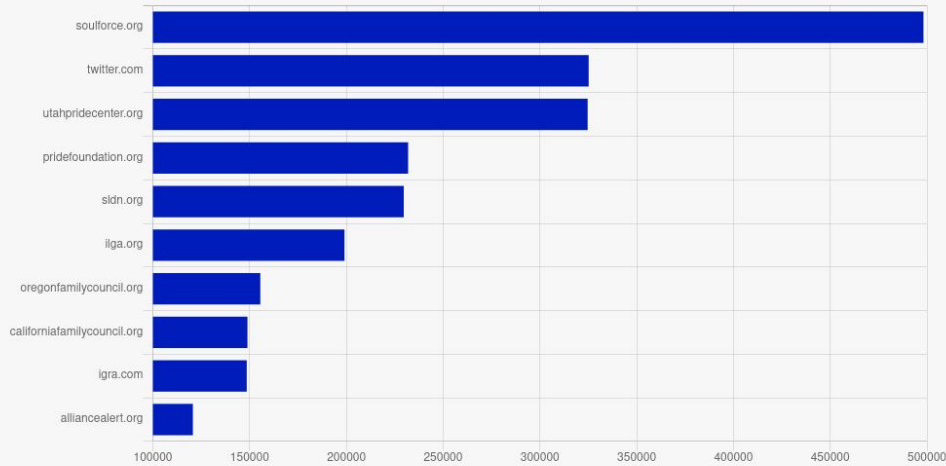
crawl_date	source	target	count
20180327152758	porco.ga	porco.ga	12081
20180327153659	susi.ml	susi.ml	12081
20180327152801	porco.gq	porco.gq	12081
20130408205620	facebook.com	facebook.com	11851
20130408203658	facebook.com	facebook.com	11254
20160616014151	tubepornstars.com	tubepornstars.com	9290
20170609194110	annakooliman.com	annakooliman.com	8961
20170609194142	squarespace.com	squarespace.com	8904
20130409013244	tubepornstars.com	tubepornstars.com	8806
20130408194055	facebook.com	facebook.com	8118
20130408194040	facebook.com	facebook.com	8063
20160629165957	facebook.com	facebook.com	7758
20130408190950	facebook.com	facebook.com	7720
20130408194787	facebook.com	facebook.com	7676

[Download Preview Data](#)

## U.S. LGBTQ Web Collection: Domain frequency

[Learn More: ARCH Documentation](#)

### Domains



This graph shows the top 10 domains and the number of times that they occur in this collection.

### Dataset(s)

A CSV with the following columns: domain and count.

**File name:** domain-frequency.csv.gz

**File size:** 170.4 KB

**Result count:** 23,910 lines

**Date completed:** 2022-10-13

**Checksum(s):** md5:e618bb22efe9f7ff2a85d28798d5b45f



[Download](#)

### Preview

domain	count
soulforce.org	498191
twitter.com	325181
utahpridecenter.org	324651
pridefoundation.org	231888
sidn.org	229624
ilga.org	198982
oregonfamilycouncil.org	155503
californiamfamilycouncil.org	148862
lgra.com	148490
alliancealert.org	120648
gaygames.net	104389
narth.com	99862
gaypatriot.net	93066
...	...

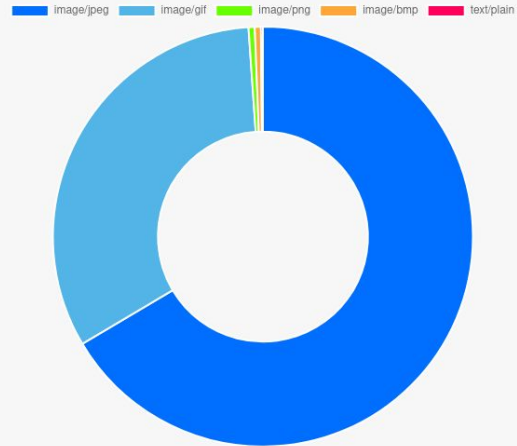
[Download Preview Data](#)



## Canada U15: Extract image information

Learn More: [ARCH Documentation](#)

### File Format Distribution



This graph shows the distribution of the various file formats in the web archive collection.

### Dataset(s)

A CSV with the following columns: crawl date, last modified date, URL of the image, filename, image extension, MIME type as provided by the web server, MIME type as detected by Apache TIKA, image width, image height, image MD5 hash and image SHA1 hash.

**File name:** image-information.csv.gz  
**File size:** 1.6 MB  
**Result count:** 24,705 lines  
**Date completed:** 2022-04-28  
**Checksum(s):** md5:191494a2a761f203f1ec4123e751b334



Download

 Open in Colab



# My Research Page

The screenshot shows a Jupyter Notebook interface with a dark theme. The notebook title is "image-information.ipynb". The menu bar includes "File", "Edit", "View", "Insert", "Runtime", "Tools", and "Help". The toolbar shows options for "+ Code", "+ Text", and "Copy to Drive".

## Image Information Dataset Exploration

We're going to take a look at a few examples of how we can explore the Image Information dataset.

The first thing we need to do is enter the URL for our Image Information dataset in the cell below. You can get this by right clicking the Download icon, and selecting "Copy Link".

```
dataset: "https://webdata.archive-it.org/ait/files/download/ARCHIVEIT-14462/ImageInformationExtraction/image-information.csv.gz?access=FIRYZOPIJQCZFRACASNMRYA7SJC55XSRR"
```

[Show code](#)

```
https://webdata.archive-it.org/ait/files/download/ARCHIVEIT-14462/ImageInformationExtraction/image-information.csv.gz?access=FIRYZOPIJQCZFRACASNMRYA7SJC55XSRR
```

## pandas

Next, we'll setup our environment so we can load our Image Information dataset into [pandas](#) DataFrames. If you're unfamiliar with DataFrames, but you've worked with spreadsheets before, you should feel comfortable pretty quick.

```
[ ] import pandas as pd
```

## Data Table Display

Colab includes an extension that renders pandas DataFrames into interactive displays that can be filtered, sorted, and explored dynamically. This can be very useful for taking a look at what each DataFrame provides, and doing some initial filtering!

Data table display for pandas dataframes can be enabled by running:

```
%load_ext google.colab.data_table
```

and disabled by running

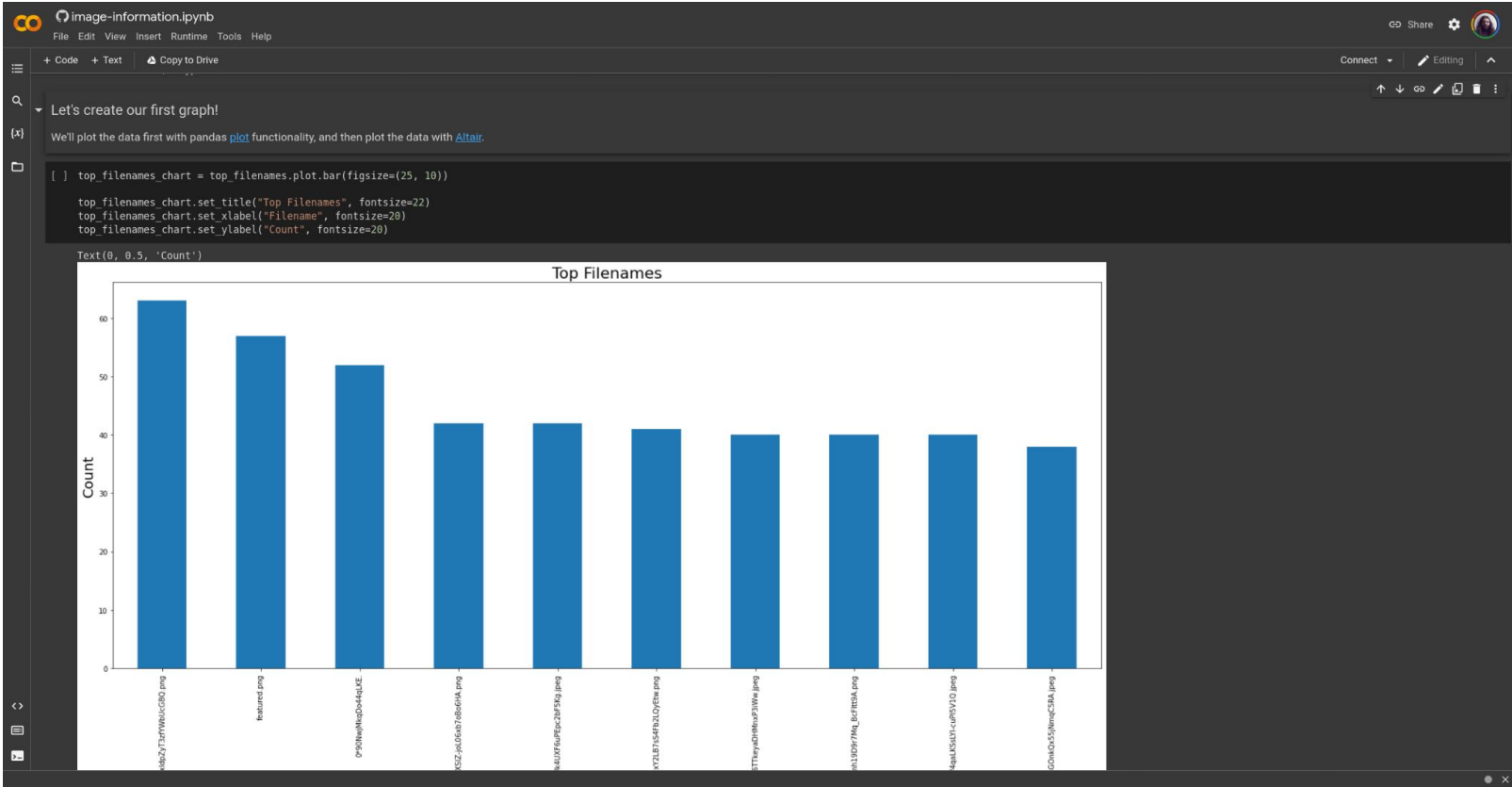
```
%unload_ext google.colab.data_table
```

```
[ ] %load_ext google.colab.data_table
```

## Loading our ARCH Dataset as a DataFrame



# BRAIN CLUB



## Creating order from the mess: web archive derivative datasets and notebooks

Nick Ruest <sup>a</sup>, Samantha Fritz <sup>b</sup> and Ian Milligan <sup>b</sup>

<sup>a</sup>York University, Toronto, Canada; <sup>b</sup>University of Waterloo, Waterloo, Canada

### ABSTRACT

For a quarter-century, memory institutions have been preserving web-based content. These web archives have been collected and stored in ARC and WARC (W/ARC) file formats and will form a basis for contemporary histories. Yet, these formats present significant challenges to researchers who wish to access and use web archival data. This is primarily due to the nature of collecting, storing, and providing access to these multifaceted digital objects. In other words, web archives are messy. Applying traditional archival methods of description to digital-born collections is complicated due to issues of provenance, original order, and scale. However, we believe that archival description offers a practical starting point for thinking about access. This paper argues a robust finding aid must extend beyond basic collection-level description to allow for more meaningful interactions with web archives. As such, we propose a reimagining of a traditional finding-aid model into a three-level mode of description to include computational methods, the generation of derivative datasets, and interactive code-rich notebooks. These three factors combine to ultimately contribute to the expanded access and use of web archives.

### KEYWORDS

Web archives; big data; notebooks; finding aids; data science

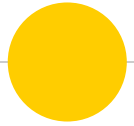
### Introduction

Since 1996, organizations such as the Internet Archive have been collecting web archives with an eye to making digital heritage accessible. A quarter-century later, attention is turning to the problem of analysis: what can we do to make the petabytes of information — over 100PB in the case of the Internet Archive — usable.

Web archives are messy. Collected using a variety of platforms, web resources are aggregated in standardized file formats: the ARC file, followed by the successor WARC format (proposed as an ISO standard in 2009, and a standard since 2017). Working with these files is difficult for users. Resources are co-mingled, meaning the text of a webpage is interspersed with multi-faceted digital objects and data such as HTML, CSS, JavaScript, binary encoded images, videos, documents, or legacy file formats. Not only do WARCs contain a diverse range of digital objects, but the order also ultimately eludes human readership. Since crawlers often work in parallel, files from dozens or hundreds of websites are woven together. While web archives are traditionally accessed through a Wayback Machine replay instance, this only allows one-page-at-a-time browsing,

**CONTACT** Nick Ruest  [ruestn@yorku.ca](mailto:ruestn@yorku.ca)

© 2022 The Author(s). Published by Informa UK Limited, trading as Taylor & Francis Group  
This is an Open Access article distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivatives License (<http://creativecommons.org/licenses/by-nc-nd/4.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is properly cited, and is not altered, transformed, or built upon in any way.



# Supporting Research

Archives Unleashed Cohorts

## AU Cohorts

- Purpose:
  - Facilitate research engagement with computational use of large web/digital archives through use of the ARCH platform in funded scholarly research projects
- Scope:
  - 10 total research teams (5 each year)
- Research support:
  - Funding (~1\$0K USD), technical support, bi-weekly calls
  - 2 in person events





## AU Cohort Projects

---

- ◉ AWAC2 Analysing Web Archives of the COVID Crisis through the IIPC Novel Coronavirus dataset
- ◉ Everything Old is New Again: A Comparative Analysis of Feminist Media Tactics between the 2nd- to 4th Waves
- ◉ Mapping and tracking the development of online commenting systems on news websites between 1996–2021
- ◉ Crisis Communication in the Niagara Region during the COVID-19 Pandemic
- ◉ Viral health misinformation from Geocities to COVID-19





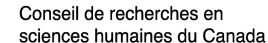
## AU Cohort Projects

---

- ◉ Latin American Women's Rights Movements: Tracing Online Presence through Language, Time and Space
- ◉ Historicizing Aughts-Era Mormon Mommy Blogging Media Landscapes
- ◉ Web Archiving and the Saskatchewan COVID Archive: Expanding Coverage to Capture Social Media, Medical Misinformation, and Radicalization
- ◉ Querying Queer Web Archives
- ◉ Using Web Archives for Mapping the Use of Cultural Practices in Postconflict Societies and During Reconciliation Processes

# Acknowledgements of Institutional Support

ARCH is a joint project between [Internet Archive](#) and [Archives Unleashed](#). This work is primarily supported by the [Andrew W. Mellon Foundation](#). Other financial and in-kind support has come from the [Social Sciences and Humanities Research Council](#), [Compute Canada](#), [York University Libraries](#), [Start Smart Labs](#), and the [Faculty of Arts](#) and [David R. Cheriton School of Computer Science](#) at the [University of Waterloo](#).





# Thanks!

*Any questions ?*

Connect with our project teams:

- [arch@archive.org](mailto:arch@archive.org)
- [archivesunleashed@gmail.com](mailto:archivesunleashed@gmail.com)
- [jefferson@archive.org](mailto:jefferson@archive.org)
- [nick@archivesunleashed.org](mailto:nick@archivesunleashed.org)