

# All Our Yesterdays: A toolkit to explore web archives in Colab

Tim Ribaric

Sam Langdon



# Motivation?

- WARC files are difficult to work with directly (especially if they are gigabytes in size)
- No standard desktop tool is suitable to explore them
- Often you don't want to look at the whole archive anyway



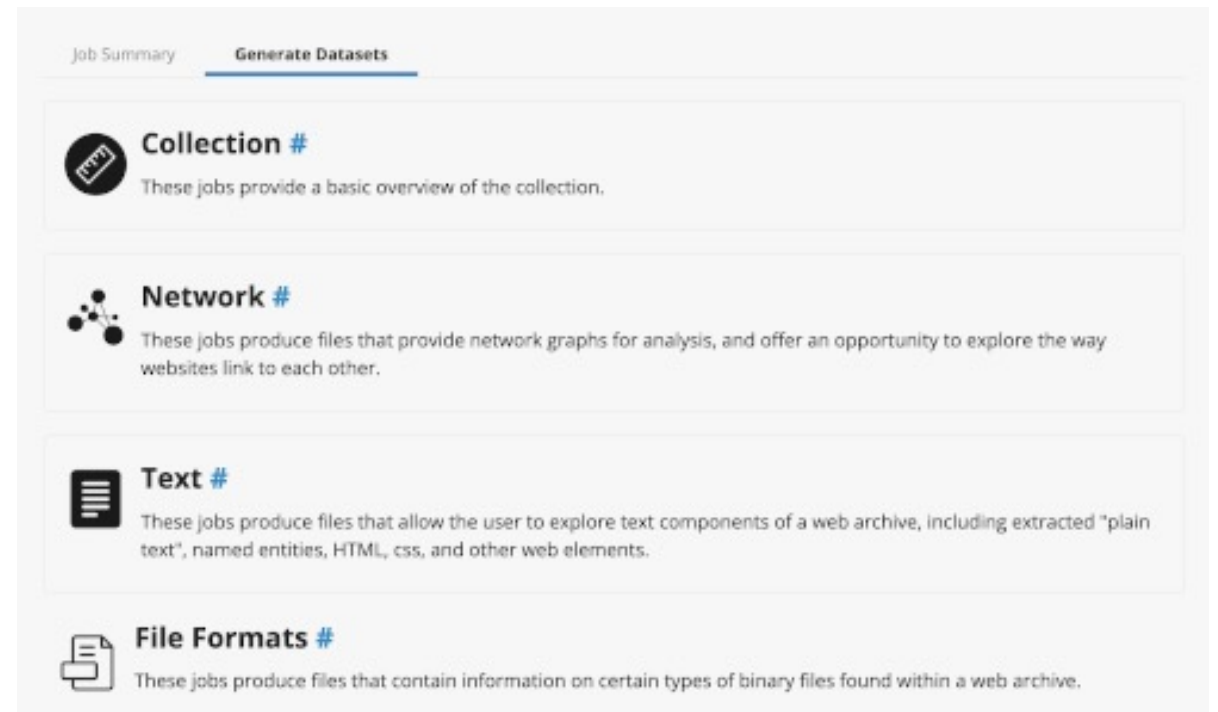
# Motivation?

- Looking to make a toolkit that scaffolds up an environment that will create the derivatives and plot out basic analysis steps
- Build on the new additions to the ARCH tool that is being developed as part of ARCHIVE-IT



# Motivation?

- One of the things Archives Unleashed has been working on
- An integrated derivative generator
- Used these derivatives to perform some in-depth analysis



# Motivation













- What do you do with your derivative? Particularly if you don't have much background / experience working with CSV Files
- Google Colab is a great place to start

## ARCH\_Data\_Explore

Notebooks and datasets for the Archives Unleashed Cohort Grant for the **Covid-19 In Niagara** proje

More details on project site: [https://brockdsl.github.io/archives\\_unleashed/](https://brockdsl.github.io/archives_unleashed/)

### Notebook listing

ARCH Data Exploration	<a href="#">Notebook</a>	 <a href="#">Open in Colab</a>
COMM 4P35 Tutorial	<a href="#">Notebook</a>	 <a href="#">Open in Colab</a>
Hackfest notebook	<a href="#">Notebook</a>	 <a href="#">Open in Colab</a>
Muni Data Export	<a href="#">Notebook</a>	 <a href="#">Open in Colab</a>
Prep Domain Data	<a href="#">Notebook</a>	 <a href="#">Open in Colab</a>
Twitter Data Export	<a href="#">Notebook</a>	 <a href="#">Open in Colab</a>
Municipal Data Similarity	<a href="#">Notebook</a>	 <a href="#">Open in Colab</a>
Another example of Municipal Data Similarity using SpaCy	<a href="#">Notebook</a>	 <a href="#">Open in Colab</a>
Municipal Data Similarity using TF-IDF	<a href="#">Notebook</a>	 <a href="#">Open in Colab</a>
Content size of pages over time	<a href="#">Notebook</a>	 <a href="#">Open in Colab</a>
Frequency of page updates over time	<a href="#">Notebook</a>	 <a href="#">Open in Colab</a>
Word frequency visualized	<a href="#">Notebook</a>	 <a href="#">Open in Colab</a>

# Demo

Derivative Generation

# Demo

Analysis of a derivative

# Roadmap

- Archives Unleashed Grant for 'Niagara COVID Archive'
- Prototype of notebooks created for analysis
- Match of Mind Grant to develop AOYTK software
- ...
- Further Match of Mind Grant to do User Testing on user notebooks



# More info

- Project is available on GitHub:
  - <https://brockdsl.github.io/AOYTK/>
- Looking for collaborators to test? [dsl@brocku.ca](mailto:dsl@brocku.ca)

