# What's in a Web Archive Collection? Summarization and Discovery of Archived Webpages

**Dr. Michele C. Weigle**

Professor, Department of Computer Science
Web Science and Digital Libraries (WS-DL) Group
Old Dominion University

Slides: https://bit.ly/weigle-nlm-2022

# Featuring work done in collaboration with

## Web Science and Digital Libraries (WS-DL) Research Group at ODU

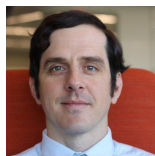Dr. Michael L. Nelson

### PhD Students and Alumni
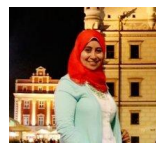
Shawn Jones
(PhD 2021)

Alexander Nwala
(PhD 2020)
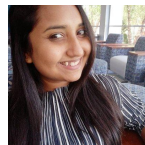
Sawood Alam
(PhD 2020)

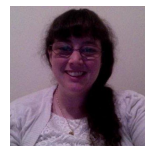Mat Kelly
(PhD 2019)

Yasmin AlNoamany
(PhD 2016)

Kritika Garg

Himarsha Jayanetti

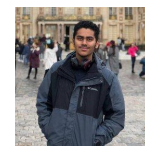### MS Students and Alumni

Lesley Frew

Abigail Mabe
(MS 2021)

Dhruv Patel
(MS 2021)

Maheedhar Gunnam
(MS 2018)

Surbhi Shankar
(MS 2017)

### Funders

NATIONAL ENDOWMENT FOR THE Humanities

NSF

INSTITUTE of Museum and Library SERVICES

IIPC netpreserve.org

# The Web holds our stories

ESPN.com          March 13, 2020



CNN.com          March 14, 2020

# But webpages can change or disappear

ESPN.com          November 1, 2022

CNN.com          November 1, 2022

# Maybe it's archived?



[CNN.com archived on March 13, 2020](#)

More than just saving a screenshot, web archives strive to capture and allow replay of the *entire contents of a web page*, including its source HTML and embedded images, stylesheets, and JavaScript source code.



Michele C. Weigle, "On the Importance of Web Archiving", *SSRC Parameters*, Sep 2018, [https://items.ssrc.org/parameters/on-the-importance-of-web-archiving/](https://items.ssrc.org/parameters/on-the-importance-of-web-archiving/)

# Web archives are essential for studying recent history and culture



Ian Milligan (University of Waterloo), https://twitter.com/NetPreserve/status/1141321443373920256

# There are several public web archives

# How do webpages get archived?          Anatomy of a Web Page



CSS

JavaScript

https://www.cs.odu.edu/~mweigle/
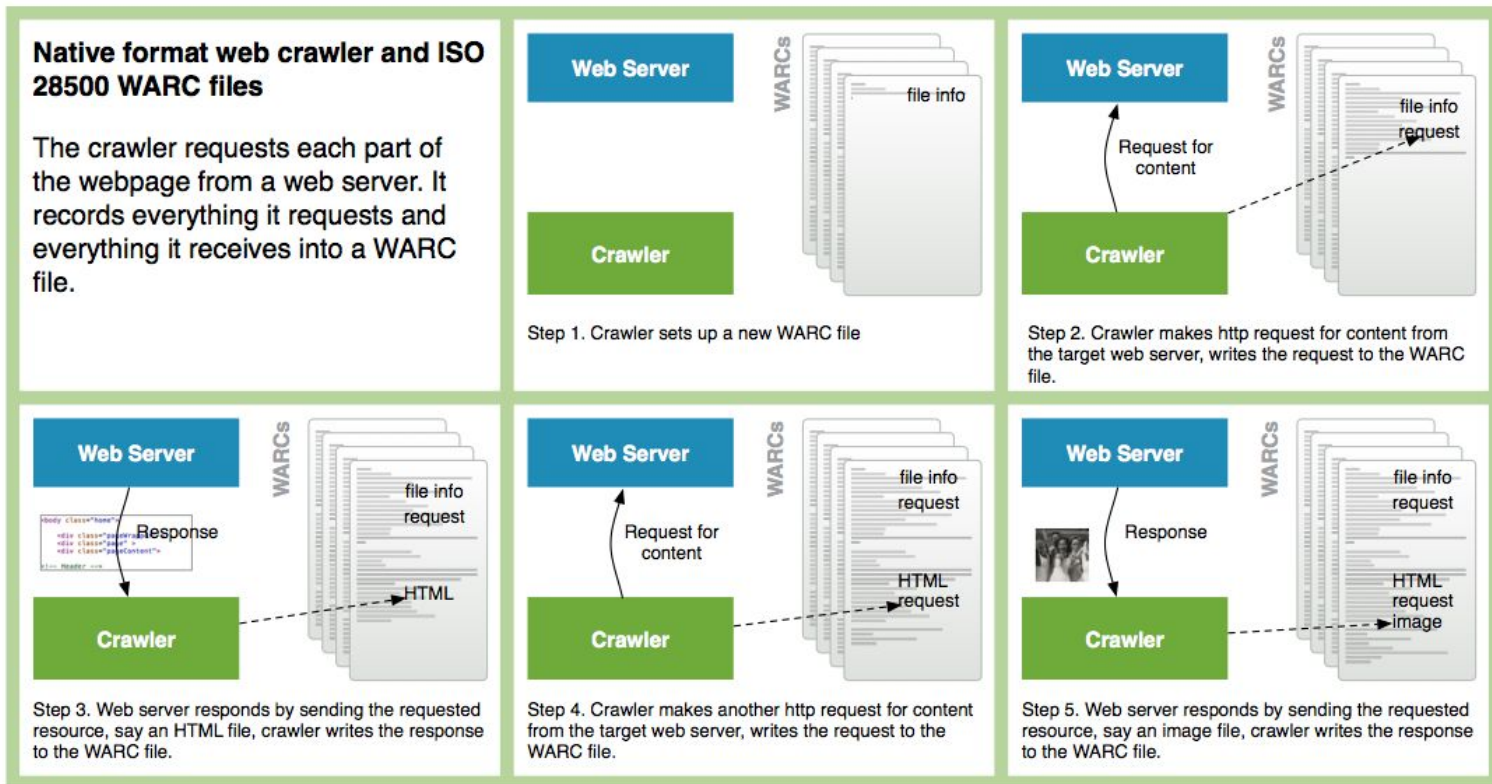
# How do webpages get archived?     Web Crawling Process

# How do webpages get archived?    WARC: Web ARChive file format



**Native format web crawler and ISO 28500 WARC files**

The crawler requests each part of the webpage from a web server. It records everything it requests and everything it receives into a WARC file.

Step 1. Crawler sets up a new WARC file

Step 2. Crawler makes http request for content from the target web server, writes the request to the WARC file.

Step 3. Web server responds by sending the requested resource, say an HTML file, crawler writes the response to the WARC file.

Step 4. Crawler makes another http request for content from the target web server, writes the request to the WARC file.

Step 5. Web server responds by sending the requested resource, say an image file, crawler writes the response to the WARC file.

https://www.hanzo.co/blog/iso-28500-warc

# On-Demand Web Archiving

**ARCHIVE-IT**

https://web.archive.org

**Save Page Now**

https:// 

SAVE PAGE

**Conifer** by RHIZOME.org

**New Capture**

URL to capture

| Add to collection | Public | | | | |
|---|---|---|---|---|---|
| Select browser | Use Current Browser | | | | |
| ▸ Session settings | **browser** | **version** | **release** | **OS** | **capabilities** |
| | Chrome | v76 | 2019-08-05 | linux | flash, autopilot |
| | Firefox | v68 | 2019-07-09 | linux | flash, autopilot |

https://conifer.rhizome.org

**Add Seeds**

Enter one seed URL per line below to add them to this collection.

✓ https://craftycocktails.com/pages/mixology-101/

Access: Private

Frequency: One-Time

Seed Type:
Standard
✓ One Page
One Page Plus External Links (One Page+)
Standard Plus External Links (Standard+)

Cancel

Add Seeds

**archive.today** webpage capture

https://archive.today

My url is alive and I want to archive its content

http://www.domain.com/url          save

www.archive-it.org

# Webrecorder and Conifer



Mode Menu / Capture Size Meter

Browser Selector

URL

Toggle Autopilot Beta Panel

https://guide.conifer.rhizome.org/



**Ilya Kreymer**
@IlyaKreymer

Thanks @machawk1 for the generous mention! 😌

A well-deserved award! Several years before Webrecorder, WARCreate really paved the way in providing an alternative, user-driven approach to small-scale web archiving!

**NDSA** @NDSA2 · Jul 28, 2020
New NDSA blog series: Catching up with past NDSA Innovation Awards Winners! First post features Mat Kelly ndsa.org//2020/07/28/ca… #DigiPres

3:22 PM · Jul 28, 2020 · Twitter Web App

https://twitter.com/IlyaKreymer/status/1288193177279541248

**Mat Kelly** and Michele C. Weigle, "WARCreate - Create Wayback-Consumable WARC Files from Any Webpage," In JCDL 2012, https://www.cs.odu.edu/~mweigle/papers/kelly-jcdl12.pdf

# Query web archives for https://www.nlm.nih.gov/hmd/

https://web.archive.org/



https://archive.today/

# Mementos for https://www.nlm.nih.gov/hmd/

https://web.archive.org/web/19990000000000*/https://www.nlm.nih.gov/hmd/

https://archive.ph/https://www.nlm.nih.gov/hmd/

# Navigating the Wayback Machine



https://web.archive.org/web/19990302100723/https://www.nlm.nih.gov/hmd/

# Querying Multiple Web Archives Using Memento Time Travel



https://timetravel.mementoweb.org/

Herbert Van de Sompel, Michael L. Nelson, and Robert D. Sanderson, "HTTP Framework for Time-Based Access to Resource States -- Memento", RFC 7089, https://www.rfc-editor.org/rfc/rfc7089.html

# Querying Multiple Web Archives Using MemGator

`curl https://memgator.cs.odu.edu/timemap/link/https://www.nlm.nih.gov/hmd/`

https://memgator.cs.odu.edu/

```
mweigle-laptop:[~]% curl https://memgator.cs.odu.edu/timemap/link/https://www.nlm.nih.gov/hmd/
<https://www.nlm.nih.gov/hmd/>; rel="original",
<https://memgator.cs.odu.edu/timemap/link/https://www.nlm.nih.gov/hmd/>; rel="self"; type="application
/link-format",
<https://webarchive.loc.gov/all/19981205095201/http://www.nlm.nih.gov/hmd/>; rel="first memento"; date
time="Sat, 05 Dec 1998 09:52:01 GMT",
```

```
<https://web.archive.org/web/19990127152411/http://www.nlm.nih.gov:80/hmd>;
rel="memento"; datetime="Wed, 27 Jan 1999 15:24:11 GMT",
<https://webarchive.loc.gov/all/19990209110725/http://www.nlm.nih.gov/hmd/>
; rel="memento"; datetime="Tue, 09 Feb 1999 11:07:25 GMT",
```

```
<https://webarchive.loc.gov/all/19990302020841/http://www.nlm.nih.gov/hmd/>; rel="memento"; datetime="
Tue, 02 Mar 1999 02:08:41 GMT",
<https://web.archive.org/web/19990302020841/http://www.nlm.nih.gov:80/hmd/>; rel="memento"; datetime="
Tue, 02 Mar 1999 02:08:41 GMT",
<https://webarchive.loc.gov/all/19990302034201/http://www.nlm.nih.gov/hmd/>; rel="memento"; datetime="
Tue, 02 Mar 1999 03:42:01 GMT",
```

**Sawood Alam**, Michael Nelson. MemGator - A Portable Concurrent Memento Aggregator: Cross-Platform CLI and Server Binaries in Go. In JCDL 2016.
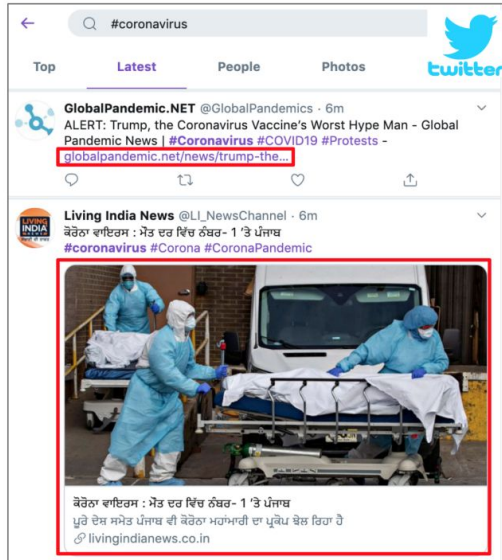
# Web archivists can create curated collections



https://archive-it.org/collections/4887

# Bootstrapping Collections from Social Media







**Alexander Nwala**, Michele Weigle, and Michael Nelson, "Bootstrapping Web Archive Collections from Social Media," In ACM Hypertext (HT) 2018, https://www.cs.odu.edu/~mln/pubs/ht-2018/hypertext-2018-nwala-bootstrapping.pdf

**Alexander Nwala**, Michele Weigle, and Michael Nelson, "Using Micro-collections in Social Media to Generate Seeds for Web Archive Collections," In JCDL 2019, https://arxiv.org/abs/1905.12220.

**Alexander Nwala**, Michele Weigle, and Michael Nelson, "Garbage, Glitter, or Gold: Assigning Multi-dimensional Quality Scores to Social Media Seeds for Web Archive Collections," In JCDL 2021, https://arxiv.org/abs/2107.02680.

# Archive-It is a popular platform for developing and accessing collections



https://archive-it.org/collections/4887

# Each seed can have multiple captures, or mementos



Collection 6435:
Northern Illinois University
Collected by Northern Illinois University

23.18%
of collections
have the behavior
Seeds Early,
Seed Mementos
Continuously

Shawn M. Jones, Alexander Nwala, Michele C. Weigle, and Michael L. Nelson, "The Many Shapes of Archive-It," In iPres 2018, https://arxiv.org/abs/1806.06878.

# Full-text search can be tricky

If term appears in multiple mementos for a seed, which one should be returned?



Sites for this collection are listed below. Narrow your results at left, or enter a search query below to find a site, specific URL or to search the text of archived webpages.

`"masks required"` [Search] [Clear]

The following results were found for the term(s): "masks required"

- **757 Sites** were found.
- 3278 result(s) for "masks required" within the page text.

| Sites | Search Page Text |

Page 1 of 164 (3,278 Total Results)     [Next Page ▶]

Sort By: Best Match

**Novel Coronavirus (COVID-19) | Department of Health**
URL: https://coronavirus.health.ny.gov/home
This text was captured on Dec 23, 2021 Show All Captures
Required in Indoor Public Places Effective Dec. 13, masks must be worn in all indoor public places unless businesses or venues implement a vaccine requirement.
*Content: text/html Size: 47 KB*
More Results from coronavirus.health.ny.gov

Dec 09, 2021

Dec 16, 2021

Dec 23, 2021

Dec 30, 2021

# Search for added phrases



Last memento before the addition: https://wayback.archive-it.org/4887/20220203045021/https://coronavirus.health.ny.gov/home
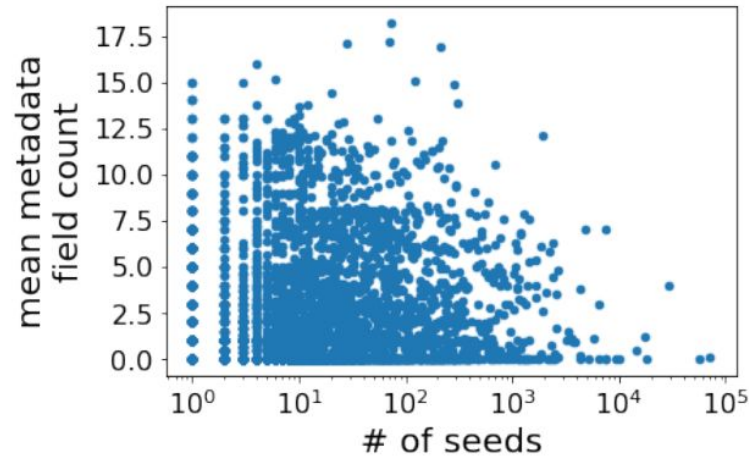First memento after the addition:   https://wayback.archive-it.org/4887/20220211212050/https://coronavirus.health.ny.gov/home

**Lesley Frew**,  Michael L. Nelson, and Michele C. Weigle, *Work in progress*

# And, facets are only useful when curators add metadata



*more seeds -> less metadata*



**Shawn M. Jones**, **Alexander Nwala**, Michele C. Weigle, and Michael L. Nelson, "The Many Shapes of Archive-It," In iPres 2018, https://arxiv.org/abs/1806.06878.

**Shawn M. Jones**, "Improving Collection Understanding for Web Archives with Storytelling: Shining Light Into Dark and Stormy Archives", PhD Dissertation, ODU, 2021, https://digitalcommons.odu.edu/computerscience_etds/131/

# NLM's COVID-19 Collection is an excellent resource

And has been carefully curated,
so several facets are available
for filtering



mweigle@odu.edu (@weiglemc) | ws-dl.cs.odu.edu (@WebSciDL) | November 17, 2022 | National Library of Medicine History Talks

25

# The National Library of Medicine Global Health Events web archive, coronavirus disease (COVID-19) pandemic collecting

Susan L. Speaker, PhD, Christie Moffatt, MSLS

## ABSTRACT

Since January 30, 2020, when the World Health Organization declared the SARS CoV-2 disease (COVID-19) t̶
emergency of international concern, the National Library of Medicine's (NLM's) Web Collecting and Archiving̶
been collecting a broad range of web-based content about the emerging pandemic for preservation in an Interne̶
other Global Health Events web collections, this content will have enduring value as a multifaceted historical r̶
and understanding of this event. This article describes the scope of the COVID-19 project; some of the content̶
websites, blogs, and social media; collecting criteria and methods; and related COVID-19 collecting efforts by̶
growing collection—2,500 items as of June 30, 2020—chronicles the many facets of the pandemic: epidemiolo̶
research; disease control measures and resistance to them; effects of the pandemic on health care institutions an̶
commerce, and many aspects of social life; effects for especially vulnerable groups; role of health disparities in̶
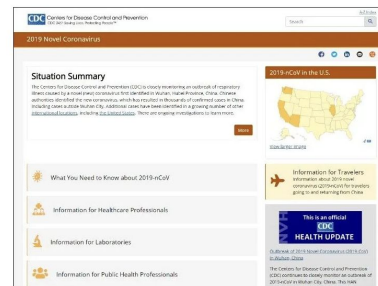mortality; and recognition of racism as a public health emergency.

Susan Speaker and Christie Moffatt, "The National Library of Medicine Global Health Events web archive, coronavirus disease (COVID-19) pandemic collecting", *Journal of the Medical Library Association,* Vol 108, No. 4 (2020), Oct 2020, https://jmla.mlanet.org/ojs/jmla/article/view/1090

Christie Moffatt, "COVID-19 Web Collecting: Reflections at One Year", Jan 2021, https://circulatingnow.nlm.nih.gov/2021/01/28/covid-19-web-collecting-reflections-at-one-year/

## COVID-19 WEB COLLECTING: REFLECTIONS AT ONE YEAR

🗓 January 28, 2021    🗂 About Us, Archives & Manuscripts, Collections, News    💬 8 comments

*By Christie Moffatt ~*

U.S. Centers for Disease Control and Prevention captured January 30, 2020. Archived pages available at:
https://wayback.archive-it.org/4887/*/https://www.cdc.gov/coronavirus/2019-nCoV/

One year ago, on January 30, 2020, the World Health Organization declared the current coronavirus disease (COVID-19) pandemic a Public Health Emergency of International Concern. The disease was yet to be officially named and there were 5 positive cases in the United States. With this designation, following National Library of Medicine Collection Development Guidelines, the NLM Web Collecting and Archiving Working Group began collecting web and social media to document the emerging global pandemic. While we had no idea for how long the pandemic would last, nor the extent of its impact, we began this effort with a belief that web content would be a significant part of the primary historical record for future research and understanding of this historic time. As the nature of web content is ephemeral, future access to the web and social media documenting the pandemic—its impact, our reactions,

# New Mexico Department of Health

March 14, 2020



September 21, 2020



September 21, 2021



March 15, 2022

November 2022
(live)

# How has a single webpage changed over time?



**Abigail Mabe**, **Dhruv Patel**, **Maheedhar Gunnam**, **Surbhi Shankar**, **Mat Kelly**, **Sawood Alam**, Michael L. Nelson, Mic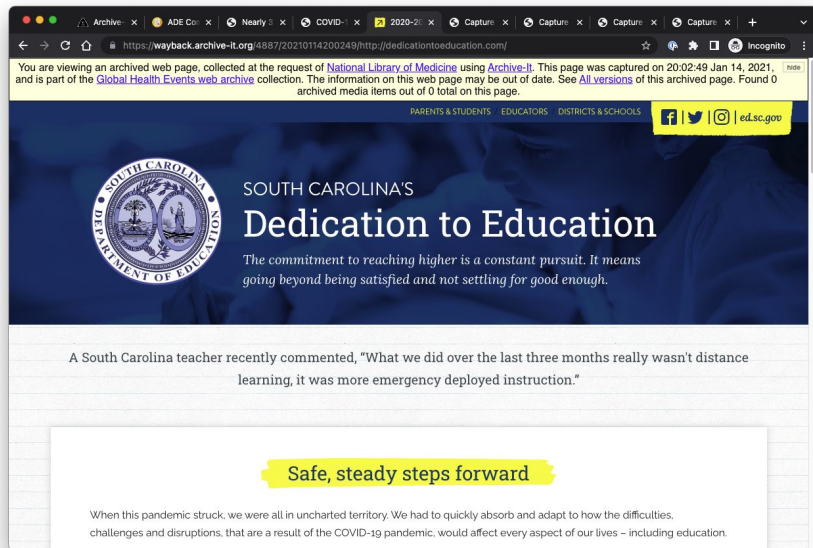hele C. Weigle, "Visualizing Webpage Changes Over Time", https://arxiv.org/abs/2006.02487, https://ws-dl.blogspot.com/2020/05/2020-05-21-visualizing-webpage-changes.html

# Getting Michiganders to take COVID-19 vaccine: 'My trust just isn't there'

Published Jul 28, 2020
First archived Sep 4, 2020



https://wayback.archive-it.org/4887/20200904135409/https://www.bridgemi.com/michigan-health-watch/getting-michiganders-take-covid-19-vaccine-my-trust-just-isnt-there

# South Carolina Dept of Education, Dedication to Education Page



January 14, 2021



*today*
redirects to SC DOE's COVID Newsroom,
https://ed.sc.gov/newsroom/covid-19-coronavirus-and-south-carolina-schools/

https://wayback.archive-it.org/4887/*/http://dedicationtoeducation.com/

# Twitter Hashtags in Collection

#2019nCoV

#COVID19

#DontBeASpreader

#IStayHomeFor

#MyCOVIDVaccine

#SafeHands

#SleeveUp

#SongsOfComfort

#StayAtHome

#TheMoment

#VaccinEquity

#WearAMask

#flattenthecurve

#healthworkers

#millionmaskchallenge

#SolidarityNotStigma

#coronavirus

#coronacomic

https://archive-it.org/collections/4887?q=hashtag&sort=f_sort_title.asc&fc=websiteGroup%3ACoronavirus+disease+%28COVID-19%29+outbreak&fc=meta_Type%3ASocial+media

# #IStayHomeFor on Twitter (April 2020)



April 26, 2020

https://wayback.archive-it.org/4887/20200426034124/https://twitter.com/hashtag/IStayHomeFor/

# #SongsOfComfort on Twitter (March 2020)



March 19, 2020

https://wayback.archive-it.org/4887/20200319102952/https://twitter.com/hashtag/SongsOfComfort/

# Pro-Tip: Archive-It mementos are also available via IA Wayback

https://wayback.archive-it.org/4887/20200319102952/https://twitter.com/hashtag/SongsOfComfort/

https://web.archive.org/web/20200319102952/https://twitter.com/hashtag/SongsOfComfort/

# #COVID19 on Twitter (Mar 2020 - May 2022)



March 4, 2020



May 3, 2022

https://wayback.archive-it.org/4887/*/https://twitter.com/hashtag/COVID19?src=hashtag/

# Issues with Archiving Twitter



June 4, 2020

Kritika Garg and Himarsha Jayanetti, "Twitter Was Already Difficult To Archive, Now It's Worse!", July 2020, https://ws-dl.blogspot.com/2020/07/2020-07-15-twitter-was-already.html

Kritika Garg, Himarsha R. Jayanetti, Sawood Alam, Michele C. Weigle, and Michael L. Nelson, "Replaying Archived Twitter: When your bird is broken, will it bring you down?," In JCDL 2021, pp. 160-169, https://arxiv.org/abs/2108.12092

# COVID-19 Web Archive at Archive-It (154 collections)



https://covid19.archive-it.org/

# From the Web to Collection to Story



*collection* - thematic sample from the web

*story* - arranged sample from a collection

The Web

Archive-It Collections

Story

**Yasmin AlNoamany,** Michele C. Weigle, and Michael L. Nelson, "Generating Stories From Archived Collections", In ACM WebSci, 2017, https://doi.org/10.1145/3091478.3091508

# Storify Example



http://web.archive.org/web/20171214031158/https://storify.com/ait_stories/2823spst0s/

**Yasmin AlNoamany**, Michele C. Weigle, and Michael L. Nelson, "Characteristics of social media stories", *IJDL,* September 2016, https://doi.org/10.1007/s00799-016-0185-3

# We've developed a suite of tools for collection storytelling

From 23,376 mementos

To a sample of 36 mementos, visualized as social cards, phrases, and images



Dark and Stormy Archives (DSA) Toolkit

**Shawn M. Jones, Himarsha R. Jayanetti,** Alex Osborne, Paul Koerbin, Martin Klein, Michele C. Weigle, Michael L. Nelson, "The DSA Toolkit Shines Light Into Dark and Stormy Archives", *The Code4Lib Journal*, Issue 53, 2022, https://journal.code4lib.org/articles/16441

# From Collections to Stories with the DSA Toolkit

# Story based on IIPC's COVID Collection (generated April 2020)



**striking image**

**most frequent entities:**
china, wuhan, cdc, japan, chinese

**most frequent sumgrams:**
covid 19, public health, the centers for disease control and prevention, the world health organization

**36 social cards, each representing a memento**

https://oduwsdl.github.io/dsa-puddles/stories/shari/2020/04/22/archive-it_collection_13529___novel_coronavirus_(covid-19)/

# Archive-It Utilities

Python library

Given an Archive-It collection, obtain collection metadata, seed list, and seed metadata.

```
In [3]: aic.get_collection_name()
Out[3]: 'Social Media'

In [4]: aic.get_collectedby()
Out[4]: 'Willamette University'

In [5]: aic.get_description()
Out[5]: 'Social media content created by Willamette University.'

In [6]: aic.get_collection_uri()
Out[6]: 'https://archive-it.org/collections/5728'

In [7]: aic.get_archived_since()
Out[7]: 'Apr, 2015'

In [8]: aic.is_private()
Out[8]: False

In [9]: len(aic.list_seed_uris())
Out[9]: 113

In [10]: aic.list_seed_uris()[0]
Out[10]: 'http://blog.willamette.edu/mba/'

In [11]: seed_url = aic.list_seed_uris()[0]

In [12]: aic.get_seed_metadata(seed_url)
Out[12]:
{'collection_web_pages': [{'title': 'Willamette MBA Blog',
    'description': ['Blog for the Willamette University Atkinson Graduate School of Management']}]]
```

Shawn M. Jones, Alexander Nwala, Michele C. Weigle, and Michael L. Nelson, "The Many Shapes of Archive-It," In iPres 2018, https://arxiv.org/abs/1806.06878.

# Hypercane provides composable operators

| | |
|---|---|
| **sample** | Produce a list of exemplars with an existing algorithm. |
| **cluster** | Divide a collection into meaningful sub-collections. |
| **score** | Rank mementos in a cluster according to certain criteria. |
| **filter** | Remove mementos from consideration. |
| **order** | Sort mementos based on the preceding score operation. |

# MementoEmbed generates a social card from a memento

<em>

http://mementoembed.ws-dl.cs.odu.edu/

**Delaware concerned over low numbers of minority populations getting vaccine**

Preserved by ARCHIVE-IT.ORG

Member of the Collection Global Health Events web archive

information about the archive

information from the memento content

So far just 5% of the people who have received the vaccine in Delaware are Black, and just 2% are Hispanic.

6ABC.COM @ 2021-06-23T17:55:45Z

information about the memento

Other Versions || Current version

other mementos for this page

current (live) version of this page

https://wayback.archive-it.org/4887/20210623175545/https://6abc.com/delaware-covid-19-vaccine-underserved-communities-black-americans-hispanic/10273580/

**Shawn Jones**, A Preview of MementoEmbed: Embeddable Surrogates for Archived Web Pages, https://ws-dl.blogspot.com/2018/08/2018-08-01-preview-of-mementoembed.html,

# Raintale can produce stories in several formats

https://oduwsdl.github.io/raintale/



HTML



Markdown



Twitter Thread

**S. M. Jones**, M. Klein, M. C. Weigle, and M. L. Nelson. "MementoEmbed and Raintale for Web Archive Storytelling", 2020, https://arxiv.org/abs/2008.00137

DSA Puddles  Visualizations of the Dark and Stormy Archives Project

2022, NOV 08

StoryGraph Biggest Story 2022-11-08 -- election day (14), the 2020 election (11), secretary state (9), two years (7), key races (7)

2022, NOV 07

StoryGraph Biggest Story 2022-11-07 -- new york (15), president joe biden (13), election day (12), gov ron (11), a campaign event (11)

2022, NOV 06

StoryGraph Biggest Story 2022-11-06 -- president joe biden (11), mehmet oz (11), new york (8), two years ago (7), election day (7)

Browse By Post Category

| All |   | Story Generated by an expert archivist |   | storygraph |   | Early Experiments with Raintale |   | CIKM 2019 study examples |
| Stories Generated By Algorithms |   | Wikipedia References |

2022 midterm elections (12), mehmet oz (9), president joe biden's (8), election day (8), john fetterman (7)

web archive

2020 election (17), president joe biden (16), election day (10), two years (9), speaker nancy pelosi

https://oduwsdl.github.io/dsa-puddles/

https://oduwsdl.github.io/dsa-puddles/categories/Stories_Generated_By_Algorithms/

# NLM Collections as a Stories



Archive-It Collection 9280: Opioid Epidemic web archive
A collection curated by National Library of Medicine

A short summary story of 20 documents automatically selected by AINoamany's Algorithm

https://oduwsdl.github.io/dsa-puddles/stories/archivei
t_collections/2022/11/02/archive-it_collection_9280/

https://oduwsdl.github.io/dsa-puddles/stories/archivei
t_collections/2022/11/02/archive-it_collection_9626/

https://oduwsdl.github.io/dsa-puddles/stories/archivei
t_collections/2022/11/03/archive-it_collection_7219/

# Challenge: How do users discover appropriate collections?
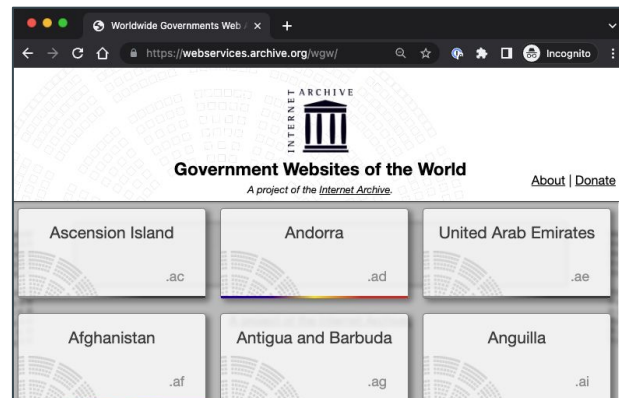


https://oduwsdl.github.io/dsa-puddles/stories/shari/2020/04/22/archive-it_collection_13529___novel_coronavirus_(covid-19)/

https://covid19.archive-it.org/

https://web.archive.org

https://webservices.archive.org/wgw

# Challenge: How can users work with collections as research data?



Archives
Unleashed

https://archivesunleashed.org/

ARCH (Archives Research Compute Hub)
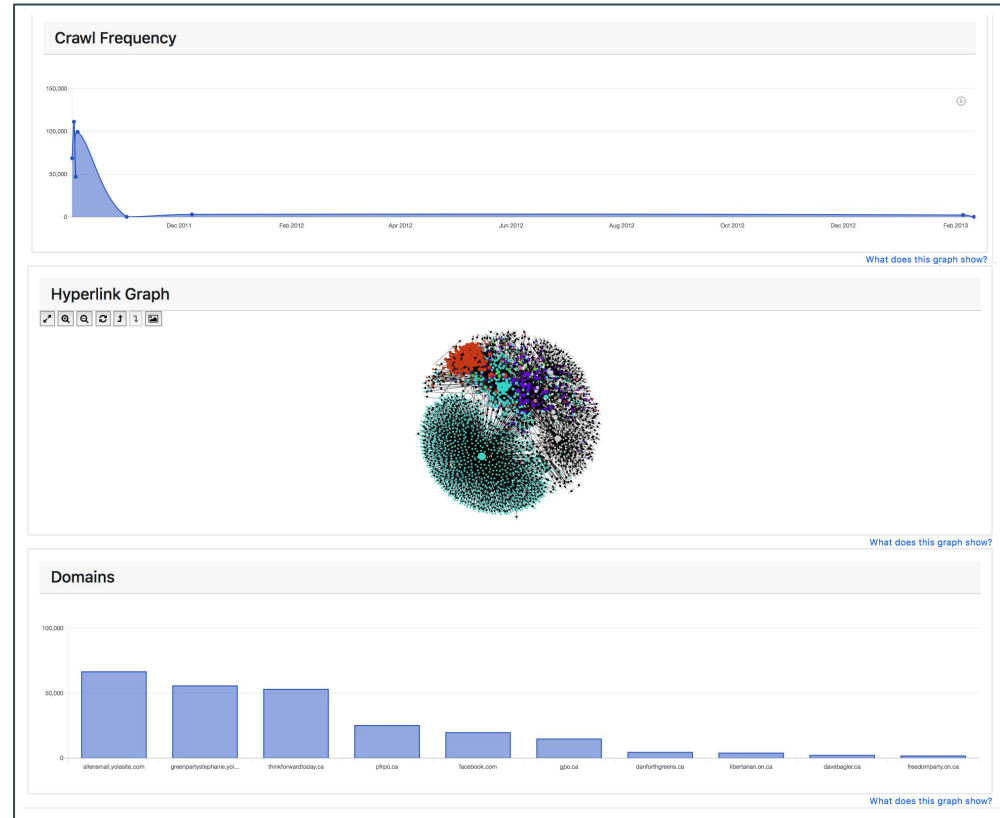https://webservices.archive.org/pages/arch

| Gephi 3.11 MB | Raw Network 1.45 MB | Domains 8.21 KB | Web Page Text 1.37 GB | Text by Domains 21.6 KB |

# Challenge: How can users work with collections as research data?



https://archivesunleashed.org/

Applications of Web Archive Research with the Archives Unleashed Cohort Program, http://blog.archive.org/2022/03/21/library-as-laboratory-recap-applications-of-web-archive-research-with-the-archive-unleashed-cohort-program/



ARCH (Archives Research Compute Hub)
https://webservices.archive.org/pages/arch

Supporting Computational Use of Web Collections, http://blog.archive.org/2022/03/07/library-as-laboratory-recap-supporting-computational-use-of-web-collections/

# Challenge: How do users discover specific mementos?

View unavailable page
- Deletion - page
- Site unavailable on live web
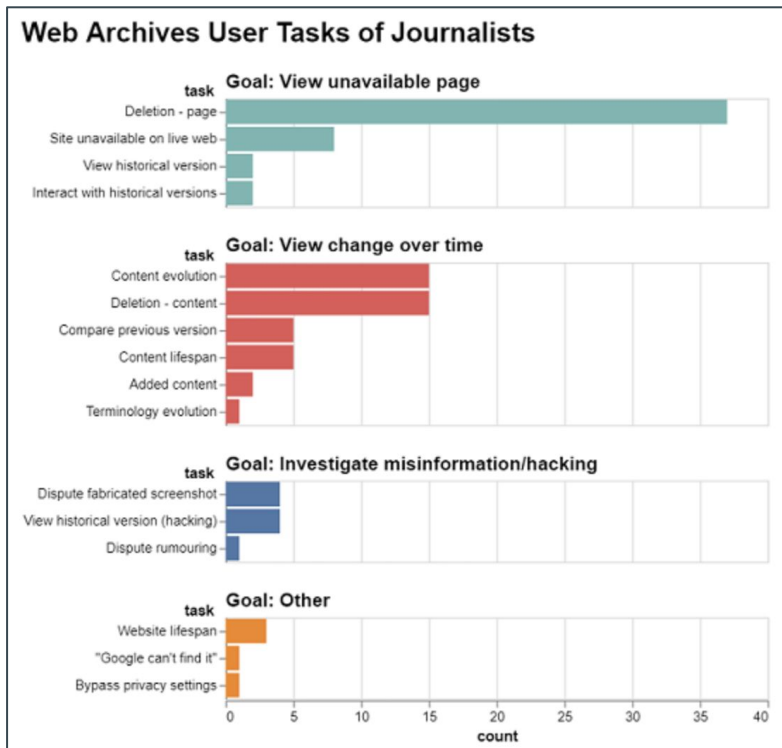
View change over time
- Content evolution
- Deletion - content

Investigate misinformation/hacking
- Dispute fabricated screenshot
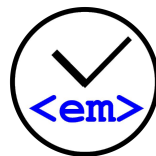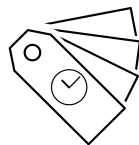- View historical version

Other
- Website lifespan



Lesley Frew, Web Archiving in Popular Media II: User Tasks of Journalists, https://ws-dl.blogspot.com/2022/08/2022-08-04-web-archiving-in-popular.html

# Conclusion

- Web archives contain traces of our history and culture.

- Web archive collections are thematic groups of archived web pages that can help with discovery.

- The DSA Toolkit can produce overviews of web archive collections, especially useful for large collections.

- Challenges related to discovery in web archives remain.

Slides: https://bit.ly/weigle-nlm-2022          Tools and Code: https://github.com/oduwsdl