

# On Tool Building and Evaluation of the Archived Web

Mat Kelly

Old Dominion University  
Web Science & Digital Libraries Research Group  
Department of Computer Science  
Norfolk, Virginia USA  
[mkelly@cs.odu.edu](mailto:mkelly@cs.odu.edu)



Seminar, Penn State University  
February 13, 2019



# Who I Am

- PhD Candidate (ADB) of Computer Science
- Defending Dissertation in 2019
- Floridian moving progressively Northward

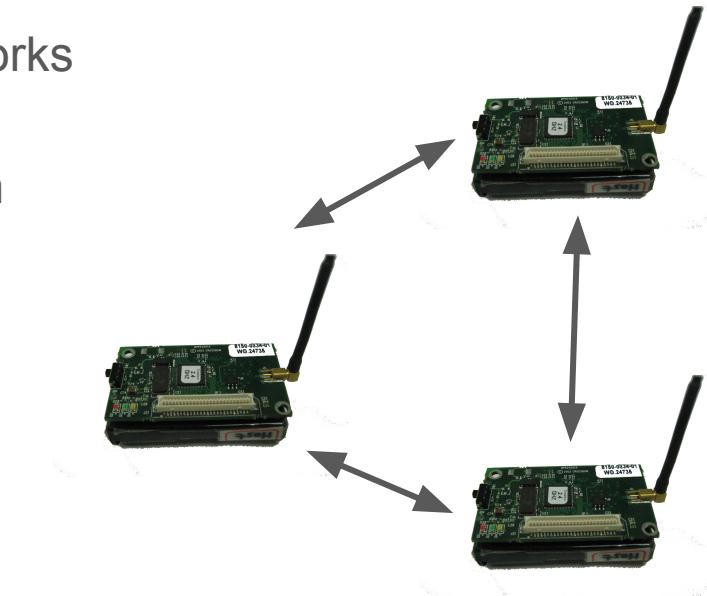


# My Research Topic

- Personal, Private, and Public Web archiving
- Technical perspective involving standards
- Tool Builder to support research
  - Often with a solution seeking a problem

# The Origin Topic

- Started off researching wireless sensor networks
- Focus: distributed emergency detection
- Exploratory NesC trilateration implementation



<https://github.com/machawk1/alert-codebase>

# Shift Topics & Labs with My PhD Advisor

(Dr. Michele C. Weigle)

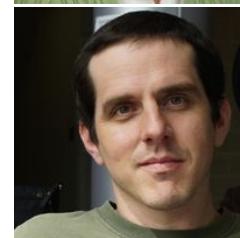
Wireless Sensor Networks



Networking (iNeTs)  
Research Group



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



Web Archiving

# Web Archives? Like Internet Archive?

- Saving pages on the Web of today
  - For exploration and research later
- The Archived Web: A Culturally significant resource
- The Internet Archive (IA) started saving the Web in 1996

But there are other institutional, public archiving efforts beyond IA

# Digital History on the Web

The screenshot shows the Penn State College of Information Sciences and Technology website. At the top, there's a navigation bar with links for Faculty Search, IST Students, Prospective Students, Alumni, Directory, and Search. A yellow "DONATE" button is also present. The main header features the Penn State logo and the text "College of Information Sciences and Technology". Below the header, a banner promotes "Managing Information, Powering Intelligence". The main content area has tabs for "THE COLLEGE", "THE EDUCATION", "THE RESEARCH", and "THE APPLICATION". A prominent graphic shows a rocket launching from clouds, with text about accepting applications for David Rusenko Entrepreneur Scholarships by April 1. Below this, there are sections for "NEWS" and "EVENTS". The news section includes images of people at events and a caption about the Future Forum. The events section lists "Training - Personal Safety in Today's Times" and "Startup Week".

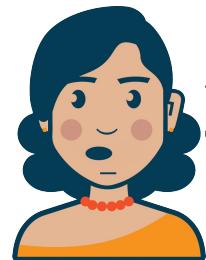
Now

This screenshot shows the Penn State IST website from February 2003, as captured by the Wayback Machine. The page title is "Welcome to The School of Info". The header includes the Penn State logo and the text "IST School of Information Sciences and Technology". The main menu links to Contacts, About IST, Dean's Message, Academic Advising, Information for Students, Research, IST Solutions Institute, and Partnerships. On the left, there's a sidebar with a search function and links for Request Information (Prospective Students, Prospective Graduate Students, Prospective Faculty, Employers, Corporate Partners, Alumni, Press Inquiries, Contact the Dean, General Inquiries) and Focus on Research (IST's Faculty, Recent Faculty Publications, Recent Faculty Grants, An Infrastructure for Innovation in Information Computing, A Field Study of Individual Differences in the Shaping of Gender and IT, ITR (GEO/SE) Ontologies in Architecture, Architecture to Support Investigation of Linked Health-Environment Information). The right side features a large image of the Information Sciences and Technology Building, a "View Live Web Cam" link, and a "Message from the Dean" from James B. Thomas, Dean, dean@ist.psu.edu. The footer includes sections for "In the Spotlight", "IST Stats at a Glance", "IST Solutions Institute", and "news@ist.psu.edu".

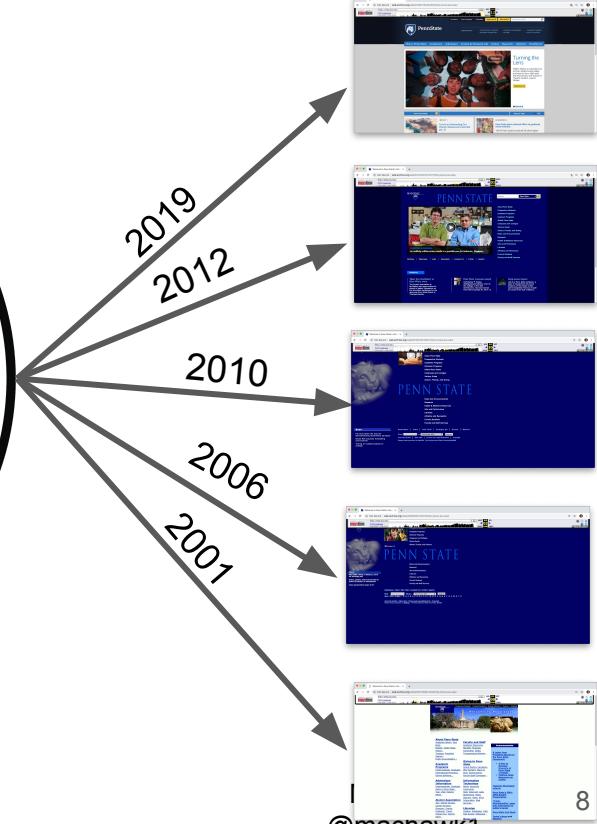
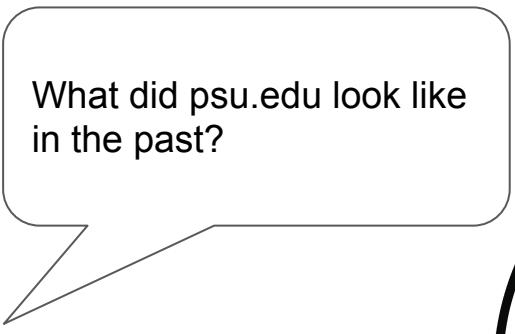
February 2003

Mat Kelly  
@machawk1

# PSU.edu of the past



What did psu.edu look like  
in the past?



# Multiple archival efforts (3 of many)





# First brush with Web Archiving beyond IA

- Tasked to revamp and maintain ArchiveFacebook
  - Mozilla Firefox extension/add-on
- Provided mechanism to allow preservation of user's Facebook contents
- Created browser-accessible cache of FB web pages

“Replay” UI

The screenshot shows the Mozilla Firefox interface with the ArchiveFB extension installed. The extension has a toolbar icon and a context menu. The context menu is open, showing options like "Archive All" (Ctrl+Shift+A), "Show in Sidebar" (Alt+K), and checkboxes for "Profile Wall", "Limit Number of Unroll", "Unroll Until End", "Friends", "Photos", "Notes", and "Events". Below the menu is a sidebar window titled "archivefb" containing a "Full Text Search" bar and a list of archive entries: "Archive (7/14/2011)", "Archive (7/14/2011)", "Archive (7/15/2011)", "Archive (7/15/2011)", and "Archive (7/15/2011)".

“Preservation” UI

The screenshot shows a Facebook profile page for "Mat Kelly" in a Mozilla Firefox window. The page includes sections for Wall, Info, Photos, Notes, Friends, and Events. A sidebar on the right lists "People You Know" and "Sponsored" posts. At the bottom, there's a "Master App Dev" section and an "Archiving Info Section". An "archivefb" sidebar is overlaid on the bottom left, showing a list of archive entries: "Wall", "Info", "Photos (70)", "Notes", and "Friends". The main content area shows Mat Kelly's profile picture, status update, and recent activity feed.

# Data Liberation vs. WYSIWYG



## Facebook Native Profile Download

Mat Kelly

Sex: Male  
Birthday: 11/16/1982  
Relationship Status: Married - Melissa Kelly  
Family: Jennifer Kelly Price (sister)  
Melissa Kelly (sister)  
Michele Glaser Kelly (mother)  
Jill Craver (cousin)  
Elle Craves (cousin)  
Craver Is Jesus (cousin)  
Steve Glaser (cousin)  
Kevin Glaser (uncle)  
Sharon Robbins Glaser (aunt)  
Carol Bartol (aunt)  
Eileen Kelly (grandmother)  
Joyce Baker (aunt)  
Kelly Baker (cousin)  
Brian Kelly (uncle)  
Renee Kelly Scarzafava (cousin)  
Rebecca Stage (cousin)  
Email: me@matkelly.com  
Facebook Profile: <http://www.facebook.com/profile.php?>

**Profile**

- Wall
- Photos
- Friends
- Notes
- Messages



## Archive Facebook Archiving Session Result

Mat Kelly - Mozilla Firefox

File Edit View History Bookmarks ArchiveFB Tools Help

resource://archivefb/data/20110719230547/index.html

Archive (7/19/2011) Archive (7/19/2011)

facebook

Mat Kelly

Mobile Applications Developer/Programmer at NASA through SSAT Studied Computer Science at Old Dominion University Lives in Virginia Beach, Virginia Married to Melissa Kelly From Labelle, Florida Born on November 16, 1982 Add languages you know Edit Profile

People You May Know

- Adalberto Gonzalez Jr.
- Tina Carter

Sponsored

State Farm Jessica Hester

When everyone forgets how to drive, Calliope Jessica Hester in Virginia Beach For a quote: 757-481-1840

Great Not Big

Advice on running a software development firm from someone who's been there.

GREAT NOT BIG

Like · 42 Like · Chat (30)

Master App Dev! 1 Chat (30)

Archiving Info Section

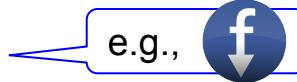
http://www.facebook.com/?ref=logo

# Research Software Beyond my Use Case

- Rapid prototyping
- Public releases of software
- Open source, permissively licensed (GPL or MIT)
- Rationales for Tool Build:
  - Data generation for further experimentation
  - Medium melding and merging (e.g., **live** & **archived** Web)
  - Exploration on the dynamics of previously unpreserved

# Lessons Learned

- Site-specific scrapers are fragile
- Little guidance on the Web on Archival Tool Building
- Testing was ad hoc and laborious (moving target) but effective
- Created Framework for MS Thesis
  - Made these sort of tools more robust and adaptive



AN EXTENSIBLE FRAMEWORK FOR CREATING  
PERSONAL ARCHIVES OF WEB RESOURCES  
REQUIRING AUTHENTICATION

by

Matthew Ryan Kelly  
B.S. June 2006, University of Florida

A Thesis Submitted to the Faculty of  
Old Dominion University in Partial Fulfillment of the  
Requirements for the Degree of

MASTER OF SCIENCE  
COMPUTER SCIENCE

OLD DOMINION UNIVERSITY  
August 2012

Approved by:

Michele C. Weigle (Director)

Michael L. Nelson (Member)

Yaohang Li (Member)

# Site-Agnostic Preservation

- Preserve everything you see!
- Created files that adhere to standard ISO28500 (Web ARChive) format
- Enable individuals to preserve any Web page from their browser

[github.com/machawk1/warcreate](https://github.com/machawk1/warcreate)

Mat Kelly and Michele C. Weigle, "WARCreate - Create Wayback-Consumable WARC Files from Any Webpage," In *Proceedings of the ACM/IEEE Joint Conference on Digital Libraries (JCDL)*. June 2012

## WARCreate - Create Wayback-Consumable WARC Files from Any Webpage

Mat Kelly  
Department of Computer Science  
Old Dominion University  
Norfolk, Virginia  
mkelly@cs.odu.edu

Michele C. Weigle  
Department of Computer Science  
Old Dominion University  
Norfolk, Virginia  
mweigle@cs.odu.edu

### ABSTRACT

The Internet Archive's Wayback Machine is the most common way that typical users interact with web archives. The Internet Archive uses the Heritrix web crawler to transform pages on the publicly available web into Web ARChive (WARC) files, which can then be accessed using the Wayback Machine. Because Heritrix can only access the publicly available web, many personal pages (e.g., password-protected pages, social media pages) cannot be easily archived in the standard WARC format. We have created a Google Chrome extension, WARCreate, that allows a user to create a WARC file from any webpage. Using this tool, content that might have been otherwise lost in time can be archived in a standard format by any user. This tool provides a way for casual users to easily create archives of personal online content. This is one of the first steps in resolving issues of "long term storage, maintenance, and access of personal digital assets that have emotional, intellectual, and historical value to individuals" [3].

### Categories and Subject Descriptors

H.3.4 [Information Storage and Retrieval]: Systems and Software; H.3.7 [Digital Libraries]: Personal Web Archiving

### General Terms

Design

### Keywords

Personal Web Archiving, WAR, WARC, Internet Archive



### 1. INTRODUCTION

ACM/IEEE - CS Joint Conference on Digital Libraries

The Internet Archive, along with many libraries and institutions, has done a remarkable job at archiving the public web. But in recent years, the web has become a home for a significant amount of original user-generated content, such as that posted on social media sites. Users are becoming increasingly aware of the need for personal web archiving [4, 5]. Unfortunately, this content is largely unavailable to standard web archives because it lives behind the "walled garden" of authentication and is part of the "deep

Copyright is held by the author/owner(s).  
*JCDL '12*, June 10–14, 2012, Washington, DC, USA.  
ACM 978-1-4503-1154-0/12/06.

on the current webpage, the user's icon in the address bar and WARC button (see Figure 1). The resources (including external scripts, CSS and images) and HTTP headers normally used by the web browser to generate a webpage and adds metadata (the *warinfo* records) to generate a WARC file that conforms to the standard's specification (Figure 2). Adherence to the specification allows the WARC to be read by Wayback.

When the compilation of the WARC file is complete, the file is downloaded to the local file system. The browser ex-

<sup>1</sup>

<https://webarchive.jira.com/wiki/display/Heritrix/Heritrix>

<sup>2</sup>Archived on July 25, 2011

<http://machawk1.com/warcreate>



# WARCreate for Google Chrome

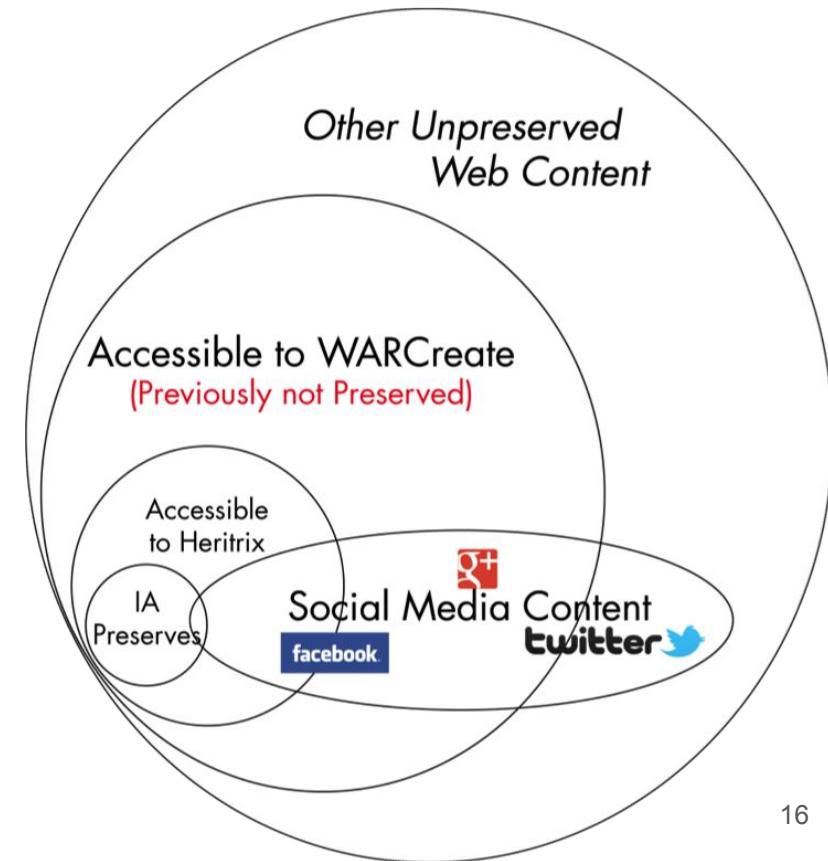
- Facilitate preservation through familiar viewport (the browser)
- Extension for Google Chrome
  - Predated WebExtensions standard API
- Easy usage:
  - One-click, current webpage → WARC
- Acts as a “buffer” until commanded to create WARC

[github.com/machawk1/warcreate](https://github.com/machawk1/warcreate)



# Archiving the Previously Unarchivable

- Target audience are for users that won't go to CLI
- Leveraging browser medium was novel and facilitated consistency





# Initial limitations

- Interacting with the File System was Limited
  - This was pre-HTML 5 File API
- Initial idea was to have Server-Side replay to also mitigate file limitations
  - This spun off “Web Archiving Integration Layer”
- As File APIs evolved, a “Server” was no longer needed for WARC generation

# High Level Overview of WARC format

Concatenated records consisting of:

- WARC Headers
- WARC Payload
  - HTTP Headers
  - HTTP Payload

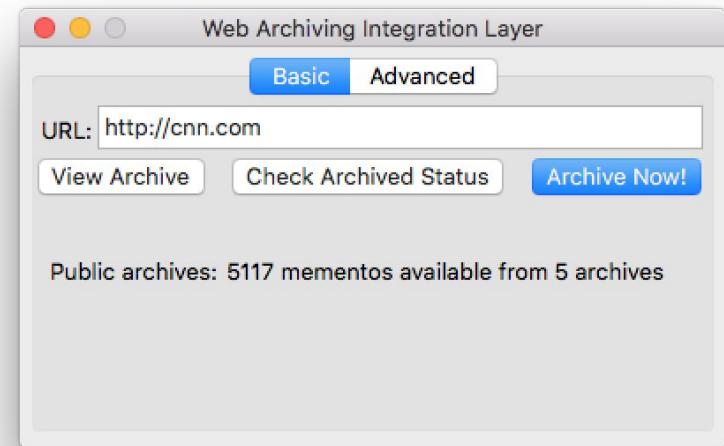
Information on HTTP request, response,  
DNS, metadata, general resources

```
20160907003819654.warc UNREGISTERED  
20160907003819654.warc x  
52 WARC/1.0  
53 WARC-Type: response  
54 WARC-Target-URI: http://ipwb.example.com/  
55 WARC-Date: 2016-09-07T00:38:19Z  
56 WARC-Record-ID: <urn:uuid:1e3907a9-2e5c-9981-6a92-964a465d998e>  
57 Content-Type: application/http; msgtype=response  
58 Content-Length: 800  
59  
60 HTTP/1.1 200 OK  
61 Host: ipwb.example.com  
62 Connection: close  
63 Content-Type: text/html; charset=UTF-8  
64 Content-Length: 666  
65  
66 <html><head>  
67 <title>InterPlanetary Wayback</title>  
68 <link rel="stylesheet" type="text/css" href="style.css">  
69 </head>  
70 <body>  
71 <h1>This is site for Space Dog</h1>  
72   
73 <p>InterPlanetary Wayback (ipwb) facilitates permanence and collabo  
74  
75 </body></html>  
76  
77  
78 WARC/1.0  
79 WARC-Type: request  
80 WARC-Target-URI: http://ipwb.example.com/style.css  
81 WARC-Date: 2016-09-07T00:38:19Z  
82 WARC-Concurrent-To: <urn:uuid:2d315cc1-a34d-3945-c5d9-ab4c7ac13fe6>  
83 WARC-Record-ID: <urn:uuid:5a1491a6-f5be-d75e-25bd-6650c69a7182>  
Line 91, Column 14 Tab Size: 4 Plain Text
```



# Web Archiving Integration Layer (WAIL)

- Written in Python, compiled to native application
  - initially OS X, Windows, and Linux
- Bundled and preconfigured “Institutional Grade” archiving tools
  - Heritrix (archival grade Web crawler)
  - OpenWayback (Web archive replay system)
- Again, simple interfaces to facilitate usage



# Software that uses WARCs

Writers/Crawlers



...

Readers/Replay Engines



**pywb**



...

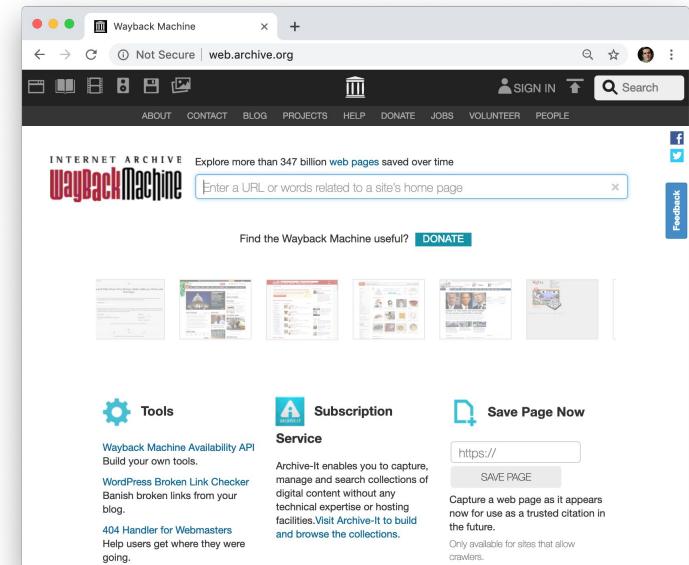
# IA's Web Archives, stored in WARCs, use same tools



# WayBack Machine

OpenWayback

# HERITRIX



A screenshot of the Wayback Machine website ([web.archive.org](https://web.archive.org)). The page features a header with the Wayback Machine logo and navigation links for About, Contact, Blog, Projects, Help, Donate, Jobs, Volunteer, and People. Below the header is a search bar and a message encouraging users to explore over 347 billion web pages. The main content area includes a "Find the Wayback Machine useful? DONATE" button and several thumbnail previews of archived web pages. At the bottom, there are sections for Tools (Wayback Machine Availability API, WordPress Broken Link Checker, 404 Handler for Webmasters), Subscription Service (Archive-It), and a "Save Page Now" feature.

# Studying the Archived Web Beyond Tools and Formats

# Study of Archiving Difficulties

- An initial examination of large Web archives
  - cf.live Web
- Which things are hard to preserve?

Mat Kelly, Justin F. Brunelle, Michele C. Weigle, and Michael L. Nelson, “On the Change in Archivability of Websites Over Time,” In Proceedings of the International Conference on Theory and Practice of Digital Libraries (TPDL). September 2013, pp. 35-47

## On the Change in Archivability of Websites Over Time

Mat Kelly, Justin F. Brunelle, Michele C. Weigle, and Michael L. Nelson  
Old Dominion University, Department of Computer Science  
Norfolk VA, 23529, USA  
{mkelly,jbrunelle,mweigle,mln}@cs.odu.edu

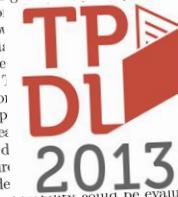
**Abstract.** As web technologies evolve, web archivists work to keep up so that our digital history is preserved. Recent advances in web technologies have introduced client-side executed scripts that load data without a referential identifier or that require user interaction (e.g., content loading when the page has scrolled). These advances have made automating methods for capturing web pages more difficult. Because of the evolving schemes of publishing web pages along with the progressive capability of web preservation tools, the *archivability* of pages on the web has varied over time. In this paper we show that the archivability of a web page can be deduced from the type of page being archived, which aligns with that page's accessibility in respect to dynamic content. We show concrete examples of when these technologies were introduced by referencing memoranda of pages that have persisted through a long evolution of available technologies. Identifying these reasons for the inability of these web pages to be archived in the past in respect to accessibility serves as a guide for ensuring that content that has longevity is published using good practice methods that make it available for preservation.

**Keywords:** Web Archiving, Digital Preservation

### 1 Introduction

The web has gone through a gradient yet demarcated series of phases in which websites were static. A web page to respond to user input more usable. Ajax is the ability to perform updates to the web page without a full page reload. This phase in the progression of the web has also progressed but in a less linear manner. The fluidity of user interaction with the web has increased over time.

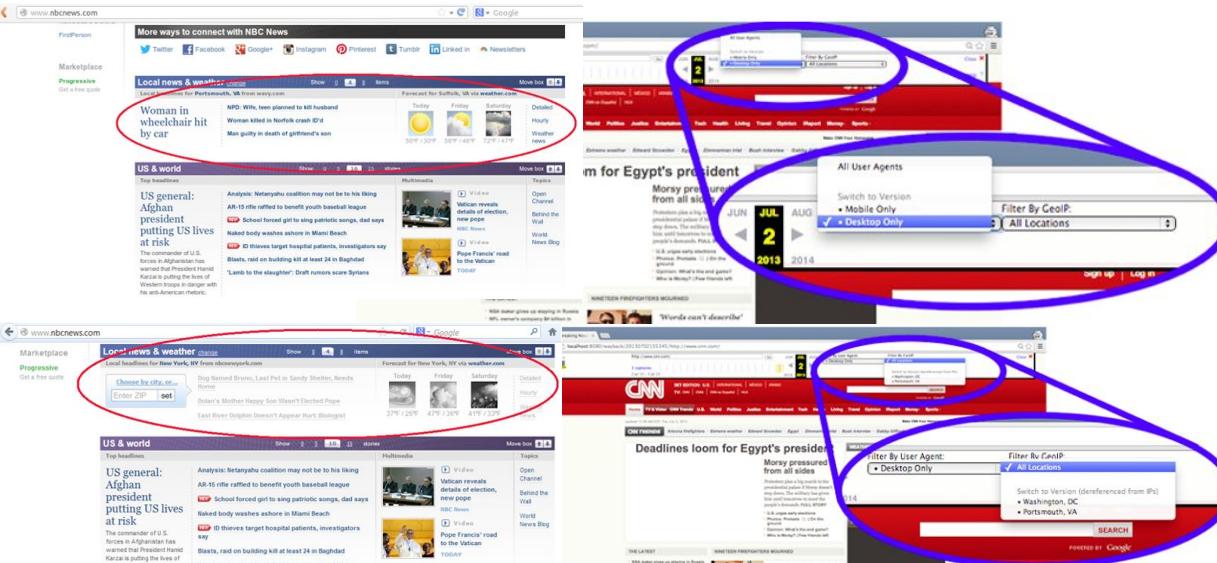
A large amount of the difficulty in archiving the web is due to the insufficient ability to capture the state of the web page. When a page is loaded, the crawler needs to execute the client-side scripts to determine the final state of the page. It should follow that the archivability could be evaluated using a consistent replay medium. The medium used to archive (normally a web crawler) is frequently different from the medium used to replay the archive (henceforth, the *web browser*, the predominant means of



Mat Kelly  
@machawk1

# Personalized Pages in the Archives

- An initial examination of large Web archives cf. live Web
- Some preserved things are personalized



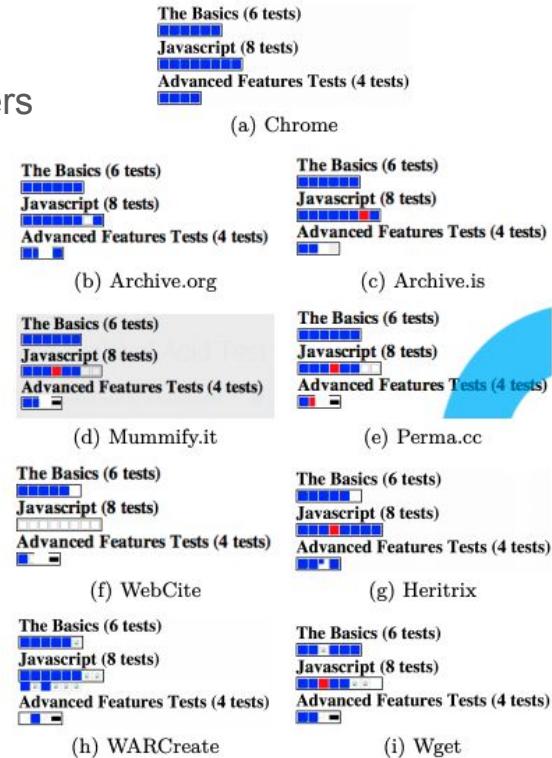
This figure shows a screenshot of the D-Lib Magazine website. The main content area displays an article titled "A Method for Identifying Personalized Representations in Web Archives" by Mat Kelly, Justin F. Brunelle, Michele C. Weigle, and Michael L. Nelson. The article is dated November/December 2013, Volume 19, Number 11/12, and includes a Table of Contents. The page features a blue header with the magazine logo and navigation links. The article text discusses methods for identifying personalized representations in web archives, mentioning the use of user-agent strings and context negotiation. It also includes several figures and tables related to their research findings.

Mat Kelly, Justin F. Brunelle, Michele C. Weigle and Michael L. Nelson, "A Method for Identifying Personalized Representations in the Archives," D-Lib Magazine, 19(11/12), 2013.

# Existing Tools' Capabilities

Punchline:

- Preservation tools lag in capability cf. Web browsers
- How well do archiving tools perform?



## The Archival Acid Test: Evaluating Archive Performance on Advanced HTML and JavaScript

Mat Kelly, Michael L. Nelson, and Michele C. Weigle  
Old Dominion University  
Department of Computer Science  
Norfolk, Virginia 23529 USA  
{mkelly,mln,mweigle}@cs.odu.edu

### ABSTRACT

When preserving web pages, archival crawlers sometimes produce a result that varies from what an end-user expects. To quantitatively evaluate the degree to which an archival crawler is capable of comprehensively reproducing a web page from the live web into the archives, the crawlers' capabilities must be tested. In this paper, we propose a set of metrics to evaluate the capability of archival crawlers and other preservation tools to pass the Acid Test concept. For a variety of web preservation tools, we examine previous captures within web archives and note the features that produce incomplete or unexpected results. From there, we design the test to produce a quantitative measure of how well each tool performs its task.

### Categories and Subject Descriptors

H.3.7 [Online Information Services]: Digital Libraries and Archives

### General Terms

Experimentation, Standardization, Web

### Keywords

Web Crawler, Web Archiving, Digital

### 1. INTRODUCTION

Since much of our cultural discourse is now online, web archiving is used for preserving web pages so they can be re-experienced at a later date. Web archiving tools capture the live web in a manner similar to web crawlers (crawlers) and preserve the page's data and contextual information about the page so it can be re-experienced. These "archival crawlers" take different approaches in digital preservation and thus their capability and scope vary.

Because archival crawlers attempt to duplicate what a user would see if he accessed the page on the live web, variance from what was seen and what would have been seen compromises the integrity of the archive. The functional difference between archival crawlers and web browsers causes this sort of unavoidable discrepancy in the archives, but it is difficult to evaluate how good of a job the crawlers did if the information no longer exists on the live web. By comparing what sort of web content is inaccurately represented or missing from the web archives, it would be useful to evaluate the completeness of archival crawlers (in respect to that of web browsers that implement the latest technologies) to determine what might be missing from the functional repertoire.

Web browsers exhibited this deviation between each other in the early days of Web Standards. A series of "Acid Tests" that implemented the Web Standards allowed each browser to visually and functionally render a web page and produce an evaluation of how well the browser conformed to the standards. In much the same way, we have created an "Archival Acid Test" to implement features of web browsers in a web page. While all standards-compliant browsers will correctly render the Acid Test, the difference can be seen in how each browser handles the page when the difference can be experienced. The difference can be experienced by comparing to what the page looked like in a variety of web browsers.

Heritrix paved the way for Internet Archive (IA) to utilize their open source Heritrix to create ARC and WARC files from web crawls while capturing all resources necessary to replay a web page [2]. Other tools have since added WARC creation functionality [3, 4, 5]. Multiple software platforms exist that can replay WARC's but IA's Wayback Machine (and its open source counterpart<sup>1</sup>) is the de facto standard.

Multiple services exist that allow users to submit URLs for preservation. IA recently began offering a "Save Page Now" feature co-located with their web archive browsing interface.

<https://github.com/iipc/openwayback>

978-1-4799-5569-5/14/\$31.00 ©2014 IEEE.



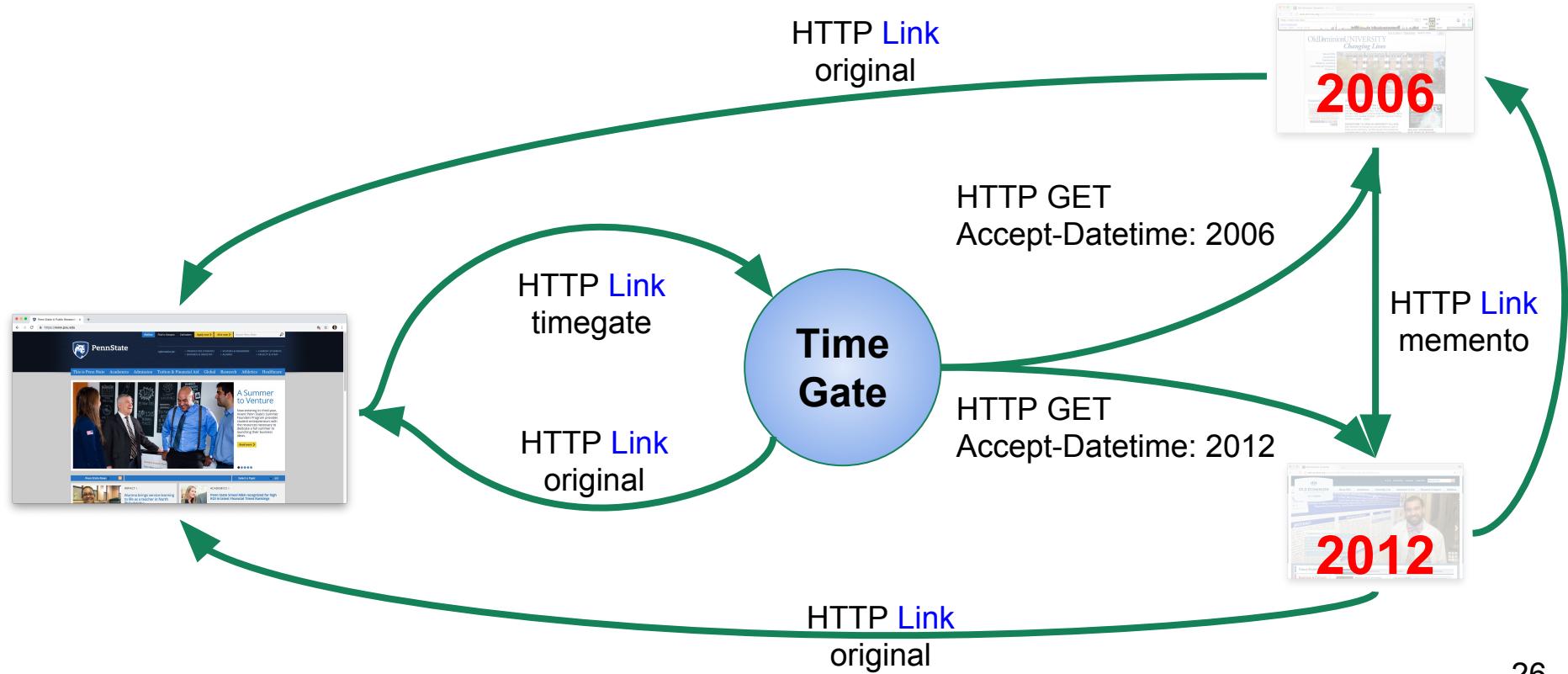
Digital Libraries 2014  
London, UK  
September 2014

looks in a variety of ways. It is used by institutions such as the Web ARCHive (WARC) to communicate payload, metadata, and other information in a single or an extensibly defined set of WARC files.

Heritrix paved the way for Internet Archive (IA) to utilize their open source Heritrix to create ARC and WARC files from web crawls while capturing all resources necessary to replay a web page [2]. Other tools have since added WARC creation functionality [3, 4, 5]. Multiple software platforms exist that can replay WARC's but IA's Wayback Machine (and its open source counterpart<sup>1</sup>) is the de facto standard.

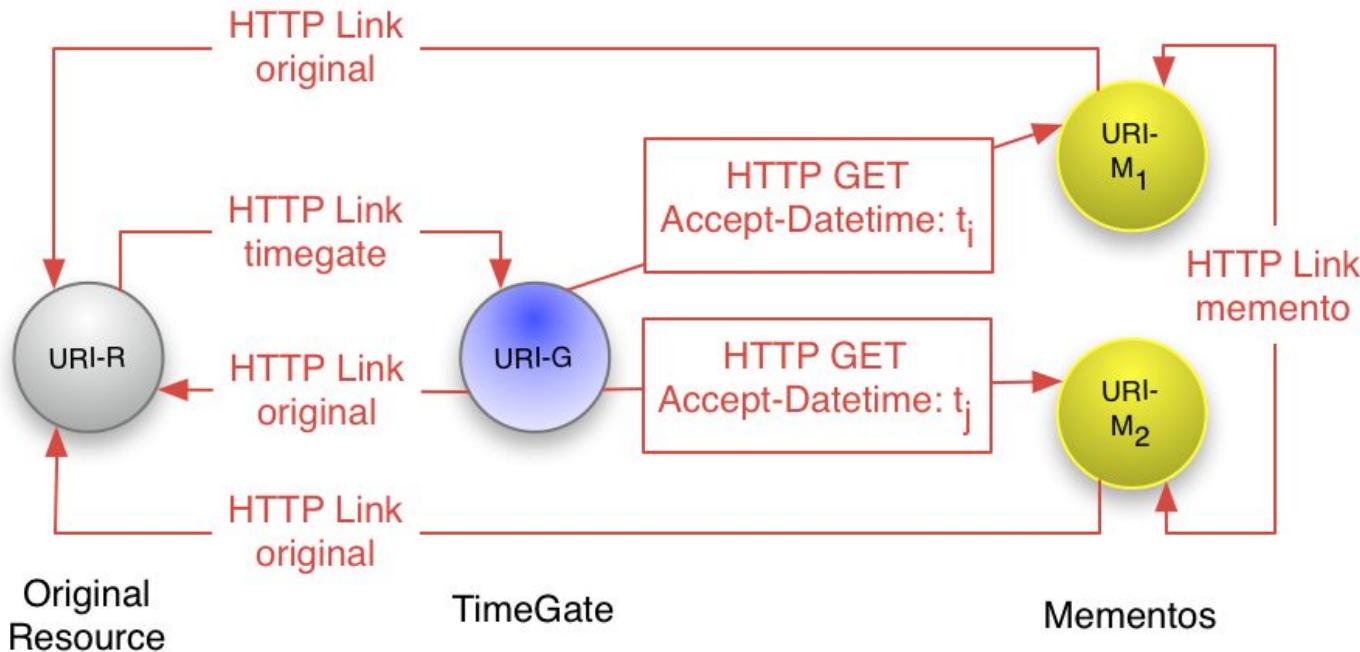
Multiple services exist that allow users to submit URLs for preservation. IA recently began offering a "Save Page Now" feature co-located with their web archive browsing interface.

# Representations can be **Linked** in time





# Memento provides relations

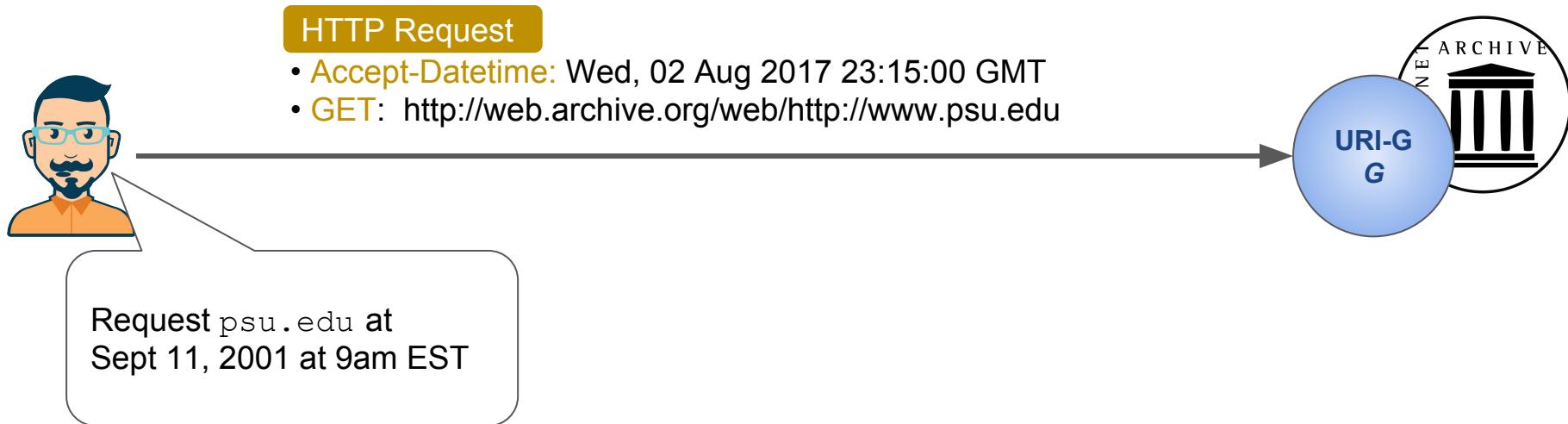


Memento Guide: Introduction. <http://www.mementoweb.org/guide/quick-intro/>, January 2015.

\* H. Van de Sompel et al. *HTTP Framework for Time-Based Access to Resource States – Memento*. IETF RFC 7089, December 2013.

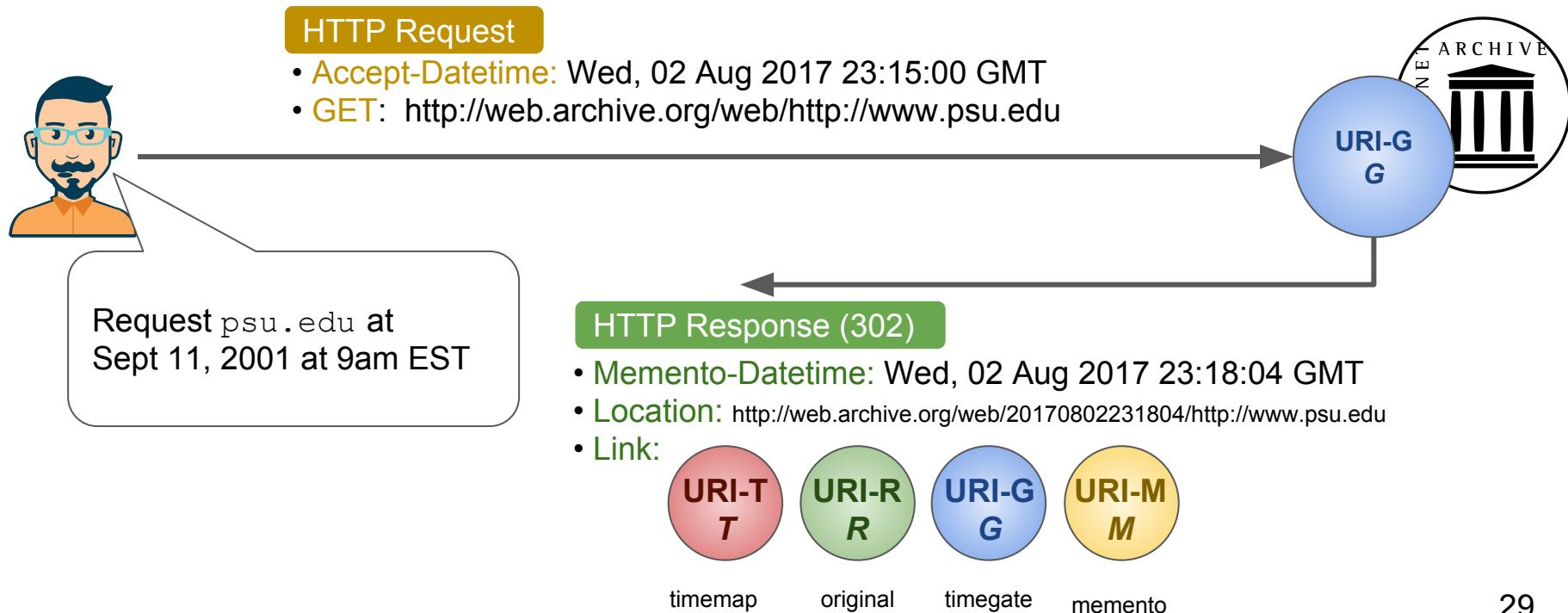


# Memento Request Example

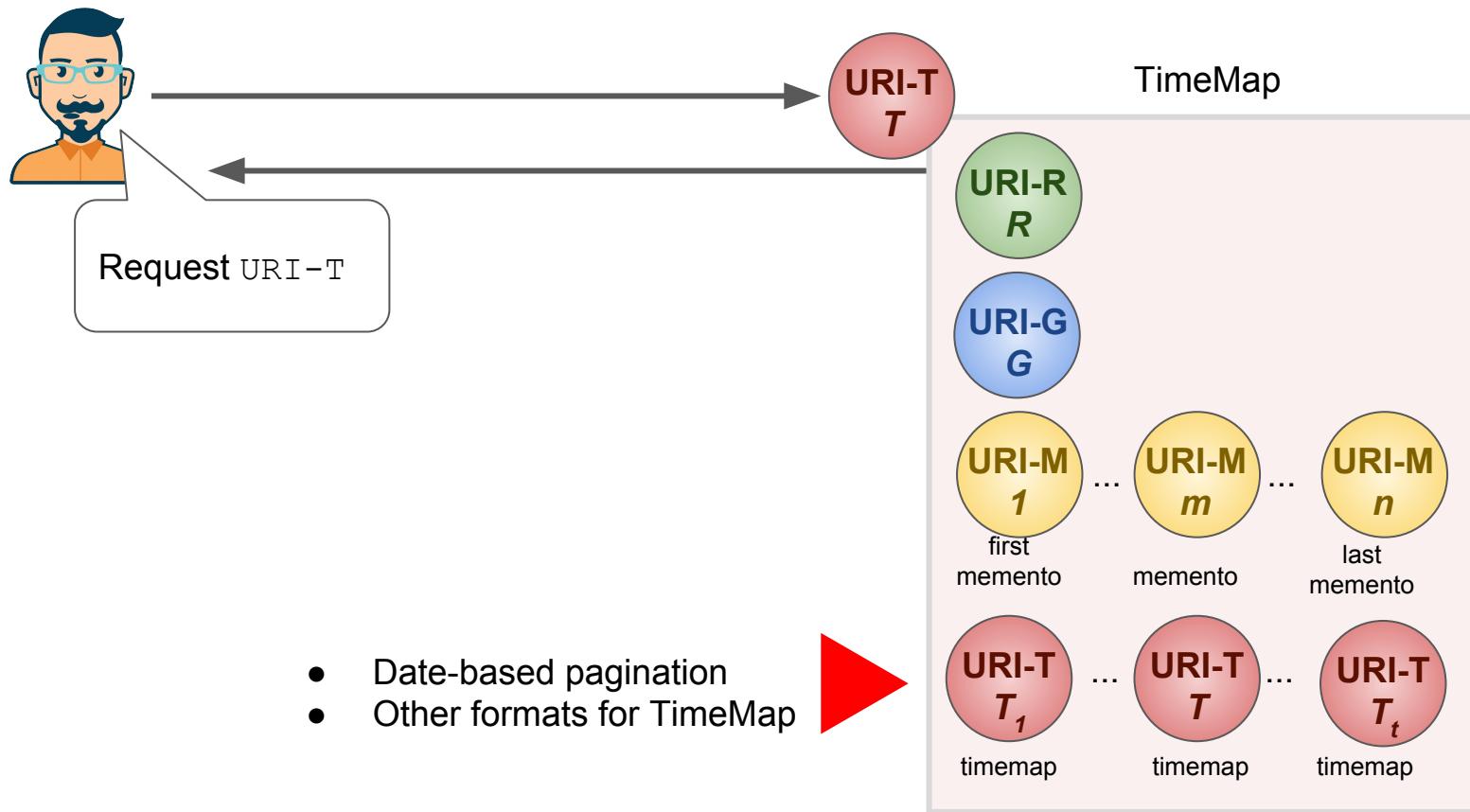




# Memento Request Example



# Background: Dereferencing a TimeMap at URI-T



# Memento “Damage” Metric

- Not all missing resource are created equal
- Multiple studies establishing metric
- Evaluated through Mechanical Turk

The collage includes:

- A white paper titled "Not All Mementos Are Created Equal: Measuring The Impact Of Missing Resources" by Justin F. Brunelle, Mat Kelly, Hany SalahEldeen, Michele C. Weigle, and Michael L. Nelson. It's published in the International Journal of Digital Libraries (Int J Digit Libr) 2015, 16:283–301. DOI: 10.1007/s00799-015-0150-6.
- A CrossMark logo indicating the document is part of a permanent record.
- A red book cover for "INTERNATIONAL JOURNAL OF Digital Libraries" from Springer.
- A white paper titled "Not all mementos are created equal: measuring the impact of missing resources" by the same authors, with a note that it was received on March 3, 2014, revised on April 22, 2015, accepted on April 22, 2015, and published online on May 6, 2015.
- A small image of the Big Ben clock tower in London.
- Text snippets from the papers discussing the damage rating algorithm and its application to the Internet Archive.
- Footnotes at the bottom of the white paper detailing author information and acknowledgments.
- Footnotes at the bottom of the red book cover page.
- Text at the bottom right: "Mat Kelly @machawk1".

# Impact of Missing Resources

**xkcd** A WEBCOMIC OF ROMANCE, SARCASM, MATH, AND LANGUAGE.

YOU CAN GET A BIG BLUE PIECE OF PAPER! IT USES SPARE WORDS TO TELL YOU WHY YOU WILL NOT GO TO SPACE TODAY.

ARCHIVE WHAT IF? BLAG STORE ABOUT

YEAR IN REVIEW

WE GO LIVE TO OUR 2015 YEAR IN REVIEW. THAT IS, IN 2015, I DONT SEE AN AURORA BOREALIS. I- WHATT?

THE NORTHERN LIGHTS PROBABLY WOULDNT FINALLY BE THE YEAR, BUT IT DIDNT HAPPEN.

OH, IHL...WHY ABOUT THE REST OF THE YEAR? I WANT TO CLEAN UP ANY BIG HEAD STORIES? DAYUM, TONI.

WELL, THAT WAS 2015 YEAR IN REVIEW. ILL BE DOING MY OWN CLEANING UP ILL BE OUTSIDE.

| < < PREV RANDOM NEXT > > |

PERMANENT LINK TO THIS COMIC: [HTTP://XKCD.COM/1302/](http://xkcd.com/1302/)  
IMAGE URL (FOR HOTLINKING/EMBEDDING): [HTTP://IMGS.XKCD.COM/COMICS/YEAR\\_IN REVIEW.PNG](http://IMGS.XKCD.COM/COMICS/YEAR_IN REVIEW.PNG)

SEARCH COMIC TITLES AND TRANSCRIPTS:  SEARCH

RSS FEED - ATOM FEED

COMICS I ENJOY:

- THREE WORD PHRASE, OLAf (new),
- SMBc, DILBERT, COT, A SOFTER WORLD, BUTTERSCAF, PERRY BIBLE
- FELLOWSHIP, QUESTIONABLE CONTENT, BUTTERCUP FESTIVAL

WARNING: THIS COMIC OCCASIONALLY CONTAINS STRONG LANGUAGE (WHICH MAY BE UNSUITABLE FOR CHILDREN), UNUSUAL HUMOR (WHICH MAY BE UNSUITABLE FOR ADULTS), AND ADVANCED MATHEMATICS (WHICH MAY BE UNSUITABLE FOR LIBERAL-ARTS MAJORS).

THIS WORK IS LICENSED UNDER A CREATIVE COMMONS ATTRIBUTION-NONCOMMERCIAL 2.5 LICENSE.

THIS MEANS YOU'RE FREE TO COPY AND SHARE THESE COMICS (BUT NOT TO SELL THEM). MORE DETAILS.

**xkcd** A WEBCOMIC OF ROMANCE, SARCASM, MATH, AND LANGUAGE.

CHRISTMAS PRESENTS FOR DATA

YEAR IN REVIEW

WE GO LIVE TO OUR 2015 YEAR IN REVIEW. THAT IS, IN 2015, I DONT SEE AN AURORA BOREALIS. I- WHATT?

THE NORTHERN LIGHTS PROBABLY WOULDNT FINALLY BE THE YEAR, BUT IT DIDNT HAPPEN.

OH, IHL...WHY ABOUT THE REST OF THE YEAR? I WANT TO CLEAN UP ANY BIG HEAD STORIES? DAYUM, TONI.

WELL, THAT WAS 2015 YEAR IN REVIEW. ILL BE DOING MY OWN CLEANING UP ILL BE OUTSIDE.

| < < PREV RANDOM NEXT > > |

PERMANENT LINK TO THIS COMIC: [HTTP://XKCD.COM/1302/](http://xkcd.com/1302/)  
IMAGE URL (FOR HOTLINKING/EMBEDDING): [HTTP://IMGS.XKCD.COM/COMICS/YEAR\\_IN REVIEW.PNG](http://IMGS.XKCD.COM/COMICS/YEAR_IN REVIEW.PNG)

SEARCH COMIC TITLES AND TRANSCRIPTS:  SEARCH

RSS FEED - ATOM FEED

COMICS I ENJOY:

- THREE WORD PHRASE, OLAf (new),
- SMBc, DILBERT, COT, A SOFTER WORLD, BUTTERSCAF, PERRY BIBLE
- FELLOWSHIP, QUESTIONABLE CONTENT, BUTTERCUP FESTIVAL

WARNING: THIS COMIC OCCASIONALLY CONTAINS STRONG LANGUAGE (WHICH MAY BE UNSUITABLE FOR CHILDREN), UNUSUAL HUMOR (WHICH MAY BE UNSUITABLE FOR ADULTS), AND ADVANCED MATHEMATICS (WHICH MAY BE UNSUITABLE FOR LIBERAL-ARTS MAJORS).

THIS WORK IS LICENSED UNDER A CREATIVE COMMONS ATTRIBUTION-NONCOMMERCIAL 2.5 LICENSE.

THIS MEANS YOU'RE FREE TO COPY AND SHARE THESE COMICS (BUT NOT TO SELL THEM). MORE DETAILS.

**xkcd** A WEBCOMIC OF ROMANCE, SARCASM, MATH, AND LANGUAGE.

CHRISTMAS SHOPPING

YEAR IN REVIEW

WE GO LIVE TO OUR 2015 YEAR IN REVIEW. THAT IS, IN 2015, I DONT SEE AN AURORA BOREALIS. I- WHATT?

THE NORTHERN LIGHTS PROBABLY WOULDNT FINALLY BE THE YEAR, BUT IT DIDNT HAPPEN.

OH, IHL...WHY ABOUT THE REST OF THE YEAR? I WANT TO CLEAN UP ANY BIG HEAD STORIES? DAYUM, TONI.

WELL, THAT WAS 2015 YEAR IN REVIEW. ILL BE DOING MY OWN CLEANING UP ILL BE OUTSIDE.

| < < PREV RANDOM NEXT > > |

PERMANENT LINK TO THIS COMIC: [HTTP://XKCD.COM/1302/](http://xkcd.com/1302/)  
IMAGE URL (FOR HOTLINKING/EMBEDDING): [HTTP://IMGS.XKCD.COM/COMICS/YEAR\\_IN REVIEW.PNG](http://IMGS.XKCD.COM/COMICS/YEAR_IN REVIEW.PNG)

SEARCH COMIC TITLES AND TRANSCRIPTS:  SEARCH

RSS FEED - ATOM FEED

COMICS I ENJOY:

- THREE WORD PHRASE, OLAf (new),
- SMBc, DILBERT, COT, A SOFTER WORLD, BUTTERSCAF, PERRY BIBLE
- FELLOWSHIP, QUESTIONABLE CONTENT, BUTTERCUP FESTIVAL

WARNING: THIS COMIC OCCASIONALLY CONTAINS STRONG LANGUAGE (WHICH MAY BE UNSUITABLE FOR CHILDREN), UNUSUAL HUMOR (WHICH MAY BE UNSUITABLE FOR ADULTS), AND ADVANCED MATHEMATICS (WHICH MAY BE UNSUITABLE FOR LIBERAL-ARTS MAJORS).

THIS WORK IS LICENSED UNDER A CREATIVE COMMONS ATTRIBUTION-NONCOMMERCIAL 2.5 LICENSE.

THIS MEANS YOU'RE FREE TO COPY AND SHARE THESE COMICS (BUT NOT TO SELL THEM). MORE DETAILS.

 $m_0$ 
 $m_1$ 
 $m_2$

# Damage in the eyes of Mechanical Turkers

- $m_0$ : live Web
- $m_1$ : comic removed
- $m_2$ : two logo images removed

The turkers selected  $m_0$  as the preferred memento 81% of the time, and more consistently for larger  $\Delta M_m$  values.

| $\Delta M_m$ | Splits |      |      |      |      |      | Total |
|--------------|--------|------|------|------|------|------|-------|
|              | 5-0    | 4-1  | 3-2  | 2-3  | 1-4  | 0-5  |       |
| 1.0          |        |      |      |      |      |      | 0.00  |
| 0.9          |        |      |      |      |      |      | 0.00  |
| 0.8          | 4      |      |      |      |      |      | 0.07  |
| 0.7          |        |      |      |      |      |      | 0.00  |
| 0.6          |        |      |      |      |      |      | 0.00  |
| 0.5          | 1      | 1    |      |      |      |      | 0.04  |
| 0.4          |        |      |      |      |      |      | 0.00  |
| 0.3          | 15     | 5    |      |      |      |      | 0.36  |
| 0.2          | 2      |      |      |      |      |      | 0.04  |
| 0.1          | 5      | 4    | 4    | 2    |      | 1    | 0.29  |
| 0.0          | 5      | 3    | 1    | 3    |      |      | 0.22  |
| Total        | 0.58   | 0.23 | 0.09 | 0.09 | 0.00 | 0.02 | 1.0   |

Justin F. Brunelle, Mat Kelly, Hany SalahEldeen, Michele C. Weigle, and Michael L. Nelson, “Not All Mementos Are Created Equal: Measuring the Impact of Missing Resources,” International Journal on Digital Libraries (IJD), 16(3), pp. 283-301.

# Impact of JavaScript on Archivability

- Missing JavaScript has big ramifications
  - Content Complexity (CC) measure
  - URIs shared over Twitter and from Archive-It collection
  - Evaluated WebCitation, wget, and the Heritrix
  - 4.2% of the Twitter collection is perfectly archived by all of these tools
  - 34.2% of the Archive-It collection is perfectly archived.



# Missing resources a direct result of JS usage



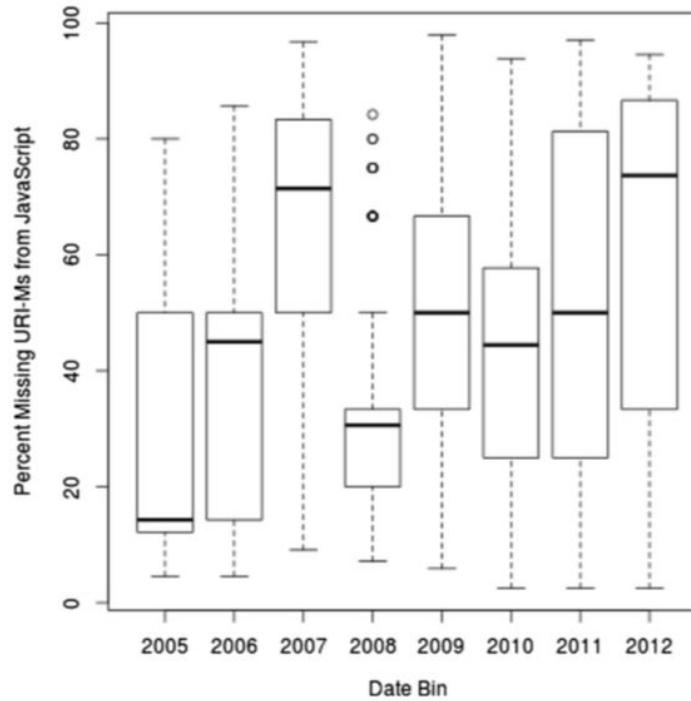
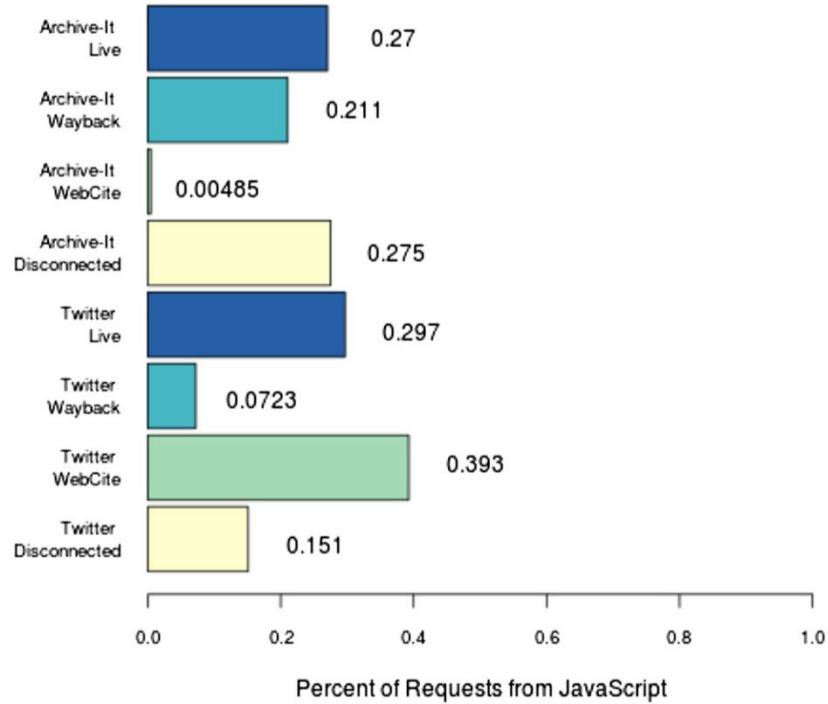
[cmt.com](http://cmt.com)

Impact of JS evident in missing  
resources

Mat Kelly  
@machawk1

35

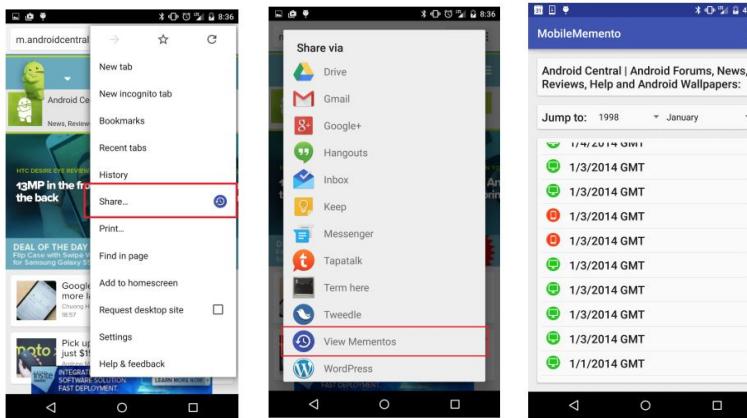
# Problem was getting worse



Justin F. Brunelle, Mat Kelly, Michele C. Weigle and Michael L. Nelson, "The Impact of JavaScript on Archivability," International Journal on Digital Libraries (IJD), 17(2), pp. 95-117.

# Mobile Mink

- Android Application
  - Few web archiving offerings in this realm
- Leveraged Native “Share” feature for archival lookup
- Identified and associated archived mobile representations



## Mobile Mink: Merging Mobile and Desktop Archived Webs

Wesley Jordan<sup>1</sup>, Mat Kelly<sup>2</sup>, Justin F. Brunelle<sup>2,3</sup>, Laura Vobrak<sup>1</sup>, Michele C. Weigle<sup>2</sup>, and Michael L. Nelson<sup>2</sup>  
<sup>1</sup> New Horizons Regional Education Center Governor's School for Science and Technology  
<sup>2</sup> Old Dominion University, Department of Computer Science  
<sup>3</sup> The MITRE Corporation

### ABSTRACT

We describe the mobile app *Mobile Mink* which extends Mink, a browser extension that integrates the live and archived web. Mobile Mink discovers mobile and desktop URLs and provides the user an aggregated TimeMap of both mobile and desktop mementos. Mobile Mink also allows users to submit mobile and desktop URLs for archiving at the Internet Archive and Archive.today. Mobile Mink helps to increase the archival coverage of the growing mobile web.

### Categories and Subject Descriptors

H.3.7 [Online Information Services]: Digital Libraries

### General Terms

Design; Experimentation; Measurement

### Keywords

Web Archiving; Digital Preservation; Memento; TimeMaps

### 1. INTRODUCTION

Mink [4] is a browser extension for Google Chrome that more closely integrates the past and present web. Mink uses the Memento framework [8] to present archived versions of visited pages to the user, allowing the users to seamlessly navigate between the archived and live web.

Memento is a framework that standardizes Web archive access and terminology. Live web resources are identified by URI-R. Archived versions of URI-Rs are called *mementos* and are identified by URI-M. Memento TimeMaps are machine-readable lists of mementos (at the level of single-archives or aggregation-of-archives) sorted by archival date.

While Mink works well in the traditional, desktop-oriented web, the mobile web continues to be less prominent in the archives. This phenomenon persists even as mobile devices grow in power, use, and ubiquity and the mobile web continues to grow and become more prevalent [9]. Because of

their prevalence on the web, it is increasingly important to archive mobile resources and representations. However, because mobile resources are not always directly linked from their desktop counterparts, it is difficult for crawlers to find pages in the mobile web [2].

Mobile Mink is a mobile application that – in the same way Mink integrated the past and present desktop webs – bridges the mobile and desktop webs. Mobile Mink uses URI permutations to discover mobile and desktop versions of the same resource. Mobile Mink provides the user an aggregate TimeMap of mobile and desktop mementos, and provides the opportunity to submit the mobile and desktop URI-Rs to the Save Page Now service at the Internet Archive [6] and Archive.today [1].

### 2. AGGREGATE TIMEMAPS

Mobile Mink is an Android application that is currently in development and will be released for download in the Google Play app store. Much like its desktop browser parent, Mobile Mink offers a TimeMap of mementos that allows the user to navigate between the past and present webs. Mobile Mink also allows the user to submit mobile and desktop URI-Rs to be archived by archival services.

When using a web browser native to the Android operating system, the user is presented with an expandable menu in the top right of the browser window (called a “view as list”). Selecting this sign opens a menu of options, one of which is the option to “Share” the page (Figure 1(a)).

Mobile Mink adds the option to “View Mementos” of the currently viewed page to the list of sharing options (Figure 1(b)).

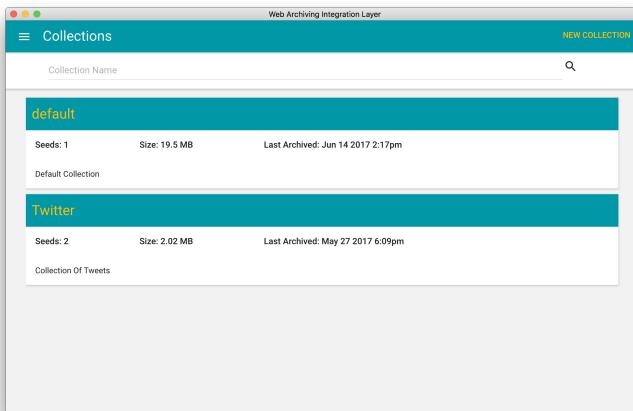
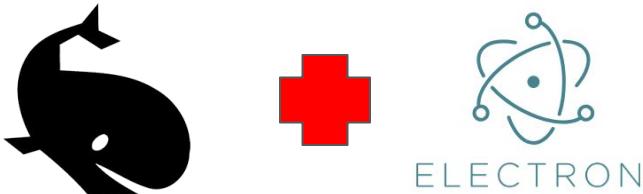
Selecting the option of viewing mementos begins the process of discovering mobile and desktop URLs of the current URL-R. First, Mobile Mink identifies the URL-R of the currently viewed page. Mobile Mink identifies the URL-R as either a desktop URL or a mobile URL. Second, if the URL is a desktop URL, Mobile Mink translates the URL to a mobile URL; if the URL is a mobile URL, Mobile Mink translates the URL to a desktop URL. We use the same URI modifications as in Schneider and McCown’s work [7] and test for the mobile URL’s existence on the live web (i.e., returns an HTTP 200 response) in the archives (returns a TimeMap of cardinality  $> 0$  from the Memento aggregator).

Note that our previous research demonstrated that differentiating between the mobile and desktop versions of a page can be difficult if the same URL is used to identify the mobile and desktop representations, and only content-negotiation based on the user-agent is used by the server to

Permit the making of digital or hard copies of part or all of this work for personal or for profit or commercial advantage, and that copies bear this notice and the full citation on the first page. Copyright for third-party components of this work must be honored. For all other uses, contact the owner/author(s). Copyright is held by the author(s).  
ICWS '15, June 21–25, 2015, Knoxville, Tennessee, USA.  
ACM 978-1-4503-3594-2/15/06. \$10.275696.  
<https://doi.org/10.1145/2756940.2756956>

# WAIL reimagined

- Archive from the desktop with higher fidelity than conventional archiving tools



**WAIL: Collection-Based Personal Web Archiving**

John A. Berlin, Mat Kelly, Michael L. Nelson, Michele C. Weigle  
(jberlin,mklly,mln,mweigle)@cs.odu.edu

**ABSTRACT**

Web Archiving Integration Layer (WAIL) is a desktop application written in Python that integrates Heritrix and OpenWayback. In this work we recreate and extend WAIL from the ground up to facilitate collection-based personal Web archiving. Our new iteration of the software, WAIL-Electron, leverages native Web technologies (e.g., JavaScript, Chromium) using Electron to open new potential for Web archiving by individuals in a stand-alone cross-platform native application. By replacing OpenWayback with PyWB, we provide a novel means for personal Web archivists to curate collections of their captures from their own personal computer rather than relying on an external archival Web service. As extended features we also provide the ability for a user to monitor and automatically archive Twitter users' feeds, even those requiring authentication, as well as provide a reference implementation for integrating a browser-based preservation tool into an OS native application.

**KEYWORDS**  
Personal Web Archiving

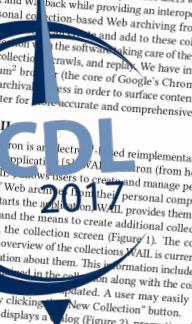
**ACM Reference format:**  
John A. Berlin, Mat Kelly, Michael L. Nelson, Michele C. Weigle. 2017. WAIL: Collection-Based Personal Web Archiving. In *Proceedings of Joint Conference on Digital Libraries, Toronto, Ontario, Canada, June 2017 (JCDL '17)*. 2 pages.  
DOI: 10.XXXXXXXX

**1 INTRODUCTION**  
Subscription-based Web archiving services like Archive-It allow users with limited technical knowledge to create and replay personalized collections of Web archives. Archive-It provides its users with a simple interface to create collections and to launch comprehensive archival crawls. Similar to Archive-It is Webrecoorder<sup>1</sup>, which allows any user to register for the service and provides them with the ability to create and manage personalized collections of Web pages. But unlike Archive-It, Webrecoorder requires its user to manually drive the preservation process or upload content for replay, only providing its users up to five gigabytes of storage. Individuals that wish to freely (*gratis et libere*) archive Web pages without arbitrary restrictions beyond the limitations of their personal computers using institutional grade tools must setup an archival Web crawler (e.g., Heritrix) and replay system (e.g., Wayback), time consuming and technical tasks potentially beyond the individual's skill level. Even if a user is able to successfully set up these tools, they must also configure the crawls via Heritrix and come up with their own means of associating the Web archives to each other for collection-based replay via Wayback.

We have developed a tool that provides users with access to both Heritrix and Wayback while providing an interoperable mechanism for personal collection-based Web archiving from their personal computers. We can crawl and add to these collections through the WAIL interface, taking care of the details in managing crawls, and replay. We have integrated a native Chromium<sup>2</sup> browser (the core of Google's Chrome Web browser) into the archival process in order to surface content specific to sites like Twitter for more accurate and comprehensive preservation.

**2 WAIL**  
WAIL-Electron is an Electron-based reimplementing of the original WAIL application [5]. WAIL-Electron (from here on referred to as WAIL) allows users to create and manage personalized collections of Web archives from their personal computers. When a user first starts the application, it provides them with a default collection and the means to create additional collections straight away from the collection screen (Figure 1). The collection view displays an overview of the collections. WAIL is currently managing two collections in the collection view along with the collection's size and last archived date. A user may easily create a new collection by clicking the "New Collection" button.

Doing so displays a dialog (Figure 2), prompting the user for a collection name, title, and description. These values are propagated to the WAIL interface and are viewable when replaying the collection through Wayback. When viewing a collection, WAIL displays



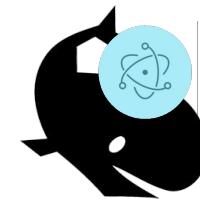
JCDL  
2017

Permission is granted to make digital or hard copies of part or all of this work for personal or classroom use is granted without prior permission or fee. If copies are not made or distributed for profit or commercial advantage and the copyright holder bears this notice and the full citation for the original publication in this proceedings, for all other use, contact the owner/authors(s).  
© 2017 Copyright held by the owner/authors. XXX-YYYY-ZZ-AAA/BB/CC... \$15.00  
DOI: 10.XXXXXXXX

<sup>1</sup><https://webrecoorder.net/>

<sup>2</sup><https://www.chromium.org/>

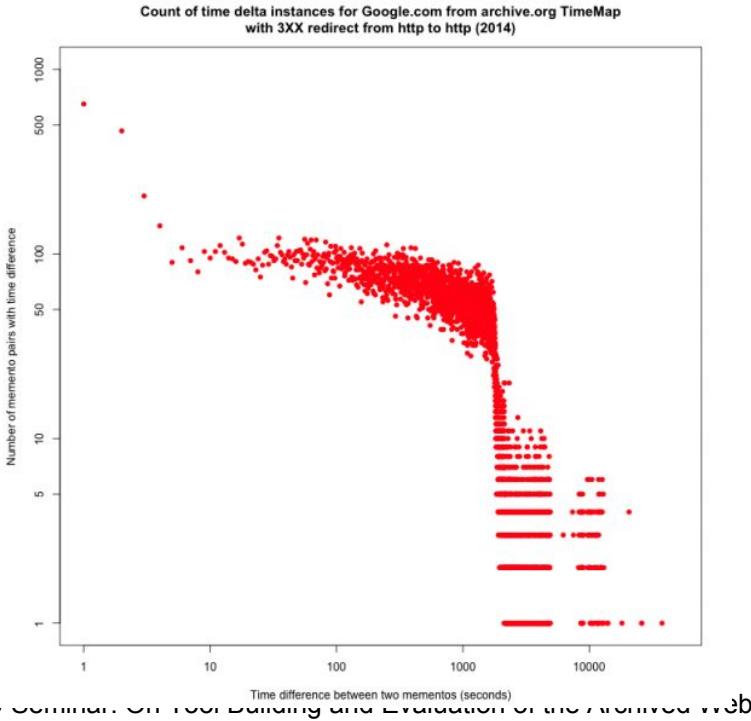
<sup>3</sup><http://electron.atom.io/>



|                              |   |   |
|------------------------------|---|---|
| <b>Code</b>                  | Python  | HTML, JavaScript, Electron  |
| <b>Archival organization</b> | Single Collection   | Multiple Collections  |
| <b>User interface</b>        | System Native   | Material  |
| <b>Archival Crawler</b>      | Heritrix  | Heritrix, node-warc   |
| <b>Archival Replay</b>       | OpenWayback   | pywb  |
| <b>Release</b>               | macOS, Windows  | macOS, Windows, Linux   |
| <b>Source</b>                | <a href="https://github.com/machawk1/wail">github.com/machawk1/wail</a> | <a href="https://github.com/n0tan3rd/wail">github.com/n0tan3rd/wail</a> |

# Conicalization's effects over time

URI coalescence considered harmful for archives



## Impact of URI Canonicalization on Memento Count

Mat Kelly, Lulwah M. Alkwai, Sawood Alam,  
Michael L. Nelson, and Michele C. Weigle  
Old Dominion University  
Department of Computer Science  
Norfolk, Virginia, USA  
{mkelly, lalkwai, salam, mln, mwweigle}@cs.odu.edu

Herbert Van de Sompel  
Los Alamos National Laboratory  
Los Alamos, New Mexico, USA  
herberty@lanl.gov

### ABSTRACT

Memento TimeMaps [5] list identifiers for archived web captures (URI-Ms). When some URI-Ms are dereferenced, they point to a different URI-M instead of a unique representation at the datatype. This suggests that confidently obtaining an accurate count quantifying the number of non-reducing captures for a Original Resource URI (URI-R) is not possible using a TimeMap alone and that the magnitude of a TimeMap is not equivalent to the number of representations it identifies. This work represents an abbreviated version of the full technical report describing this phenomena in depth [3]. For google.com we found that 84.9% of the URI-Ms in a TimeMap result in an HTTP redirect when dereferenced. The full study applies this technique to seven other URI-Rs of large Web sites and 14 academic institutions. Using a ratio metric for the number of URI-Ms without redirects to those requiring a redirect when dereferenced, five of the large Web sites' and two of the academic institutions' TimeMaps had a ratio of less than one, indicating that more than half of the URI-Ms in these TimeMaps result in redirects when dereferenced.

### 1 INTRODUCTION

Web archives return TimeMaps with a list of URI-Ms for HTTP transactions observed at archiving time. TimeMaps have generally been used as a compact representation of a URI-R present in a TimeMap. However, TimeMaps may include URI-Ms for archiving representations, redirections, and errors [2]. For example, if the URI-Ms for http://simon.com produce an HTTP 302 redirect to another URI-M in the TimeMap that returns a status code of OK, TimeMaps do not explicitly return a "count" field to indicate the number of mementos listed in the TimeMap that produce a non-redirection HTTP status code when dereferenced. The heuristic of determining how many captures represented by URI-Ms in a TimeMap cannot be completed without dereferencing.

Dereferencing in a Web archive can be attributed to a variety of canonicalization rules [3]. Preserving and replaying these redirects allows an archive to accurately reproduce the HTTP transactions that would have occurred when the URI being accessed resided on the live Web. Because of the potential for redirection, the heuristic of counting URI-Ms with reduction rules of "memento" is an inaccurate means of determining the number of unique representations inferred from a TimeMap. We further emphasize the distinction per the Memento specification that the identifiers for mementos

Table 1: Google over time (abbreviated), bucketed by year, based on IA mementos extracted from the TimeMap.  $M_{T_M}$  is the memento count based solely on the data in the TimeMap, while  $M_{RC}$  is the count based on exclusion of redirects when dereferencing.  $DI$  is the ratio of non-reducing mementos to redirecting mementos.

(URI-Ms) in a TimeMap are identifiers for archived HTTP transactions (e.g., transmission of HTTP 2XX, 3XX, 4XX, etc.) rather than identifiers for representations. Because the number of URI-Ms in a TimeMap not necessarily resolve to unique mementos when archival redirects are followed, we examined the mementos from contemporary TimeMaps to evaluate the patterns and schemes used in Memento canonicalization. Through this, we identify the difference between the number of mementos available as compared to the TimeMap through native "rel" counting heuristics to the temporally unique mementos identified once these mementos are dereferenced.

### 2 BACKGROUND AND RELATED WORK

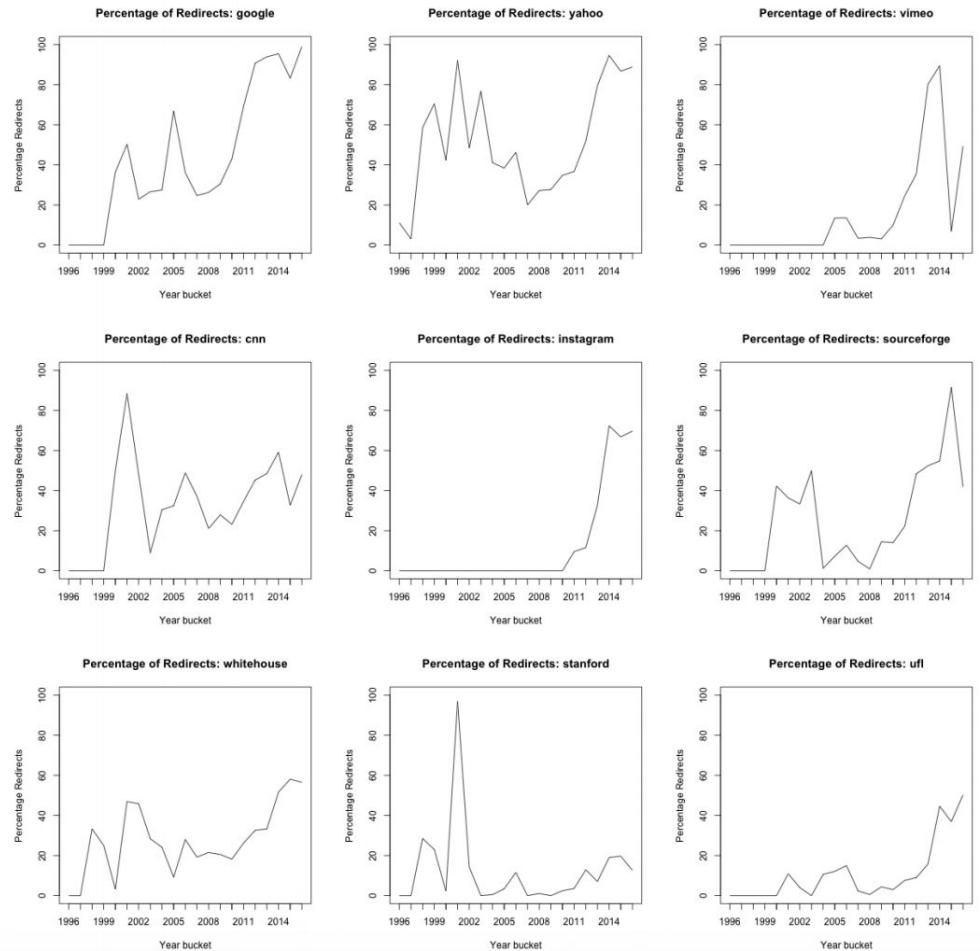
URI canonicalization associates differently formatted URLs with their after-the-fact clustering of URLs that likely refer to the same resource. As URI schemes from a Web site change over time, canonicalization is critical for retaining a cohesive, comprehensive listing of the mementos available for a Web page.

AlSum et al. [1] analyzed memento redirection patterns relating to HTTP redirects to supply the user with the correct memento when a redirect is encountered in the archives. They introduced the notion of "URI-stability" to give a quantitative measure of the presence of HTTP 3XX status codes that result when URI-Ms in TimeMaps are dereferenced.

# Many URI-Ms are actually redirects

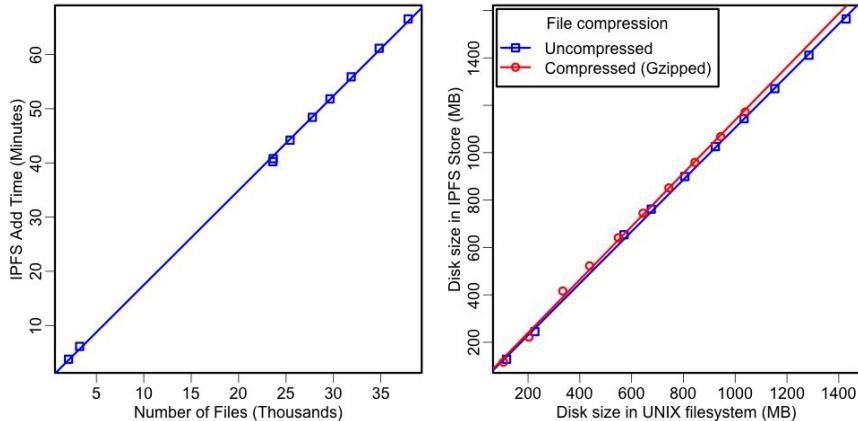
| host        | % 3XX | % 200 | $M_{TM}$ | $DI$  | host      | % 3XX | % 200 | $M_{TM}$ | $DI$  |
|-------------|-------|-------|----------|-------|-----------|-------|-------|----------|-------|
| google      | 84.89 | 15.11 | 695,525  | 0.178 | stanford  | 62.14 | 37.84 | 19,309   | 0.609 |
| yahoo       | 88.16 | 11.83 | 418,896  | 0.134 | princeton | 60.10 | 39.88 | 9,355    | 0.663 |
| sourceforge | 73.34 | 26.63 | 31,408   | 0.363 | columbia  | 48.01 | 51.88 | 9,882    | 1.082 |
| instagram   | 67.32 | 32.65 | 55,228   | 0.485 | harvard   | 33.91 | 65.96 | 7,699    | 1.948 |
| vimeo       | 57.04 | 42.94 | 199,262  | 0.752 | caltech   | 33.13 | 66.86 | 5,474    | 2.017 |
| cnn         | 49.97 | 50.01 | 87,148   | 1.001 | mit       | 26.57 | 73.24 | 6,379    | 2.763 |
| wikipedia   | 44.62 | 55.19 | 25,973   | 1.240 | gatech    | 26.03 | 73.94 | 3,907    | 2.841 |
| whitehouse  | 44.57 | 55.24 | 26,006   | 1.243 | ufl       | 24.76 | 75.23 | 4,927    | 3.038 |
|             |       |       |          |       | vt        | 23.07 | 76.92 | 4,061    | 3.334 |
|             |       |       |          |       | lsu       | 15.06 | 84.93 | 2,974    | 5.638 |
|             |       |       |          |       | nsu       | 13.82 | 86.00 | 1,208    | 6.233 |
|             |       |       |          |       | odu       | 9.727 | 90.27 | 1,727    | 9.279 |
|             |       |       |          |       | tcc       | 5.429 | 94.57 | 884      | 17.41 |

Mat Kelly, Lulwah M. Alkwai, Michael L. Nelson, Michele C. Weigle, and Herbert Van de Sompel, “Impact of URI Canonicalization on Memento Count,” Technical Report arXiv:1703.03302, 2017.



# InterPlanetary Wayback (ipwb)

- Personal archives are more resilient when propagated.
- How much does it cost to have resilient personal archives?



**InterPlanetary Wayback: The Permanent Web Archive**

Sawood Alam, Matt Kelly, and Michael L. Nelson  
Old Dominion University, Department of Computer Science, Norfolk VA, 23529, USA  
(salam,mkelly,mnln)@cs.odu.edu

**ABSTRACT**

To facilitate preservation and collaboration in web archives, we built InterPlanetary Wayback to disseminate the contents of WARC files into the IPFS network. IPFS is a peer-to-peer content addressable system that inherently supports de-duplication and facilitates efficient storage. We split the header and body of WARC response records to build efficient indexing and retrieval at the time of replay. From a 1.0 GB sample, we find that the original WARC contains 21,994 mementos and can be indexed and retrieved in less than 10 minutes. We also show that it is feasible to store and propagate large amounts of data in IPFS, such as the 1.0 GB sample, in less than 66 minutes.

**JCDL 2016**  
**NEWARK**

**TPDL 2016 HANNOVER**

**Keywords:** Web Archives, Memento, Peer-to-Peer, IPFS

**InterPlanetary Wayback: Peer-To-Peer Performance of Web Archives**

Mat Kelly, Sawood Alam, Michael L. Nelson, and Michele C. Weigle  
Old Dominion University, Department of Computer Science, Norfolk VA, 23529, USA  
(mkelly,salam,mnln,mweigle)@cs.odu.edu

**Abstract.** We have integrated Web ARCHive (WARC) files with the peer-to-peer content addressable InterPlanetary File System (IPFS), to allow the payload content of web archives to be easily propagated. We provide an archival replay system extended from ipwb to fetch the WARC content from IPFS and re-assemble the originally archived HTTP responses for replay. From a 1.0 GB sample Archive-It collection of WARCs containing 21,994 mementos, we show that extracting and indexing takes 66.6 minutes response content of WARCs containing IPFS lookup hashes takes multiple lines for readability

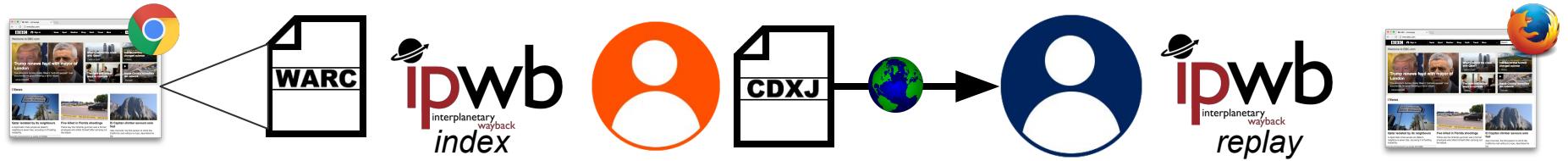
ipwb



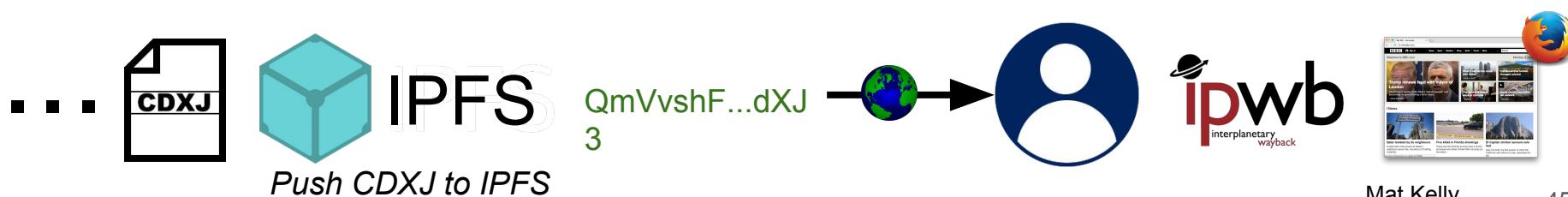
- Persistence of archived web data dependent on resilience of organization and availability of data
- Remove massive redundancy in web archive files of exact duplicate content
- Determine feasibility of pushing WARCes into IPFS

# ipwb Base Dynamics

- IPWB CDXJs may be transferred for our users' replay



- CDXJ-by-hash recursive fetch/replay
  - Share hash of CDXJ then `$ ipwb replay hash` to replicate experience





WARC Store

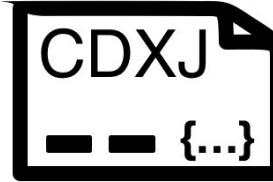


IPFS Store

Indexer

extract HTTP  
headers+payloads

1



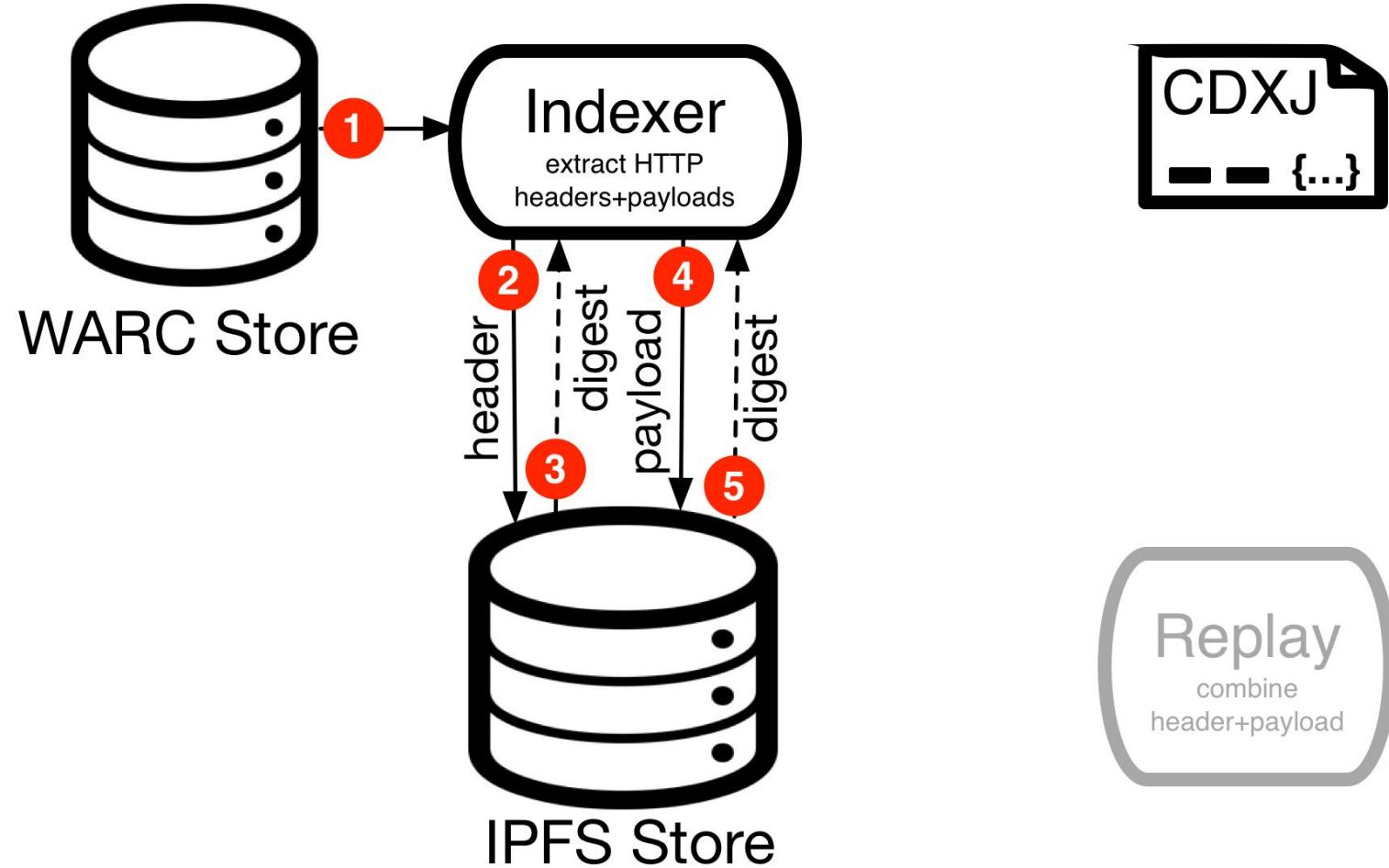
CDXJ

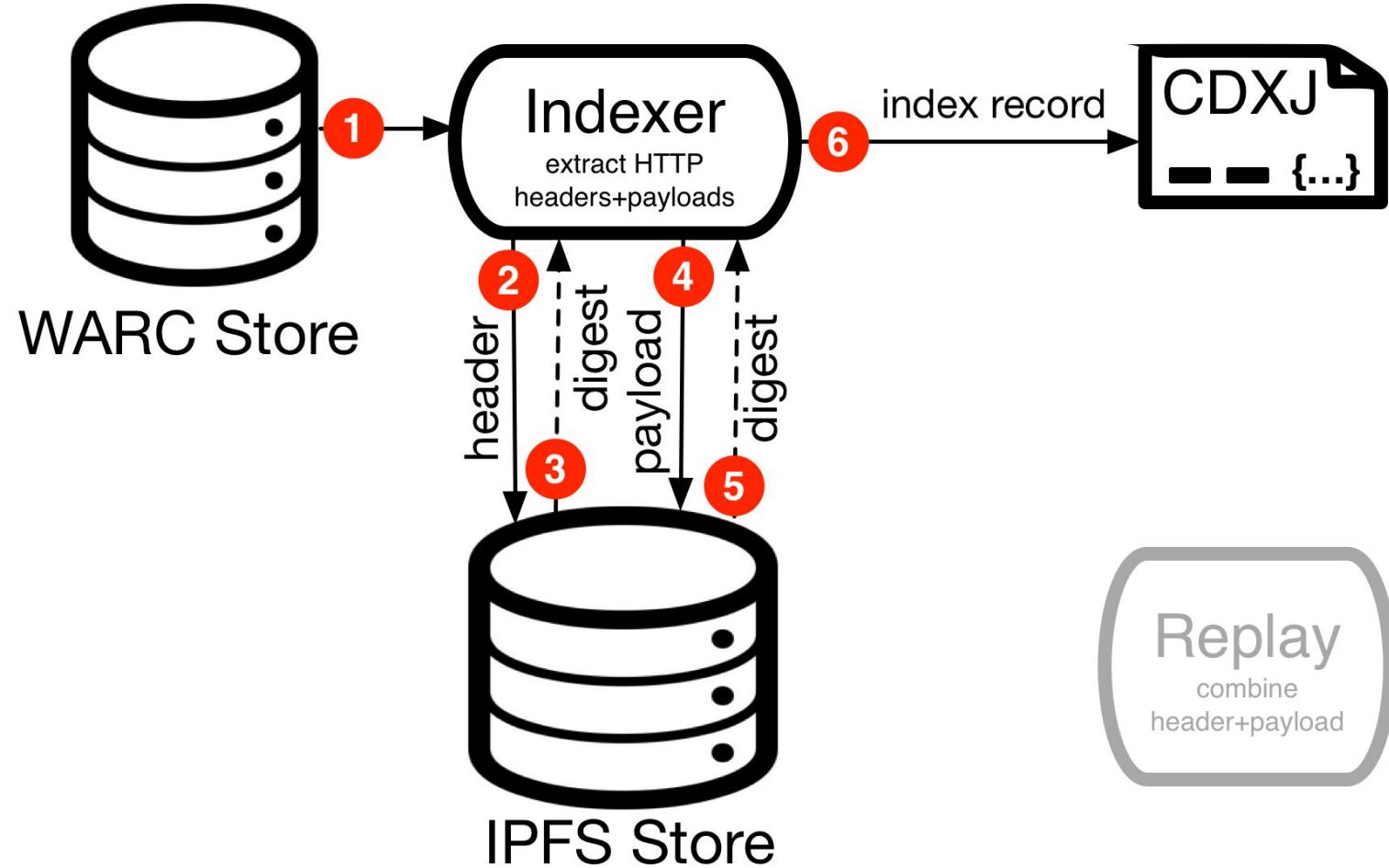
...  
---

Replay

combine  
header+payload





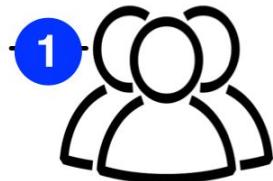
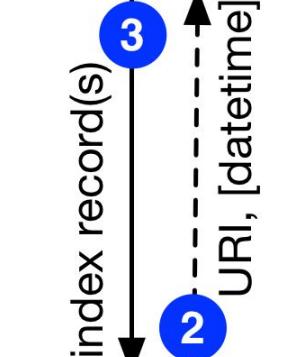
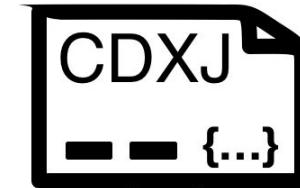




WARC Store

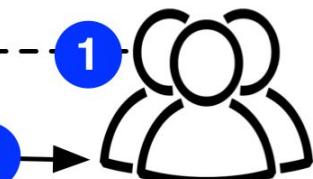
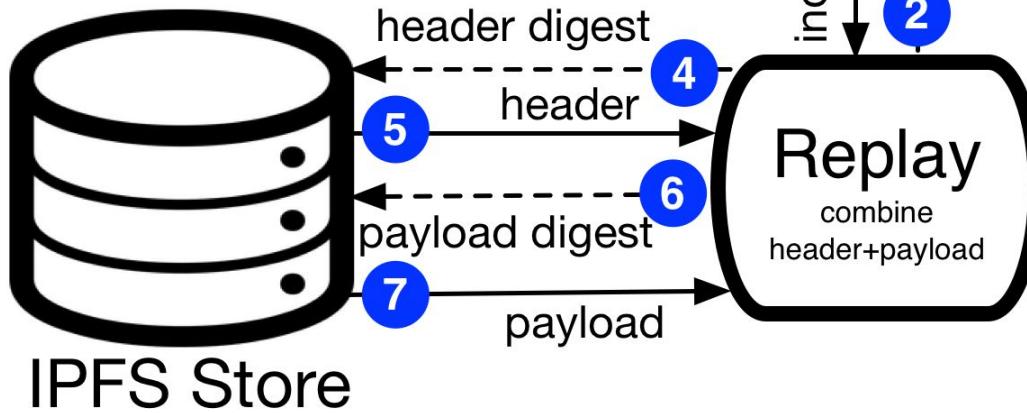


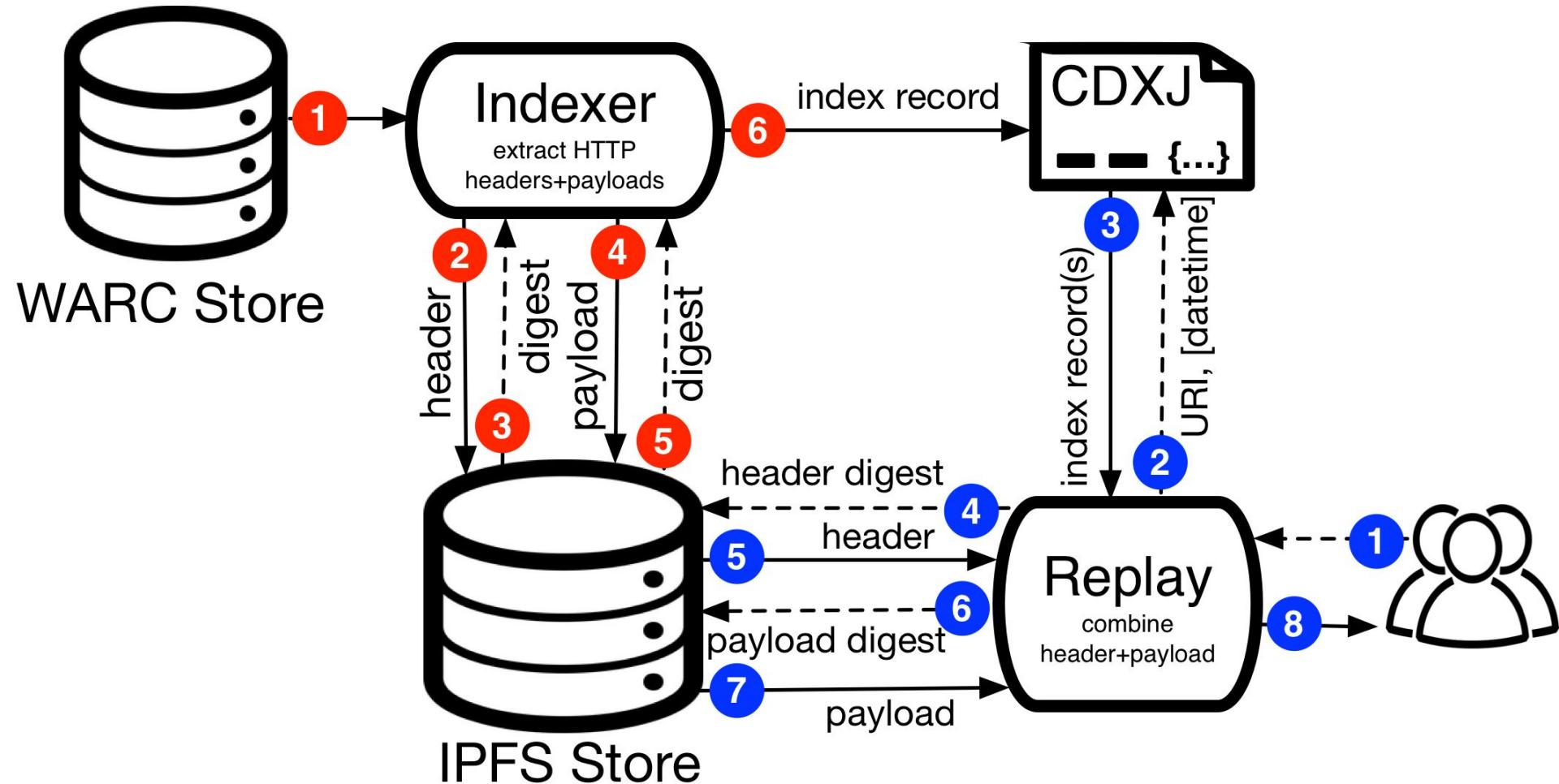
IPFS Store





WARC Store





# WARC Record extraction to CDXJ

```
20160907003819654.warc * UNREGISTERED
20160907003819654.warc
52 WARC/1.0
53 WARC-Type: response
54 WARC-Target-URI: http://ipwb.example.com/
55 WARC-Date: 2016-09-07T00:38:19Z
56 WARC-Record-ID: <urn:uuid:1e3907a9-2e5c-9981-6a92-964a465d998e>
57 Content-Type: application/http; msgtype=response
58 Content-Length: 800
59
60 HTTP/1.1 200 OK
61 Host: ipwb.example.com
62 Connection: close
63 Content-Type: text/html; charset=UTF-8
64 Content-Length: 666
65
66 <html><head>
67 <title>InterPlanetary Wayback</title>
68 <link rel="stylesheet" type="text/css" href="style.css">
69 </head>
70 <body>
71 <h1>This is site for Space Dog</h1>
72 
73 <p>InterPlanetary Wayback (ipwb) facilitates permanence and collaboration i
74
75 </body></html>
76
77 WARC/1.0
78 WARC-Type: response
79 WARC-Target-URI: http://ipwb.example.com/style.css
80 WARC-Date: 2016-09-07T00:38:19Z
81 WARC-Record-ID: <urn:uuid:2502e65a-b0bd-70c6-8799-8687029a71f4>
82 Content-Type: application/http; msgtype=response
83 Content-Length: 144
84
85 HTTP/1.1 200 OK
86 Host: ipwb.example.com
87 Connection: close
88 Content-Type: text/css; charset=UTF-8
89 Content-Length: 19
90
91
92 img {width: 250px;}
```

ipwb.example.com)/ 20160905022013 {"locator": "urn:ipfs/QmcN9eWwRF73dZj5BgT4x8jeEcFrurX1hot8QwCbMi9PB/Qmczh9YnB4H1ptPeqxcaTZA4aMmuNUswTLTwzXntvbp9sL", "mime\_type": "text/html", "status\_code": "200"}  
ipwb.example.com)/style.css 20160905022013 {"locator": "urn:ipfs/QmU1k71bT6ibZBSdxBL35cQXwoVtih8cTB4CXfrjyMfZxE/QmbvUAo9U31wSdvARjvbPeVBTAwCjN1kyPhQ4ho3n8TAZo", "mime\_type": "text/css", "status\_code": "200"}  
ipwb.example.com)/ipwb.png 20160905022013 {"locator": "urn:ipfs/QmTjfMixFGvbP4nwFoq3tNYDPW6gC99i5njrqsXSw6QRvHa/QmYMKZbnk53kuPJirahJHGevCCy2afLyePRdX38TukFUwd", "mime\_type": "image/png", "status\_code": "200"}  
ipwb.example.com)/fileduration.png 20160905022013 {"locator": "urn:ipfs/QmaCj6LNngxwqxaLmfplxCyxcwDt2Uzqf8gCG6bVyQppYC/QmdgtMcGprTF8bqv7ytgMwtoi5BhRxfuvBjD6Vj2U7ohz1", "mime\_type": "image/png", "status\_code": "200"}  
ipwb.example.com)/filesize.png 20160905022013 {"locator": "urn:ipfs/QmNPjrSVY31oGDooMiA18ZDNHfkLnEg3j5gRj1dFdrqmS4/Qmb4heB8PU58nkWt6w5tBgmfpelTKuU7iuxg9tFdoPsF1B", "mime\_type": "image/png", "status\_code": "200"}}

# IPFS multihashes in IPWB CDXJs

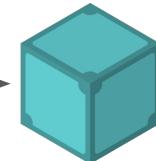
```
com,example)/index.html 20170301192639 {"locator":  
"urn:ipfs/QmPdyY6Pm66iWtGpTc7PqK11hvsnYSKMVL57G69RiNjGcm/QmNZ6mKS  
SAXAmXEocQj5gT4y4kdcr5D2C173ubWJ6PSKEZ", "mime_type": "text/html",  
"status_code": "200"}  
com,example)/images/frog.png 20170301192639 {"locator":  
"urn:ipfs/QmUeko8zM7Xanwz6F9GtRH4rLAI4Poj3EMECGsci3BRQfs/QmPhMnX74c  
wqx2xgj9d3N3gTra8CzafXwSbUwU8xagMfqR", "mime_type": "image/png"  
"status_code": "200"}
```



# Content Addressing



<http://foo.com/spaceDog.jpg>



**IPFS**

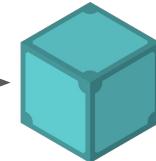


QmZAD4xeeNeYF3TmwWgBXypLKTiCGwGRMXHW7MtheWKtw4

====



<http://example.org/yuri.jpg>



**IPFS**



QmZAD4xeeNeYF3TmwWgBXypLKTiCGwGRMXHW7MtheWKtw4

\$ ipfs cat QmZAD4xeeNeYF3TmwWgBXypLKTiCGwGRMXHW7MtheWKtw4 > doge.jpg



# Methodology - IPWB WARC indexing

- warc-response record body extracted into temp files
  - HTTP header and entity body (payload) separated
  - Response metadata (e.g., datetime) retained
- temp files pushed into IPFS via locally running daemon
  - Two IPFS hashes (for header and payload) returned
- CDXJ record created representing warc-response contents
  - Contains URI-R, archived HTTP status, encoded IPFS hashes

# Methodology - Replaying Archives

- Extension of pywb API to read CDXJ files
- On encountering IPFS URN, fetch `warc-response temp` files from IPFS using local daemon
  - This may occur on a separate machine using a separate daemon
- With WARC contents fetched, replay contents using pywb where the locator value in the CDXJ is used to dereference the temp files pulled from IPFS

# CDXJ in ipwb

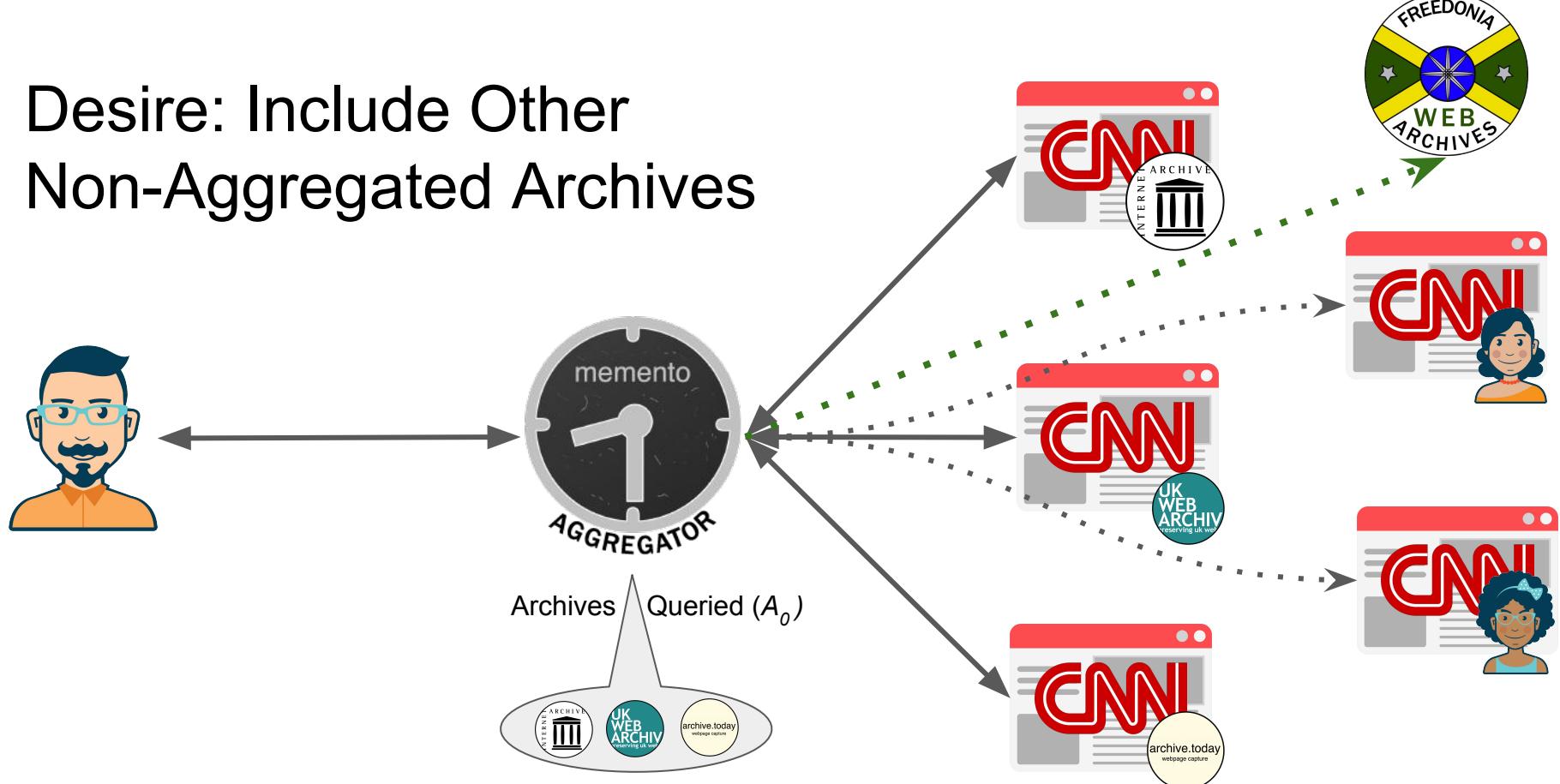
```
1 SURT_URI DATETIME {  
2     "id": "WARC-Record-ID",  
3     "url": "ORIGINAL_URI",  
4     "status": "3-DIGIT_HTTP_STATUS",  
5     "mime": "Content-Type",  
6     "locator": "urn:ipfs/HEADER_DIGEST/PAYLOAD_DIGEST"  
7 }
```

# A Framework for Aggregating Private and Public Web Archives

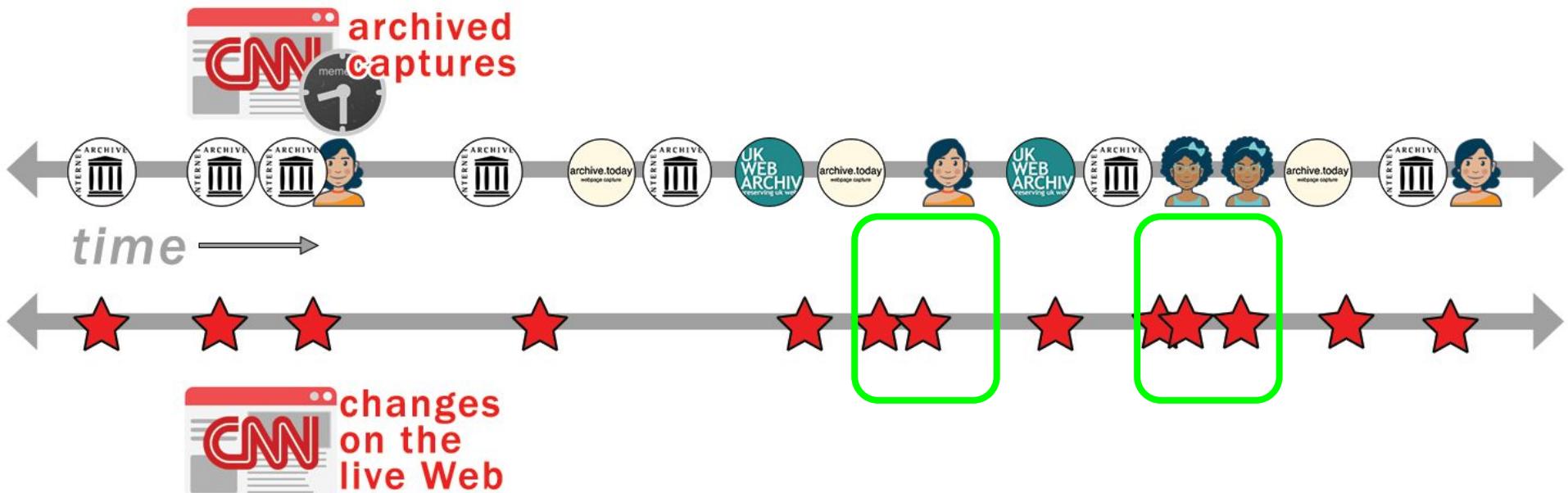
- Mitigate outstanding issues in Web archiving beyond public scope.
  - High-level of dissertation topic
  - Introduced the “Memento Framework”

# A Framework for Aggregating Private and Public Web Archives

# Desire: Include Other Non-Aggregated Archives

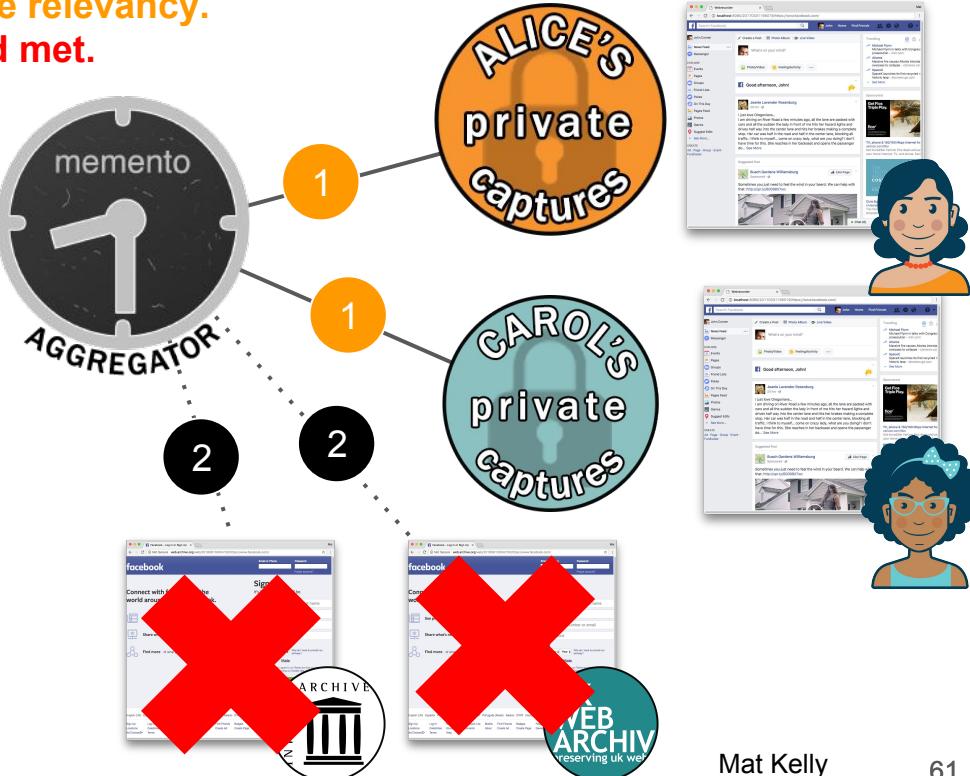
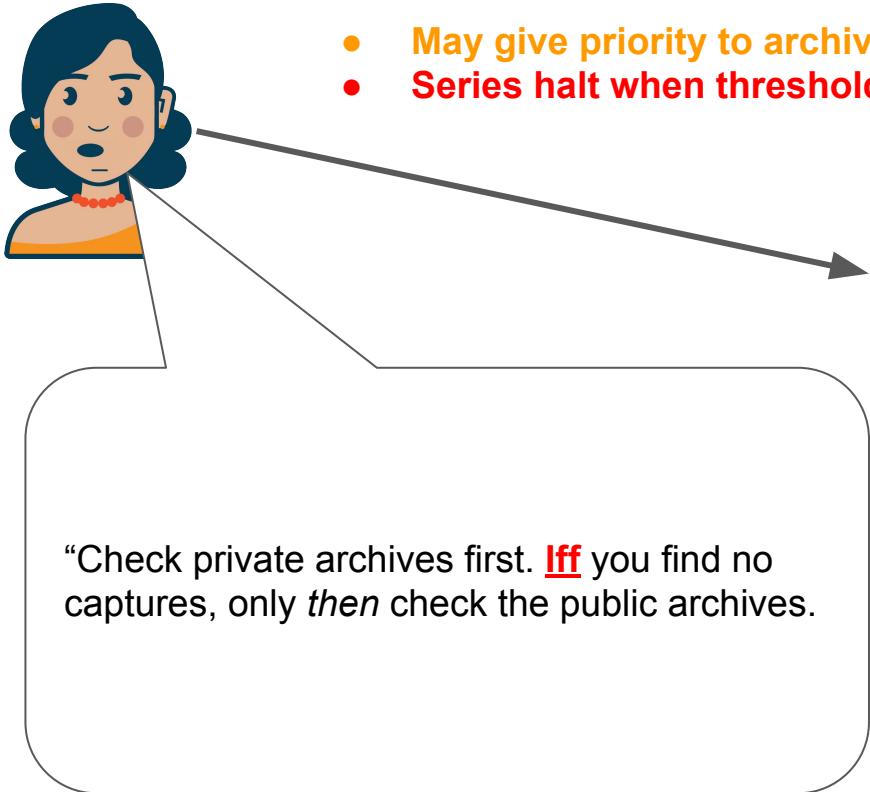


# Archiving More Archives Provides a Better Picture of the Web

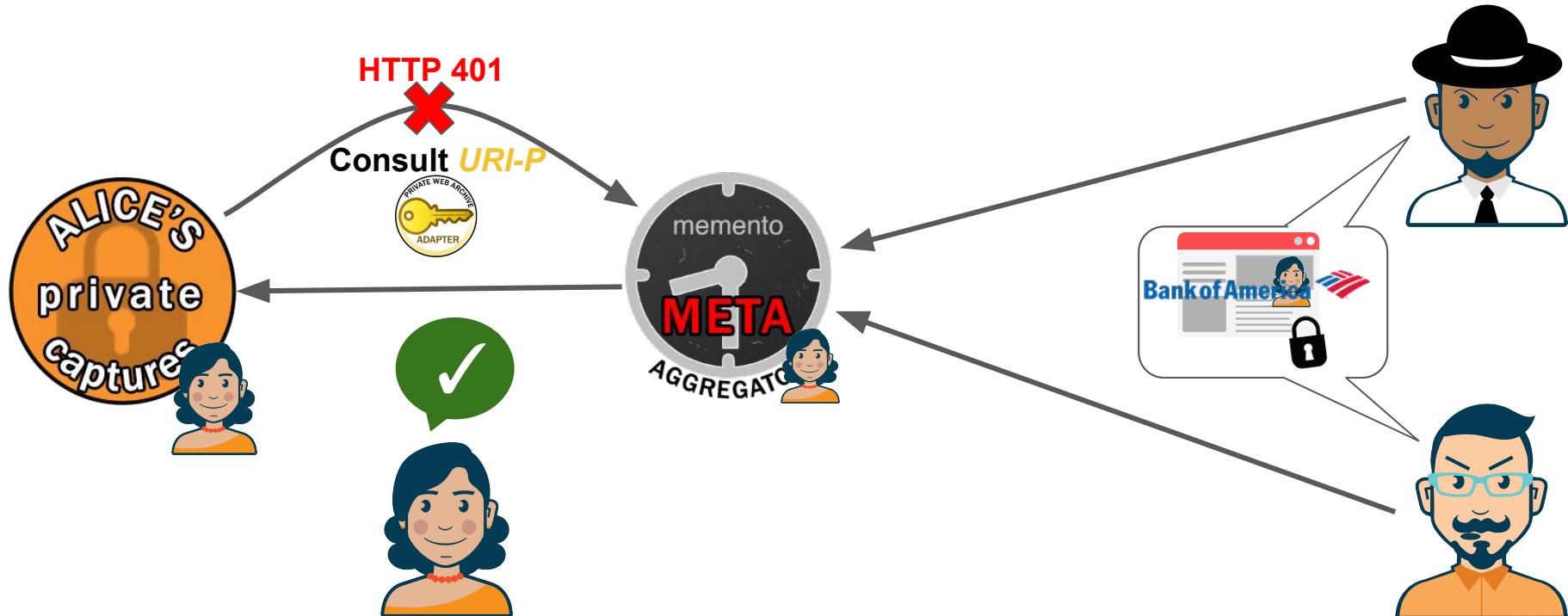


# Query Precedence & Short-Circuiting

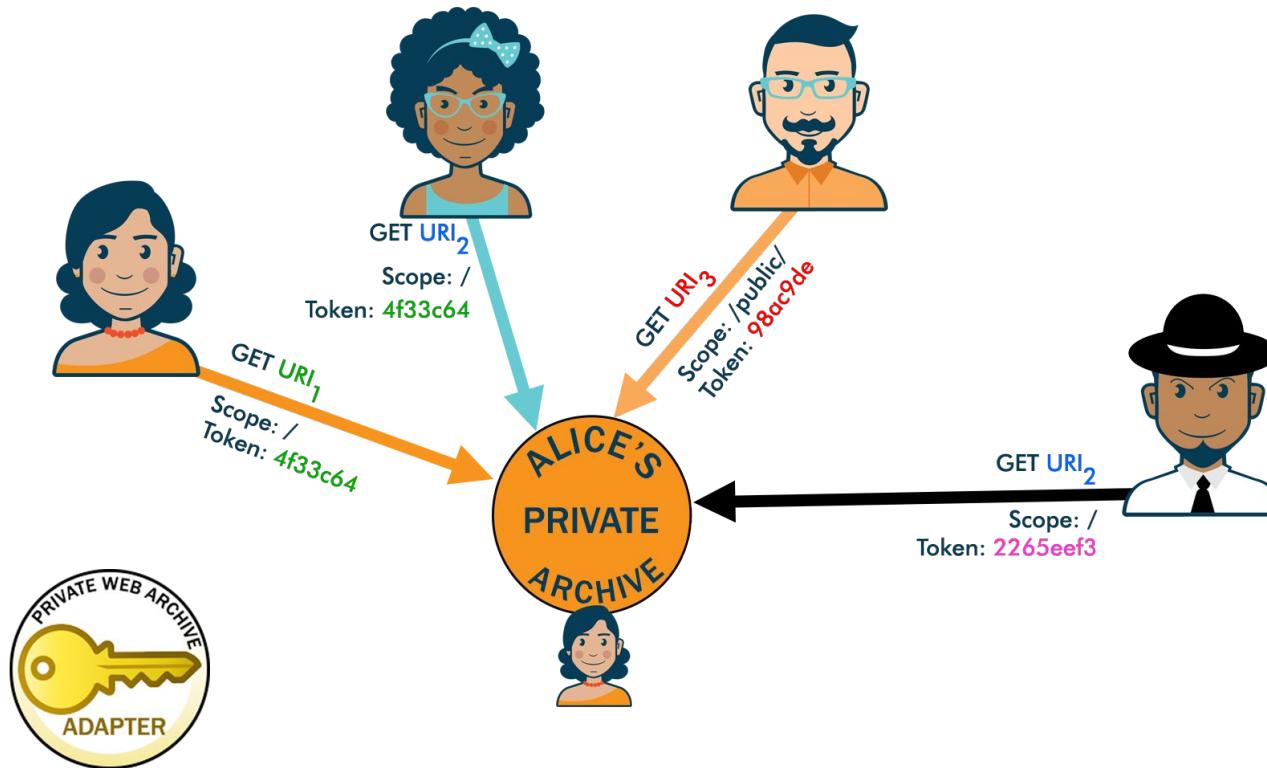
- May give priority to archive relevancy.
- Series halt when threshold met.



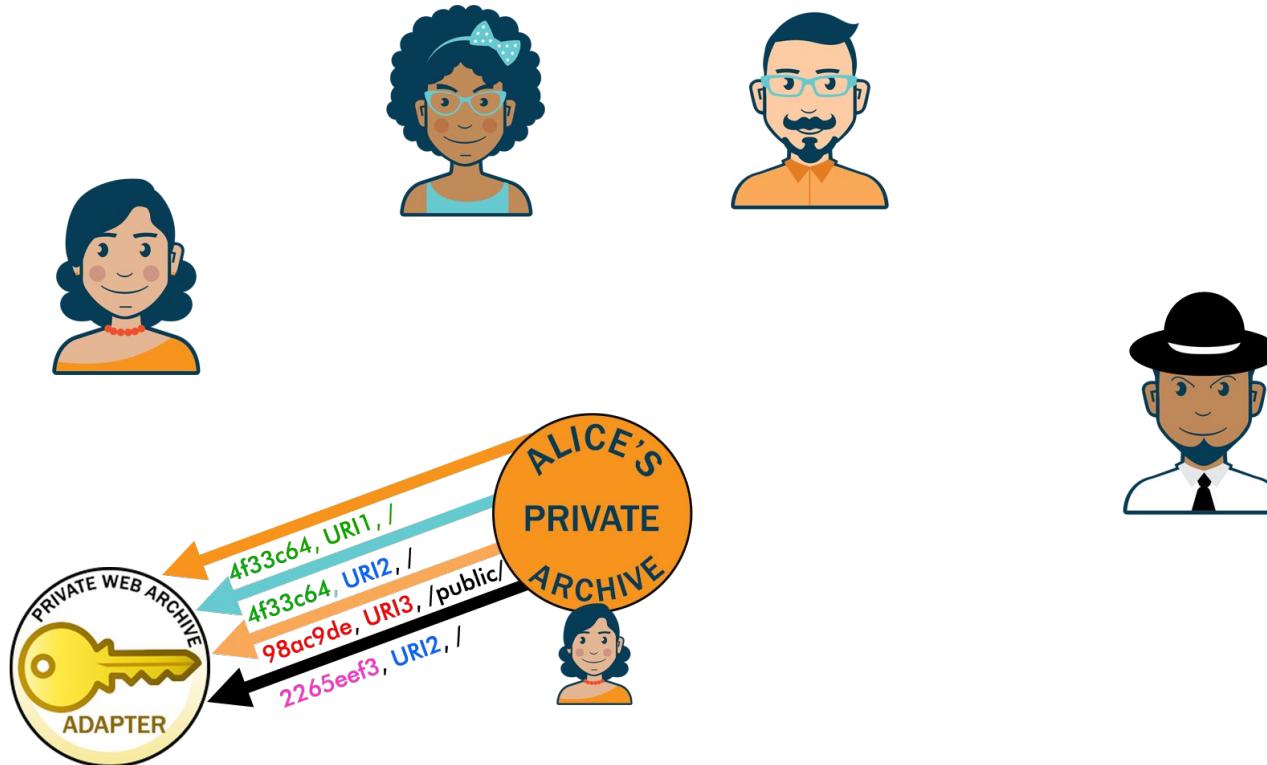
# Aggregation with Access Regulation



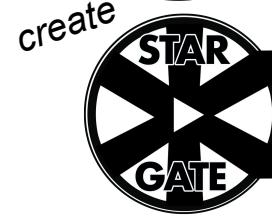
# OAuth2-based tokenization patterns



...with offloading of the procedure from the archives



# Evaluation Through Implementation



Extend for client-side archival specification

Exhibit features of an MMA

Regulate access to Private Web archives

Facilitate archival negotiation in more dimensions

The screenshot shows a web application interface for managing mementos. At the top, there's a navigation bar with three circles, a left arrow, a right arrow, and a URL field containing <https://my.fancydomain.com/somePage.html>. To the right of the URL is a user icon with the number '1' and a vertical ellipsis menu.

The main content area displays a list titled "569 Mementos Available". It includes a logo for "MINK INTEGRATING the WEBS". Below the title are filter options: "VIEW BY: Dropdown" (selected), "Columns" (selected), "VizMethodFoo", and "VizMethodBar".

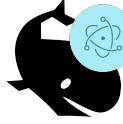
The data table has four columns:

| 2005 | 5   | January  | 9         |
|------|-----|----------|-----------|
| 2006 | 6   | February | 22        |
| 2007 | 36  | April    | 33        |
| 2008 | 42  | November | 15        |
| 2009 | 57  |          | 1st       |
| 2010 | 3   |          | 2nd       |
| 2011 | 2   |          | 9th       |
| 2012 | 0   |          | 22        |
| 2013 | 79  |          | 18th      |
| 2014 | 81  |          | 30th      |
| 2015 | 99  |          | 1         |
| 2016 | 156 |          | 5         |
| 2017 | 3   |          | 09:30 GMT |
|      |     |          | 11:06 GMT |
|      |     |          | 12:58 GMT |
|      |     |          | 20:06 GMT |
|      |     |          | 20:08 GMT |

To the right of the table is a "Sources" sidebar with checkboxes for "Internet Archive" (checked), "Archive.is" (unchecked), "Local Archive1" (unchecked), "Local MMA1" (checked), and "Remote MA1" (unchecked). There are also edit and delete icons for the sources.

# On Tool Building and Evaluation of the Archived Web

## Open Source Web Archiving Tools



## GitHub

machawk1/warcreate machawk1/wail

n0tan3rd/wail

machawk1/mink

oduwsdl/ipwb

More details of studies:

- Measuring the Impact of Missing Resources ([Conf](#), [Journal](#))
- Impact of URI Canonicalization ([Conf](#), [arXiv](#))
- Archivability Over Time ([Conf](#), [arXiv](#))
- Impact of JS on Archivability ([Journal](#))
- Personalization in Web Archives ([article](#))

Mat Kelly

<https://www.cs.odu.edu/~mkelly/>  
February 13, 2019