

A Framework for Aggregating Private and Public Web Archives

Mat Kelly

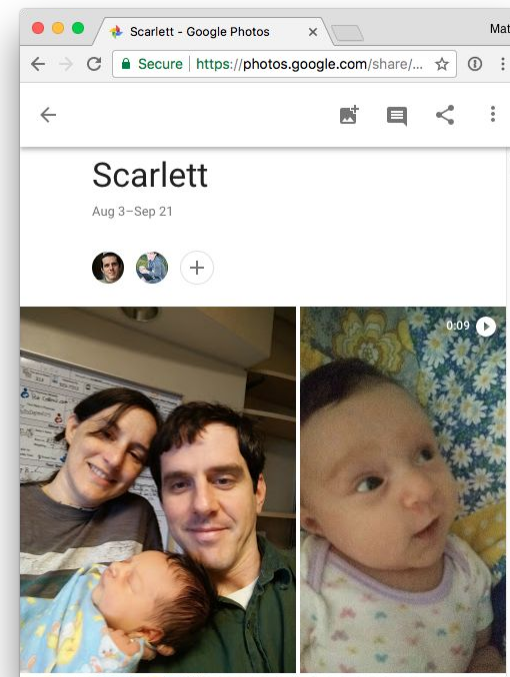
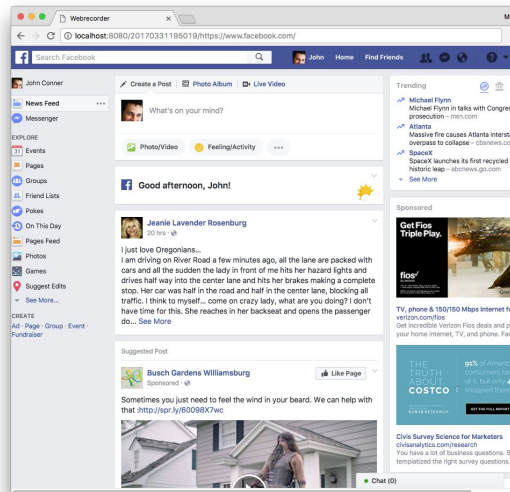
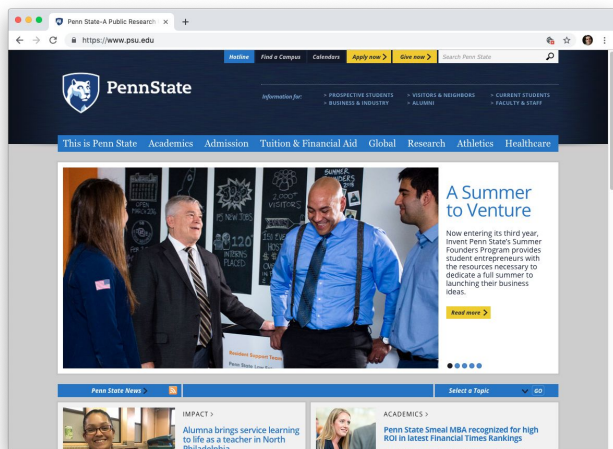
Old Dominion University
Web Science & Digital Libraries Research Group
Department of Computer Science
Norfolk, Virginia USA
`mkelly@cs.odu.edu`



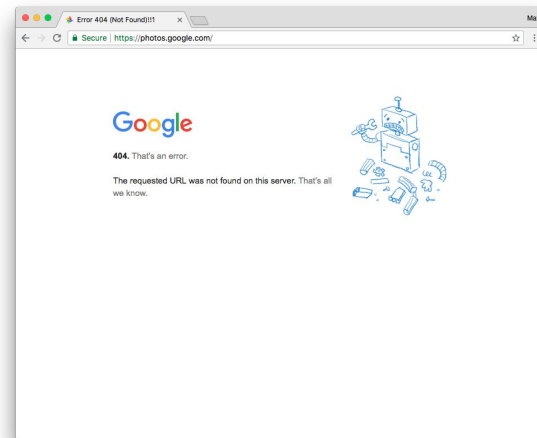
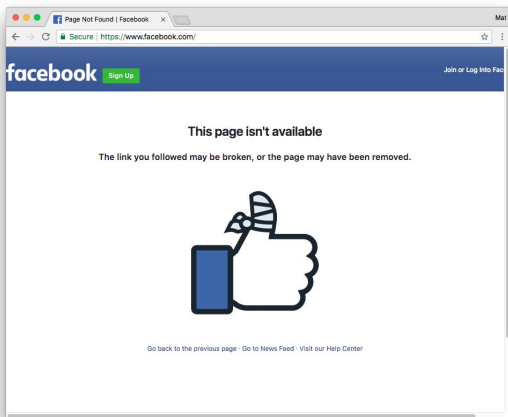
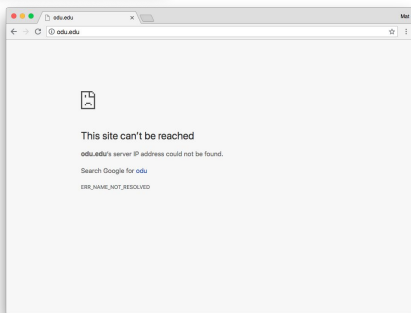
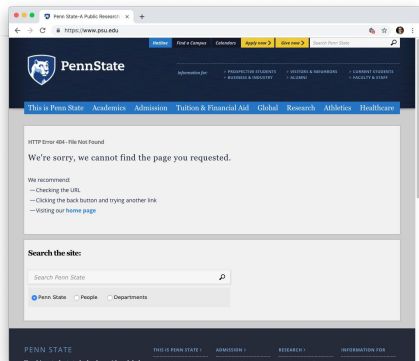
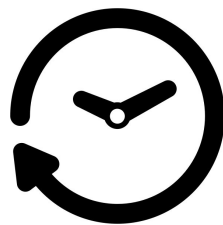
Seminar, Penn State University
February 14, 2019



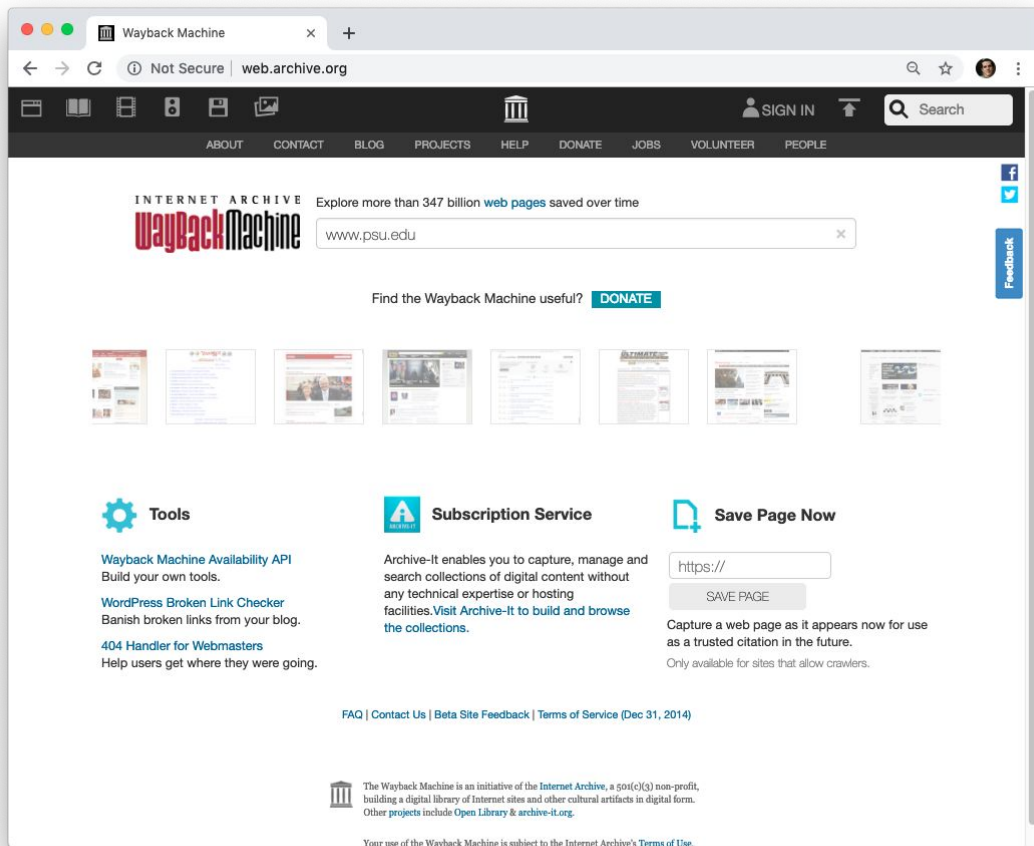
The Web



The Web is Ephemeral

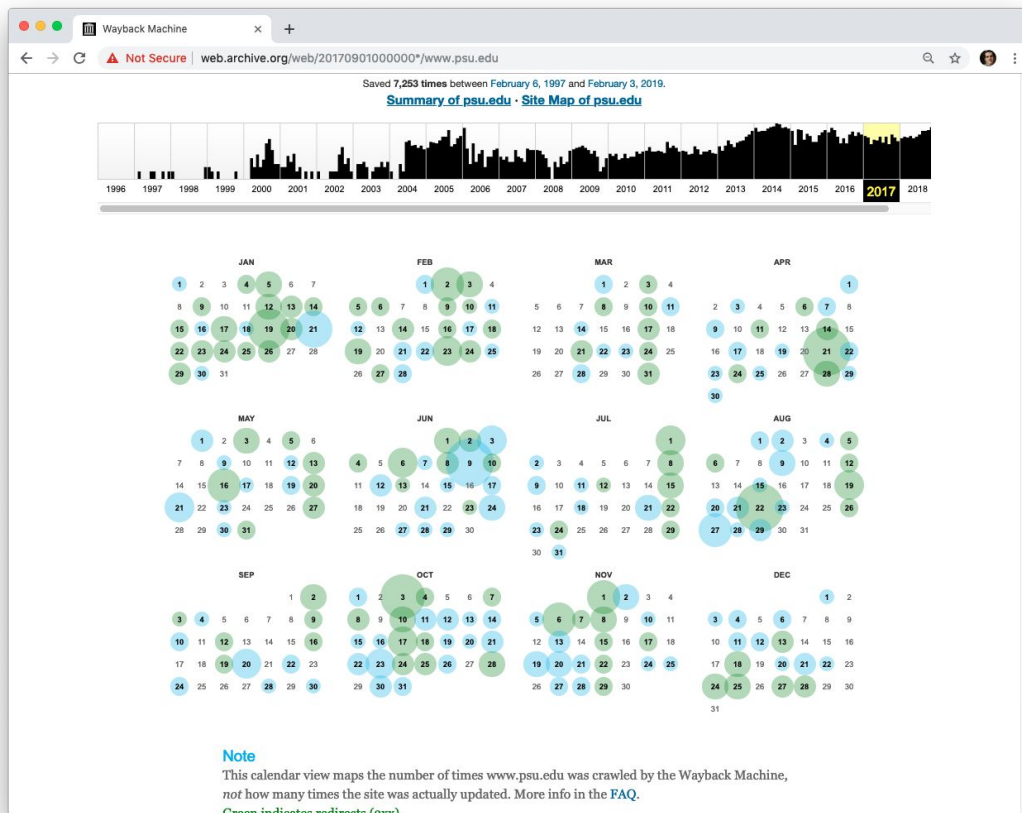


Web Archives to the Rescue: Typical Access



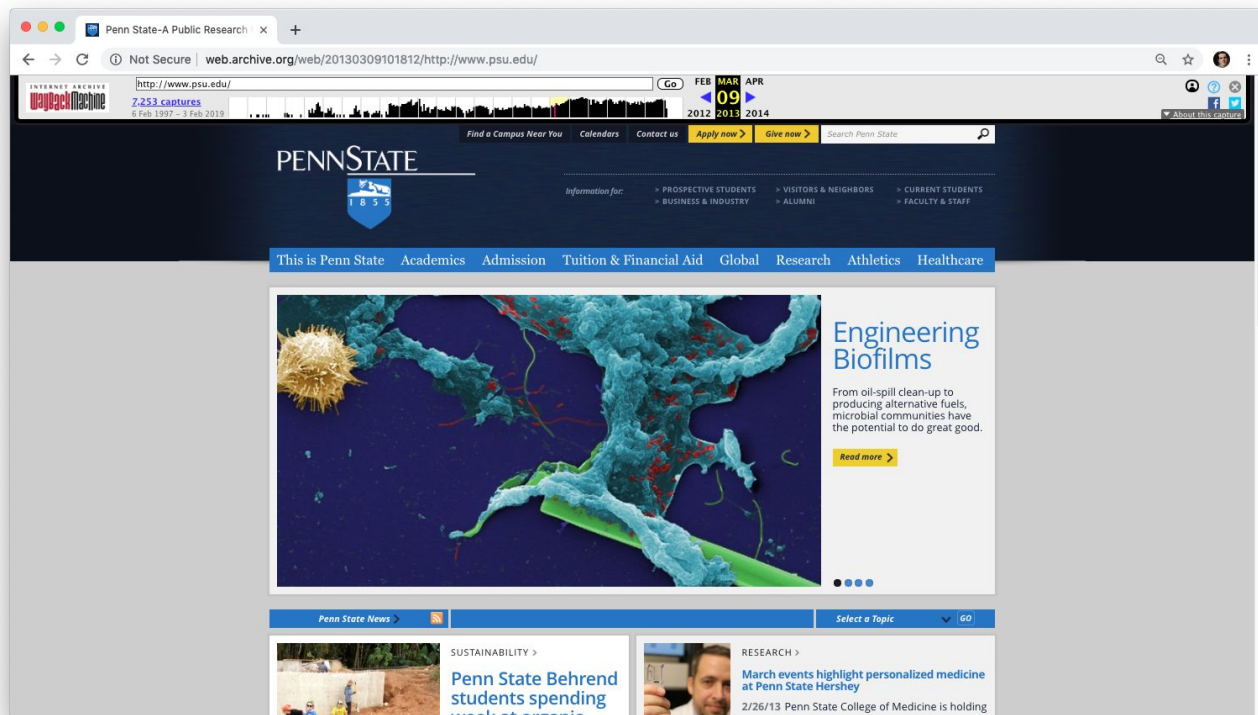
1. Go to `archive.org` in your browser
2. Enter the URL you want to see in the past in the form field
3. Submit your query

Web Archives to the Rescue: Typical Access

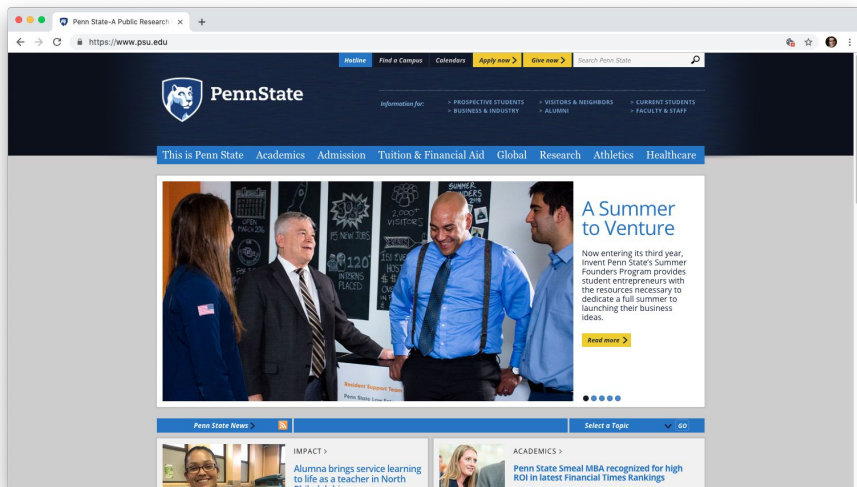


4. Locate the capture on the calendar or histogram view
5. Select the year/capture for the day
6. Repeat until you find the closest date and time

7. Finally, view the capture

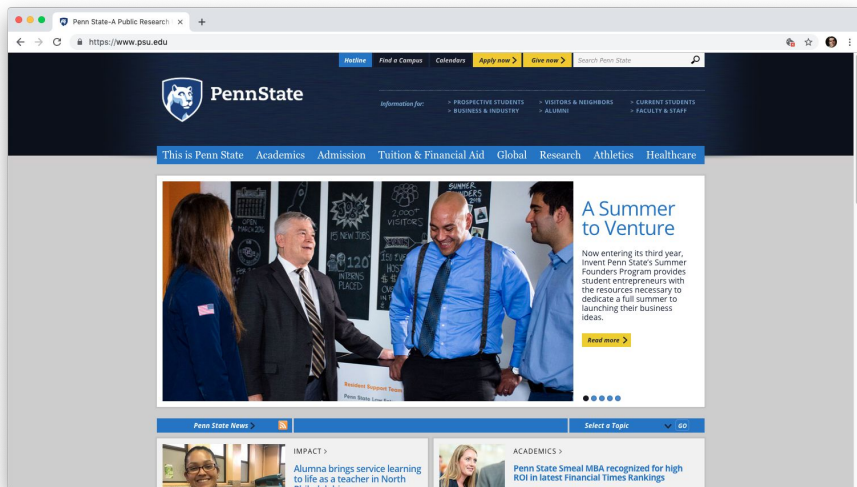


Web Archiving - Live Web psu.edu

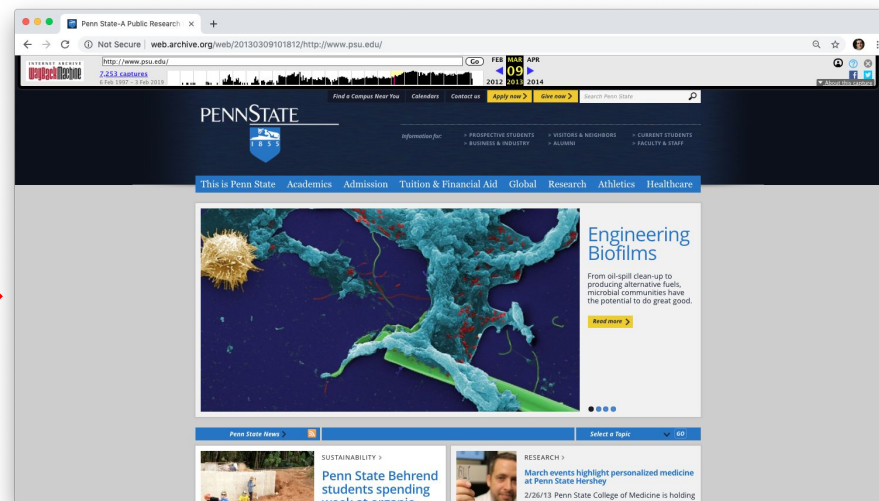


Now

Web Archiving - Archival Capture



Now



March 9, 2013

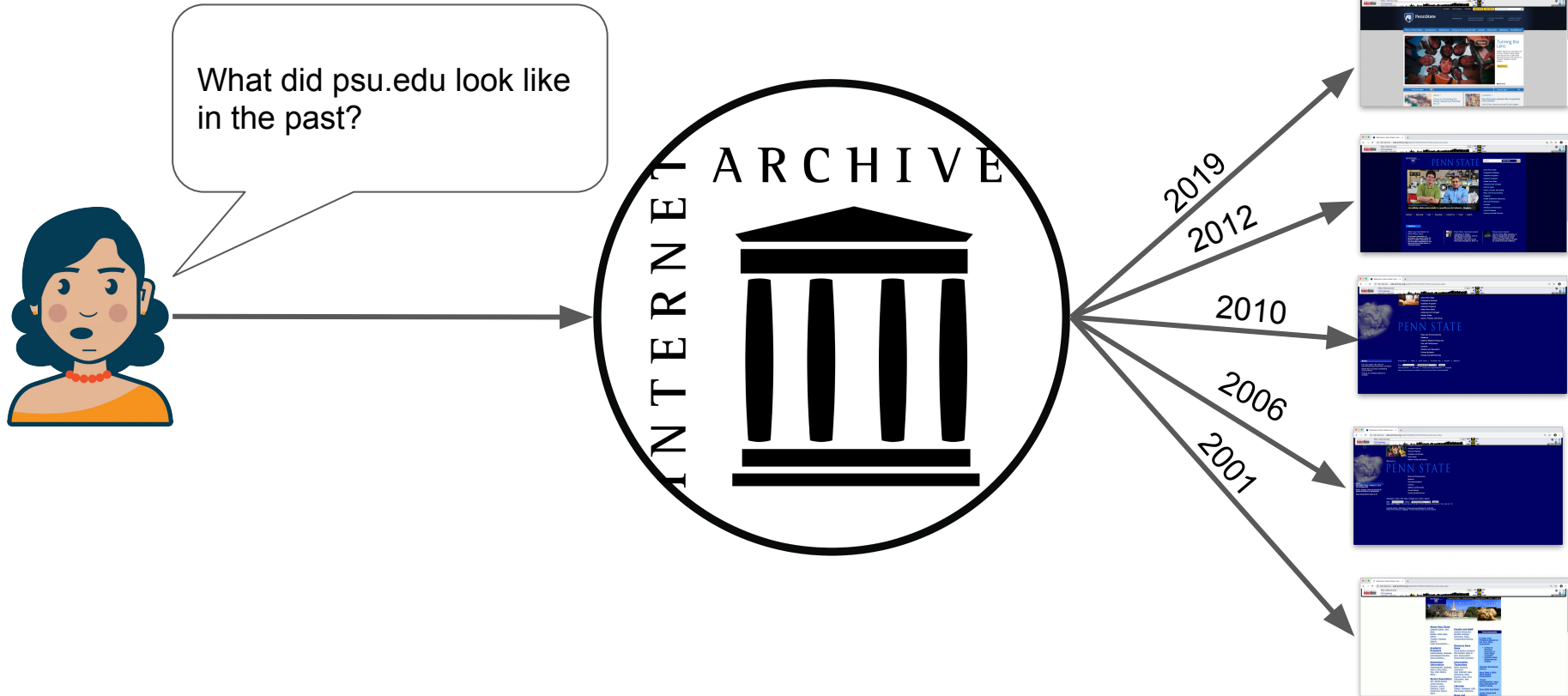
Web Archiving

Associate live Web URIs



With their archived representations

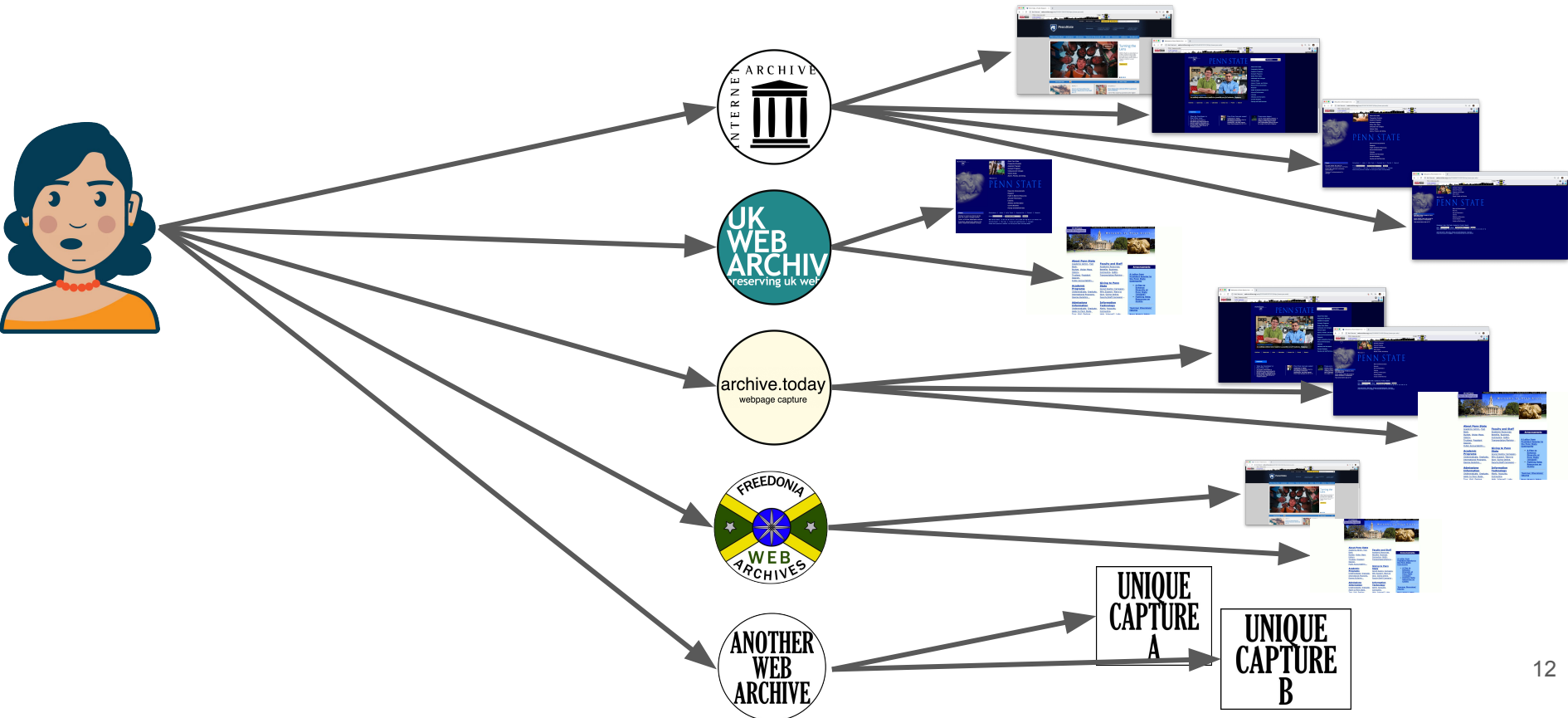
Web Archives provides access to the **Web** that *was*



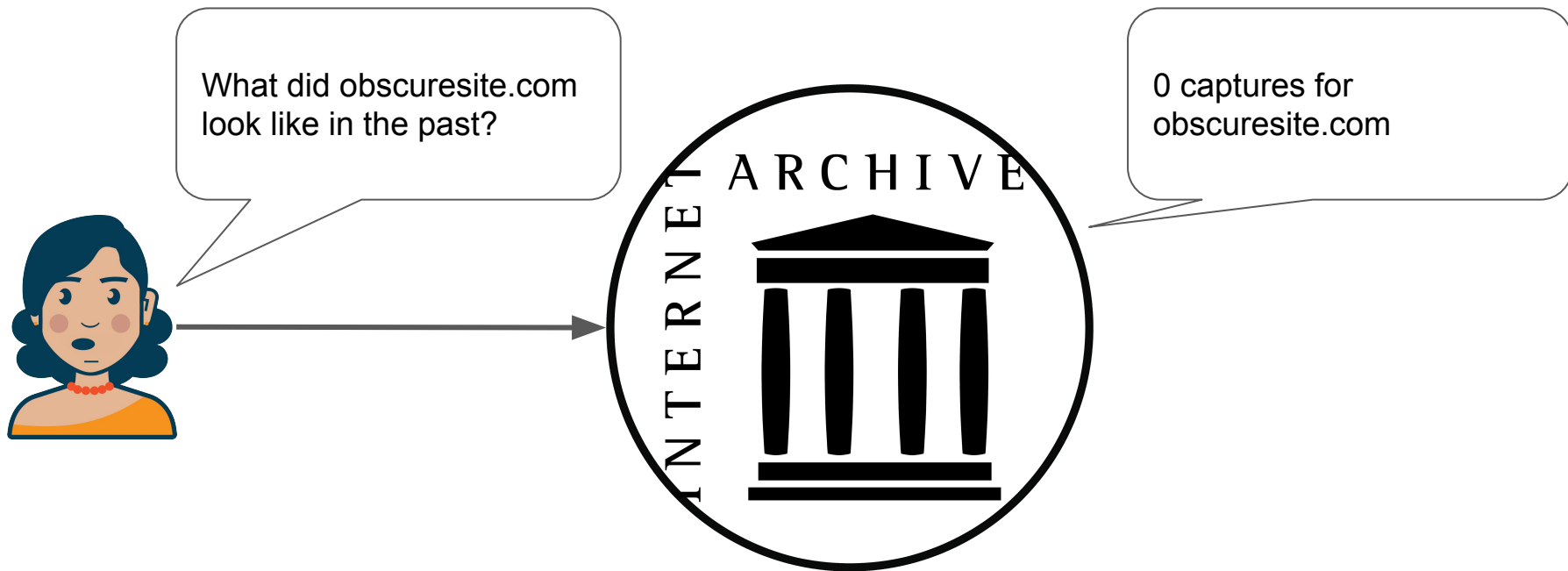
Multiple archival efforts (3 of many)



More archives produces a more comprehensive picture



Even then, not everything is preserved



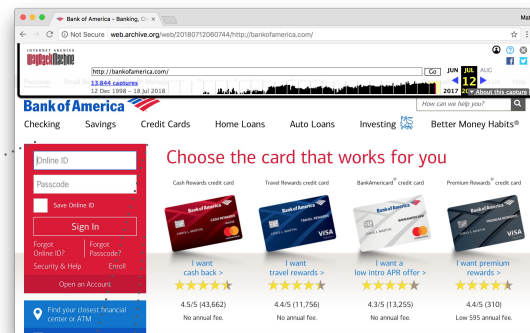
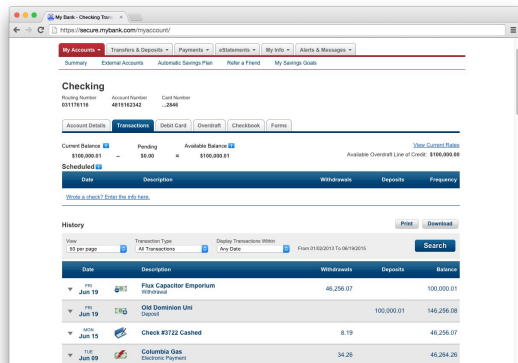
User sees on live Web may not be what is captured

What did facebook.com look like in the past?

The diagram illustrates the concept of web archiving. A woman's head is shown with a speech bubble asking "What did facebook.com look like in the past?". An arrow points from her head to a circular logo with a classical building facade and the text "INTERNET ARCHIVE". To the left is a screenshot of the Facebook homepage as seen through a Webrecorder browser window. To the right is a screenshot of the Facebook "Sign Up" page as seen through a web.archive.org browser window.

...And oftentimes that is for the best

Have you preserved my online banking?

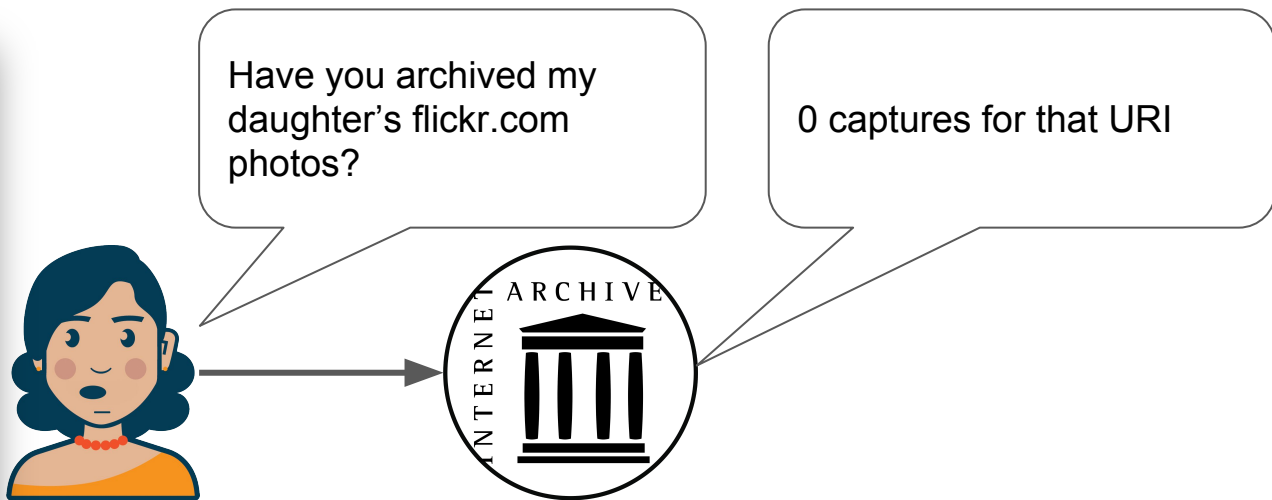
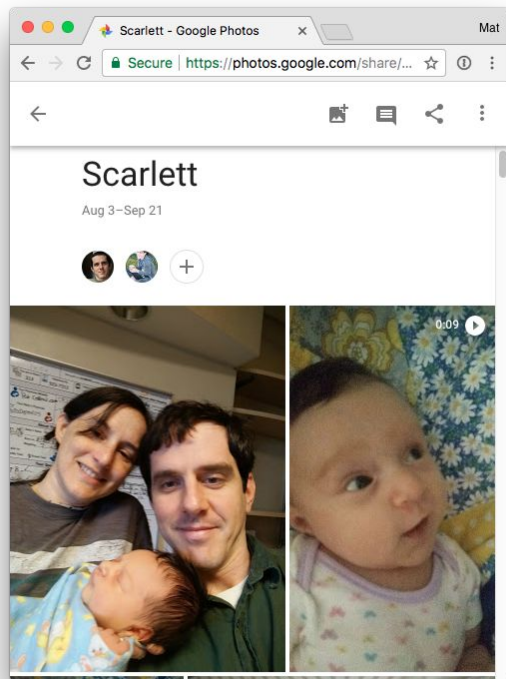


Online ID

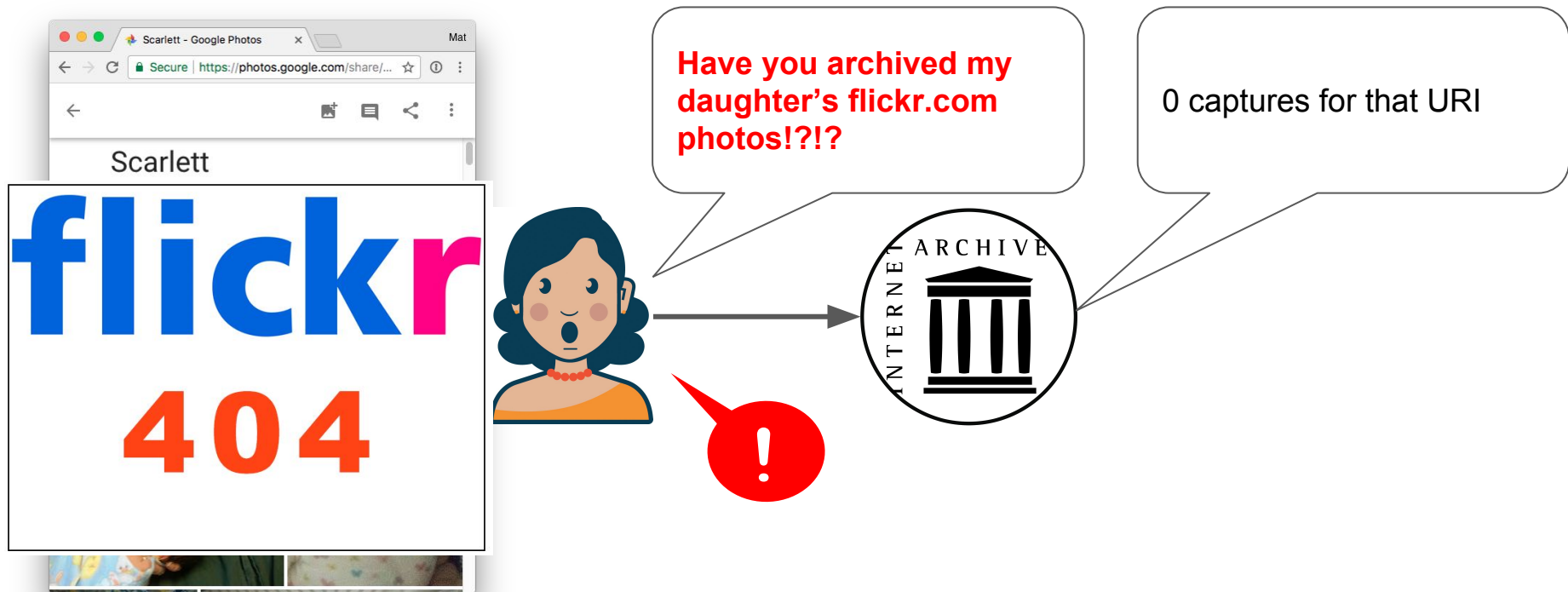
Passcode



Other times, we may want our content archived



...especially when it has disappeared

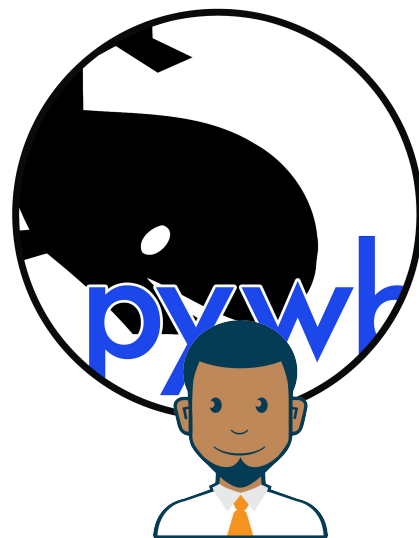


“Save this, but only for me.”

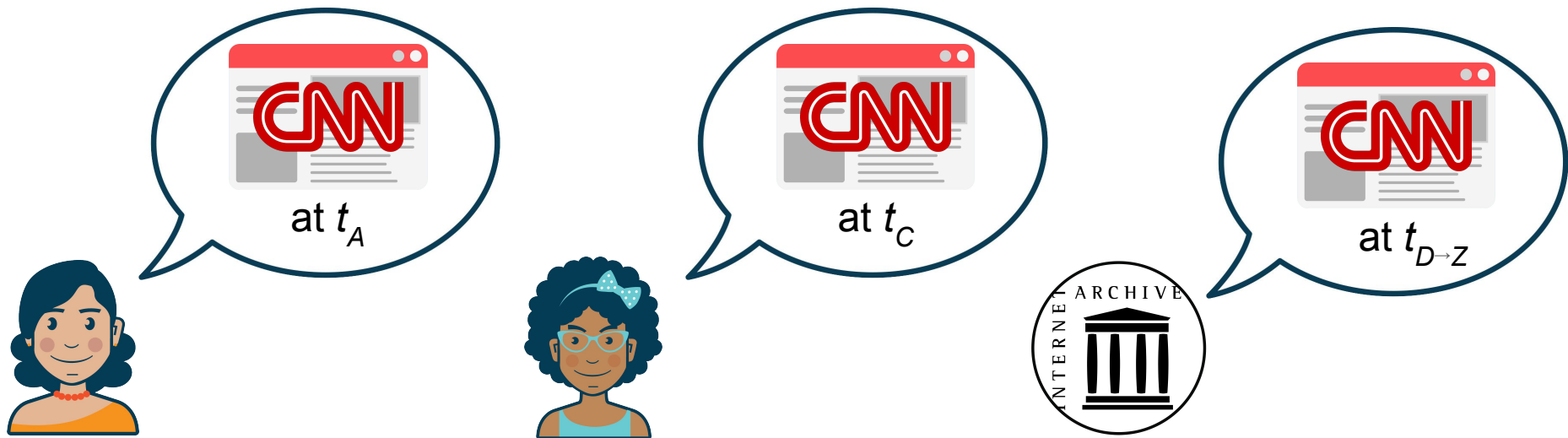
- **Screenshots** of Web pages are insufficient
 - Not interactive/representative, do not integrate, lose context otherwise provided in metadata
- Large-scale archives' tools are open source
- Individuals can archive, but there are still technical barriers



Individuals, Too, Can Archive The Web



Captures from Institutional and Personal Sources

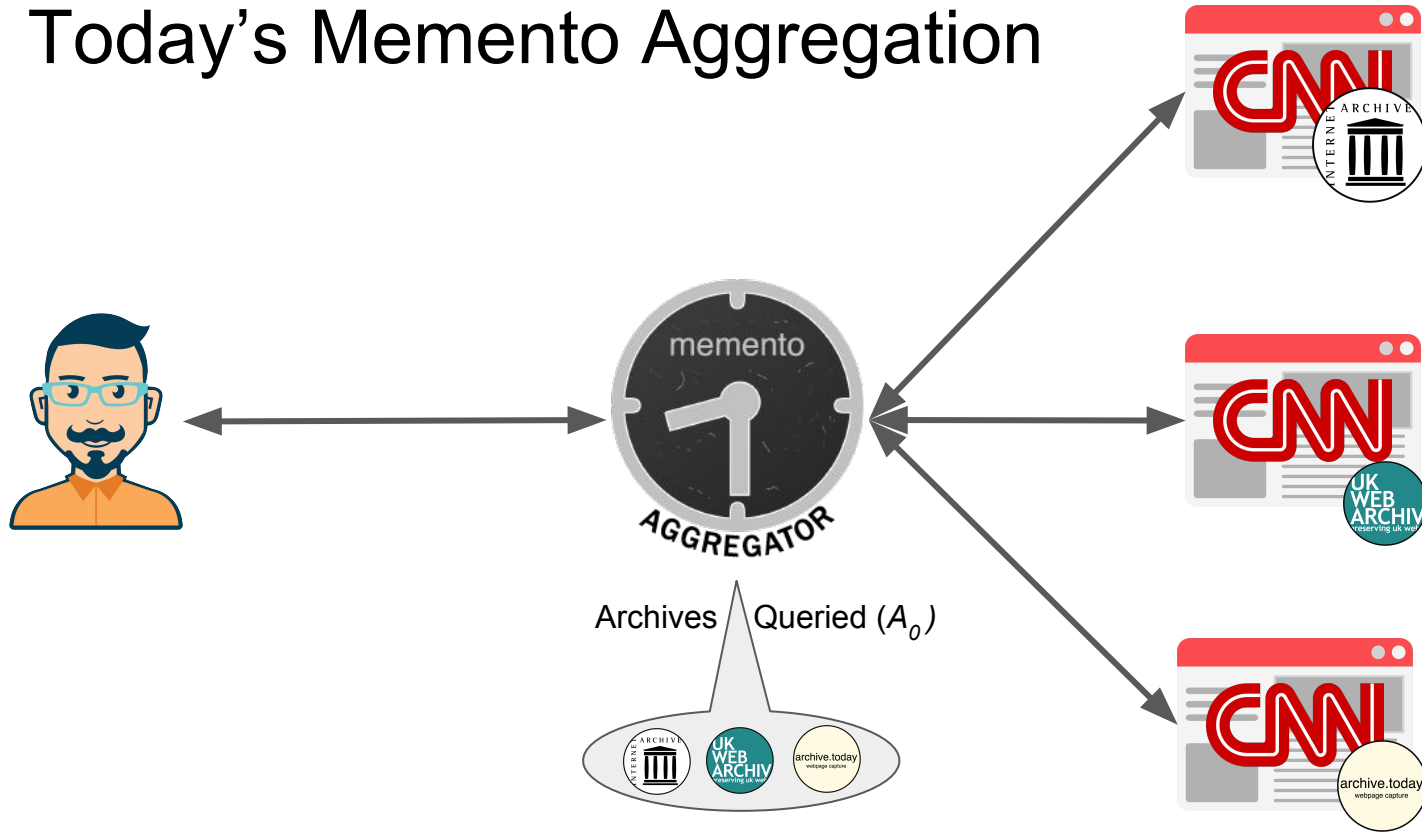


Memento Facilitates this Aggregation

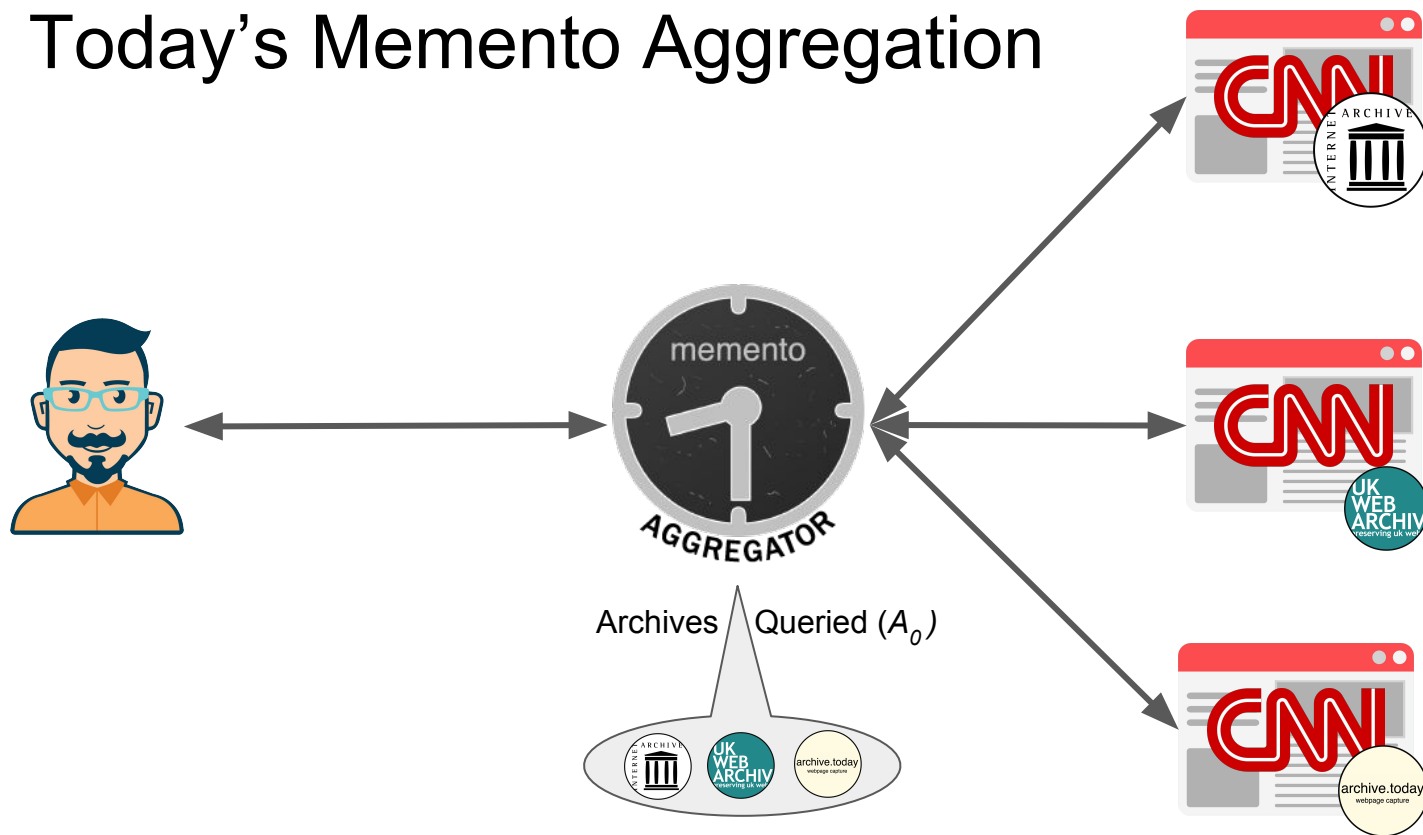
RFC7089



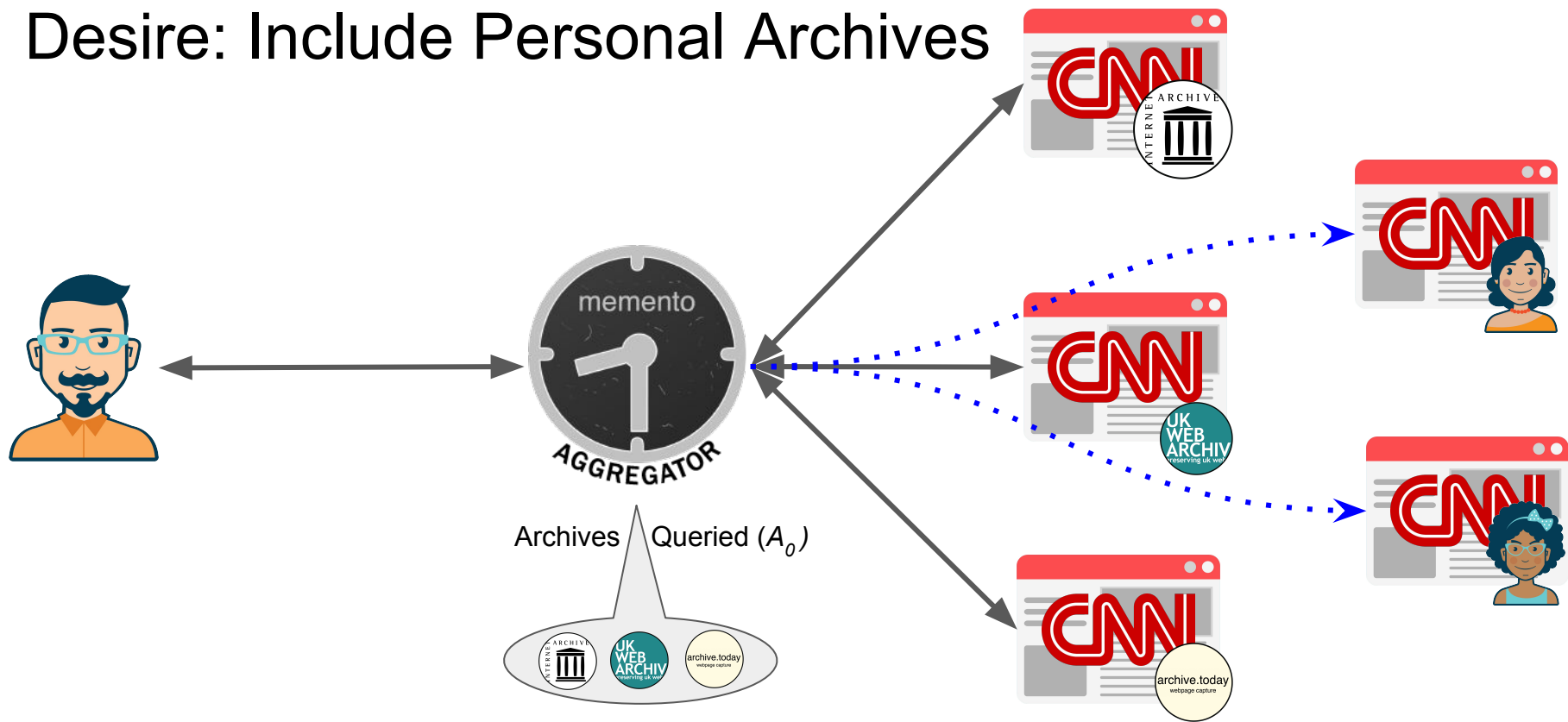
Today's Memento Aggregation



Today's Memento Aggregation



Desire: Include Personal Archives

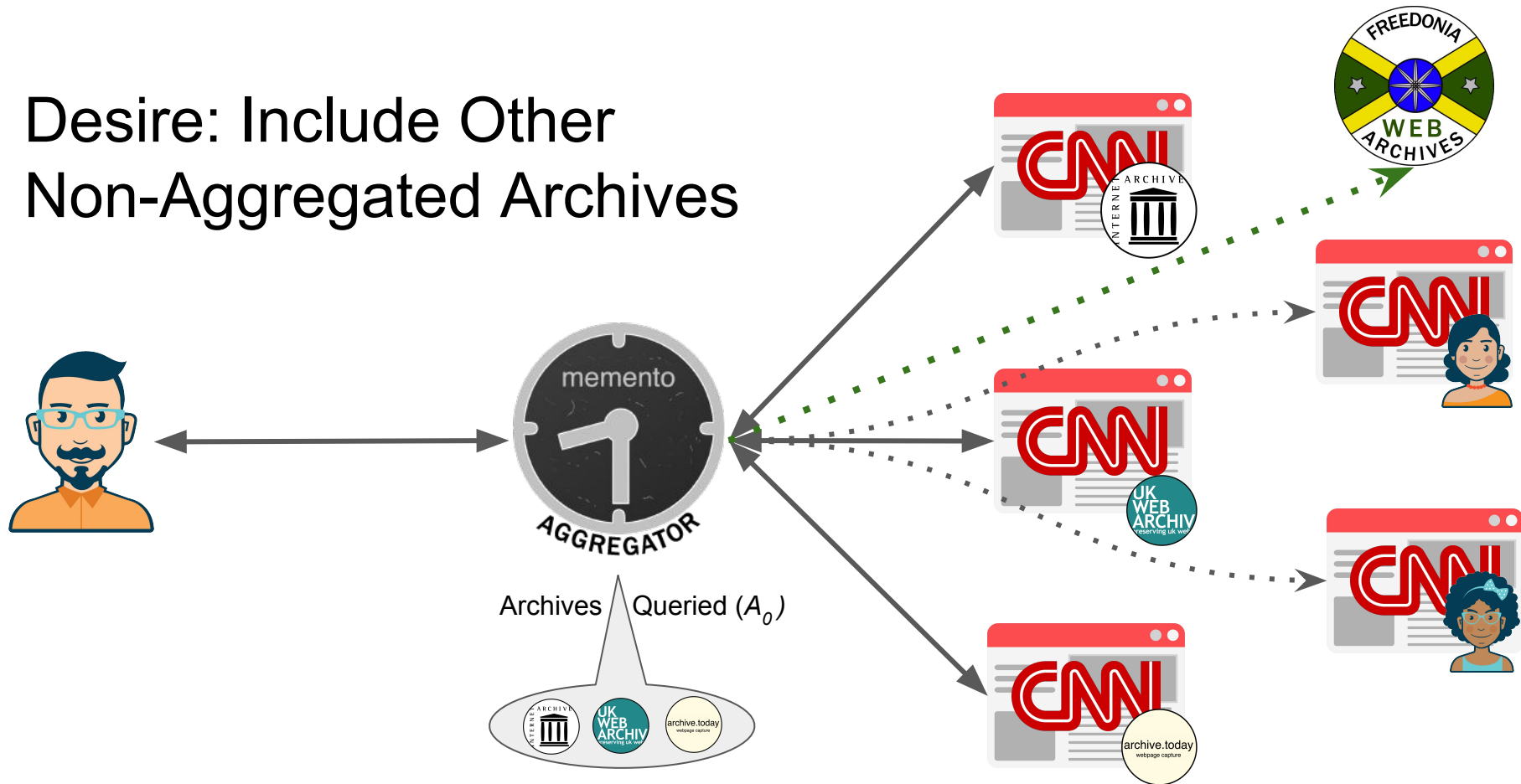


A Framework for Aggregating Public and Private Web Archives

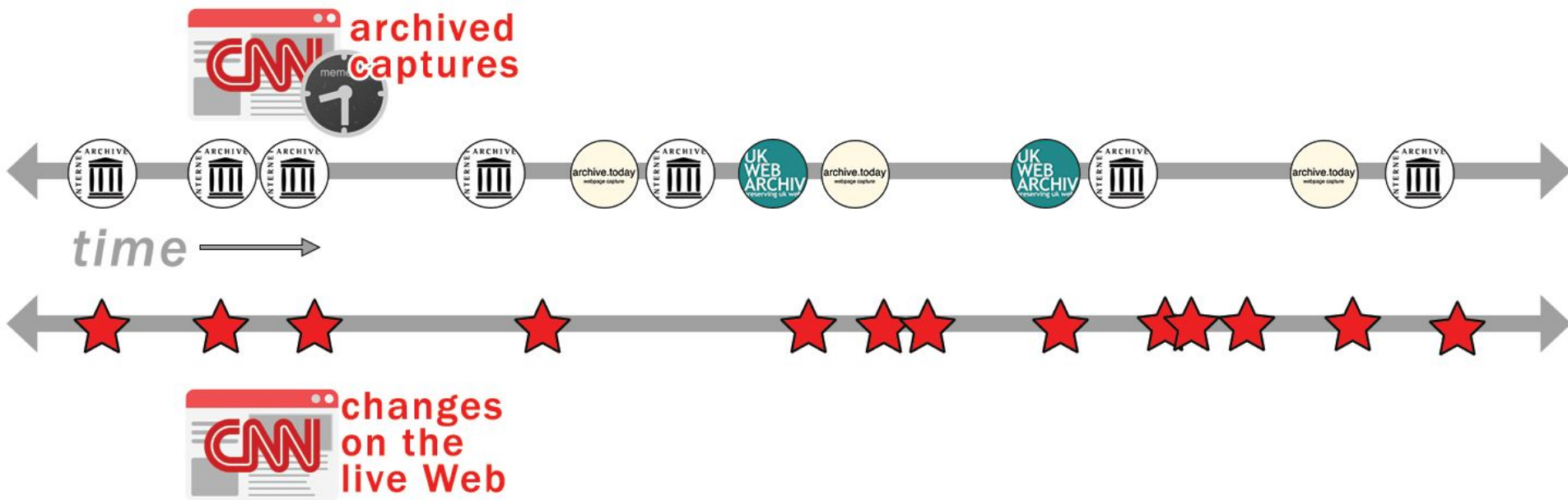
February 14, 2019

Mat Kelly

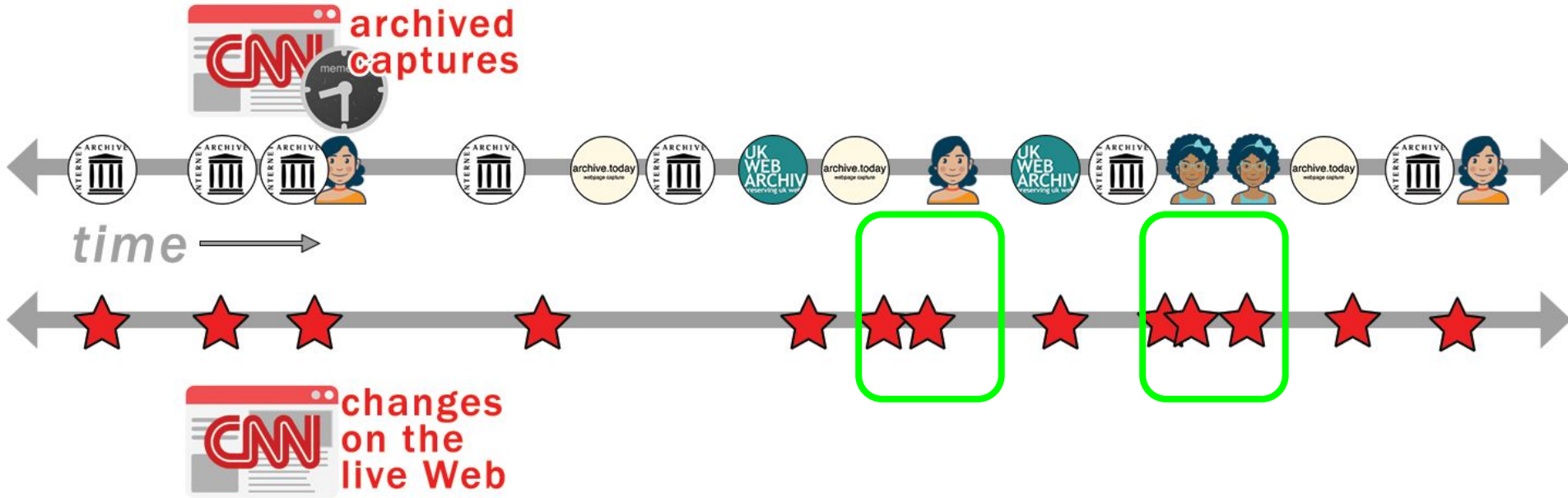
Desire: Include Other Non-Aggregated Archives



Rapidly Changing Pages May Not Be Comprehensively Captured



Archiving More Archives Provides a Better Picture of the Web



Research Questions

RQ1: What sort of content is difficult to capture and replay for preservation from the perspective of a Web browser?

RQ2: How do Web browser APIs compare in potential functionality to the capabilities of archival crawlers?

RQ3: What issues exist for capturing and replaying content behind authentication?

RQ4: How can content that was captured behind authentication signal to Web archive replay systems that it requires special handling?

RQ5: How can Memento aggregators indicate that private Web archive content requires special handling to be replayed, despite being aggregated with publicly available Web archive content?

RQ6: What kinds of access control do users who create private Web archives need to regulate access to their archives?

Research Questions

RQ1: What sort of **content is difficult to capture** and replay for preservation from the perspective of a Web browser?

RQ2: How do **Web browser APIs compare** in potential functionality to the capabilities of archival crawlers?

RQ3: What issues exist for capturing and replaying **content behind authentication**?

RQ4: How can **content** that was captured behind authentication **signal** to Web archive replay systems that it **requires special handling**?

RQ5: How can Memento **aggregators indicate** that private Web archive content requires **special handling** to be replayed, despite being aggregated with publicly available Web archive content?

RQ6: What kinds of access control do users who create private Web archives need to **regulate access** to their archives?

Outline

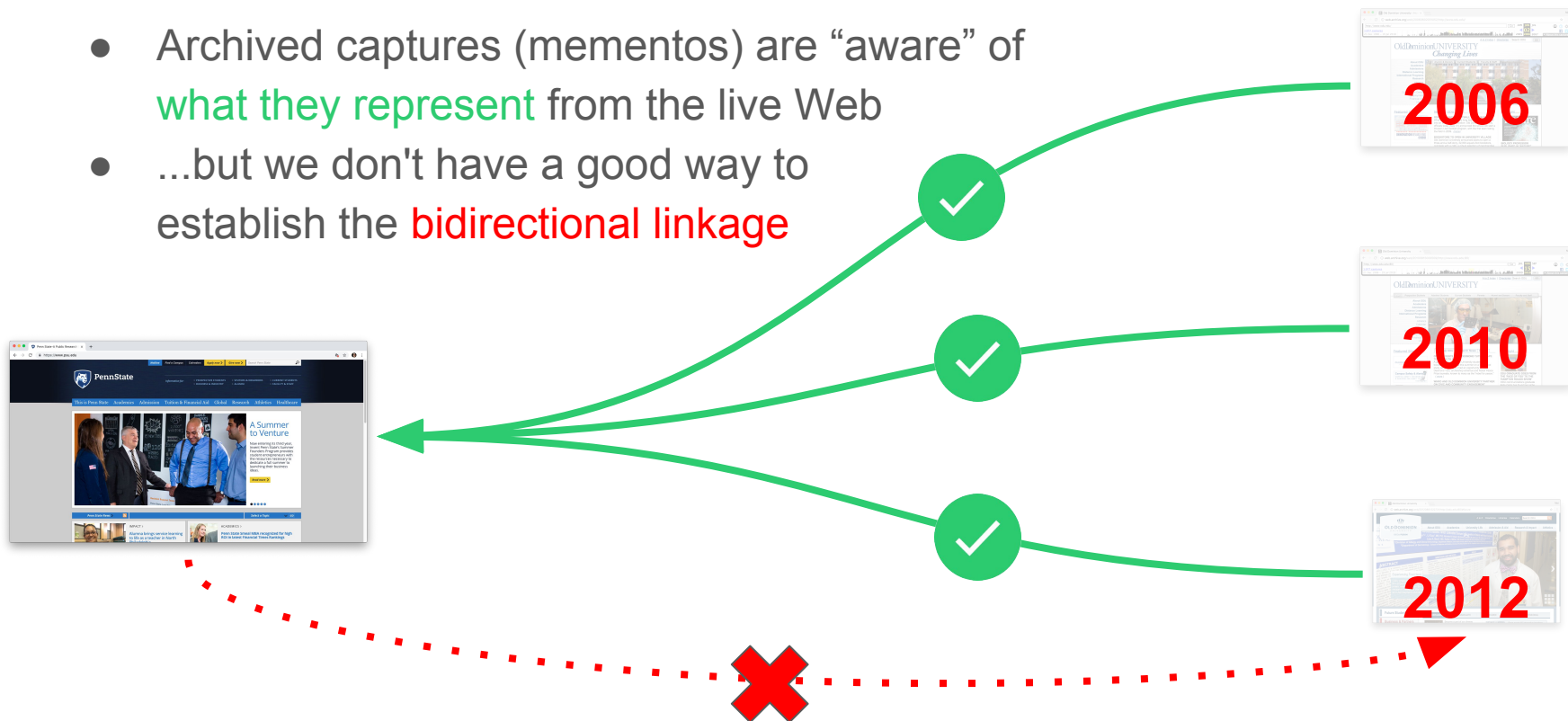
- Introduction/Motivation
- Background
- Preliminary Research
- Proposed Framework
- Evaluation Plan

Outline

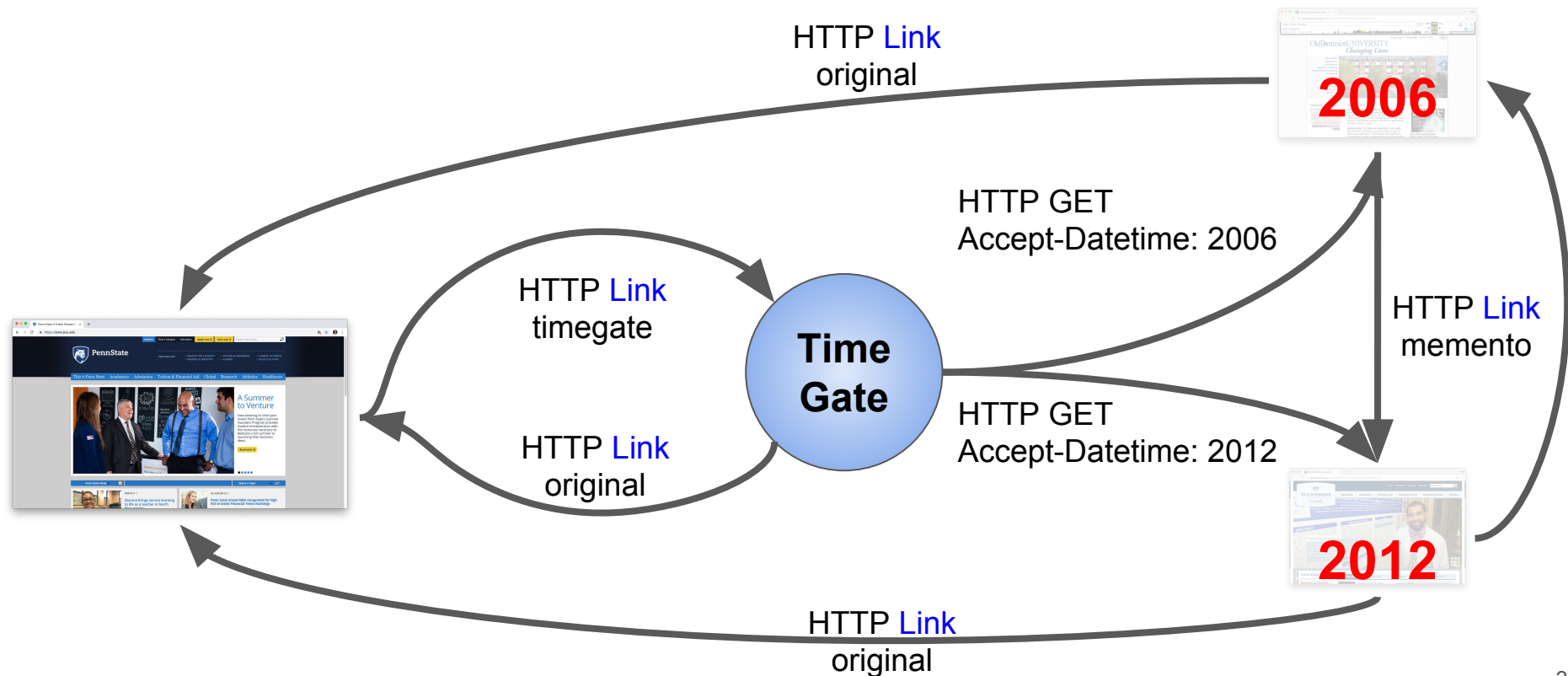
- Introduction/Motivation
- **Background**
- Preliminary Research
- Proposed Framework
- Evaluation Plan

Needed Association of Live-to-Archived Web

- Archived captures (mementos) are “aware” of **what they represent** from the live Web
- ...but we don't have a good way to establish the **bidirectional linkage**

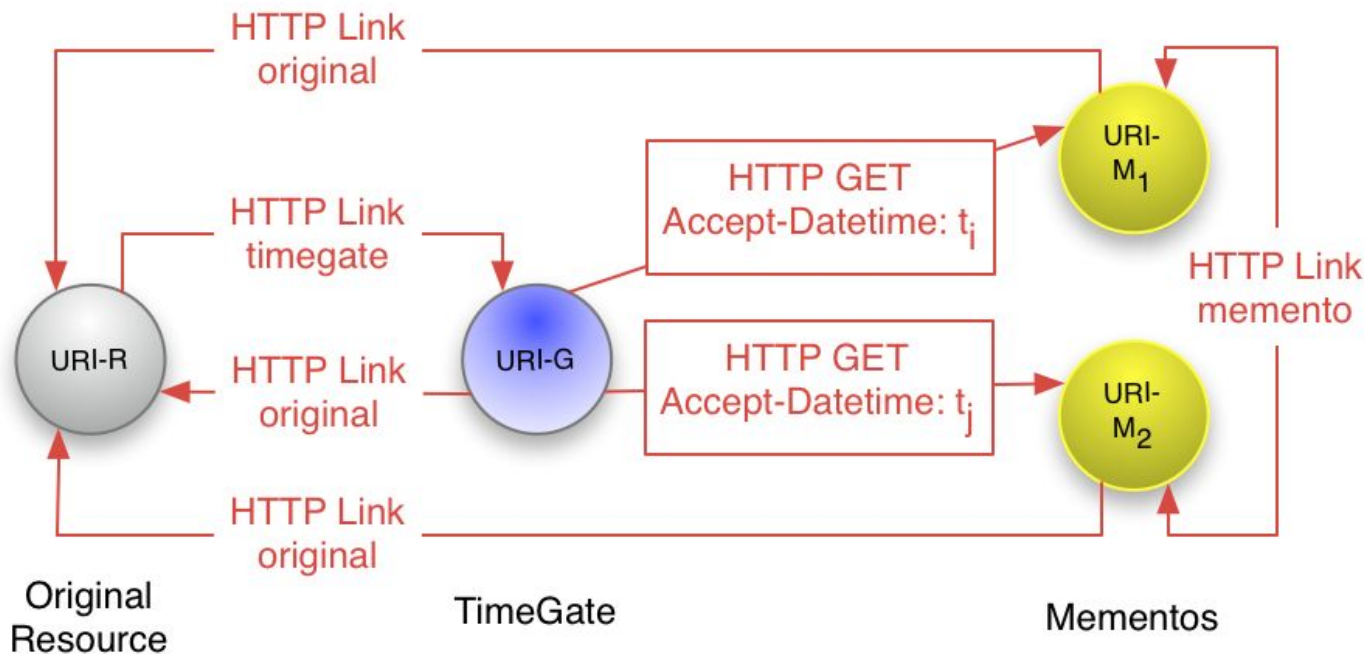


Representations can be **Linked** in time





Background: Memento



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.



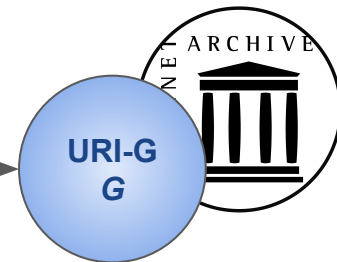
Background: Memento Request Example

HTTP Request

- **Accept-Datetime:** Wed, 02 Aug 2017 23:15:00 GMT
- **GET:** <http://web.archive.org/web/http://www.cnn.com>



Request `cnn.com` at Sept
11, 2001 at 9am EST





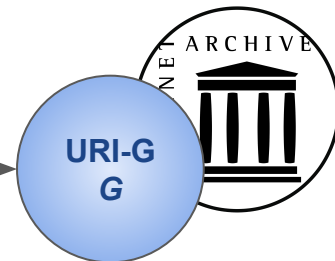
Background: Memento Request Example



HTTP Request

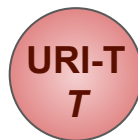
- **Accept-Datetime:** Wed, 02 Aug 2017 23:15:00 GMT
- **GET:** <http://web.archive.org/web/http://www.cnn.com>

Request `cnn.com` at Sept 11, 2001 at 9am EST



HTTP Response (302)

- **Memento-Datetime:** Wed, 02 Aug 2017 23:18:04 GMT
- **Location:** <http://web.archive.org/web/20170802231804/http://www.cnn.com/>
- **Link:**



timemap



original



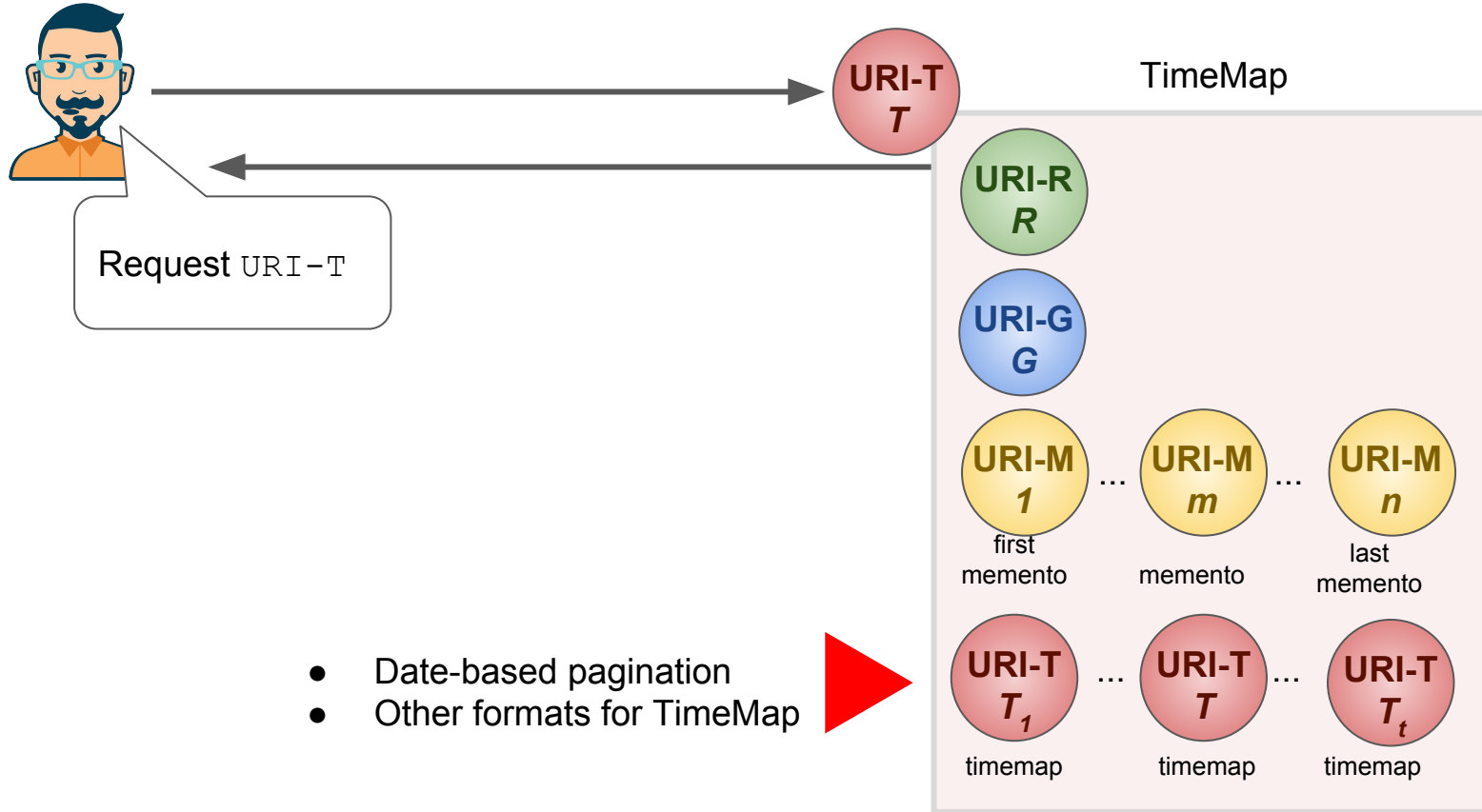
timegate



memento

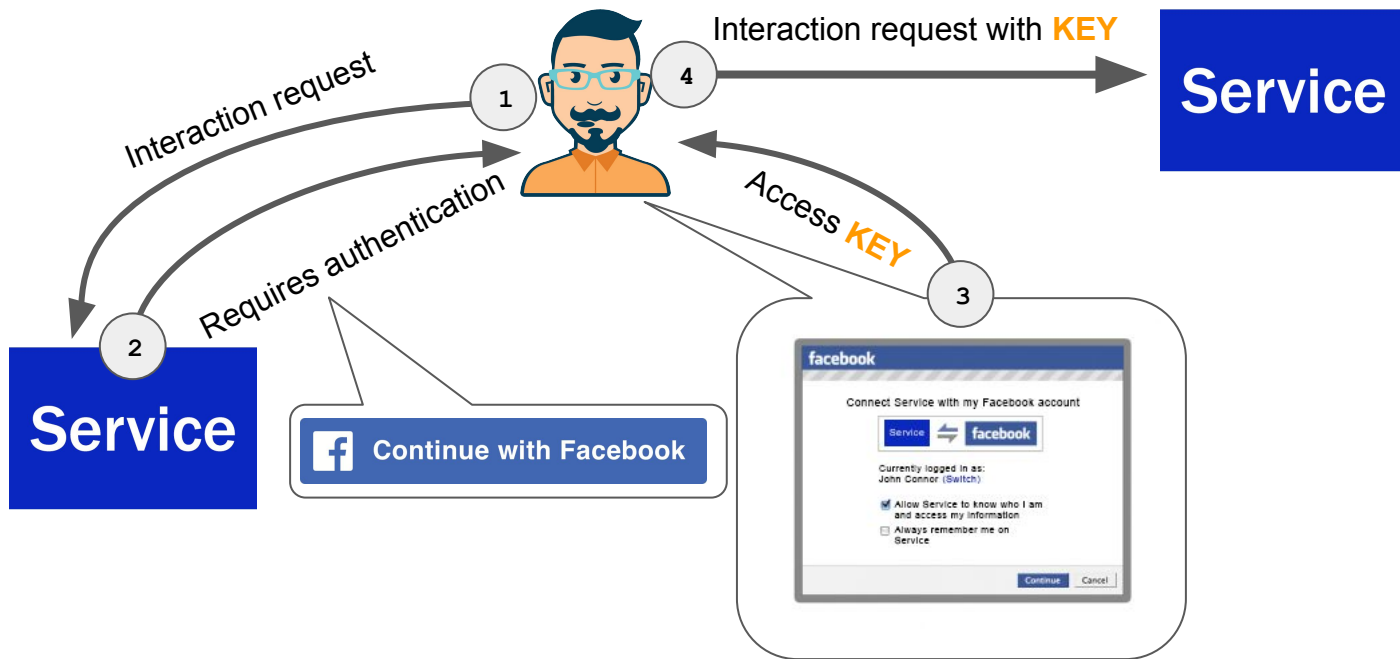


Background: Dereferencing a TimeMap at URI-T



Role-based delegation and authentication

A familiar paradigm used for authentication on the live Web





Background - Privacy and Security

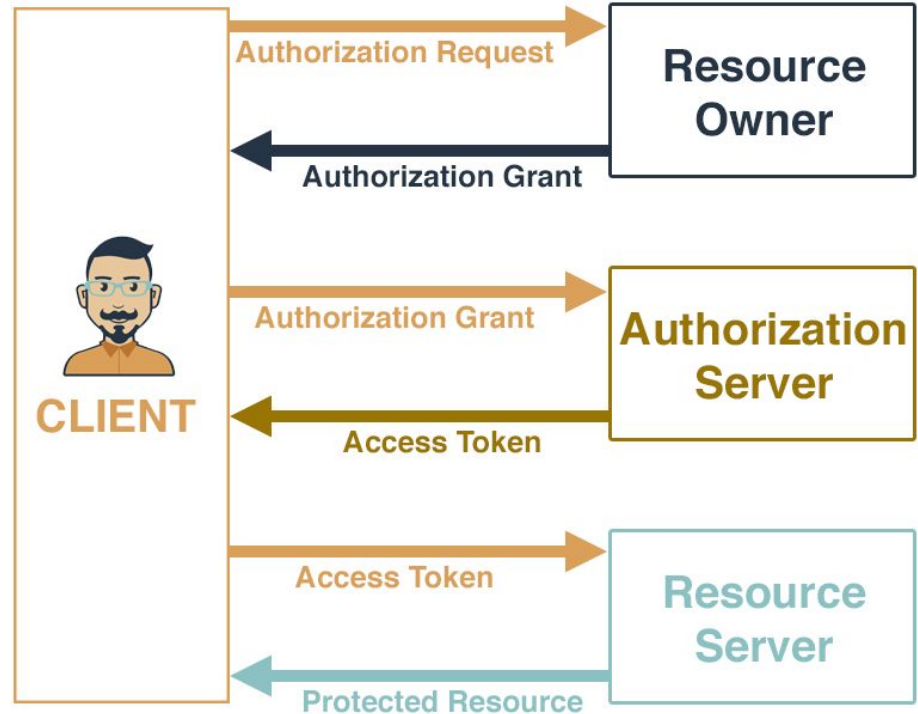
- Web users question trusting institutions to preserve private Web contents¹
- OAuth 2.0² facilitates authentication cohesion of entities

RQ3: What issues exist for capturing and replaying content behind authentication?

RQ6: What kinds of access control do users who create private Web archives need to regulate access to their archives?

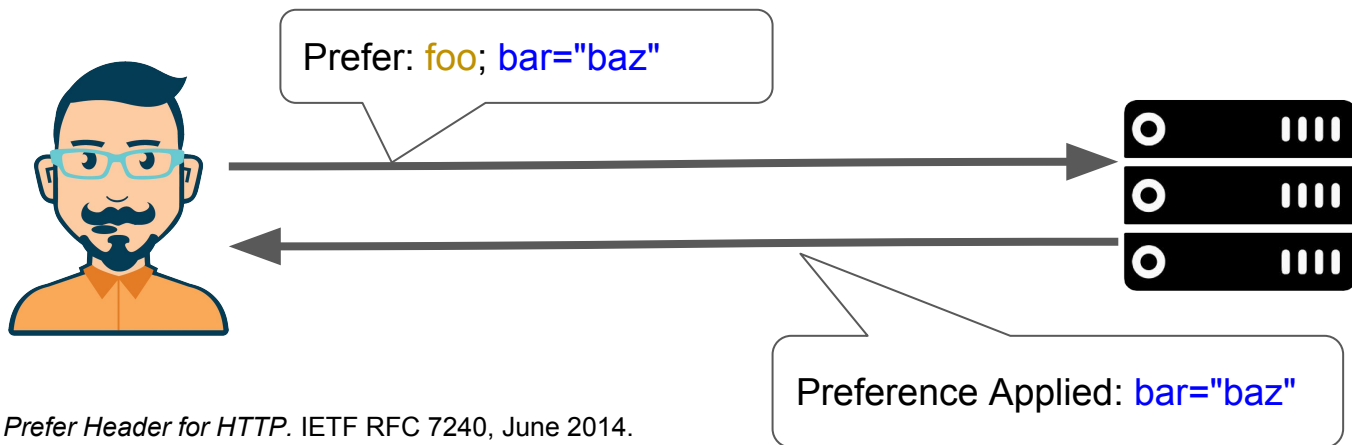
¹ Marshall and Shipman., "On the Institutional Archiving of Social Media", JCDL 2012

² D. Hardt. *The OAuth 2.0 Authorization Framework*. IETF RFC 6749, October 2012.



HTTP Prefer

- HTTP negotiation already available via Accept-* headers
- *Prefer* syntax provide mechanism for client to specify preferences
 - ...with which servers may not comply

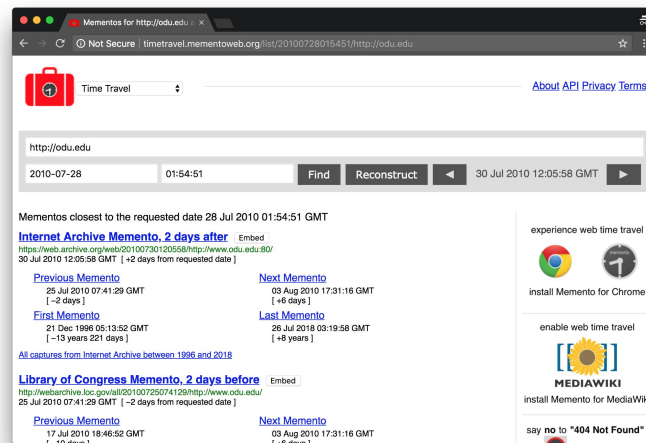
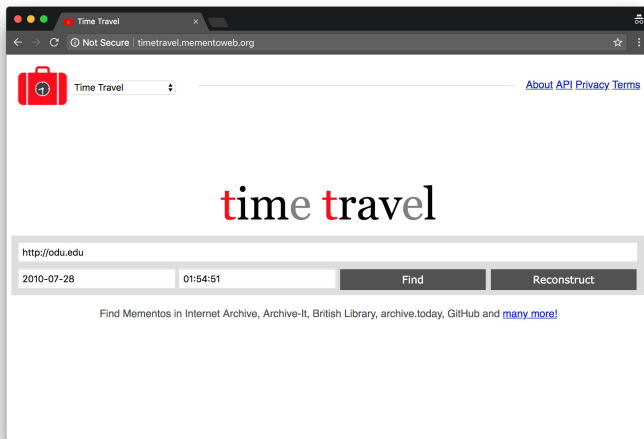


* J. Snell. *Prefer Header for HTTP*. IETF RFC 7240, June 2014.

Memento Aggregation State of the Art



Memento Aggregation - MementoWeb

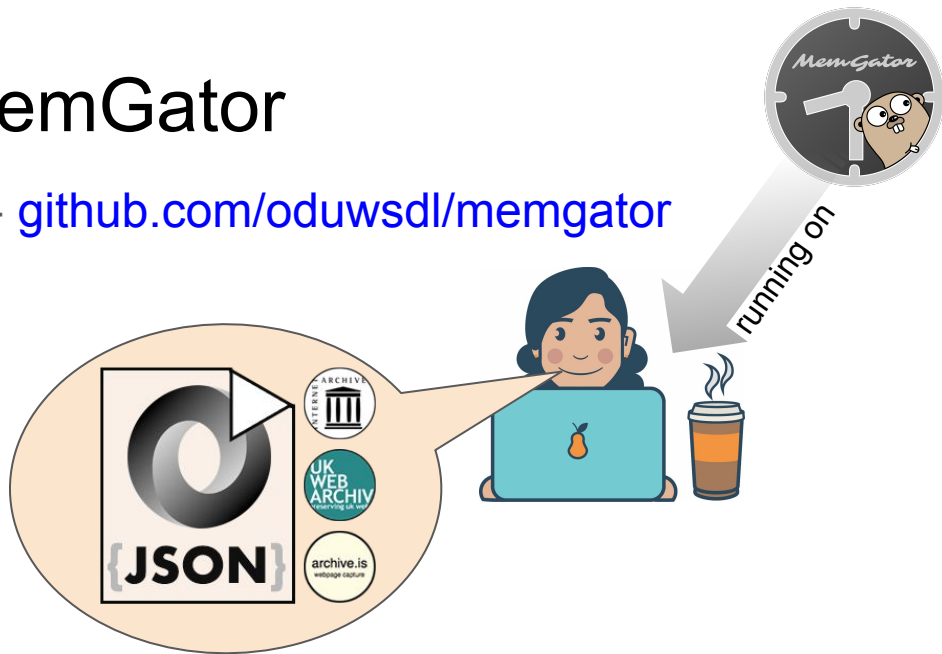


Also available via CLI:

```
$ curl http://timetravel.mementoweb.org/timemap/link/http://odu.edu
```

Memento Aggregation - MemGator

- Open Source Memento Aggregator - github.com/oduwsdl/memgator
- Easy personal/local deployment
- Specify archive list on launch
 - Easily configurable **JSON** →
 - Use default collection if not specified
- TimeMap Formats:
 - Link
 - **JSON**
 - **CDXJ**



* Alam and Nelson, “MemGator - A Portable Concurrent Memento Aggregator: Cross-Platform CLI and Server Binaries in Go”, JCDL 2016

A Framework for Aggregating Public and Private Web Archives

February 14, 2019

Mat Kelly

CDXJ: An Alternative TimeMap Format

```
<http://matkelly.com>; rel="original",
<http://localhost:1208/timemap/link/http://matkelly.com>;
rel="self"; type="application/link-format",
<http://web.archive.org/web/20060514123511/http://www.mat
kelly.com:80/>; rel="first memento"; datetime="Sun, 14
May 2006 12:35:11 GMT",
<http://web.archive.org/web/20060516213852/http://www.mat
kelly.com/>; rel="memento"; datetime="Tue, 16 May 2006
21:38:52 GMT",
...
<http://web.archive.org/web/20180128152125/http://matkell
y.com>; rel="memento"; datetime="Sun, 28 Jan 2018
15:21:25 GMT",
<http://web.archive.org/web/20180319141920/http://matkell
y.com/>; rel="last memento"; datetime="Mon, 19 Mar 2018
14:19:20 GMT",
<http://localhost:1208/timemap/link/http://matkelly.com>;
rel="timemap"; type="application/link-format",
<http://localhost:1208/timemap/json/http://matkelly.com>;
rel="timemap"; type="application/json",
<http://localhost:1208/timemap/cdxj/http://matkelly.com>;
rel="timemap"; type="application/cdxj+ors",
<http://localhost:1208/timegate/http://matkelly.com>;
rel="timegate"
```

```
!context ["http://tools.ietf.org/html/rfc7089"]
!id {"uri": "http://localhost:1208/timemap/cdxj/http://matkelly.com"}
!keys ["memento_datetime_YYYYMMDDhhmmss"]
!meta {"original_uri": "http://matkelly.com"}
!meta {"timegate_uri":
"http://localhost:1208/timegate/http://matkelly.com"}
!meta {"timemap_uri": {"link_format":
"http://localhost:1208/timemap/link/http://matkelly.com", "json_format":
"http://localhost:1208/timemap/json/http://matkelly.com", "cdxj_format":
"http://localhost:1208/timemap/cdxj/http://matkelly.com"}}}
20060514123511 {"uri":
"http://web.archive.org/web/20060514123511/http://www.matkelly.com:80/",
"rel": "first memento", "datetime": "Sun, 14 May 2006 12:35:11 GMT"}
20060516213852 {"uri":
"http://web.archive.org/web/20060516213852/http://www.matkelly.com/",
"rel": "memento", "datetime": "Tue, 16 May 2006 21:38:52 GMT"}
...
20180128152125 {"uri":
"http://web.archive.org/web/20180128152125/http://matkelly.com", "rel":
"memento", "datetime": "Sun, 28 Jan 2018 15:21:25 GMT"}
20180319141920 {"uri":
"http://web.archive.org/web/20180319141920/http://matkelly.com/", "rel":
"last memento", "datetime": "Mon, 19 Mar 2018 14:19:20 GMT"}
```

Link (RFC 7089) TimeMap

Original URI (URI-R)

Other TimeMaps (URI-Ts)

TimeGate (URI-G)

Relative Relations

CDXJ TimeMap

See Alam, [“CDXJ: An Object Resource Stream Serialization Format”](#), 2015

CDXJ: An Alternative TimeMap Format

MAIN POINTS

Link, CDXJ, and JSON TimeMaps:

Multiple formats to express same information

Link syntax is not expandable:

They were meant to be displayed in a constrained environment
(in HTTP Link headers)

Link (RFC 7089) TimeMap

CDXJ TimeMap

Original URI (URI-R)

Other TimeMaps (URI-Ts)

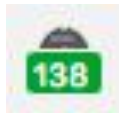
TimeGate (URI-G)

Relative Relations



: visual user interaction with aggregators

- Bridges gap between live and archived Webs
- Leverages Memento aggregator's capability, returns TimeMaps
- Indicates # of captures for a URI while you browse
- Provides navigation of mementos while browsing live Web
- Single-click submission of URI-R to multiple Web archives



8673 mementos available.

Archive Page To...

List mementos by: Dropdown Drilldown

Year	Month	Day	Time
1997	5	Jan	19 1st
1998	3	Feb	34 2nd
1999	5	Mar	30 3rd
2000	118	Apr	29 4th
2001	31	May	41 5th
2002	42	Jun	53 6th
2003	27	Jul	54 7th
2004	177	Aug	101 8th
2005	515	Sep	23 9th
2006	251	Oct	27 10th
2007	107	Nov	47 11th
2008	80	Dec	57 12th
2009	153		13th
2010	119		14th
2011	259		15th
2012	194		16th
2013	671		17th
2014	2058		18th
2015	968		19th
2016	1035		20th
2017	660		22nd
2018	1103		23rd
2019	112		24th

This is Penn

A Summer to Venture

Now entering its third year, Invent Penn State's Summer Founders Program provides student entrepreneurs with the resources necessary to dedicate a full summer to launching their business ideas.

Read more >

Outline

- Introduction/Motivation
- Background
- **Preliminary Research**
- Proposed Framework
- Evaluation Plan

Preliminary Research

1. JCDL 2012 - WARCreate
2. TPDL 2013 - Change in Archivability
3. DLib 2013 - Method for Identifying
4. JCDL 2014 - Mink
5. JCDL 2014 - Archival Acid Test
6. JCDL 2014 - Not All Mementos
7. IJDL 2015 - Impact of JavaScript
8. IJDL 2015 - Not All Mementos
9. JCDL 2015 - Mobile Mink
10. JCDL 2016 - InterPlanetary Wayback (ipwb)
11. TPDL 2016 - ipwb extended
12. JCDL 2017 - WAIL Electron
13. JCDL 2017 - ServiceWorker Replay
14. JCDL 2017 - Impact of Canonicalization
15. JCDL 2018 - ArchiveNow
16. JCDL 2018 - Replay Banners
17. JCDL 2018 - A Framework...

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 into 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 tools that allows a user to “long term storage, maintenance, and access of personal digital assets that have emotional, intellectual, and historical value to individuals” [3].

Preserve everything you see!

Categories and Subject Descriptors

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

General Terms

Personal Web Archiving, WARC, Browser, Wayback Machine, Internet Archive

INTRODUCTION

The Internet Archive, along with web archives in other languages, has become an important resource for web archives. It is 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

web” [1]. Our goal is to allow users, once past authentication, to generate their own archives that can be browse-able in a user-friendly manner.

The Internet Archive's Wayback Machine is the most well-known interface for accessing web archives. The archived pages are stored in the standard Web ARChive (WARC) format [2] and are generated by the Heritrix¹ crawler. Unfortunately, Heritrix is limited to crawling only publicly accessible pages, so many personal pages (*e.g.*, password-protected pages, social media pages) cannot be easily archived. In addition, the pages that are located on user-owned servers (the version archived at Internet Archive is the one that the Heritrix crawler (run from San Francisco) sees. For example, the most recently available version² of <http://www.craigslist.org> redirects to <http://sfbay.craigslist.org>.

To allow a user to use of the standard WARC format to archive personal web archives, we have developed a tool to allow a user to archive any page, edit its metadata, and submit it to an instance of the Wayback Machine (from here on referred to as Wayback).

2. WARCREATE

WARCreate³ is an extension for the Google Chrome web browser that allows a user to generate a WARC file from the current webpage. In addition to creating a valid WARC that can be viewed in Wayback, the extension provides options for how the WARC file is generated. For example, two different users see different content at <http://facebook.com>, and the extension allows the user to generate a WARC file that contains the content as seen by the user. The extension also allows the user to generate a WARC file that contains the content as seen by the user.

To create a WARC file from the current webpage, the user clicks on the browser extension's icon in the address bar and the extension generates the WARC file (see Figure 1). The extension also allows the user to generate a WARC file that contains the content as seen by the user. The extension also allows the user to generate a WARC file that contains the content as seen by the user. The extension also allows the user to generate a WARC file that contains the content as seen by the user.

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

Preliminary Research

1. JCDL 2012 - WARCreate
2. TPDL 2013 - Change in Archivability
3. DLib 2013 - Method for Identifying
4. JCDL 2014 - Mink
5. JCDL 2014 - Archival Acid Test
6. JCDL 2014 - Not All Mementos
7. IJDL 2015 - Impact of JavaScript
8. IJDL 2015 - Not All Mementos
9. JCDL 2015 - Mobile Mink
10. JCDL 2016 - InterPlanetary Wayback (ipwb)
11. TPDL 2016 - ipwb extended
12. JCDL 2017 - WAIL Electron
13. JCDL 2017 - ServiceWorker Replay
14. JCDL 2017 - Impact of Canonicalization
15. JCDL 2018 - ArchiveNow
16. JCDL 2018 - Replay Banners
17. JCDL 2018 - A Framework...

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 page types that have presented challenges to the preservation 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.

Which things are hard to preserve?

Keywords: Web Archiving, Digital Preservation

RQ1: What sort of content is difficult to capture and replay for preservation from the perspective of a Web browser?

Adoption of JavaScript allowed the components on a web page to respond to users' actions or be manipulated in ways that made the page more usable. Ajax [9] combines multiple web technologies to give web pages the ability to perform operations asynchronously. The adoption of Ajax by web developers facilitated the fluidity of user interaction on the web. Through each phase in the progression of the web, the ability to preserve the content displayed to the user has also progressed but in a less linear trend.

A large amount of the difficulty in web archiving stems from the crawler's insufficient ability to capture content related to JavaScript. Because JavaScript is executed on the client side (i.e., within the browser after the page has loaded), it should follow that the archivability could be evaluated using a consistent replay medium. The medium used to archive (normally a web crawler tailored for archiving, e.g., Heritrix [21]) is frequently different from the medium used to replay the archive (henceforth, the *web browser*, the predominant means of

Preliminary Research

1. JCDL 2012 - WARCreate
2. TPDL 2013 - Change in Archivability
3. DLib 2013 - Method for Identifying
4. JCDL 2014 - Mink
5. JCDL 2014 - Archival Acid Test
6. JCDL 2014 - Not All Mementos
7. IJDL 2015 - Impact of JavaScript
8. IJDL 2015 - Not All Mementos
9. JCDL 2015 - Mobile Mink
10. JCDL 2016 - InterPlanetary Wayback (ipwb)
11. TPDL 2016 - ipwb extended
12. JCDL 2017 - WAIL Electron
13. JCDL 2017 - ServiceWorker Replay
14. JCDL 2017 - Impact of Canonicalization
15. JCDL 2018 - ArchiveNow
16. JCDL 2018 - Replay Banners
17. JCDL 2018 - A Framework...

A Framework for Aggregating Public and Private Web Archives

February 14, 2019

Mat Kelly



D-Lib Magazine

November/December 2013
Volume 19, Number 11/12
[Table of Contents](#)

A Method for Identifying Personalized Representations in Web Archives

Mat Kelly, Justin F. Brunelle, Michele C. Weigle, and Michael L. Nelson
Old Dominion University
mkelly,jbrunelle,mweigle,min@cs.odu.edu

[doi:10.1045/november2013.kelly](#)

[Printer-friendly Version](#)

Abstract

Web resources are becoming increasingly personalized – two different users clicking on the same link at the same time can see content customized for each individual user. These changes result in multiple representations of a resource that cannot be canonicalized in Web archives. We identify characteristics of this problem by presenting a potential solution to generalize personalized representations in archives. We also present our proof-of-concept prototype that analyzes WARC (Web ARChive) format files, inserts metadata establishing relationships, and provides archive users the ability to navigate on the additional dimension of environment variables in a modified Wayback Machine.

Introduction

Personalized Web resources offer different representations [1] to different users based on the user-agent string and other values in the HTTP request headers, GeoIP, and other environmental factors. This means Web crawlers capturing content for archives may receive representations based on the crawl environment which will differ from the representations returned to the interactive users. In summary, what we archive is increasingly different from what we as interactive users experience.

Some preserved things are personalized

Web archives have long been criticized for not capturing the full range of the Web. With the increasing prevalence of mobile browsers on the Web (30% - 54% of users have mobile representations [23]), it is becoming important to capture these mobile representations of resources.

Mobile pages often contain links to additional resources instead of embedded text and often reduce the number of images embedded in the page [3]. For example, the mobile representation of <http://espn.go.com/> contains a section on ESPN Videos, while the desktop representation does not. When <http://espn.go.com/> (the "original resource", identified by URI-R), is accessed, it redirects to <http://m.espn.go.com/>, effectively giving two separate but related URI-R values that go into the archive.

Because of the differences in content, the desktop and mobile representations of a resource are not interchangeable. The desktop representation is presented to the user. To quantify the differences, the desktop representation contains 201 links, while the mobile representation contains only 58 links. These link sets are mutually exclusive, with the mobile representation linking to specific resources (such as box-scores and gamecasts) while the desktop representation links to higher-level resources (such as narratives that include box-scores and may have links to gamecasts). A user may review news articles or other content on a mobile device and be unable to recall the article in an archive. To capture and record the complete set of content at <http://espn.go.com/>, each of these different representations, both mobile and desktop, need to be stored in Web archives.

One way to capture both representations is to use a crawler that can crawl the mobile Web by setting its user-agent string to a mobile browser. This can potentially lead to multiple representations of the same content being

RQ2: How do Web browser APIs compare in potential functionality to the capabilities of archival crawlers?

Therefore, it is no longer sufficient to only navigate archives in the temporal dimension; archives must also provide users the opportunity to understand how a representation was recorded and

In this work, we explore the issue of personalized representations in Web archives, propose a framework to solve this problem, and present a proof-of-concept prototype that integrates personalized representations into the existing Wayback Machine. We use two different methods for creating archived representations called mementos (identified by URI-M) to a canonical representation. This prototype extends the description of mementos from only "when" they were archived (temporal dimension) to "where" and "how" (GeoIP and browser environments). Users can then browse between mementos based on temporal or environmental dimensions.

Personalized, Anonymous Representations

Dynamic and personalized representations of Web 2.0 resources that are generated by technologies such as JavaScript can differ greatly depending on several factors. For example, some sites attempt to provide alternate representations by interpreting the user-agent portion of the HTTP GET headers and use content negotiation to determine which representation to return.

We ran a pair of limited crawls of the cnn.com front page with Heritrix 3.1 and then accessed the mementos captured by Heritrix with a desktop Mac and an Android phone. The first crawl captured the cnn.com front page and specified a desktop version of the Mozilla browser as the user-agent: string in the header string, as seen in Figure 1. The resulting Web ARChive (WARC) file [24] is viewed in a local installation of the Wayback Machine [22] and is shown in Figures 3(a) and 3(c).

The second crawl captured the cnn.com front page and specified an iPhone version of the Mozilla browser as the user-agent: string in the header, as seen in Figure 2. The resulting WARC, as viewed in the Wayback Machine, is shown in Figures 3(b) and 3(d). The mobile and desktop representations differ in archives, but their relationship as permutations of each other is neither recorded nor seen by users; a user of the Wayback Machine may not understand how these representations are generated since they are identified by the same URI-R. We refer to these differing representations of the same URI-R built with differing environments as *personalized representations* of the resource R.

The headers in Figures 1 and 2 reference the user-agent: string with <http://yourdomain.com>, which is a place holder for the URI for whom the crawl is being executed. For example, a crawl originating from Old Dominion University's Computer Science department would read <http://www.cs.odu.edu/>.

```
WARC/1.0
WARC-Type: request
WARC-Target-URI: http://www.cnn.com/
WARC-Date: 2013-03-03T02:16:57.000Z
```


Preliminary Research

1. JCDL 2012 - WARCreate
2. TPDL 2013 - Change in Archivability
3. DLib 2013 - Method for Identifying
4. JCDL 2014 - Mink
5. JCDL 2014 - Archival Acid Test
6. JCDL 2014 - Not All Mementos
7. IJDL 2015 - Impact of JavaScript
8. IJDL 2015 - Not All Mementos
9. JCDL 2015 - Mobile Mink
10. JCDL 2016 - InterPlanetary Wayback (ipwb)
11. TPDL 2016 - ipwb extended
12. JCDL 2017 - WAIL Electron
13. JCDL 2017 - ServiceWorker Replay
14. JCDL 2017 - Impact of Canonicalization
15. JCDL 2018 - ArchiveNow
16. JCDL 2018 - Replay Banners
17. JCDL 2018 - A Framework...



Mink: Integrating the Live and Archived Web Viewing Experience Using Web Browsers and Memento

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

We describe Mink, a new web browser extension that provides a different model for integration of the live and archived web. While a user browses the live web, Mink actively queries the archives and reports other instances of the page in the archives without requiring active querying by the user. Further, by querying the archives dynamically and asynchronously, a user can view the extent to which the currently viewed page on the live web has been archived and proactively submit a request to various archives using an overlay on the live web page and a simple interface.

web. We have developed a new browser extension, Mink², that instead uses an unobtrusive alert model to remind the user about the past. This model allows the user to quickly poll through the mementos available while maintaining the paradigm of relying on what is returned by the server to determine whether the user stays in the past or returns to the present. The additional feature of allowing the user to seamlessly jump from the past to the present while maintaining a quick return to the past makes Mink's approach unique.

Categories and Subject Descriptors

H.3.7 [Online Information Services]: Web-based services and Archives

1. INTRODUCTION

To better integrate the past and live web, implementations of the Memento framework [1] provide the facilities to query the archives (using URI and HTTP Accept-Datetime headers as parameters) to provide resources on the past web.

RQ2: How do Web browser APIs compare in potential functionality to the capabilities of archival crawlers?

RQ5: How can Memento aggregators indicate that private Web archive content requires special handling to be replayed, despite being aggregated with publicly available Web archive content?

We chose the Google Chrome browser extension environment due to the browser's popularity, but the logic is simple enough to be ported to other browsers.

When a user loads a page, Mink queries a meta-search engine for Memento URIs and expects a response. While processing the response, Mink displays a "spinning" animation at the bottom right of the browser viewport and provides a "TimeMap" is paginated with a reference to a subsequent TimeMap, a button is provided to the user to invoke the iterative fetching of the TimeMap. Once the TimeMap is fully expanded, a badge is displayed in the bottom right corner of the browser viewport, indicating how many mementos are available for the current page. This badge is updated dynamically as the user observes how well pages are archived without needing to commit to browsing the archived web nor to proactively submit a request to the archives to receive this archival metadata about the live web.

Once a user has accessed an archived page using Mink, the interface provides an additional button that allows the user to return to the live web with a single click for easy compar-

²Named for Minkowski Space

³Available at <https://github.com/machawk1/mink>

Preliminary Research

1. JCDL 2012 - WARCreate
2. TPDL 2013 - Change in Archivability
3. DLib 2013 - Method for Identifying
4. JCDL 2014 - Mink
5. JCDL 2014 - Archival Acid Test
6. JCDL 2014 - Not All Mementos
7. IJDL 2015 - Impact of JavaScript
8. IJDL 2015 - Not All Mementos
9. JCDL 2015 - Mobile Mink
10. JCDL 2016 - InterPlanetary Wayback (ipwb)
11. TPDL 2016 - ipwb extended
12. JCDL 2017 - WAIL Electron
13. JCDL 2017 - ServiceWorker Replay
14. JCDL 2017 - Impact of Canonicalization
15. JCDL 2018 - ArchiveNow
16. JCDL 2018 - Replay Banners
17. JCDL 2018 - A Framework...

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 evaluated. In this paper, we propose a set of metrics to evaluate the capability of archival crawlers and other preservation tools using 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, Verification

Keywords

Web Crawler, Web Archiving, Digital Preservation

RQ1: What sort of content is difficult to capture and replay for preservation from the perspective of a Web browser?

Because archival crawlers attempt to duplicate what a user would see if he accessed the page on the live web, variance from what is preserved 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 crawler did if the information no longer exists on the live web. By examining what sort of web content is inaccurately represented or missing from the web archives, it would be useful to evaluate the capability of archival crawlers (in respect to that of web browsers that implement the latest technologies) to determine what might be missing from their functional repertoire.

Web browsers exhibited this deviation between each other in the early days of Web Standards. A series of "Acid Tests" followed each browser to evaluate 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 live page, this is not always the case when the archived version of the page is rendered. This difference can be used to highlight the features that archival crawlers are lacking compared to web browsers and thus emphasize the deviations that will occur in web archives compared to what the original page would have displayed.

Web archiving is to capture web pages so they can be "replayed" at a later date. Web archiving tools access these pages on the live web in a manner similar to tools used by search engines (crawlers) and preserve the pages in a format that allows the data and contextual information about the crawl to be re-experienced. These "archival crawlers" take different approaches in digital preservation and thus their capability and scope vary.

The development of a standard format for the preservation of web pages is a goal of many tools in a variety of formats. An ISO standard format utilized by institutional and personal web archivists alike is the Web Archive (WARC) format [1]. WARC files allow HTTP communication that occurred during a crawl as well as payload, metadata and other archival features to be encoded 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 WARCs but IA's Wayback Machine (and its open source counterpart¹) is the de facto standard.

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

Preliminary Research

1. JCDL 2012 - WARCreate
2. TPDL 2013 - Change in Archivability
3. DLib 2013 - Method for Identifying
4. JCDL 2014 - Mink
5. JCDL 2014 - Archival Acid Test
6. JCDL 2014 - Not All Mementos
7. IJDL 2015 - Impact of JavaScript
8. IJDL 2015 - Not All Mementos
9. JCDL 2015 - Mobile Mink
10. JCDL 2016 - InterPlanetary Wayback (ipwb)
11. TPDL 2016 - ipwb extended
12. JCDL 2017 - WAIL Electron
13. JCDL 2017 - ServiceWorker Replay
14. JCDL 2017 - Impact of Canonicalization
15. JCDL 2018 - ArchiveNow
16. JCDL 2018 - Replay Banners
17. JCDL 2018 - A Framework...

A Framework for Aggregating Public and Private Web Archives
February 14, 2019
Mat Kelly

Not All Mementos Were Created Equal: Measuring The Impact Of Missing Resources

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

ABSTRACT

Web archives do not capture every resource on every page that they attempt to archive. This results in archived pages missing a portion of their embedded resources. These embedded resources have varying historic, utility, and importance values. The proportion of missing embedded resources does not provide an accurate measure of their impact on the Web page; some embedded resources are more important to the utility of a page than others. We propose a method to measure the relative value of embedded resources and assign a damage rating to archived pages as a way to evaluate archival success. In this paper, we show that Web users' perceptions of damage are not accurately estimated by the proportion of missing embedded resources. The proportion of resource damage and a human solution. We propose a damage rating algorithm that provides closer alignment to Web user perception, providing an overall improved agreement with users on memento damage by 17% and an improvement by 31% if the mementos are not similarly damaged. We use our algorithm to measure damage in the Internet Archive, showing that it is getting better at mitigating damage over time (going from 0.16 in 1998 to 0.13 in 2013). However, we show that a greater number of important embedded resources (2.05 per memento on average) are missing

Categories and Subject Descriptors

General Terms

Design, Experimentation, Measurement

Keywords

Web Archiving, Digital Preservation, HTTP, TimeMaps



1. INTRODUCTION

Web archives are valuable cultural repositories that capture and store Web content. Users make use of archives like the Internet Archive [16, 25] to retrieve archived material [11, 14] for a variety of purposes and in a variety of ways [3]. However, the resources being requested by Web users may not be complete; embedded resources are sometimes missing from an archived Web page [4]. Missing embedded resources return a non-200 HTTP status (e.g., 404, 503) when their URL is dereferenced.

Large images are often more important to an archived page's utility than small images. Similarly, stylesheets that format visible content are more important to the representation of the page than stylesheets without significant formatting. We propose a method to assess the relative importance of embedded resources in the archives.

Throughout this paper we use Memento Framework terminology. Memento [26] is a framework that allows web users to browse in the temporal dimension by aggregating the different versions of a single page.

Original (or live web) resources are identified by URI-R, and 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.

Web user satisfaction (i.e., the utility of mementos). Using resource importance, we can estimate the perceived damage of a missing embedded resource of the memento is missing (e.g., a main image or video essential to the user's understanding of the page), or the missing embedded resource is a spacer image or a small button logo that contributes little to the memento's utility for the user. We propose a method of weighting embedded resources in a memento according to importance. We show that this is an

improved damage rating over an unweighted count of missing embedded resources. We use the unweighted measure of damage as the proportion of missing embedded resources to all requested resources (M_m) and compare it to our algorithm's calculation of damage (D_m).

Third and finally, we measure damage in the Internet

Not all missing resources are created equal

RQ1: What sort of content is difficult to capture and replay for preservation from the perspective of a Web browser?

BEST STUDENT PAPER AWARD
at JCDL 2014

Preliminary Research

1. JCDL 2012 - WARCreate
2. TPDL 2013 - Change in Archivability
3. DLib 2013 - Method for Identifying
4. JCDL 2014 - Mink
5. JCDL 2014 - Archival Acid Test
6. JCDL 2014 - Not All Mementos
7. IJDL 2015 - Impact of JavaScript
8. IJDL 2015 - Not All Mementos
9. JCDL 2015 - Mobile Mink
10. JCDL 2016 - InterPlanetary Wayback (ipwb)
11. TPDL 2016 - ipwb extended
12. JCDL 2017 - WAIL Electron
13. JCDL 2017 - ServiceWorker Replay
14. JCDL 2017 - Impact of Canonicalization
15. JCDL 2018 - ArchiveNow
16. JCDL 2018 - Replay Banners
17. JCDL 2018 - A Framework...

A Framework for Aggregating Public and Private Web Archives
February 14, 2019
Mat Kelly

The impact of JavaScript on archivability

Justin F. Brunelle · Mat Kelly · Michele C. Weigle ·
Michael L. Nelson

Received: 7 November 2013 / Revised: 12 January 2015 / Accepted: 14 January 2015 / Published online: 25 January 2015
© Springer-Verlag Berlin Heidelberg 2015

Abstract As web technologies evolve, web archivists work to adapt so that digital history is preserved. Recent advances in web technologies have introduced client-side executed scripts (Ajax) that, for example, load data without a change in top level Universal Resource Identifier (URI) or require user interaction (e.g., content loading via Ajax when the user scrolls). Supporting a complete and accurate digital history of the web requires a more robust effort to understand why mementos (archived versions of live resources) in today's archives vary in completeness and sometimes pull content from the live web, we present a study of web resources and archival tools. We used a collection of URIs shared over Twitter and a collection of URIs curated by Archive-It in our investigation. We created local archived versions of the URIs from the Twitter and Archive-It sets using WebCite, wget, and the Heritrix crawler. We found that only

12.0 % from 2005 to 2012. We also show that JavaScript is responsible for 33.2 % more missing resources in 2012 than in 2005. This shows that JavaScript is responsible for an increasing proportion of the embedded resources unsuccessfully loaded by mementos. JavaScript is also responsible for 52.7 % of all missing embedded resources in our study.

J. F. Brunelle (✉) · M. Kelly · M. C. Weigle · M. L. Nelson
Department of Computer Science, Old Dominion University,
Norfolk, VA 23529, USA
e-mail: jbrunelle@cs.odu.edu

M. Kelly
e-mail: mkelly@cs.odu.edu

M. C. Weigle
e-mail: mweigle@cs.odu.edu

M. L. Nelson
e-mail: mln@cs.odu.edu

Missing JavaScript has big ramifications

1 Introduction
How well can we archive the web? This is a question that is becoming increasingly important and more difficult to answer. Additionally, this question has significant impact on

RQ1: What sort of content is difficult to capture and replay for preservation from the perspective of a Web browser?

The web has gone through a gradient of changes fueled by early websites were relatively static and the description of web technologies has made the pages personalized and more interactive.

JavaScript, which executes on the client, provides additional features for the web user, enabling or increasing interactivity, client-side state changes, and personalized representations. These additional features offer an enhanced browsing experience for the user.

JavaScript has enabled a wide-scale migration from web pages to web applications. This migration continued with the introduction of Ajax (first introduced in 2005 [28]), which combined multiple technologies to give web pages the ability to perform asynchronous client-server interactions after the HTML is loaded. The first wide-scale implementation of Ajax was in Google Maps in 2005, but Ajax was officially added as a standard in 2006 [70]. While archival tools per-

Preliminary Research

1. JCDL 2012 - WARCreate
2. TPDL 2013 - Change in Archivability
3. DLib 2013 - Method for Identifying
4. JCDL 2014 - Mink
5. JCDL 2014 - Archival Acid Test
6. JCDL 2014 - Not All Mementos
7. IJDL 2015 - Impact of JavaScript
8. IJDL 2015 - Not All Mementos
9. JCDL 2015 - Mobile Mink
10. JCDL 2016 - InterPlanetary Wayback (ipwb)
11. TPDL 2016 - ipwb extended
12. JCDL 2017 - WAIL Electron
13. JCDL 2017 - ServiceWorker Replay
14. JCDL 2017 - Impact of Canonicalization
15. JCDL 2018 - ArchiveNow
16. JCDL 2018 - Replay Banners
17. JCDL 2018 - A Framework...

A Framework for Aggregating Public and Private Web Archives
February 14, 2019
Mat Kelly

Not all mementos are created equal: measuring the impact of missing resources

Justin F. Brunelle¹ · Mat Kelly¹ · Hany SalahEldeen¹ · Michele C. Weigle¹ · Michael L. Nelson¹

Received: 3 December 2014 / Revised: 22 April 2015 / Accepted: 22 April 2015 / Published online: 6 May 2015
© Springer-Verlag Berlin Heidelberg 2015

Abstract Web archives do not always capture every resource on every page that they attempt to archive. This results in archived pages missing a portion of their embedded resources. These embedded resources have varying historic, utility, and importance values. The proportion of missing embedded resources on a page can vary significantly based on their impact on the Web page; some embedded resources are more important than others. We propose a method to estimate the impact of missing embedded resources and assign a damage rating to archived pages as a way to evaluate archival success. In this paper, we show that web users' perceptions of damage are not accurately estimated by the proportion of missing embedded resources. In fact, the proportion of missing embedded resources is a less accurate estimate of resource damage than a random selection. We propose a damage rating algorithm that provides

an improved estimate of damage. Web pages with a damage rating of 0.16 in 1998 to 0.13 in 2013). However, we show that a greater number of important embedded resources (2.05 per memento on average) are missing over time. Alternatively, the damage in WebCite is increasing over time (going from 0.375 in 2007 to 0.475 in 2014), while the missing embedded resources are decreasing over time (going from 0.375 in 2007 to 0.475 in 2014). Finally, we investigate the impact of JavaScript on the resources, showing that a crawler that can store Web content (e.g., 404, 503) when their URI is dereferenced.

✉ Justin F. Brunelle
jbrunelle@cs.odu.edu
Mat Kelly
mkelly@cs.odu.edu
Hany SalahEldeen
hany@cs.odu.edu
Michele C. Weigle
mweigle@cs.odu.edu
Michael L. Nelson
mln@cs.odu.edu

¹ Department of Computer Science, Old Dominion University, Norfolk, VA 23529, USA

Is the metric for missing resources applicable across Web?

Web archives do not always capture every resource on every page that they attempt to archive. This results in archived pages missing a portion of their embedded resources. These embedded resources have varying historic, utility, and importance values. The proportion of missing embedded resources on a page can vary significantly based on their impact on the Web page; some embedded resources are more important than others. We propose a method to estimate the impact of missing embedded resources and assign a damage rating to archived pages as a way to evaluate archival success. In this paper, we show that web users' perceptions of damage are not accurately estimated by the proportion of missing embedded resources. In fact, the proportion of missing embedded resources is a less accurate estimate of resource damage than a random selection. We propose a damage rating algorithm that provides

Keywords Web architecture · Web archiving · Digital preservation · Memento damage

RQ1: What sort of content is difficult to capture and replay for preservation from the perspective of a Web browser?

1 Introduction

Archivists work to ensure archives are as complete—and as high quality—as possible. Through identifying sources of missing content or archival difficulties, archivists can address archival challenges by taking steps to adjust processes or to fill in gaps in archive collections. Reyes et al. identified current efforts within several archives to assess their archival collections [4]. Of the archivists sampled, 61 % confirmed that their goal is to assess the quality of every Web page captured, 43 % assess quality

Preliminary Research

1. JCDL 2012 - WARCreate
2. TPDL 2013 - Change in Archivability
3. DLib 2013 - Method for Identifying
4. JCDL 2014 - Mink
5. JCDL 2014 - Archival Acid Test
6. JCDL 2014 - Not All Mementos
7. IJDL 2015 - Impact of JavaScript
8. IJDL 2015 - Not All Mementos
9. JCDL 2015 - Mobile Mink
10. JCDL 2016 - InterPlanetary Wayback (ipwb)
11. TPDL 2016 - ipwb extended
12. JCDL 2017 - WAIL Electron
13. JCDL 2017 - ServiceWorker Replay
14. JCDL 2017 - Impact of Canonicalization
15. JCDL 2018 - ArchiveNow
16. JCDL 2018 - Replay Banners
17. JCDL 2018 - A Framework...

Mobile Mink: Merging Mobile and Desktop Archived Webs

Wesley Jordan¹, Mat Kelly², Justin F. Brunelle^{2,3}, Laura Vobrak¹, Michele C. Weigle²,
and Michael L. Nelson²

¹ New Horizons Regional Education Center Governor's School for Science and Technology

² Old Dominion University, Department of Computer Science

³ 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 URIs and provides the user an aggregated TimeMap of both mobile and desktop mementos. Mobile Mink also allows users to submit mobile and desktop URIs 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, Languages, Management

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 the current web. Mementos are snapshots of web pages archived by URI-Rs. 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 popularity, and the mobile web continues to become more prevalent [9].

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].

Recoupled mobile and desktop archived Webs

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 resources 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 called a “view as” menu. This menu offers a variety of options, one of which is the option to “Share” the page (Figure 1(a)). The “Share” option opens a list of sharing options (Figure 1(b)).

Selecting the option of viewing mementos begins the process of discovering mobile and desktop URIs of the current URI-R. First, Mobile Mink identifies the URI-R of the currently viewed page. Mobile Mink identifies the URI-R as either a desktop URI or a mobile URI. Second, if the URI is a desktop URI, Mobile Mink translates the URI to a mobile URI.

Mobile Mink translates the URI to a mobile URI by applying the same URI modifications as the live web (i.e., returns an HTTP 200 response) and in the archives (returns a TimeMap of mementos > 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 URI is used to identify the mobile and desktop representations, and only content-negotiation based on the user-agent is used by the server to

RQ2: How do Web browser APIs compare in potential functionality to the capabilities of archival crawlers?

Permission is granted to reproduce copies of part or all of this work for personal or classroom use is granted without fee provided that copies are made or distributed for profit or commercial advantage, and that copies bear this notice and the citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner(s). Copyright is held by the author(s).

JCDL'15, June 21–25, 2015, Knoxville, Tennessee, USA.

ACM 978-1-4503-3594-2/15/06.

http://dx.doi.org/10.1145/2756406.2756956.



BEST POSTER AWARD
at JCDL 2015

A Framework for Aggregating Public and Private Web Archives

February 14, 2019

Mat Kelly

Preliminary Research

Sawood Alam, Mat Kelly, and Michael L. Nelson
Old Dominion University, Department of Computer Science, Norfolk VA, 23529, USA
{salam,mkelly,mln}@cs.odu.edu

1. JCDL 2012 - WARCreate
2. TPDL 2013 - Change in Archivability
3. DLib 2013 - Method for Identifying
4. JCDL 2014 - Mink
5. JCDL 2014 - Archival Acid Test
6. JCDL 2014 - Not All Mementos
7. IJDL 2015 - Impact of JavaScript
8. IJDL 2015 - Not All Mementos
9. JCDL 2015 - Mobile Mink
10. JCDL 2016 - InterPlanetary Wayback (ipwb)
11. TPDL 2016 - ipwb extended
12. JCDL 2017 - WAIL Electron
13. JCDL 2017 - ServiceWorker Replay
14. JCDL 2017 - Impact of Canonicalization
15. JCDL 2018 - ArchiveNow
16. JCDL 2018 - Replay Banners
17. JCDL 2018 - A Framework...

ABSTRACT

To facilitate permanence 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 file system that inherently allows deduplication and facilitates opt-in replication. We split the header and payload of WARC response records before disseminating into IPFS to leverage the deduplication, build a CDXJ index, and combine them at the time of replay. From a 1.0 GB sample Archive-It collection of WARCs containing 21,994 mementos, we found that on an average, 570 files can be indexed and disseminated into IPFS per minute. We also found that in our naive prototype implementation, replay took on an average 370 milliseconds per request.

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

Figure 1: A single-line CDXJ record template, shown on multiple lines for readability

about WARC records within IPFS (i.e., the content digest needed for lookup in IPFS).

IPFS is a content addressable peer-to-peer distributed file system [2]. By extracting the HTTP response body (henceforth “payload”) from the records within a WARC file, IPFS allows our prototype to generate a signature uniquely representative of this content. This payload can then be pushed into the IPFS system and retrieved at a later date when the URI-EM is queried. Content addressability allows the user to retrieve the content in the content in the network.

1. INTRODUCTION

The recently created InterPlanetary File System (IPFS) [2] is showing the potential to facilitate data persistence through a peer-to-peer network. In this paper we describe how we use this paper we describe how we use InterPlanetary Wayback (ipwb), that partitions, indexes, and deploys the payloads of archival data into the IPFS peer-to-peer “permanent web” to facilitate permanent and redundant preservation and replay.

The Web Archive (WARC) format is an ISO standard² to store live web archive content in a concatenated record-based file. IA’s web crawler, Heritrix [3], generates WARC files to be read and the content re-experienced in an archival replay system. OpenWayback³ (written in Java) and pywb⁴ (written in Python) are two such replay systems. We leverage

Personal archives are more resilient when propagated

2. IMPLEMENTATION

CDXJ is a text-based file format that we utilize to store indexes of the archived content. Each line in the CDXJ file holds one index record. The line begins with a SURT URI⁵ and datetime followed by a single-line JSON block that stores reference to the content and other arbitrary metadata (Figure 1). We utilize the last field in a CDXJ record (a JSON

RQ4: How can content that was captured behind authentication signal to Web archive replay systems that it requires special handling?

RQ6: What kinds of access control do users who create private Web archives need to regulate access to their archives?

²<https://github.com/oduwsdl/ipwb>

³<https://github.com/odw/openwayback>

for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

JCDL '16 June 19-23, 2016, Newark, NJ, USA
© 2016 Copyright held by the owner/author(s).
ACM ISBN 978-1-4503-4229-2/16/06...\$15.00
DOI: <https://dx.doi.org/10.1145/2910896.2925467>

block in the URI-R is crawled, and a UUID to identify a memento. The two digests that are used to locate the contents of the memento are the “original URI” and the “payload digest”.

In designing ipwb, it was critical to consider the HTTP header returned at crawl time separately from the HTTP response body. The HTTP response header’s content will change with every capture, as the datetime returned from a server is temporally dependent. Compare this to the response body, which very often contains the same content on every replay. In the HTTP response body, the HTTP response body is pushed to IPFS, every IPFS hash would be unique, nullifying the potential for deduplication. In our design decision, ipwb only re-experiences the content of the response body, not the header, which is a design decision that is not considered the WARC request record upon replay. While including request records may be useful in the future (for

instance, to take into account the user-agent originally used to view the live website), WARC content is currently fully replayable without preserving the request records.

⁵http://crawler.archive.org/articles/user_manual/glossary.html#surt

⁶<https://www.w3.org/TR/uri-clarification/>

Preliminary Research

1. JCDL 2012 - WARCreate
2. TPDL 2013 - Change in Archivability
3. DLib 2013 - Method for Identifying
4. JCDL 2014 - Mink
5. JCDL 2014 - Archival Acid Test
6. JCDL 2014 - Not All Mementos
7. IJDL 2015 - Impact of JavaScript
8. IJDL 2015 - Not All Mementos
9. JCDL 2015 - Mobile Mink
10. JCDL 2016 - InterPlanetary Wayback (ipwb)
11. TPDL 2016 - ipwb extended
12. JCDL 2017 - WAIL Electron
13. JCDL 2017 - ServiceWorker Replay
14. JCDL 2017 - Impact of Canonicalization
15. JCDL 2018 - ArchiveNow
16. JCDL 2018 - Replay Banners
17. JCDL 2018 - A Framework...

A Framework for Aggregating Public and Private Web Archives
February 14, 2019
Mat Kelly

InterPlanetary Wayback: Peer-To-Peer Permanence 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, mln, 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 also 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 the HTTP response content of WARCs containing IPFS lookup hashes takes

How much does it cost to have
resilient personal archives?

1 Motivation

RQ4: How can content that was captured behind authentication signal to Web archive replay systems that it requires special handling?

RQ6: What kinds of access control do users who create private Web archives need to regulate access to their archives?

The Web ARChive (WARC) format is an ISO standard [4] to store live web archive content in a concatenated record-based file. IA's web crawler, Heritrix [7], generates WARC files to be read and the content re-experienced in an archival

¹ <https://github.com/oduwsdl/ipwb>

Preliminary Research

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

1. JCDL 2012 - WARCreate
2. TPDL 2013 - Change in Archivability
3. DLib 2013 - Method for Identifying
4. JCDL 2014 - Mink
5. JCDL 2014 - Archival Acid Test
6. JCDL 2014 - Not All Mementos
7. IJDL 2015 - Impact of JavaScript
8. IJDL 2015 - Not All Mementos
9. JCDL 2015 - Mobile Mink
10. JCDL 2016 - InterPlanetary Wayback (ipwb)
11. TPDL 2016 - ipwb extended
12. JCDL 2017 - WAIL Electron
13. JCDL 2017 - ServiceWorker Replay
14. JCDL 2017 - Impact of Canonicalization
15. JCDL 2018 - ArchiveNow
16. JCDL 2018 - Replay Banners
17. JCDL 2018 - A Framework...

A Framework for Aggregating Public and Private Web Archives
February 14, 2019
Mat Kelly

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. 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.XXX/XXXX

1 INTRODUCTION

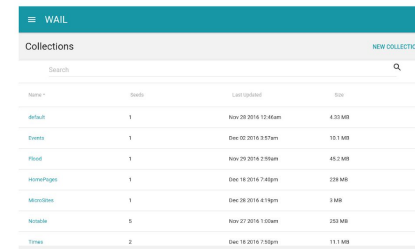
Subscription-based Web archiving services like Archive-it allow users with limited technical knowledge to create and replay personal collections of Web pages. However, these services are limited with a single seed and a single crawl. Similar to Archive-It is Webrecorder¹, which allows any user to create and manage personalized collections of Web archives. But unlike Archive-It, Webrecorder requires its user to manually drive the preservation process or upload content for replay while only providing its users up to five gigabytes of storage. Individuals that wish to freely (*gratis* and *libre*) 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

¹<https://webrecorder.io/>

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

JCDL'17, Toronto, Ontario, Canada

© 2017 Copyright held by the owner/author(s). XXX-YYYY-ZZ-AAA/BB/CC...\$15.00
DOI: 10.XXX/XXXX



Name	Seeds	Last Updated	Size
default	1	Nov 20 2016 13:04pm	6.53 MB
Events	1	Nov 20 2016 3:57pm	10.1 MB
Fixed	1	Nov 20 2016 7:04pm	45.2 MB
HeritrixPages	1	Nov 16 2016 7:40pm	228 MB
Microsites	1	Nov 16 2016 4:19pm	3 MB
Nicole	5	Nov 27 2016 1:04pm	253 MB
Times	2	Nov 16 2016 7:04pm	11.1 MB

Figure 1: Collections screen

Archive from the desktop
With higher fidelity than institutions

skill level. Even if a user is able to successfully set up these tools, their own means of associating the Web archives to each other for access to both Heritrix and Wayback while providing an interoperable mechanism for personal collection-based Web archiving from their personal computers. Users can create and add to these collections through WAIL-Election with the software taking care of the details in managing the collections, crawls, and replay. We have integrated a native Chromium² browser (the core of Google's Chrome Web browser) into the archival process in order to surface content specific to sites

RQ2: How do Web browser APIs compare in potential functionality to the capabilities of archival crawlers?

2. WAIL allows users to create and manage personalized collections of Web archives from their personal computers. When a user first starts the application, WAIL 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 and information about them. This information includes the number of seeds contained in the collection along with the collection's size and the last time it was updated. 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

²<https://www.chromium.org/>

³<http://electron.atom.io/>

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

Preliminary Research

1. JCDL 2012 - WARCreate
2. TPDL 2013 - Change in Archivability
3. DLib 2013 - Method for Identifying
4. JCDL 2014 - Mink
5. JCDL 2014 - Archival Acid Test
6. JCDL 2014 - Not All Mementos
7. IJDL 2015 - Impact of JavaScript
8. IJDL 2015 - Not All Mementos
9. JCDL 2015 - Mobile Mink
10. JCDL 2016 - InterPlanetary Wayback (ipwb)
11. TPDL 2016 - ipwb extended
12. JCDL 2017 - WAIL Electron
13. JCDL 2017 - ServiceWorker Replay
14. JCDL 2017 - Impact of Canonicalization
15. JCDL 2018 - ArchiveNow
16. JCDL 2018 - Replay Banners
17. JCDL 2018 - A Framework...

ABSTRACT

We use the ServiceWorker (SW) web API to intercept HTTP requests for embedded resources and reconstruct Composite Mementos without the need for conventional URL rewriting typically performed by web archives. URL rewriting is a problem for archival replay systems, especially for URLs constructed by JavaScript; frequently resulting in incorrect URI references. By intercepting requests on the client using SW, we are able to strategically reroute instead of rewrite. Our implementation moves rewriting to clients, saving servers' computing resources and allowing servers to return responses more quickly. Our experiments show that retrieving the original instead of rewritten pages from the archive reduces time overhead by 35.66% and data overhead by 19.68%. Our system prevents Composite Mementos from leaking the live web while being easy to distribute and maintain.

CCS CONCEPTS

•Information systems → World Wide Web;

KEYWORDS

ServiceWorker, Memento, Composite Memento, Web Archive, Archival Replay

ACM Reference format:

Sawood Alam, Mat Kelly, Michele C. Weigle, and Michael L. Nelson. 2019. Client-side Reconstruction of Composite Mementos Using ServiceWorker. In *Conference on Digital Libraries, Toronto, Ontario, Canada, June 2019*.

RQ2: How do Web browser APIs compare in potential functionality to the capabilities of archival crawlers?

1 INTRODUCTION

ServiceWorker (SW) is a new client-side web API [11] that can be used to intercept all the network requests for embedded resources originating from web pages in its scope. A Composite Memento [2] is an archived HTML page along

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

JCDL'17, Toronto, Ontario, Canada
© 2017 Copyright held by the owner/author(s). Publication rights licensed to ACM. 000-0000-00-00/00/00...\$15.00
DOI: 00.000/000.0



Figure 1: Live Ad Zombie Leaks into Archived Page. The image shows a screenshot of a CNN.com page from an archived version. It features a large article on the left, a sidebar with various links and images, and a prominent advertisement for Obama's campaign on the right. The page layout and content are typical of a news website from that era.

with all the embedded resources that are necessary to render the page correctly. Web archival replay systems rewrite embedded resources to point to their archival versions (e.g., a reference to `example.net/logo.png` is changed to `example.net/logo.png`).

We use SW API to reconstruct Composite Mementos from the originally captured data without any such URL rewriting. By intercepting requests on the client-side we are essentially rerouting instead of rewriting. Rerouting is an effective mechanism to block live web leakage, or “zombies” that might happen after executing potential JavaScript (JS), otherwise not discoverable by static analysis. For example, in Figure 1

the page was archived on September 16th, 2008, but when the page is replayed, the ad from the time of the 2012 presidential candidates [5]. Client-side rerouting also saves the content on the client side, such as to include archival banners, hence, there is no need to send extra data with each HTTP response. Client-side solutions such as Memento for Chrome¹ involve installing a browser add-on, which limits the adoption by users and adds the burden of maintaining the add-on while only available for Google Chrome users. Our exploratory technique works well when SW is supported. However, a server-side fallback is necessary for production usage to avoid the risk of zombies and broken references when SW is not supported.

Our experiments show that retrieving the original instead of rewritten pages from the Internet Archive (IA) reduces time overhead by 35.66% and data overhead by 19.68%. Our system prevents Composite Mementos from zombies while being easy to distribute and maintain. It is a lightweight and portable system that can be used with any Memento server such as a web archive or a Memento aggregator.

¹<http://bit.ly/memento-for-chrome>

Preliminary Research

1. JCDL 2012 - WARCreate
2. TPDL 2013 - Change in Archivability
3. DLib 2013 - Method for Identifying
4. JCDL 2014 - Mink
5. JCDL 2014 - Archival Acid Test
6. JCDL 2014 - Not All Mementos
7. IJDL 2015 - Impact of JavaScript
8. IJDL 2015 - Not All Mementos
9. JCDL 2015 - Mobile Mink
10. JCDL 2016 - InterPlanetary Wayback (ipwb)
11. TPDL 2016 - ipwb extended
12. JCDL 2017 - WAIL Electron
13. JCDL 2017 - ServiceWorker Replay
14. JCDL 2017 - Impact of Canonicalization
15. JCDL 2018 - ArchiveNow
16. JCDL 2018 - Replay Banners
17. JCDL 2018 - A Framework...

A Framework for Aggregating Public and Private Web Archives
February 14, 2019
Mat Kelly

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,mh,mweigle}@cs.odu.edu

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

ABSTRACT

Memento TimeMaps [5] list identifiers for archival web captures (URI-Ms). When some URI-Ms are dereferenced, they redirect to a different URI-M instead of a unique representation at the datetime. This suggests that confidently obtaining an accurate count quantifying the number of non-forwarding captures for an 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 13 academic institutions. Using a ratio of the number of URI-Ms that result in a redirect to the number of URI-Ms in the TimeMap, we found that the 'large web sites' and two of the thirteen 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 the HTTP transactions observed at archival time. TimeMaps have generally been used as a count of the number of representations of a URI-R. However, some URI-Ms may include URI-Ms for archived representations that redirect another URI-M in the TimeMap that returns a HTTP Status OK.

TimeMaps do not explicitly return a "count" value to indicate the number of mementos listed in the TimeMap that produce a non-redirecting HTTP status code when dereferenced. The heuristic of determining how many captures are represented by URI-Ms in a TimeMap cannot be completed without

Redirection rules can be attributed to the Web architecture of canonicalization rules [3]. Preserving and replaying these redirection rules 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 relation values 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

year	M_{TM}	M_{RC}	DI
2006	735	483	1.917
2007	1,055	842	3.953
2008	1,376	894	1.855
2009	6,074	4,335	2.493
2010	9,326	6,530	2.335
2011	20,634	9,279	0.817
2012	102,533	16,240	0.188
2013	228,405	25,203	0.124
2014	164,865	22,738	0.160
2015	17,978	11,286	1.686
2016	120,520	5,505	0.049

Table 1: Google over time (abbreviated), bucketed by year, based on IA mementos extracted from the TimeMap. M_{TM} is the memento count based solely on the data in the TimeMap, M_{RC} is the count based on aclusion of redirects when dereferenced, and DI is the ratio of M_{RC} to M_{TM} .

URI coalescence considered harmful for archives

(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.

Based on the number of URI-Ms in a TimeMap not necessarily resolving to unique mementos when archival redirects are followed, we examined the mementos from contemporary TimeMaps to determine the extent to which they identify the difference between the number of mementos available as reported by the TimeMap through naive "rel" counting heuristics to the temporally unique mementos identified once these mementos are dereferenced.

2 BACKGROUND AND RELATED WORK

URI canonicalization associates differently formatted URIs [4] and allows for the fact clustering of URIs that likely reference the same resources. As URI schemes from a Web archive can be attributed to the Web architecture of canonicalization rules [3]. Preserving and replaying these redirection rules 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 relation values 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

BEST POSTER AWARD
at JCDL 2017

1. JCDL 2012 - WARCreate
2. TPDL 2013 - Change in Archivability
3. DLib 2013 - Method for Identifying
4. JCDL 2014 - Mink
5. JCDL 2014 - Archival Acid Test
6. JCDL 2014 - Not All Mementos
7. IJDL 2015 - Impact of JavaScript
8. IJDL 2015 - Not All Mementos
9. JCDL 2015 - Mobile Mink
10. JCDL 2016 - InterPlanetary Wayback (ipwb)
11. TPDL 2016 - ipwb extended
12. JCDL 2017 - WAIL Electron
13. JCDL 2017 - ServiceWorker Replay
14. JCDL 2017 - Impact of Canonicalization
15. **JCDL 2018 - ArchiveNow**
16. JCDL 2018 - Replay Banners
17. JCDL 2018 - A Framework...

ArchiveNow: Simplified, Extensible, Multi-Archive Preservation

Mohamed Aturban, Mat Kelly, Sawood Alam, John A. Berlin,
Michael L. Nelson, and Michele C. Weigle
Old Dominion University
Department of Computer Science
Norfolk, Virginia, USA
{maturban,mkelly,salam,jberlin,mln,mweigle}@cs.odu.edu

ArchiveNow is a Python module for preserving web pages in on-demand web archives. This module allows a user to submit a URI of a web page for archiving at several configured web archives. Once the web page is captured, *ArchiveNow* provides the user with links to the archived copies of the web page. *ArchiveNow* is initially configured to use four archives but is easily configurable to add or remove other archives. In addition to pushing web pages to public archives, *ArchiveNow*, through the use of *Wget* and *Squidward*, allows users to generate local WARC files, enabling them to create their own personal and private archives.

```
% archivenow -all --ccapi_key=7e..3f http://money.cnn.com/2018/01/27/technology/future/spacex-falcon-heavy-everything-you-need-to-know/index.html
{
  "uri": "http://money.cnn.com/2018/01/27/technology/future/spacex-falcon-heavy-everything-you-need-to-know/index.html",
  "request_datetime": "20180129094723",
  "mentos": {
    "archive.org": "https://web.archive.org/web/20180129094728/http://money.cnn.com/2018/01/27/technology/future/spacex-falcon-heavy-everything-you-need-to-know/index.html",
    "archive.is": "https://archive.is/hr41is",
    "archive.org": "https://web.archive.org/web/20180129094728/http://money.cnn.com/2018/01/27/technology/future/spacex-falcon-heavy-everything-you-need-to-know/index.html",
    "archive.org": "https://web.archive.org/web/20180129094728/http://money.cnn.com/2018/01/27/technology/future/spacex-falcon-heavy-everything-you-need-to-know/index.html"
  }
}
```

• Information systems → Digital libraries and archives; World Wide Web: **Create & Submit**

Web Archiving, Memento, WARC and local W

ACM Reference Format:
Mohamed Ataburn, Mat Kelly, Sawood Alam, John A. Berlin, Michael L. Nelson, and Michele C. Weale. 2018. ArchiveNow: Simplified, Extensible, Multi-Archive Preservation. In *JCDL '18: The 18th ACM/IEEE Joint Conference on Digital Libraries*, June 3–7, 2018, Fort Worth, TX, USA. ACM, New York, NY, USA, 2 pages. <https://doi.org/10.1145/3197026.3203880>

RQ1: What sort of content is difficult to capture and replay for preservation from the perspective of a Web browser?

RQ1: What sort of content is preserved from the perspective of the user?

that notifies a user of any available archived copies for a viewed page and suggests to archive the page in three archives. Welsh [10] developed several tools intended to archive news-related resources. For example, Welsh's *Savemynews* (www.savemynews.com) saves web pages in two archives. Users of this service are required to create accounts. In addition to



Permission is granted to reproduce copies of part or all of this work for personal or classroom use without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Requests for reproduction for other than personal or classroom use should be directed to the owner/author(s).

JCDL '18, June 3-7, 2018, Fort Worth, TX, USA
© 2018 Copyright held by the owner/author(s).
ACM ISBN 978-1-4503-5178-2/18/06.
<https://doi.org/10.1145/3197026.3203880>

BEST POSTER AWARD at JCDL 2018

Difficult to capture and replace content of a Web browser?

Preliminary Research

1. JCDL 2012 - WARCreate
2. TPDL 2013 - Change in Archivability
3. DLib 2013 - Method for Identifying
4. JCDL 2014 - Mink
5. JCDL 2014 - Archival Acid Test
6. JCDL 2014 - Not All Mementos
7. IJDL 2015 - Impact of JavaScript
8. IJDL 2015 - Not All Mementos
9. JCDL 2015 - Mobile Mink
10. JCDL 2016 - InterPlanetary Wayback (ipwb)
11. TPDL 2016 - ipwb extended
12. JCDL 2017 - WAIL Electron
13. JCDL 2017 - ServiceWorker Replay
14. JCDL 2017 - Impact of Canonicalization
15. JCDL 2018 - ArchiveNow
16. JCDL 2018 - Replay Banners
17. JCDL 2018 - A Framework...

Unobtrusive and Extensible Archival Replay Banners Using Custom Elements

Sawood Alam, Mat Kelly, Michele C. Weigle, and Michael L. Nelson

Old Dominion University
Department of Computer Science
Norfolk, Virginia, USA

{salam,mkelly,mweigle,mln}@cs.odu.edu

ABSTRACT

We compare and contrast three different ways to implement an archival replay banner. We propose an implementation that utilizes *Custom Elements* and adds some unique behaviors, not common in existing archival replay systems, to enhance the user experience. Our approach has a minimal user interface footprint and resource overhead while still providing rich interactivity and extended on-demand provenance information about the archived resources.

CCS CONCEPTS

• Information systems → Digital libraries and archives; • Human-centered computing → User interface design;

KEYWORDS

Memento; Archival; Custom Elements; Archival Replay

ACM Reference Format:

Sawood Alam, Mat Kelly, Michele C. Weigle, and Michael L. Nelson. 2018. Unobtrusive and Extensible Archival Replay Banners Using Custom Elements. In *JCDL '18: The 18th ACM/IEEE Joint Conference on Digital Libraries, June 3–7, 2018, Fort Worth, TX, USA*. ACM, New York, NY, USA, 2 pages. <https://doi.org/10.1145/3197026.3203881>

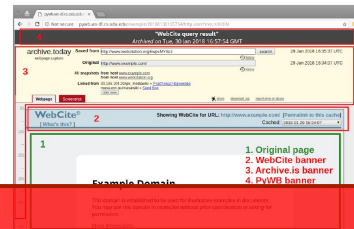
1 MOTIVATION

Web archival replay systems express that a user is interacting with a *memento* (an archived representation of a resource) by adding an archival banner. Archival banners provide metadata about both the *memento* and the original resource. They also allow users to interactively replicate the live web experience when viewing a *memento*, by providing a visual representation of the original resource and not the live web. Any component injection in the page makes it different from the original and may consume additional screen real estate. We illustrate this in Figure 1 by archiving example.com in Memento.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored.

For all other uses, contact the owner/author(s).
JCDL '18, June 3–7, 2018, Fort Worth, TX, USA
© 2018 Copyright held by the owner/author(s).

ACM ISBN 978-1-4503-5178-2/18/06.
<https://doi.org/10.1145/3197026.3203881>



Enable archival navigation on replay to be more extensible

by using a custom test, and allows drag-and-drop repositioning of the banner itself when not needed. In the on-demand extended mode it provides a set of interactive visualizations and provenance information that are customizable by the archive.

2 METHODOLOGY

There are three primary ways to serve an archival banner with an archived web page that shares the rendering space with the *memento*. Browser toolbars (e.g., the now defunct *MementoFox*) and archival emulators (e.g., *NetCapsule*) are out of scope of this work.

Inline Plain HTML Banners – This is the simplest and most common method. It involves inserting the banner directly into the archived *HTML*. While simple, it poses some serious issues, such as conflicts with the style of the *memento* or hiding important elements of the page, such as the header of the site.

IFrame Banners – This method involves creating a separate *HTML* element for the banner, which is then loaded into an *iframe* within the *memento* page, or 2) making the banner document as the outer page and serving the *memento* inside an *iframe*. For example, *WebCite* uses the first approach while many archives, such as the *Portuguese Web Archive*, use the second.

PyWB, a popular web archival replay system, uses the second approach by default, but allows using plain inline *HTML* banners. *Iframes* provide full document isolation, both style and origin. Therefore, *iframe* banners do not conflict with the position

Preliminary Research

1. JCDL 2012 - WARCreate
2. TPDL 2013 - Change in Archivability
3. DLib 2013 - Method for Identifying
4. JCDL 2014 - Mink
5. JCDL 2014 - Archival Acid Test
6. JCDL 2014 - Not All Mementos
7. IJDL 2015 - Impact of JavaScript
8. IJDL 2015 - Not All Mementos
9. JCDL 2015 - Mobile Mink
10. JCDL 2016 - InterPlanetary Wayback (ipwb)
11. TPDL 2016 - ipwb extended
12. JCDL 2017 - WAIL Electron
13. JCDL 2017 - ServiceWorker Replay
14. JCDL 2017 - Impact of Canonicalization
15. JCDL 2018 - ArchiveNow
16. JCDL 2018 - Replay Banners
17. JCDL 2018 - A Framework...



BEST PAPER AWARD FINALIST
at JCDL 2018

A Framework for Aggregating Private and Public Web Archives

Mat Kelly
Old Dominion University
Norfolk, Virginia, USA
mkelly@cs.odu.edu

Michael L. Nelson
Old Dominion University
Norfolk, Virginia, USA
mln@cs.odu.edu

Michele C. Weigle
Old Dominion University
Norfolk, Virginia, USA
mweigle@cs.odu.edu

ABSTRACT

Personal and private Web archives are proliferating due to the increase in the tools to create them and the realization that Internet Archive and other public Web archives are unable to capture personalized (e.g., Facebook) and private (e.g., banking) Web pages. We introduce a framework to mitigate issues of aggregation in private, personal, and public Web archives without compromising potentially sensitive information contained in private captures. We amend Memento syntax and semantics to allow TimeMap enrichment to incorporate additional attributes to be expressed in private or personal captures. We then describe a framework that can be used to provide a method to give the user further in the negotiation of an archival captures in dimensions beyond time. We introduce a framework for archival querying precedence and short-circuiting of captures when aggregating private and personal Web archive captures with those from public Web archives through Memento. Negotiation of this sort is novel to Web archiving and allows for the more seamless

inappropriate (e.g., requires a specific user's credentials) for these crawlers and systems to preserve. For this reason and enabled by the recent influx of personal Web archiving tools, such as WARCreate, WAIL, and Webrecorder.io, individuals are preserving live Web content and personal Web archives are proliferating [20].

Personal and private captures, or mementos, of the Web, particularly those preserving content that requires authentication on the live Web, have potential privacy ramifications if shared or made publicly accessible after being preserved [22]. Given the privacy issues, strategically regulating access to these personal and private Web archives is a necessary and prudent way to address and mitigate privacy considerations to the aggregate view of the Web. Web archives provide a more comprehensive picture of the Web while mitigating privacy violations.

This work has four primary contributions to Web archiving:

Archival Query Precedence and Short-circuiting: Allow

Storing and replaying content

TimeMap Link Enrichment Provides additional, more descriptive attributes to URIs for more efficient querying and interaction (Section 4).

Multidimensional navigation Facilitates the negotiation of content captured behind URIs for URIs in both temporal and other dimensions (Sections 5 and 6.1).

Archive replay systems that i

Memento aggregation using OAuth (Section 6.2).

1.1 Solutions Beyond Institutions

regulators indicate that private information, such as time sensitive verification

special handling to be (Figure 1a)

mementos

dated with publicly available

archives simply preserve the facebook.com again page (Figure 1b).

Both captures are representative of facebook.com, and they may have even been captured at the same time. Users may be hesitant to share their mementos of facebook.com for other personal or private control do users who create as

those archives can be restricted

regulate access to their

susceptible to disappearing without an institution's backing. Maintaining backups of archived content is unwieldy, requires diligence or automation, and is still at the mercy of hardware failures. While

Preliminary research aggregating private and public Web archives

RQ3: What issues exist for capturing and replaying content behind authentication?

RQ4: How can content that was captured behind authentication signal to Web archive replay systems that it requires special handling?

RQ5: How can Memento aggregators indicate that private Web archive content requires special handling to be replayed, despite being aggregated with publicly available Web archive content?

RQ6: What kinds of access control do users who create private Web archives need to regulate access to their archives?

Outline

- Introduction/Motivation
- Background
- Preliminary Research
- **Proposed Framework**
- Evaluation Plan

Proposed Framework

(for aggregating private and public Web archives)

Proposed Framework

- Archival negotiation beyond time
- Query precedence & short-circuiting
- Mementities

PROPOSED FRAMEWORK

Archival Negotiation Beyond Time

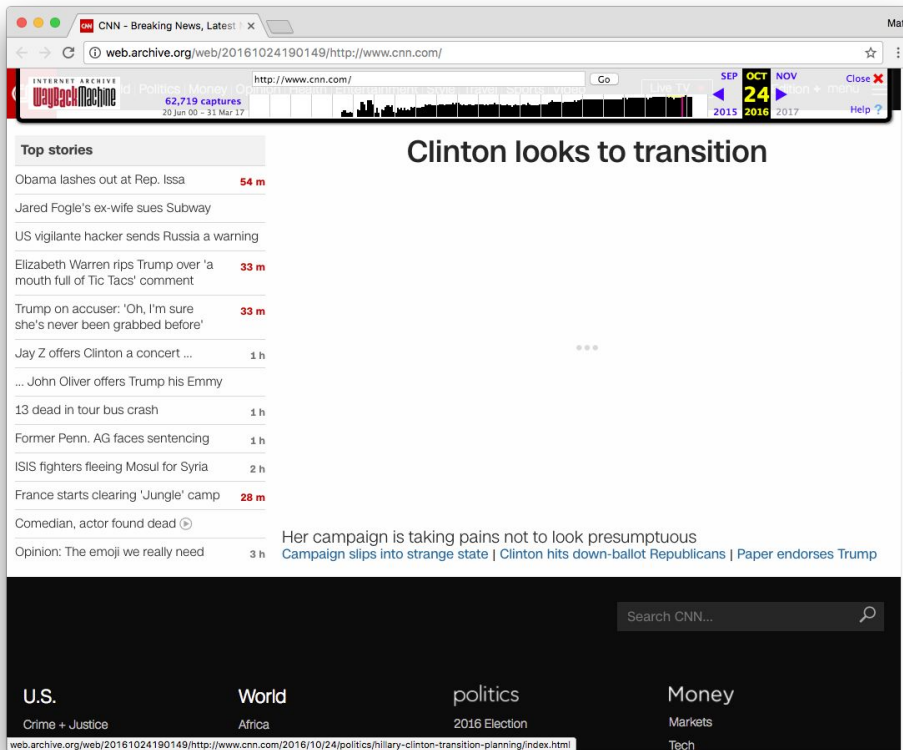
More Expressive TimeMaps

- Memento Quality (e.g., Damage)¹
- How Many Captures?²
- How Many Are Identical?^{2,3}
- Other Attributes of Mementos...

¹ Brunelle *et al.*, JCDL 2014, IJDL 2015

² Kelly *et al.*, JCDL 2017

³ AISum and Nelson, ECIR 2014



Additional TimeMap Attributes

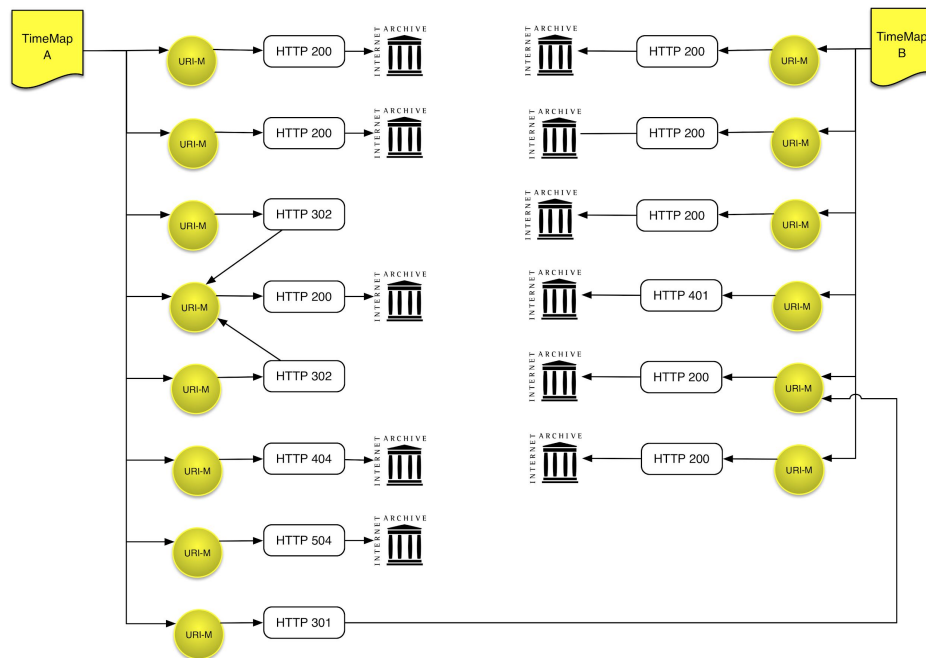
Content-based Attributes

Derived Attributes

Access Attributes

TimeMap Enrichment: Content-Based Attributes

- Status Code¹
- Content-Digest
 - In WARC & CDX
 - Not all archives expose CDX
- Would allow more info about mementos without requiring comprehensive dereferencing



¹ Kelly *et al.*, “Impact of URI Canonicalization on Memento Count”, JCDL 2017, arXiv 1703.03302

TimeMap Enrichment: Derived Attributes

- Thumbnails (e.g, via SimHash)¹
 - Calculation based on root memento's HTML
- Memento Damage (JCDL 2014, IJDL)²
 - Requires dereferencing embedded resources



apple.com, many duplicate mementos!

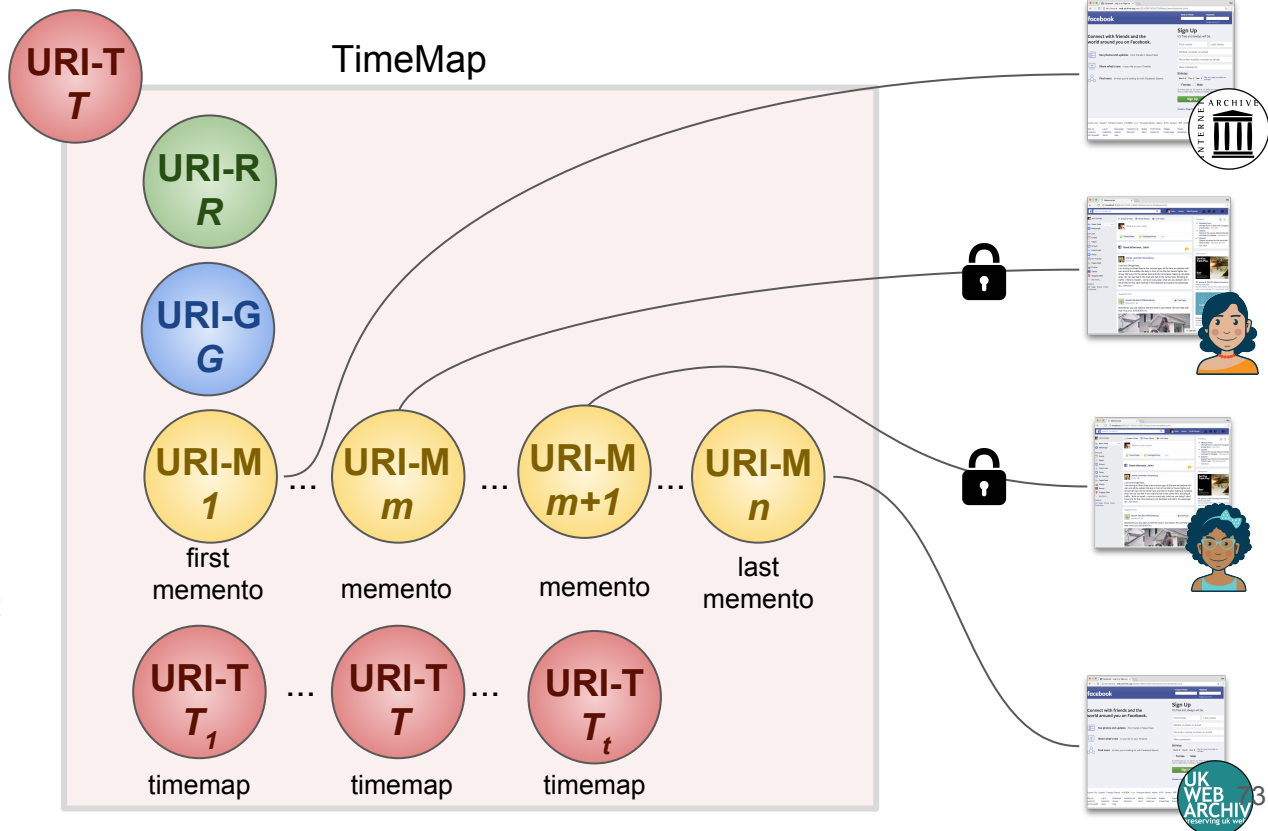
¹ AISum and Nelson, Thumbnail Summarization Techniques for Web Archives, ECIR 2014, pp. 299-310.

² Brunelle *et al.*, “The Impact of JavaScript on Archivability,” IJDL, 17(2), pp. 95-117. January 2016.

TimeMap Enrichment: Access Attributes

How to distinguish
Private captures
from
Public captures
in a TimeMap?

RQ5: How can Memento aggregators indicate that private Web archive content requires special handling to be replayed, despite being aggregated with publicly available Web archive content?



TimeMap Enrichment - in a CDXJ TimeMap

Line breaks added for clarity, CDXJ records occupy a single line

```
19981212013921 {  
  "uri": "http://localhost:8080/20101116060516/http://facebook.com/",  
  "rel": "memento",  
  "datetime": "Tue, 16 Nov 2010 06:05:16 GMT",  
  "status_code": 200,  
  "digest": "sha1:LK26DRRQJ4WATC6LBVF3B3Z4P2CP5ZZ7",  
  "damage": 0.24,  
  "simhash": "6551110622422153488",  
  "content-language": "en-US",  
  "access": {  
    "type": "Blake2b",  
    "token": "c6ed419e74907d220c69858614d86...ef0a3a88a41"  
  }  
}
```

Content-based attributes

Derived Attributes

Access Attributes

TimeMap
+
Enrichment with Additional Attributes

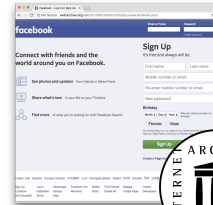
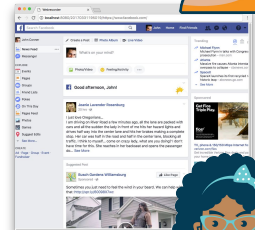
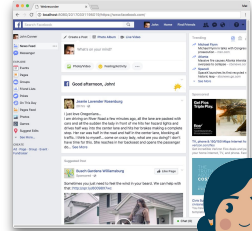
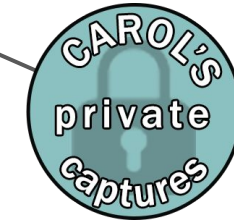
“StarMap”

PROPOSED FRAMEWORK

Query Precedence
- and -
Short Circuiting

Query Precedence

- More control of querying in series and parallel



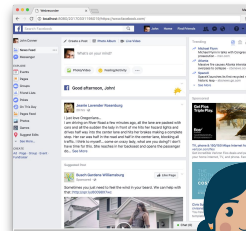
“Check my archive first, then Carol’s, then all public archives.”

Query Precedence

- More control of querying in series and parallel



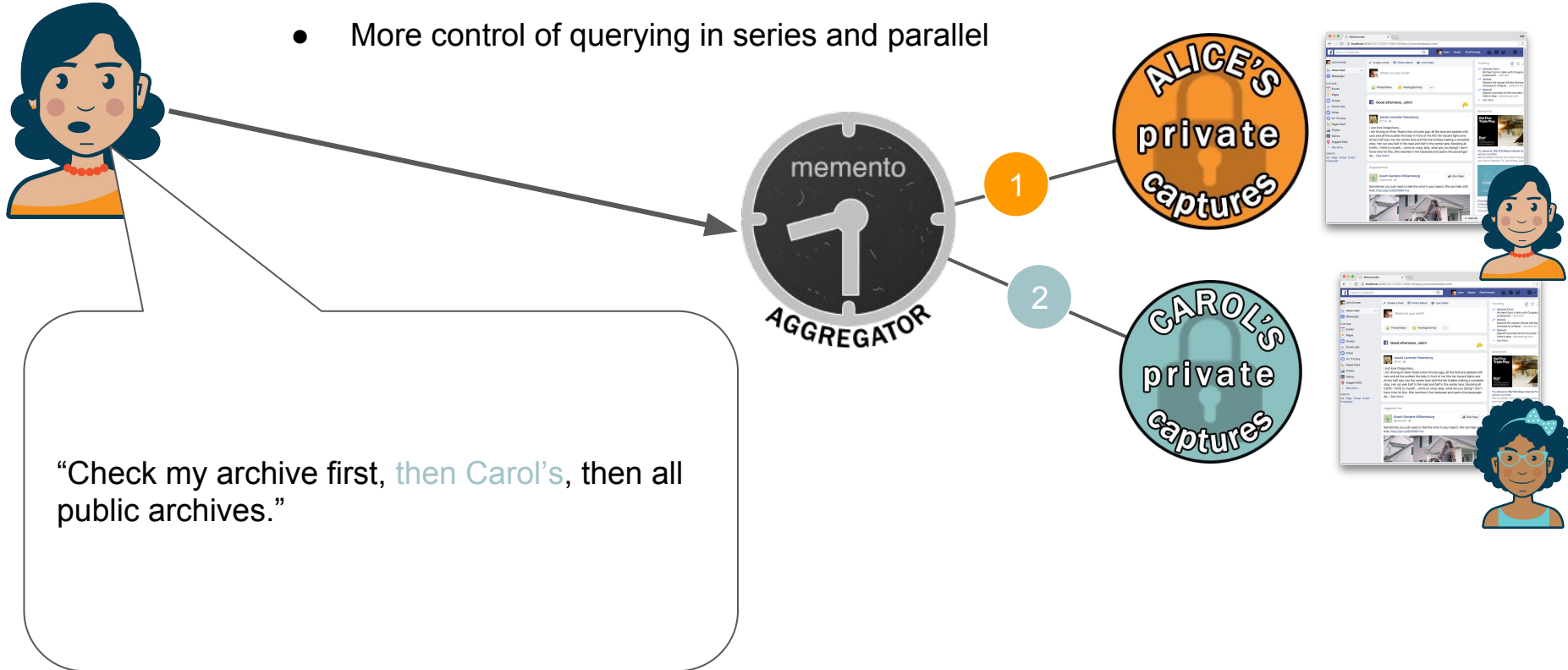
1



“Check **my archive first**, then Carol’s, then all public archives.”

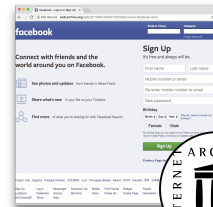
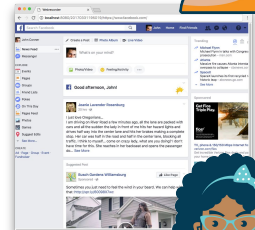
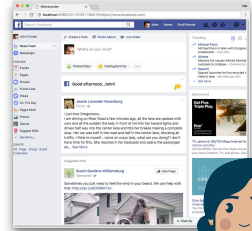
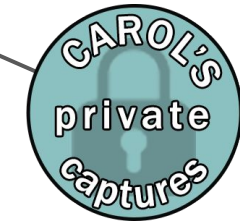
Query Precedence

- More control of querying in series and parallel



Query Precedence

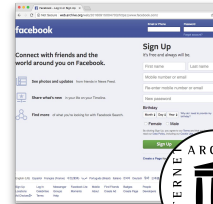
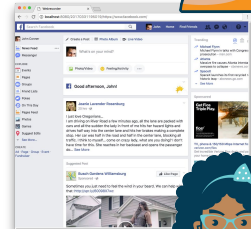
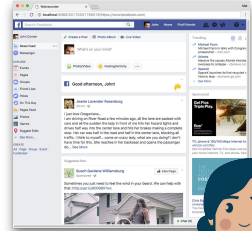
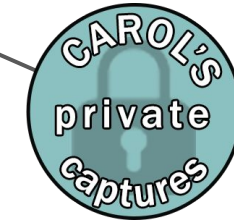
- More control of querying in series and parallel



“Check my archive first, then Carol’s, then all public archives.”

Query Short-Circuiting

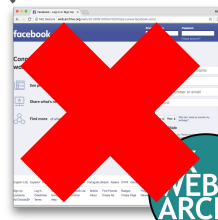
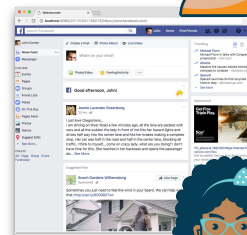
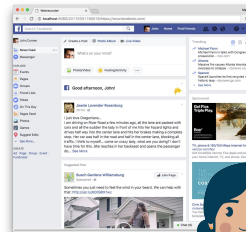
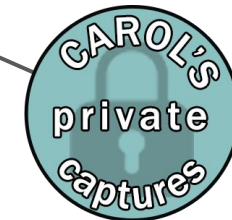
- May give priority to archive relevancy.
- Series halt when threshold met.



“Check private archives first. **Iff** you find no captures, only *then* check the public archives.

Query Short-Circuiting

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



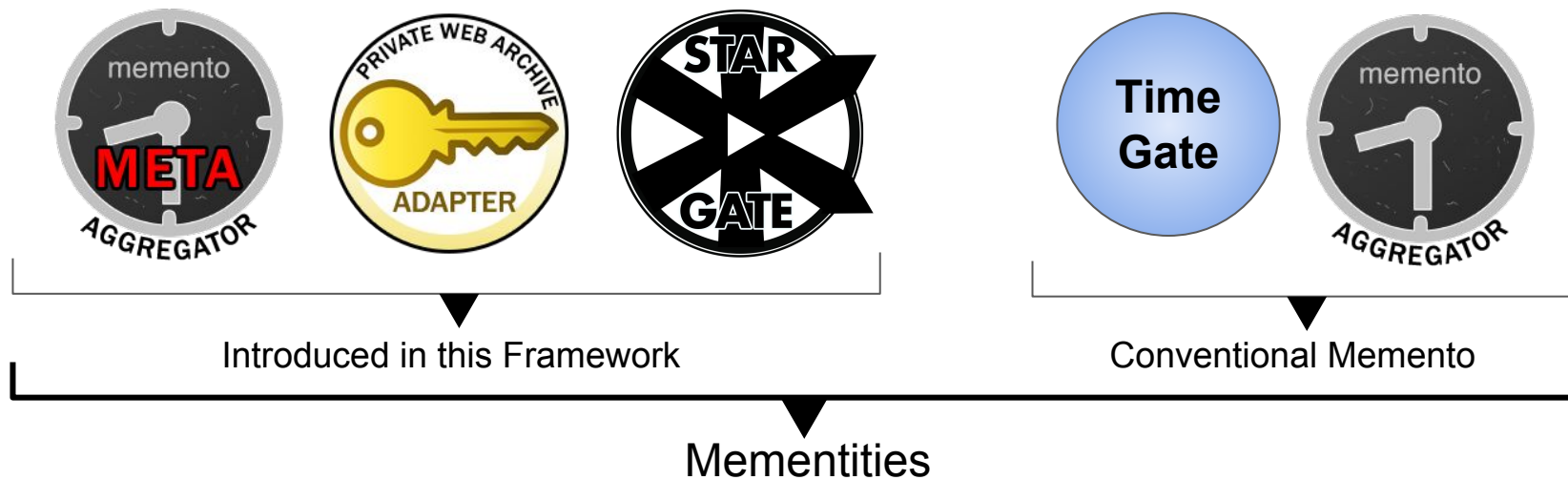
“Check private archives first. **Iff** you find no captures, only *then* check the public archives.

PROPOSED FRAMEWORK

Mementies

Mementities

- Memento + Entity (*entity* term already overused)

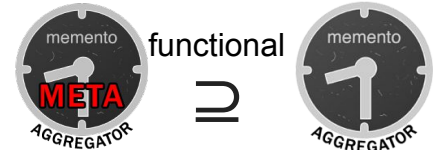
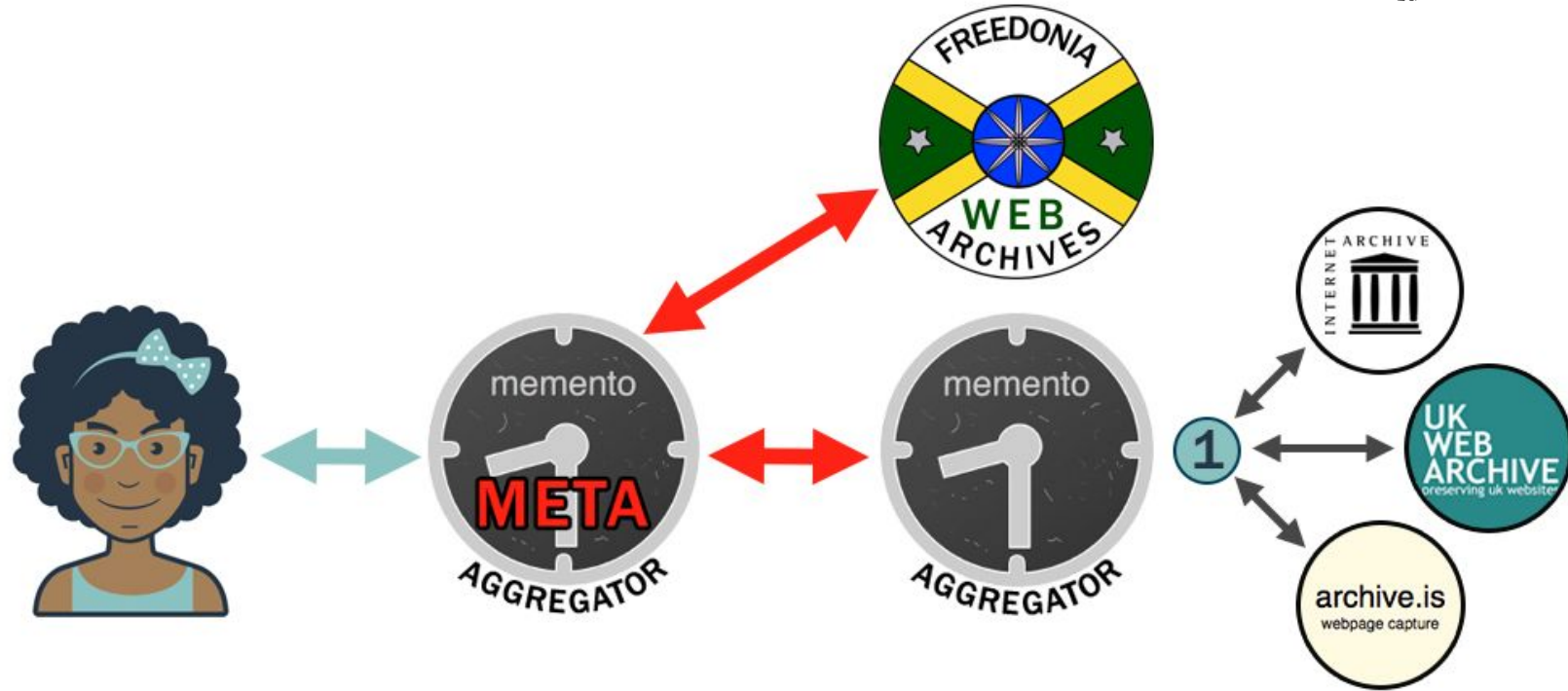


PROPOSED FRAMEWORK

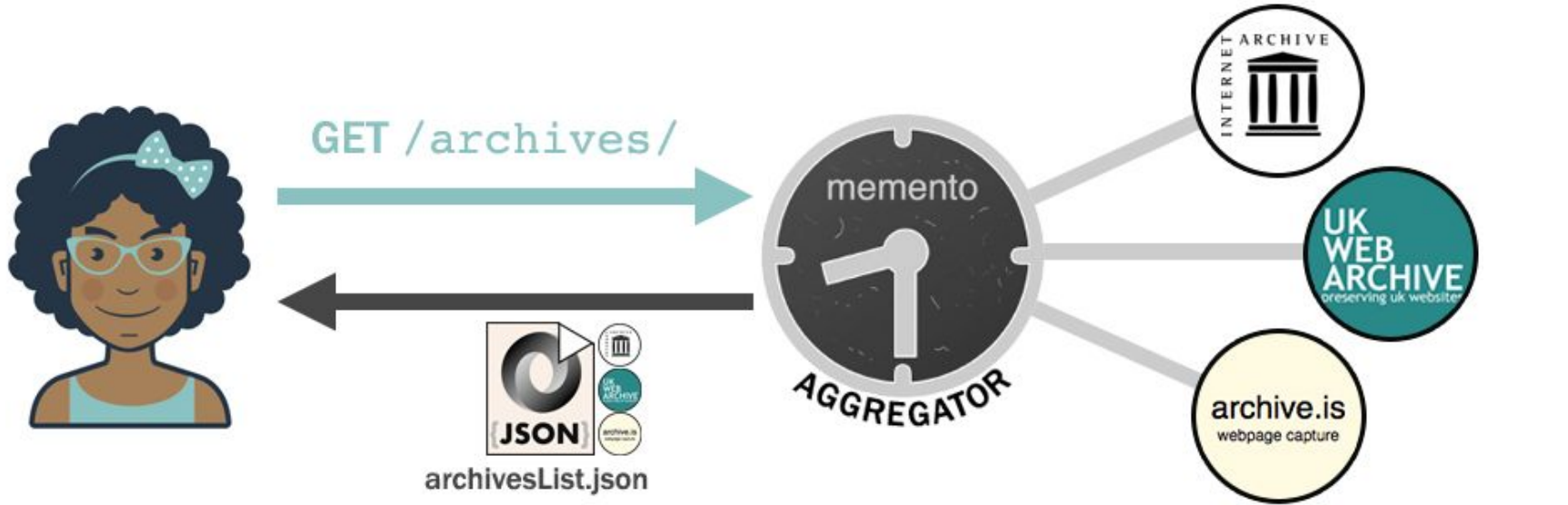
Mementities



Memento Meta-Aggregator (MMA)



MMA: Archive Selection



MMA: User-Driven Archival Specification



MMA Aggregation sources

MMA_{α} :

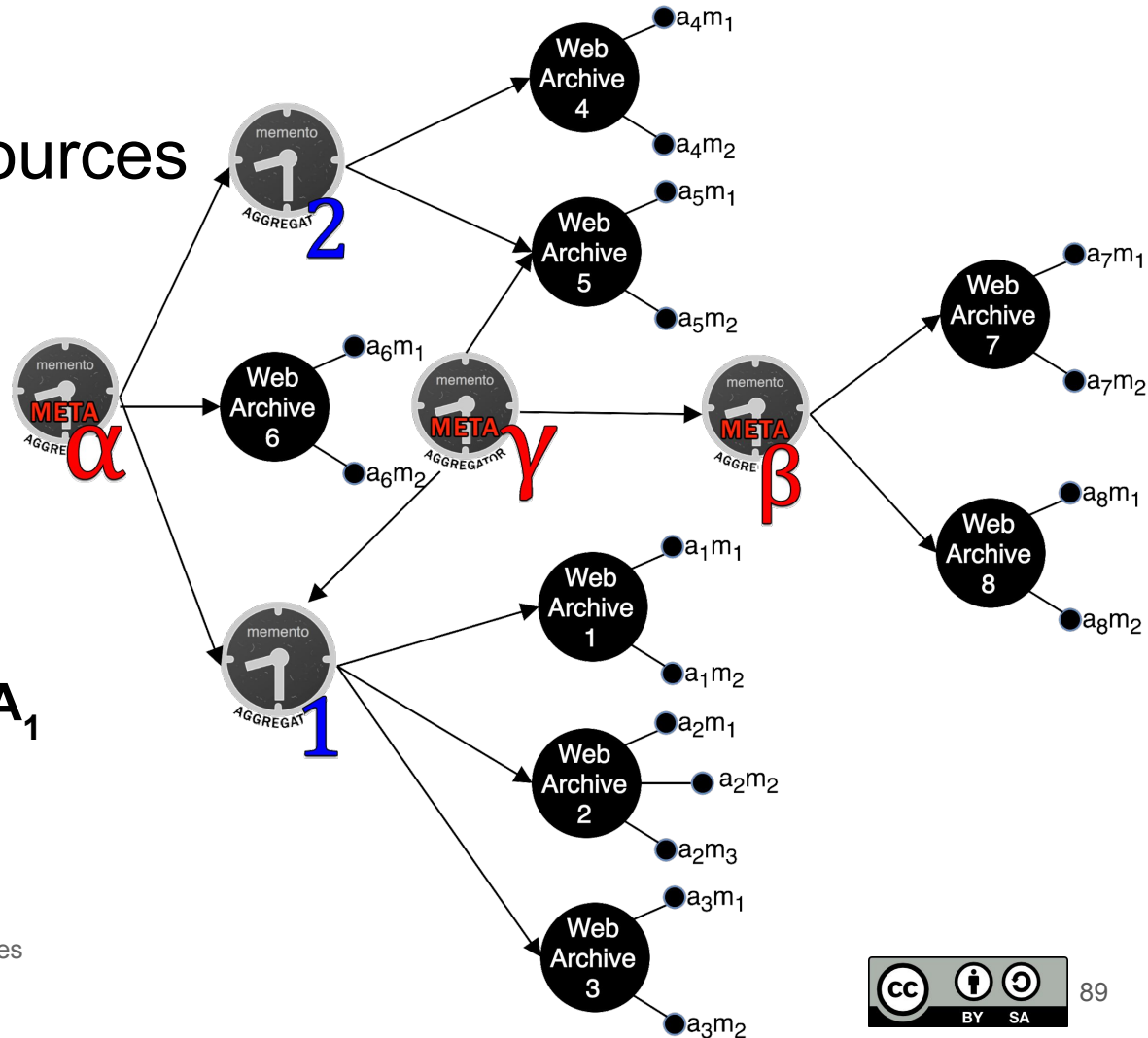
from MA_2 , MA_1 and WA_6

MMA_{β} :

from WA_7 and WA_8

MMA_{γ} :

from MMA_{β} , MA_5 , and WA_1

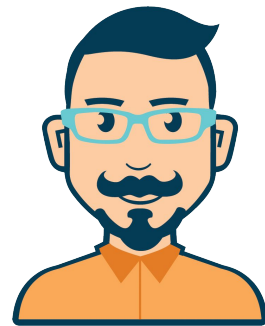


MMA Dynamics By-Example

- Personal Archive Aggregation
- MMA Chaining
- Client-Side Aggregation Preference



ALICE



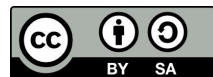
BOB



CAROL



MALCOLM



MMA Dynamics - Personal Archive Aggregation



bbc

homepage

Public videos

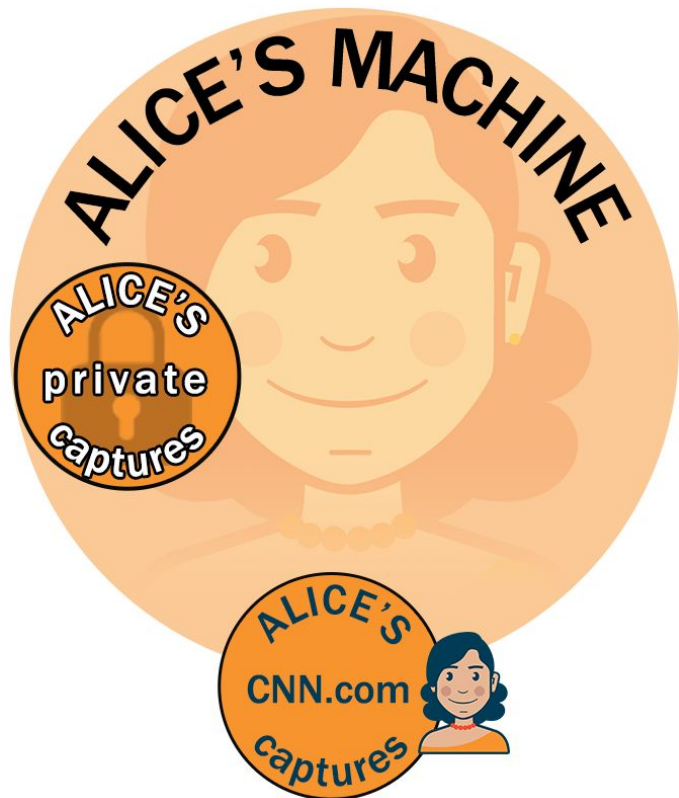


FB

bank

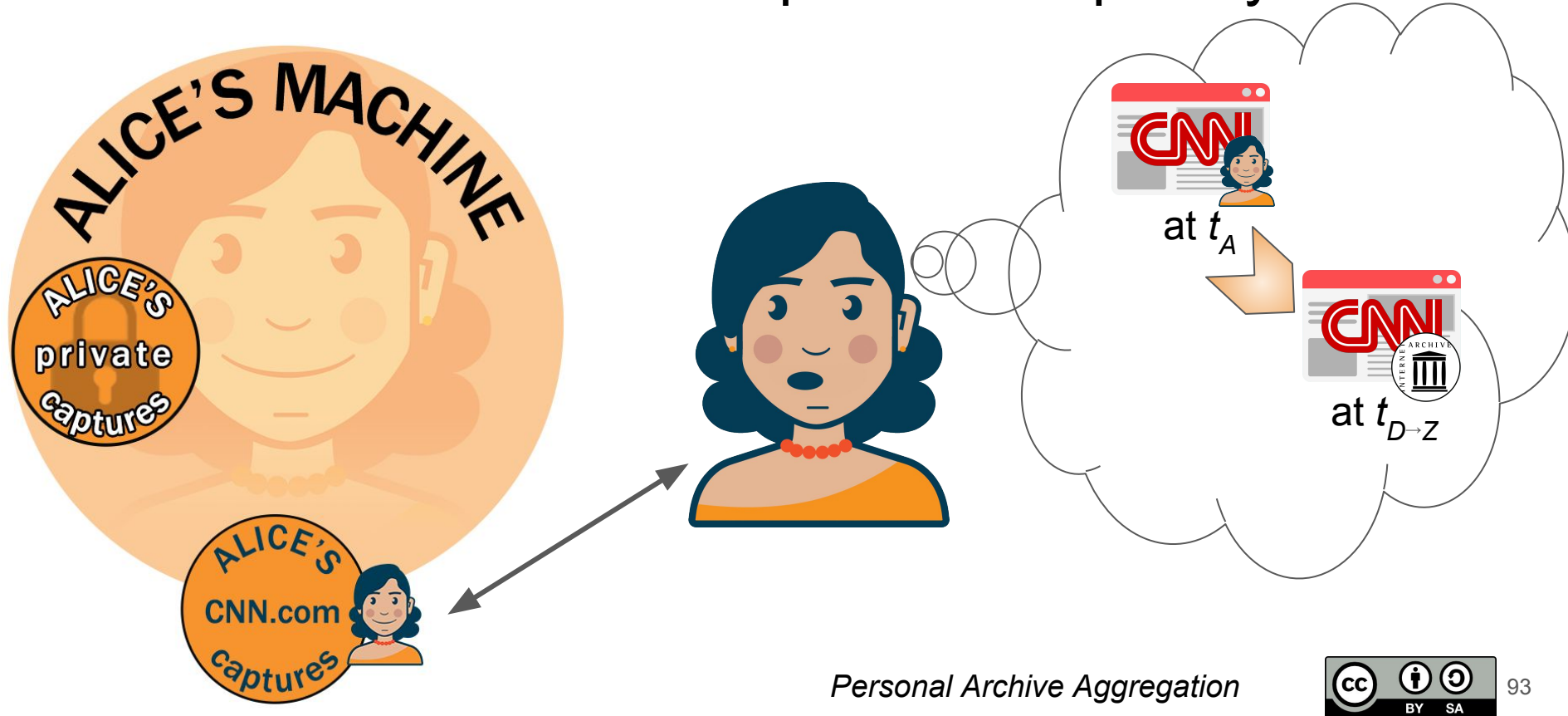
flickr

Alice Saves the Web

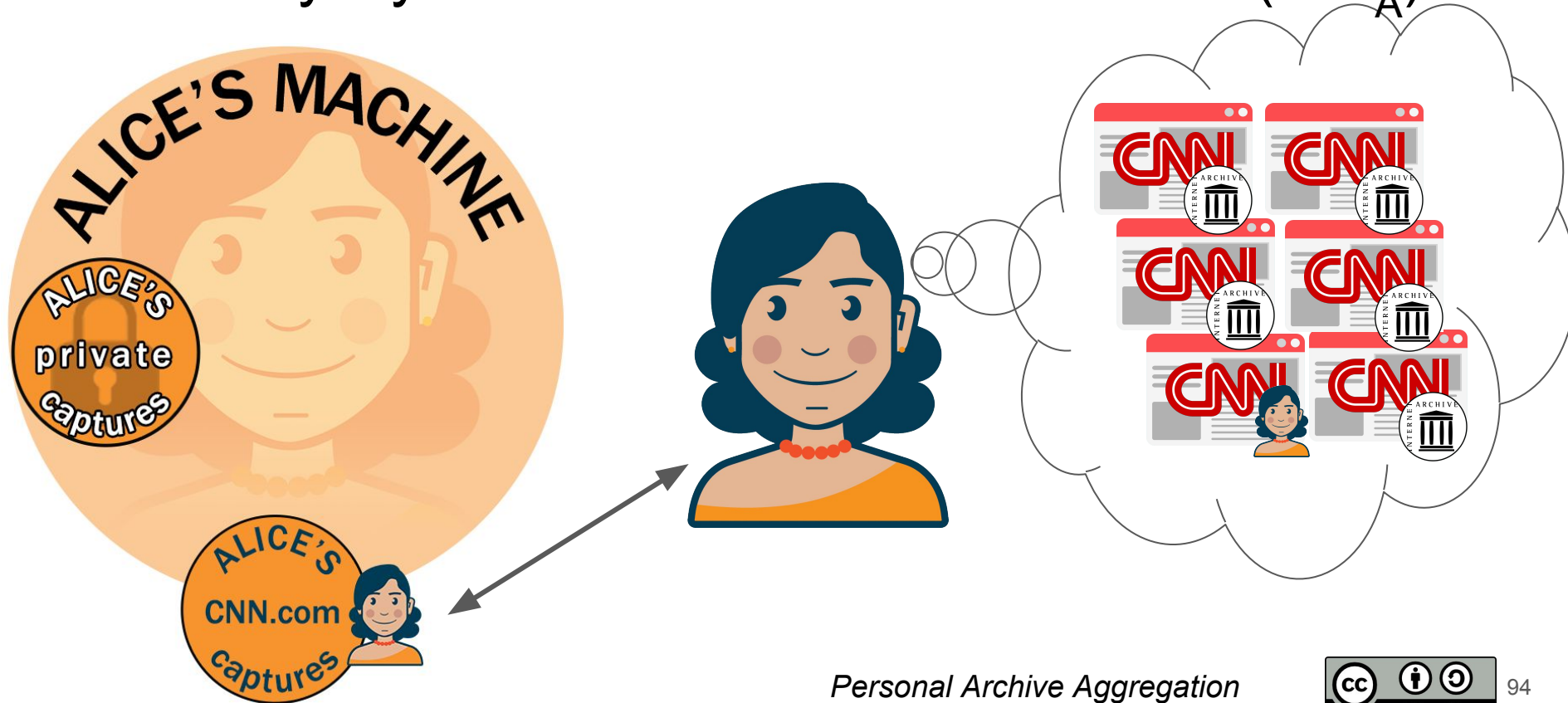


Personal Archive Aggregation

Alice Wants to See Her Captures Temporally Inline

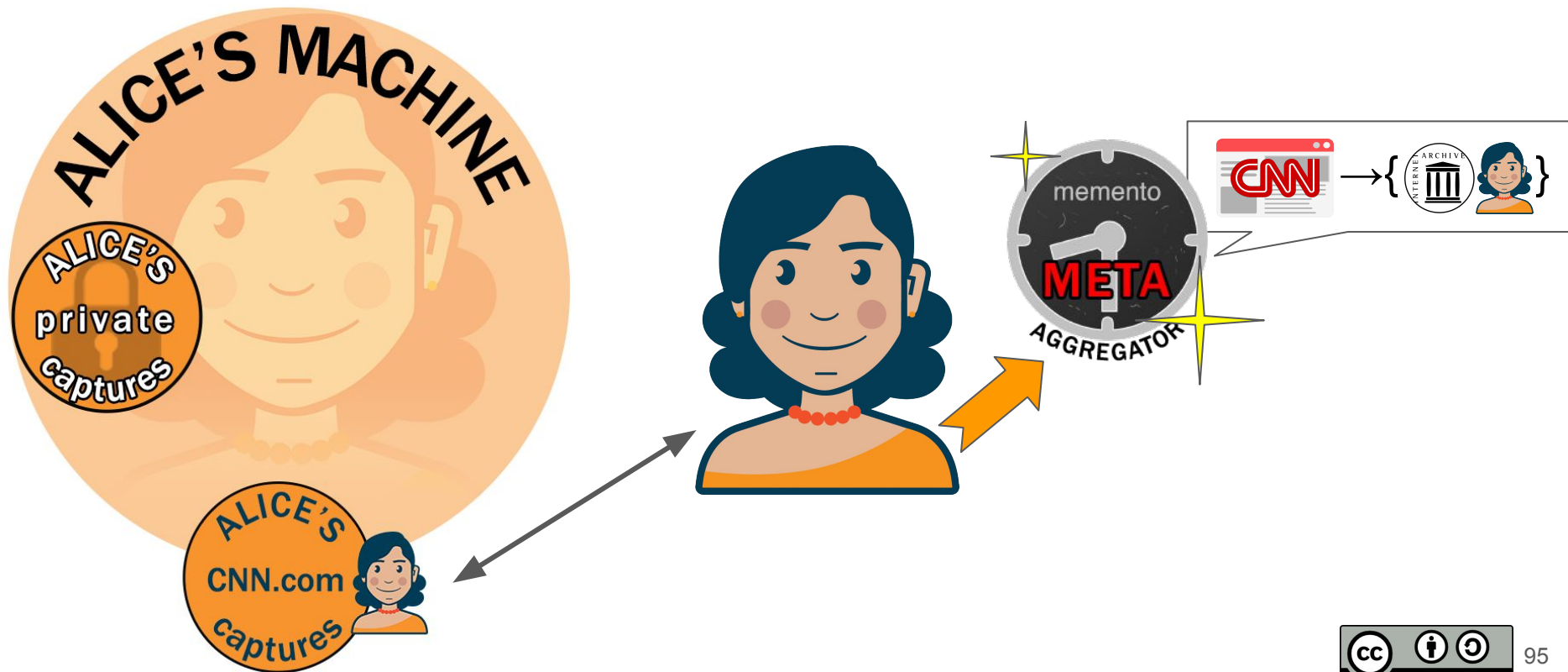


Mementity Dynamics - Alice & Her Archives (WA_A)

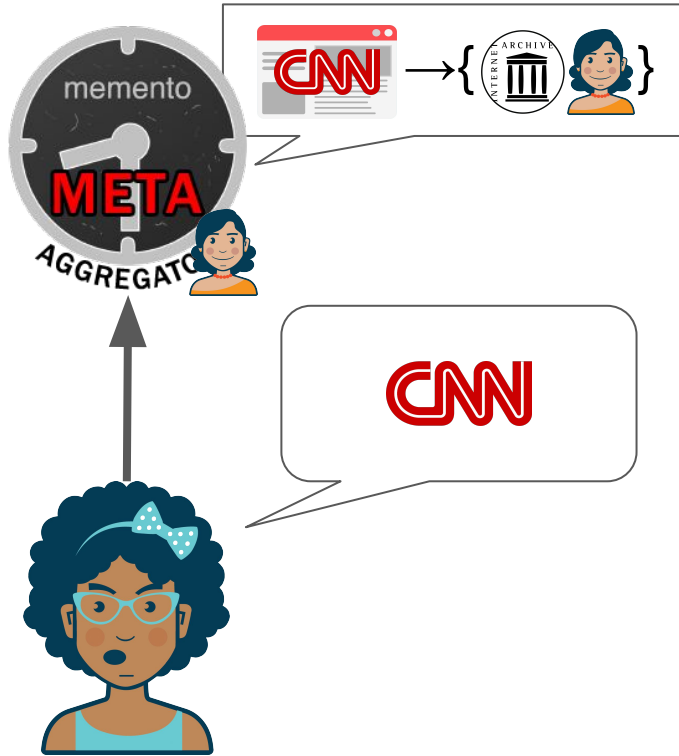


Personal Archive Aggregation

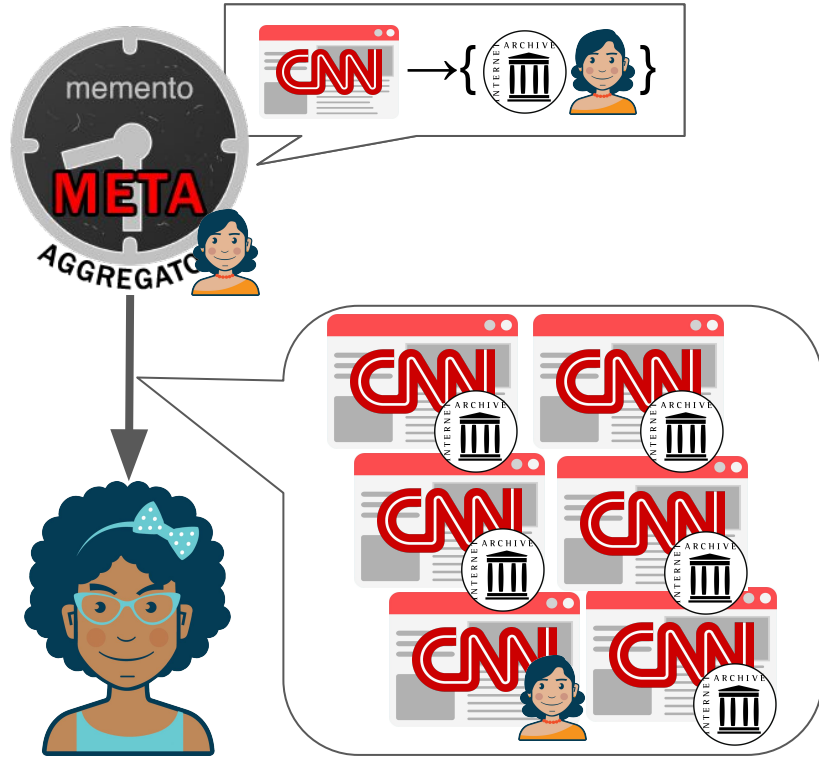
Alice Deploys MMA_A



Carol Asks MMA_A for CNN

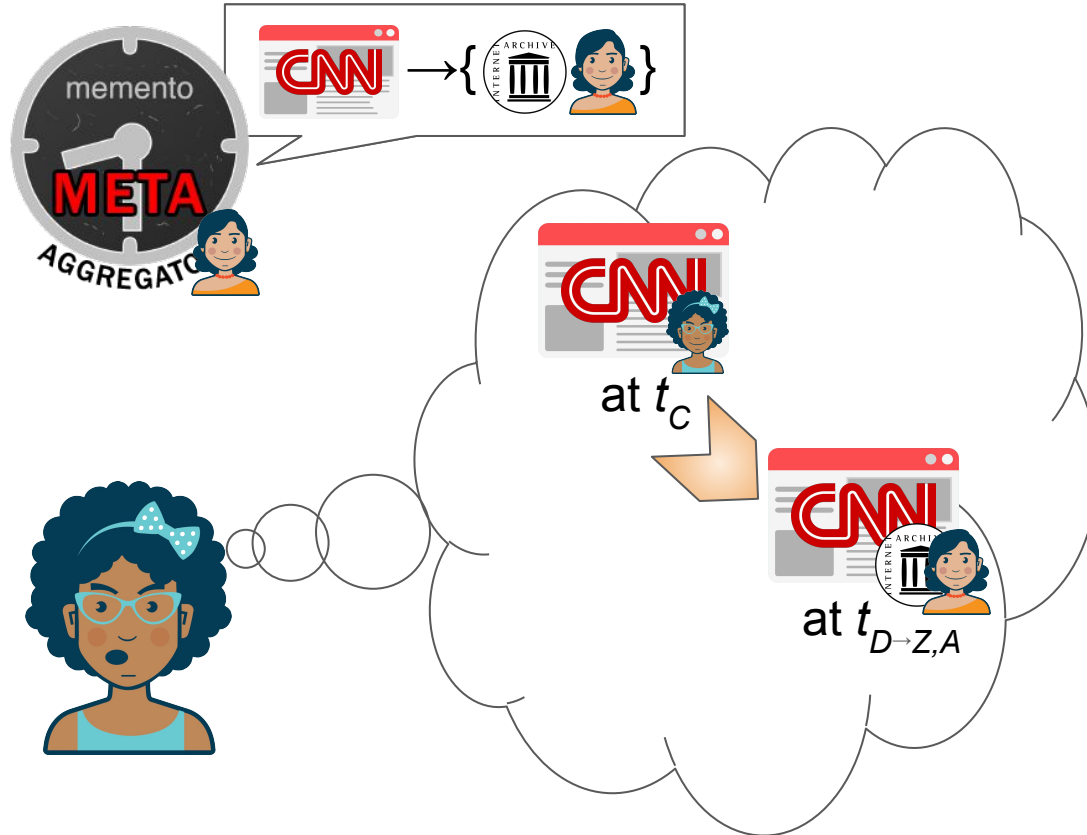


MMA_A returns CNN Memento $\{M_A, M_{IA}\}$

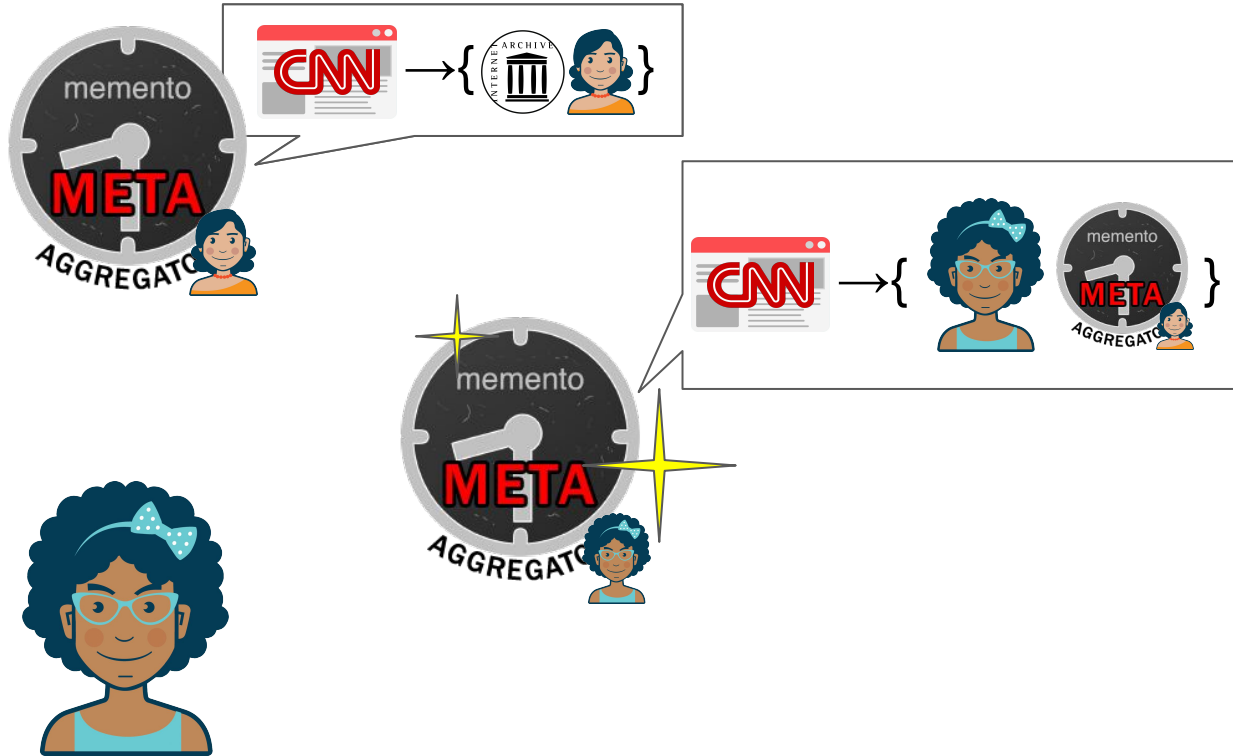


Carol Wants to Aggregate Her Own Captures

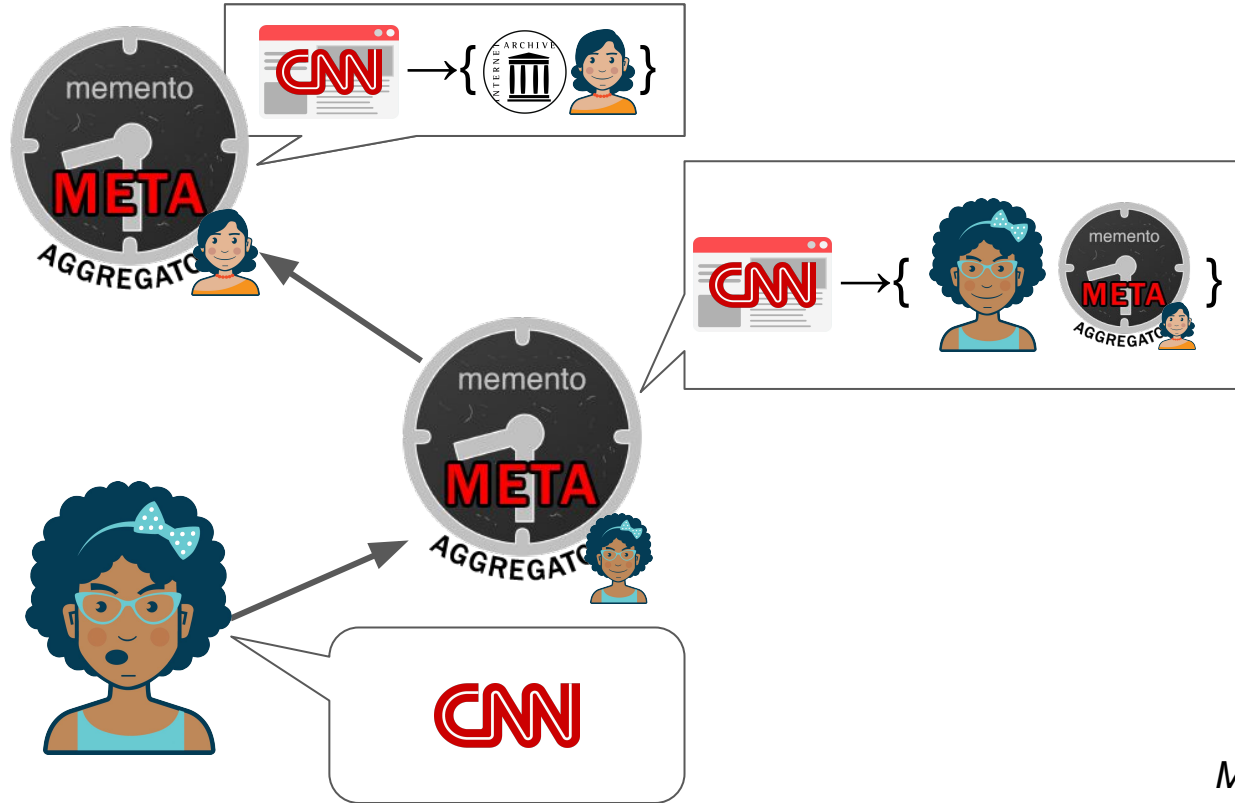
CNN(M(WA_C))



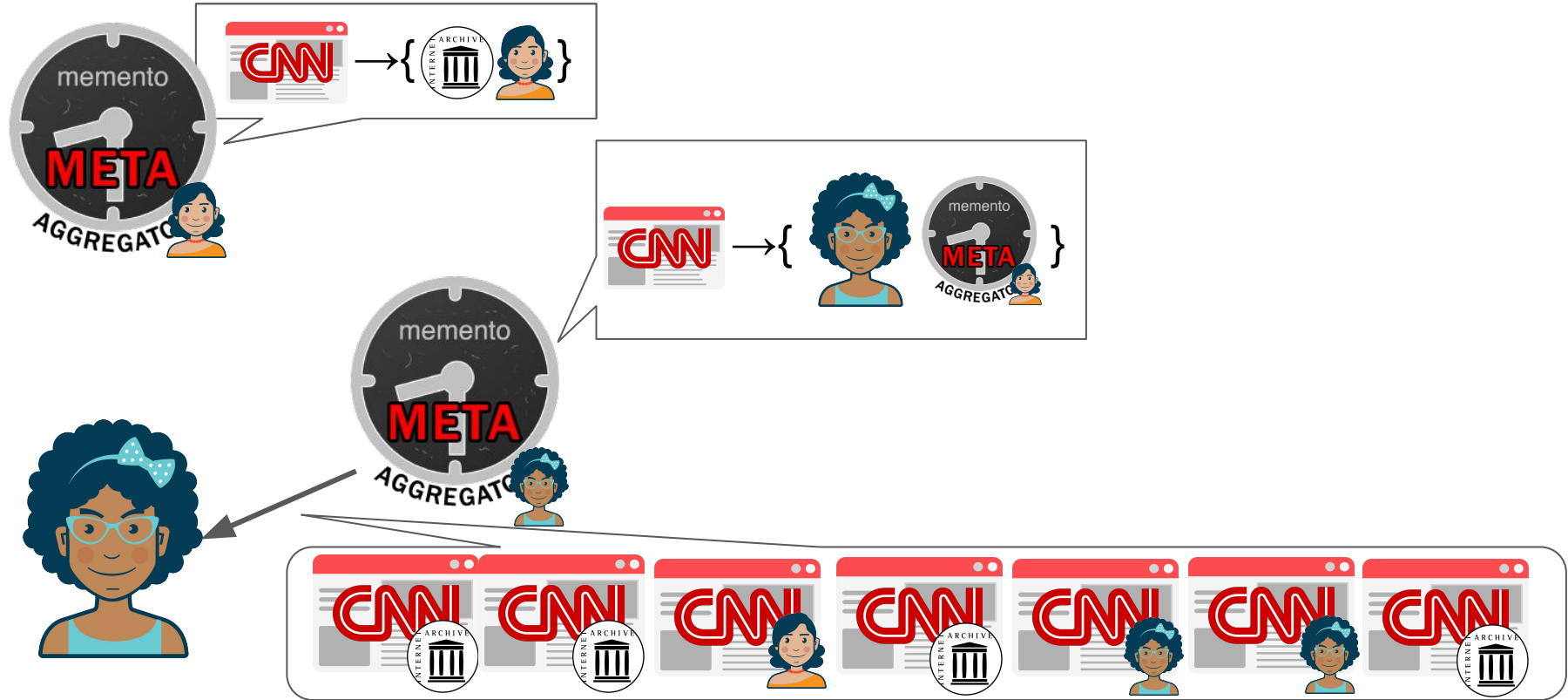
Carol Creates MMA_C to Access WA_C and MMA_A



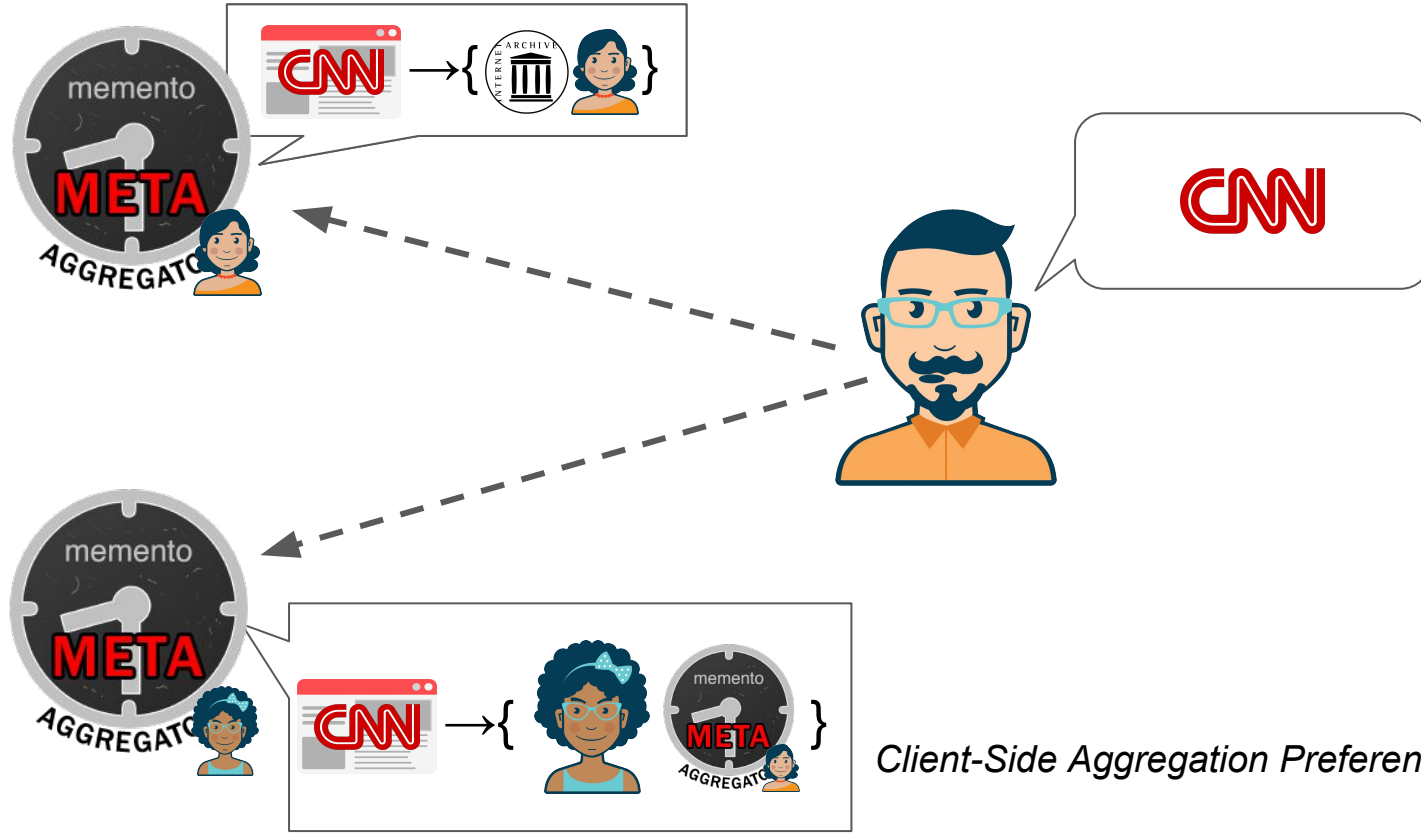
Carol Asks MMA_C For CNN



MMA_A returns CNN Memento $\{M_A, M_{IA}, M_C\}$



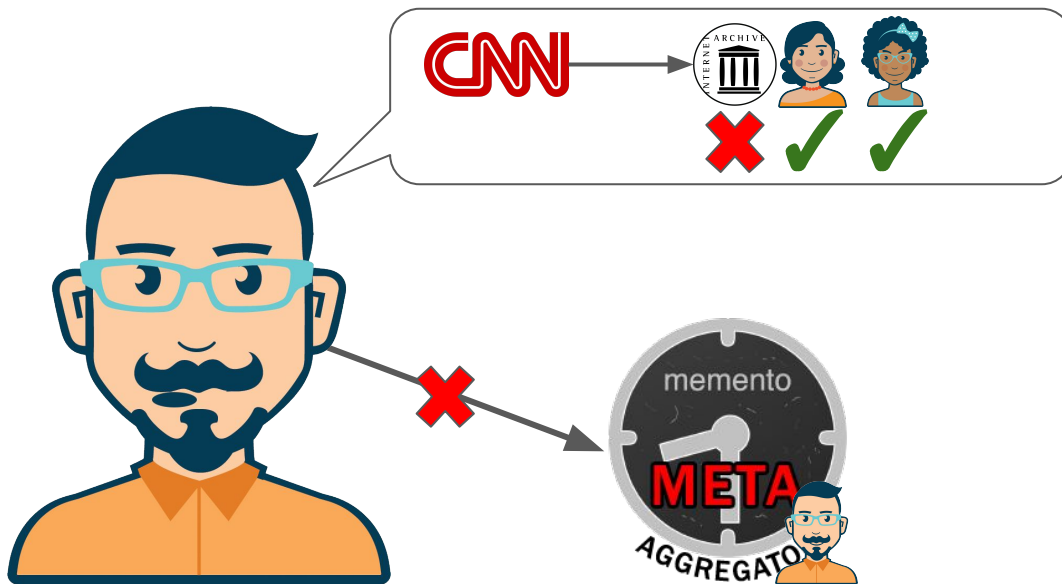
Bob May Request $M(\text{CNN})$ From MMA_A or MMA_C



Client-Side Aggregation Preference

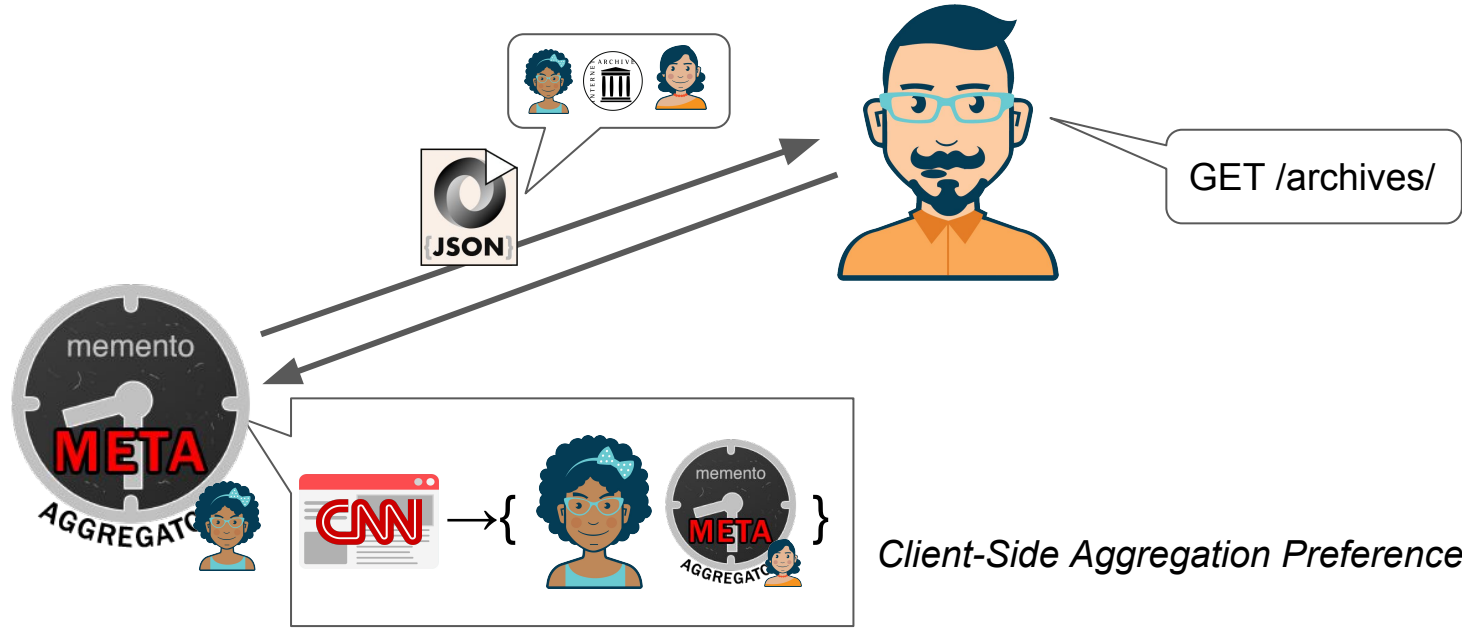
Bob Prefers to Exclude IA Captures

...and does not want to setup his own MMA

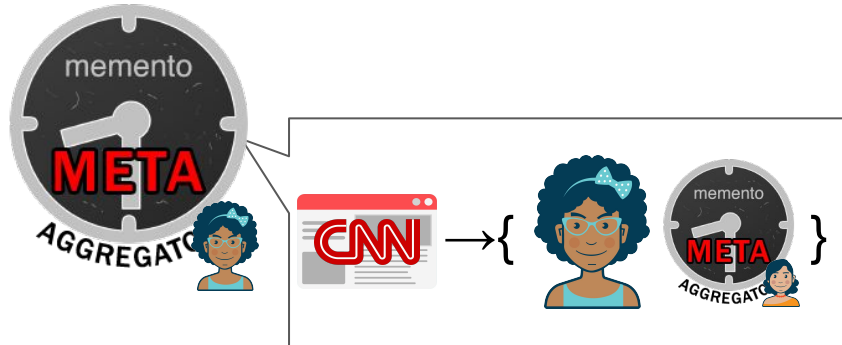
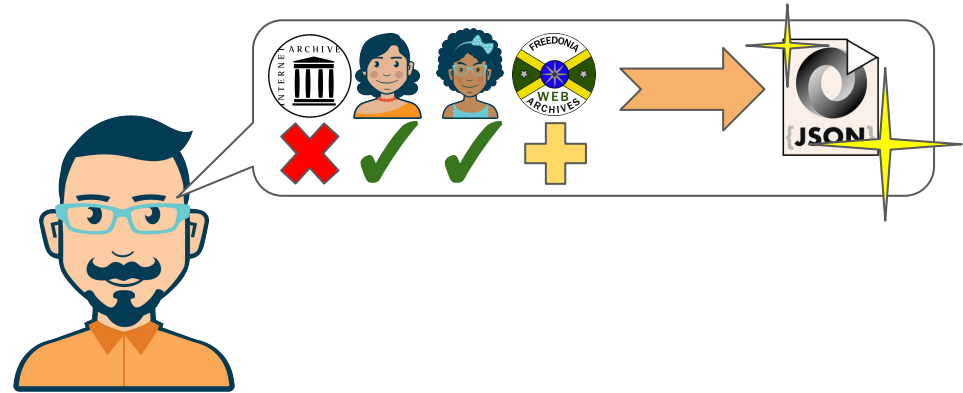


Client-Side Aggregation Preference

Bob Requests Supported Archives

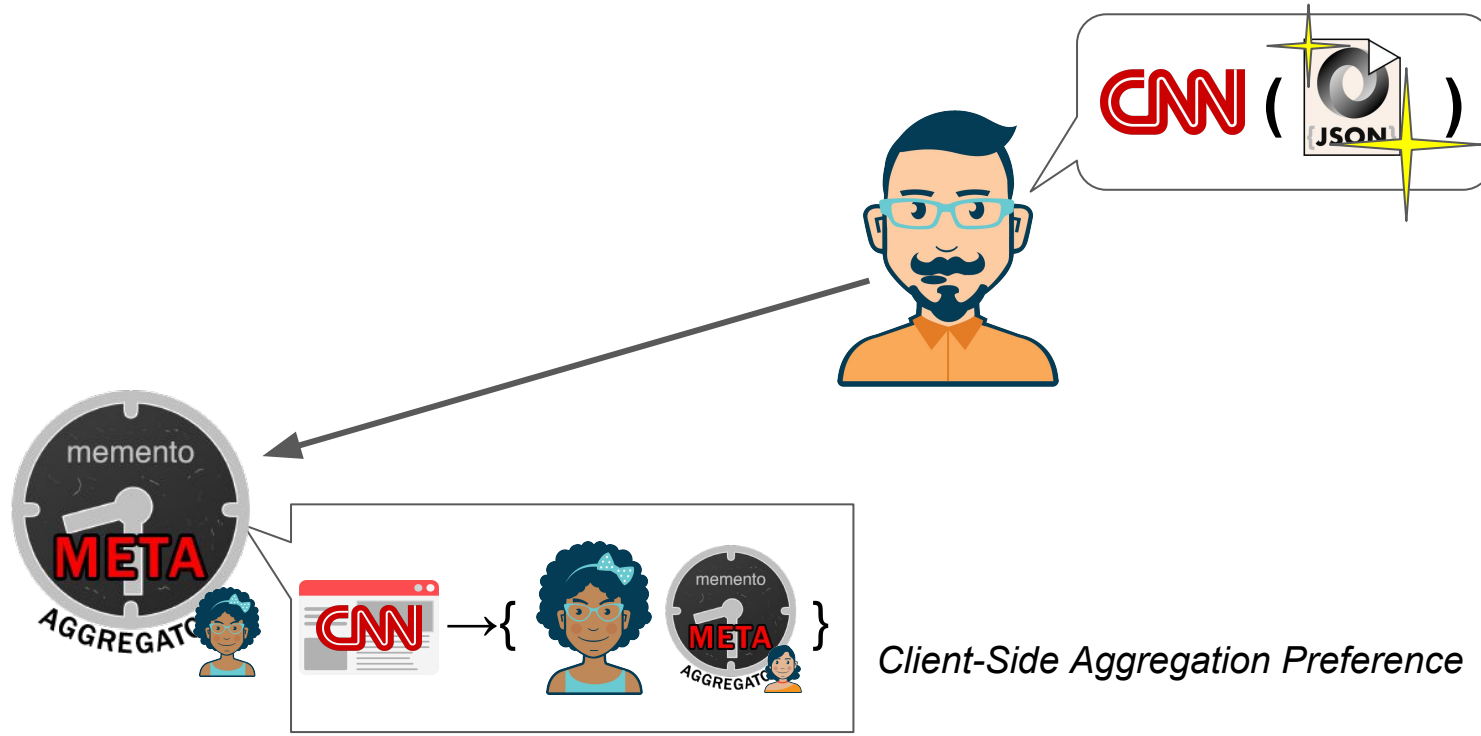


Bob Customizes the Set in the JSON

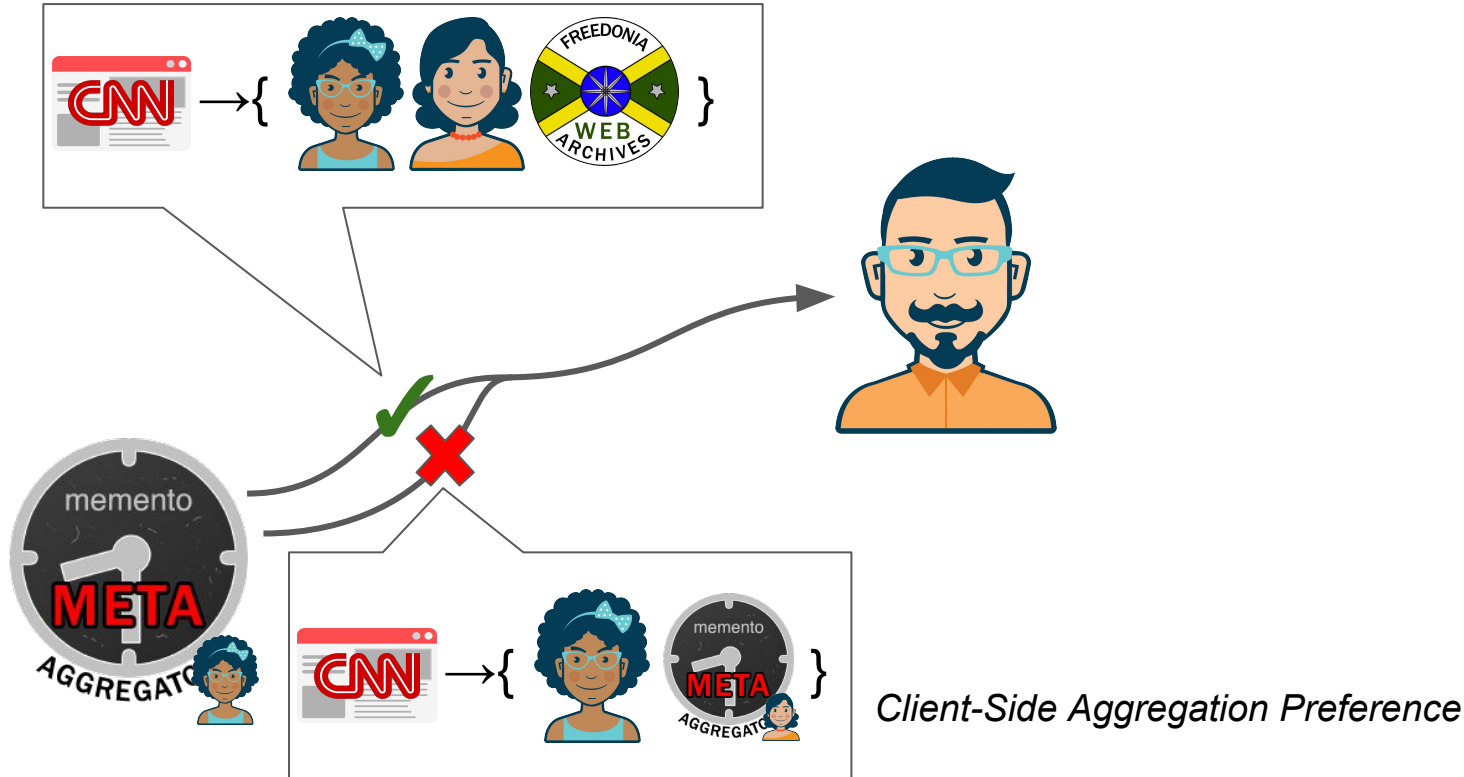


Client-Side Aggregation Preference

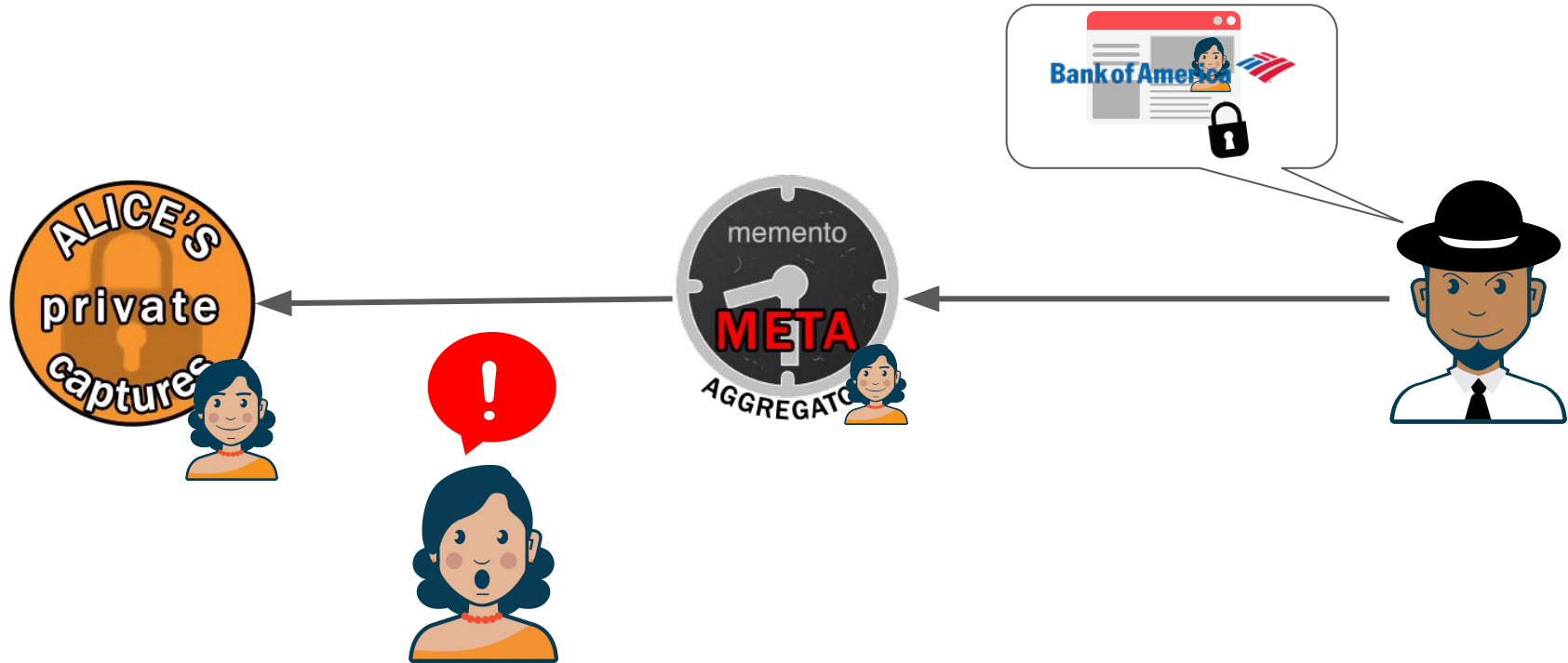
Bob Requests CNN for His Custom Set



MMA Complies or Ignores Preference



Hooray, Aggregation!



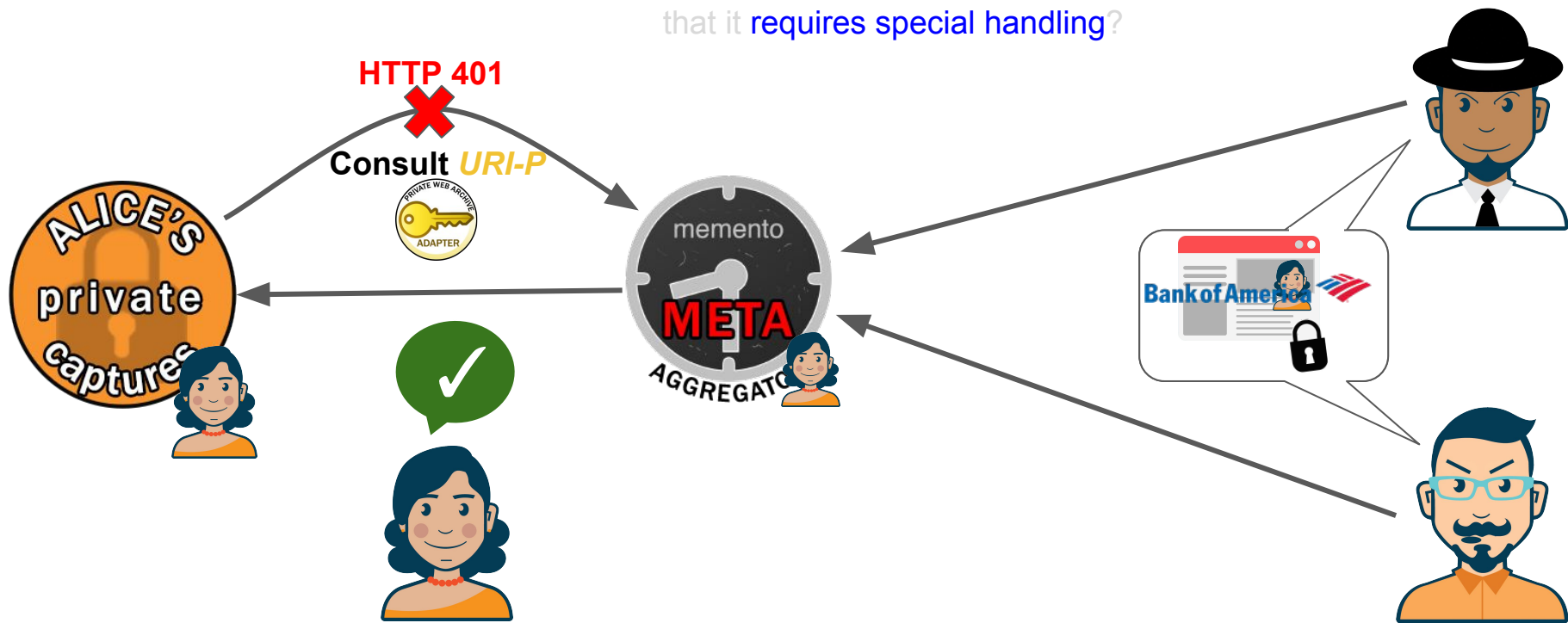
PROPOSED FRAMEWORK

Mementities



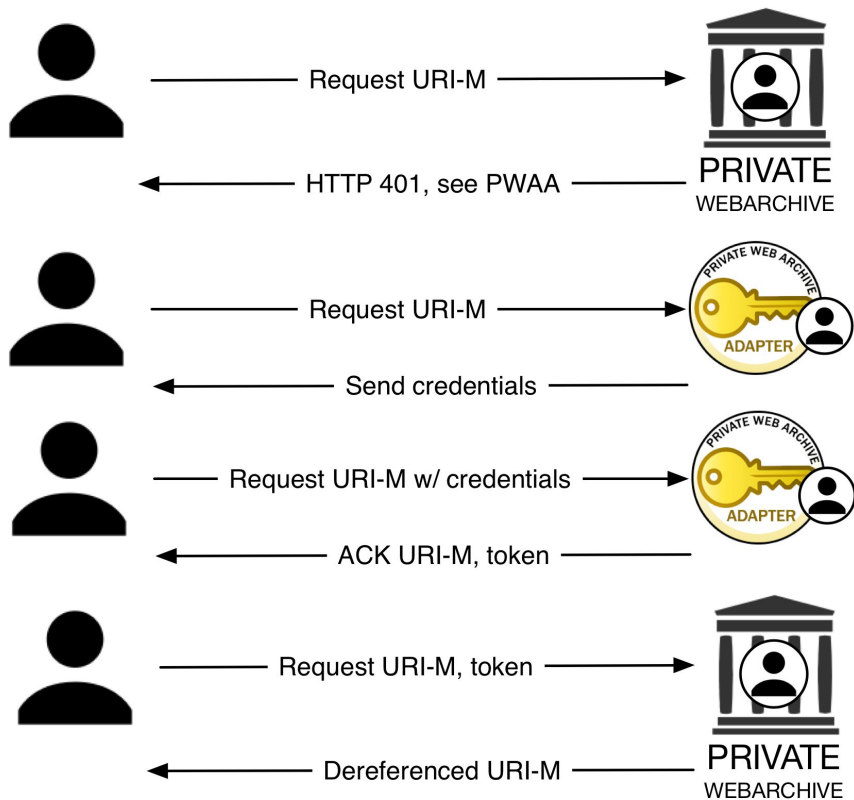
Hooray, Aggregation!

RQ4: How can **content** that was captured behind authentication **signal** to Web archive replay systems that it **requires special handling**?





Private Web Archive Adapter (PWAA)



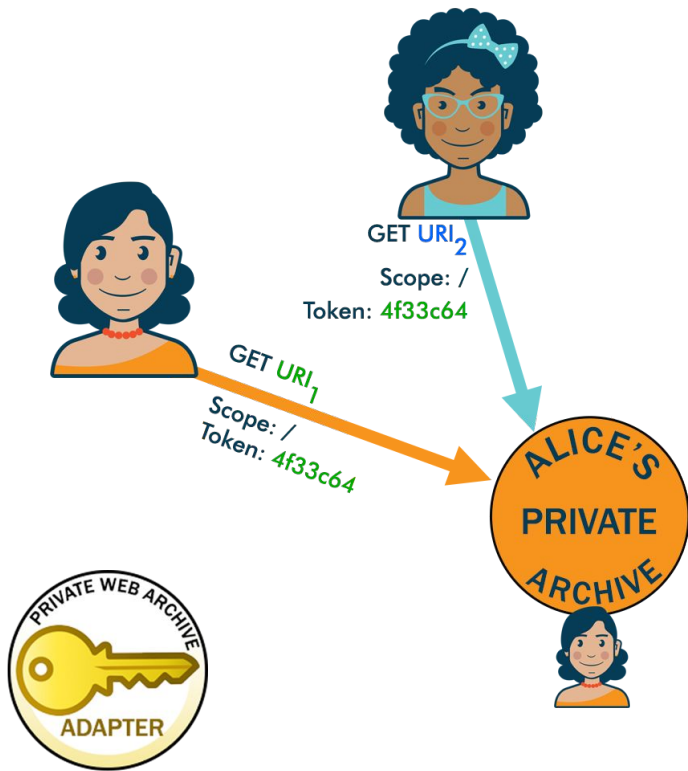
- Auth Layer for to encourage Private Web archive aggregation
- Typical OAuth 2.0 flow
- Auth role cohesive to PWAA
- Persistent access through tokenization

RQ6: What kinds of access control do users who create private Web archives need to **regulate access** to their archives?

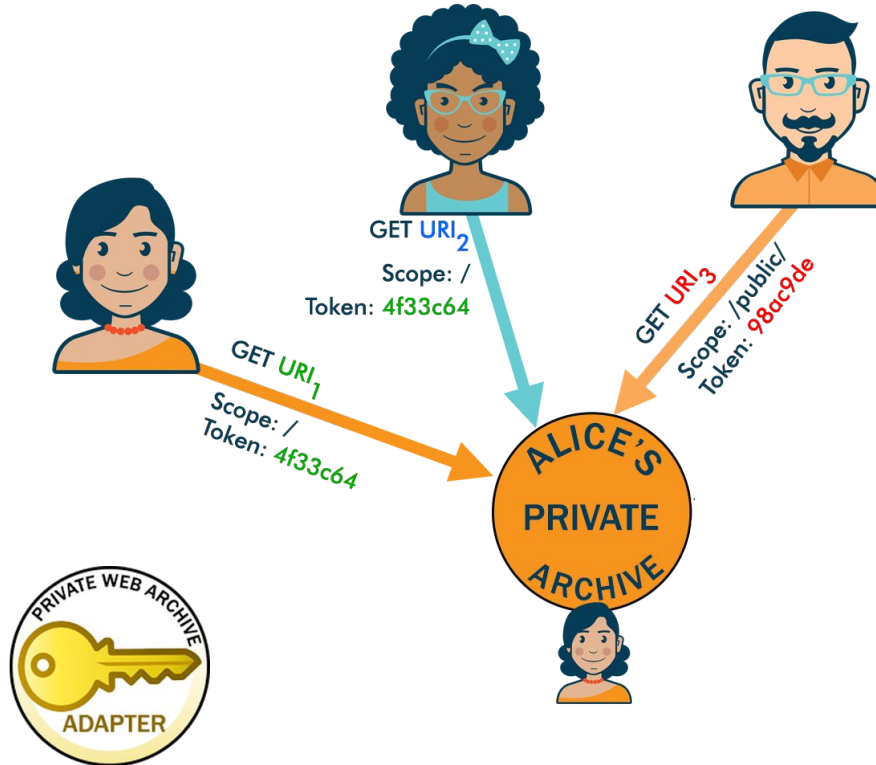


PWAA - Sharing Tokens

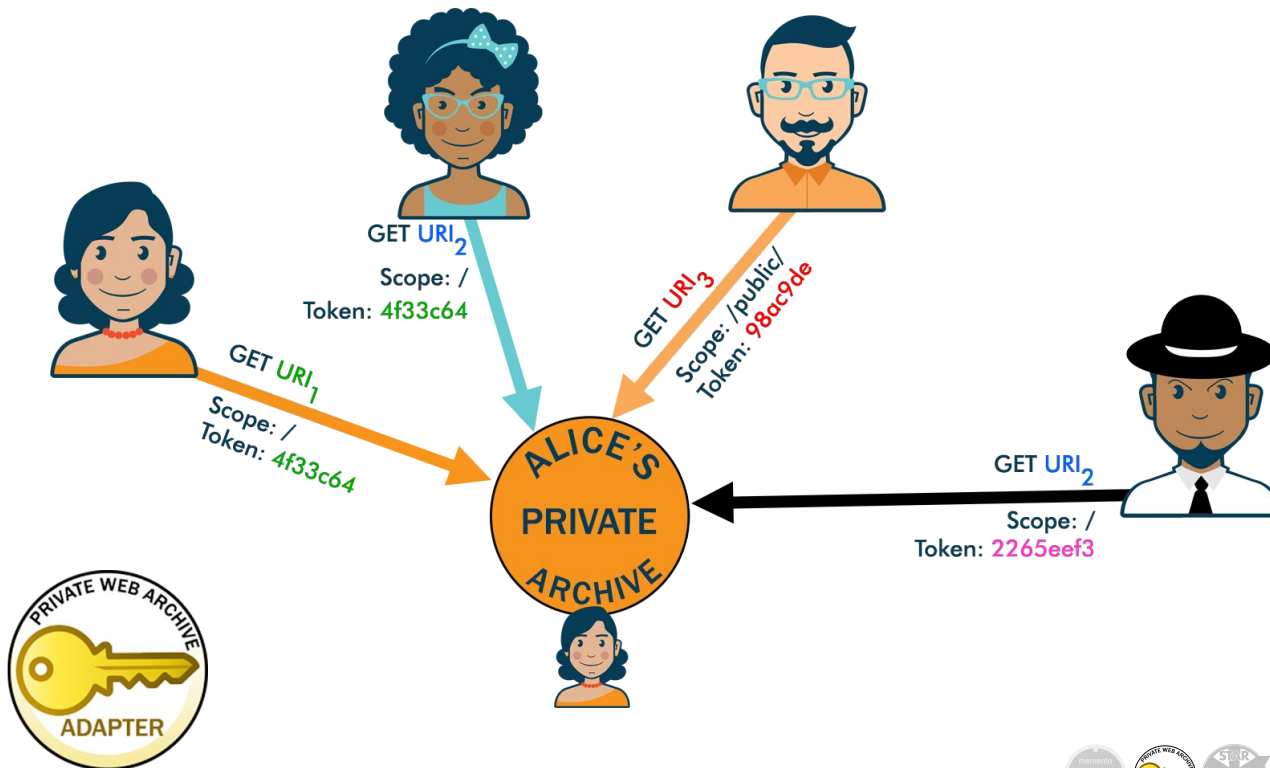
RQ6: What kinds of access control do users who create private Web archives need to regulate access to their archives?



PWAA - Previously Authorized

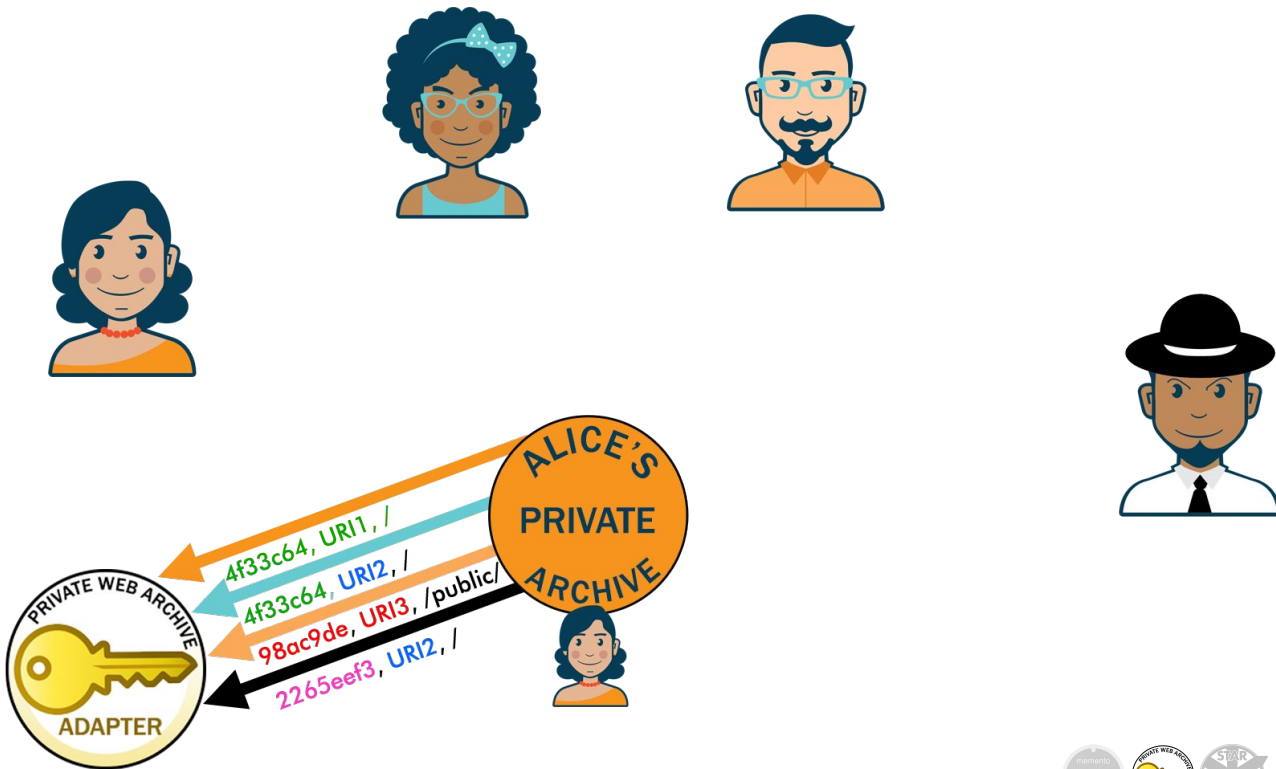


PWAA - Unauthorized Request

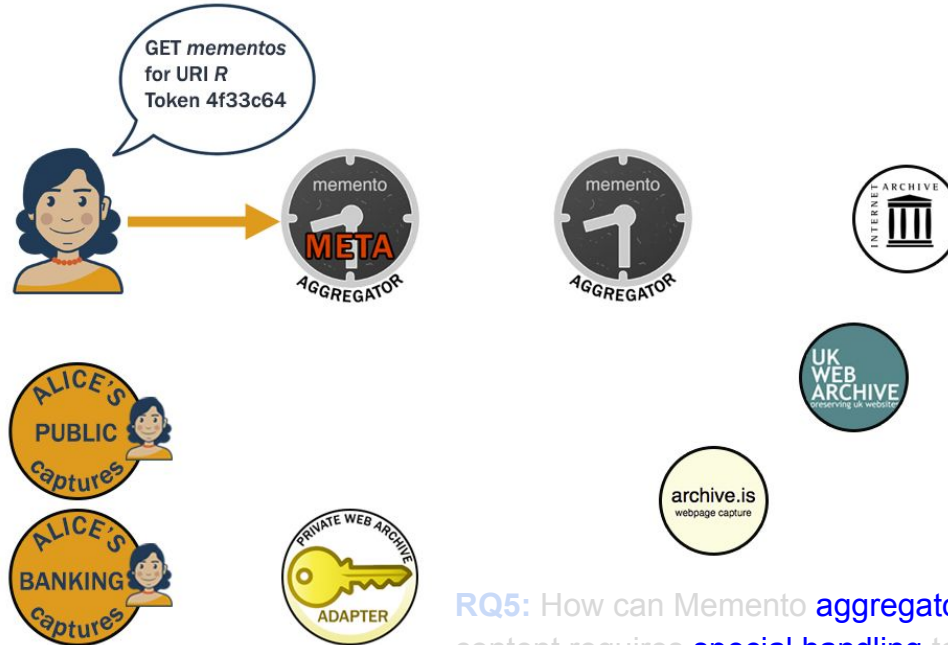


PWAA - Sharing Tokens

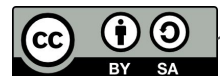
RQ6: What kinds of access control do users who create private Web archives need to regulate access to their archives?



Alice Passes Associative Token to MMA

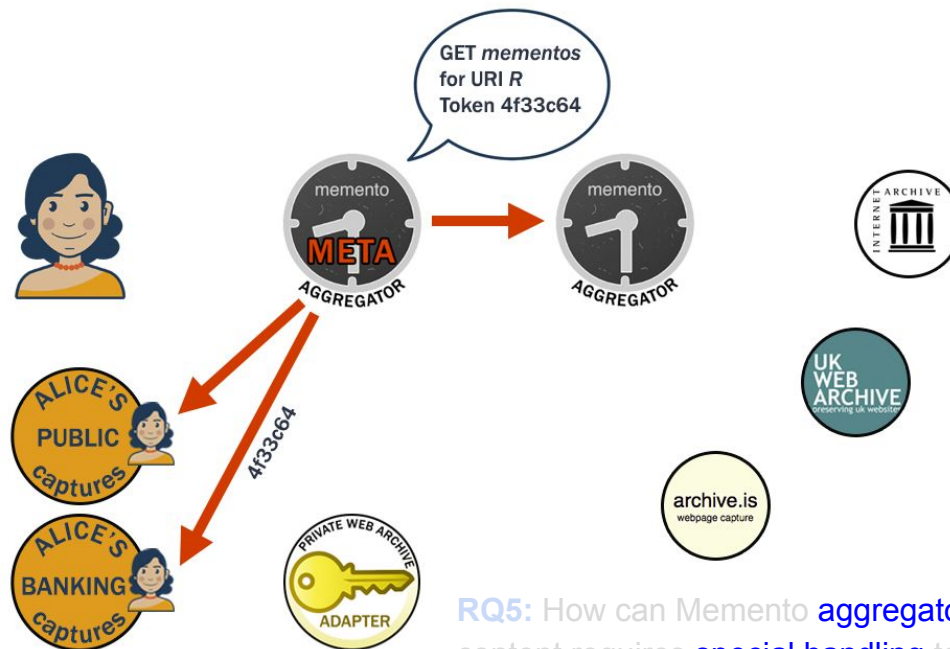


RQ5: How can Memento aggregators indicate that private Web archive content requires special handling to be replayed, despite being aggregated with publicly available Web archive content?



MMA requests URI-R...

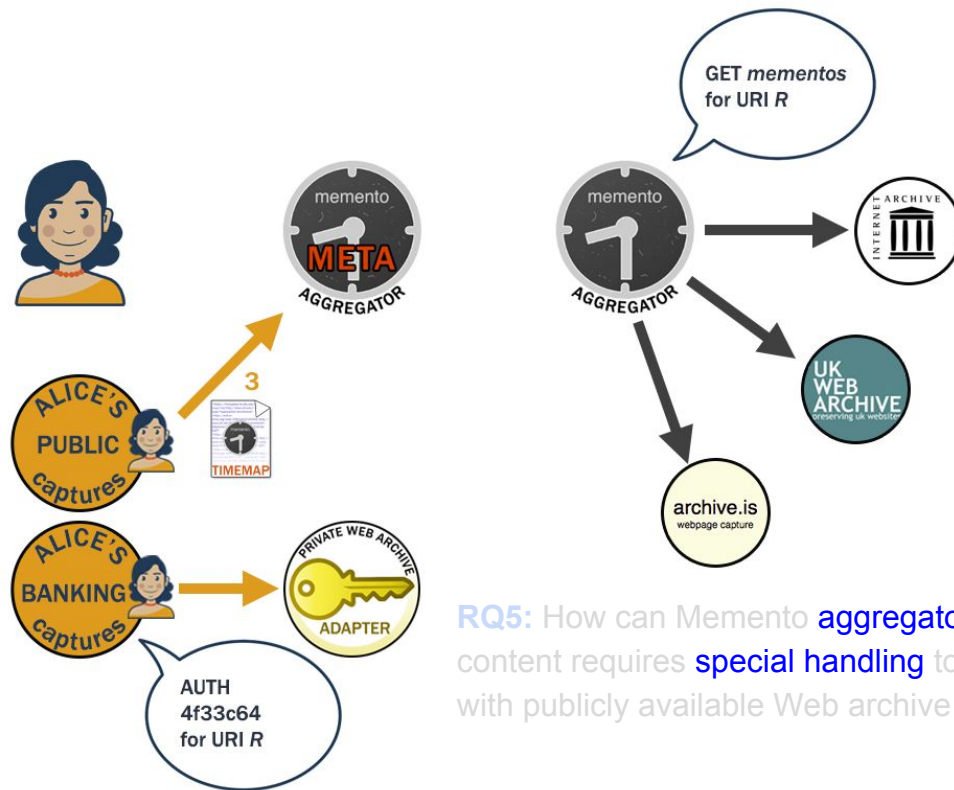
...relays token where applicable



RQ5: How can Memento aggregators indicate that private Web archive content requires special handling to be replayed, despite being aggregated with publicly available Web archive content?



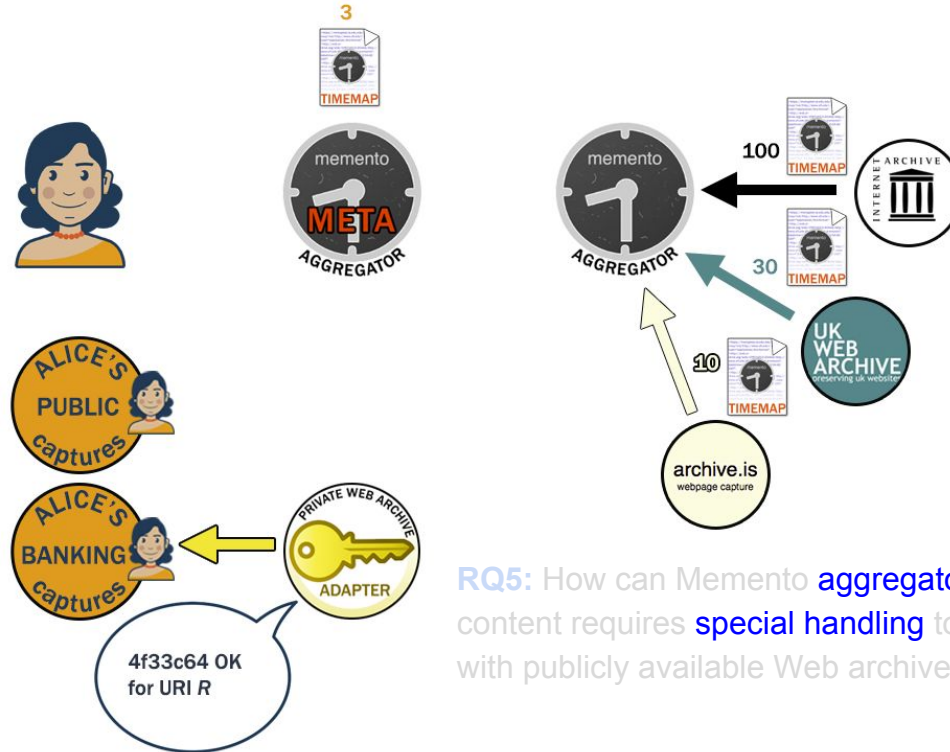
Private Archive Validates with PWAA



RQ5: How can Memento aggregators indicate that private Web archive content requires special handling to be replayed, despite being aggregated with publicly available Web archive content?



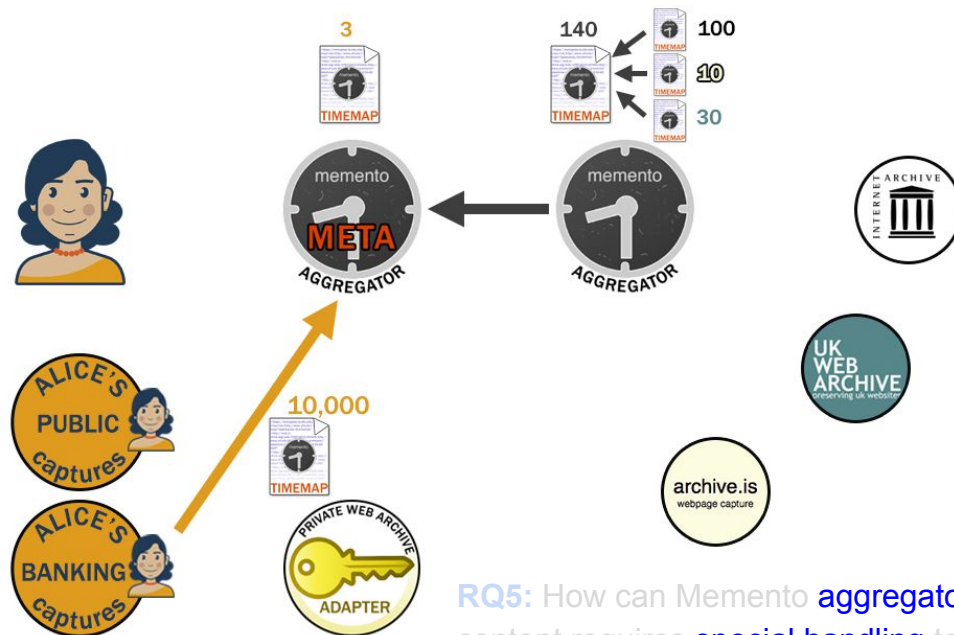
PWAA Confirms Token



RQ5: How can Memento aggregators indicate that private Web archive content requires special handling to be replayed, despite being aggregated with publicly available Web archive content?



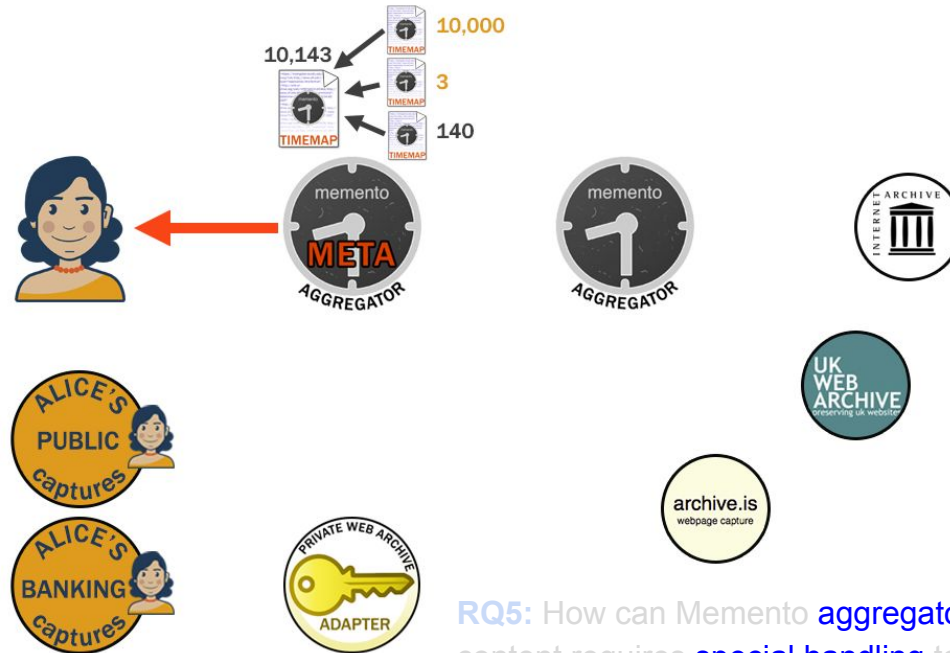
Private Archive Returns Captures



RQ5: How can Memento aggregators indicate that private Web archive content requires special handling to be replayed, despite being aggregated with publicly available Web archive content?



MMA Aggregates, Associates Token

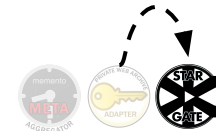


RQ5: How can Memento aggregators indicate that private Web archive content requires special handling to be replayed, despite being aggregated with publicly available Web archive content?

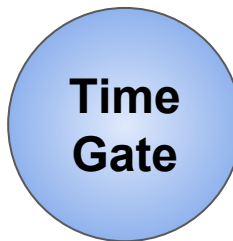


PROPOSED FRAMEWORK

Mementities



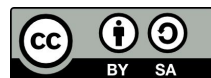
StarGate



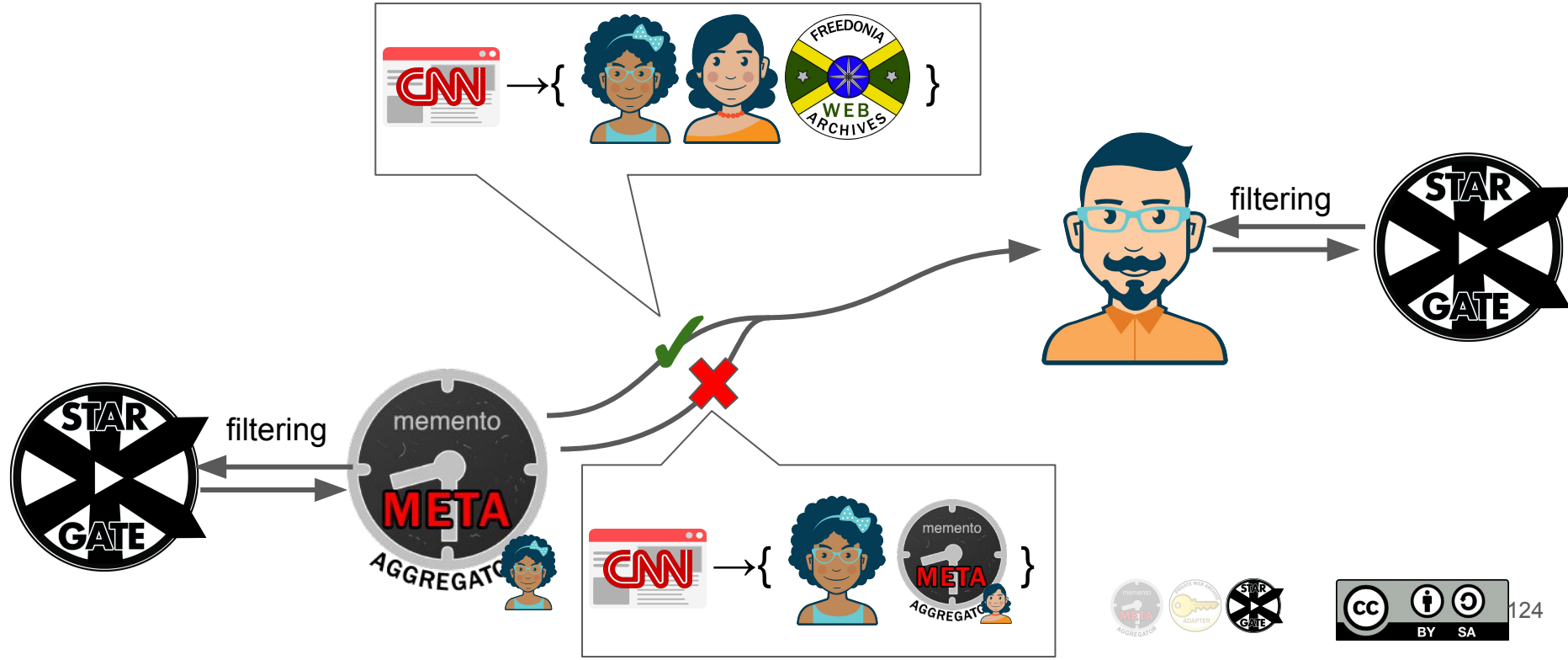
functional
 \subseteq



- Content negotiation in Web archives **beyond time**
- “Star” ~ wildcard (*) → any dimension of negotiation
- Allow for queries like: *Only show me mementos...*
 - That are not redirects (*content-based attribute* HTTP Status \neq 3XX)
 - Of a sufficient quality (*derived attribute* Memento Damage < 0.4)
 - Are from personal Web archives (*access attribute* indicate Facebook.com memento is not a login page)

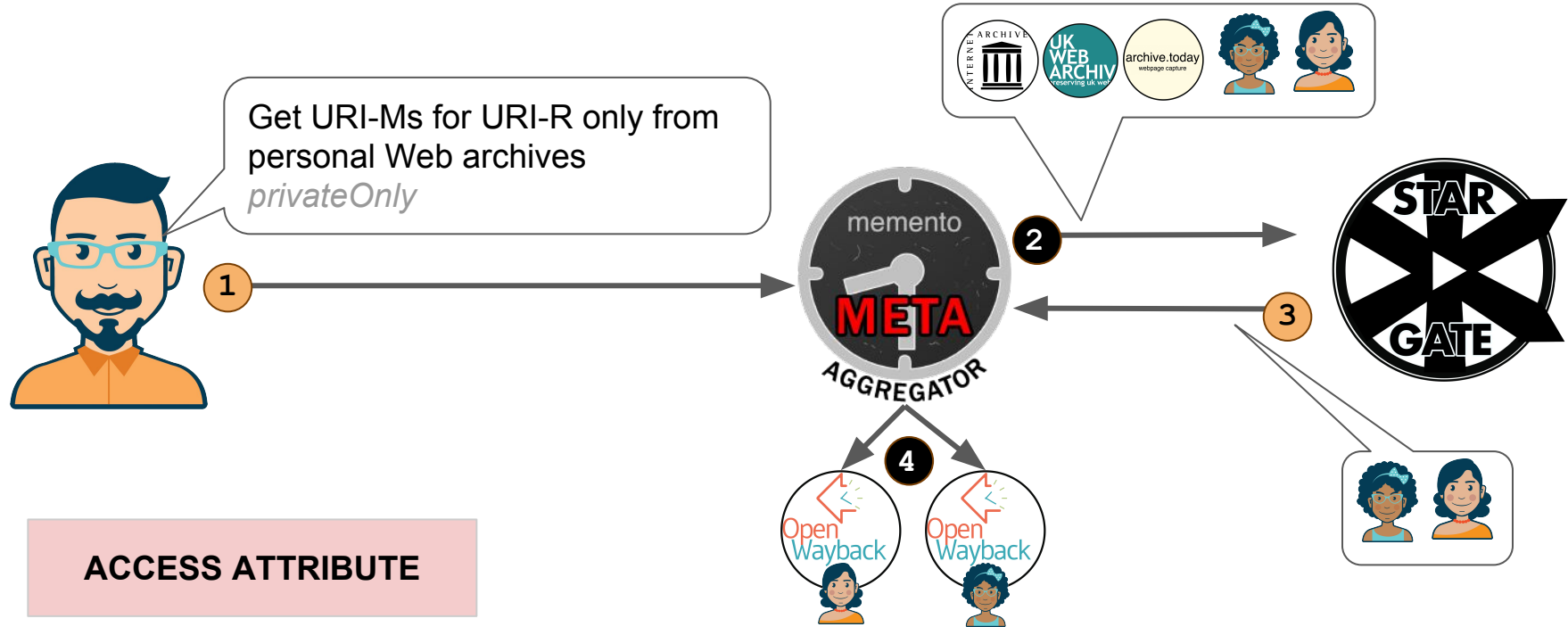


Implicit Filtering via MMA or Directly (a la TG)



Negotiation in the Privacy Dimension

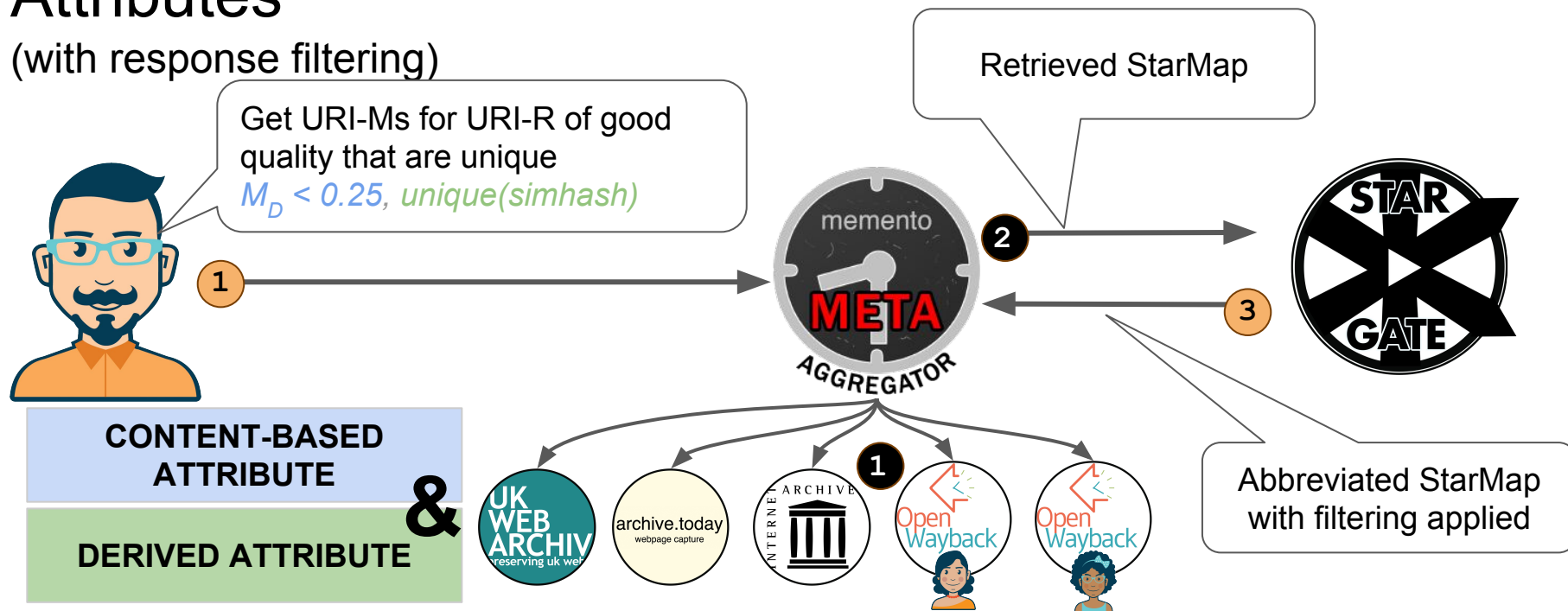
(via short circuiting)



ACCESS ATTRIBUTE

Negotiation on Content-Based or Derived Attributes

(with response filtering)



A Framework for Aggregating Public and Private Web Archives

February 14, 2019

Mat Kelly

Outline

- Introduction/Motivation
- Background
- Preliminary Research
- Proposed Framework
- **Evaluation Plan**

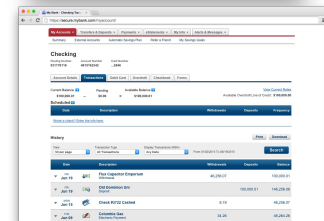
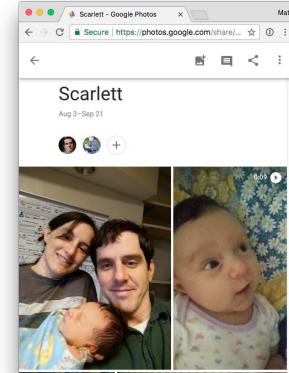
Framework Evaluation

- Evaluation of mementity design decisions
- Costs of more expressive TimeMaps (StarMaps) and Link header enrichment
- Evaluation through implementation



Evaluation of Mementity Design Decisions

- Effectiveness in resolving initial use cases and access patterns
- *“It was there yesterday, where did it go?”*
- *“Save this, but only for me.”*
- *“I want to share this but control who can see it.”*



A Framework for Aggregating Public and Private Web Archives

February 14, 2019

Mat Kelly

EVALUATION: [Design/Practicality](#) • Enrichment costs • Implementation



Costs of more expressive TimeMaps (StarMaps) and Link header enrichment

- Computational:
 - Mostly server-side, potential to further leverage client
- Temporal
 - Required on variant generation
- Spatial
 - Permutation variant storage
- Access
 - Variant negotiation

Evaluation Through Implementation

extend



Extend for client-side archival specification

extend



Exhibit features of an MMA

create



Regulate access to Private Web archives

create



Facilitate archival negotiation in more dimensions

569 Mementos Available

Sources

- ✓ Internet Archive
- Archive.is
- Local Archive1
- ✓ Local MMA1
- Remote MA1

VIEW BY: Dropdown Columns VizMethodFoo VizMethodBar

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

A Framework for Aggregating Private and Public Web Archives

Mat Kelly

Old Dominion University
Web Science & Digital Libraries Research Group
Department of Computer Science
Norfolk, Virginia USA
mkelly@cs.odu.edu



Seminar, Penn State University
February 14, 2019



Backup Slides

Research Questions

RQ1: What sort of content is difficult to capture and replay for preservation from the perspective of a Web browser?

RQ2: How do Web browser APIs compare in potential functionality to the capabilities of archival crawlers?

RQ3: What issues exist for capturing and replaying content behind authentication?

RQ4: How can content that was captured behind authentication signal to Web archive replay systems that it requires special handling?

RQ5: How can Memento aggregators indicate that private Web archive content requires special handling to be replayed, despite being aggregated with publicly available Web archive content?

RQ6: What kinds of access control do users who create private Web archives need to regulate access to their archives?

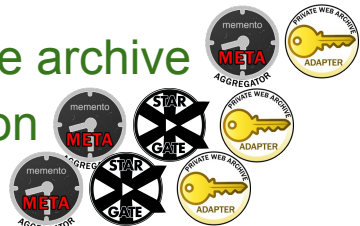
User Access Patterns

- Pattern 1: Single archive access
- Pattern 2: Aggregation of multiple Web archives

Pre-existing archival usage

Contribution beyond proposal

- Pattern 3: Aggregator chaining 
- Pattern 4: Aggregation with authentication  
- Pattern 5: Aggregation including a hybrid public-private archive
- Pattern 6: Aggregation with filtering via MMA interaction
- Pattern 7: Aggregation with filtering via SG interaction



CDXJ: An Alternative TimeMap Format

```
<http://matkelly.com>; rel="original",
<http://localhost:1208/timemap/link/http://matkelly.com>;
rel="self"; type="application/link-format",
<http://web.archive.org/web/20060514123511/http://www.mat
kelly.com:80/>; rel="first memento"; datetime="Sun, 14
May 2006 12:35:11 GMT",
<http://web.archive.org/web/20060516213852/http://www.mat
kelly.com/>; rel="memento"; datetime="Tue, 16 May 2006
21:38:52 GMT",
...
<http://web.archive.org/web/20180128152125/http://matkell
y.com>; rel="memento"; datetime="Sun, 28 Jan 2018
15:21:25 GMT",
<http://web.archive.org/web/20180319141920/http://matkell
y.com/>; rel="last memento"; datetime="Mon, 19 Mar 2018
14:19:20 GMT",
<http://localhost:1208/timemap/link/http://matkelly.com>;
rel="timemap"; type="application/link-format",
<http://localhost:1208/timemap/json/http://matkelly.com>;
rel="timemap"; type="application/json",
<http://localhost:1208/timemap/cdxj/http://matkelly.com>;
rel="timemap"; type="application/cdxj+ors",
<http://localhost:1208/timegate/http://matkelly.com>;
rel="timegate"
```

```
!context ["http://tools.ietf.org/html/rfc7089"]
!id {"uri": "http://localhost:1208/timemap/cdxj/http://matkelly.com"}
!keys ["memento_datetime_YYYYMMDDhhmmss"]
!meta {"original_uri": "http://matkelly.com"}
!meta {"timegate_uri":
"http://localhost:1208/timegate/http://matkelly.com"}
!meta {"timemap_uri": {"link_format":
"http://localhost:1208/timemap/link/http://matkelly.com", "json_format":
"http://localhost:1208/timemap/json/http://matkelly.com", "cdxj_format":
"http://localhost:1208/timemap/cdxj/http://matkelly.com"}}}
20060514123511 {"uri":
"http://web.archive.org/web/20060514123511/http://www.matkelly.com:80/",
"rel": "first memento", "datetime": "Sun, 14 May 2006 12:35:11 GMT"}
20060516213852 {"uri":
"http://web.archive.org/web/20060516213852/http://www.matkelly.com/",
"rel": "memento", "datetime": "Tue, 16 May 2006 21:38:52 GMT"}
...
20180128152125 {"uri":
"http://web.archive.org/web/20180128152125/http://matkelly.com", "rel":
"memento", "datetime": "Sun, 28 Jan 2018 15:21:25 GMT"}
20180319141920 {"uri":
"http://web.archive.org/web/20180319141920/http://matkelly.com/", "rel":
"last memento", "datetime": "Mon, 19 Mar 2018 14:19:20 GMT"}
```

Link (RFC 7089) TimeMap

Original URI (URI-R)

Other TimeMaps (URI-Ts)

TimeGate (URI-G)

Relative Relations

CDXJ TimeMap

See Alam, [“CDXJ: An Object Resource Stream Serialization Format”](#), 2015

Private & Public Archives May Differ for the Same URI



Should Public Archives *Really* Capture the Private Web?

