# Save *My* Web Things For Me

## and Others

# but Not All

## Mat Kelly, PhD

Assistant Professor, Information Science
Drexel University
College of Computing & Informatics
mkelly@drexel.edu
https://matkelly.com

Master's of Information Webinar Series
February 4, 2020

Drexel UNIVERSITY

# Drexel CCI'S DCM Major

- Specialization in Digital Content Management for Masters in Information Science degree
- Courses
  - 3 foundation
  - 5 core
  - 6 electives
  - Capstone (2 quarters)
- Sample Courses in Major:
  - INFO633 Information Visualization
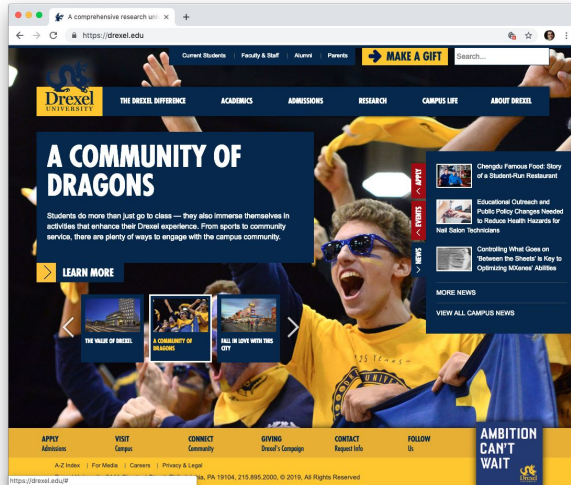  - INFO654 Enterprise Content Management (being taught by Dr. Kelly, Spring 2020)
  - INFO 676 Applied Ontology

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

2

# The Web

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
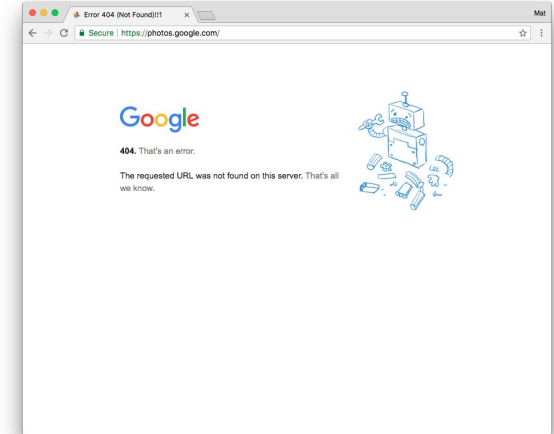Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

3

# The Web...is Ephemeral

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

4

# The Web...can be preserved



Internet Archive's interface to "Save Page Now"



IA's result of preserving facebook.com

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
**Computing &**
**Informatics**

5

# The Web...can be preserved*



≠

What you see on the live web

IA's result of preserving facebook.com

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

6

# Archiving Institutions' Efforts



S. Alam et al., Web Archive Profiling Through CDX Summarization, International Journal on Digital Libraries, 17(3), pp. 223--238, 2016.

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

7

# Institutions' Efforts...likely won't have your things

Private social media pages

Bank account ledger

Online photo storage

- Niche Sites
- Corporate *intra*nets
- etc.

...and it might be best they didn't.

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

8

# The Divide between the Live & Archived Web

For what they *do* capture:

- What (URI?) do these captures represent?
- When were they captured?
- Can the entities in the archive give us more context?

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

9

# Needed a Standard Way for Archives to Express this Association

What did drexel.edu look like in the past?



INTERNET ARCHIVE

2019
2012
2008
2004
1997

# Memento Bridges the Divide



TimeMap for https://drexel.edu/cci/

original: https://drexel.edu/cci/
Memento1: March 14, 2014 @ $M_3$
Memento2: March 18, 2016 @ $M_1$
Memento3: February 3, 2020 @ $M_2$

TimeMap for https://drexel.edu

original: https://drexel.edu/
Memento1: August 15, 2007 @ $M_4$
...

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

11

# Web Archiving Components<sup>(two of them)</sup>

**Preservation**

- Saving Live Web Content Offline in a Standard, Interoperable Format (WARC)
  - Includes HTML, CSS, JavaScript, Flash, MP4, etc. and resources within each

**Access**

- Associating a request for a URI in time to a WARC
- Re-assembling Pages (no easy task)

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

12

# Preservation

- Capture an HTML page at a URI
- Capture all resource representations embedded within a page
- Associate the what (URI) and when (time) of the captured resource representations



CCI homepage ~now

CCI homepage, "recently" archived

Web Archiving Components
# Preservation

# ISO 28500:2017
## WARC file format

- Concatenated "records"
- WARC header provides metadata and provenance info
- HTTP header blocks preserved as Part of the transaction

WARC response header

HTTP header block

HTTP payload block

```
52  WARC/1.0
53  WARC-Type: response
54  WARC-Target-URI: http://ipwb.example.com/
55  WARC-Date: 2016-09-07T00:38:19Z
56  WARC-Record-ID: <urn:uuid:1e3907a9-2e5c-9981-6a92-964a465d998e>
57  Content-Type: application/http; msgtype=response
58  Content-Length: 800

60  HTTP/1.1 200 OK
61  Host: ipwb.example.com
62  Connection: close
63  Content-Type: text/html; charset=UTF-8
64  Content-Length: 666

66  <html><head>
67  <title>InterPlanetary Wayback</title>
68  <link rel="stylesheet" type="text/css" href="style.css">
69  </head>
70  <body>
71  <h1>This is site for Space Dog</h1>
72  <img src="yuri.jpg">
73  <p>InterPlanetary Wayback (ipwb) facilitates permanence and collabo
74
75  </body></html>
76
77
78  WARC/1.0
79  WARC-Type: request
80  WARC-Target-URI: http://ipwb.example.com/style.css
81  WARC-Date: 2016-09-07T00:38:19Z
82  WARC-Concurrent-To: <urn:uuid:2d315cc1-a34d-3945-c5d9-ab4c7ac13fe6>
83  WARC-Record-ID: <urn:uuid:5a1491a6-f5be-d75e-25bd-6650c69a7182>
```
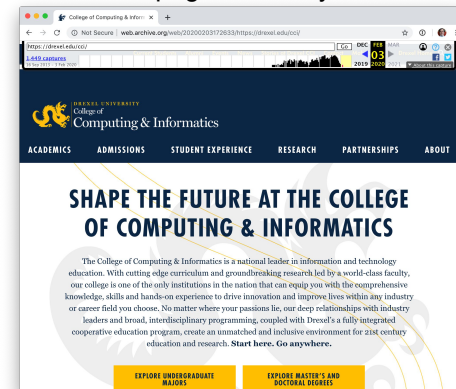
Line 91, Column 14    Tab Size: 4    Plain Text

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

14

# Access



- WARCs are indexed for URI→ record lookup
- HTML content is served with URIs rewritten to "point" back into the archive
  - Inclusive of embedded resources and "external links" pointing to other archived web pages
- Further lookups result with subsequent rewriting until no further HTTP requests are made

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

15

Web Archiving Components
# Access

Experience the past representation online, e.g.,

- http://web.archive.org/web/20130319235928/http://www.cnn.com/
    - Cnn.com at Tue, 19 Mar 2013 23:59:28 GMT
- http://archive.ph/9tEp6
    - cnn.com at Wed, 10 Apr 2013 11:33:55 GMT

URIs are opaque*, semantics should not be inferred

e.g., **20130319** does not necessarily mean the capture represents the page on **March 19, 2013**

* https://www.w3.org/DesignIssues/Axioms.html#opaque

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

16

# Backtracking - Save My Things

**Mat Kelly, PhD**
https://matkelly.com

Save My Web Things For Me and Others but Not All
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

17

# Preserving Content That *Was* Behind Authentication

- By **reference** (URI) vs. **By-Value** (representation)*
- Proxy Approaches
  - MITM: warcprox, Webrecorder (link these, screenshots)
- Browser-based approach (by-value)
  - **WARCreate**



**WARCreate**
Google Chrome Extension

Archive facebook.com!

\* Previously funded by grant, *"Archive What I See Now"*

NATIONAL ENDOWMENT FOR THE HUMANITIES

*#HK-50181-14*

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

18

# Archival Replay

- Authentication mechanism is not exhibited once a page is preserved
  - For example, entering credentials for facebook.com does not send the request to facebook.com, which might not exist in the future.
- By value stores the resource representation
- This content might contain personally identifiable or sensitive information
  - sharing the captures to others needs strategic consideration.
- URI-M alone for access could make the memento publicly accessible (bad)

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

19

# Sharing Captures



Check out my capture at
http://archive.alice.com/web/202002031205/https://facebook.com

That probably should not be publicly available...

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

20

# Memento TimeMaps

- Coalescence URI-Ms by the original URI (URI-R) they represented on the live web
  - **https://cnn.com**
  - **http://www.cnn.com**
  - **https://cnn.com/index.php**
  - **http://www2.cnn.com:80/homepage.html**

  **Various identifiers (original URI, URI-Rs) of the CNN homepage over time**

- ...might all represent the same site in the past
- "Canonicalization" helps to normalize these variants for retrieval

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

21

# Cross-archive Association

**Live web "original URI" (URI-R)**

<**https://cnn.com**>; rel="original",
<https://memgator.cs.odu.edu/timemap/link/https://cnn.com>; rel="self"; type="application/link-format",
<https://arquivo.pt/wayback/20000620180259mp_/http://cnn.com/>; rel="first memento"; datetime="Tue, 20 Jun 2000 18:02:59 GMT",
<http://web.archive.org/web/20000621011731/http://cnn.com:80/>; rel="memento"; datetime="Wed, 21 Jun 2000 01:17:31 GMT",
<http://wayback.vefsafn.is/wayback/20000621140928/http://cnn.com/>; rel="memento"; datetime="Wed, 21 Jun 2000 14:09:28 GMT",
<http://web.archive.org/web/20000621140928/http://cnn.com:80/>; rel="memento"; datetime="Wed, 21 Jun 2000 14:09:28 GMT",
<https://arquivo.pt/wayback/20000706192838mp_/http://cnn.com/>; rel="memento"; datetime="Thu, 06 Jul 2000 19:28:38 GMT",
<http://wayback.vefsafn.is/wayback/20000706192838/http://cnn.com/>; rel="memento"; datetime="Thu, 06 Jul 2000 19:28:38 GMT",
<http://web.archive.org/web/20000706192838/http://cnn.com:80/>; rel="memento"; datetime="Thu, 06 Jul 2000 19:28:38 GMT",
<http://wayback.archive-it.org/all/20150601215659/http://cnn.com/>; rel="memento"; datetime="Mon, 01 Jun 2015 21:56:59 GMT",
...

INTERNET ARCHIVE
WayBackMachine

ARQUIVO.PT

ARCHIVE-IT

Vefsafn.is

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

22

# Memento Aggregation

- A web archive is unlikely to provide links to other archives' captures (mementos)
- Aggregation services send a request for a **URI-R** (live web URI) to multiple archives and "aggregate" the results



Show me captures for **cnn.com**

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

23

# Who Runs the Aggregators?

- Los Alamos National Laboratory hosts their Memento aggregator
  - e.g., http://timetravel.mementoweb.org/timemap/link/http://cnn.com
- The set of archives it uses is unchangeable by the querying user
- When new archives come about, to be included:
  - The archive's **MUST** be Memento compliant
  - The archive's endpoints **MUST** be manually added by someone at LANL
- The captures on *your machine* will never be included in these results.

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

24

# Example TimeMap from a Memento Aggregator

```
<http://drexel.edu>; rel="original",
<http://some.aggregator.com/timemap/link/http://drexel.edu>; rel="self"; type="application/link-format",
<http://web.archive.org/web/20060514123511/http://drexel.edu:80/>; rel="first memento"; datetime="Sun, 14 May 2006
12:35:11 GMT",
<http://web.archive.org/web/20060516213852/http://www.drexel.edu/>; rel="memento"; datetime="Tue, 16 May 2006 21:38:52
GMT",
...
<http://web.archive.org/web/20180128152125/http://drexel.edu>; rel="memento"; datetime="Sun, 28 Jan 2018 15:21:25 GMT",
<http://web.archive.org/web/20180319141920/http://drexel.edu/>; rel="last memento"; datetime="Mon, 19 Mar 2018 14:19:20
GMT",
<http://some.aggregator.com/timemap/link/http://drexel.edu>; rel="timemap"; type="application/link-format",
<http://some.aggregator.com/timemap/json/http://drexel.edu>; rel="timemap"; type="application/json",
<http://some.aggregator.com/timemap/cdxj/http://drexel.edu>; rel="timemap"; type="application/cdxj+ors",
<http://some.aggregator.com/timegate/http://drexel.edu>; rel="timegate"
```
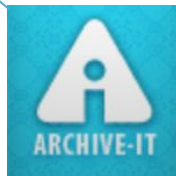
Original URI (URI-R) | Other TimeMaps (URI-Ts) | TimeGate (URI-G) | Relative Relations

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

25

# MemGator

- *Open Source* Memento aggregator software
- Individuals can easily deploy an aggregator *onto their system*.
- Allows users to define which archives are queried
  - provides a comprehensive set of defaults

running on

S. Alam and M. Nelson, MemGator - A Portable Concurrent Memento Aggregator: Cross-Platform CLI and Server Binaries in Go, Proceedings of JCDL 2016, pp. 243-244, 2016.

https://github.com/oduwsdl/memgator

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

26

# A Framework for Aggregating Private and Public Web Archives

- **Aggregation with private web archives**
- **Client-side archive specification**

- **Authentication layer to systematically interface with private Web archives**

- **Archival negotiation in dimensions beyond time**

more information:
M. Kelly et al., "A Framework for Aggregating Private and Public Web Archives," In Proceedings of the ACM/IEEE Joint Conference on Digital Libraries (JCDL), 2018, pp. 273-282.

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

27

# Memento Meta Aggregator

- Enhanced MemGator to enable more powerful client-side querying
  - Includes client specifying set of archives to aggregate
- Empowers querying client to:
  - optionally specify the set of archives aggregated, beyond what might be done by default
  - Specify archival precedence
  - Interface with mementos that require access beyond simply dereferencing URI-M
  - ...and more.

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

28

# Personal Archive Aggregation: Alice Saves the Web

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

29

# Alice Wants to See Her Captures Temporally Inline



at $t_A$

at $t_{D \to Z}$

*Personal Archive Aggregation*

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

30

# Mementity Dynamics - Alice & Her Archives



*Personal Archive Aggregation*

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

31

# Alice Deploys MMA

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

32

# Carol Asks MMA for CNN

# MMA Returns CNN Memento

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

34

# Carol Wants to Aggregate Her Own Captures

**Mat Kelly, PhD**
https://matkelly.com

Sav~~~ r Me and Others but Not All
Master's of I~~~ ebinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

35

# Carol Spawn Her Own MMA to Access Alice's Web Archive and Alice's MMA

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

36

# Carol Asks Her MMA For CNN

**Mat Kelly, PhD**
https://matkelly.com

o Things For Me and Others but Not All
information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

37

# Carol's MMA Returns CNN Mementos from the Sources She Specified

# Hooray, Aggregation!

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

39

# Regulating Access to Private Captures

- Some form of access control should be used to regulate private web archives to prevent information leakage and exposure of sensitive information.
- Private Web Archive Adapter
  - Enables systematic (OAuth2) based access control to private web archives
- Delegates authentication/authorization to a separate "entity"
  - archives remain functionally cohesive
  - decouples archives from authentication role to facilitate interoperability

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

40

# Dereferencing Beyond Simple URI Access

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

41

# Adaptation of OAuth2 to Web Archives

Request URI-M →

← HTTP 401, see PWAA

**PRIVATE WEBARCHIVE**

Request URI-M →

← Send credentials

Request URI-M w/ credentials →

← ACK URI-M, token

Request URI-M, token →

← Dereferenced URI-M

**PRIVATE WEBARCHIVE**

- Auth Layer for to encourage Private Web archive aggregation

- Typical OAuth 2.0 flow

- Auth role cohesive to PWAA

- Persistent access through tokenization

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

42

# Using Auth Tokens for Private Web Archive Access

Token grants can be:

- Shared
- Disavowed
- Limit scope



GET URI$_2$
Scope: /
Token: 4f33c64

GET URI$_3$
Scope: /public/
Token: 98ac9de

GET URI$_1$
Scope: /
Token: 4f33c64

ALICE'S PRIVATE ARCHIVE

GET URI$_2$
Scope: /
Token: 2265eef3

PRIVATE WEB ARCHIVE
ADAPTER

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing & Informatics

43

# Using Auth Tokens for Private Web Archive Access

Token grants can be:

- Shared
- Disavowed
- Limit scope



4f33c64, URI1, /
4f33c64, URI2, /
98ac9de, URI3, /public/
2265eef3, URI2, /

ALICE'S PRIVATE ARCHIVE

PRIVATE WEB ARCHIVE
ADAPTER

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

# Richer Querying of Web Archives Makes Them More Useful

Memento introduces negotiation in **time**, also useful would be negotiation:

- private/public captures, only, first, if-and-only-if conditional
- Memento quality above a threshold*
  - Crawlers sometimes miss necessary resource representations!
- Only mementos that are not redirects†
  - Over 80% of the google.com URI-Ms at Internet Archive are HTTP 3XX redirects!

* J. Brunelle et al., "Not All Mementos Are Created Equal: Measuring the Impact of Missing Resources," International Journal on Digital Libraries (IJDL), 16(3), pp. 283-301, 2015.

† M. Kelly et al., "Impact of URI Canonicalization on Memento Count," Technical Report arXiv:1703.03302, 2017.

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

45

# There exists a need in querying web archives to...

- Filter results on the server side (cf. making users sort through verbose results)
- Specify the sources used in aggregation (e.g., include *my* captures)
- Allow aggregators to optionally comply
- Allow clients/users of web archives to express these PREFERENCES

- The **StarGate** "concepts" facilitates
*content negotiation with web archives in dimensions beyond time*

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

46

# Despite Your Sources, Use A Subset

Get URI-Ms for URI-R only from personal Web archives
*privateOnly*

①

②

③

④

**ACCESS ATTRIBUTE**

# Many Nearly Identical

- Temporal adjacent captures might be similar or identical
- Finding a summary of representations in time can be expensive
- Generating screenshots of all would be wasteful

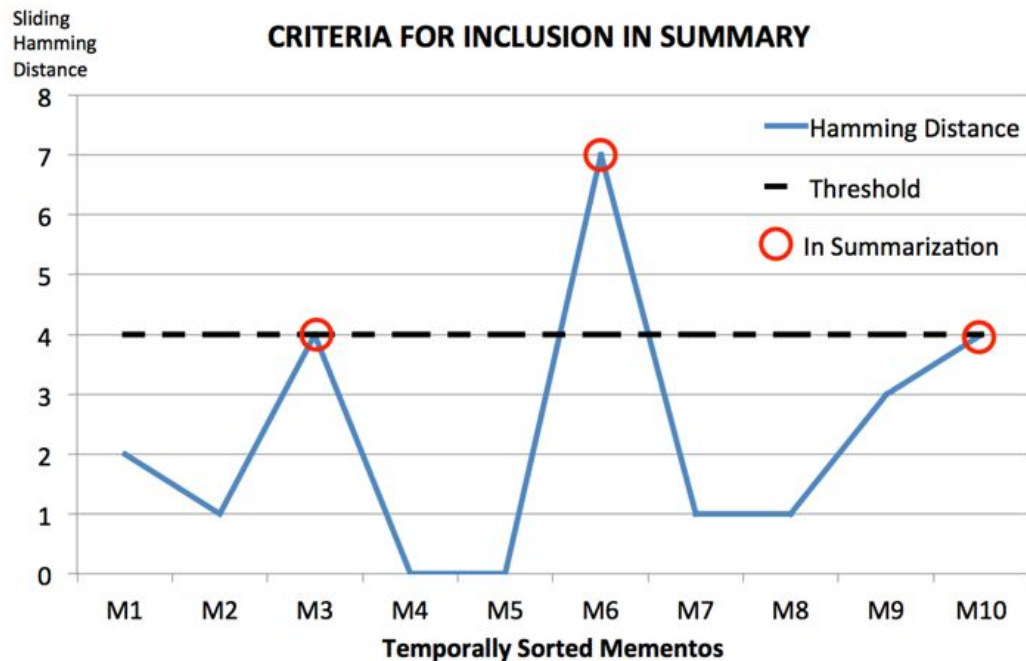**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

48

# Identifying Similarity by Calculating Hamming Distance

**HAMMING DISTANCE**

| | | |
|---|---|---|
| HTML of apple.com March 3, 2008 | c39f0abc...b9 | N/A pivot |
| HTML of apple.com March 5, 2008 | c39d0abc...c9 | 2 |
| HTML of apple.com April 12, 2008 | c39d0abc...b9 | 1 |
| HTML of apple.com October 4, 2008 | c770ad1b...b9 | 7 |

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

49

# Use dissimilarity from pivot as basis for inclusion



Sliding Hamming Distance

**CRITERIA FOR INCLUSION IN SUMMARY**

Legend:
- Hamming Distance
- Threshold
- In Summarization

**Temporally Sorted Mementos**

X-axis: M1, M2, M3, M4, M5, M6, M7, M8, M9, M10

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

50

# StarGate

- Allows for content negotiation with archives in dimensions beyond time
  - e.g., Bob asks for only captures from the set of archives that are of a sufficient quality of capture.

Get URI-Ms for URI-R of good quality that are unique
$M_D < 0.25$, unique(simhash)

Retrieved StarMap

Abbreviated StarMap with filtering applied

CONTENT-BASED ATTRIBUTE

&

DERIVED ATTRIBUTE

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

51

# Summation

**SAVE MY WEB THINGS FOR ME**

Preserve content behind authentication

**...AND OTHERS**

Use an Aggregator for others' access

**...BUT NOT ALL**

Use a PWAA for Access Control

...and further relevant topics beyond the title!

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

52

# Open-Ended Research Efforts

- Efficient querying in numerous dimensions to mitigate the combinatorial explosion of permutations
- Make software easy, intuitive, native, self-hosted
  - Users that preserve might not be technical, this should not be a barrier
- Facilitate resilience -- allow users to have captures **safely** distributed
  - See InterPlanetary Wayback https://github.com/oduwsdl/ipwb
- Introduce more sophisticated, stronger ways to regulate access

**Mat Kelly, PhD**
https://matkelly.com

**Save My Web Things For Me and Others but Not All**
Master's of Information Webinar Series - February 4, 2020

DREXEL UNIVERSITY
College of
Computing &
Informatics

53