

# Where Did the Web Archive Go?

(Published in TPDL 2021)

Mohamed Aturban<sup>1</sup>   Michael L. Nelson<sup>2</sup>   Michele C. Weigle<sup>2</sup>

<sup>1</sup>Columbia College, Columbia MO

<sup>2</sup>Old Dominion University, Norfolk VA

Presented by – Yasith Jayawardana (@yasithdev)

Old Dominion University, Norfolk VA

CS895 Web Archiving Forensics, Fall 2022

# Outline

1 Recap on the Memento Protocol

2 Testing the Fixity of Mementos

3 Analysis

# Headers in the Memento Protocol

An HTTP protocol extension for **temporal** content negotiation of web resources<sup>1</sup>

## Request Headers

**Accept-Datetime:** <approximate datetime of the memento>

---

<sup>1</sup><https://www.slideshare.net/hvdsomp/memento-101>

# Headers in the Memento Protocol

An HTTP protocol extension for **temporal** content negotiation of web resources<sup>1</sup>

## Request Headers

**Accept-Datetime:** <approximate datetime of the memento>

## Response Headers

**Memento-Datetime:** <exact datetime of the memento>

**Vary:** Accept-Datetime, ...

**Link:** <URL>; rel="<original/timegate/timemap/memento>", ...

**Location:** <URI-M of the memento>

---

<sup>1</sup><https://www.slideshare.net/hvdsomp/memento-101>

# Headers in the Memento Protocol

An HTTP protocol extension for **temporal** content negotiation of web resources<sup>1</sup>

## Request Headers

**Accept-Datetime:** <approximate datetime of the memento>

## Response Headers

**Memento-Datetime:** <exact datetime of the memento>

**Vary:** Accept-Datetime, ...

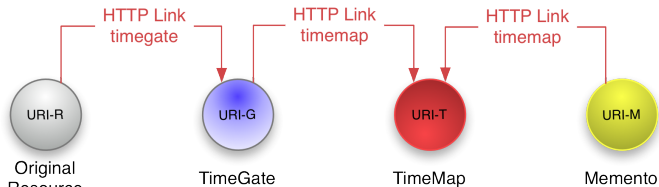
**Link:** <URL>; rel="original/timegate/timemap/memento", ...

**Location:** <URI-M of the memento>

- Datetime should comply with **RFC 9110** section 5.6.7
- For **Link** header entries
  - If rel="memento" ⇒ datetime="" is mandatory
  - If rel="timemap" ⇒ type="" is recommended

<sup>1</sup><https://www.slideshare.net/hvdsomp/memento-101>

# Using the Memento Protocol



**Figure:** Components of the Memento Protocol ([mementoweb.org](http://mementoweb.org))

# Using the Memento Protocol

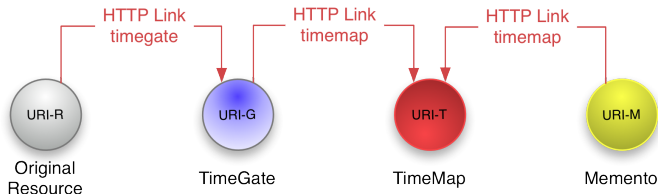


Figure: Components of the Memento Protocol (mementoweb.org)

## Requesting with Accept-Datetime Headers

(req) **URI-R** → (res) Link: **URI-G**  
(req) **URI-G** → (res) Link: **URI-R**, **URI-T**, Location: temporally closest **URI-M**  
(req) **URI-T** → (res) list of **URI-M**  
(req) **URI-M** → (res) Link: **URI-R**, **URI-T**, **URI-G**

## Example URI-M

<http://www.collectionscanada.gc.ca/webarchives/20060208075019/http://www.cdc.gov/>

# Outline

- 1 Recap on the Memento Protocol
- 2 Testing the Fixity of Mementos
- 3 Analysis



# Not even archived resources are permanent<sup>2</sup>

The web archives might...

**Go Defunct** (funding issues, etc.)



Figure: Now Defunct Web Archives

- Internet Memory Foundation (formerly European Archive Foundation)
- PRONI (now in archive-it.org)

---

<sup>2</sup>[https://en.wikipedia.org/wiki/Wikipedia:List\\_of\\_web\\_archives\\_on\\_Wikipedia](https://en.wikipedia.org/wiki/Wikipedia:List_of_web_archives_on_Wikipedia)

# Not even archived resources are permanent<sup>2</sup>

The web archives might...

## Go Defunct (funding issues, etc.)



Figure: Now Defunct Web Archives

- Internet Memory Foundation (formerly European Archive Foundation)
- PRONI (now in archive-it.org)

## Change Infrastructure

```
$ curl -i https://perma-archives.org
HTTP/2 301
date: Mon, 24 Oct 2022 03:03:52 GMT
location: https://perma.cc/
```

Figure: perma-archives.org is now perma.cc

- perma-archives.org  
→ perma.cc
- collectionscanada.gc.ca  
→ bac-lac.gc.ca  
→ library-archives.canada.ca

<sup>2</sup>[https://en.wikipedia.org/wiki/Wikipedia:List\\_of\\_web\\_archives\\_on\\_Wikipedia](https://en.wikipedia.org/wiki/Wikipedia:List_of_web_archives_on_Wikipedia)

# 16,627 Mementos, 3,698 URI-Rs, 17 Public Web Archives

Collected from 2017-11 to 2019-01

**Table 1: A set of 17 public web archives.**

Archive URI	Archive Name	Purpose
<a href="http://swap.stanford.edu">swap.stanford.edu</a>	Stanford Web Archive Portal	General
<a href="http://web.archive.org">web.archive.org</a>	The Internet Archive	General and on-demand
<a href="http://archive.bibalex.org">archive.bibalex.org</a>	Bibliotheca Alexandrina's Internet Archive	National
<a href="http://arquivo.pt">arquivo.pt</a>	The Portuguese Web Archive (PWA)	National
<a href="http://collectionscanada.gc.ca">collectionscanada.gc.ca</a>	Library and Archives Canada	National
<a href="http://digar.ee">digar.ee</a>	The Estonian Web Archive	National
<a href="http://nationalarchives.gov.uk">nationalarchives.gov.uk</a>	The National Archives	National
<a href="http://vefsafn.is">vefsafn.is</a>	The Icelandic Web Archive	National
<a href="http://webarchive.loc.gov">webarchive.loc.gov</a>	Library of Congress Web Archives	National
<a href="http://webarchive.org.uk">webarchive.org.uk</a>	The UK Web Archive (UKWA)	National
<a href="http://webarchive.proni.gov.uk">webarchive.proni.gov.uk</a>	Public Record Office of Northern Ireland (PRONI)	National
<a href="http://webharvest.gov">webharvest.gov</a>	Congressional & Federal Government Web Harvests	National
<a href="http://archive-it.org">archive-it.org</a>	Archive-It - Web Archiving Services for Libraries and Archives	On-demand
<a href="http://archive.is">archive.is</a>	Archive.is	On-demand
<a href="http://perma.cc">perma.cc</a>	Perma.cc	On-demand
<a href="http://webcitation.org">webcitation.org</a>	WebCite	On-demand
<a href="http://europarchive.org">europarchive.org</a>	The European Archive	Organizational

**Figure:** Public web archives used in the study (source: "Collecting 16K archived web pages from 17 public web archives", Aturban et al., 2019)

# Outline

- 1 Recap on the Memento Protocol
- 2 Testing the Fixity of Mementos
- 3 Analysis

## 4 Web Archives (1,981 Mementos) Changed Base URIs

With no redirection to their new base URIs

- **Library and Archives Canada** (351 Mementos)  
*49 changed, 2 missing*
- **National Library of Ireland** (979 Mementos)  
*192 changed, 0 missing*
- **Public Record Office of Northern Ireland** (469 Mementos)  
*114 changed, 0 missing*
- **Perma.cc** (182 Mementos)  
*164 changed, 18 missing*

### Overall Statistics

- **Total:** 1,981
- **Changed:** 537 (27.11%)
- **Missing:** 20 (1.01%)

# Breakdown of Changed (537) and Missing (20) Mementos

## Criteria – “Memento-Datetime” Header / HTTP Status / URI-R Changes

Table 1: Web archive changes based on how mementos changed. The number of missing mementos is shown in **bold**.

Original archive → New archive	Same Memento-Datetimes?	Same status codes?	Same URI-Rs?	URI-Ms
collectionscanada.gc.ca → bac-lac.gc.ca	Yes	Yes	Yes	302
	NO	Yes	Yes	<b>28</b>
	NO	Yes	NO	<b>18</b>
	NO	NO	Yes	<b>1</b>
	NO	NO	NO	<b>2</b>
europarchive.org/NLI → internetmemory.org/NLI	Yes	Yes	Yes	979
internetmemory.org/NLI → archive-it.org	Yes	Yes	Yes	787
	Yes	NO	Yes	<b>1</b>
	Yes	NO	NO	<b>2</b>
	NO	Yes	Yes	<b>184</b>
	NO	Yes	NO	<b>5</b>
proni.gov.uk → archive-it.org	Yes	Yes	Yes	355
	Yes	NO	Yes	<b>2</b>
	NO	Yes	Yes	<b>106</b>
	NO	Yes	NO	<b>6</b>
perma-archives.org → perma.cc	NO	Yes	Yes	<b>164</b>
	NO	NO	NO	<b>18</b>

## Why did “Memento-Datetime” Headers Change?

The new archive may have post-processed the original WARC files

# Library and Archives Canada

351 mementos, 49 changed, 2 missing

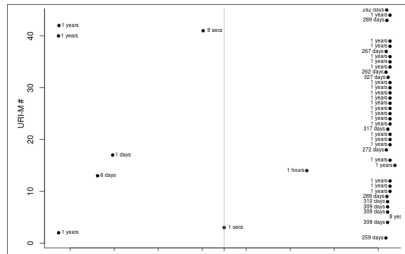
- Old mementos incorrectly redirected to the main webpage
- New mementos were accessible from webarchive.bac-lac.gc.ca
- Some new mementos didn't match old status codes / headers

```
http://www.collectionscanada.gc.ca/webarchives/20070220181041/http://www.berlin.gc.ca/ (302)
http://www.collectionscanada.gc.ca/webarchives/20070220181041/http://www.dfait-maeci.gc.ca/canadaeuropa/germany/ (302)
http://www.collectionscanada.gc.ca/webarchives/20070220181204/http://www.dfait-maeci.gc.ca/canadaeuropa/germany/ (302)
http://www.collectionscanada.gc.ca/webarchives/20070220181204/http://www.international.gc.ca/global/errors/404.asp?404%3Bhttp://www.dfait-maeci.gc.ca/canadaeuropa/germany/ (404)
```

Fig. 4: The HTTP status codes of the URI-M <http://www.collectionscanada.gc.ca/webarchives/20070220181041/http://www.berlin.gc.ca/> from the original archive.

```
http://webarchive.bac-lac.gc.ca:8080/wayback/20070220181041/http://www.berlin.gc.ca/ (Redirect by JavaScript (JS))
http://webarchive.bac-lac.gc.ca:8080/wayback/20070220181041/http://www.dfait-maeci.gc.ca/canadaeuropa/germany/ (Redirect by JS)
http://webarchive.bac-lac.gc.ca:8080/wayback/20070220181204/http://www.international.gc.ca/global/errors/404.asp?404%3Bhttp://www.dfait-maeci.gc.ca/canadaeuropa/germany/ (Redirect by JS)
http://webarchive.bac-lac.gc.ca:8080/wayback/20071115025620/http://www.international.gc.ca/canada-europa/germany/ (302)
http://webarchive.bac-lac.gc.ca:8080/wayback/20071115023828/http://www.international.gc.ca/canada-europa/germany/ (200)
```

Fig. 5: The HTTP status codes of the URI-M <http://webarchive.bac-lac.gc.ca:8080/wayback/20070220181041/http://www.berlin.gc.ca/> from the new archive.

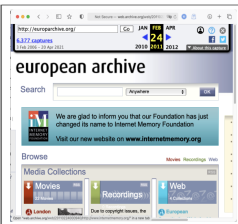


# National Library of Ireland Changed Infrastructure Twice

- (1) [europarchive.org/NLI](http://europarchive.org/NLI) → [internetmemory.org/NLI](http://internetmemory.org/NLI)
- (2) [internetmemory.org/NLI](http://internetmemory.org/NLI) → [wayback.archive-it.org/10702](http://wayback.archive-it.org/10702)



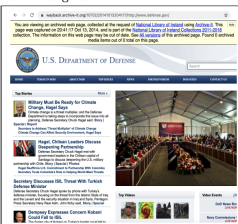
(a) In [europarchive.org](http://europarchive.org)



(b) European Archive announcing their name change to IMF



(c) In [internetmemory.org](http://internetmemory.org)



(d) In [archive-it.org](http://archive-it.org)



979 mementos, 192 changed, 0 missing

- ```
$ curl --head --location --silent http://wayback.archive-it.org
/10702/20121221162201/http://bbc.co.uk/news/ | egrep -i "(HTTP/|
location:|^Memento-Datetime)"
```

Location: /10702/20121221163248/http://www.bbc.co.uk/news/  
HTTP/1.1 200 OK  
Memento-Datetime: Fri, 21 Dec 2012 16:32:48 GMT

```
$ curl --head --silent http://wayback.archive-it.org
/10702/20121223031837/http://www2008.org/
```

Date: Thu, 05 Sep 2019 08:28:27 GMT

A scatter plot showing the distribution of UPRIM# values for 1000 random samples across various time scales. The y-axis is labeled 'UPRIM#' and ranges from 0 to 150. The x-axis is labeled 'time scale' and ranges from -115 days to +115 days, with a central vertical line at 0. The data points are labeled with their corresponding time scales, such as '5 mins', '1 day', '10 mins', etc. The distribution is roughly symmetric around the 0 time scale, with a higher density of points between -10 and +10 minutes.

Fig. 8: Difference between the Memento-Datetimes for URI-Ms from [internetmemory.org](http://internetmemory.org) and the corresponding URI-Ms from [archive-it.org](http://archive-it.org).

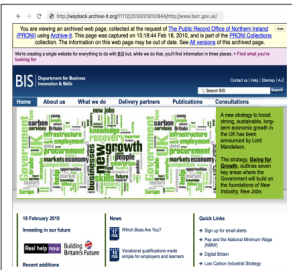
# Public Record Office of Northern Ireland (PRONI)

469 mementos, 114 changed, 0 missing

- webarchive.proni.gov.uk → archive-it.org/collections/11112



(a) In [webarchive.proni.gov.uk](http://webarchive.proni.gov.uk)



(b) In [archive-it.org](http://archive-it.org)

```
$ curl --head --location http://webarchive.proni.gov.uk/20100218151844/  
http://www.berr.gov.uk/
```

```
HTTP/1.1 302 Found  
Cache-Control: no-cache  
Content-Length: 0  
Location: https://webarchive.proni.gov.uk/20100218151844/http://www.berr.  
gov.uk/  
HTTP/2 404  
date: Fri, 20 Sep 2019 08:13:45 GMT  
server: Apache/2.4.18 (Ubuntu)  
content-type: text/html; charset=iso-8859-1
```

Fig. 12: The HTTP status codes of the URI-M <http://webarchive.proni.gov.uk/20100218151844/http://www.berr.gov.uk/> from the original archive [webarchive.proni.gov.uk](http://webarchive.proni.gov.uk/).

```
# With mod_rewrite  
RewriteEngine on  
RewriteRule "^/(\\d{14})/(.+)"  
http://wayback.archive-it.org/11112/$1/$2 [L,R=301]
```

Fig. 14: The Apache `mod_rewrite` rules that can be used to handle redirects from [webarchive.proni.gov.uk](http://webarchive.proni.gov.uk) to [archive-it.org](http://archive-it.org).

- `perma-archives.org/warc/<datetime>` → `perma.cc/<id>`
- Since 2020-02 perma.cc supports the Memento protocol
- New mementos don't have the 14-digit memento-datetime in the URI-M  
e.g., `https://perma.cc/T8U2-994F`

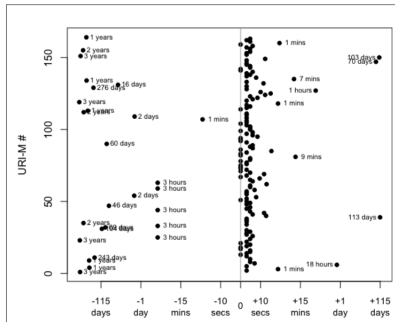


Fig. 16: Difference between the Memento-Datetimes for the long-form Perma.cc URI-Ms and the corresponding short-form URI-Ms.

## Investigation

- Goal – detect variations and changes in replaying individual mementos
- Sample – 16,627 mementos between 2017-11-01 and 2019-01-01 from 17 public web archives

## Findings

- 4 web archives changed their base URIs and left no *machine-readable* method to locate their new base URIs
- 537 of 1,981 mementos from these web archives were impacted
  - 517 rediscovered with changes to their Memento-Datetime, URI-R, or HTTP status code
  - 20 could not be found at all
- If a Web archive changes its base URL, setting rewrite rules may help to avoid broken URI-Ms