# Right HTML, Wrong JSON: Challenges in Replaying Archived Webpages Built with Client-Side Rendering

**Michele C. Weigle**, Michael L. Nelson

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

Sawood Alam, Mark Graham

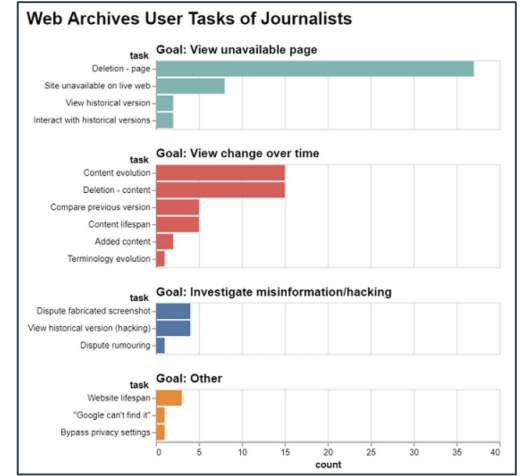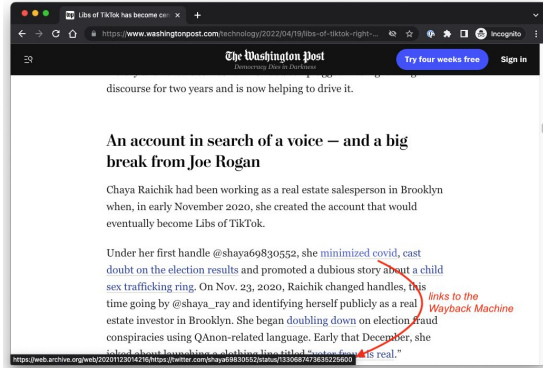Wayback Machine
Internet Archive

# Web archives have been used for journalistic and legal evidence

Taylor Lorenz, "Meet the woman behind Libs of TikTok", *Washington Post*, April 2022,
https://www.washingtonpost.com/technology/2022/04/19/libs-of-tiktok-right-wing-media/



*Last Week Tonight*, Mar 18, 2018

"... conference addressing, understanding and preventing homosexuality."





**Web Archives User Tasks of Journalists**

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

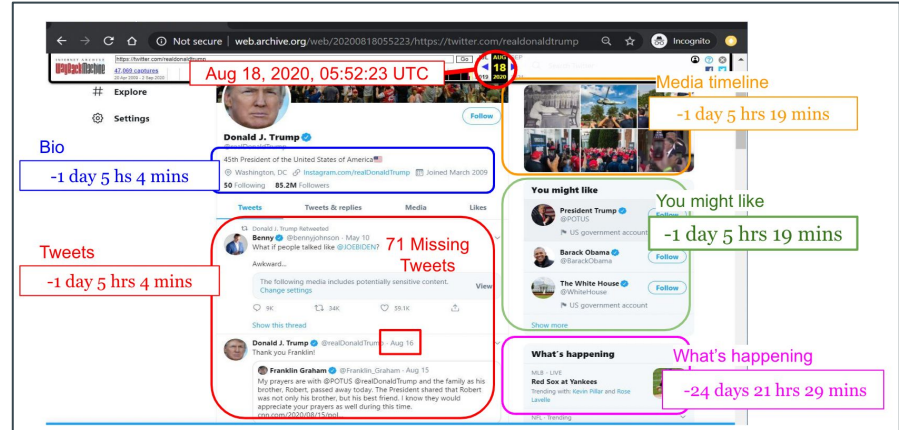

**Telewizja Polska USA, Inc. v. Echostar Satellite Corp.** (2004)

Nicholas Taylor and Joe Dugan, "Evidentiary Use of the Temporal Web", Mar 2022,
https://nullhandle.org/pdf/2023-03-02_evidentiary_use_of_the_temporal_web.pdf#page=74
Nicholas Taylor, "A Brief Primer on Using Web-Archived Evidence", June 2023,
https://nullhandle.org/blog/2023-06-07-a-brief-primer-on-using-web-archived-evidence.html

# However, webpages replayed from web archives are not always exactly what was observed



Root Memento-Datetime: 2004-12-09T19:09:26

−15 hours

missing

+9 hours

+9 months

S. Ainsworth, M.L. Nelson, H. Van de Sompel, "Only One Out of Five Archived Web Pages Existed as Presented", ACM Hypertext 2015, https://ws-dl.blogspot.com/2015/12/2015-12-08-evaluating-temporal.html

K. Garg, H.R. Jayanetti, S. Alam, M.C. Weigle, and M.L. Nelson, "Replaying Archived Twitter: When your bird is broken, will it bring you down?," JCDL 2021, https://arxiv.org/abs/2108.12092



Aug 18, 2020, 05:52:23 UTC

Media timeline
−1 day 5 hrs 19 mins

Bio
−1 day 5 hs 4 mins

You might like
−1 day 5 hrs 19 mins

Tweets
−1 day 5 hs 4 mins

71 Missing Tweets

What's happening
−24 days 21 hrs 29 mins



On 2018-02-28          On 2018-07-08          On 2018-08-25

*same memento (capture), replayed on different dates*

M. Aturban, M. Klein, H. Van de Sompel, S. Alam, M.L. Nelson, and M.C. Weigle, "Hashes are not suitable to verify fixity of the public archived web", *PLOS ONE*, June 2023, https://doi.org/10.1371/journal.pone.0286879

# We found temporal violations in some CNN.com mementos



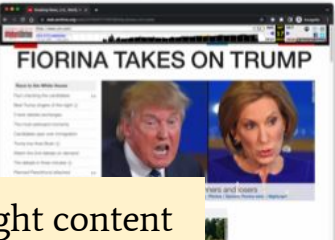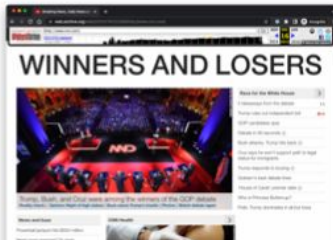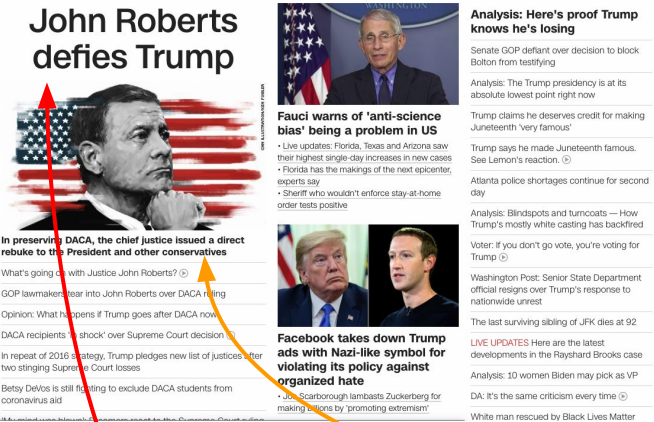| Sep 17, 2015 | Oct 14, 2015 | Dec 16, 2015 | Apr 20, 2016 | Jul 29, 2016 |

wrong content

right content

# CNN.com's front page uses client-side rendering, loading many JSON zone-manager.izl files that contain HTML code

Snippet from
https://wayback.archive-it.org/4887/20200618234850if_/https://www.cnn.com/data/ocs/section/index.html:**homepage1-zone-1**/views/zones/common/zone-manager.izl, which builds the "Hero" headlines.

```
{"izlData":{"cards":[]},"html":"<section class=\"zn zn-homepage1-zone-1 zn-left-flui
zn-left-fluid-share zn-has-multiple-containers zn-has-3-containers\" data-eq-pts=\"x
1100\" id=\"homepage1-zone-1\" data-vr-zone=\"zone-0-0\" data-zone-label=\"Hero\" da
class=\"l-container zn-top__label\"></div><div class=\"l-container zn-top__banner\">
class=\"zn-top__background\"></div></div><div class=\"l-container\"><div class=\"zn_
zn__column--idx-0\"><ul class=\"cn cn-list-hierarchical-xs cn--idx-0 cn-container_17
data-layout=\"list-hierarchical-xs\" data-vr-zone=\"home-top-col1\"><li><article cla
cd--vertical cd--has-siblings cd--has-media cd--media__image cd--has-banner\"
data-vr-contentbox=\"/2020/06/18/politics/john-roberts-supreme-court-daca-legacy/index.html\" data-eq-pts=\"xsmall: 0, small: 300,
medium: 460, large: 780, full16x9: 1100\" data-section-name=\"politics\"><a
href=\"/2020/06/18/politics/john-roberts-supreme-court-daca-legacy/index.html\" class=\"link-banner\" ><h2 class=\"banner-text
screaming-banner-text banner-text-size--char-26\" data-analytics=\"_list-hierarchical-xs_article_\">John Roberts defies
Trump</h2></a><div class=\"cd__wrapper\" data-analytics=\"_list-hierarchical-xs_article_\"><div class=\"media\"><a
href=\"/2020/06/18/politics/john-roberts-supreme-court-daca-legacy/index.html\" >
…
<span class=\"cd__headline-text\"><strong>In preserving DACA, the chief justice issued a direct rebuke to the President and other
conservatives</strong></span><span
```

# Each zone-manager.izl file builds a separate section of the CNN.com front page



homepage1-zone-1

(Hero)

https://www.cnn.com/data/ocs/section/index.html:**homepage1-zone-1**/views/zones/common/zone-manager.izl

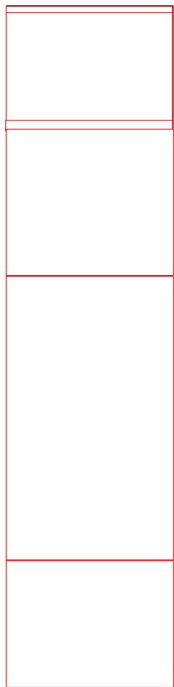# Each zone-manager.izl file builds a separate section of the CNN.com front page



homepage2-zone-1

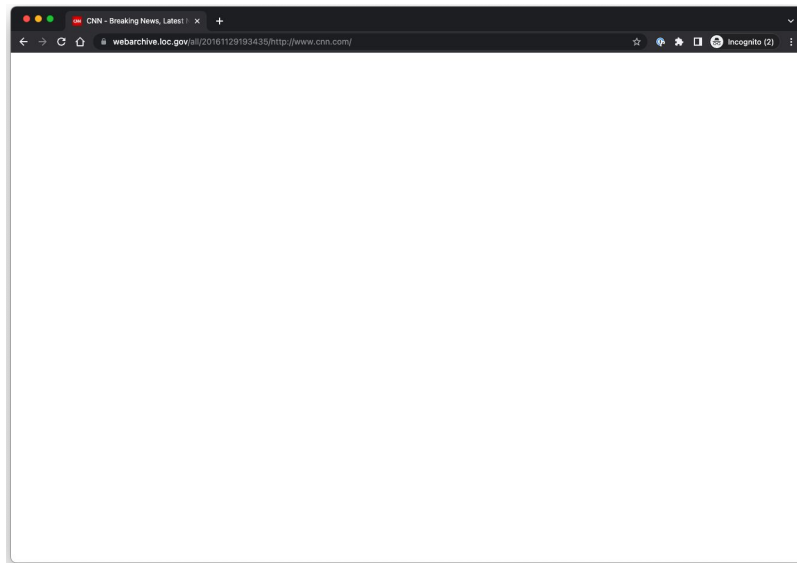(News and buzz, Privacy and tech, Life during the pandemic, CNN Business)

https://www.cnn.com/data/ocs/section/index.html:**homepage2-zone-1**/views/zones/common/zone-manager.izl

# Without executing JavaScript, conventional web archive crawlers can't archive the content (just the placeholder)



Archived CNN.com, Library of Congress (uses conventional crawler), Nov 29, 2016
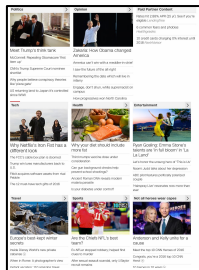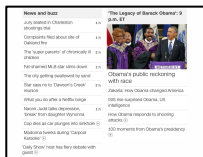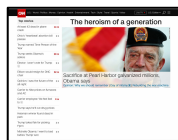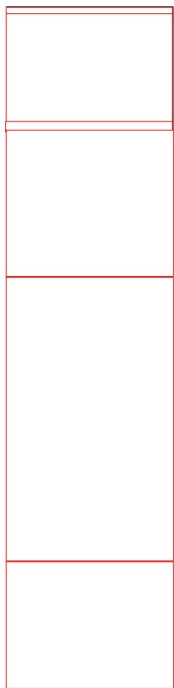


Justin F. Brunelle, Michele C. Weigle, and Michael L. Nelson, "Archival Crawlers and JavaScript: Discover More Stuff but Crawl More Slowly," JCDL 2017.
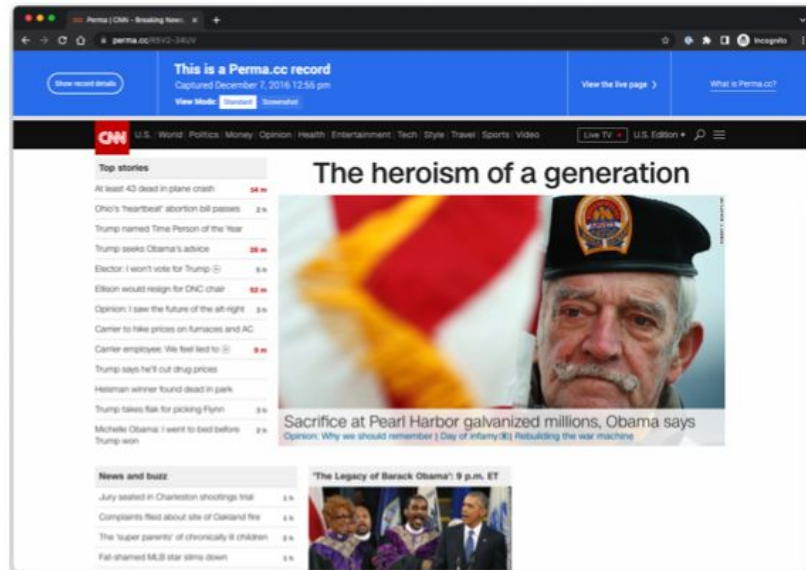
https://webarchive.loc.gov/all/20161129193435/http://www.cnn.com/

# Browser-based crawlers can capture both the placeholder and the content



Archived CNN.com, Perma.cc (uses browser-based crawler), Dec 7, 2016



https://perma.cc/R5V2-34UV

# Mixing the results of conventional and browser-based crawling can cause temporal violations

Archive-It

- 2020 - subscribers can use either Heritrix (conventional crawler) or Brozzler (browser-based crawler)
- subscribers can upload captures created with other tools (like Conifer, a browser-based crawler)

IA Wayback Machine

- Jun 2016 - incorporating captures from perma.cc
- Oct 2019 - Save Page Now service begins using Brozzler
- Large crawls still use Heritrix

https://blogs.harvard.edu/perma/2016/06/27/perma-v0-67-collaboration-internet-archive/
https://blog.archive.org/2019/10/23/the-wayback-machines-save-page-now-is-new-and-improved/

# When did CNN start using client-side rendering (CSR)?

- Used Internet Archive's CDX API to identify 200,000 CNN.com mementos with HTTP 200 OK since 2010


- Used selenium-wire headless browser to replay over 1100 CNN.com mementos from Internet Archive's Wayback Machine (IAWM) between 2013-2020 and logged requests for zone-manager files

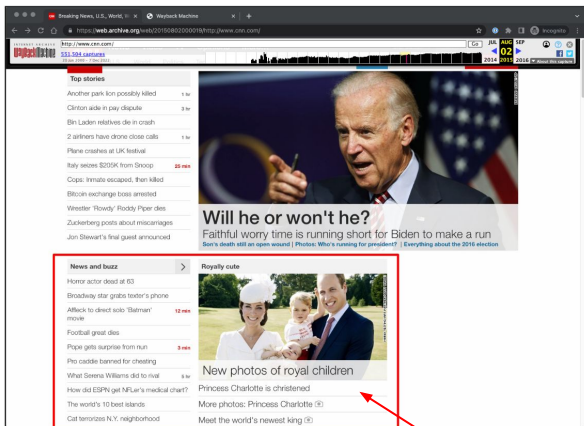https://github.com/internetarchive/wayback/blob/master/wayback-cdx-server/README.md
https://pypi.org/project/selenium-wire/

# Apr 2015: Everything but Hero zone requested via CSR

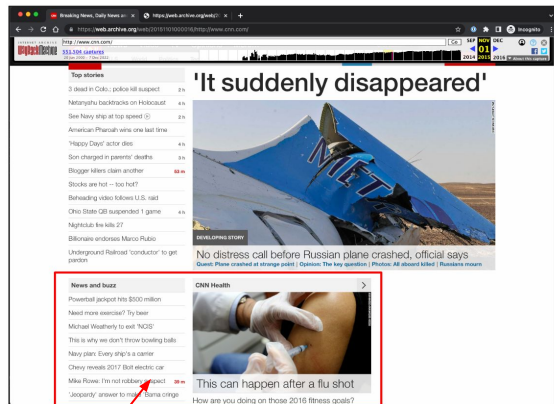| Date | Zone Content Delivery |
|---|---|
| Feb 18, 2015 | content in base HTML divided into zones |
| Apr 24, 2015 | content other than Hero zone requested via CSR in `zone-manager.html` files |
| Sep 17, 2015 | Hero zone (homepage1-zone-1) *sometimes* requested via CSR |
| Oct 18, 2016 | zone-manager format changed to `.izl.json` |
| Nov 1, 2016 | all zones requested via CSR |
| Jan 31, 2017 | zone-manager extension changed to `.izl` |

# Hero story in base HTML, but temporal violations in 2nd level content



Aug 2, 2015



Oct 1, 2015



Nov 1, 2015

from Jul 10, 2015

from Jan 6, 2016
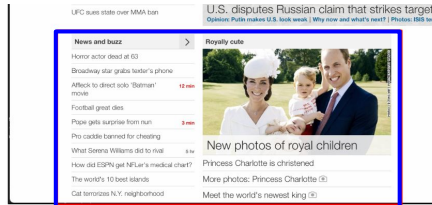
https://web.archive.org/web/20150802000019/http://www.cnn.com/

https://web.archive.org/web/20151101000016/http://www.cnn.com/

https://web.archive.org/web/20151001000018/https://www.cnn.com/

# Every replay between July 10, 2015 and Oct 8, 2015 will use same homepage2-zone-1



mementos for zone-manager files (captured via browser-based crawler)

mementos for CNN.com base HTML

# Every replay between Oct 8, 2015 and Jan 6, 2016 will use same homepage2-zone-1



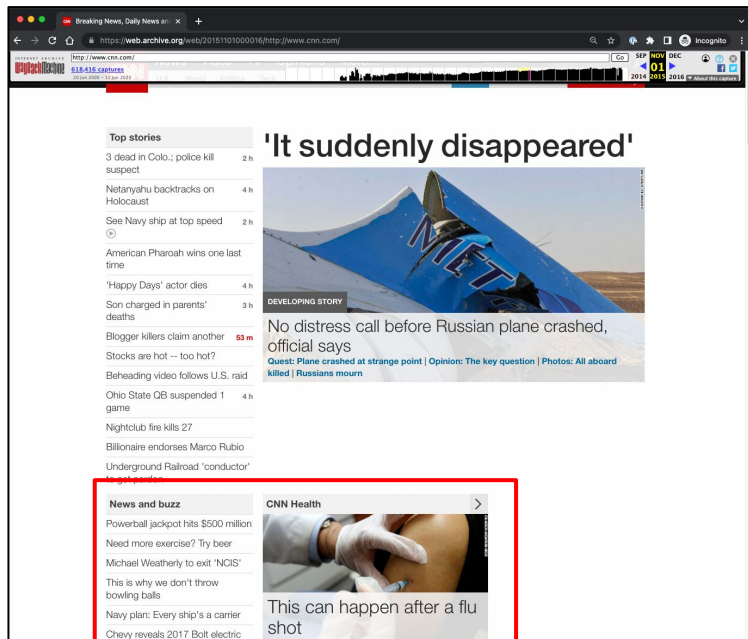mementos for zone-manager files (captured via browser-based crawler)
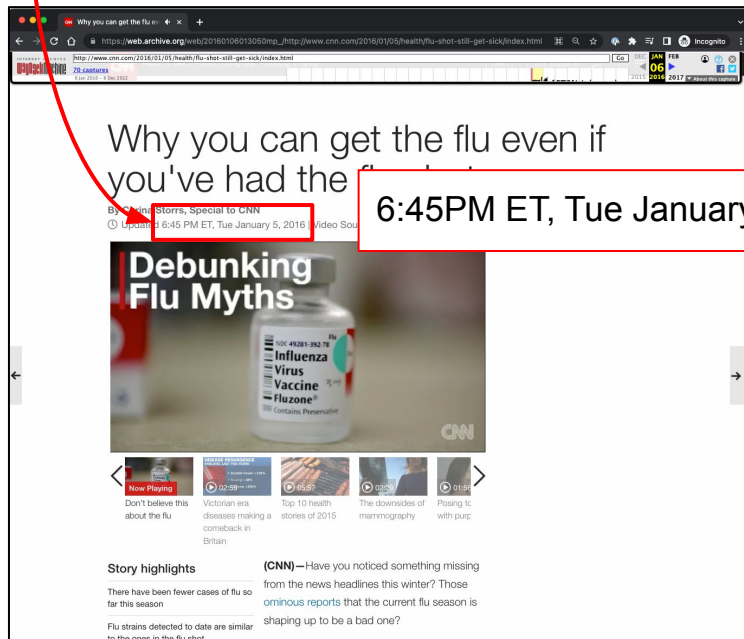
mementos for CNN.com base HTML

# Temporal violations are especially hard to detect on CNN.com

No date displayed on front page, even though individual articles do contain the date.
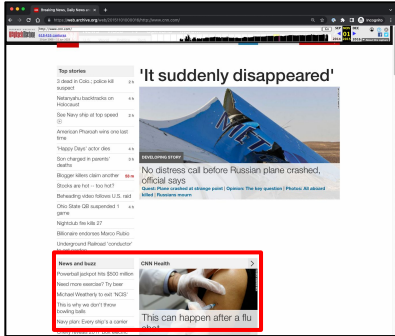
Nov 1, 2015



6:45PM ET, Tue January 5, 2016

https://web.archive.org/web/20151101000016/http://www.cnn.com/

https://web.archive.org/web/20160106013050mp_/http://www.cnn.com/2016/01/05/health/flu-shot-still-get-sick/index.html

# IAWM's "About this capture" lists temporal differences



https://web.archive.org/web/20151101000016/http://www.cnn.com/

# But "About this capture" does not include zone-manager files



zone-manager.html files should appear as +2 months 6 days

# Sep 2015 - Nov 2016: Hero zone sometimes requested via CSR

We downloaded raw HTML (using "id_") of 1316 mementos between 2010-2022 to investigate when Hero story delivered via base HTML vs. CSR.

| Date | Zone Content Delivery |
|------|----------------------|
| Feb 18, 2015 | content in base HTML divided into zones |
| Apr 24, 2015 | content other than Hero zone requested via CSR in `zone-manager.html` files |
| Sep 17, 2015 | Hero zone (homepage1-zone-1) *sometimes* requested via CSR |
| Oct 18, 2016 | zone-manager format changed to `.izl.json` |
| Nov 1, 2016 | all zones requested via CSR |
| Jan 31, 2017 | zone-manager extension changed to `.izl` |

https://web.archive.org/web/20130806040521/http://faq.web.archive.org/page-without-wayback-code/

# Hero story temporal violations

| Sep 17, 2015 | Oct 14, 2015 | Dec 16, 2015 | Apr 20, 2016 | Jul 29, 2016 |
|---|---|---|---|---|

wrong content

right content

# Hero stories *sometimes* in the base HTML



Oct 14, 2015

https://web.archive.org/web/20151014120016/https://www.cnn.com/

If the base HTML contains a `<section>` with
`id="homepage1-zone-1"`, then Hero content is in the HTML.

```
<section class="zn zn-homepage1-zone-1
zn-left-fluid-right
    -stack zn--idx-0 t-light zn-loaded zn-left-fluid-
    right-stack zn-has-multiple-containers zn-has-3-
    containers" data-eq-pts="xsmall: 0,
    medium: 460, large: 780, full16x9: 1100"
    id="homepage1-zone-1" data-vr-zone="zone-0-0"
    data-zone-label="Hero" data-containers="3">
[...]
<h2 class="banner-text js-screaming-banner-text
screaming-banner-text">CLINTON CONFIDENT, POLISHED</h2>
[...]
```

# Hero stories *sometimes* loaded via CSR

If there is no `id="homepage1-zone-1"` in the base HTML, then Hero content is in the zone-manager file.

```
[...]
<h2 class="banner-text banner-text--maximized
banner-text-size--char-23">'Fear is not
strength'</h2>
[...]
```

Content loaded from
https://web.archive.org/web/**20160729**003156/http://www.cnn.com/data/ocs/section/index.html:**homepage1-zone-1**/views/zones/common/**zone-manager.html**



Oct 14, 2015

https://web.archive.org/web/**20151014**000014/https://www.cnn.com/

verified by examining raw HTML of 1069 mementos between Apr 24, 2015 - Nov 1, 2016 and comparing to the requests logged by the headless browser

# Hero story violation if requested via CSR and no hero memento



Internet Archive, homepage1-zone-1

# But, each day had mementos with the Hero story in the HTML



Internet Archive, homepage1-zone-1

homepage1-zone-1

hero via CSR

potential mitigation: redact those mementos that load homepage1-zone-1 via CSR during this time period

cnn.com

2015-09  2015-11  2016-01  2016-03  2016-05  2016-07  2016-09  2016-11

# Nov 1, 2016: All content requested via CSR

| Date | Zone Content Delivery |
| --- | --- |
| Feb 18, 2015 | content in base HTML divided into zones |
| Apr 24, 2015 | content other than Hero zone requested via CSR in `zone-manager.html` files |
| Sep 17, 2015 | Hero zone (homepage1-zone-1) *sometimes* requested via CSR |
| Oct 18, 2016 | zone-manager format changed to `.izl.json` |
| Nov 1, 2016 | all zones requested via CSR |
| Jan 31, 2017 | zone-manager extension changed to `.izl` |

# Starting in 2017, good coverage of mementos of the top 3 CNN zones



Internet Archive

# Contributions of homepage1-zone-1 zone-manager files from browser-based AIT collections, perma.cc, and Save Page Now



CNN.com - Archive-It collection 7678, Jul 2016 - Sep 2019

perma.cc

Live Web Proxy Crawls, Save Page Now

https://web.archive.org/web/collections/20170901000000*/https://www.cnn.com/data/ocs/section/index.html:**homepage1-zone-1**/views/zones/common/zone-manager.izl

# How many days would have missing content if IAWM limited the allowed distance between replayed HTML and JSON?



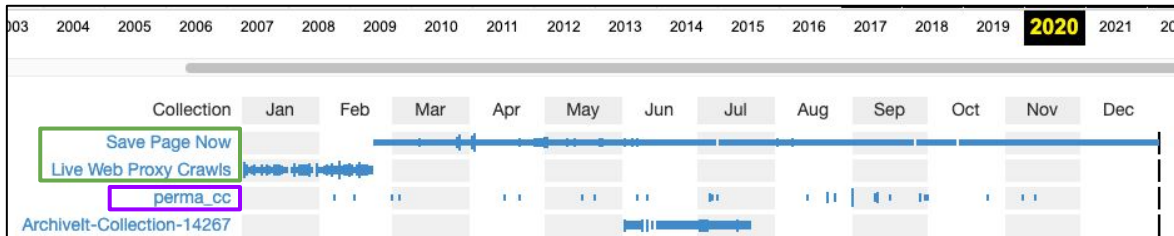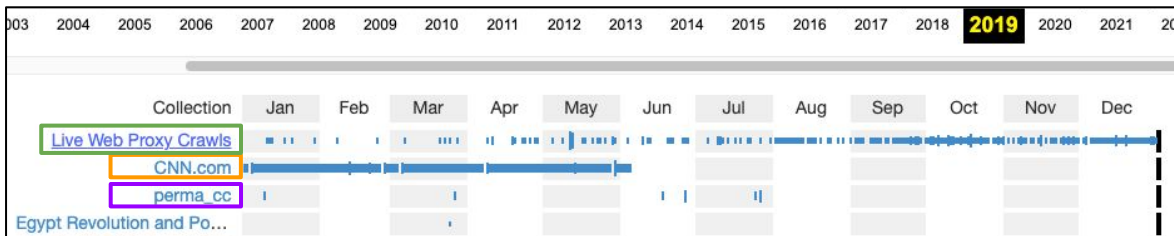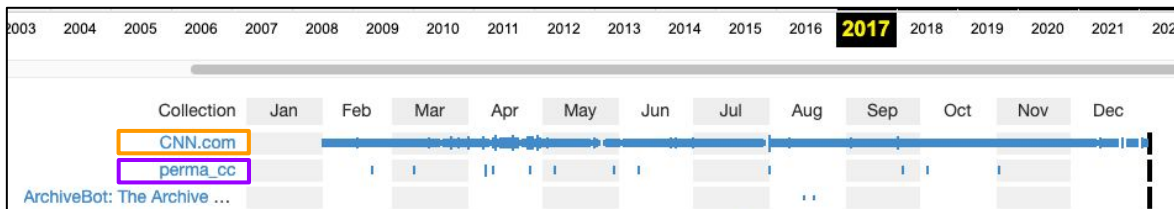| Date Range | Total | | Threshold | Affected | |
|---|---|---|---|---|---|
| | Mementos | Days | | Mementos | Days |
| Apr 24, 2015– | 788 | 29 | 1 hr | 786 | 29 |
| May 23, 2015 | | | 2 hrs | 785 | 29 |
| | | | 6 hrs | 780 | 28 |
| (none) | | | 24 hrs | 754 | 27 |
| | | | 48 hrs | 724 | 26 |
| May 23, 2015– | 17,017 | 425 | 1 hr | 16,863 | 425 |
| Jul 21, 2016 | | | 2 hrs | 16,738 | 425 |
| | | | 6 hrs | 16,244 | 410 |
| (sparse) | | | 24 hrs | 15,142 | 370 |
| | | | 48 hrs | 14,178 | 342 |
| Jul 21, 2016– | 154,917 | 2256 | 1 hr | 95,009 | 1826 |
| Sep 23, 2022 | | | 2 hrs | 67,992 | 1425 |
| | | | 6 hrs | 13,893 | 190 |
| (dense) | | | 24 hrs | 1812 | 37 |
| | | | 48 hrs | 526 | 8 |

VS.

# 48 hrs: Only 3/29 days would have 2nd level and lower content



| Date Range | Total Mementos | Days | Threshold | Affected Mementos | Days |
|---|---|---|---|---|---|
| Apr 24, 2015– | 788 | 29 | 1 hr | 786 | 29 |
| May 23, 2015 | | | 2 hrs | 785 | 29 |
| | | | 6 hrs | 780 | 28 |
| (none) | | | 24 hrs | 754 | 27 |
| | | | 48 hrs | 724 | 26 |
| May 23, 2015– | 17,017 | 425 | 1 hr | 16,863 | 425 |
| Jul 21, 2016 | | | 2 hrs | 16,738 | 425 |
| | | | 6 hrs | 16,244 | 410 |
| (sparse) | | | 24 hrs | 15,142 | 370 |
| | | | 48 hrs | 14,178 | 342 |
| Jul 21, 2016– | 154,917 | 2256 | 1 hr | 95,009 | 1826 |
| Sep 23, 2022 | | | 2 hrs | 67,992 | 1425 |
| | | | 6 hrs | 13,893 | 190 |
| (dense) | | | 24 hrs | 1812 | 37 |
| | | | 48 hrs | 526 | 8 |



VS.

# 48 hrs: Only 83/425 days would have 2nd level and lower content



| Date Range | Total | | | Affected | |
|---|---|---|---|---|---|
| | Mementos | Days | Threshold | Mementos | Days |
| Apr 24, 2015– | 788 | 29 | 1 hr | 786 | 29 |
| May 23, 2015 | | | 2 hrs | 785 | 29 |
| | | | 6 hrs | 780 | 28 |
| (none) | | | 24 hrs | 754 | 27 |
| | | | 48 hrs | 724 | 26 |
| May 23, 2015– | 17,017 | 425 | 1 hr | 16,863 | 425 |
| Jul 21, 2016 | | | 2 hrs | 16,738 | 425 |
| | | | 6 hrs | 16,244 | 410 |
| (sparse) | | | 24 hrs | 15,142 | 370 |
| | | | 48 hrs | 14,178 | 342 |
| Jul 21, 2016– | 154,917 | 2256 | 1 hr | 95,009 | 1826 |
| Sep 23, 2022 | | | 2 hrs | 67,992 | 1425 |
| | | | 6 hrs | 13,893 | 190 |
| (dense) | | | 24 hrs | 1812 | 37 |
| | | | 48 hrs | 526 | 8 |

vs.

# 48 hrs: Only 8 days would be *missing* 2nd level and lower content



| Date Range | Total | | Threshold | Affected | |
| --- | Mementos | Days | | Mementos | Days |
| --- | --- | --- | --- | --- | --- |
| Apr 24, 2015– | 788 | 29 | 1 hr | 786 | 29 |
| May 23, 2015 | | | 2 hrs | 785 | 29 |
| | | | 6 hrs | 780 | 28 |
| (none) | | | 24 hrs | 754 | 27 |
| | | | 48 hrs | 724 | 26 |
| May 23, 2015– | 17,017 | 425 | 1 hr | 16,863 | 425 |
| Jul 21, 2016 | | | 2 hrs | 16,738 | 425 |
| | | | 6 hrs | 16,244 | 410 |
| (sparse) | | | 24 hrs | 15,142 | 370 |
| | | | 48 hrs | 14,178 | 342 |
| Jul 21, 2016– | 154,917 | 2256 | 1 hr | 95,009 | 1826 |
| Sep 23, 2022 | | | 2 hrs | 67,992 | 1425 |
| | | | 6 hrs | 13,893 | 190 |
| (dense) | | | 24 hrs | 1812 | 37 |
| | | | 48 hrs | 526 | 8 |

VS.

# Archive-It collections can be sparse, could have continuing issues

Mar 27, 2020

Archive-It collection 4887
National Library of Medicine's Global Health Events
https://archive-it.org/collections/4887

Internet Archive Wayback Machine (IAWM)



*Hero story content from Jun 18, 2020*

https://wayback.archive-it.org/4887/**20200327123616**/https://www.cnn.com/

https://web.archive.org/web/**20200327123616**/https://www.cnn.com/
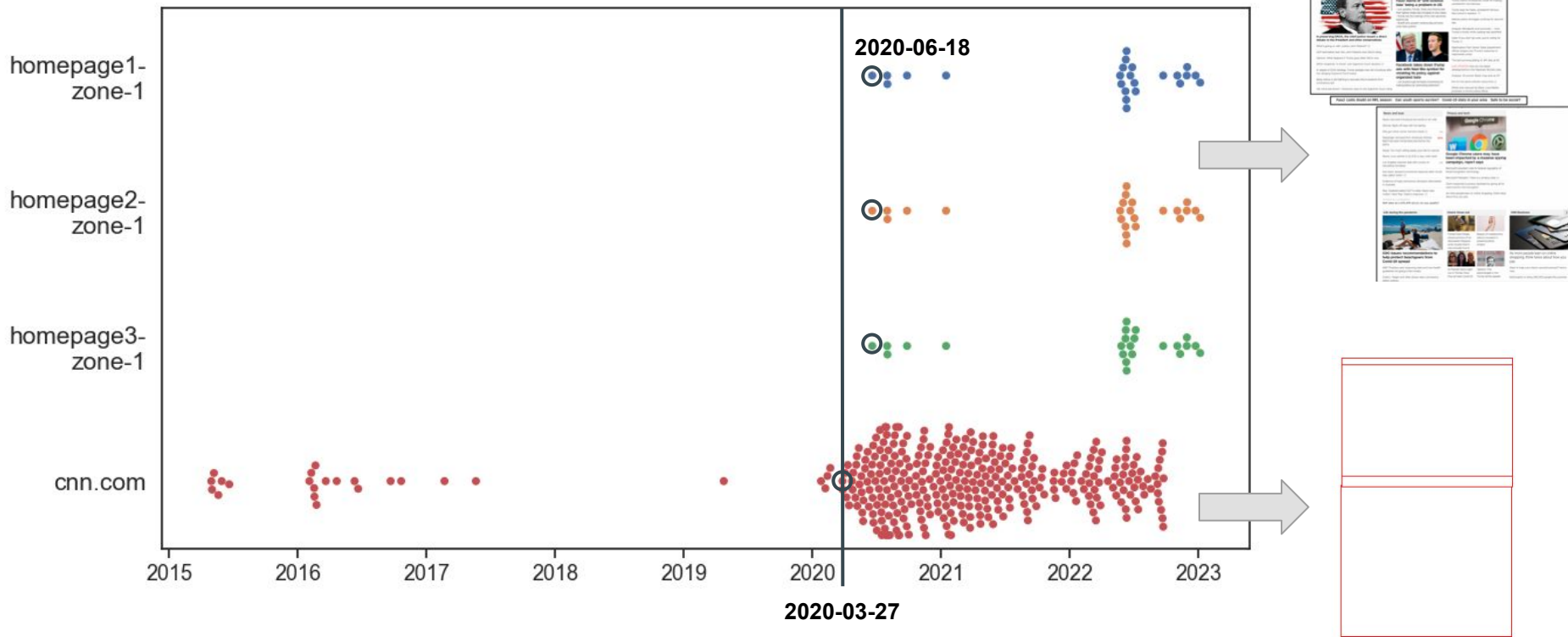
# March 27 replay contains placeholder from March 27 and content from June 18



Archive-It collection 4887

2020-06-18

2020-03-27

# Recommendations

- Use browser-based crawling for websites with heavy CSR, like CNN.com
- Develop crawling methods that can detect when CSR being used, nominate sites for browser-based crawls
- Be aware of implications of mixing mementos from conventional and browser-based crawls, audit replay of previous captures
- Develop replay methods that can highlight temporal violations on the page

- Potential mitigations
  - redact mementos between Sep 2015-Nov 2016 with temporal violations in the Hero story
  - increase the MIME types/file extensions included in "About this Capture"
  - include alerts for (or don't load) JSON files with a temporal gap greater than some threshold
  - allow Archive-It collections to borrow certain resources (like JSON) from other collections

# Summary

- Since 2015, CNN.com has used client-side rendering (CSR) to build its front page

- CNN.com's CSR model conflicts with standard web archiving models (i.e., load content from closest memento)

- For webpages built with CSR (not only CNN.com), the mix of captures from conventional (Heritrix) and browser-based crawlers can result significant temporal violations upon replay

- We found over 15,000 mementos of CNN.com in the Internet Archive's Wayback Machine with temporal violations greater than 2 days

- We have documented the issue, but more work should be done in detecting these types of problems in real-time