# Archiving Digital Marketing

## Examining Preservation of Dynamic Content on the Web Through the Lens of Online Advertisements

Christopher Rauch, Alex H. Poole, Travis Reid, Michele C. Weigle, Michael L. Nelson, Faryaneh Poursardar, and **Mat Kelly**

mkelly@drexel.edu • @machawk1 • @machawk1@digipres.club

September 19, 2024

A COLLABORATION OF

DREXEL UNIVERSITY
College of
Computing &
Informatics

**&** OLD DOMINION UNIVERSITY®

SUPPORTED BY

INSTITUTE of
Museum and Library
SERVICES

Slides: https://bit.ly/ipres2024

# Why Preserve Online Advertisements

Advertisements are not just promotional tools but reflect and shape cultural and societal norms over time.



Source: Wikimedia Commons

# Why Preserve Online Advertisements

ONLINE ADVERTISEMENTS

**ARE DOCUMENTARY EVIDENCE**

for understanding mid-1990s history onward
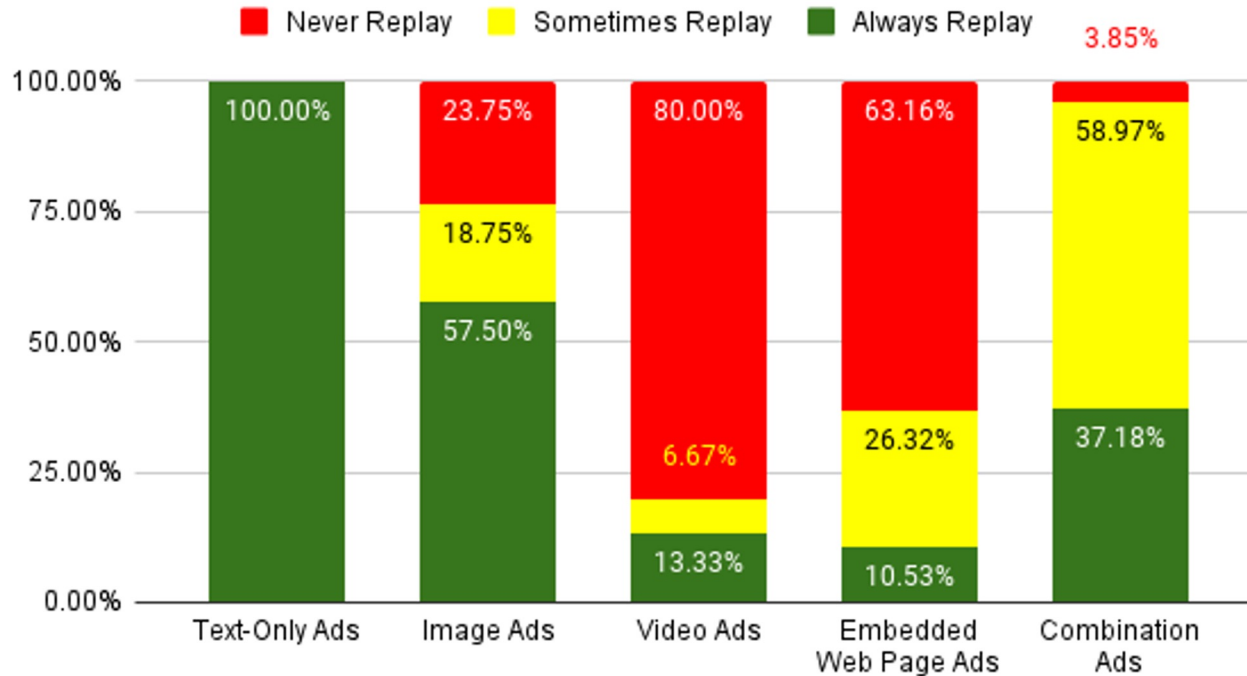
# Web Ads are a Rich Source for Historians

Archived ads enable scholars to:

- **EXPLORE** new questions
- **REVISIT** existing findings
- **ANALYZE** legal and social advertising practices.

# Challenges: Ephemerality



Never Replay  Sometimes Replay  Always Replay

| | Never Replay | Sometimes Replay | Always Replay |
|---|---|---|---|
| Text-Only Ads | | | 100.00% |
| Image Ads | 23.75% | 18.75% | 57.50% |
| Video Ads | 80.00% | 6.67% | 13.33% |
| Embedded Web Page Ads | 63.16% | 26.32% | 10.53% |
| Combination Ads | 3.85% | 58.97% | 37.18% |

Replaying Ads in Containing Web Page

# Challenges: Institutional Resistance



Scheme    Random Value    Subdomain    Domain

https://e76308bcf1c30aa4c853507f4b382285.safeframe.googlesyndication.com/safeframe/1-0-40/html/container.html

Path

# Enhanced User Interface Tools



https://github.com/internetarchive/brozzler

# Increased Support for Dynamic Content



Starfield Wiki page (https://starfield.fandom.com/wiki/Starfield)

# Toolbox Approach

```
[mrk335@CCI-63JQLNF9 Squidwarc % ./run-crawler.sh -c conf.json
Running Crawl From Config File conf.json
With great power comes great responsibility!
Squidwarc is not responsible for ill behaved user supplied scripts!

Crawler Operating In page-only mode
Crawler Will Be Preserving 1 Seeds
Crawler Generated WARCs Will Be Placed At /private/tmp/Squidwarc
Crawler Will Be Generating WARC Files Using the filenamified url
Crawler Navigating To https://ipres2024.pubpub.org/
Crawler Navigated To https://ipres2024.pubpub.org/
Running user script
Crawler Generating WARC
Crawler Has 0 Seeds Left To Crawl
Crawler shutting down. Have nice day :)
mrk335@CCI-63JQLNF9 Squidwarc % 
```

Squidwarc

https://github.com/N0taN3rd/Squidwarc

# Tools Comparison


REPLAY WEBPAGE

ReplayWeb.page's URL Search

Our Tool For Displaying Ads From A WARC file

Wayback Machine's URL Search

https://github.com/webrecorder/replayweb.page

# Study Results


Number of Ads Per Category

# Study Results



Chart legend: Never Replay (red), Sometimes Replay (yellow), Always Replay (green)

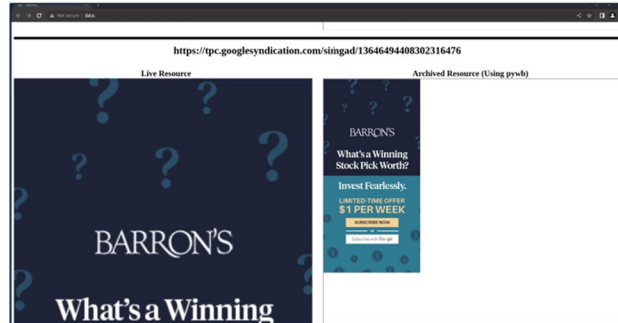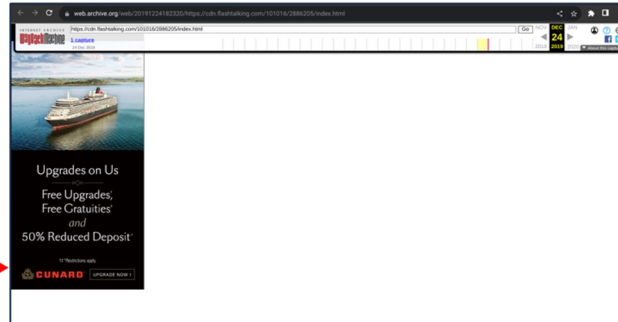| Category | Always Replay | Sometimes Replay | Never Replay |
|---|---|---|---|
| Text-Only Ads | 100.00% | | |
| Image Ads | 57.50% | 18.75% | 23.75% |
| Video Ads | 13.33% | 6.67% | 80.00% |
| Embedded Web Page Ads | 10.53% | 26.32% | 63.16% |
| Combination Ads | 37.18% | 58.97% | 3.85% |

Replaying Ads in Containing Web Page

# Pursue Partnerships and Collaboration

INTERNATIONAL CONFERENCE
ON DIGITAL PRESERVATION
16 - 20 SEPTEMBER 2024
Ghent & Flanders, Belgium

iPRES 2024

Archiving Digital Marketing • @machawk1    https://bit.ly/ipres2024    13

# Future Work





10 Steps to Personas

World Knowledge — Public Web Text (~$10^{14}$ tokens) — Compress — Represented by distributed carriers — Persona Hub (1 billion personas) ~$10^{10}$ tokens — Decompress — Generate texts with their knowledge — World Knowledge — Public Web Text (~$10^{14}$ tokens)
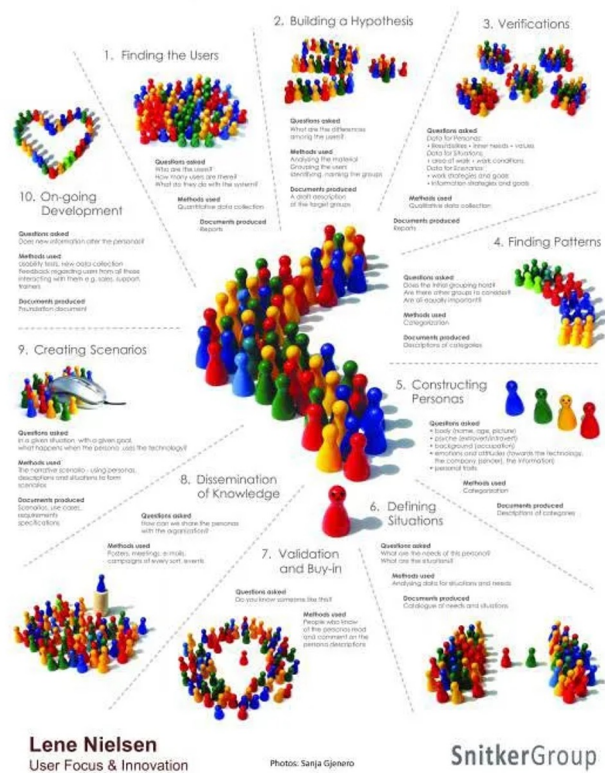
(Delétang et al., 2023; Ge et al., 2024), https://github.com/tencent-ailab/persona-hub

Lene Nielsen.

# In Summary

- Online ads are cultural artifacts, preservation is essentials
- Modern ads are highly personalized and dynamically loaded
- Web Ads tech is complex and the ads themselves are ephemeral
- Text/img ads are easy to capture and reply but majority use more complex techs
- Encoding personas into crawler can help to surface and preserve resources crawlers will never encounter

iPRES 2024
INTERNATIONAL CONFERENCE ON DIGITAL PRESERVATION
16 - 20 SEPTEMBER 2024
Ghent & Flanders, Belgium

Archiving Digital Marketing
Mat Kelly • @machawk1@digipres.club • mkelly@drexel.edu
https://bit.ly/ipres2024

INSTITUTE of Museum and Library SERVICES