# "What You See No One Saw"

**Mat Kelly**, Alexander H. Poole, Michele C. Weigle, Michael L. Nelson, Travis Reid, Christopher Rauch, and Hyung Wook Choi

**Drexel University CCI** & **Old Dominion University WS-DL**
mkelly@drexel.edu – @machawk1

**IIPC Web Archiving Conference (WAC)**
Oslo, Norway
April 10, 2025

slides:
bit.ly/iipcwac2025

# The Past Web Saved Is Not The Web That Was



ARCHIVAL WEB CRAWLER

WEB USER

WEB USER

WEB USER

WEB USER

WEB USER

WEB ARCHIVE USER

"What You See No One Saw"
▸ Slides: bit.ly/iipcwac2025

**Mat Kelly**
@machawk1

IIPC Web Archiving Conference (WAC) 2025
April 10, 2025
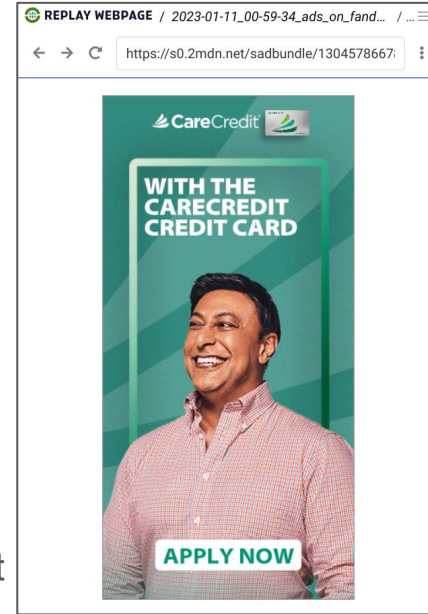
2

# Crawler Preserve A False Web (Premise)

- Archival crawlers preserve a version of the web inconsistent with web users' experience, a web that actually wasn't
- Customization, personalization based on user history is not *canonical*
- Crawlers (rightfully) see a clean/agnostic version of web sites, devoid of any individuals' experience, PII
- Ergo, what crawlers preserve is a version of the web inconsistent with what a user would have seen at that time
- False history? Nature of experience

ARCHIVAL
WEB CRAWLER

≠

WEB USER

"What You See No One Saw"
‣ Slides: bit.ly/iipcwac2025

**Mat Kelly**
@machawk1

IIPC Web Archiving Conference (WAC) 2025
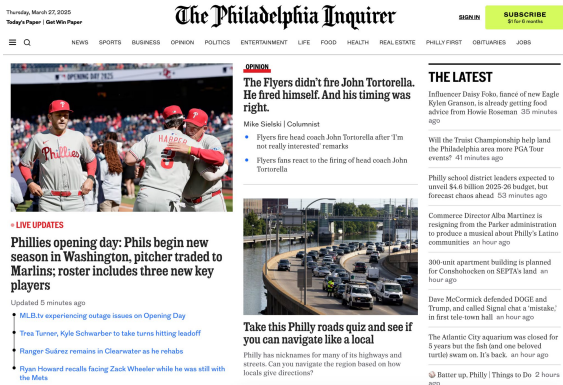April 10, 2025

3

# A Valid Perspective, Just Not A Web Users'



- Crawler's perspective *is* a valid representation
- There is no *one true* representation of a personalized web
- Web Ads?
  - Annoying but useful for study
  - Hyper-personalized, distinguishing factor from generic representation
  - Google's Manifest V3 / Ad blocker drama* means Chrome users are returning to an ad-ridden web
  - As with analog advertisements, web ads represent zeitgeist in retrospect



\* {
https://www.eff.org/deeplinks/2021/12/chrome-users-beware-manifest-v3-deceitful-and-threatening
https://developer.chrome.com/docs/extensions/develop/migrate/mv2-deprecation-timeline
https://blog.mozilla.org/addons/2024/03/13/manifest-v3-manifest-v2-march-2024-update/
https://developer.chrome.com/blog/resuming-the-transition-to-mv3/

"What You See No One Saw"
▸ Slides: bit.ly/iipcwac2025

**Mat Kelly**
@machawk1

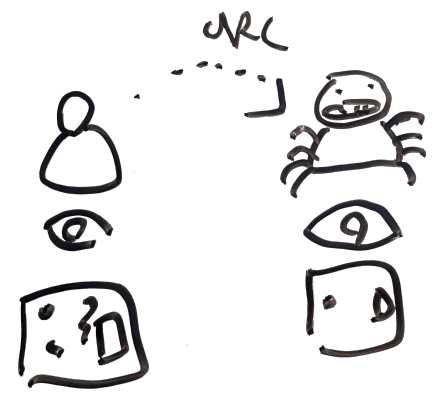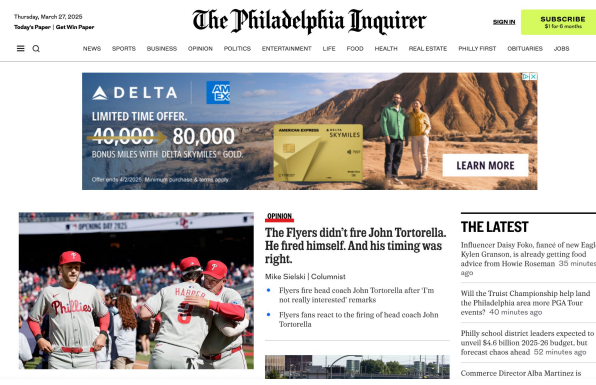IIPC Web Archiving Conference (WAC) 2025
April 10, 2025

4

# Can We Save The Web We See From Our Perspective?

- Repurpose user's daily driver profile as crawler basis
- Permutate attributes of a user to represent a "persona", producing a web experience closer to that of an actual web user cf. crawler
- Avoid clean slate crawling and delegation to a user-agnostic crawler
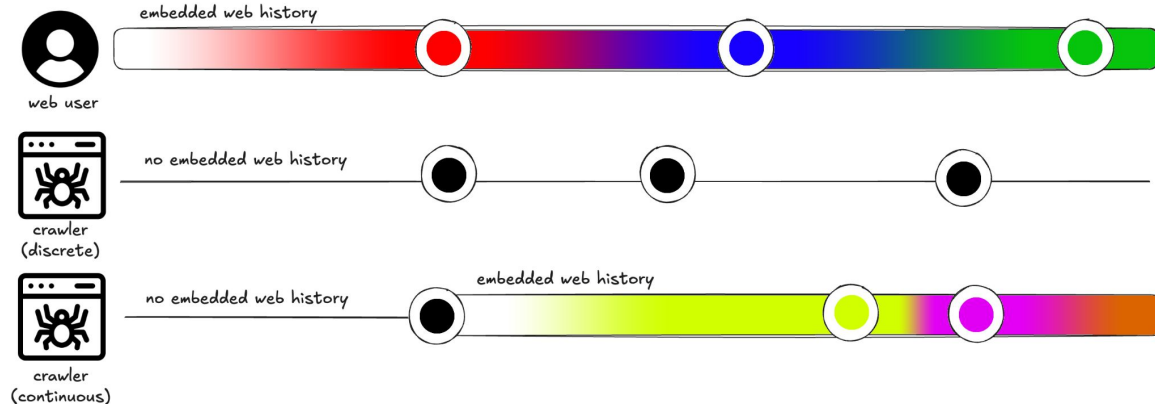- Scale?



vs

"What You See No One Saw"
▸ Slides: bit.ly/iipcwac2025

**Mat Kelly**
@machawk1

IIPC Web Archiving Conference (WAC) 2025
April 10, 2025

5

# Crawling w/ a Web History + Discrete vs. Continuous

- Want to either reuse browser user profile or extract feature (e.g., cookies) to be used as the basis for what is served at archive time
  - What else is contained in this profile?
  - Is reuse possible/feasible for web archiving? What are Selenium's capabilities? Other headless crawlers

"What You See No One Saw"
▸ Slides: bit.ly/iipcwac2025

**Mat Kelly**
@machawk1

IIPC Web Archiving Conference (WAC) 2025
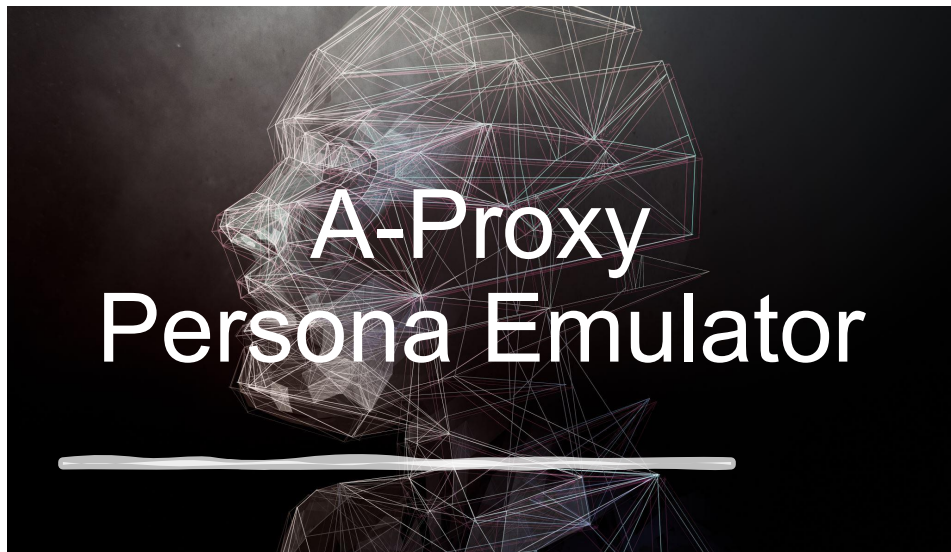April 10, 2025

6

Prior Technical Work
# Leveraging Perspective-Based Crawling

- WARCreate - browser extension that archives by-value (cf. URI as basis)
  - Manifest V3 caveat (webRequest)
- Warcprox - save representations as they come over the wire
- Ad Blockers
  - Are users seeing the ad-ridden, true representation of the web?

"What You See No One Saw"
▸ Slides: bit.ly/iipcwac2025

**Mat Kelly**
@machawk1

IIPC Web Archiving Conference (WAC) 2025
April 10, 2025

7

# WIP: Persona-based Web Archiving

- Rapidly prototyped crawl director
- Side-load Chrome profile with Selenium WebDriver
- UI for user to specify crawl profile attributes
- Based on Andy Jackson's Sliver
  - https://github.com/anjackson/sliver



A-Proxy
Persona Emulator

https://github.com/savingads/a-proxy

"What You See No One Saw"
▸ Slides: bit.ly/iipcwac2025

**Mat Kelly**
@machawk1

IIPC Web Archiving Conference (WAC) 2025
April 10, 2025

8

# A-Proxy

**Dashboard**
**Personas**

WEB BROWSING
**Browse As**
**Archived Pages**

**Settings**

Enter URL | Preview | Archive

## Welcome to A-Proxy
Web Testing with Geolocation and Language Simulation

### About A-Proxy

A-Proxy is a tool that allows you to test websites with different geolocation and language settings by using VPN connections and a customized browser setup. This is useful for testing how websites behave for users in different countries and with different language preferences.

With A-Proxy, you can:

- Connect to VPN servers in different countries
- Simulate different browser language settings
- Override geolocation data in the browser
- Create and manage user personas with demographic, psychographic, behavioral, and contextual data
- Take screenshots of websites with the simulated settings

Go to Dashboard | Manage Personas

### VPN Status

**VPN Running:** False

VPN is not currently running. Start it from the dashboard to simulate different locations.

Start VPN

## Your Browser Information

### 🌐 Location & Language

| | |
|---|---|
| Geolocation | 39.6815, -74.2389 |
| Language | en-US |
| Time Zone | America/New_York |

### 🖥 Device Information

| | |
|---|---|
| Platform | Win32 |
| Screen Resolution | 1920 x 1080 |
| Device Pixel Ratio | 1 |
| Touch Support | No |

### ⚙ Browser Details

| | |
|---|---|
| User Agent | Mozilla/5.0 (Windows NT 10.0; Win64; ... |
| Browser Name | Firefox |
| Cookies Enabled | Yes |
| Do Not Track | No |

### 📶 Connection Information

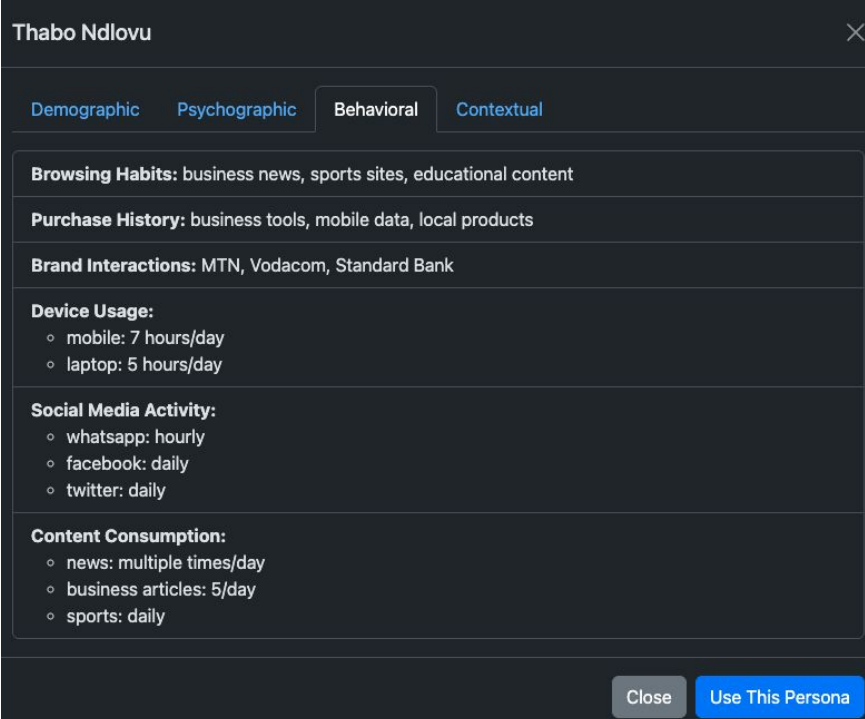| | |
|---|---|
| Online Status | Online |
| Connection Type | Not available |
| Effective Connection Type | Not available |
| Downlink Speed | Not available |

### 📍 Your Location


Your browser location

# Parameters of Perspective / Personalization

- Demographic
  - ex: Location, Language
- Psychographic
  - ex: interests, values
- Behavioral
  - ex: browsing habits, social media activity
- Contextual
  - ex: time of day, weather, browser



Thabo Ndlovu                                                    ✕

Demographic    Psychographic    **Behavioral**    Contextual

**Browsing Habits:** business news, sports sites, educational content

**Purchase History:** business tools, mobile data, local products

**Brand Interactions:** MTN, Vodacom, Standard Bank

**Device Usage:**
  ○ mobile: 7 hours/day
  ○ laptop: 5 hours/day

**Social Media Activity:**
  ○ whatsapp: hourly
  ○ facebook: daily
  ○ twitter: daily

**Content Consumption:**
  ○ news: multiple times/day
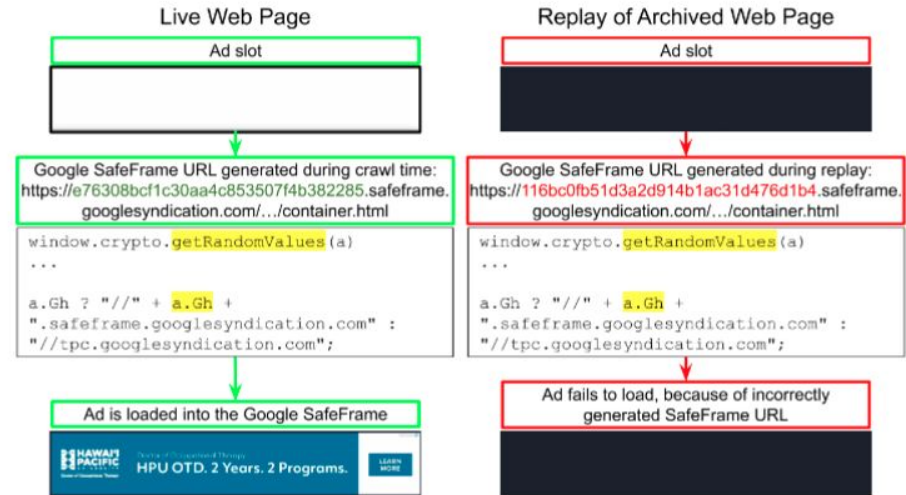  ○ business articles: 5/day
  ○ sports: daily

Close    Use This Persona

"What You See No One Saw"
▸ Slides: bit.ly/iipcwac2025

**Mat Kelly**
@machawk1

IIPC Web Archiving Conference (WAC) 2025
April 10, 2025

10

# Initial Observations and Experiments

- Examine existing crawler-derived WARCs, observe personalization traits
- Enumerate attributes like location, user-agent, language for rudimentary personalization
- Use case: capture delta of web ads
  - Difficult due to randomization, dynamic
  - Requires replay amendment
  - Web ads are more than just images
    - e.g., video, combo, interactive



See arXiv:2502.01525, 2025

"What You See No One Saw"
▸ Slides: bit.ly/iipcwac2025

**Mat Kelly**
@machawk1

IIPC Web Archiving Conference (WAC) 2025
April 10, 2025

11

# "What You See Now One Saw"



**WEB USER**

**WEB ARCHIVE USER**

Project supported by
**INSTITUTE of Museum and Library SERVICES**

- Sufficient personalized context is lost when delegating to a crawl by URI
- That which we consider the historical web was captured through the lens of a perspective-agnostic crawler
- Project underway, interpolating personas, gathering data, building A-Proxy
- Ongoing dev work, data at
  <u>https://github.com/savingads</u>

<span style="color:red">**See our recent tech report on archiving web ads!  arXiv:2502.01525, 2025. →**</span>

**Slides: bit.ly/iipcwac2025**



### Archiving and Replaying Current Web Advertisements: Challenges and Opportunities

TRAVIS REID, Old Dominion University, USA
ALEX H. POOLE, Drexel University, USA
HYUNG WOOK CHOI, Drexel University, USA
CHRISTOPHER RAUCH, Drexel University, USA
MAT KELLY, Drexel University, USA
MICHAEL L. NELSON, Old Dominion University, USA
MICHELE C. WEIGLE, Old Dominion University, USA

Although web advertisements represent an inimitable part of digital cultural heritage, serious archiving and replay challenges persist. To explore these challenges, we created a dataset of 279 archived ads. We encountered five problems in archiving and replaying them. For one, prior to August 2023, Internet Archive's Save Page Now service excluded not only well-known ad services' ads, but also URLs with ad related file and directory names. Although after August 2023, Save Page Now still blocked the archiving of ads loaded on a web page, it permitted the archiving of an ad's resources if the user directly archived the URL(s) associated with the ad. Second, Brozzler's incompatibility with Chrome prevented ads from being archived. Third, during crawling and replay sessions, Google's and Amazon's ad scripts generated URLs with different random values. This precluded archived ads' replay. Updating replay systems' fuzzy matching approach should enable the replay of these ads. Fourth, when loading Flashtalking web page ads outside of ad iframes, the ad script requested a non-existent URL. This prevented the replay of ad resources. But as was the case with Google and Amazon ads, updating replay systems' fuzzy matching approach should enable Flashtalking ads' replay. Finally, successful replay of ads loaded in iframes with the src attribute of "about:blank" depended upon a given browser's service worker implementation. A Chromium bug stopped service workers from accessing resources inside of this type of iframe, which in turn prevented replay. Replacing the "about:blank" value for the iframe's src attribute with a blob URL before an ad was loaded solved this problem. Resolving these replay problems will improve the replay of ads and other dynamically loaded embedded

#### 1    INTRODUCTION

Brewster Kahle, founder of the Internet Archive, [...]
of valuable scientific, cultural and historical inf[...]
characterized the web in similar terms, but also s[...]
for the study of almost every possible aspect of th[...]
scholars, however, web content has been hemor[...]

Whether impelled by legal obligation, business[...]
and/or historical research, web archiving involve[...]
to content [6, 14, 51, 53]. Web archives may be u[...]
about the period in which the archived content [...]

Because the web depends upon advertising re[...]
dynamic content. Just as physical ephemera in li[...]

Authors' addresses: Travis Reid, Department of Computer S[...]
Alex H. Poole, Department of Information Science, Drexel U[...]
Department of Information Science, Drexel University, Phila[...]
Information Science, Drexel University, Philadelphia, PA, 1[...]
Drexel University, Philadelphia, PA, 19104, USA, mkelly@d[...]
University, Norfolk, VA, 23529, USA, mln@cs.odu.edu; Mic[...]
Norfolk, VA, 23529, USA, mweigle@cs.odu.edu.