# Saving the Web

**Mat Kelly**
mkelly@drexel.edu
Assistant Professor, Information Science
Drexel University College of Computing and Informatics

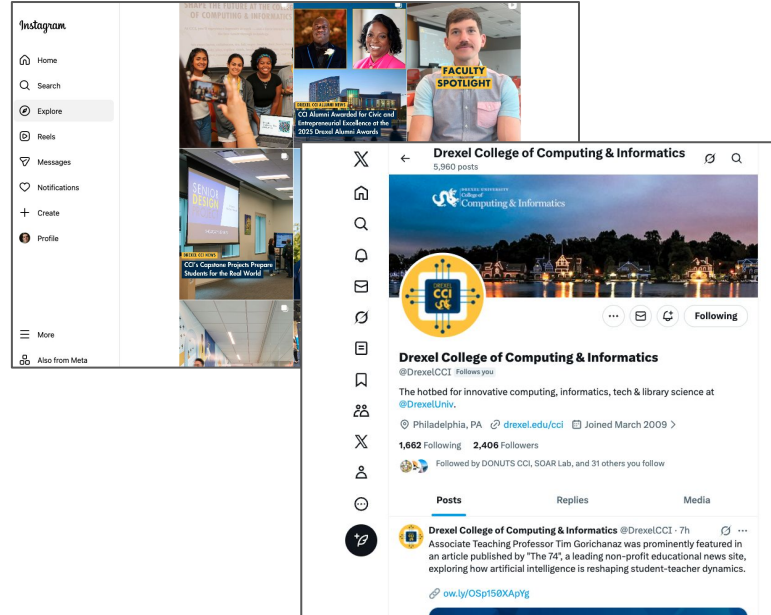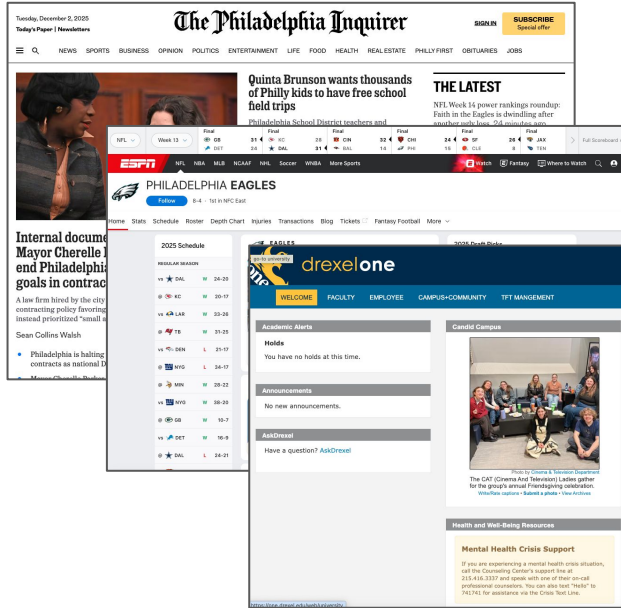**INFO591 Week 10**
December 4, 2025

Slides: https://bit.ly/info591

**https://matkelly.com**

# The Web

"Saving the Web"
▸ Slides: bit.ly/info591
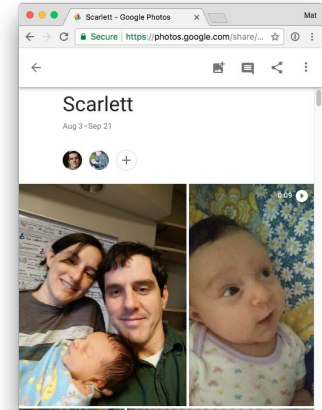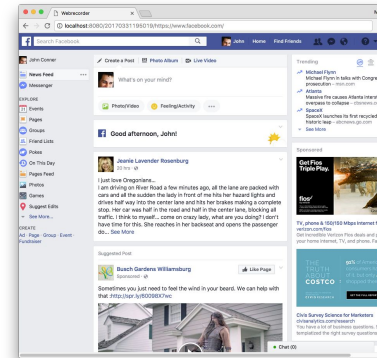
**Mat Kelly, Asst. Prof., Drexel IS**
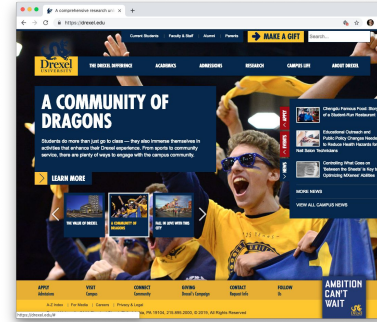@machawk1 • https://matkelly.com
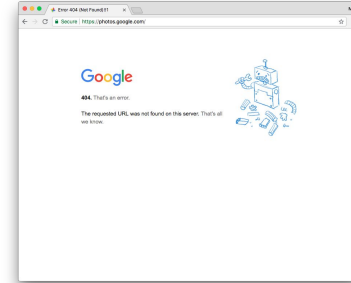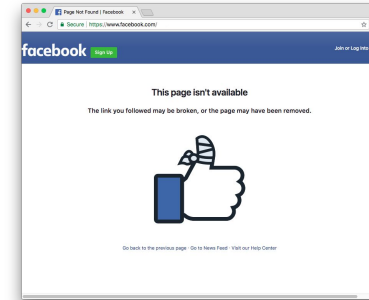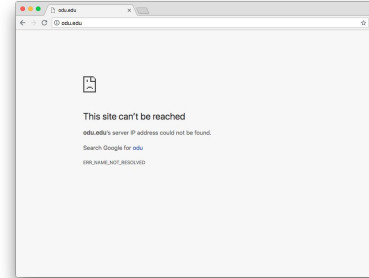
INFO591 Week 10
December 4, 2025

2

# The Web's Importance

- Information access
- Personal interfacing (e.g., social media)
- Cultural artifact
- Personal knowledge reference

"Saving the Web"
▸ Slides: bit.ly/info591

**Mat Kelly, Asst. Prof., Drexel IS**
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025
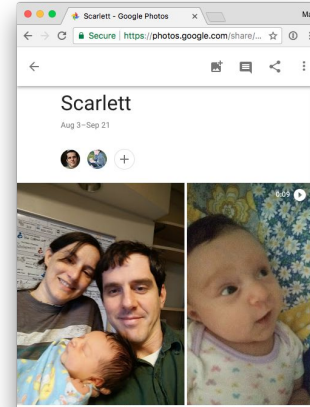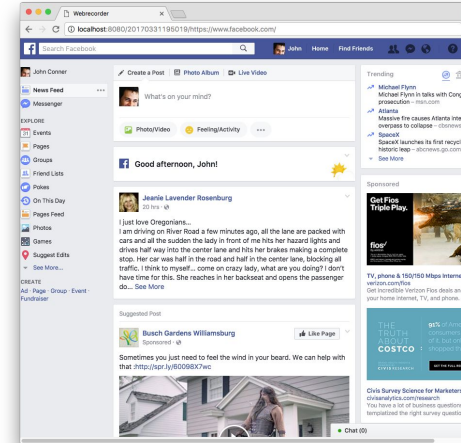
3

# The Web is Ephemeral

- No guarantee that information or services available today will be available tomorrow
- We are reliant on this service
  - Inclusive of mobile apps
- Links rot quickly
- User-generated content is especially fragile
  - Apps, social media posts, and comments vanish as services evolve

"Saving the Web"
▸ Slides: bit.ly/info591

**Mat Kelly, Asst. Prof., Drexel IS**
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025

4

# *Our* Web

- Born digital data
- Personalized Content
- Private Content
- Unique content
- Typically not preserved
  - Should it be?

"Saving the Web"
▸ Slides: bit.ly/info591

**Mat Kelly, Asst. Prof., Drexel IS**
@machawk1 • https://matkelly.com

INFO591 Week 10
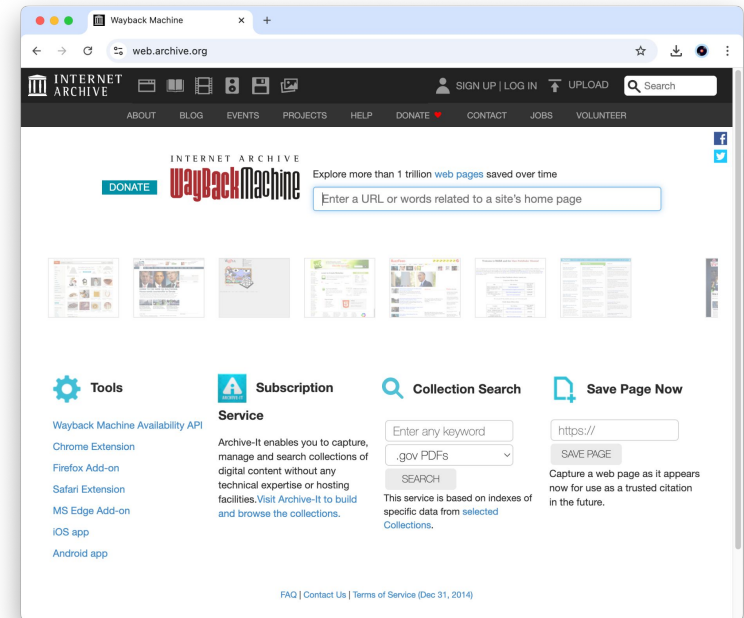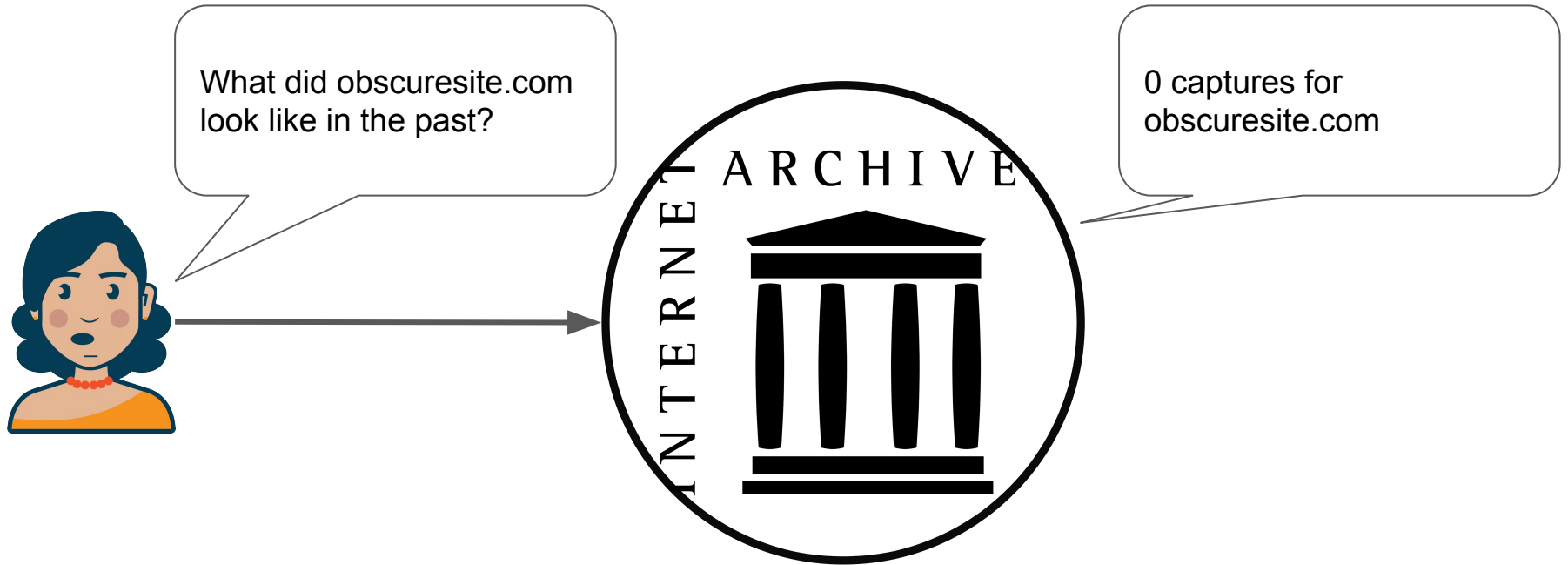December 4, 2025

5

# But Can We Save It?

# The Internet Archive



- Captures the web at 150 TB / day; currently 175 PB of web history (Dec 2025)
- Also: music recordings, radio broadcasts, podcasts, software, images, etc.

"Saving the Web"
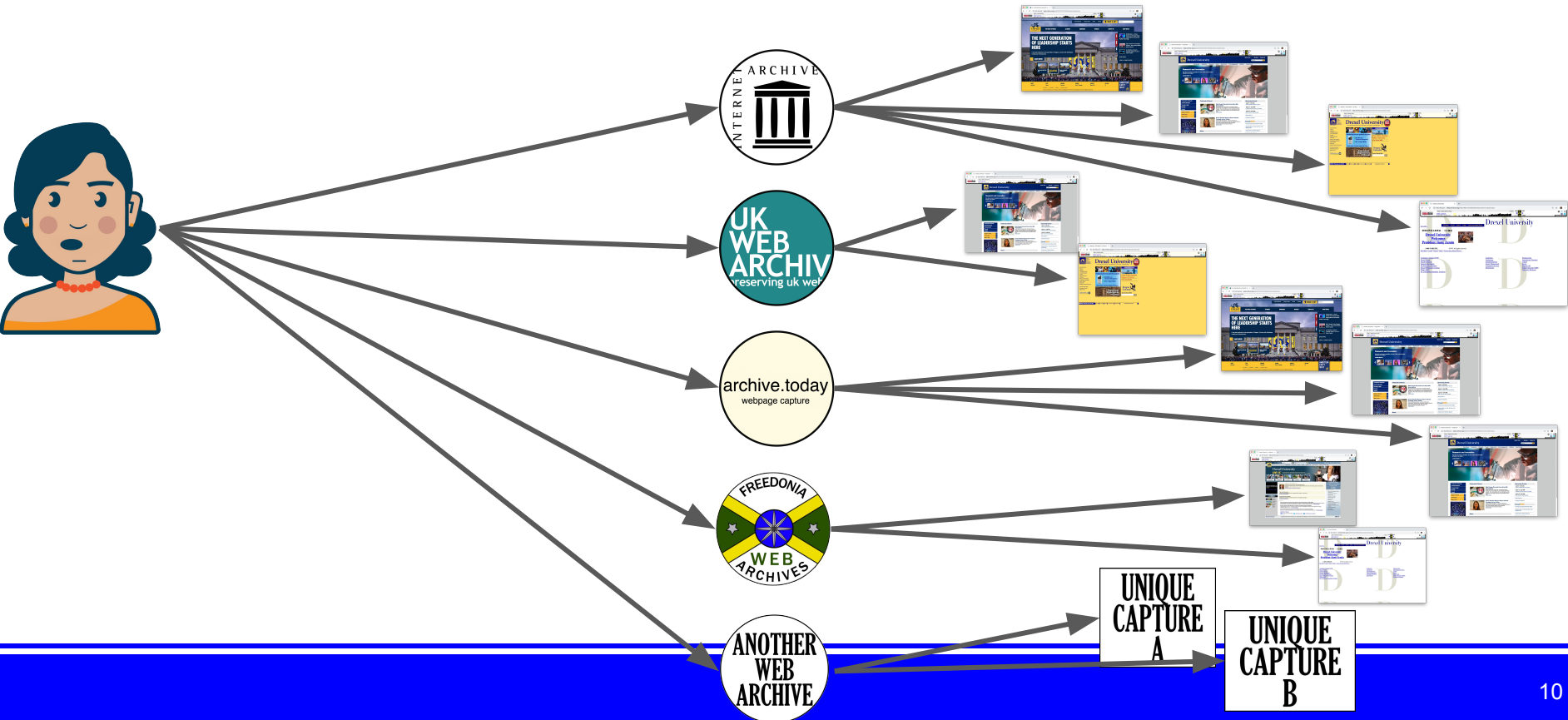▸ Slides: bit.ly/info591

Mat Kelly, Asst. Prof., Drexel IS
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025

# Not everything is preserved



What did obscuresite.com look like in the past?

0 captures for obscuresite.com

"Saving the Web"
▸ Slides: bit.ly/info591

Mat Kelly, Asst. Prof., Drexel IS
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025

8

# Multiple archival efforts (3 of many)

"Saving the Web"
▸ Slides: bit.ly/info591

Mat Kelly, Asst. Prof., Drexel IS
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025

9

# More archives produces a more comprehensive picture

# Web Archiving



Associate live Web URIs

With their archived representations

# Web Archiving: not just "Save As…"

- Metadata
- Provenance
- Standards
- Big Data
- Composite Resources (cf. analog/conventional archiving)
- Technical Nuances
  - Indexing
  - Deduplication
  - Resilience
  - Handling Dynamicism
- Preservation is a race against time — the moment something appears online, it can begin disappearing

"Saving the Web"
▸ Slides: bit.ly/info591

Mat Kelly, Asst. Prof., Drexel IS
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025      12

# Web Archiving Metadata

- The how, when and why web content was collected, preserved, and made accessible
- High-level descriptors
- Fields in crawls
- Derivative metadata formats
- Types:
    - Technical (format, checksums)
    - Descriptive (URL, pub date, site title)
    - Preservation (relationship w/ other captures, date/time, crawler used)
    - Administrative (curatorial decision, access restrictions)

"Saving the Web"
‣ Slides: bit.ly/info591

**Mat Kelly, Asst. Prof., Drexel IS**
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025     13

# Web Archiving Provenance

- Provenance documents where archived content came from, how it was collected, and when.
- Establishes the authenticity and integrity of archived web resources
    - Was it really as depicted? Who did it? When?
- Often realized through metadata
- Essential for reproducibility
    - Others can understand or repeat the capture process
- Reveals biases or gaps
    - what was included, excluded, or technically unreachable
- Clarifies chain of custody when content is transformed, migrated, or redacted over time

"Saving the Web"
▸ Slides: bit.ly/info591

Mat Kelly, Asst. Prof., Drexel IS
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025

14

# Web Archiving Standards

- WARC: ISO 28500:2017
  - Storage format for web archives
  - Container that contains transactional information and metadata
- Memento: IETF RFC7089
  - Provides mechanisms for interacting with the past web

"Saving the Web"
‣ Slides: bit.ly/info591

**Mat Kelly, Asst. Prof., Drexel IS**
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025    15

# Web Archiving Standards: WARC

- ISO-standardized file format (ISO 28500) used for storing web crawls and related digital content
- Uses records (e.g., response, request, metadata, revisit, resource) to represent different capture types
- Preserves original HTTP request/response pairs, supporting authenticity and replay
- Designed for long-term digital preservation: stable, extensible, and format-agnostic

"Saving the Web"
▸ Slides: bit.ly/info591

**Mat Kelly, Asst. Prof., Drexel IS**
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025

16

# WARC Anatomy

WARC response record

Warc-response header

HTTP resp header

HTTP resp payload

WARCs also contain:
- HTTP requests
- warc-info
- warc-metadata records
- etc.

# Web Archiving Standards: Memento

- Standardizes semantics and syntax for time on the web
- Content negotiation in the temporal dimension
- IETF RFC7089 (standard)

**Original Resource**
Resource that exists or used to exist; we are interested in a past state of it

```
Link: <URI-R>;type="original"
        points to Original Resource
```

**Memento**
Resource that encapsulates a past state of the Original Resource

```
Link: <URI-M>;type="memento"
        points to Memento
```

**TimeGate**
Resource that "decides", based on a given datetime, which is the temporally best Memento for an Original Resource

```
Link: <URI-TG>;type="timegate"
        points to TimeGate
```

"Saving the Web"
▸ Slides: bit.ly/info591

**Mat Kelly, Asst. Prof., Drexel IS**
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025     18

# Background: Memento Request Example

**HTTP Request**
- Accept-Datetime: Wed, 02 Aug 2017 23:15:00 GMT
- GET: http://web.archive.org/web/http://www.cnn.com

URI-G
*G*

Request `cnn.com` at Sept 11, 2001 at 9am EST

"Saving the Web"
‣ Slides: bit.ly/info591

Mat Kelly, Asst. Prof., Drexel IS
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025

19

# Background: Memento Request Example

**HTTP Request**
- Accept-Datetime: Wed, 02 Aug 2017 23:15:00 GMT
- GET: http://web.archive.org/web/http://www.cnn.com

Request `cnn.com` at Sept 11, 2001 at 9am EST

URI-G
*G*

**HTTP Response (302)**
- Memento-Datetime: Wed, 02 Aug 2017 23:18:04 GMT
- Location: http://web.archive.org/web/20170802231804/http://www.cnn.com/
- Link:

URI-T
*T*

URI-R
*R*

URI-G
*G*

URI-M
*M*

timemap

original

timegate

memento

"Saving the Web"
▸ Slides: bit.ly/info591

Mat Kelly, Asst. Prof., Drexel IS
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025

20

# Web Archiving (Big) Data

- Entire sites, documents, media, files
- Massive, complex, heterogeneous datasets, often measured in petabytes.
- Volume (billions of URLs harvested, longitudinal captures over the years)
- Variety (structured, semi-structured, unstructured content; media, scripts)
- Velocity (continuous capture of rapidly changing web content)

"Saving the Web"
▸ Slides: bit.ly/info591

**Mat Kelly, Asst. Prof., Drexel IS**
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025

21

# Composite Resources

- HTML → CSS, images, media → deeper embedded resources
- Content may appear dynamically based on user interaction

**Browser**

↓

**HTML Document**
(index.html, etc.)

CSS  JS  IMG  FONT

"Saving the Web"
▸ Slides: bit.ly/info591

**Mat Kelly, Asst. Prof., Drexel IS**
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025

22

# Indexing WARCs

- Fast lookup of a URL w/o scanning terabytes of WARC data
- Reply systems (e.g., Wayback, pywb) for reconstructing a page quickly
- Efficient queries over billions of archived documents
- Without indexing, random access is impossible
- Single WARC file can hold thousands of individual resources
- Index formats:
    - CDX (classic, well-defined fields)
    - CDXJ (modern, more flexible, supports extra metadata fields)
    - Database indexes (performance at scale, full-text search and complex queries)

"Saving the Web"
▸ Slides: bit.ly/info591

**Mat Kelly, Asst. Prof., Drexel IS**
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025

23

# Index Examples

```
edu,drexel)/ 19970626040823 http://www.drexel.edu:80/ text/html 200 NA4IS3YV4BOPOE3KRG6VMNQLFFTZMWJG 1814
edu,drexel)/ 20220214060114 https://www.drexel.edu/ text/html 301 3I42H3S6NNFQ2MSVX7XZKYAYSCX5QBYJ 459
edu,drexel)/ 20230709040752 https://drexel.edu/ text/html 200 XMW7AUSKNSWO2WMOML3TH2I3HVRVVHGD 30966
edu,drexel)/ 20251005193011 https://drexel.edu/ text/html 200 H5GRERCRQSZOWAEUV5APC7U67W2ZO7LQ 32444
```

datetime          content-type          WARC locator ID

original URI captured

Sort-friendly URI Reordering Transformation (SURT)          Offset for record in WARC

HTTP status code

"Saving the Web"
▸ Slides: bit.ly/info591

**Mat Kelly, Asst. Prof., Drexel IS**
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025

# Deduplication

- Identify and handle duplicate or near-duplicate content
- Storage reduced, replay accurate, efficient indexing, avoid unnecessary redundancy
- Faster crawling, more efficient indexing, improved data quality with reduced noise, better preservation integrity
- Can be performed at crawl time or during replay
- Types:
  - URL-level (same page at different URLs)
  - Payload (identical content at different hosts)
  - Temporal (repeated captures unchanged over time)
  - Structural (template pages w/ minimal diffs)



Screenshots of apple.com of the past

"Saving the Web"
▸ Slides: bit.ly/info591

Mat Kelly, Asst. Prof., Drexel IS
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025

25

# Resilience of Web Archives

- Wholesale or piecemeal copies of archives with means to check integrity
- How can you verify what is being purported as history is truly representative?
- Does authenticity matter if used for personal reference?



https://github.com/oduwsdl/ipwb



http://foo.com/spaceDog.jpg → IPFS → QmZAD4xeeNeYF3TmwWgBXypLKTiCGwGRMXHW7MtheWKtw4

===

QmZAD4xeeNeYF3TmwWgBXypLKTiCGwGRMXHW7MtheWKtw4

http://example.org/yuri.jpg → IPFS →

"Saving the Web"
▸ Slides: bit.ly/info591

**Mat Kelly, Asst. Prof., Drexel IS**
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025          26

# Handling Dynamicism

- The delta between a crawler and a browser causes resources to be missed
- JavaScript-driven, dynamic sites are hard to capture.

doi:10.1007/978-3-642-40501-3_5     doi:10.1007/s00799-015-0140-8     doi:10.1109/JCDL.2014.6970146

# Web Archive Aggregation



Archives Queried ($A_0$)

"Saving the Web"
▸ Slides: bit.ly/info591

**Mat Kelly, Asst. Prof., Drexel IS**
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025

28

# Where Mat Fits

# You Are Responsible for Saving What's Important

"Saving the Web"
▸ Slides: bit.ly/info591

**Mat Kelly, Asst. Prof., Drexel IS**
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025

30

# Most Web Archiving Tools Cater to Archiving the Public Web

- The face of society, Zeitgeist
- Easier to vet – captures from multiple institutions give legitimacy of the record
- Institutional self-interest
    - e.g., Stanford University Archives, UK Web Archive

"Saving the Web"
▸ Slides: bit.ly/info591

**Mat Kelly, Asst. Prof., Drexel IS**
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025

31

# The Fallacious Key-Based Approach



- A URL is not enough
- Personalized Representations
- Dynamic web sites
- Context lost



facebook.com capture with browser-based archiving tools



facebook.com capture from Internet Archive

"Saving the Web"
▸ Slides: bit.ly/info591

**Mat Kelly, Asst. Prof., Drexel IS**
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025

32

# 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

"Saving the Web"
▸ Slides: bit.ly/info591

Mat Kelly, Asst. Prof., Drexel IS
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025

33

# 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**

"Saving the Web"
▸ Slides: bit.ly/info591

Mat Kelly, Asst. Prof., Drexel IS
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025

34

# Crawlers 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

"Saving the Web"
▸ Slides: bit.ly/info591

Mat Kelly, Asst. Prof., Drexel IS
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025

35

# By-Value Archiving with Posthoc Metadata Ascription

- Most true to the original form

- Capture complex content

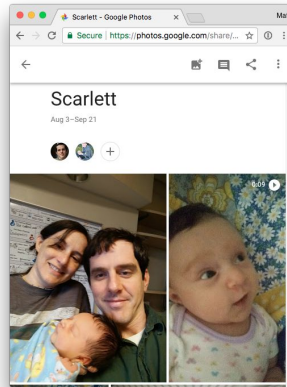- Personalized representations archivable

- Easier to fabricate content (not a good thing)
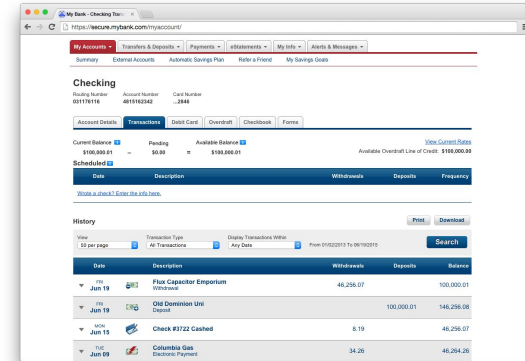
    - Requires a means of vetting authenticity

"Archive What I See Now: Bringing Institutional Web Archiving Tools to the Individual Researcher", National Endowment for the Humanities Digital Humanities Implementation Grant (HK-50181-14)

"Saving the Web"
▸ Slides: bit.ly/info591

**Mat Kelly, Asst. Prof., Drexel IS**
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025

36

# Archiving Private and Personal Web Content

- The things we care about may not belong in an institutional web archiving
- Bearer of our own important information
- Stewardship, ownership, copyright
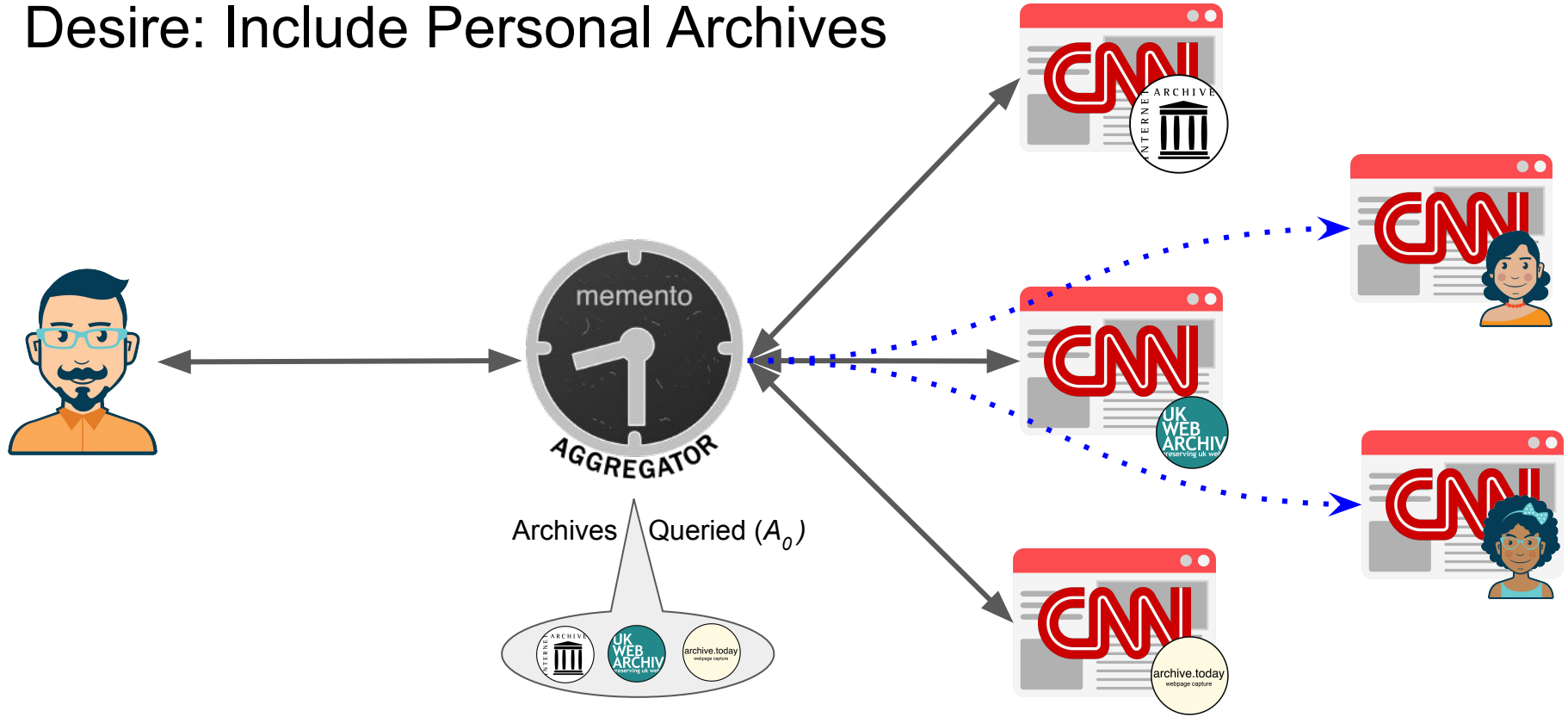- Integration into conventional public web archive corpora


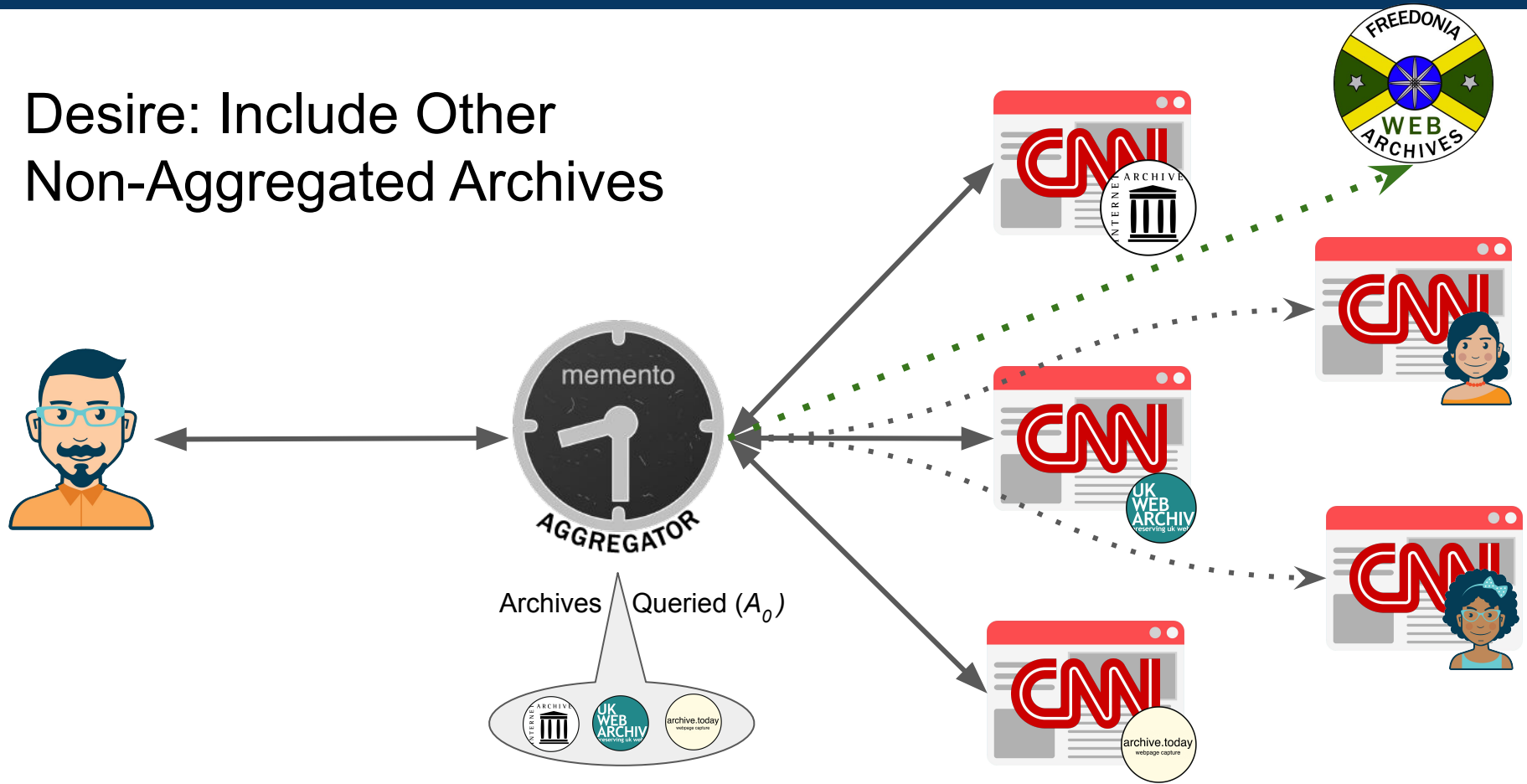Historical capture of Google Photos


Historical capture of Online Banking Interface

"Saving the Web"
▸ Slides: bit.ly/info591

Mat Kelly, Asst. Prof., Drexel IS
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025

37

# Desire: Include Personal Archives



Archives Queried ($A_0$)

"Saving the Web"
▸ Slides: bit.ly/info591

**Mat Kelly, Asst. Prof., Drexel IS**
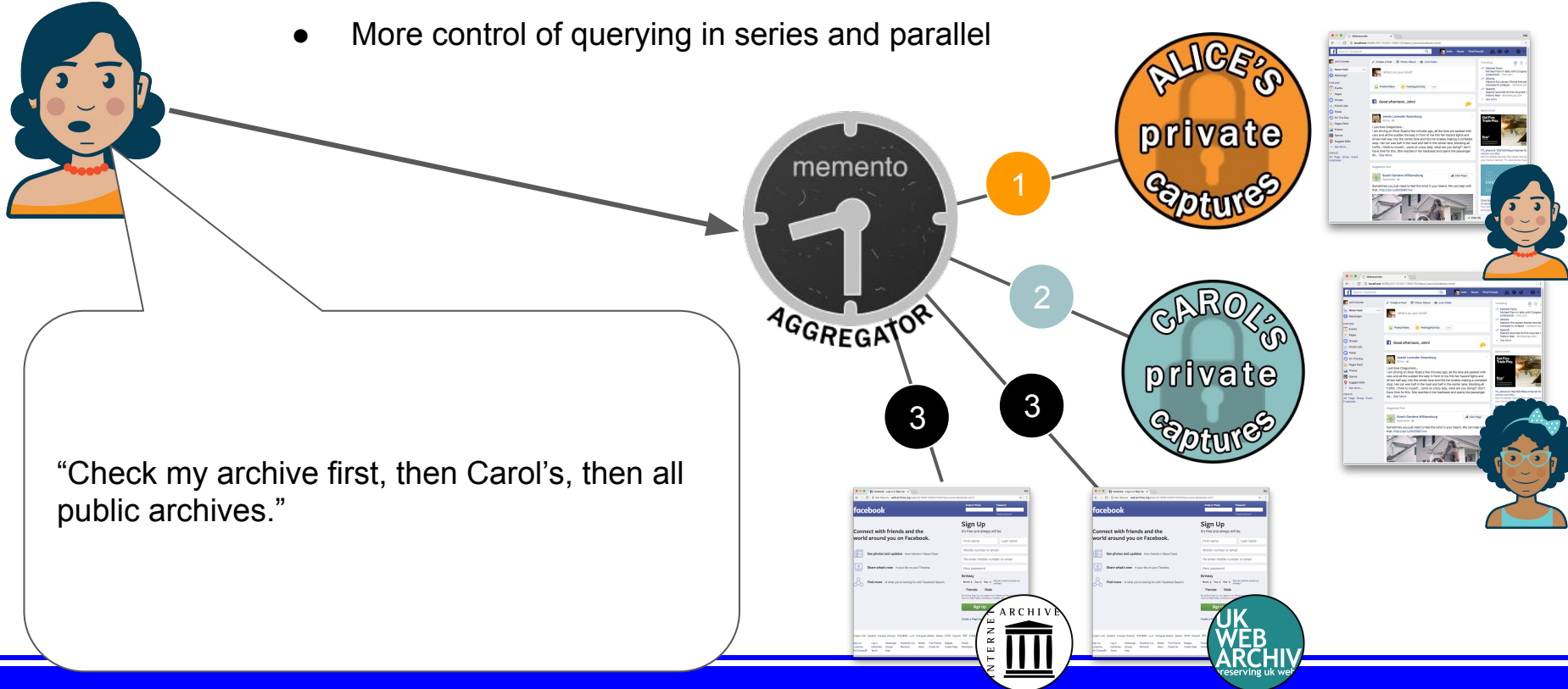@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025

38

# Desire: Include Other Non-Aggregated Archives

Archives Queried ($A_0$)

"Saving the Web"
▸ Slides: bit.ly/info591

**Mat Kelly, Asst. Prof., Drexel IS**
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025

39

# Query Precedence



- More control of querying in series and parallel

"Check my archive first, then Carol's, then all public archives."
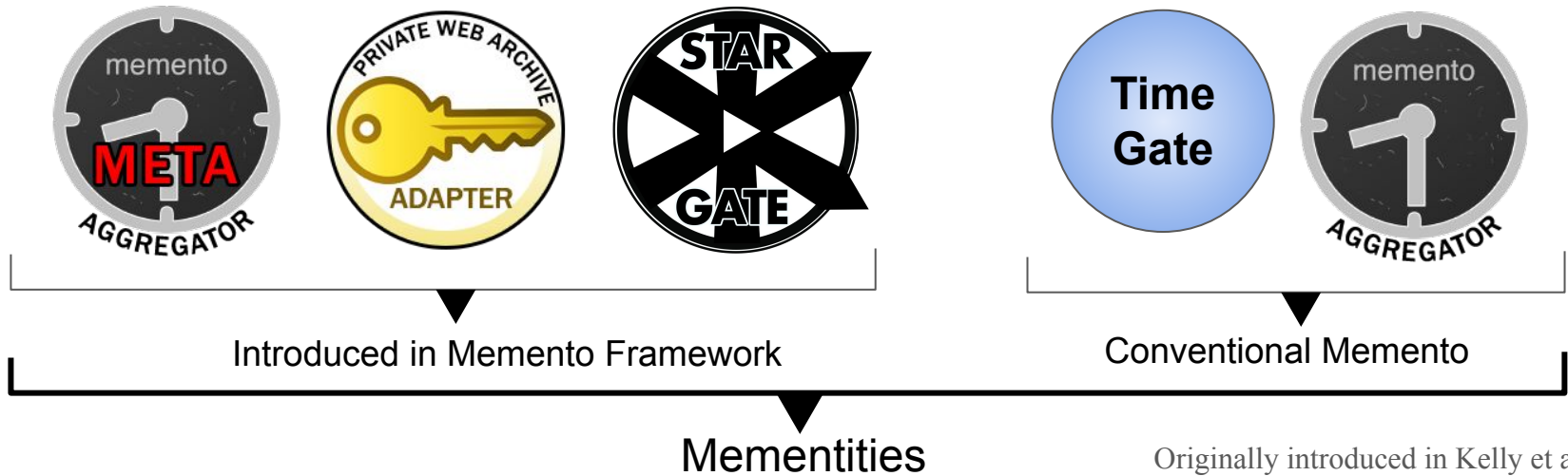
# Query Short-Circuiting

- May give priority to archive relevancy.
- **Series halt when threshold met.**

"Check private archives first. **Iff** you find no captures, only *then* check the public archives.

# Mementities

- Memento + Entity (*entity* term already overused in web parlance)

Introduced in Memento Framework

Conventional Memento

Mementities

Originally introduced in Kelly et al. 2018

"Saving the Web"
▸ Slides: bit.ly/info591

**Mat Kelly, Asst. Prof., Drexel IS**
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025          42

# Future Access

- Should these private photos/bank statements be accessible to anyone?
    - Strawman answer: no, but how can this be systematically implemented to be both sophisticated for security and privacy and usable by the non-technical
- How do we balance ability to access in the future with privacy and security?
- Should a page's public/private access now persist into the future?
- Reducing accessibility now can be at odds with any accessibility in the future

"Saving the Web"
‣ Slides: bit.ly/info591

**Mat Kelly, Asst. Prof., Drexel IS**
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025
43

# 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



URL

"Saving the Web"
▸ Slides: bit.ly/info591

Mat Kelly, Asst. Prof., Drexel IS
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025

44

# Outstanding Problems Under Investigation

*A few of many*

- Web Advertisements are intentionally lost but are important
  - AI-driven personas to surface deep web content; reusing browser profile to inform crawls
- Private/Personal web archive resilience through replication with security
- Resilience of remote resources through sophisticated linked data structures
- User-informed preference for information retrieval of the past in dimensions beyond time
  - e.g., topical, quality, ad hoc supplied lambda functions
- *Better* defining fuzziness of public/private and personal/institutional captures to better align with non-boolean reality



IPNS => https://web.archive.org/dweb/timemap/https://example.com/

First            Latest



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

Retrieved StarMap

memento
META
AGGREGATOR

STAR GATE

CONTENT-BASED ATTRIBUTE
&
DERIVED ATTRIBUTE

UK WEB ARCHIVE
archive.today
Open Wayback
Open Wayback

Abbreviated StarMap with filtering applied

-30-

"Saving the Web"
▸ Slides: bit.ly/info591

Mat Kelly, Asst. Prof., Drexel IS
@machawk1 • https://matkelly.com

INFO591 Week 10
December 4, 2025