# Archive What I See Now
## Personal Web Archiving with WARCs

Michele C. Weigle, Michael L. Nelson, **Mat Kelly**, and John Berlin
Web Science and Digital Libraries (WS-DL) Research Group
Old Dominion University
ws-dl.cs.odu.edu • @WebSciDL

**@machawk1**

# Web Archiving Tools for Web Users

Standard Web archiving tools are difficult for non IT experts.

"Save Page As" is not suitable for archiving purposes.

Pages are behind authentication.

Pages change quickly, but current state needs archiving.

**ARCHIVE**

**WHAT I SEE**

**NOW**

http://bit.ly/iipcWAC2017

# Why?

- Allow non-technical users to locally create+replay own archives
- Preserve the previously unpreserved

**more archives → more better**

http://bit.ly/iipcWAC2017

@machawk1

NATIONAL ENDOWMENT FOR THE
Humanities

# CREATION

# +

# ACCESS

**of personal and private web archives**

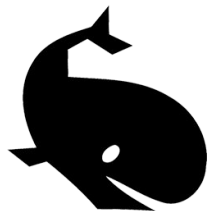http://bit.ly/iipcWAC2017

**@machawk1**

NATIONAL ENDOWMENT FOR THE
HUMANITIES

# Goals: Advance Development of 3 Tools

## WARCreate
Create a WARC from what you see in your browser

## Web Archiving Integration Layer (WAIL)
Replay the WARC using software of your desktop
Your captures never leaves your machine

## Mink
See how your captures temporally integrate with institutions'
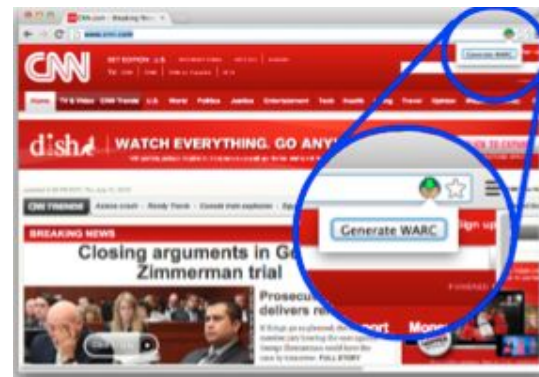Submit new URIs to Web archives (was to-WAIL in scope?)

NATIONAL ENDOWMENT FOR THE
Humanities

**WARCreate**

# WARCreate

- Google Chrome browser extension
- Save WARC files from your browser
- No credentials pass through 3$^{rd}$ party
- Heavily leverages Chrome webRequest API
- Built in '12, APIs and libraries have evolved!

http://bit.ly/iipcWAC2017

@machawk1

NATIONAL ENDOWMENT FOR THE
Humanities
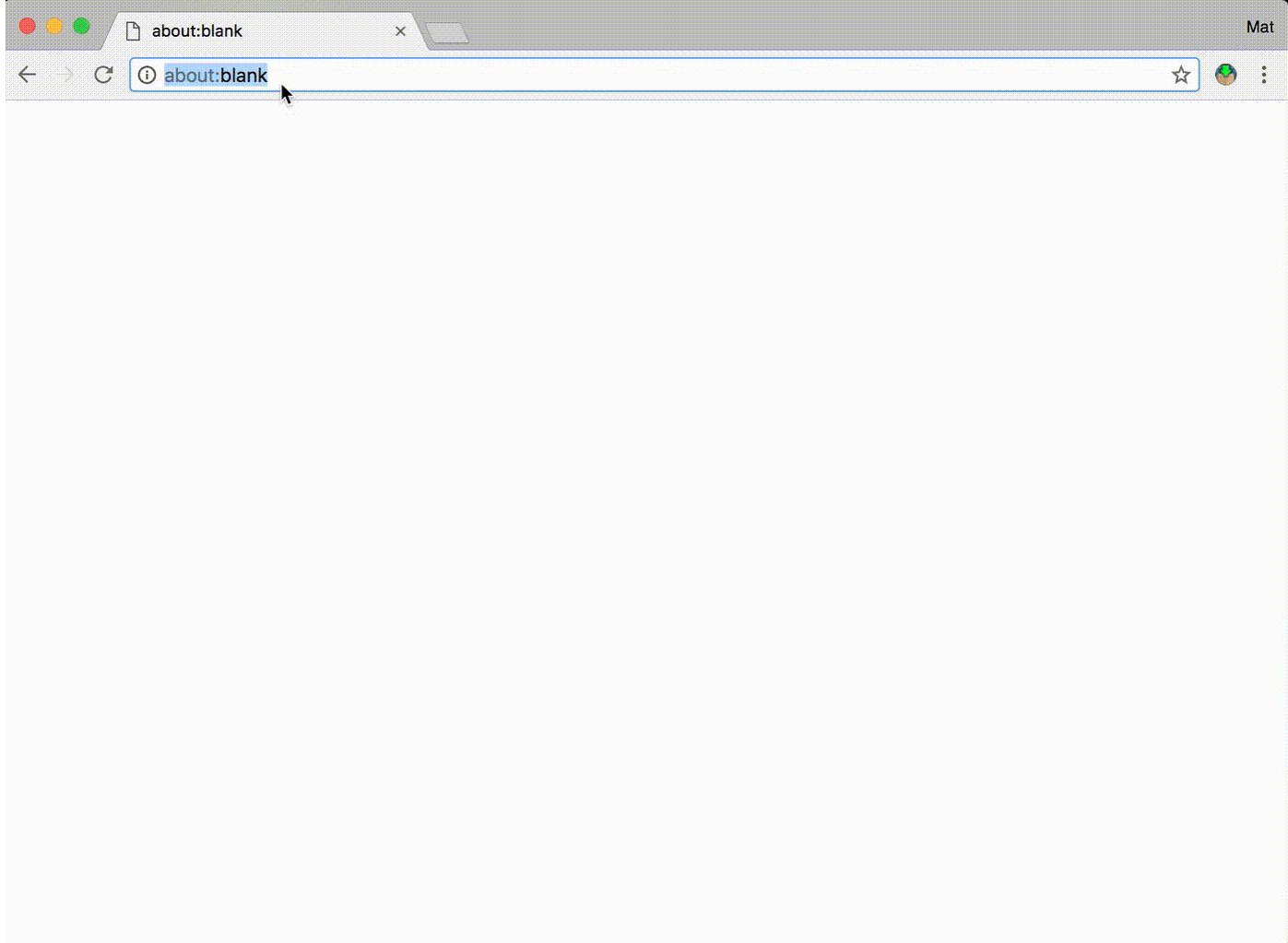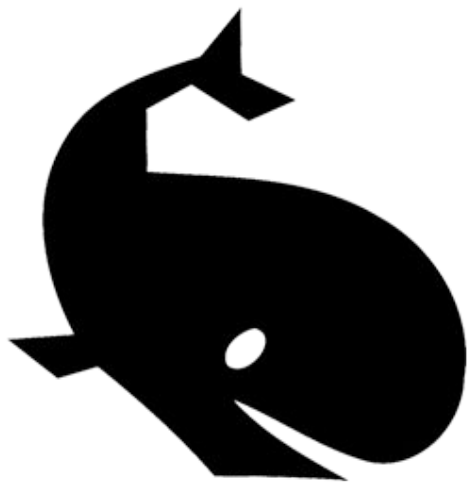HD-51670-13 • HK-50181-14

# WARCreate - Recent Advancements

- Three New Modes for Browser-Based Preservation
  - Record Mode - retain buffer as you browse
  - Countdown Mode - preserve reloading page on an interval
  - Event Mode - preserve page when it's automatically reloaded
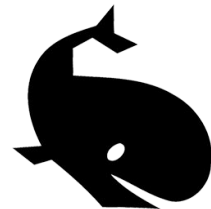- Save to local Web archive (e.g., WAIL)

http://bit.ly/iipcWAC2017

@machawk1

**Web Archiving Integration Layer (WAIL)**
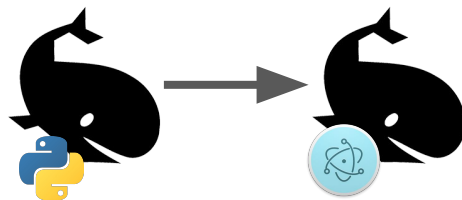
# Web Archiving Integration Layer (WAIL)

- Stand-alone desktop application
- Collection-based Web Archiving
- Includes Heritrix for crawling, OpenWayback for Replay
- Python scripts compiled to OS-native binaries (.app, .exe)

- What to do with WARCs?

- See: [How WAIL came about, "Lipstick or Ham"](#)

NATIONAL ENDOWMENT FOR THE
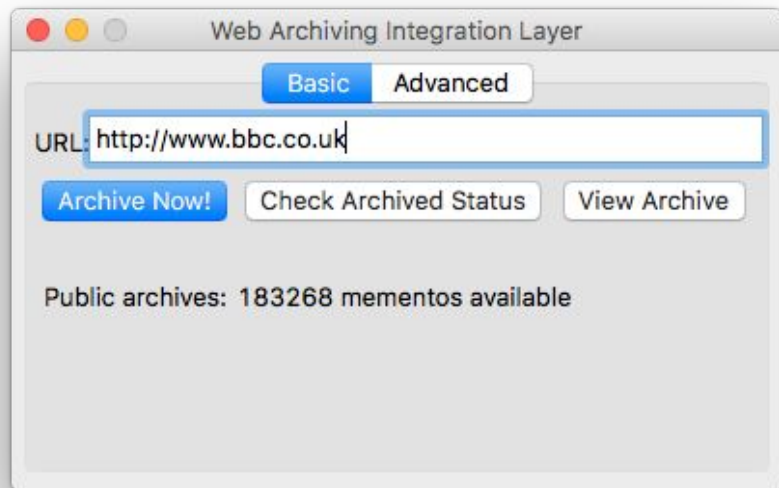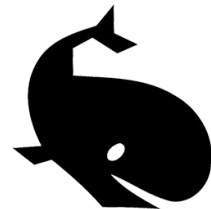HUMANITIES
HD-51670-13 • HK-50181-14
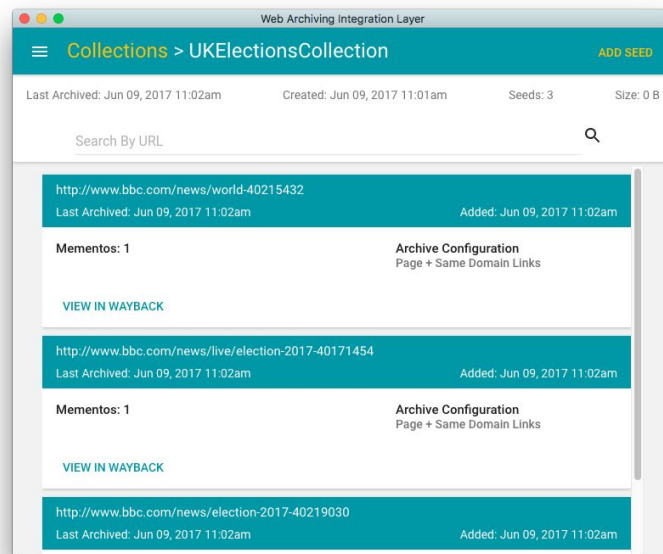
# WAIL - Recent Advancements

- New User Interface
- Ported from Python to Electron
  - Now using Web technologies to archive the Web
- Single archive to collection-based archiving
- OpenWayback to pywb
- Twitter integration

WAIL-Electron Feature Walk-through

NATIONAL ENDOWMENT FOR THE
Humanities
HD-51670-13 • HK-50181-14

# WAIL - New User Interface



**Original one-click interface**



**New collection-based interface**

http://bit.ly/iipcWAC2017

**@machawk1**

NATIONAL ENDOWMENT FOR THE
**Humanities**
HD-51670-13 • HK-50181-14

Collections > UKElectionsCollection

ADD SEED

Last Archived: Jun 09, 2017 4:02pm          Created: Jun 09, 2017 4:01pm          Seeds: 3          Size: 0 B

Search By URL

http://www.bbc.com/news/world-40215432

Last Archived: Jun 09, 2017 4:02pm                                    Added: Jun 09, 2017 4:02pm

Mementos: 1                                                           Archive Configuration
                                                                      Page + Same Domain Links

VIEW IN WAYBACK

http://www.bbc.com/news/live/election-2017-40171454

Last Archived: Jun 09, 2017 4:02pm                                    Added: Jun 09, 2017 4:02pm

Mementos: 1                                                           Archive Configuration
                                                                      Page + Same Domain Links

VIEW IN WAYBACK

http://www.bbc.com/news/election-2017-40219030

Last Archived: Jun 09, 2017 4:02pm                                    Added: Jun 09, 2017 4:02pm

Mementos: 1                                                           Archive Configuration
                                                                      Page + Same Domain Links

# Web Archiving Integration Layer

## Crawls

| Crawl URL(s) | Status | For Collection | Timestamp | Discovered | Queued | Downloaded | Actions |
|---|---|---|---|---|---|---|---|
| http://www.netpreserve.org/general-assembly/2017/overview | Ended | default | Jun 09 2017 12:03am | 147 | 117 | 30 | ⋮ |
| http://www.netpreserve.org/general-assembly/2017/overview | Ended | default | Jun 08 2017 11:40pm | 167 | 134 | 31 | ⋮ |
| http://www.cs.odu.edu/~mweigle/ | Ended | Radon | May 18 2017 12:05am | 44 | 34 | 9 | ⋮ |
| http://www.cs.odu.edu/~mln/ | Ended | Radon | May 17 2017 9:33pm | 27 | 0 | 27 | ⋮ |
| http://www.cs.odu.edu/~mln/ | Ended | Radon | May 17 2017 9:33pm | 27 | 0 | 27 | ⋮ |
| http://www.cs.odu.edu/~mkelly/ | Ended | Radon | May 17 2017 9:50pm | 244 | 128 | 117 | ⋮ |
| http://www.cs.odu.edu/~melly/ | Ended | Radon | May 17 2017 9:29pm | 4 | 0 | 4 | ⋮ |
| http://www.cs.odu.edu/~mkelly/ | Ended | Test2 | May 16 2017 8:20pm | 0 | 0 | 0 | ⋮ |

RESCAN JOB DIRECTORY ≡✓

LAUNCH WEB UI ⬆

IIPC Web Archiving Conference 2017
June 15, 2017
London, UK

http://bit.ly/iipcWAC2017

@machawk1

NATIONAL ENDOWMENT FOR THE
Humanities
HD-51670-13 • HK-50181-14

Web Archiving Integration Layer

Archive Twitter

Info
Monitoring @netpreserve's Timeline For 5
minutes. Tweets Added To
UKElectionsCollection!

Monitoring & Archiving Configuration

Required

ScreenName

NEXT

Info

Monitoring @netpreserve's Timeline For 5
minutes. Tweets Added To
UKElectionsCollection!

IIPC Web Archiving Conference 2017
June 15, 2017
London, UK

http://bit.ly/iipcWAC2017

@machawk1

NATIONAL ENDOWMENT FOR THE
Humanities
HD-51670-13 • HK-50181-14

memento

MINK
INTEGRATING the WEBS

# Mink

- Google Chrome browser extension
- Indicates archival capture count as you browse
- Quickly submit URI to multiple archives from UI
- From Mink(owski Space)

http://bit.ly/iipcWAC2017

@machawk1

NATIONAL ENDOWMENT FOR THE
Humanities
HD-51670-13 • HK-50181-14

# Mink - Recent Advancements

- Enhance interface
  - Add number of archived pages to icon at bottom of page
  - Allow users to set preferences on how to view large set of mementos
- Communication with user-specified (or local) archive in additional to aggregated institutional archives' results

# Mink - Previous Interface



➢ Interface affected by page CSS

➢ Obtrusive on the viewport by default

➢ Haphazard, inconsistent animations

NATIONAL ENDOWMENT FOR THE
Humanities
HD-51670-13 • HK-50181-14

# Mink - User Interface Revamp

# Mink - User Interface Revamp



➤ Interface-on-demand

➤ Shadow DOM, no CSS intrusion

➤ More consistent, intuitive Miller columns for many captures

HD-51670-13 • HK-50181-14

# Mink - User Interface Revamp

HD-51670-13 • HK-50181-14

# Mink - Communication with Local Archives

HD-51670-13 • HK-50181-14

Mink usage GIF, also available at:
https://youtu.be/bGjxofpTgv4

http://bit.ly/iipcWAC2017

# Tools' Integration

IIPC Web Archiving Conference 2017
June 15, 2017
London, UK

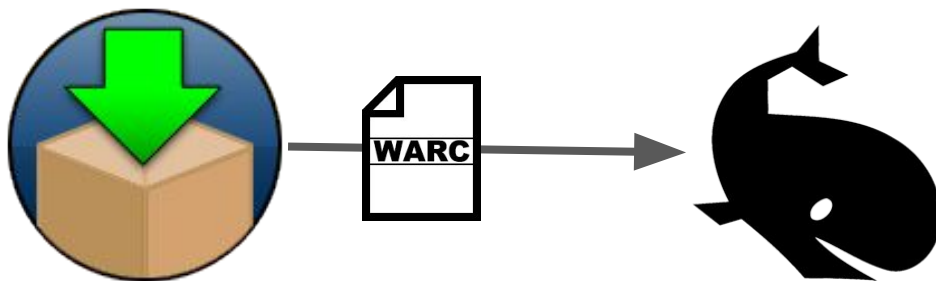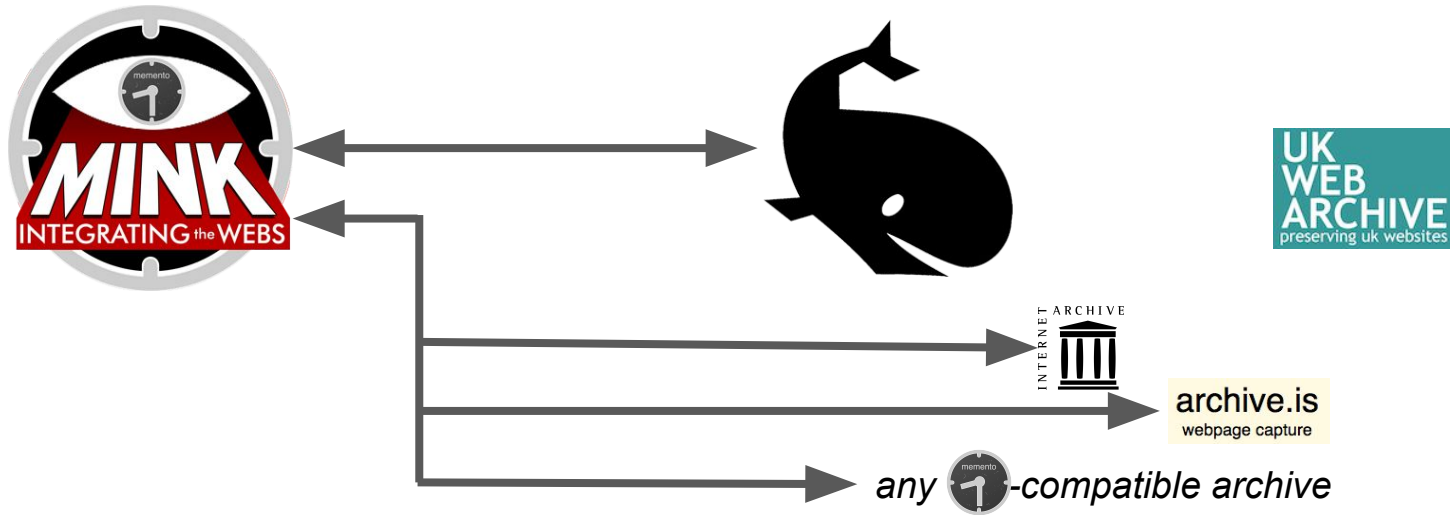http://bit.ly/iipcWAC2017

@machawk1

NATIONAL ENDOWMENT FOR THE
Humanities
HD-51670-13 • HK-50181-14

# Tools' Integration: WARCreate→WAIL

- Save WARC directly to local archive (by reference [easier]
  - By-value integration feasibility being investigated a la WASAPI
- Automatically indexed and replayable

http://bit.ly/iipcWAC2017

**@machawk1**

NATIONAL ENDOWMENT FOR THE

**Humanities**

**HD-51670-13 • HK-50181-14**

# Tools' Integration: Mink→WAIL
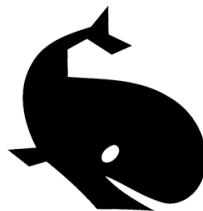
# Some Future Work

- Decouple Mink from external Memento aggregator
  - Client-side customizable set of archives instead
- WARC replay using browser extensions/apps
- Further integration with other archiving tools in WAIL
  - Re-add Memgator Memento aggregator (removed from Electron version)
- Firefox version of tools
  - XUL→WebExtensions
  - Decouple from Chrome APIs
- Integration with InterPlanetary Wayback (speaking about later today)

NATIONAL ENDOWMENT FOR THE
Humanities
HD-51670-13 • HK-50181-14

# Acknowledgements

- NEH Grant #s [HD-51670-13]{.link} • [HK-50181-14]{.link}
- Dr. Liza Potts and WIDE Research Center at Michigan State University
- ODU SEES Travel Grant

http://bit.ly/iipcWAC2017

**@machawk1**

NATIONAL ENDOWMENT FOR THE
**Humanities**
**HD-51670-13 • HK-50181-14**

# Archive What I See Now
## Personal Web Archiving with WARCs

Michele C. Weigle, Michael L. Nelson, **Mat Kelly**, and John Berlin
Web Science and Digital Libraries (WS-DL) Research Group
Old Dominion University
ws-dl.cs.odu.edu

**@machawk1**

http://bit.ly/iipcWAC2017