# Web Archiving Workflows

This training session was developed in partnership by the International Internet Preservation Consortium (IIPC) and the Digital Preservation Coalition (DPC)
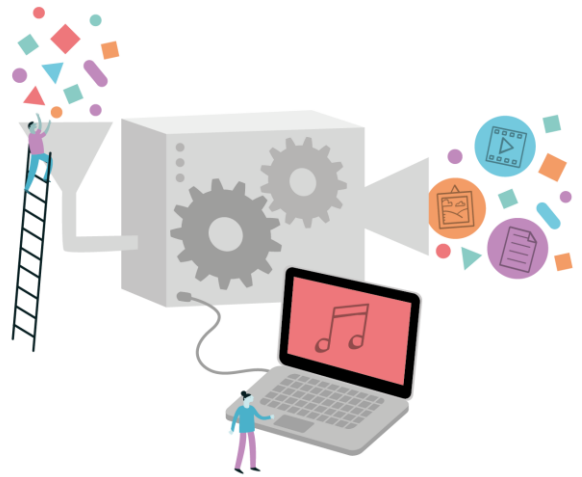
1

## Workflows We Will Cover

**Capture:** live web content is downloaded & stored

**Preserve:** downloaded files are checked, converted to a stable file type if necessary, and looked after over time

**Playback:** the archived web content is accessed through a tool that allows users to interact with it like the original

2

# Capture: Downloading and Storing Live Web Content
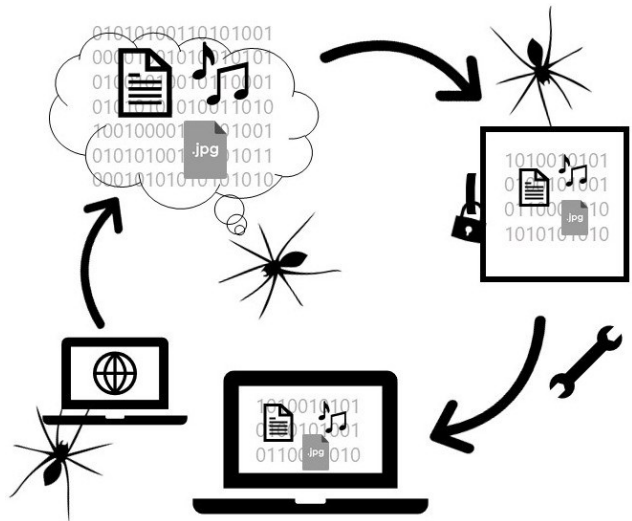
3

# Intro to Crawling

Crawler or "Spider"

Code & files needed to reproduce original website

Playback Tool
(like Wayback Machine)

4

# Basic Crawl

- A tool ('crawler') systematically browses the web
- Uses a set of parameters, or defined scope (e.g. from a seed list)
- Downloads code, images, documents, and other files
  - Whatever is essential to reproducing the web content as similarly to original form as possible

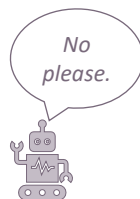Also referred to as "harvesting", "web crawling", "spidering", or "an internet bot"

5

# Crawling: Parameters

- **Seed URLs** or **URIs**: starting point(s) for web crawler; the crawler follows links out from this initial URL or set of URLs
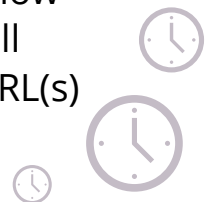
- **Robots.txt**: a file included in web content that instructs a crawler not to capture that content or to capture only parts of it

*No please.*

- **Crawl Depth** or "**Hops**": number of links away from the seed URL/URI the crawler will capture

- **Crawl Frequency**: how often the crawler will capture the same URL(s)

6

## Other Capture Methods

### Dynamic Capture

- Tools built to capture interactive or complex content
  - Ex. videos & other media
  - Ex. social media and other platform-based web
  - Ex. complex JavaScript
- Tools like Webrecorder/Conifer, Internet Archive's Brozzler

### API Harvesting

- Only available for web resources that provide an API
- An API allows authenticated users to extract data directly from the platform through the web
- Works for the modern "platformized" web
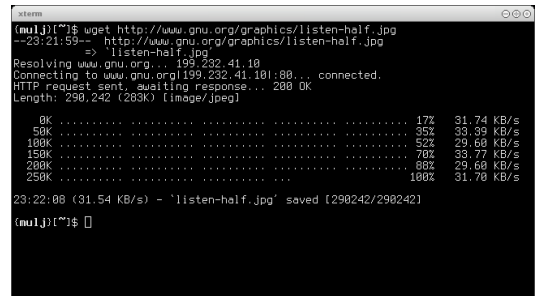- Tools for using APIs: Twarc, Social Feed Manager, others

7

## Tools to Support Capture of Web Archives



8

# Tools: GNU wget

- Command line tool that downloads files from the web
- Runs on Unix and Mac OS, but also has a Windows version
- Supports HTTP, HTTPS, and FTP
- Operates continuously in the background
- Usable on slow or unstable networks
- Allows scoping & configuration

- Supports writing to a WARC file
- Free and open source under GNU General Public License

```
xterm                                                                    ⊖⊝⊘
(mulj)[~]$ wget http://www.gnu.org/graphics/listen-half.jpg
--23:21:59--  http://www.gnu.org/graphics/listen-half.jpg
           => `listen-half.jpg'
Resolving www.gnu.org... 199.232.41.10
Connecting to www.gnu.org|199.232.41.10|:80... connected.
HTTP request sent, awaiting response... 200 OK
Length: 290,242 (283K) [image/jpeg]

     0K .......... .......... .......... .......... ..........  17%  31.74 KB/s
    50K .......... .......... .......... .......... ..........  35%  33.39 KB/s
   100K .......... .......... .......... .......... ..........  52%  29.60 KB/s
   150K .......... .......... .......... .......... ..........  70%  33.77 KB/s
   200K .......... .......... .......... .......... ..........  88%  29.60 KB/s
   250K .......... .......... .......... ...               100%  31.70 KB/s

23:22:08 (31.54 KB/s) - `listen-half.jpg' saved [290242/290242]

(mulj)[~]$ []
```

9

# Tools: Heritrix (+ Umbra)

**Heritrix**
- From the Internet Archive
- Web crawler that downloads websites and embedded media
- Suitable for large collections
- Available for Windows and Unix-like environments
- Supports configurable scoping and deduplication
- Supports writing to a WARC file
- Less effective at triggering and capturing client side script

**Heritrix + Umbra**
- Browser automation tool that runs alongside Heritrix
- Mimics the way a browser would access a page
- Executes client side scripts so previously undetectable URLs can be accessed
- Supports the capture of JavaScript
- Allows for dynamic scrolling
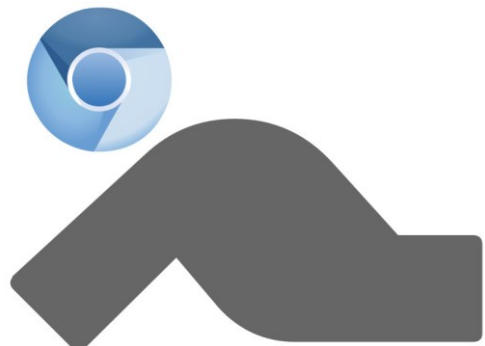
10

## Tools: Heritrix-Based Curator Tools

- Examples:
  - Archive-It
  - Web Curator Tool
  - NetArchive Suite
- Run Heritrix with user interfaces that make it easier to manage collections
- Many used by IIPC members
- Some curator tools are subscription-based

11

## Tools: Conifer

- "Browser" + "Crawler" = "Brozzler"
- Captures HTTP traffic as it loads
- Uses a real browser to fetch pages and embedded URLs, and to extract links
- Implements a tool called youtube-dl to improve media capture
- Requires: Python 3.4 or later, RethinkDB deployment, Chromium or Google Chrome version 64 or higher

12

## Tools: Webrecorder/Conifer

- User-driven capture rather than automated crawler
- Focus on dynamic web content (embedded video and JavaScript)
- Simple to use interface
- Captures page by page – can be labour intensive!
- Can be structured by Collection and capture session
- Captures can be downloaded as WARC files

**Conifer** = Hosted Service from Rhizome
- Up to 5GB of free storage
- Some use-cases and integrations may require additional support or storage that requires a fee
- Quick Start Guide: https://guide.conifer.rhizome.org/

**Webrecorder** = Desktop App
- Same functionality as Conifer but on local desktop
- Slower, but only limited by local storage…

13

## Tools: Social Feed Manager

- Open source tool created by George Washington University Libraries
- Harvests data from Twitter, Flickr, Sina Weibo, and Tumblr
- Captures data through platform APIs
- Captures linked URLs and embedded media
- Supports the curation and management of archived collections

14

# Tools: Social Media Download

- Available for Twitter, Facebook (limited), Google, and others
- Function in Settings
- Only permitted for the account owner
- Good practice for institutions with one or more public-facing social media accounts
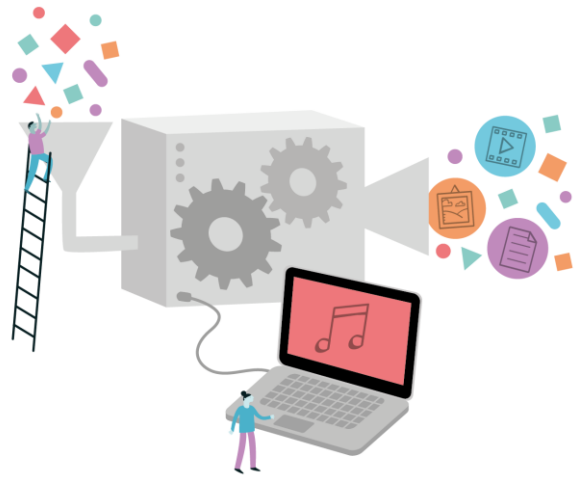- Good practice for personal digital archiving

15

# Capture: MOAR TOOLS!

16

# Preserve:
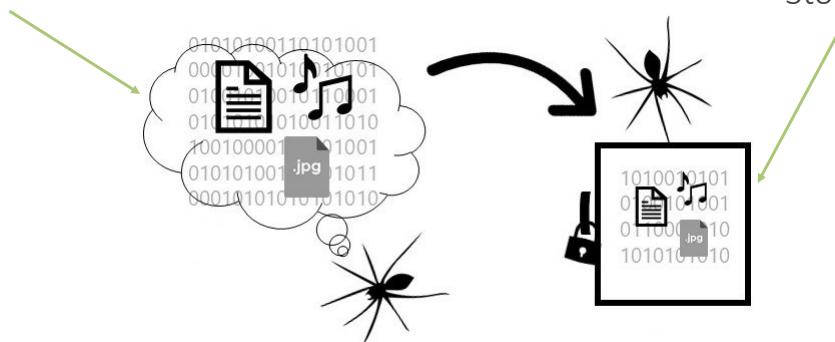# What Happens to Captured Content?

17

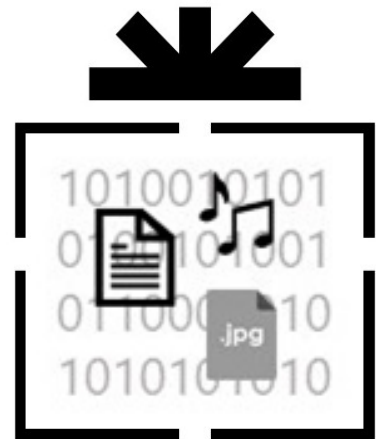# Capture to Preservation

Captured Content
or
"The Crawl"

Web Archive
Preserved in Safe
Storage

18

## Introduction to WARC

- WARC (Web ARChive)
- WARC is a wrapper for archived web objects developed by the IIPC
- Tools that write to WARC create files with the extension .warc
- A WARC file can be ingested into a digital preservation system
- WARC was preceded by the ARC (.arc) format

19

## WARC Standard

- File format standard
- ISO 28500:2017 (formerly ISO 28500:2009)
- Packages together multiple files of different types from a web crawl or capture
- Maintains and describes relationships between web pages or related content
- Accommodates different forms of metadata
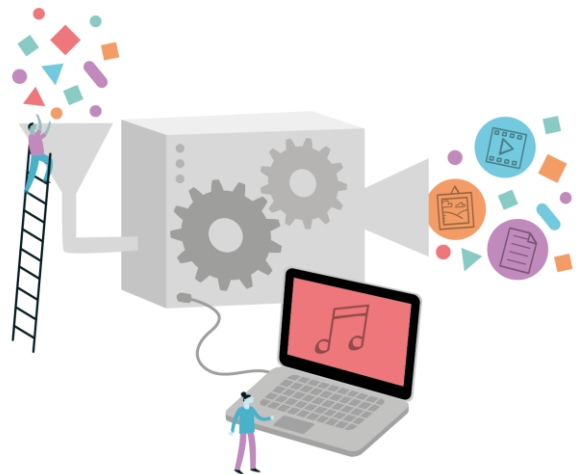- Requires special access tools or viewers

20

## Preserve: Actions

- Quality assurance
  - Successfully completed capture?
  - Capture content complete?
  - Sensitive data review
- "Patching" any issues
- Generating metadata
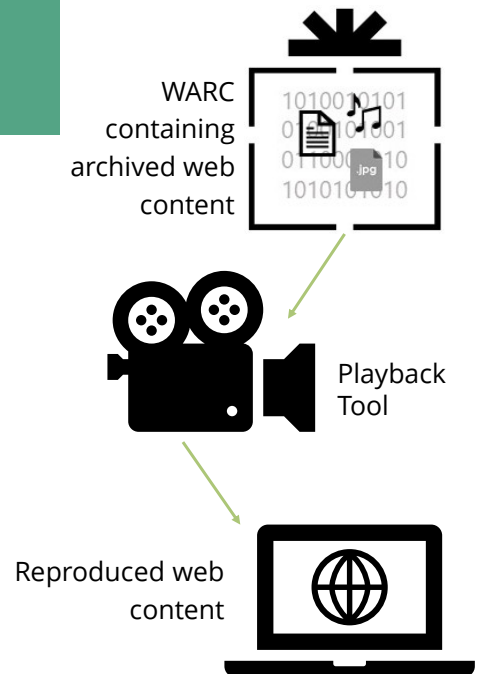- Transfer to archival storage

21

## Playback: Providing Access to Archived Web Content

22

## How Do We Provide Access?

- Playback Tools required
- Designed to render archived web content
- Read and display WARC files
- Examples:
  - Wayback Machine
  - ReplayWeb (dynamic, interactive)
  - Third-party service platforms

WARC containing archived web content

Playback Tool

Reproduced web content

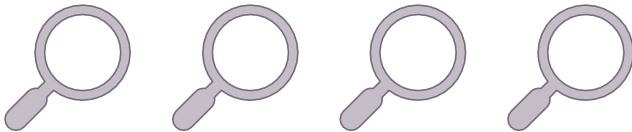23

## Playback: Facilitating Use

24

# Indexing and Search

**Indexing**

- Allows the search and retrieval of archived web content
- Enables search based on metadata fields, including keywords
- Required to enable access & re-use of web content

**Full-text Search Index**

- Indexed to allow end users to search broader range of keywords or phrases
- Enhances digital preservation planning
  - UK Web Archive uses full-text search capability to search tags to track the birth and death of specific features like HTML elements

25

# Banner or User Notice

- To designate the viewed content as archived to avoid confusion with the live web

National Records of Scotland

You are viewing an archived web page captured at 1:03:51 Sep 03, 2017, which is part of the National Records of Scotland Web Archive. The information on this web page may out of date. See all captures of this archived web page. We do not use cookies but some may be left in your browser from archived websites. Find out more about cookies.
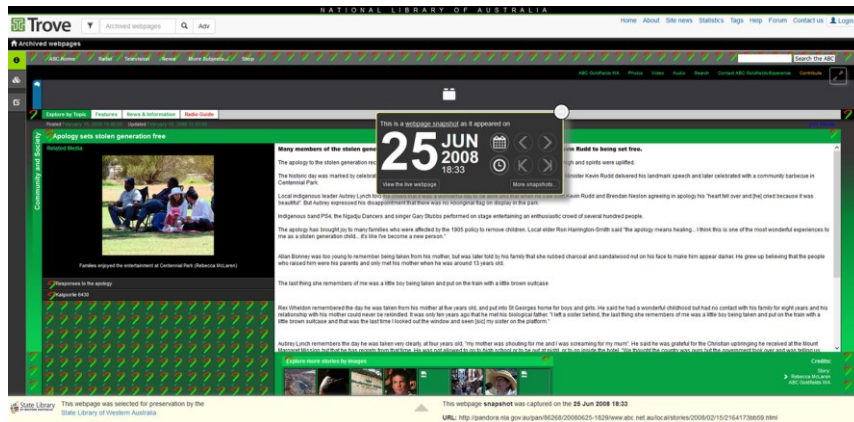
hide

https://webarchive.nrscotland.gov.uk/#!/

26

## Date and Time of Capture

- To compare with concurrent information, such as major events or other publications



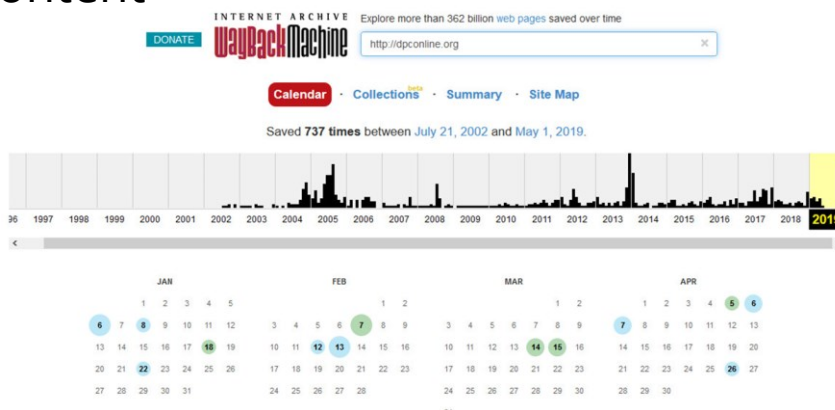https://webarchive.nla.gov.au/awa/20080625083334/http://pandora.nla.gov.au/pan/86268/20080625-1829/www.abc.net.au/local/stories/2008/02/15/2164173bb59.html

27

## Timeline Navigation

- To show the timeline of captures for collections of web content



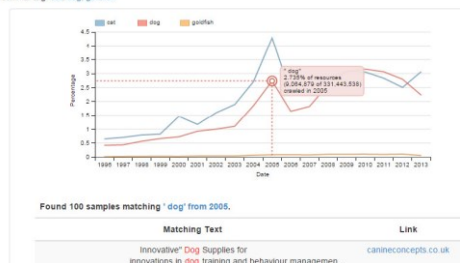https://web.archive.org/web/*/http://dpconline.org

28

## Web Content as Data

- Users may wish to analyze web content or social media data for trends over time or across the web
- The UK Web Archive's SHINE historical search engine is a prototype web-based tool for analyzing trends in web content



https://www.webarchive.org.uk/shine

29

## Playback:
## Tools to Replay
## Web Archives



30

## Wayback Machine & OpenWayback

**Wayback Machine**

- Developed by the Internet Archive
- Used to "play back" archived web content contained in a WARC file in an end user's browser
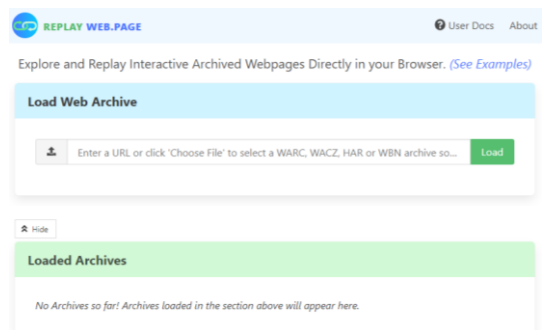- Open source software to query and access archived web content

**OpenWayback**

- Shared development project to address common requirements and improve testing

INTERNET ARCHIVE
WayBackMachine

IIPC netpreserve.org
INTERNATIONAL INTERNET PRESERVATION CONSORTIUM

31

## ReplayWeb.page

- Developed by the Webrecorder project and complements Conifer and Webrecorder
- Available as a browser-based replay tool or an downloadable app
- Browser-based replay tool can be used offline once it is loaded
- Supports WARC file types (.warc, .warc.gz)

REPLAY WEB.PAGE                                    User Docs   About

Explore and Replay Interactive Archived Webpages Directly in your Browser. *(See Examples)*

**Load Web Archive**

Enter a URL or click 'Choose File' to select a WARC, WACZ, HAR or WBN archive so...   Load

Hide

**Loaded Archives**

*No Archives so far! Archives loaded in the section above will appear here.*

https://replayweb.page/

32

16

## Other Tools and Services

- When working with archived web content as data:
  - ArchiveSpark
  - Archives Unleashed Toolkit
- Third-Party Services
  - Archive-It
  - MirrorWeb
  - Hanzo
- For creating Collaborative Collections
  - Momento
  - COBWEB
  - UNT Nomination Tool

33

## Group Discussion

Questions to discuss:

1. Do you currently have web archiving process in place?
2. What tools do you use/have you tested?
3. What would you like to add to your web archiving programme?

34