# Websites, Web Archives, WARCs, and Archived Web Objects

**Grace Thomas**

**Digital Collections Specialist,
Library of Congress
Web Archiving Team**

Designing Storage Architectures
September 10, 2019

# SCALE

**17,440**
web archives available on loc.gov

**3.9 PB**
total size

**1.85 PB**
unique size

**9 Million**
total files on storage

**4.3 Million**
unique files on storage

**18+ Billion**
total archived
web objects

**6.2+ Billion**
unique archived
web objects

# Discussion Points

- **Lifecycle of web archives data from live website to containers on storage to public-facing, replayed content**

- **Using Wayback Machine indexes to study individual resources and groups of individual resources in the archive**

LIBRARY LIBRARY OF CONGRESS

www.xkcd.com, accessed 09/04/2019

.warc

.warc.gz - 1GB

GZIP

# BagIt Bags - 1TB

- Copy on long-term storage

- Copy on presentation storage

LIBRARY
LIBRARY
OF CONGRESS

webarchive.loc.gov/all/*/http://xkcd.com



webarchive.loc.gov/all/20180808065952/https://www.xkcd.com/

Data Lifecycle - Access

# Wayback Indexes

- Sit between the content (W/ARCs) on presentation storage and the Wayback Machine software to replay content

- Each line represents a single archived web object with related metadata

Fields:
- canonized URL
- capture timestamp
- full URL
- reported MIME type of original document
- response code
- checksum
- document length (size)
- W/ARC offset
- W/ARC filename
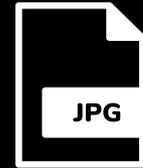
# Using Wayback Indexes Creatively

- Counting total & unique objects
- Qualifying those numbers based on:
  - HTTP response (do 404s count as objects since their records take up space?)
  - Renderable content (200s; uncorrupted)

- Number of objects per "collection" or "web archive" as defined by LC

- Grouping unique objects by MIME type
- Counting objects of certain MIME type per "collection" or "web archive"





www.loc.gov/item/lcwaN0009950/

# Thank you!

Grace Thomas

@gracehthom

grth@loc.gov
webcapture@loc.gov