**LIBRARY** LIBRARY OF CONGRESS

# Library of Congress Web Archiving: Selective Archiving at Scale

**DPC Web Archiving & Preservation Webinar January 30. 2020**

**Abbie Grotke**
**Lead Librarian, Web Archiving Team**
abgr@loc.gov
**@agrotke**

# Program Overview



https://www.loc.gov/programs/web-archiving/about-this-program/

https://www.loc.gov/websites/

http://webarchive.loc.gov/

# Where is Web Archiving organizationally within LC?

```
                          Librarian of
                           Congress
                               |
Congressional   Office of the   Library       OCIO      Copyright
Research        General         Collections and
Service         Counsel         Services Group
                               |
              Law Library    Library Services
                               |
Collection    Acquisition and   Preservation   Digital      Special      General and
Development    Bibliographic                    Services     Collections  International
Office         Access                           Directorate               Collections
                                                    |
                              ILS program      Digital          Business
                              office           Collections      Analysis Team
                                               Management &
                                               Services Division
                                                    |
                              Digitization     Digital Content
                              Services         Management
                                               Section
                                                    |
                                               Web Archiving
                                               Team
                                               (5 FTEs)
```
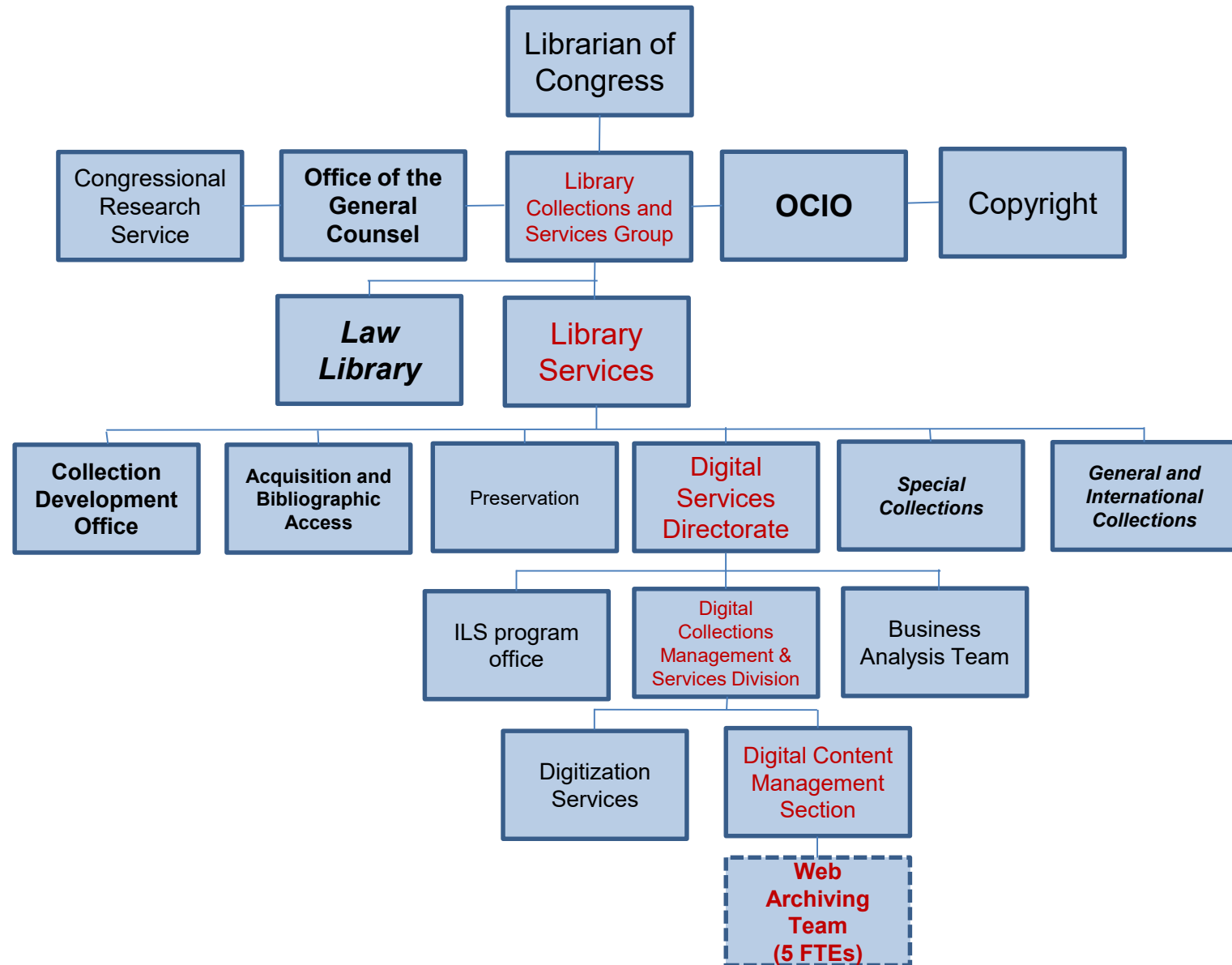
# Our Technical Approach and Tools

- Acquisition
    - Primarily outsourced crawling under contract with Internet Archive
        - IA manages and runs crawls (Heritrix and Brozzler)
        - We do QA before transferring to the Library for preservation and access
    - Limited crawling in-house
        - Heritrix
        - Some experimentation/testing of WebRecorder
    - .gov content from 1996-2001 acquired via backfile purchase from IA
    - Copies of data collected through collaborative efforts (such as End of Term project)
- Access
    - Open Wayback
    - Some testing of Pywb
    - Collections integrated into loc.gov – MODS records are searchable
    - No full-text indexing
    - Some datasets available through https://labs.loc.gov/experiments/webarchive-datasets/
- Storage and Processing
    - LC infrastructure and tools used:
        - Content Transfer Services copy to long term tape storage and a copy for access
        - Participated in recent experimentations with cloud processing of collections

# The LC Approach: Event and Thematic Collections

| | | |
|---|---|---|
| American Music Industry Web Archive | active | 04/26/2018 - ongoing |
| Author Websites Web Archive | active | 04/13/2018 - ongoing |
| Banking Industry in Southeast Asia Web Archi... | active | 05/11/2017 - ongoing |
| Brazil Cordel Literature Web Archive | active | 11/29/2011 - ongoing |
| Brazilian Presidential Election 2018 Web Arc... | active | 04/26/2018 - 01/02/2019 |
| Brozzler Supplemental Crawling | active | 06/13/2018 - ongoing |
| Business in America Web Archive | active | 03/27/2015 - ongoing |
| China-Pakistan Economic Corridor (CPEC) 2018... | active | 02/02/2018 - 12/31/2020 |
| Comics Literature and Criticism Web Archive | active | 01/09/2018 - ongoing |
| Digital Formats Web Archive | active | 08/01/2009 - ongoing |
| East European Government Ministries Web Arch... | active | 04/25/2014 - ongoing |
| Executive Branch Federal Government Web Arch... | active | 07/28/2016 - ongoing |

Library of Congress Collection Policy Statements:
http://www.loc.gov/acq/devpol/cpsstate.html
Web Archiving Supplemental Guidelines:
http://www.loc.gov/acq/devpol/webarchive.pdf

Our web archive collections are typically:

- **Thematic or subject-focused** (e.g., Authors Web Archive, LGBTQ Studies Web Archive, Web Cultures Web Archive)

OR

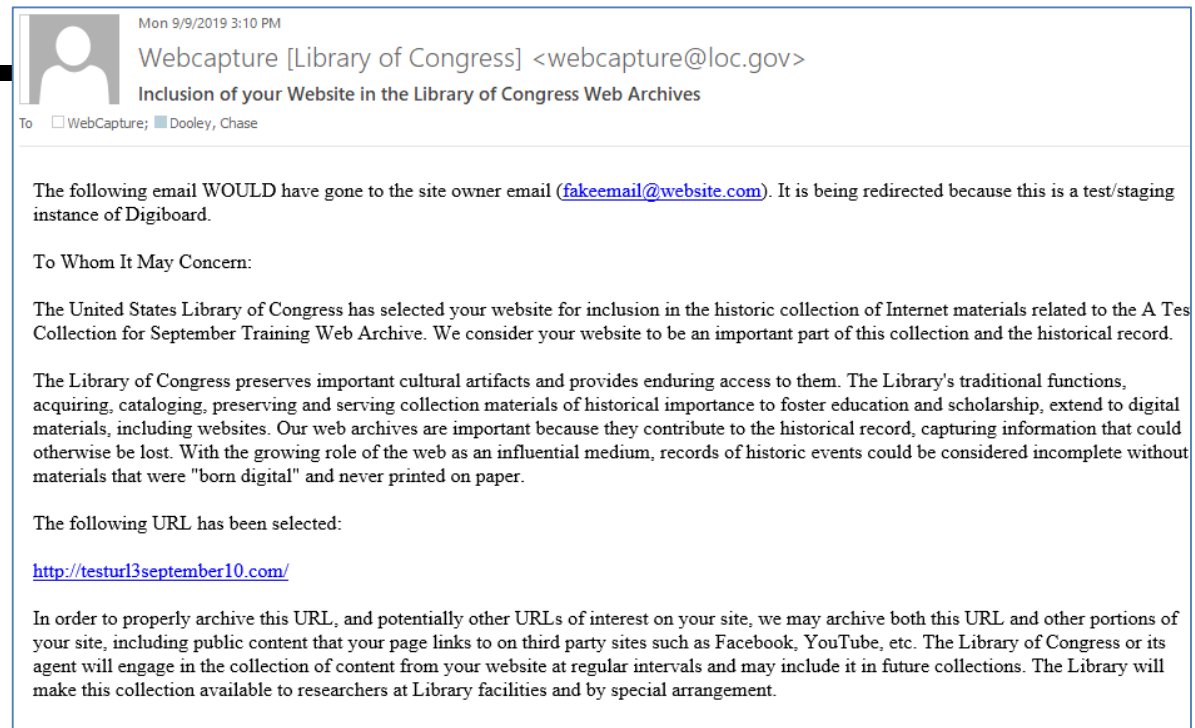- **Event-focused** (e.g.; national or foreign elections)

We manage over **150** collections

**63** active are currently ACTIVE, event-based collections

**9** are "administrative" collections to help us manage our crawls
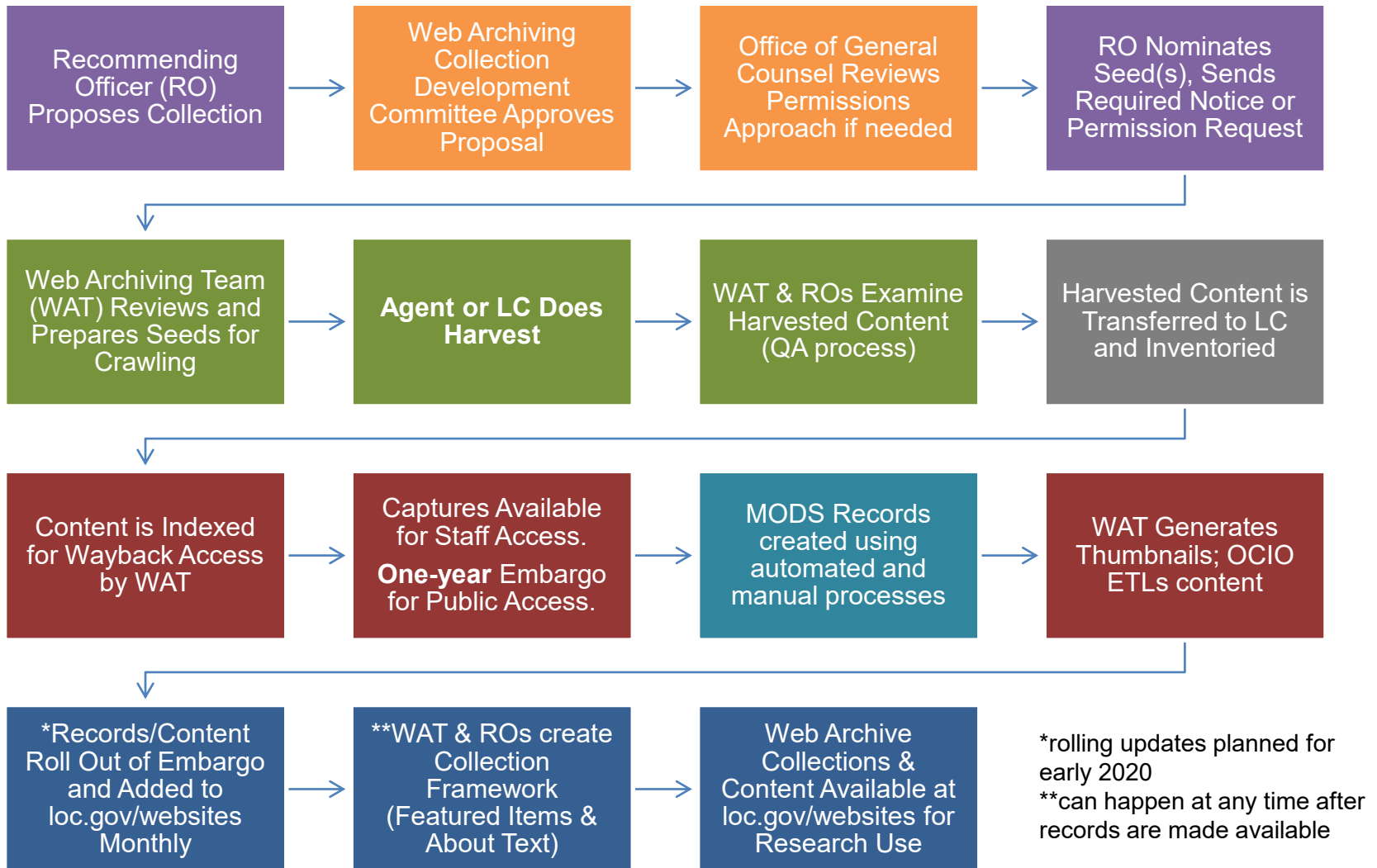
LIBRARY LIBRARY OF CONGRESS

## We follow a permissions-based approach

- Permissions/notices typically have to be sent for anything selected for web archiving
- Permissions are currently based on the **COUNTRY of PUBLICATION** and the **CATEGORY** of site.



Mon 9/9/2019 3:10 PM

Webcapture [Library of Congress] <webcapture@loc.gov>
Inclusion of your Website in the Library of Congress Web Archives

To   ☐ WebCapture; ■ Dooley, Chase

The following email WOULD have gone to the site owner email (fakeemail@website.com). It is being redirected because this is a test/staging instance of Digiboard.

To Whom It May Concern:

The United States Library of Congress has selected your website for inclusion in the historic collection of Internet materials related to the A Tes Collection for September Training Web Archive. We consider your website to be an important part of this collection and the historical record.

The Library of Congress preserves important cultural artifacts and provides enduring access to them. The Library's traditional functions, acquiring, cataloging, preserving and serving collection materials of historical importance to foster education and scholarship, extend to digital materials, including websites. Our web archives are important because they contribute to the historical record, capturing information that could otherwise be lost. With the growing role of the web as an influential medium, records of historic events could be considered incomplete without materials that were "born digital" and never printed on paper.

The following URL has been selected:

http://testurl3september10.com/

In order to properly archive this URL, and potentially other URLs of interest on your site, we may archive both this URL and other portions of your site, including public content that your page links to on third party sites such as Facebook, YouTube, etc. The Library of Congress or its agent will engage in the collection of content from your website at regular intervals and may include it in future collections. The Library will make this collection available to researchers at Library facilities and by special arrangement.

- **Notice/Notice:** We can crawl once notice is sent and make accessible after the embargo period
- **Notice/Permission:** We can crawl once the notice is sent, and make available offsite if the site owner grants permission, otherwise it is made available onsite after the embargo period.
- **Permission/Permission:** We must have explicit permission to crawl and make accessible offsite.
- **No Notice:** US Government sites and those with creative commons

# The LC Web Archiving Process

| | | | |
|---|---|---|---|
| Recommending Officer (RO) Proposes Collection | → Web Archiving Collection Development Committee Approves Proposal | → Office of General Counsel Reviews Permissions Approach if needed | → RO Nominates Seed(s), Sends Required Notice or Permission Request |

| | | | |
|---|---|---|---|
| Web Archiving Team (WAT) Reviews and Prepares Seeds for Crawling | → **Agent or LC Does Harvest** | → WAT & ROs Examine Harvested Content (QA process) | → Harvested Content is Transferred to LC and Inventoried |

| | | | |
|---|---|---|---|
| Content is Indexed for Wayback Access by WAT | → Captures Available for Staff Access. **One-year** Embargo for Public Access. | → MODS Records created using automated and manual processes | → WAT Generates Thumbnails; OCIO ETLs content |

| | | |
|---|---|---|
| *Records/Content Roll Out of Embargo and Added to loc.gov/websites Monthly | → **WAT & ROs create Collection Framework (Featured Items & About Text) | → Web Archive Collections & Content Available at loc.gov/websites for Research Use |

*rolling updates planned for early 2020
**can happen at any time after records are made available

LIBRARY
LIBRARY OF CONGRESS

# Digiboard: Our Web Archiving Workflow Tool

**digiboard**   Nominations   Collections   Review   Search   Admin   Modules ▾                    Abbie Grotke ▾

Home

## News

**Get Help and Learn More**                                                  *Brenda Ford on 04/26/2019 @ 12:26:46 PM*
**Web Archiving Office Hours** - The next Office Hours session will be Tuesday, February 18, 2020 in the Research Orientation Center Jefferson Building, LJ 139B, 11am - Noon

**Getting Started with Web Archiving** - For new-to-web-archiving staff. Next one: TBD.

**Web Archiving Interest Group** - For anyone interested! Next one: February 4, 2020, 11-12, Law Library Multipurpose Room LM-201.

**Digiboard Training Guides** - Search or browse our online help.

**Crawl deadlines for newly added nominations**                               *Abbie Grotke on 11/09/2018 @ 1:26:07 PM*
*New Weekly seeds:* anytime (crawl starts every Tuesday)
*New Monthly, Once, Semi-Yearly, and Yearly seeds:* by the 20th of the month (the crawl starts the first Wednesday of each next month)
*New Quarterly crawl seeds for next quarterly crawl:* by February 20th (for crawl starting March 4)
*New Twice-daily (for use in crawling RSS feeds only):* by the 20th of the month (the crawls happen twice daily, but a new list is provided monthly)

## My Archive

- Actions

  - Manage my Records
  - Browse in PreCrawl by Collection or crawl date
  - Export Seedlist by crawl or Collection

- Active Collections

| Identifier | Title | Status | Crawl Dates | Records | Actions |
|---|---|---|---|---|---|
| yao85y4 | Performing Arts Foundations and Institutes W... | draft | | 0 | ❶ ✎ |
| cj11e1f | Cameroonian Parliamentary Election 2020 Web ... | new | 01/13/2020 - 02/23/2020 | 0 | ❶ ✎ |
| o1vagni | Demonstrations in India 2020 | new | 02/01/2020 - 07/31/2020 | 0 | ❶ ✎ |

**Access: loc.gov/websites/collections/**

- 21,729 web archives available
- Records from over 97 collections available
- 63 collections with contextual (aka framework) material, more on the way!
- All content is embargoed for (at least) one year after capture
- Some content restricted to onsite only; the rest is available from anywhere
- Rolling updates starting soon

## Record search at loc.gov

records describing each archived item are searchable alongside other digital collections at the Library

## URL search at webarchive.loc.gov

displays the archived content (including uncatalogued sites) up through embargo date

# Amount of data in the LC Web Archives
(as of 01/23/20)



Inventory Content Size by Project

Web Archives - Total Size (file_extension)

- WebArchive
- ndnp
- VHP Full Inventory...
- VHP Contractor A...
- Copyright Card Sc...
- AFC Veterans Hist...
- Twitter Archive
- House Streaming ...
- GMD As-Delivered...
- AFC–AIP

**2.098PB**
Sum of byte_count

Web Archives - Bag Instances Total Size - Timeline

Total Content Received (TB)

Over 18 billion documents

LIBRARY
LIBRARY OF CONGRESS

# After 2 PB and 20 Years: A Few Lessons Learned

## Automate and Reuse Data

- Do things manually for awhile to help understand the process, then figure out where you can automate
- Work with all the data you have
- Automate permissions process as much as possible. Develop permissions letter templates.
- Spreadsheets to manage tasks like QA and permissions are great but only take you so far
- Crawl reports have a ton of data that is just waiting to be used better for QA at scale

## Collaborate

- Internally - get to know all of the people who might be able to help you along the way
- Externally – DPC and IIPC are a great community, everyone is willing to help!
- Collaborative collection efforts are great when your own institutional policies make it difficult to collect rapidly or without permission
  - IIPC projects
  - IA-led efforts
  - Collaborations among other partner institutions
- Everyone has something to contribute

## Train and Retrain

- Many hires still don't have extensive web archiving skills, so training off staff takes time and investment (mostly on the job)
- Attend IIPC and other conferences (Archives Unleashed events)
  - Learn new skills –python/scripting have been critical in recent years
- Train staff as needed, and think about refresher training ALL the time
  - For Recommending officers: one-on-one, Office hours, Interest Group discussion, classroom training, etc.

## Document Everything

- Documenting decisions is key: put all policy decisions in writing
- Having good, detailed FAQ on our website geared toward site owners
- Develop good training materials for infrequent web archiving staff and for volunteers/interns brought in to help
- Develop a collection proposal template that can be used to help shape a project before it gets started
- Document workflows and processes
- Log and document all decisions for actions taken on seeds

## Accept "Good Enough"

- Defined a minimal amount of data that will allow research use but not require too much human effort to catalog
- Use experts/humans to enhance descriptive records
- Don't get too far behind, catching up on backlogs can take years and decades in some cases
- The faster we got our collections out, the more engaged the subject experts became
- We have to accept that our crawls won't be 100% perfect
- Be okay with not QAing everything

## Reinvent

- Be flexible
- Always look for ways to reinvent processes or workflows to improve things
- Test new tools as you can, but know that you might be working with the old tools for longer than you'd like (until new ones improve or scale up)
- Help the international community in efforts that will help us all

# Thank you !

Abbie Grotke
[abgr@loc.gov](mailto:abgr@loc.gov)
@agrotke
webcapture@loc.gov