# Capturing our State: The California State Government Web Archive

---

**Capturing our State:**

The California State Government Web Archive

An Infopeople Hosted Webinar
April 24, 2017

Archive of the California Government Domain, CA.gov

This work is licensed under the Creative Commons Attribution-NonCommercial 4.0 International License. To view a copy of this license, visit http://creativecommons.org/licenses/by-nc/4.0/.

---

**Credits**

Kris Kasianovitz, Government Information Librarian
Stanford University Libraries
krisk11@stanford.edu

Shari Laster, Government Information Librarian
UC Santa Barbara Library
slaster@ucsb.edu

Julie Lefevre, Digital Services Librarian
Institute of Governmental Studies Library
UC Berkeley
jlefevre@library.berkeley.edu

Lucia Orlando, Arts, Humanities, Social Sciences
and Government Information Librarian
UC Santa Cruz Library
luciao@ucsc.edu

---

**LSTA Grant Funded Project**
**(thank you & disclaimer)**

INSTITUTE of Museum and Library SERVICES

California STATE LIBRARY

This project and presentation at Best Practices Exchange was supported in whole or in part by the U.S. Institute of Museum and Library Services under the provisions of the Library Services and Technology Act, administered in California by the State Librarian.

The opinions expressed herein do not necessarily reflect the position or policy of the U.S. Institute of Museum and Library Services or the California State Library, and no official endorsement by the U.S. Institute of Museum and Library Services or the California State Library should be inferred.

---

Infopeople, a grant project of the Califa Group, is supported in part by the U.S. Institute of Museum and Library Services under the provisions of the Library Services and Technology Act, administered in California by the State Librarian.

# Capturing our State: The California State Government Web Archive

## Agenda

- Why CA.gov?
- Introduction to web archiving
- CA.gov web archive project overview
- Building your own archive

## Outcomes

- Understand the basics of web archiving using Archive-It
- Be familiar with the scope and content of the CA.gov Web Archive
- Be able to search the CA.gov Web Archive and the WayBack Machine for state government information
- Know who to contact with questions about the project or suggest sites for inclusion in the archive
- Have ideas for web harvesting projects for your local government agencies

# Capturing our State: The California State Government Web Archive

---

**CA.gov archive:**
**a collaborative effort!**

- Stephen Abrams, CDL
- Julie Lefevre, UCB Institute for Governmental Studies Library
- Kris Kasianovitz, Stanford Library
- Shari Laster, UCSB Library
- Lucia Orlando, UCSC Library
- Bill Riddle, California State Library
- Mike Smith, UCSD Library
- Joseph Yue, UCLA Library
- Michael McNeil, California State Archives

---

**Web archiving**

- Goals & principles
- Internet Archive & Archive-It
- Web crawlers

---

**Web archiving basics**

Web archiving is the process of:

collecting portions of the World Wide Web,

preserving the collections in an archival format, and then

serving the content for access and use.

Capture is usually accomplished with a *crawler* like Heritrix, which saves web pages and also follows links to additional pages.

Replay is the means of serving content for access and use: the Internet Archive's Wayback Machine, Archive-It Collections site, and Memento are good examples.
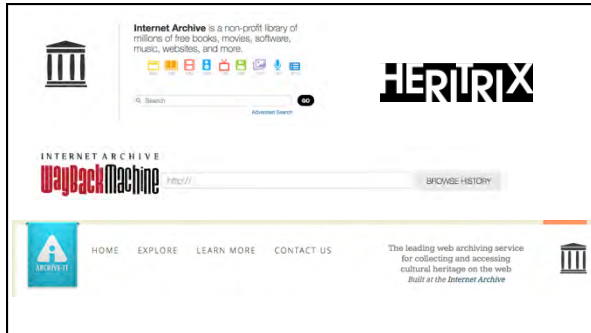
---

# Capturing our State: The California State Government Web Archive





## Crawling the web

A web crawler is a software script that fetches, analyzes, and files information from web servers at many times the speed of a human

The crawler starts with a list of URLs to visit, called the seeds

A seed can be a top level domain, like http://www.ca.gov
or a sub-directory below the top-level domain, like
http://www.ca.gov/Agencies/Secretary-of-State/

"Scoping" rules constrain the web crawl to ensure we collect what we want, and nothing more

# Capturing our State: The California State Government Web Archive

## Project overview

- Scope of project
- Work involved
- Using our content

## CA.gov archive: background

https://archive-it.org/collections/5763

- 2007 Collection established in California Digital Library's Web Archiving Service (WAS)
- 2010–2012 Minimal development
- 2013 Cooperative project of the University of California/Stanford University Government Information Librarians, with representation from California State Library
- 2015 Content migrated to Archive-It (4.2Tb of data)
- 2016 LSTA grant application
  - Grant goals:
    - comprehensive capture of the State of California web domain
    - develop quality assurance workflows
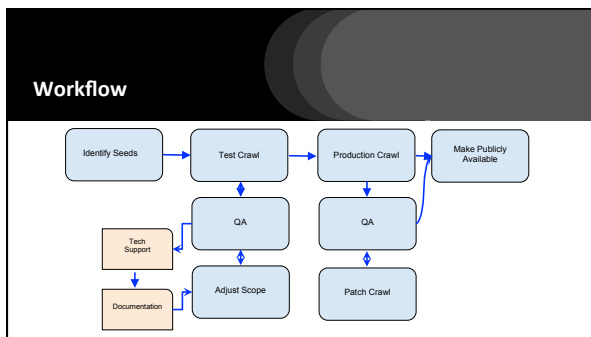    - determine a sustainable funding model

## Workflow



Identify Seeds → Test Crawl → Production Crawl → Make Publicly Available

Test Crawl → QA
QA → Tech Support
Tech Support → Documentation
QA → Adjust Scope
Production Crawl → QA
QA → Patch Crawl

# Capturing our State: The California State Government Web Archive

## Gathering our seeds

A different kind of collection development:

Identified and reviewed URLs from official state sources.

Not all seeds can be captured as entered: "Site not found" errors, robots.txt can prevent capture.

Many official state sites are missing: we add content as we find it.

Following an initial test crawl, we reviewed the data, corrected some errors, and completed our first snapshot (1.9Tb) in December 2016.

Future crawls using Archive-It will de-duplicate against the content we have collected.

---

## Managing our collection

Distributing & shifting project work is a learning process!

Check seed reports, contact website administrators to request permission to crawl blocked content

Quality assurance on captures, conduct patch crawls

Troubleshoot problems, request assistance from Archive-It to solve issues

Communicating process and progress within the group & broader communities

---

## Our collection

This year so far:

413 active seeds (URLs)

7,598,972 documents -- 2.7Tb

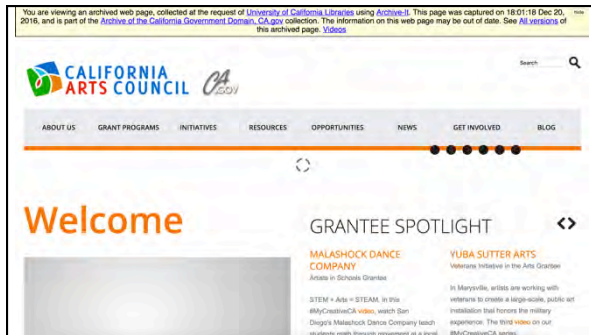This is a big collection by Archive-It standards (though not for the Internet Archive!).

---

# Capturing our State: The California State Government Web Archive

# Capturing our State: The California State Government Web Archive

# Capturing our State: The California State Government Web Archive

**Build your own collections!**

- Examples
- Using Archive-It
- Considerations

_____
_____
_____
_____
_____
_____
_____

Examples of Local Government Information Web Archives using Archive-It

_____
_____
_____
_____
_____

**Getting started**

Advantages to subscribing to Archive-It:

The Internet Archive handles the technical aspects of capture, preservation, and playback

You control what's collected, including frequency & depth of capture

Content can be exported to include in a digital archive

Community of Archive-It partners working on projects

You can contribute in other ways, too:
http://blog.archive.org/2017/01/25/see-something-save-something/

_____
_____
_____
_____
_____
_____
_____

# Capturing our State: The California State Government Web Archive

## Archive… with a plan!

Determine the scope of your collection based on your user communities.

Once you have set up & tested your crawl, schedule it to run on regular intervals.

Keep an eye on your space: Archive-It de-duplicates content, but some projects may require waiting until the next subscription year to complete an initial capture.

Be ready for the unexpected!

---

## Get in Touch

Have questions about the CA.gov web archive?

Want to suggest site(s) for inclusion?

Let us know!

email: cagovarchive@lists.berkeley.edu

---

## THANK YOU!

### Discussion and Questions

---

# Capturing our State: The California State Government Web Archive

Infopeople

Infopeople is dedicated to bringing you the best in practical library training and improving information access for the public by improving the skills of library workers. Infopeople, a grant project of the Califa Group, is supported in part by the Institute of Museum and Library Services under the provisions of the Library Services and Technology Act administered in California by the State Librarian. This material is covered by Creative Commons 4.0 Non-commercial Share Alike license.  Any use of this material should credit the funding source.