# Web Harvesting Update

March 31, 2008
Laurie Hall
Director, Library Technical Information Services
Robin Haun-Mohamed
Director, Collection Management & Preservation
U.S. Government Printing Office

---

# Web Harvesting at GPO

- Harvesting digital Federal publications since the late 1990's
- Harvesting methods
  - Manual
  - Semi-manual
- Automated harvesting tools are part of Release 2 of FDsys

## Assumptions

- GPO will continue to participate in Web harvesting efforts to obtain in-scope material for the FDLP and the Cataloging and Indexing Program as required under 44 USC.
- GPO is bound by Congressional appropriations for the Salaries & Expenses funding requirements for the FDLP and C&I Programs.
- All materials identified for inclusion in the FDLP must be brought under bibliographic control as directed by the C&I Program.

*our strategic vision in progress*

## Assumptions

- GPO does not have the authority to either give funding or gifts or to receive them.  All partnerships must represent a contribution of an equal exchange between all parties.
- Automated Web harvesting initiatives will become systematic as part of Release 2 of FDsys.
- Materials harvested under the EPA Pilot Project are being made available as staff time and processing permit.  Completion of the processing of this material will necessarily require an automated metadata extraction process that does not yet exist.

*our strategic vision in progress*

## Manual Harvesting Efforts

- Capture of known digital publications through manual identification and the saving of all associated publication files
- Monitor agency Web sites for new or updated publications
- Acquire fugitive publications through notifications in LostDocs
- Focus is on PDF files

our strategic vision in progress

## Semi-Manual Harvesting Efforts

- Use a software tool to schedule the content capture and re-harvesting of known content at known Web sites
- Used to harvest serial issues because can schedule re-harvests to acquire new issues
- Used to acquire publications in non-PDF formats

our strategic vision in progress

## Automated Web Harvesting Pilot

- Conducted in 2006 with approval of the EPA
- Two vendors crawled the EPA Web sites
- Used rules to determine if digital material was in scope of the FDLP and C&I Program
- Acquired over 200,000 files
- 14-16% of the results are not within scope
- At least 25% of the results were only partially harvested

*our strategic vision in progress*

## Type of Files Acquired

- Random sample of 1000 publications
  - 62% database results
  - 23% monographs
  - 9% Web pages
  - 3% serial issues

*our strategic vision in progress*

U.S. GOVERNMENT PRINTING OFFICE I KEEPING AMERICA INFORMED

GPO

## Processing Issues

- Must focus on Government information products, or publications as identified under Title 44
- 25% of files are out of scope
- 14-16% of files are incomplete
- Staffing limitations

our strategic vision in progress



U.S. GOVERNMENT PRINTING OFFICE I KEEPING AMERICA INFORMED

GPO

## Sample of 300 Publications

- To estimate the amount of time and the staffing implications to process all the results
- Tested two mechanisms for making the publications accessible
  - Brief bibliographic records in the CGP
  - Browse table on GPO Access
- Reviewing comments received from the depository community

our strategic vision in progress

## Sample of 300 Results

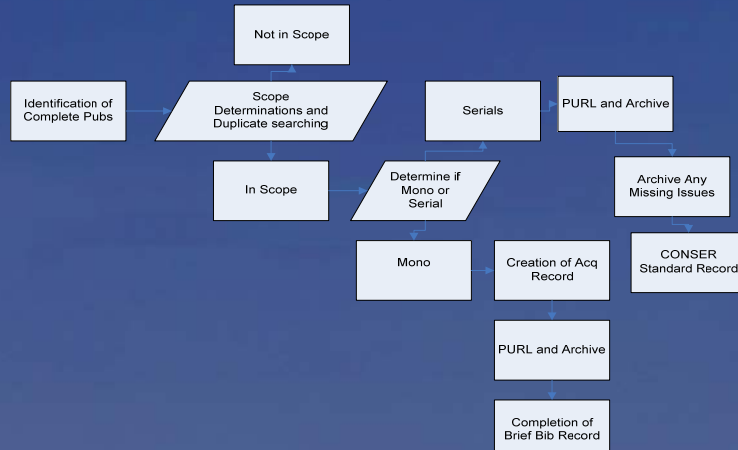| Determination | Percentage of the Sample |
|---|---|
| Already Cataloged as an electronic title | 18.5% |
| Previously distributed in tangible format | 3% |
| Not within scope | 2% |
| New publication | 62% |

our strategic vision in progress

## Processing Times for the Sample

| Processing Step | Average Time |
|---|---|
| Identification of a Complete Publication | 2 minutes |
| Scope Determination and Search for Duplicates | 17 minutes |
| Creation of Brief Bibliographic Record | 30 minutes |
| Creation of CONSER Standard Record | 2 hours 30 minutes |
| Add PURL to Publications Distributed in Tangible Format | 7.5 minutes |
| Creation of Browse Table | 4 hours total to create the entire browse table |

our strategic vision in progress

# EPA Processing Workflow

```
                          ┌──────────────┐
                          │ Not in Scope │
                          └──────────────┘
                                 ↑
┌─────────────────┐   ┌────────────────────┐   ┌──────────┐   ┌──────────────────┐
│ Identification of│──▶│ Scope              │   │ Serials  │   │ PURL and Archive │
│ Complete Pubs    │   │ Determinations and │   └──────────┘   └──────────────────┘
└─────────────────┘   │ Duplicate searching│                          │
                      └────────────────────┘                          ▼
                                 │          ┌──────────────┐   ┌──────────────────┐
                          ┌──────────┐      │ Determine if │   │ Archive Any      │
                          │ In Scope │─────▶│ Mono or      │   │ Missing Issues   │
                          └──────────┘      │ Serial       │   └──────────────────┘
                                            └──────────────┘           │
                                                    │                  ▼
                                          ┌──────┐  ┌──────────────┐  ┌──────────────────┐
                                          │ Mono │─▶│ Creation of Acq│ │ CONSER           │
                                          └──────┘  │ Record        │ │ Standard Record  │
                                                    └──────────────┘  └──────────────────┘
                                                           │
                                                           ▼
                                                    ┌──────────────────┐
                                                    │ PURL and Archive │
                                                    └──────────────────┘
                                                           │
                                                           ▼
                                                    ┌──────────────────┐
                                                    │ Completion of    │
                                                    │ Brief Bib Record │
                                                    └──────────────────┘
```

---

# Search for Complete Publications

- Vendors organized the results differently
- Each file harvested must be examined and compared to the live copy, if available
- Complete publications must be located among:

Database Results



Web Pages

# Partially Harvested Publications

# Does the Publication Belong?

- Determine if publication is within scope of the FDLP and/or the C&I Program
- Some of the questions considered:
  - Is it published by a US Government agency?
  - Is the information in it covered by copyright?
  - Who is the copyright holder?
  - Is the primary source of funding for it government money?
  - Does it contain data that may violate a citizen's privacy?

## Is the Publication a Return Visitor?

- Determine if it has previously been distributed
  - Search the legacy database
  - Search the CGP
- No further action taken if publication has an EL only record in the CGP
- All publications that have no records or were distributed in tangible format move forward

## A Fork in the Path

- Different processes for monographs, serial issues, and special materials
- Monographs: brief bibliographic records created by non-librarians; created directly in the ILS
- Serial issues: CONSER standard records created by cataloging librarians; created in OCLC

**Special Materials**

- Special materials include:
  - News releases
  - Transmittals
  - Forms
  - Announcements
- Included in List of Special Materials in the Monthly Catalog
- Investigating methods to provide bibliographic access to this material



**Special Material**

# Monograph: 1st Step

- Begin brief bibliographic record in ILS
- Includes basic bibliographic information and the SuDoc stem
- Not accessible in CGP
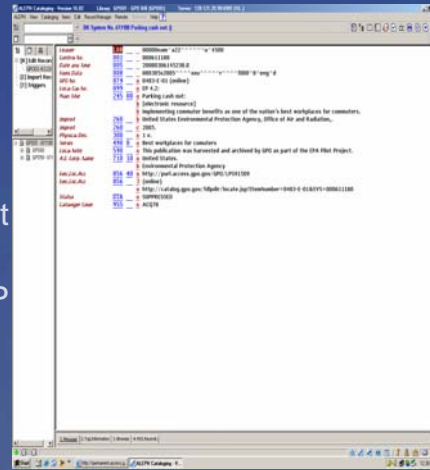- Created by staff in Content Acquisitions

# Monograph: 2nd Step

- Archive publication on the Permanent Server
- Create PURL to issuing agency's site or the archived copy
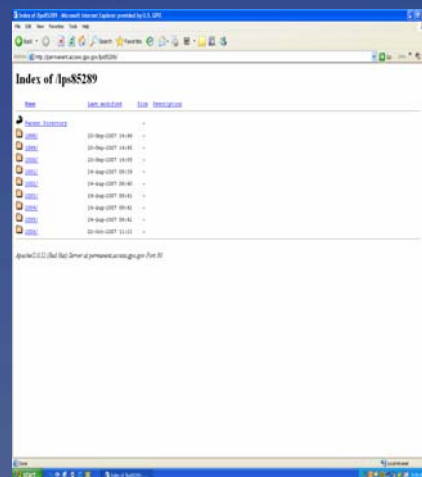- Add PURL to bibliographic record

# Monograph: Final Step

- Complete brief bibliographic record
- Validate SuDoc and item numbers
- Record indicates is part of EPA Pilot Project
- Now accessible in CGP
- Finalized by staff in Library Technical Information Support

*our strategic vision in progress*

# Serial Issue: 1st Step

- Archive issue on the Permanent Server
- Harvest and archive additional issues
- Create PURL to issuing agency's site or the archived copy

*our strategic vision in progress*

# Serial Issue: 2nd Step

- Issuing agency contacted for information, if needed
  - Is the serial still being published
  - Frequency
- Can take time to get a response
- New SuDoc class and item number created, if needed

our strategic vision in progress

---

# Serial Issue: Final Step

- Create CONSER standard record in OCLC
- Load record into the ILS
- Record in ILS indicates is part of the EPA Pilot Project
- Now accessible in CGP



our strategic vision in progress

## Related Projects: Metadata Extraction

- 2 year project with Old Dominion University to use automated metadata extraction software tools to create metadata for digital publications
- Using 1,000 publications from the EPA Pilot Project
- Currently developing software rules and designing templates

## Related Projects: Special Material

- Demonstration project to examine depository participation in harvesting activities
- Assist GPO in creating brief bibliographic records for items in the special materials category
- Opportunity for 5 depository librarians
- More information on FDLP-L as plans develop

## Related Project: Partially Harvested Publications

- Demonstration project to examine depository participation in harvesting activities
- Locate and harvest all the parts of partially harvested publications from EPA Pilot Project
- Complete the harvest of 150 publications
- 3 month project

## Related Project: Partially Harvested Publications

- Opportunity for 5 depository librarians to participate
- Must have time to devote to project from June to the beginning of September
- Must have FTP ability so GPO can retrieve harvested publications
- Call for volunteers will be posted to FDLP-L after Council

# Web Harvesting

## Web Harvesting Update for the Depository Library Council

**BACKGROUND**
GPO has been harvesting digital Federal publications since the late 1990's. These publications are acquired using two harvesting techniques: manual and semi-manual. Manual harvesting is the capture of known digital publications through manual identification and the saving of all associated publication files, while semi-manual uses a software tool to schedule the content capture and re-harvesting of known content at known Web sites.

In 2006, GPO, with the approval of the U.S. Environmental Protection Agency (EPA), undertook an automated pilot project to harvest materials from EPA Web sites. Using two different vendors, over 200,000 files were acquired. Subsequent action and review of this material has focused on ways to identify efficient methods to bring the files that are within scope of the FDLP and the Cataloging and Indexing (C&I) Program under bibliographic control.

**SAMPLE PUBLICATIONS FROM THE RESULTS OF THE EPA PILOT PROJECT**
Library Services and Content Management (LSCM) staff processed a sample of 300 publications harvested during the EPA Pilot Project. The purpose of working through this sample was to determine workflow and staffing implications, as well as to estimate the amount of time that would be required to process all of the publications acquired during the EPA Pilot Project.

Two mechanisms were used for making accessible the publications found to be within scope of the FDLP. The depository community was asked to provide feedback on both. The majority of publications were made accessible through bibliographic records in the Catalog of U.S. Government Publications (CGP). Monographs were cataloged using the new brief bibliographic record format, while serials were cataloged following the CONSER standard record format. At the request of the Depository Library Council (DLC), LSCM also explored a mechanism that enables public access to Web harvested content while these publications are in the queue for brief bibliographic records, and posted a small portion of the sample to *GPO Access* using a browse table. Publications made accessible through this mechanism will be cataloged in the CGP in the future.

Of the 300 publications identified for the sample, 62% had not previously been distributed through the FDLP. Based on the average processing times, it will take 49 minutes to process a single monograph from identification in the vendors' results to the finalization of the brief bibliographic record. Average processing time for a serial title is two hours and 49 minutes, from identification through the creation of the CONSER standard record. The complete analysis of the sample is available at **<http://www.access.gpo.gov/su_docs/fdlp/harvesting/index.html>**.

ONGOING HARVESTING ACTIVITIES

GPO continues to take an active role in Web harvesting efforts. As a member of CENDI's Web Harvesting Task Group, GPO is participating in discussions with other Federal agencies that will lead to the development of a white paper on Web harvesting best practices. Other agencies participating in the Task Group include NASA, the EPA, NTIS, and the Library of Congress. CENDI is an interagency working group of senior scientific and technical information managers from thirteen Federal agencies. Any white paper or other products resulting from this interagency effort will be shared with DLC and the depository library community. This information will also be used to further develop requirements as appropriate for the Web harvesting components of the FDsys.

Manual Web harvesting efforts also continue on a daily basis. Staff in Library Technical Information Services (LTIS) are responsible for identifying, classifying, and archiving new digital publications as part of their daily duties. These publications are included in the monthly New Electronic Titles report. While the bulk of the harvesting is performed manually, GPO staff continue to work with a semi-automated harvester to acquire issues of serials and publications in non-PDF formats.

In addition, GPO staff have begun another project to provide more access to the files acquired during the EPA Pilot Project. Following a workflow established based on the results of the sample of 300 publications, staff will process 500 files in the coming months. Many of these files will either be incomplete or not within scope of the FDLP or the C&I Program, so the depository community should not expect to locate 500 additional bibliographic records in the CGP after the completion of this project.

HARVESTING ISSUES
- Most of the material acquired during the EPA Pilot Project is still available from EPA Web sites and has not been lost. Some of the material retrieved was not in-scope, not a complete publication, or both. GPO is developing a demonstration project with libraries to assist in identifying where the complete publications can be located.
- The sheer number of files identified and retrieved during the EPA Pilot Project is overwhelming. In addition to the development of automated metadata extraction tools to assist with bibliographic control over this material, GPO is developing a demonstration project for libraries to assist in the creation of brief cataloging records for a subset of the results from the EPA Pilot Project.
- Many of the files returned from the EPA crawls were parts of a database. GPO continues to focus Web harvesting efforts on Federal Government information products or publications as identified under Title 44. Other agencies may have a different role and approach, including the downloading of entire or partial Web sites.

DEPOSITORY PARTICIPATION IN HARVESTING EFFORTS

The Federal depository library community has expressed interest in participating in LSCM's efforts to process these publications. In order to determine how best to utilize the depository libraries' offer of assistance, LSCM will undertake a demonstration project

with ten depositories. The volunteer depositories will be asked to assist LSCM either in locating the missing sections of partially harvested publications or by creating brief bibliographic records for publications in the Special Materials category, which includes press releases, transmittals, notices, and forms. Additional information about these demonstration projects will be shared via FDLP-L and the FDLP Desktop.

## RELATED PROJECTS
One thousand monographs within scope of the FDLP have been identified from the EPA Pilot Project for inclusion in the Automated Metadata Extraction Project. This two-year project with the Defense Technical Information Service (DTIC) and Old Dominion University (ODU) will use automated metadata extraction software tools to create metadata for groups of electronic publications in GPO's electronic collection. GPO expects to receive the results near the end of the project.

## ASSUMPTIONS
- GPO will continue to participate in Web harvesting efforts to obtain in-scope material for the FDLP and the C&I Program as required under 44 USC.
- GPO is bound by the Congressional appropriations for the Salaries & Expenses funding requirements for the FDLP and C&I Programs.
- All materials identified for inclusion in the FDLP must be brought under bibliographic control as directed under the C&I Program.
- GPO does not have gift authority to either give funding or gifts or to receive them. All partnerships must represent a contribution of an equal exchange between all parties.
- Automated Web harvesting initiatives will become systematic as part of release 2 of the FDsys.
- Materials harvested under the EPA Pilot Project are being made available as staff time and processing permit. Completion of the processing of this material will necessarily require an automated metadata extraction process that does not yet exist.