# Handling a Digital Backlog and Analyzing Content in Archivematica

August 7, 2012
Research Forum
Society of American Archivists Annual Conference, Beyond Borders
San Diego, California

Courtney C. Mumma, MAS/MLIS, Systems Analyst and Archivematica Community Manager

archivematica

# Overview

- Action research in the Archivematica project

- Digital backlog

- Processing digital objects

- Managing transfer backlogs
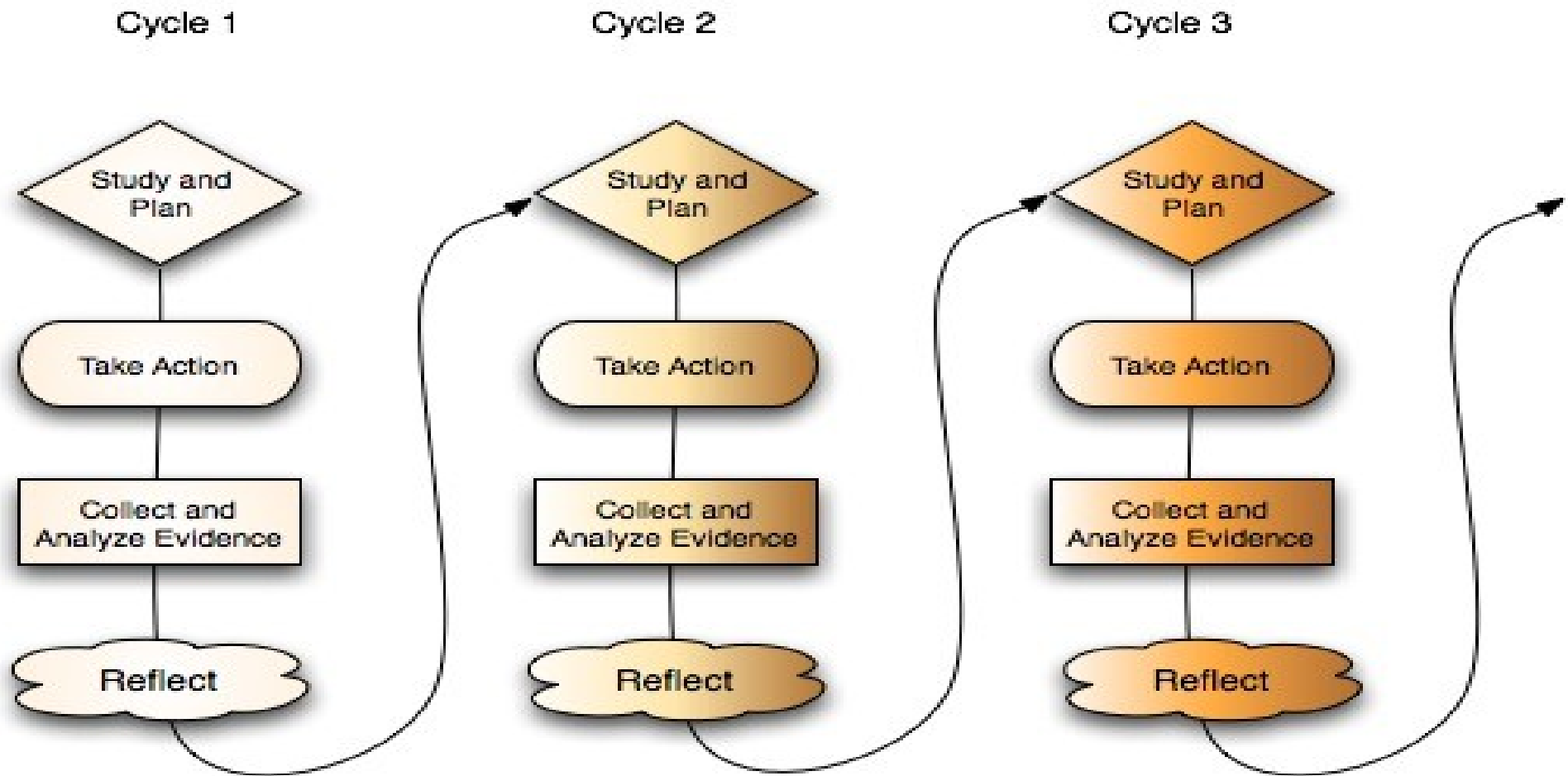
- Creating SIPs in Archivematica

# Action research

- To quote action research's instigator Kurt Lewin: "if you want truly to understand something, try to change it". This  kind of work is not simply about changing, but also improving an environment. As John Elliott says, action research is "the study of a social situation with a view to improving the quality of action within it"

(Elliott, J. (1991). Action research for educational change. Buckingham: Open University Press, p. 69).    http://en.wikiversity.org/wiki/Action_research

# Goals of action research

- The improvement of professional practice through continual learning and progressive problem solving;

- A deep understanding of practice and the development of a well specified theory of action;

- An improvement in the community in which one's practice is embedded through participatory research.

Carr, W., & Kemmis, S. (1986) Becoming critical. Lewes: Falmer Press

Cycle 1      Cycle 2      Cycle 3

Study and Plan → Take Action → Collect and Analyze Evidence → Reflect

Progressive Problem Solving with Action Research

Carr, W., & Kemmis, S. (1986) Becoming critical. Lewes: Falmer Press

# Agile development

Individuals and interactions over processes and tools

Working software over comprehensive documentation

Customer collaboration over contract negotiation

Responding to change over following a plan

Beck, Kent; et al. (2001). "Manifesto for Agile Software Development". Agile Alliance. Retrieved 14 June 2010.

# Six criteria for action research

Roles: Clarify the roles of researchers (Artefactual staff) and practioners (our clients)

Documentation: Explain the data collection approach and how data quality is managed (Observing and documenting local workflows and requirements)

Control: Explain the control measures (standards and best practices)

Usefulness: Establish the usefulness of the findings in the problem situation (assess community response)

Frameworks: Relate the action taken as well as the findings to frameworks to     support the study (build up Archivematica tools and processes)

Transferability: Explicate conditions for transfer of findings to other situations (commonalities make up default workflows in Archivematica)

Peter Axel Nielsen. "IS Action Research and Its Criteria" in Information Systems Action Research: An Applied View of Emerging Concepts (2007), p. 366

# Backlogs are inevitable

- Most of us have backlogs

- Most of us lack resources (tools, staff, funding) to prioritize and process those backlogs

# Minimum data about the backlog

- Verify transfer compliance

- Create transfer UUID and assign  file UUIDs to objects

- Verify metadata directory checksums (verifies any checksums included with the transfer)

- Assign checksums to objects

- Generate METS.xml document

- Extract packages

- Sanitize object's file and directory names

- Scan for viruses and malware

- Characterize and extract metadata

- 0.9 added transfer indexing, 1.0 accessioning as PREMIS event

# Accession metadata

* PREMIS Event = Accession

- &lt;event&gt;
- &lt;eventIdentifier&gt;
- &lt;eventIdentifierType&gt;UUID&lt;/eventIdentifierType&gt;
- &lt;eventIdentifierValue&gt;35cbe00d-d661-4174-b11a-e203f5608008&lt;/eventIdentifierValue&gt;
- &lt;/eventIdentifier&gt;
- &lt;eventType&gt;accession&lt;/eventType&gt;
- &lt;eventDateTime&gt;2012-03-14&lt;/eventDateTime&gt;
- &lt;eventDetail&gt;accession#2012-029&lt;/eventDetail&gt;
- &lt;eventOutcomeInformation&gt;
- &lt;eventOutcome&gt;&lt;/eventOutcome&gt;
- &lt;eventOutcomeDetail&gt;
- &lt;eventOutcomeDetailNote&gt;&lt;/eventOutcomeDetailNote&gt;
- &lt;/eventOutcomeDetail&gt;
- &lt;/eventOutcomeInformation&gt;
- &lt;linkingAgentIdentifier&gt;
- &lt;linkingAgentIdentifierType&gt;archivist&lt;/linkingAgentIdentifierType&gt;
- &lt;linkingAgentIdentifierValue&gt;Courtney Mumma&lt;/linkingAgentIdentifierValue&gt;
- &lt;/linkingAgentIdentifier&gt;
- &lt;/event&gt;

# 0.9 Transfer workflow

```
●  →  Allow MCP access to      →  Select preconfigured    →  Enter transfer name   →  Enter accession number
       media or storage            transfer type (generic,
       where transfer is located   dspace, bagit, etc.)

       Browse to transfer       →  Browse to all          →  Start transfer        →  Create structured
       and select                  submission                                          transfer folder
                                    documentation
                                    and select

       Assign transfer UUID     →  Assign file UUIDs       →  Quarantine            →  Log directory structure
                                    and checksums

       Extract packages         →  Assign file UUIDs       →  Verify checksums      →  Scan for viruses
                                    and checksums to           included in transfer     (generate log)
                                    extracted files

                                    Reject transfer or
                                    remove infected files
                                    and continue processing

                                    Go to next page
```

```
○                    ○
│                    │
▼                    ▼
┌──────────────┐    ┌──────────────────┐   ┌──────────────────┐   ┌──────────────────┐
│              │    │ Remove infected  │   │                  │   │    Move to       │
│Reject transfer│   │files and continue│──▶│ Characterize and │──▶│ Transfer Backlog │
│              │    │   processing     │   │ extract metadata │   │    storage       │
└──────────────┘    └──────────────────┘   └──────────────────┘   └──────────────────┘
       │                    │
       ▼              ┌──────────────┐
┌──────────────┐      │Generate log of│
│  Repository  │      │removed files  │
│   repairs    │      └──────────────┘
│Infected files│
└──────────────┘                              ◆ Reject transfer,
       │                                    Continue processing as
       ▼                                      one SIP or Create SIP ◆
┌──────────────┐
│  Resubmit    │         ┌──────────────┐      ┌──────────────────┐      ┌──────────────┐
│  transfer    │         │Reject transfer│      │Continue processing│      │  Create SIP  │
└──────────────┘         │              │      │   as one SIP     │      │              │
       │                 └──────────────┘      └──────────────────┘      └──────────────┘
       ▼                        │                      │                        │
      ◉                  ┌──────────────┐      ┌──────────────────┐      ┌──────────────────┐
                         │Ask user if they│      │Go to Ingest –    │      │Index (see Transfer│
                         │are sure about  │      │   Process SIP    │      │Indexing          │
                         │rejecting       │      └──────────────────┘      │requirements)     │
                         │transfer        │              │                 └──────────────────┘
                         └──────────────┘              ▼                        │
                                │                      ◉                 ┌──────────────────┐
                         ┌──────────────┐                               │Go to Ingest –    │
                         │User confirms │                               │   Create SIP     │
                         │reject transfer│                               └──────────────────┘
                         └──────────────┘                                        │
                                │                                                ▼
                         ┌──────────────┐                                       ◉
                         │Delete transfer│
                         │from Transfer  │
                         │Backlog storage│
                         └──────────────┘
                                │
                                ▼
                               ●
```
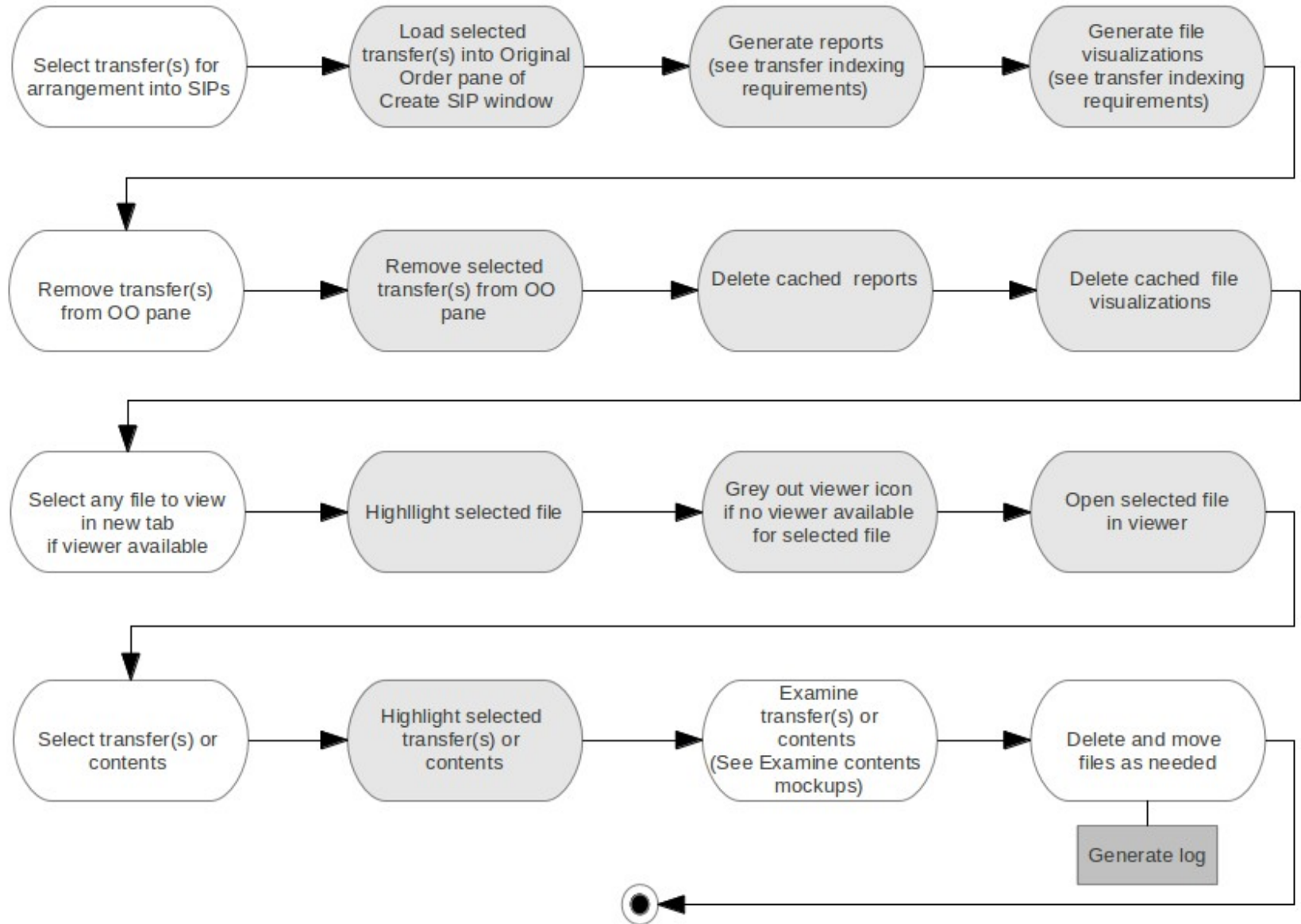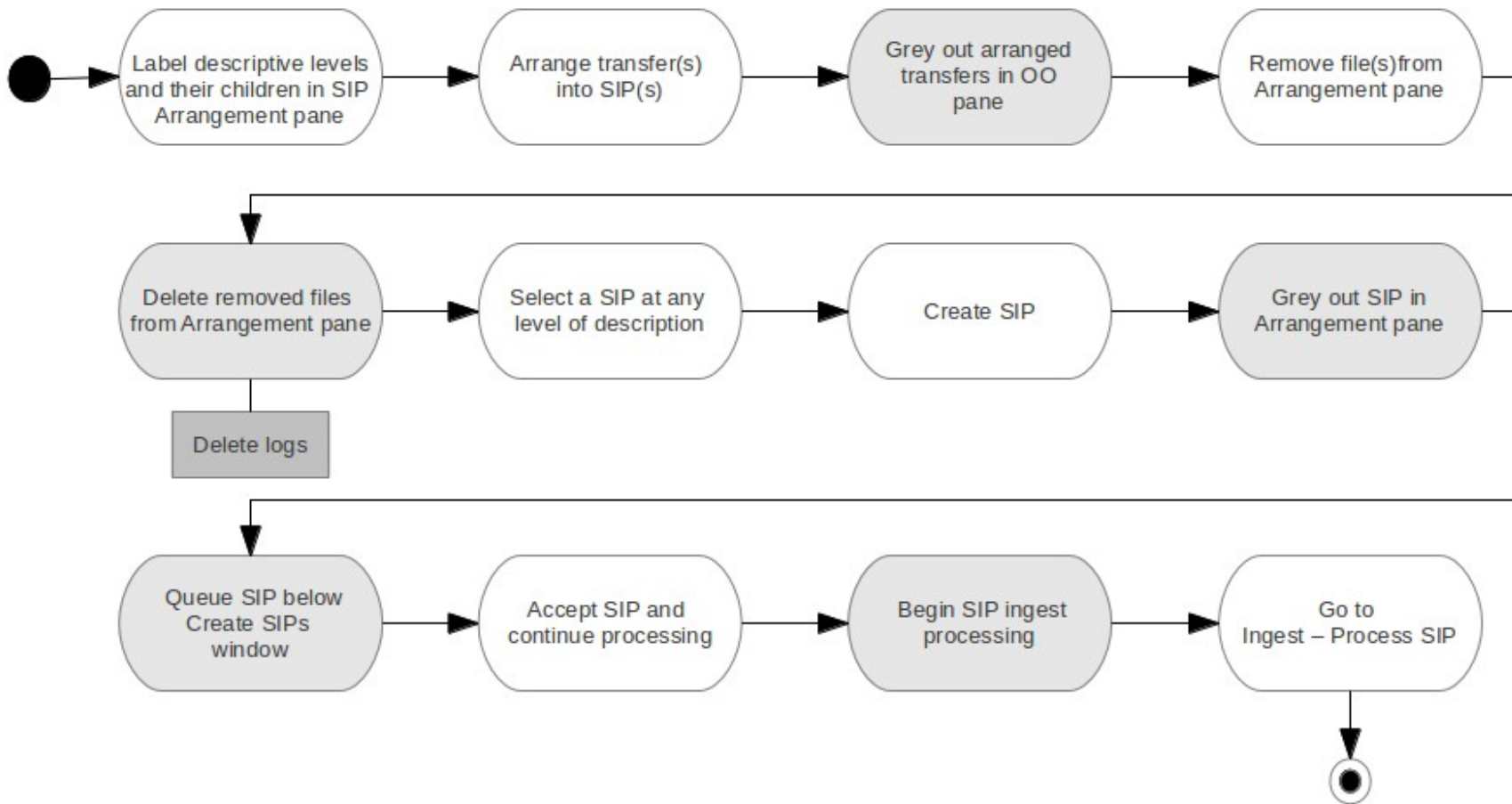
# Transfer indexing

- Full text content

- File embedded metadata

- Formats - by folder, by transfer

- Keyword & pattern matching for privacy/security sensitive information (e.g. social insurance numbers/social security numbers, credit card numbers, email addresses security keywords like 'private', 'confidential' - find or generate domain-specific taxonomies)

- Possible reports: PDFs that have not been OCR'ed, password protected / encrypted files, duplicates with their file paths

# What to do with that data

- Examine transfer(s)

- Use visualization tools and index to access transfer content

- Assign transfer(s) to appropriate intellectual levels of arrangement

- Create SIP(s) from transfer(s)

# 0.9 Ingest – Create SIP workflow

```
┌─────────────────┐      ┌─────────────────┐      ┌─────────────────┐      ┌─────────────────┐
│ Select transfer(s)│ ──▶ │ Load selected    │ ──▶ │ Generate reports │ ──▶ │ Generate file    │
│ for arrangement  │      │ transfer(s) into │      │ (see transfer    │      │ visualizations   │
│ into SIPs        │      │ Original Order   │      │ indexing         │      │ (see transfer    │
│                  │      │ pane of Create   │      │ requirements)    │      │ indexing         │
│                  │      │ SIP window       │      │                  │      │ requirements)    │
└─────────────────┘      └─────────────────┘      └─────────────────┘      └─────────────────┘
```

- Select transfer(s) for arrangement into SIPs
- Load selected transfer(s) into Original Order pane of Create SIP window
- Generate reports (see transfer indexing requirements)
- Generate file visualizations (see transfer indexing requirements)

- Remove transfer(s) from OO pane
- Remove selected transfer(s) from OO pane
- Delete cached reports
- Delete cached file visualizations

- Select any file to view in new tab if viewer available
- Highllight selected file
- Grey out viewer icon if no viewer available for selected file
- Open selected file in viewer

- Select transfer(s) or contents
- Highlight selected transfer(s) or contents
- Examine transfer(s) or contents (See Examine contents mockups)
- Delete and move files as needed

- Generate log

```
●──▶ ┌─────────────────┐   ┌─────────────────┐   ┌─────────────────┐   ┌─────────────────┐
     │ Label descriptive │   │ Arrange transfer(s)│ │ Grey out arranged │ │ Remove file(s)from│
     │ levels and their  │──▶│ into SIP(s)       │──▶│ transfers in OO  │──▶│ Arrangement pane │
     │ children in SIP   │   │                   │   │ pane             │   │                  │
     │ Arrangement pane  │   │                   │   │                  │   │                  │
     └─────────────────┘   └─────────────────┘   └─────────────────┘   └─────────────────┘
```

Label descriptive levels and their children in SIP Arrangement pane

Arrange transfer(s) into SIP(s)

Grey out arranged transfers in OO pane

Remove file(s)from Arrangement pane

Delete removed files from Arrangement pane

Select a SIP at any level of description

Create SIP

Grey out SIP in Arrangement pane

Delete logs

Queue SIP below Create SIPs window

Accept SIP and continue processing

Begin SIP ingest processing

Go to Ingest – Process SIP

http://localhost/transfer/

**@rchivematica**  Transfer | Ingest | Preservation planning | Access | Administration

Log-in

## Create SIP(s)

*Enter Accession # or Transfer Name, or use Browse to locate transfers or entire accessions*

Original Order - Select transfers for arrangement

Q search

Browse

...ment - Arrange selected records

| Transfer | Type | Size | Creation Date | Transfer Type |
|---|---|---|---|---|
| ▶ 2009_Flyers | folder | 15.6MB | 2009-01-12 | Generic |
| ▼ Cats_Parade | folder | 8.8MB | 2002-01-09 | Digitized |
| 📁 Floats | folder | 4.3MB | 2003-02-23 | |
| GiantSiamese.gif | GIF | 1.2MB | 2003-03-14 | |
| GingerRibbons.jpg | JPEG | 1.9MB | 2006-05-17 | |
| Minneke_Poes.jpg | JPEG | 1.2MB | 2008-09-21 | |
| 📁 Permits | | | | |
| 📁 Events | | | | |
| ▶ 2006_Meetings | folder | | | |
| 📁 Minutes | | | | |
| 📁 Agenda | | | | |
| ▶ Volunteers | | | | |
| 📁 ParadeCleanup | | | | |
| 📁 Cat Toss | | | | |
| 📁 Plushy donations | | | | |
| ▶ Photographs | | | | |
| 📁 2003 | | | | |

| Name | Size |
|---|---|
| AM1260 - Festival of Cats fonds | |
| 📁 Series 1 - Publicity files | 28.4 GB |
| 📁 Flyers | 12 GB |
| 📁 Press releases | 16.4 GB |
| 📁 Series 2 - Events files | 15 GB |
| 📁 Opening Party | 9.1 GB |
| 📁 Cat Parade | 5.9 GB |
| 📁 Photographs | 4 GB |
| 📁 Volunteers | 2.9 GB |
| 📁 Series 3 - Sponsor files | 14 GB |
| 📁 Series 4 - Correspondence files | 1 GB |
| 📁 Series 5 - Cat Festival Committee files | 19 GB |
| 📁 Meetings | 9 GB |
| 📁 Agenda | |

*Greyed out transfers have been added to a SIP in the arrangement pane*

*In a new tab, search indexed selections for subject terms, confidentiality statements, CC and SIN #s, etc., and create visualizations of contents*

Examine contents  🔍 ⊖

Examine contents | Add child | Add sibling | Create SIP ⊖

*Open file in viewer - greyed out if not supported in browser*
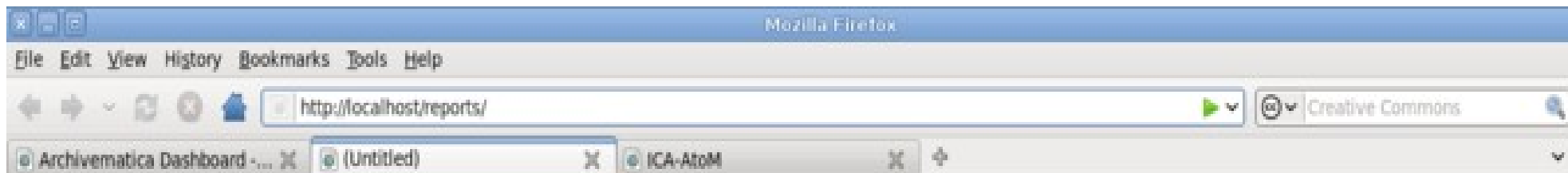
*Create SIP from selection(s)*

*Add metadata to SIP*

## SIP(s)

AM1260-S2-Photographs ✅ 📝

AM1260-S3-2-Busters-pets ✅ 📝

AM1260-S1-Flyers ✅ 📝

AM1260-S1-Press-releases ✅ 📝
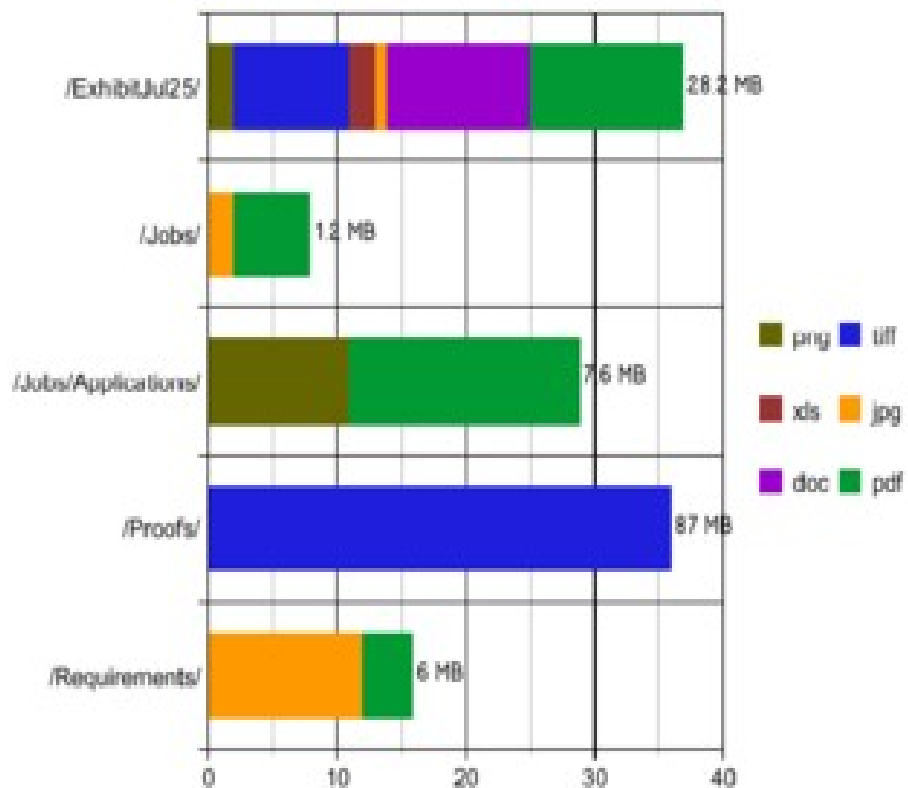
*Accept SIP and continue processing in dashboard*

# Examine contents splash page

- Information about the transfer or selected file group (number of files, size, name, uuid, accession #, and?)

- Pie graph showing file type distribution overall and bargraph showing file type by folder and ordered by size)

- Clickable links: file type opens into new tab with file browser interface of all of specified format), folders opens into new tab with file browser interface of entire folder in context of rest of transfer)

- Search box to search index (opens in new tab)

-  Report options (each opens in new tab): duplicates with full path locations, security keywords, CC numbers, SIN/SS#s, email addresses (with distribution graph), see password protected files with distribution (graph?)

File extension breakdown (by number
of files) in transferObjects/

File extensions by folder in transferObjects/
Total size of objects per folder

Archivematica Dashboard - Transfer

http://localhost/transfer/

# archivematica_

| Transfer | Ingest | Preservation planning | Access | Administration |

Log-in

## Create SIP(s)

Enter Accession # or Transfer Name, or use Browse to locate transfers or entire accessions

Original Order - Select transfers for arrangement

Q search

Browse

...ment - Arrange selected records

| Transfer | Type | Size | Creation Date | Transfer Type |
|---|---|---|---|---|
| ▶ 2009_Flyers | folder | 15.6MB | 2009-01-12 | Generic |
| ▼ Cats_Parade | folder | 8.8MB | 2002-01-09 | Digitized |
| 📁 Floats | folder | 4.3MB | 2003-02-23 | |
| GiantSiamese.gif | GIF | 1.2MB | 2003-03-14 | |
| GingerRibbons.jpg | JPEG | 1.9MB | 2006-05-17 | |
| Minneke_Poes.jpg | JPEG | 1.2MB | 2008-09-21 | |
| 📁 Permits | | | | |
| 📁 Events | | | | |
| ▶ 2006_Meetings | folder | | | |
| 📁 Minutes | | | | |
| 📁 Agenda | | | | |
| ▶ Volunteers | | | | |
| 📁 ParadeCleanup | | | | |
| 📁 Cat Toss | | | | |
| 📁 Plushy donations | | | | |
| ▶ Photographs | | | | |
| 📁 2003 | | | | |

| Name | Size |
|---|---|
| AM1260 - Festival of Cats fonds | |
| 📁 Series 1 - Publicity files | 28.4 GB |
| 📁 Flyers | 12 GB |
| 📁 Press releases | 16.4 GB |
| 📁 Series 2 - Events files | 15 GB |
| 📁 Opening Party | 9.1 GB |
| 📁 Cat Parade | 5.9 GB |
| 📁 Photographs | 4 GB |
| 📁 Volunteers | 2.9 GB |
| 📁 Series 3 - Sponsor files | 14 GB |
| 📁 Series 4 - Correspondence files | 1 GB |
| 📁 Series 5 - Cat Festival Committee files | 19 GB |
| 📁 Meetings | 9 GB |
| 📁 Agenda | 10 GB |

Greyed out transfers have been added to a SIP in the arrangement pane

In a new tab, search indexed selections for subject terms, confidentiality statements, CC and SIN #s, etc., and create visualizations of contents

**Examine contents**

🔍 ⊖

**Examine contents** | Add child | Add sibling | **Create SIP** | ⊖

Open file in viewer - greyed out if not supported in browser

Create SIP from selection(s)

Add metadata to SIP

## SIP(s)

| | | |
|---|---|---|
| AM1260-S2-Photographs | ✅ | 📝 |
| AM1260-S3-2-Busters-pets | ✅ | 📝 |
| AM1260-S1-Flyers | ✅ | 📝 |
| AM1260-S1-Press-releases | ✅ | 📝 |

Accept SIP and continue processing in dashboard

# File viewers

- This will allow the user to see individual documents in the transfer to get a better idea of their contents and technical metadata before assigning them to SIPs.

- Viewers are browser-dependent; viewer option is greyed out if viewer is not supported in browser

- Examine Contents window allows for viewing technical MD and other metadata available after Transfer microservices as well as indexing MD

@archivematica (official Twitter)
@snarkivist (me)

archivematica.org (wiki with links to documentation, downloads, user group and issues list)

courtney@artefactual.com (email)

artefactual
s y s t e m s     i n c .