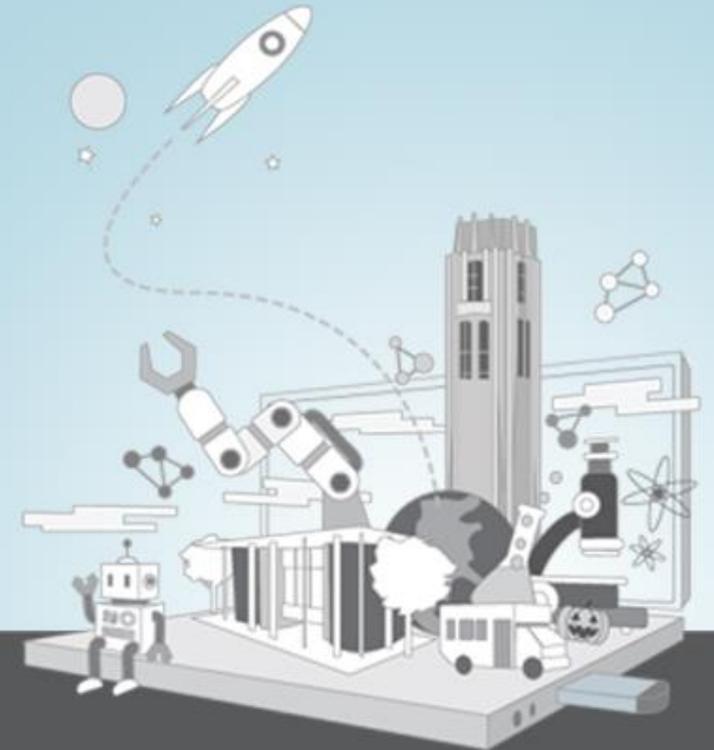# Developing software to prepare social science research data and code for sharing and preservation

Limor Peer

*Institution for Social and Policy Studies, Yale University*

2014 YALE TECHNOLOGY SUMMIT OCT. 31, EVANS HALL

# Data sharing



Source: Nature (GARY WATERS/IKON IMAGES/CORBIS)

# Using other people's data...

Lu Hugerth
@luhugerth

You know you're in trouble with #otherpeoplesdata when there's spaces on the file name

4:28 PM - 11 Dec 2013

---

Christie Bahlai retweeted

**iBartomeus** @ibartomeus 1d
Today #otherpeoplesdata problem is too many columns with derived values. Just trying to id which are raw values for now.

Details

---

**Ethan White** @ethanwhite 21h
MT @phylogenomics: @ekansa: Excel spreadsheets w/ color coding has meaning but terrible for other people to understand. #otherpeoplesdata

Details

---

**Christie Bahlai** @cbahlai 21h
Aaaannnnd: unexplained missing data. Must consult with data creator before proceeding. Will write stats code while I wait. #otherpeoplesdata

Details

---

**Christie Bahlai**
@cbahlai

MERGED CELLS IN EXCEL SPREADSHEET NOOOOOOOO! #otherpeoplesdata

11:10 AM - 6 Dec 2013

4 RETWEETS 5 FAVORITES

---

attempts were the lack

insu

on"

-cage/wp/2014/02/12/replication-in-political-science-

---

**Dr Elizabeth Sargent** @esargent184 · Sep 11
Oh no no no no no! Just received #otherpeoplesdata as a 276 page set of printed tables scanned in to a PDF

8    10

# Unusable data = lost data



Image: Shutterstock.com/Lightspring
http://slashdot.org/topic/datacenter/neglect-causes-massive-loss-of-irreplaceable-research-data

**Usable data:**

**Intelligently open**

**Independently understandable**

# Outline for today

**IF**:

Shared (and/or preserved) data may not be usable

**THEN**:

Make data usable = data curation

Project to develop curation software

- Background
- Requirements
- The software
- The architecture at Yale

## Data curation

Active and ongoing management of data through its lifecycle of interest and usefulness to scholarship, science, and education. Data curation enables data discovery and retrieval, maintains data quality, adds value, and provides for re-use over time through activities including authentication, archiving, management, preservation, and representation.

-- The University of Illinois' Graduate School of Library and Information Science

# ISPS Data Archive

# ISPS Data Archive

# Data Quality Review

**REVIEW FILES**

Assign persistent IDs * Create a citation to the study and a study level metadata record * Record file details (size, format, checksums) * Check that all files are present * Verify that content of files matches expected format * Create non-proprietary versions of the files * Implement migration strategy for file formats * Monitor bits

**REVIEW DATA**

Check for undocumented variable and value information or out of range codes * Review data for confidentiality issues
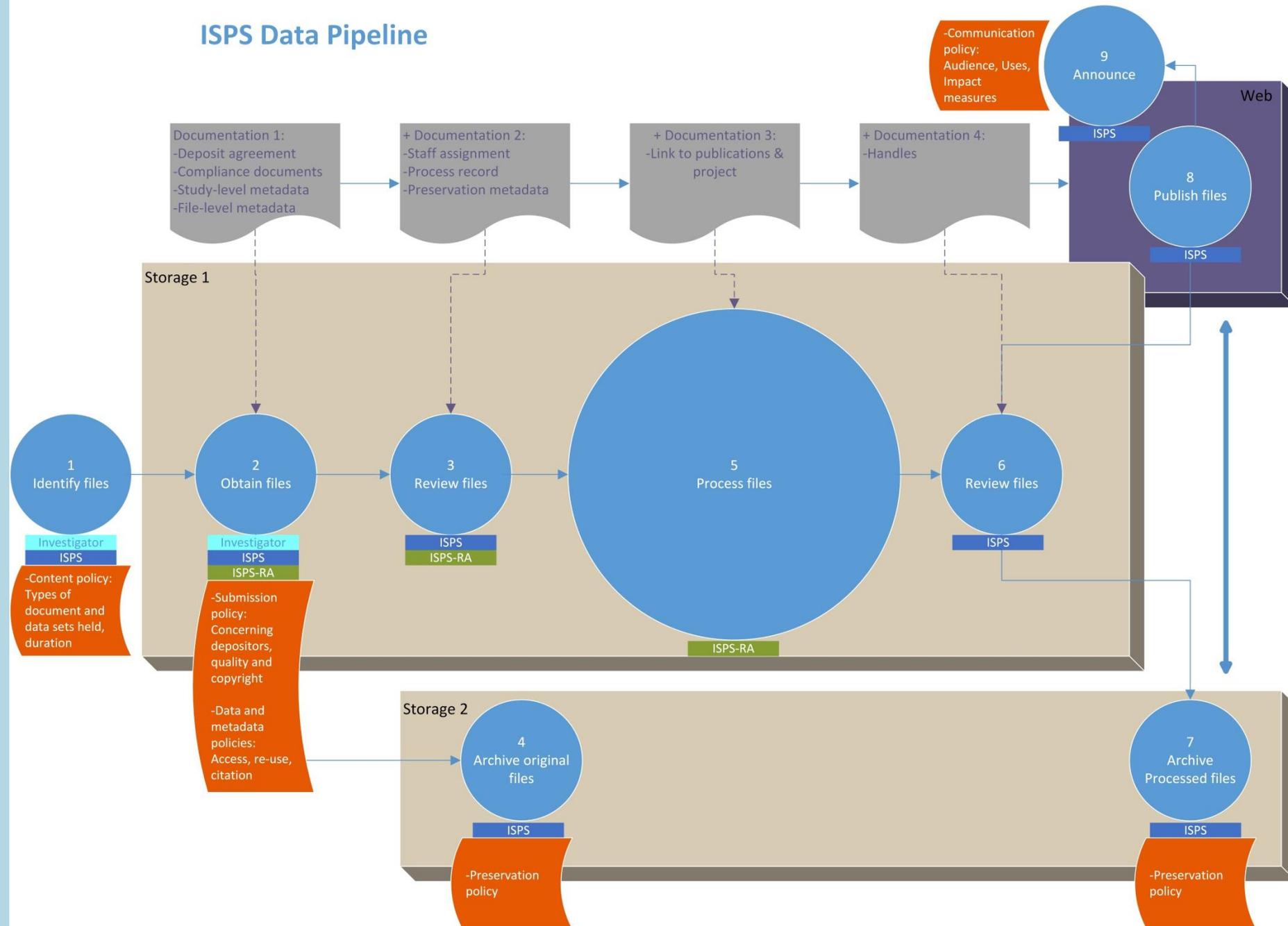
**REVIEW DOCUMENTATION**

Confirm comprehensive descriptive information for informed reuse including methodology and sampling information * Link to other research products

**REVIEW CODE**

Check and verify code for data analysis and replication

# ISPS Data Pipeline

-Communication policy: Audience, Uses, Impact measures

**9 Announce**

ISPS

Web

Documentation 1:
-Deposit agreement
-Compliance documents
-Study-level metadata
-File-level metadata

+ Documentation 2:
-Staff assignment
-Process record
-Preservation metadata

+ Documentation 3:
-Link to publications & project

+ Documentation 4:
-Handles

**8 Publish files**

ISPS

Storage 1

**1 Identify files**

Investigator
ISPS

-Content policy: Types of document and data sets held, duration

**2 Obtain files**

Investigator
ISPS
ISPS-RA

-Submission policy: Concerning depositors, quality and copyright

-Data and metadata policies: Access, re-use, citation

**3 Review files**

ISPS
ISPS-RA

**5 Process files**

ISPS-RA

**6 Review files**

ISPS

Storage 2

**4 Archive original files**

ISPS

-Preservation policy

**7 Archive Processed files**

ISPS

-Preservation policy

STRATEGIC PLAN >> 2013–2018

MORE EVIDENCE, LESS POVERTY

ipa
INNOVATIONS FOR
POVERTY ACTION

# Two Research Organizations

## The two organizations have in common…

- Similar content: Data from randomized controlled trials in the social sciences
- Similar approach to data sharing and preservation: Focus on replication, review data and code pre-publication

## Cross-fertilization…

- Build on ISPS Data Archive curation standards and practices
- Maintain key aspects of ISPS UI such as linked publications, data, and code
- Build on IPA ability to prepare data earlier in the lifecycle (e.g., pre analysis)
- Allow IPA network to access software from distributed research sites

# ISPS and IPA Requirements

- Curation workflow management (dashboard)
- Track changes to files (provenance)
- Integrate metadata production with data and code review and cleaning
- Preservation metadata and formats
- Secure storage and access
- Smooth transition to public dissemination of content
- Preference for open source solutions

# Curator software: Making data usable

A software platform that leverages the DDI Lifecycle and structures the curation workflow, including checking data for confidentiality and completeness, creating preservation formats, and reviewing and verifying code.

colectica

**Experts in social science metadata**

**Involved in DDI development**

# Features

- Web-based
- Built on DDI 3.2
- Open Source
- Builds on Existing Tools

# A Field Experiment on Legislators' Home Styles: Service versus Policy

| General | Files | Data | Collection | Quality | Related | Notes | Status | History |

## Upload Files

**+ Add files...**

Drop files here or use the button above.

File Details

# A Field Experiment on Legislators' Home Styles: Service versus Policy

| General | Files | Data | Collection | Quality | Related | Notes | Status | History |

| File | Type | Status |
| --- | --- | --- |
| Butler_et_al_JP_2012_metadata.pdf | Other | Accepted |
| Hyde_Marinov_IO_2014_observers.csv | Data | Accepted |
| Butler_et_al_JP_2012_Service_v_Policy_Limited_Dataset.dta | Data | Accepted |
| Butler_et_al_JP_2012_Service_v_Policy_Limited_Dataset.xml | Other | Accepted |
| Butler_et_al_JP_2012_JOP_Service_v_Policy_Limited_Do_File.do | Other | Accepted |
| Butler_et_al_JP_2012_Service_v_Policy_Limited_Dataset.csv | Data | Accepted |
| Butler_et_al_JP_2012_JOP_Service_v_Policy_Limited_R_File.R | Other | Accepted |

**☁ Add or Update Files**   **⬇ Download All Files**

Colectica curator software for ISPS – IPA research data repository (beta)

# Test1

| General | Files | Data | Collection | Quality | Related | Notes | Status | History |

| File | Type | Status | |
|------|------|--------|---|
| Gerber_et_al_JP_2011_README.pdf | Other | Accepted | ⊕ |
| Gerber_et_al_JP_2011_02_sub03_performmturk_analysis.do | Other | Accepted | ⊕  ⚙▾ |
| Gerber_et_al_JP_2011_mturk_replication.dta | Data | Pending | ⊕  ⚙▾ |

⚙▾
- Check Missing Labels
- Review Observation Count
- Compare Questionnaire, Codebook, and Data in Data File
- Check for Personally-Identifiable Information (PII) in Data File
- Identify Potential Errors in Data File

**⬆ Add or Update Files**   **⬇ Download All Files**

# Data curation

**Acquisition & Deposit** | **Ingest & Processing** | **Storage & Archival** | **Dissemination & Access**

**Deposit files & documentation**

Colectica Application

**Create metadata**

Colectica Repository DB for DDI Metadata

**Complete File bundle with Metadata**

**Provide preservation services**

**Provide access to files and metadata**

**Sign deposit agreement**

Colectica Application

**Build data & file set**

Colectica Application

**ITS Handle server**

**Yale Hydra Head**

**YUL FEDORA Commons**

**Yale Blacklight UI**

**YUL Dissemination server**

**ITS IIS server**

**ITS RSS file store**

**ITS SQL server**

**Web interface with Drupal nodes**

**YaleSites Dissemination server**

Blue = Curation action
Red = Applications & databases
Black = Unknown development
Green = ISPS and IPA research
Orange = ITS support / Secure access
Yellow = YUL IT support / Secure access

# Technical components & support at Yale

**ITS**

| Hardware | Windows Server (VM), 32GB RAM minimum (8 Cores), 100GB local disk for OS, applications and swap  files |
|---|---|
| Software | Colectica suite of tools, statistical software, integrated APIs |
| Storage | RSS start at 500GB, read/write/no-execute access to one or more directories |
| Application hosting | WCF application and ASP.NET MVC web application on IIS, plus a SQL Server database (10GB), a Windows Service |
| Security | Federated identification |

**Library**

| Long-term preservation | Fedora Commons / Hydra |
|---|---|
| Discovery | Blacklight |
| Persistent links | Handle service (ODAI) |

# (Target) Timeline

Project Kickoff – February 2014

Development Plan – March to April

Design + Base Platform and Basic Workflow development – May to October

Full Workflow Development – November to December

Ongoing development and maintenance – January 2015+

# Thank you!

limor.peer@yale.edu

@l_peer

In collaboration with:

**Innovations for Poverty Action**: Niall Keleher, Stephanie Wykstra

**Digital Lifecycle Research and consulting**: Ann Green

**Colectica** software company: Jeremy Iverson, Dan Smith

**Yale ITS Academic IT / Research Services**: Kiran Keshav, Themba Flowers, Paul Gluhosky

**Yale Library IT:** Michael Dula, Mike Friscia, Eric James
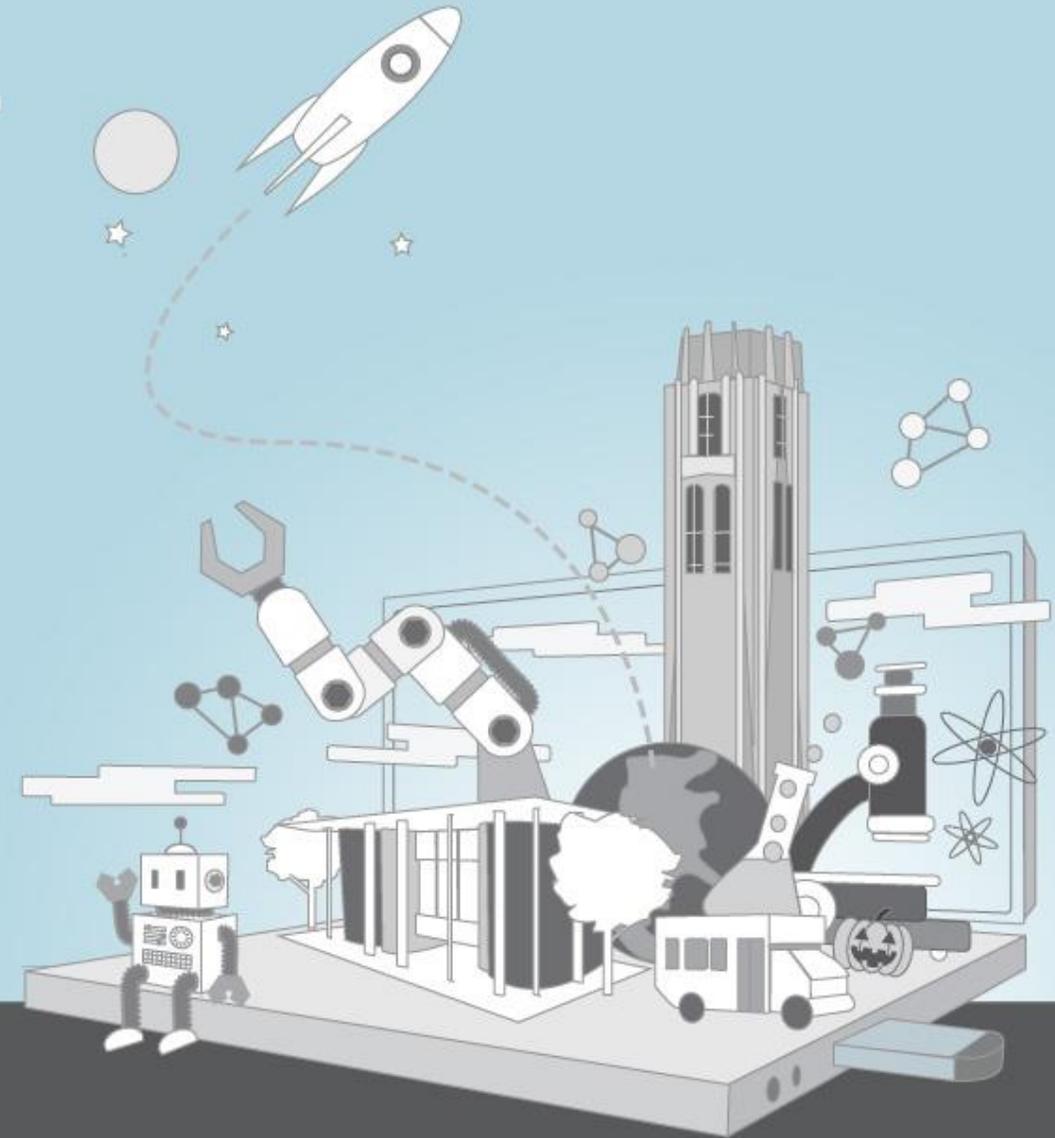
**Yale Library CSSSI**:  Michelle Hudson, Jill Parchuck

and Yale ODAI and Office of General Counsel

# 2014 YALE TECHNOLOGY SUMMIT

Friday, October 31, 2014
Evans Hall, Yale School of Management

#YaleTechSummit2014

*Presented by Yale Information Technology Services*