# New Curation Software:
# Step-by-Step Preparation
# of Social Science Data and Code
# for Publication and Preservation

Limor Peer & Stephanie Wykstra

IASSIST, June 2015

# Talk outline:

- Background: about IPA & ISPS, and our data policies/repositories

- Data curation – what we do

- More detail on code checks

- Curation workflow and partnership with Colectica

- Software demo

- Conclusion

# Background: about our groups

- Innovations for Poverty Action (IPA) works with academic researchers and partners (NGOs, governments) to carry out randomized controlled trials of poverty-alleviation programs in developing countries. Founded in 2002.
  - 300+ ongoing and completed projects.

- Institution for Social And Policy Studies (ISPS) works with academic researchers primarily in political science. Founded in 1968.
  - 100+ studies archived, many of them field experiments.

# Background: about our groups

- Both ISPS and IPA have data-sharing policies and data repositories

# IPA's data publication policy

| Level of sharing | Materials included | Allows for |
|---|---|---|
| Less than IPA-recommended | • Data and code underlying published results<br>• Readme file explaining relation between the files<br>• Minimal study-level metadata (i.e., information about the study). | • Checking that the tables in the published article can be produced by running the code<br>• limited re-use |
| IPA-recommended | *First row, plus...*<br><br>• Full set of collected and other (e.g., admin) variables, excluding personally identifiable information (PII).<br>• Data documentation including additional useful information about the study (context, notes on data-cleaning process, etc.)<br>• Survey instrument(s) | *First row, plus...*<br><br>• Fuller potential for re-use and secondary analysis<br>• Better understanding of study context – useful for systematic reviews and application to policy |
| Exceptional | *First two rows, plus...*<br><br>• Cleaning and variable construction code | *First two rows, plus...*<br><br>• Start-to-finish reproducibility |

# IPA's data repository (via Dataverse)

# ISPS Data Archive policy

**What to Archive?** Back to Top

ISPS-affiliated authors and PIs are expected to provide raw data and other information related to ISPS-supported research (e.g., instructions, treatment manuals, questionnaires, software, details of procedures, etc.). Deposits should include all data and documentation necessary to independently read and interpret the data collection. To use the ISPS Data Archive, authors and PIs are *required* to deposit the following types of files:

- Data File(s)
- Program File(s)
- Publication Citation
- Link to publication

Other types of files are *encouraged* but not required:

- Output File(s)
- Codebook(s)
- Study metadata
- Treatment Materials
- Supplementary Materials

# ISPS Data Archive

# Data curation

- "Maintaining, preserving and *adding value* to research data." (Digital Curation Center).

# Data curation

- Adding value for which purposes?
  - Reproducible research: ability to re-generate research results from data and code.
  - Data re-use for secondary analysis/ meta-analysis

- Data curation:
  - Making research materials independently usable, in order to allow reproducibility and re-use.

# Data curation for independent usability:

- Check for missing variable labels and value codes.
- Compare questionnaire and data.
- Ensure there is no personally-identifiable information (PII)
- Share key study-level metadata.
- Create open formats.
- Confirm code executes.
- Confirm code produces reported results.

- **What?**
  - ◻ Running the code and confirming that the output matches the tables in the paper.

- **Purpose?**
  - ◻ Catching discrepancies i.e., differences between output and publication tables.
  - ◻ Often, we have problems *running* the code initially and need to diagnose the issue e.g., different versions of commands, syntax errors, user-generated commands needed, missing code for some published tables.

# Data Quality Review Framework

**REVIEW FILES**

FILE://

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**

1
1 0
1 1
1 0 0
1 0 1
1 0 0 1

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

# Current workflow

**Web**

9
Announce

8
Publish files

**Processing Space**

1
Identify files

2
Obtain files

3
Review files

5
Process files

6
Review files

**Archive Space**

4
Archive original files

7
Archive processed files

# Curation tool: Workflow objectives

- Automate as many curation tasks as possible
- Technically integrate curation tasks using existing tools

And… to do this using a structured but flexible workflow from deposit to publication and preservation of data and code.

- Track changes to files and metadata

# Curation tool: Key features

- ## Leverages DDI Lifecycle
  - Machine executable, open structured format supports research transparency
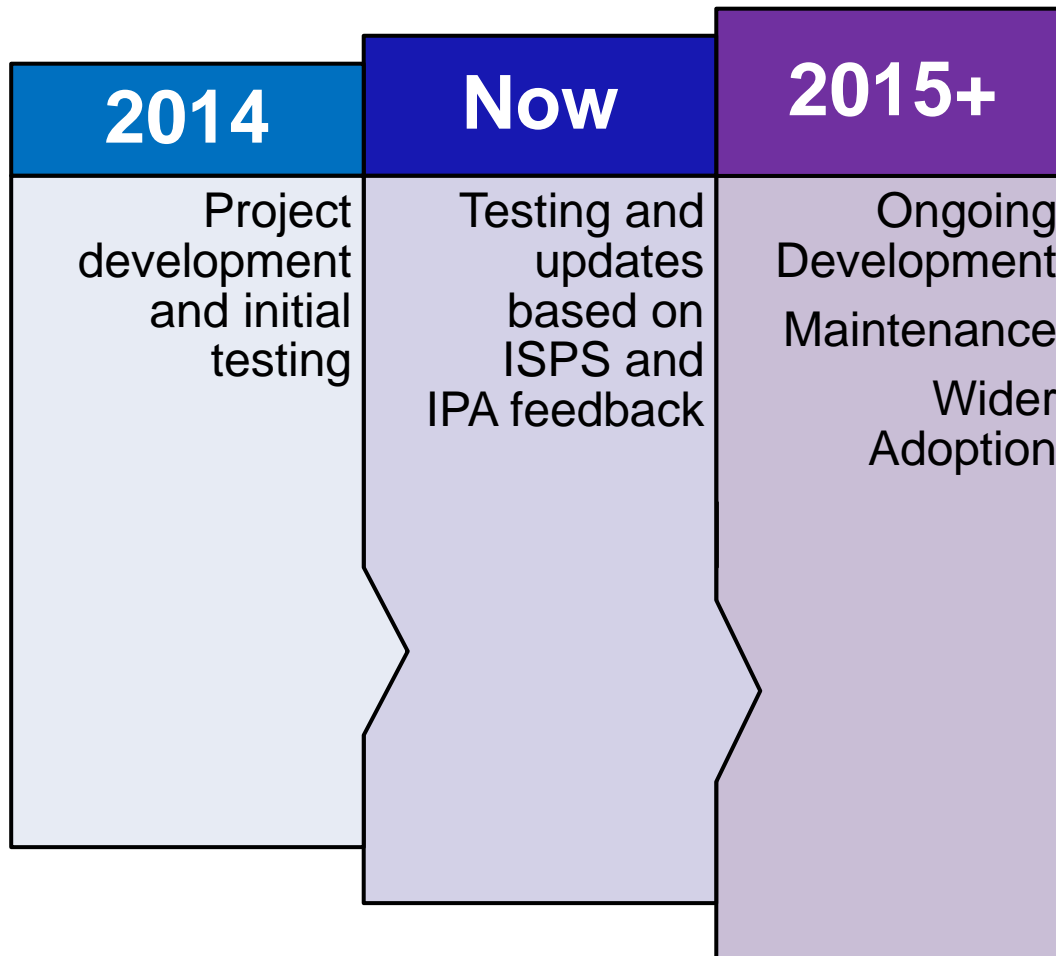  - Plays a part in review tasks – study, file & var levels

- ## Modular, open-source
  - Could be adapted to changing needs, research methods, dissemination platforms, and preservation solutions
  - Could be used by repositories, researchers, and research staff

# Timeline



| 2014 | Now | 2015+ |
|------|-----|-------|
| Project development and initial testing | Testing and updates based on ISPS and IPA feedback | Ongoing Development<br><br>Maintenance<br><br>Wider Adoption |

# Users & Roles

- **DEPOSITOR** Creates a new catalog, adds files, and submits for curation

- **ADMIN** Assigns curator

- **CURATOR** Performs review steps and requests publication.

- **APPROVER** Reviews the record, the history, and approves for publication and/or preservation.

# Demonstration: Preliminary!

# Curation tool: Summary

- ✓ Structures and tracks the curation workflow

- ✓ Helps automate parts of the data review pipeline

- ✓ Captures all metadata throughout the process

- ✓ Pushes out relevant information to pre-determined destinations

  - ✓ i.e., a user, the archive administrators, a Web based dissemination system, or preservation systems

- ✓ Can fit into repository and research workflows

# Thank you!

Limor Peer, limor.peer@yale.edu, @I_peer

Stephanie Wykstra, swykstra@poverty-action.org, @Swykstr

In collaboration with:
Jeremy Iverson & Dan Smith, Colectica
Ann Green, Digital Lifecycle Research and Consulting
Niall Keleher, Innovations for Poverty Action
Yale Academic IT / Research Services & Yale Library