# eResearch workflows for studying free and open source software development

James Howison, Andrea Wiggins, & Kevin Crowston

Syracuse University School of Information Studies

9 September 2008 ~ IFIP 2.13 - OSS 2008

# eResearch

- Scientific practices and technologies permitting distributed research collaborations using:
  - Large data sets
  - Computational resources
  - Analysis tools and workflows
  - Replicable research with provenance metadata
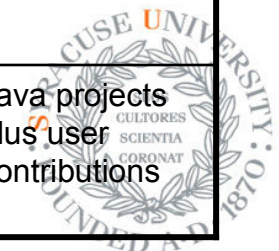- FLOSS research community is starting to move in this direction

# Current FLOSS Research Practices

- Data increasingly available in "repositories of repositories"
    - FLOSSmole
    - FLOSSmetrics
    - SRDA (Notre Dame)
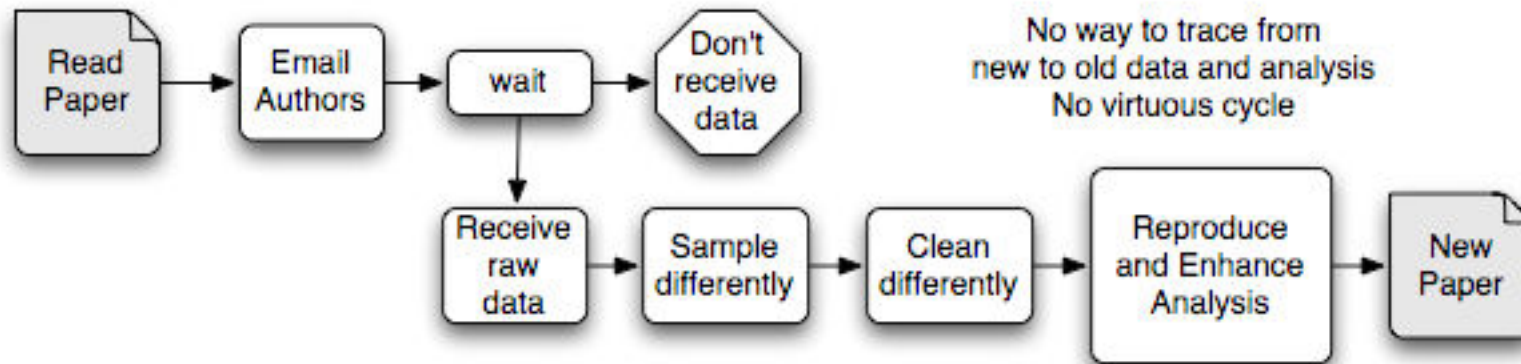- Not much sharing of analysis or methods for calculating measures; mostly bespoke scripts

# FLOSS Research Repositories

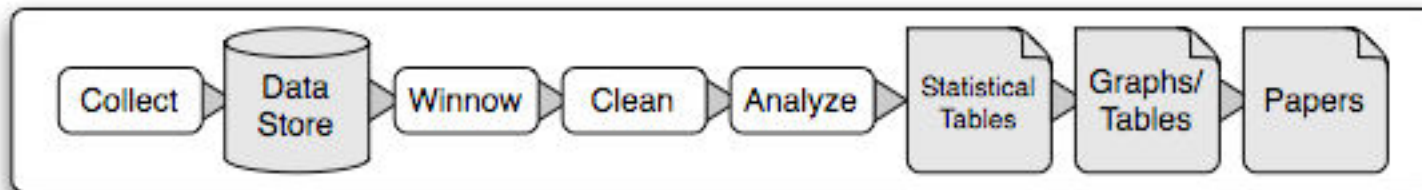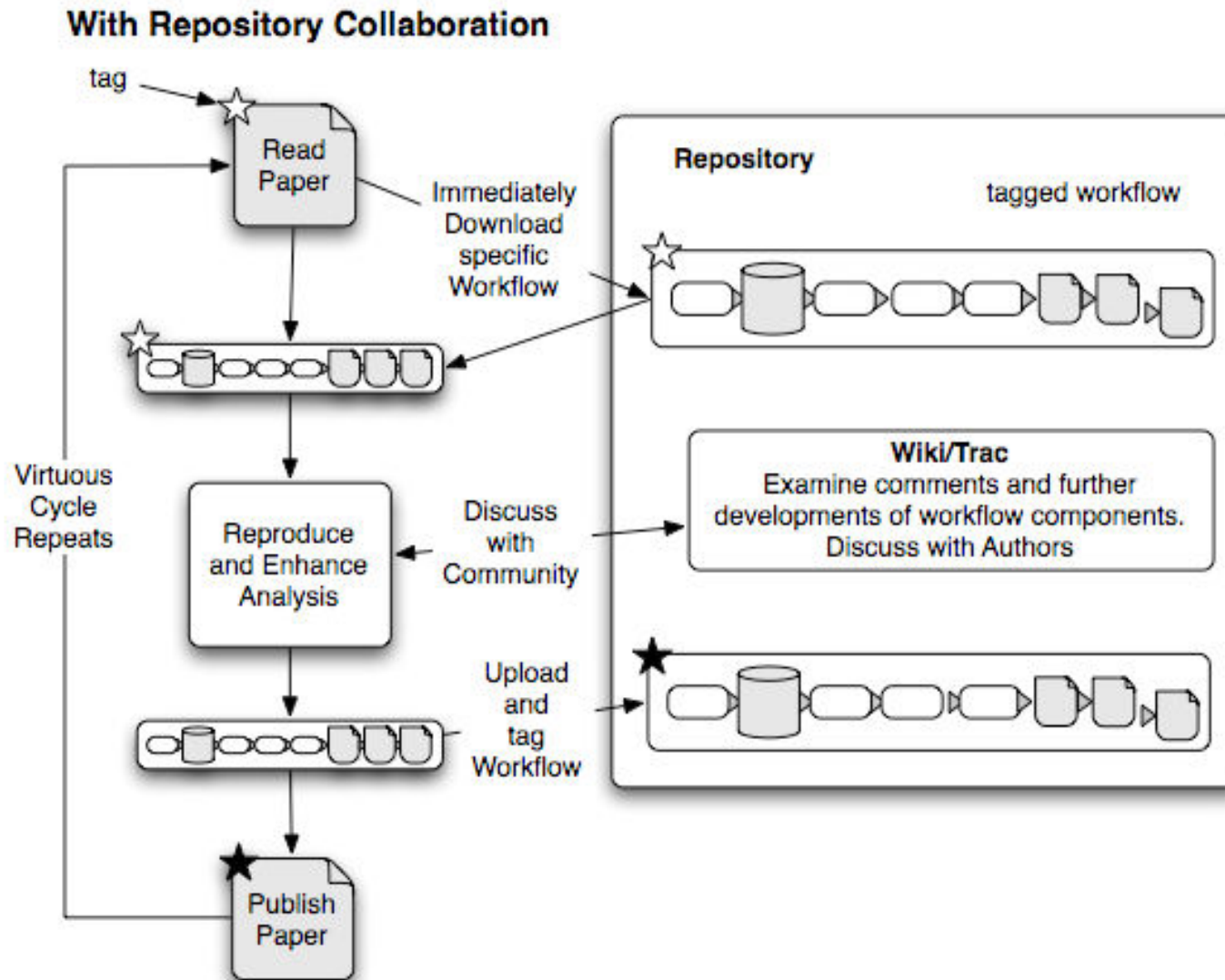| | FLOSSmole | SRDA (Notre Dame) | FLOSSmetrics & CVSAnalY | Qualoss & SQO-OSS | Source kibitzer |
|---|---|---|---|---|---|
| Project Demographics | Basic data | Basic data | Basic data & confirmed locations | | |
| Developer demographics | Memberships & roles | Memberships & roles | | | Memberships & roles |
| Communication venues | Releases; in progress: mail lists, forums, trackers | Forums & trackers | Planned/pilot: mail lists & trackers | | |
| Software venues | SVN/CVS count, packages, releases & dates | SVN/CVS count, packages, releases & dates | SVN/CVS full, size, packages, releases & dates | Planned/pilot: size & complexity metrics | Size & complexity metrics |
| Use & popularity | Downloads, pageviews, ratings, in Debian, partial: actual use | Downloads & pageviews | In Debian | | |
| Sample collected | Sourceforge, Rubyforge, ObjectWeb, Debian, freshmeat | Sourceforge | Partial: Sourceforge, ObjectWeb, Apache, GNOME | Pilot: KDE | Java projects plus user contributions |

# FLOSS Research Today

**Without Repository Collaboration**



Read Paper → Email Authors → wait → Don't receive data

wait → Receive raw data → Sample differently → Clean differently → Reproduce and Enhance Analysis → New Paper

No way to trace from new to old data and analysis
No virtuous cycle

**Abstract Workflow**

Collect → Data Store → Winnow → Clean → Analyze → Statistical Tables → Graphs/Tables → Papers

# FLOSS Research Tomorrow?

# Scientific Workflow Tools

- Tools for scientific analysis (e.g. Kepler, Taverna)
- Self-documenting analysis
  - Analysis conditions recorded at time of execution
- Steps in workflow executed by components with multiple input and output ports
- Components linked by joining input and output ports
- Supports modular analysis development, associated with easier collaboration and higher quality products
- Represented as flow diagram, stored and shared as single XML file

# Why Workflows Instead of Scripts?

- Differences and advantages over scripting languages
  - Wider accessibility
    - Programming skills helpful but not necessary to start
  - Compatibility
  - Easier integration of heterogeneous components
    - Mash up and reuse existing scripts
  - Standardized metadata
  - Out-of-the-box operation
    - Requires workflow software to execute, but little other configuration is required
- Example: BioPerl evolution from Perl scripts to libraries to workflow suites

# Taverna Workbench

- Open source stand-alone desktop application in Java
- Can also be run via "headless" server application
- Two main interface modes (v 1.71)
  - Design mode: workflow definition, automatically rendered diagram, available component types
  - Execution mode: process monitor, intermediate values, results, XML reports of status and all values
- Local and remote components are connected through input/output ports in MIME typing
  - Can also be grouped into subworkflows

# Taverna Design Mode

# Taverna Execution Mode

# Workflow Component Types

- Abstract & notification processors
- String constants: file locations, threshold values, etc.
- Beanshell: simplified Java
- Rshell: R statistical program running through Rserve
- Java widgets/shims for common operations
  - i/o, lists, text manipulation, JDBC, etc.
  - Command-line
- Web services
  - WSDL with SOAP
  - Can be "scavenged" from URLs or other workflows

# Example Workflow

- Teal: inputs
- Light blue: outputs
- Other light blue: string constant
- Green: web services
- Purple: Java shims
- Yellow: RShell
- Calculates weighted network centralization in dynamic networks, generates both numeric & graphical output

# Example Output

- Inputs
  - Sliding window size
  - Project venue
  - Date range

- Outputs
  - Time series graph of centralizations with summary stats
  - CSV of dates and centralization values for additional analysis



Network centralization for gaim-forum-users-helping-users over time

Mean = 0.34 , Standard Deviation = 0.19 Exponentially decayed edge weighting

# Benefits of Using Workflows

- Modular design yields benefits in flexibility, transparency, and ease of reuse
  - Easier to co-develop and debug components, and to integrate independent efforts
  - Can quickly change strategies with minimal adjustment to existing workflow structure
  - Can reuse existing scripts and workflow components
- Can also conduct exhaustive sensitivity testing
- Multiple ways to achieve analysis tasks

# Conclusion

- Combination of growing large-scale data sets and workflow tools present great opportunity for eResearch on FLOSS

- Work needed for eResearch infrastructure:
  - Access to data
  - Ontologies for naming data and defining relationships
  - Incorporating metadata & social science data, e.g. content analysis schemes

- Trade-offs involved in standardizing on tools to benefit collaboration, but much potential gain

# More

- Taverna demo screencast
  - Long version (24 minutes):
    floss.syr.edu/Presentations/tavernaDemoScreencast.mov
  - Short version (14 minutes):
    floss.syr.edu/Presentations/TavernaDemoRedux.m4v
- MyExperiment FLOSS group
  - www.myexperiment.org/groups/64
- 16:30 - 17:30 presentation today: Social Dynamics of FLOSS Team Communication across Channels
- Tomorrow: Workshop on Public Data about Software Development (WoPDaSD 2008)