Revealing the Detailed History of Script Outputs with Hybrid Provenance Queries

Yang Cao¹, Duc Vu², Qiwen Wang¹, Qian Zhang¹, Priyaa Ramesh³,Timothy McPhillips¹, Paolo Missier³, Bertram Ludäscher¹

¹University of Illinois, Urbana-Champaign, ²Department of Electrical and Computer Engineering, University of Illinois at Chicago, ³School of Computing Science, Newcastle University, UK



Motivation

- Data- and Workflow-Provenance are crucial for transparency and reproducibility in computational and data-driven science.
- Scientific workflow systems (Kepler, Taverna, ...) provide both **prospective provenance** (the workflow graph) and retrospective provenance (runtime observables).

Challenges

- Most computational analyses and workflows are conducted using **scripts** (Python, R, MATLAB, bash, ...) rather than workflow systems.
- Retrospective Provenance Observables, e.g., from DataONE RunManagers (file-level), ReproZip (OSlevel), or noWorkflow (Python code-level) only yield isolated fragments of the overall data lineage and processing history.
- Prospective Provenance could be used to link and contextualize fragments into a meaningful and comprehensible workflow, but scripts alone do not reveal the underlying workflow graph.
- Provenance (like other metadata) appears to be rarely actionable or immediately useful for those who are expected to provide it (provenance is "for others").

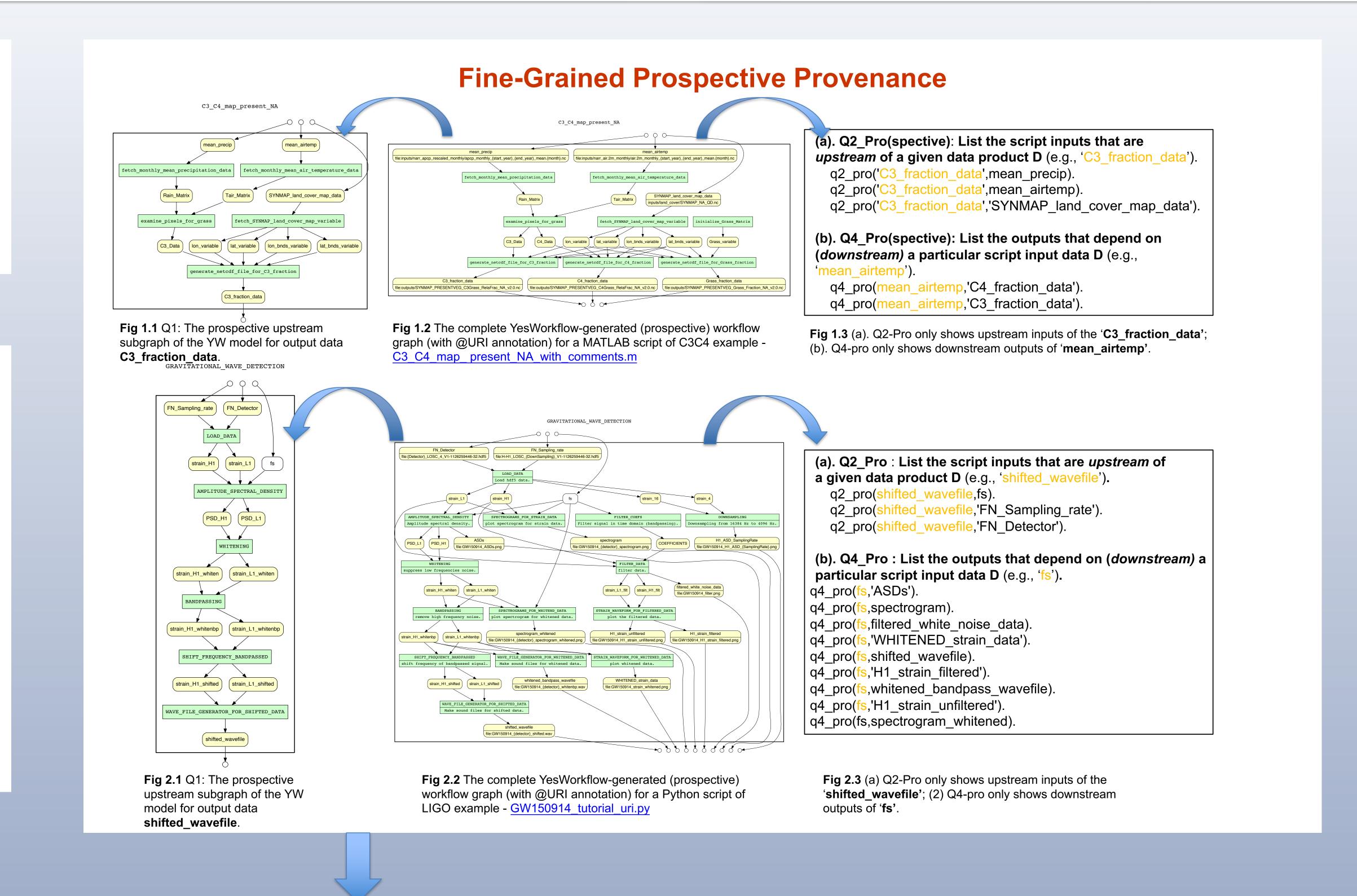
Approach

Simple YesWorkflow (YW) annotations allow users to reveal workflow (prospective provenance graph) implicit in scripts.

- Prospective provenance queries to expose and test data dependencies at the workflow level.
- Hybrid provenance queries that situate runtime observables (retrospective provenance) in the overall workflow, yielding meaningful knowledge artifacts.
- Easily share comprehensible workflow graphs and customizable provenance reports for script runs, along with data, code in scientific studies ("provenance for self').

Demo Queries

- Q1/Q3 (prospective query): Render prospective upstream/downstream subgraph of the YW model of the script for a given output data product D.
- Q2/Q4 (prospective query): List the script inputs that are upstream/downstream of a given data product D
- Q5 (hybrid query): Render retrospective graph with with concrete filename for a given output data product



Run Reconstruction

YW recon

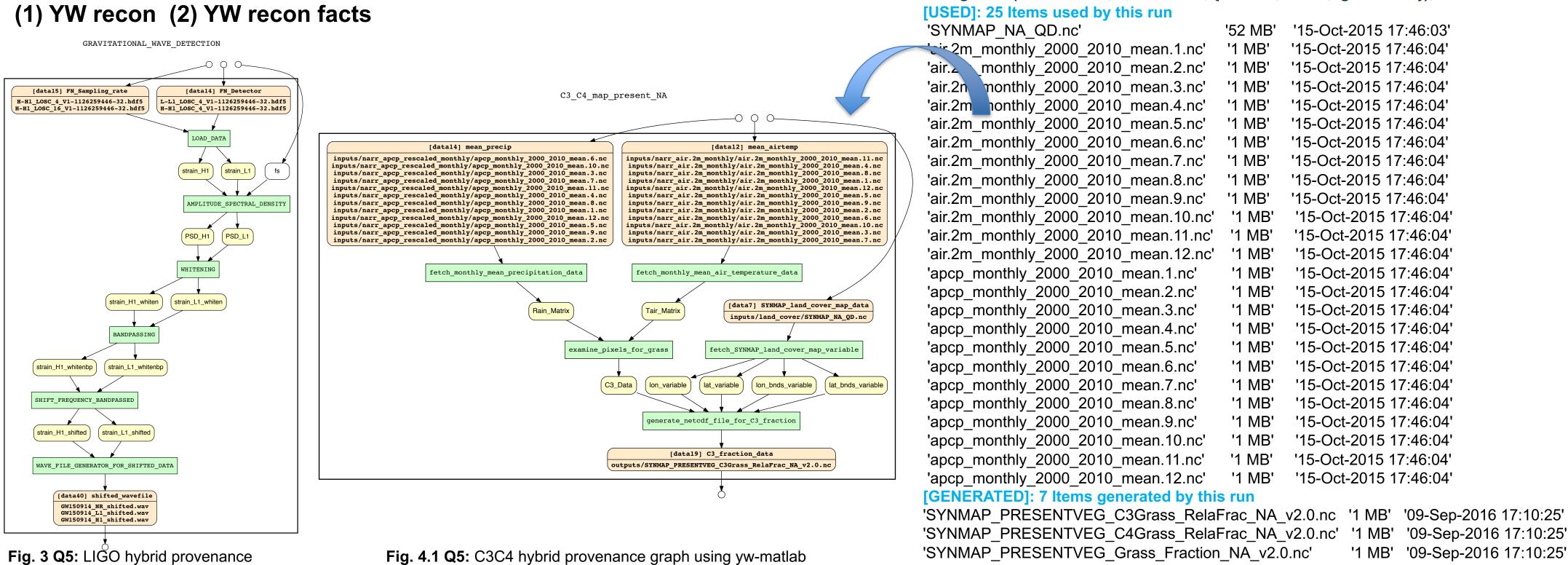
graph using yw-recon

YW recon facts

Coarse and Fine-Grained Observations of Runs

- Matlab Run Manager -> list of files input or output
- YesWorkflow -> list of files matching @URI annotations
- noWorkflow -> values assigned to variables

Hybrid Queries for Fine-Grained Retrospective Provenance



'SYNMAP NA QD.nc' '52 MB' '15-Oct-2015 17:46:03' 2m monthly 2000 2010 mean.1.nc' '1 MB' '15-Oct-2015 17:46:04' 'air.2 monthly 2000 2010 mean.2.nc' '1 MB' '15-Oct-2015 17:46:04 'air.2n monthly 2000 2010 mean.3.nc' '1 MB' '15-Oct-2015 17:46:04 'air.2m nonthly 2000 2010 mean.4.nc' '1 MB' '15-Oct-2015 17:46:04 'air.2m monthly 2000 2010 mean.5.nc' '1 MB' '15-Oct-2015 17:46:04' 'air.2m monthly 2000 2010 mean.6.nc' '1 MB' '15-Oct-2015 17:46:04' 'air.2m_monthly_2000_2010_mean.7.nc' '1 MB' '15-Oct-2015 17:46:04' 'air.2m_monthly_2000_2010_mean.8.nc' '15-Oct-2015 17:46:04' 'air.2m_monthly_2000_2010_mean.9.nc' '1 MB' '15-Oct-2015 17:46:04' 'air.2m_monthly_2000_2010_mean.10.nc' '1 MB' '15-Oct-2015 17:46:04 'air.2m monthly 2000 2010 mean.11.nc' '15-Oct-2015 17:46:04' 'air.2m_monthly_2000_2010_mean.12.nc' '1 MB' '15-Oct-2015 17:46:04' 'apcp_monthly_2000_2010_mean.1.nc' '15-Oct-2015 17:46:04' 'apcp_monthly_2000_2010_mean.2.nc' '15-Oct-2015 17:46:04' 'apcp_monthly_2000_2010_mean.3.nc' '15-Oct-2015 17:46:04 '1 MB' 'apcp_monthly_2000_2010_mean.4.nc' '15-Oct-2015 17:46:04' 'apcp_monthly_2000_2010_mean.5.nc' '15-Oct-2015 17:46:04 'apcp_monthly_2000_2010_mean.6.nc' '15-Oct-2015 17:46:04 'apcp_monthly_2000_2010_mean.7.nc' '15-Oct-2015 17:46:04 '1 MB' 'apcp_monthly_2000_2010_mean.8.nc' '15-Oct-2015 17:46:04 'apcp_monthly_2000_2010_mean.9.nc' '1 MB' '15-Oct-2015 17:46:04' 'apcp_monthly_2000_2010_mean.10.nc' '15-Oct-2015 17:46:04 'apcp_monthly_2000_2010_mean.11.nc' '1 MB' '15-Oct-2015 17:46:04' 'apcp_monthly_2000_2010_mean.12.nc' '1 MB' '15-Oct-2015 17:46:04' GENERATED]: 7 Items generated by this run

Fig. 4.2 C3C4 RunManager Screenshot

>> mgr.view('runNumber', 2, 'sections', {'details', 'used', 'generated'});

 Provenance from script runs can be revealed graphically and made actionable (e.g., to yield customizable data lineage reports) via (1) simple YW user annotations, (2) linking runtime observables (e.g. DataONE RunManager, ReproZip, noWorkflow), and (3) sharing provenance artifacts and executable queries.

Conclusions and Future Work

- Extend YW to facilitate querying log files for hybrid provenance at data level.
- Extend YW toolkit to support other (optional) workflow modeling constructs (e.g., simple control-flow to complement dataflow); to support graph pattern queries; to support project-level provenance.
- Evolve ProvONE to support project-level provenance and graph queries.

References

- Y Cao, D Vu, Q Wang, Q Zhang, P Ramesh, T McPhillips, P Missier, B Ludäscher (2016). DataONE AHM Provenance Demonstration: https://github.com/idaks/dataone-ahm-2016-
- YesWorkflow Project and Tools,
- https://github.com/yesworkflow-org
- T. McPhillips, T. Song, et al.(2015). YesWorkflow: A User-Oriented, Language-Independent Tool for Recovering Workflow Information from Scripts. Intl. Journal of Digital Curation 10, 298-313.
- T. McPhillips, S. Bowers, K. Belhajjame, B. Ludäscher (2015). Retrospective Provenance Without a Runtime Provenance Recorder. Workshop on the Theory and Practice of Provenance (TaPP).
- Cao, Y., Jones, C., Cuevas-Vicenttín, V., Jones, M.B., Ludäscher, B., McPhillips, T., Missier, et al., 2016, June. DataONE: A Data Federation with Provenance Support. Intl. Provenance and Annotation Workshop (IPAW). Springer.
- Pimentel, J.F., Dey, S., McPhillips, T., Belhajjame, K., Koop, D., Murta, L., Braganholo, V. and Ludäscher, B., 2016, June. Yin & Yang: demonstrating complementary provenance from noWorkflow & YesWorkflow. Intl. Provenance and Annotation Workshop (IPAW). Springer.

Acknowledgments. Supported by NSF awards ACI-1430508 and NSF ABI-1262458.











