# *GoldenTrail*: Retrieving the Data History that Matters from a Comprehensive Provenance Repository

Paolo Missier, Newcastle University, UK

Bertram Ludäscher, Saumen Dey, Michael Wang, Tim McPhillips, UC Davis, USA

Shawn Bowers and Michael Agun, Gonzaga University, USA

Ilkay Altintas, UC San Diego, USA

IDCC

Bristol, 6-7 Dec. 2011

Full-fledged data-mediated collaborations

IDCC '11 - P.Missier et al.

# Provenance in the experimental science lifecycle

A provenance trace is an account of the history of a data item through multiple processing steps

- Instrumental to verification and reuse of results -- Trustworthiness
- Enabler for "reproducible science" [1]

provenance trace (graph)

how did d4 come to be?
what other datasets contributed to it?
which processes were involved?

d1
i1
d2
d3
i2
d4
d5

i1 used d1 and d2

d4, d5 were generated by i2

[1] Mesirov , Jill, P. (2010). *Accessible Reproducible Research*. **Science**, 327. Retrieved from www.sciencemag.org

IDCC '11 - P.Missier et al.

4

## 2010: the DataTree Of Life summer project [2]

- Provenance *stitching*:

- Multiple, independently produced provenance traces expressed using the Open Provenance Model (OPM) can be "joined up" on shared datasets

- provided the data resides in a provenance-aware data repository.
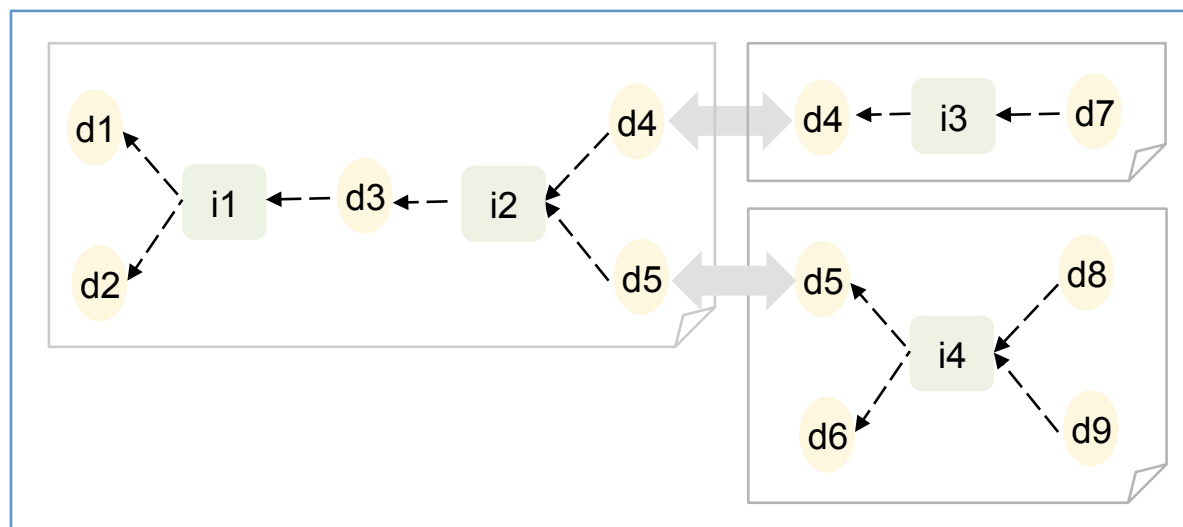


Limitations:

- automated "stitching" requires data ID mapping and provenance-aware data copy operations

- in general, it requires human intervention

[2] Missier, P., Ludascher, B., Bowers, S., Anand, M. K., Altintas, I., Dey, S., Sarkar, A., et al. (**2010**). *Linking Multiple Workflow Provenance Traces for Interoperable Collaborative Science*. Proc.s 5th Workshop on Workflows in Support of Large-Scale Science (**WORKS**).

- Experimental science is explorative and evolutionary
  - many experiments, few will succeed
  - from parameter sweeps to changes in methods

- E-science infrastructure should be able to capture the exploration process in addition to the "good" results

  - Implicit collaboration becomes "just" a special scenario

IDCC '11 - P.Missier et al.

6

- Experimental science is explorative and evolutionary
  - many experiments, few will succeed
  - from parameter sweeps to changes in methods

- E-science infrastructure should be able to capture the exploration process in addition to the "good" results

  - Implicit collaboration becomes "just" a special scenario

- Golden Data: the dataset(s) that scientists decide to share/publish
- Golden Trail: an account of how the Golden Data was obtained
  - a view over the provenance of the entire experiments history
  - describes a virtual experiment

6

# Approach: a generalized *provenance base*

## *PBase* Requirements

- Account for multiplicity of
  - workflow specifications and runs
  - workflow models
  - users

- Capture details of every execution into a persistent provenance repository

- Let scientists upload new provenance traces

- Support the provenance stitching process interactively

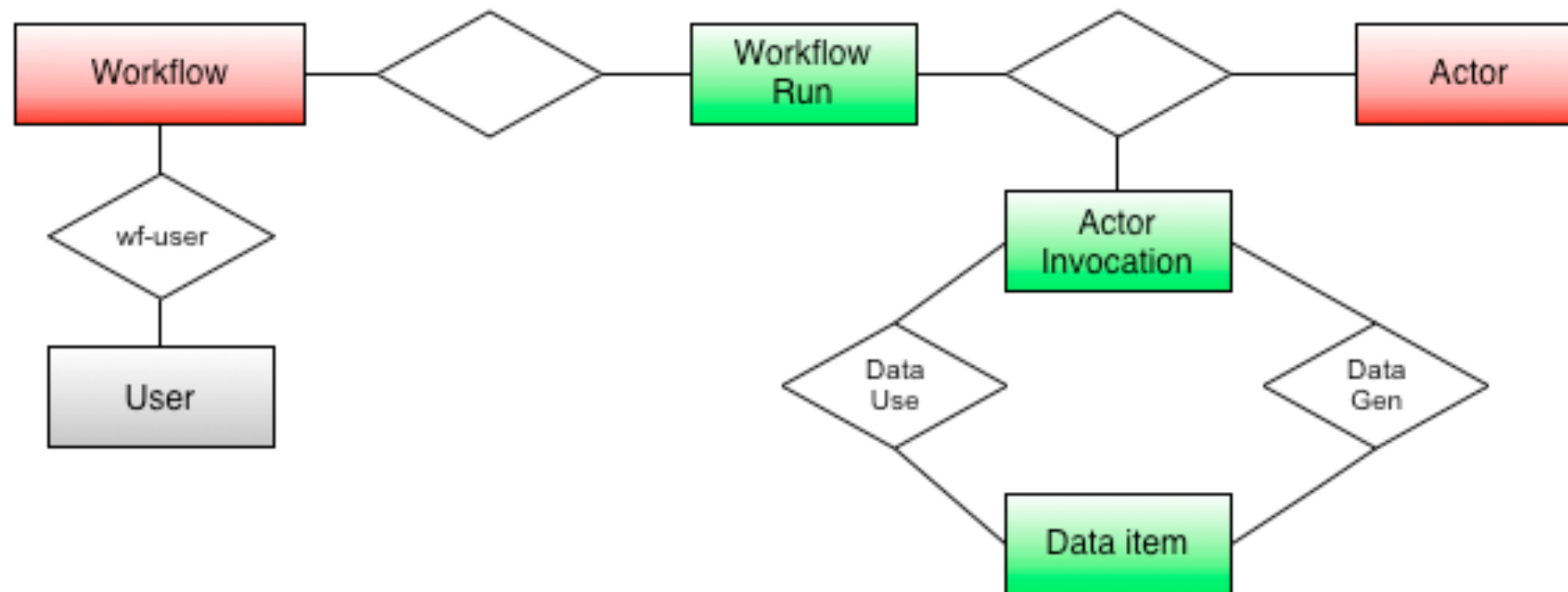- Support queries on the provenance base to compute Golden Trails

Goal:

To offer an extensible framework for building PBases

- The Open Provenance Model is adequate for describing traces of workflow execution: "trace-land"
  - to be superseded by PROV-DM, currently W3C Public Working Draft (*)

- But we also need to record workflow specifications: "workflow-land"
  - by supporting multiple heterogeneous workflow models
  - e.g. ASKALON, Galaxy, Kepler, Taverna, Pegasus, Vistrails, etc.
  - currently only Kepler (UCSD, UC Davis), Taverna (myGrid, UK) supported

- Integration with the DataONE data preservation architecture
  - Provenance base as a new type of *Member Node*

(*) FPWD as of October, 2001: http://www.w3.org/TR/2011/WD-prov-dm-20111018/
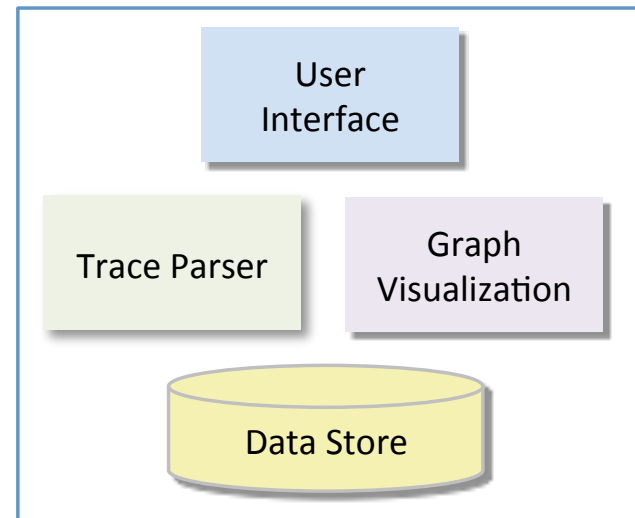
- Trace-land inspired by the OPM
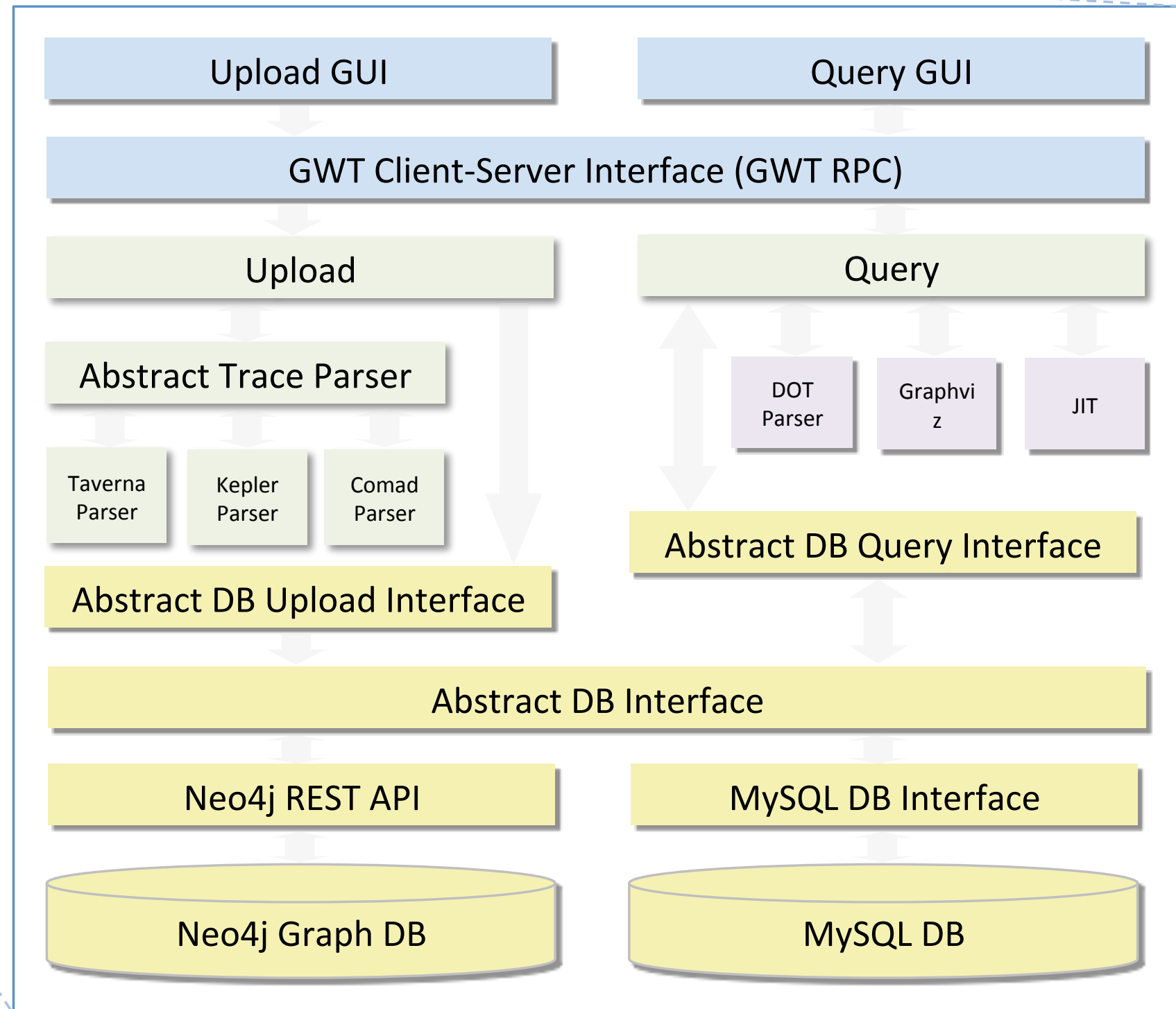
- Workflow-land inspired by Janus [1]



- Actor, a single computational step within a workflow

- Run: a single execution of an entire workflow

- Actor invocations: executions of individual steps that either Use or Generate Data Items

- Attribution: reference to users who run the workflow and thus "own" the traces.

[1] Missier, P., Sahoo, S. S., Zhao, J., Sheth, A., & Goble, C. (2010). Janus: from Workflows to Semantic Provenance and Linked Open Data. Procs. IPAW 2010. Troy, NY.

IDCC '11 - P.Missier et al.

9

User Interface

Trace Parser

Graph Visualization

Data Store

- UI: upload a new trace
- Trace Parser
  - maps native formats to D-OPM
- Graph Visualization
  - displays provenance graphs
- Data Store: provenance store

Upload GUI

Query GUI

GWT Client-Server Interface (GWT RPC)

Upload

Query

Abstract Trace Parser

DOT Parser

Graphviz

JIT

Taverna Parser

Kepler Parser

Comad Parser

Abstract DB Query Interface

Abstract DB Upload Interface

Abstract DB Interface

Neo4j REST API

MySQL DB Interface

Neo4j Graph DB

MySQL DB

- Exploit the synergy between workflow-land and trace-land

**Data-level and actor-level queries**

Ancestor / Descendant queries
(backwards / forward traversal)

Find all **Actors** that contributed to / impacted the generation of D
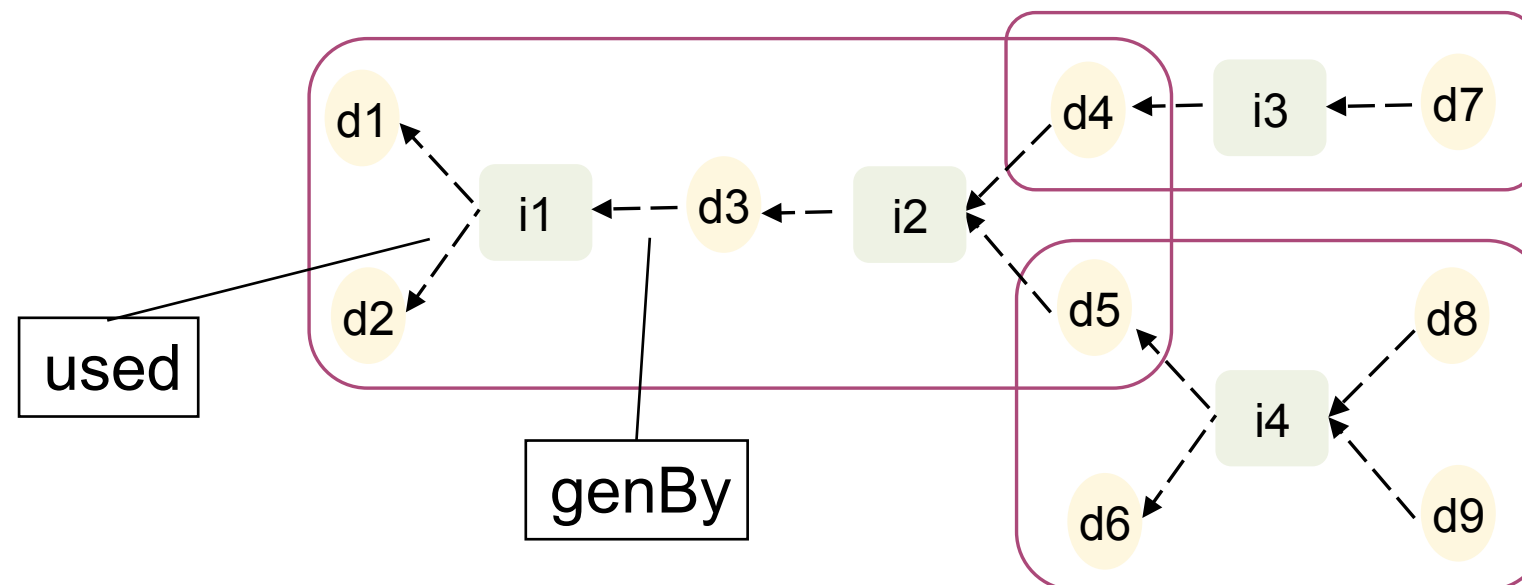
Find all **data** D'
that contributed to / impacted the generation of D

**Workflow-level queries**

Find all data that flowed through a workflow W during one run R

**User-related queries**

Find all data items used / generated on behalf of a user



used

genBy

11

**Golden-Trail**: A Provenance Repository For Storing And Retrieving Data Lineage Information

Select provenance detail level and dependency type

Filter results using conditions
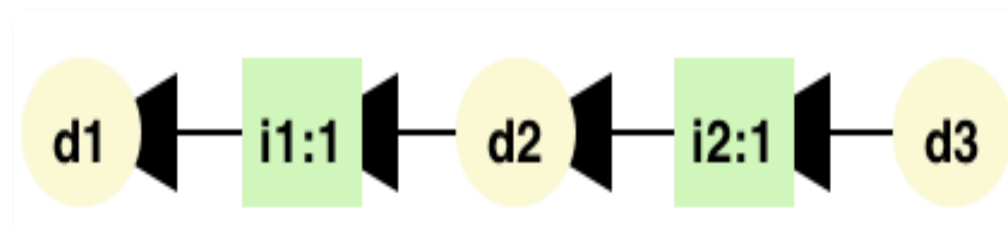
Add additional conditions

Query conditions

IDCC '11 - P.Missier et al.

In tabular format

In graphical format

IDCC '11 - P.Missier et al.

- GoldenTrail: a "Provenance Base" for workflow-related datasets
  - across users
  - across workflow models
  - across sessions

  - dedicated provenance model and query layer

- State:
  - early prototype completed (summer 2011) [1]

- Ongoing work within the DataONE project, Provenance Working Group
  - PBase to be integrated into DataONE as Member Node
  - Ongoing engagement with the scientific workflow community
    - get buy-in on the PBase idea
    - collect feedback on current prototype
    - collect additional use cases

[1] Dec. 6 2011: prototype available at: http://lore.genomecenter.ucdavis.edu:8080/GoldenApp/