# Provenance Exchange, Integration and Querying

**Marta Mattoso**

**Federal University of Rio de Janeiro, Brazil**

# Provenance
## Exchange, Integration and Querying

Contributors:

- M. David Allen, Adriane Chapman, Barbara Blaustein, Len Seligman
  [5 Getting It Together: Enabling Multi-organization Provenance Exchange]

- Anderson Marinho, Marta Mattoso, Claudia Werner, Vanessa Braganholo and Leonardo Murta
  [33 Challenges in managing implicit and abstract provenance data: experiences with ProvManager]

- Luiz M. R. Gadelha Jr., Marta Mattoso, Michael Wilde, Ian Foster
  [26 Provenance Query Patterns for Many-Task Scientific Computing]

# Importance of provenance in Science

- Interpret and reproduce data
- Understand the experiment and chain of reasoning that was used in the production of a result
- Verify that an experiment was performed according to acceptable procedures
- Identify what were the inputs to an experiment and where they came from
- Assess data quality
- Track who performed an experiment and who is responsible for its results (patents)

*Provenance is as (or more!) important as the results* (Davidson, Freire, Provenance and Workflows- SIGMOD 2008)

# Provenance along Wf levels



Databases | External Schema · Semantic WF | Workflows

Logical Schema · Abstract WF

Physical Schema · Executable WF
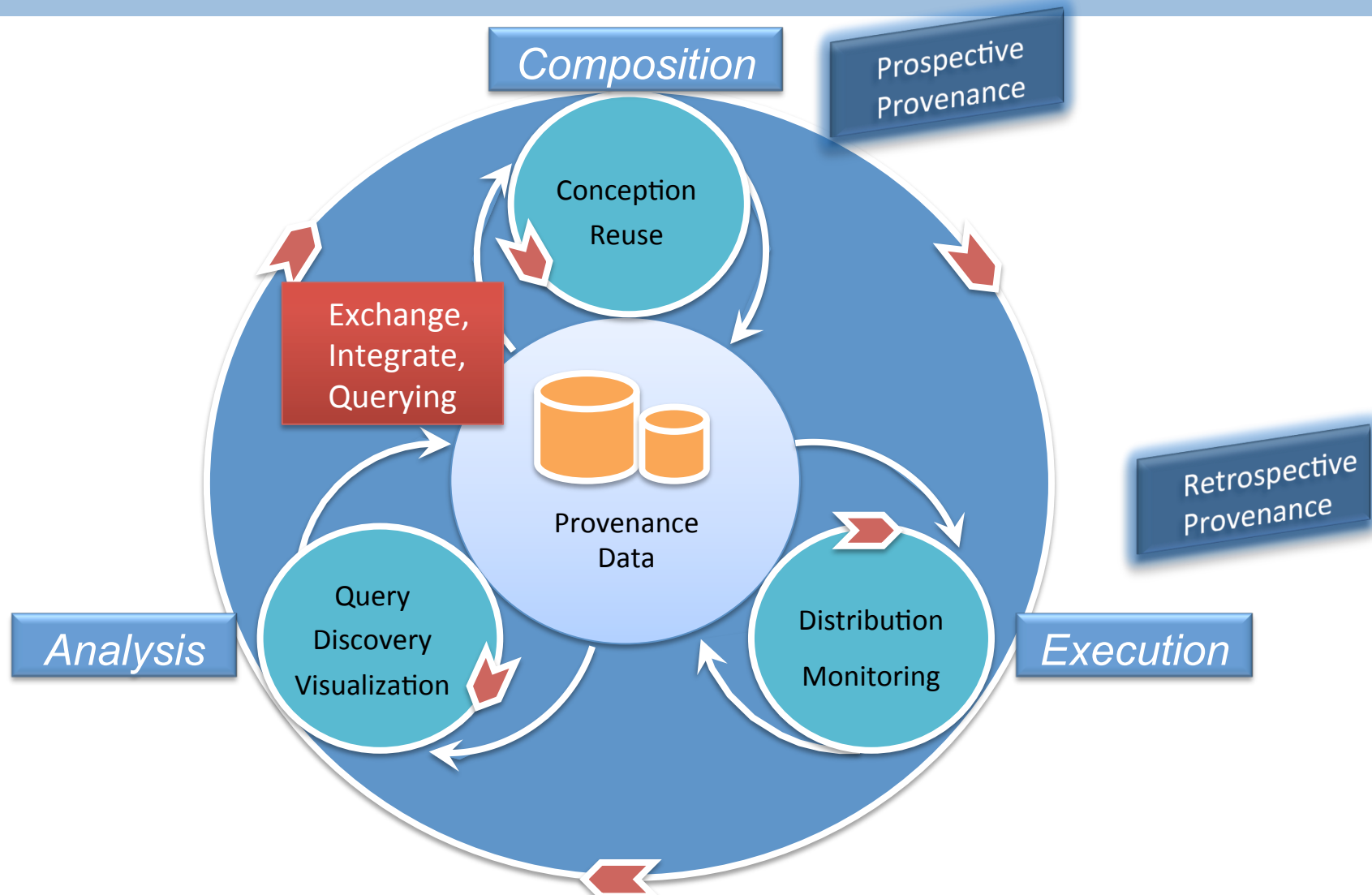
Provenance and Workflows – SIGMOD 2008 · Davidson & Freire 36

# Provenance can support analyzing scientific experiments

- Before execution:
  - What programs may be used? Is there any alternative methodology to explore?
  - Is there any dependency between activities? Which activities are mandatory?
- After execution:
  - What were the parameters used in the critical result ?
  - What were the scientific workflow activities used to obtain such result?
  - Where are the output files generated by the distributed activity A using the parameters P?
  - How many times the activity A in version V was used in the experiment E?
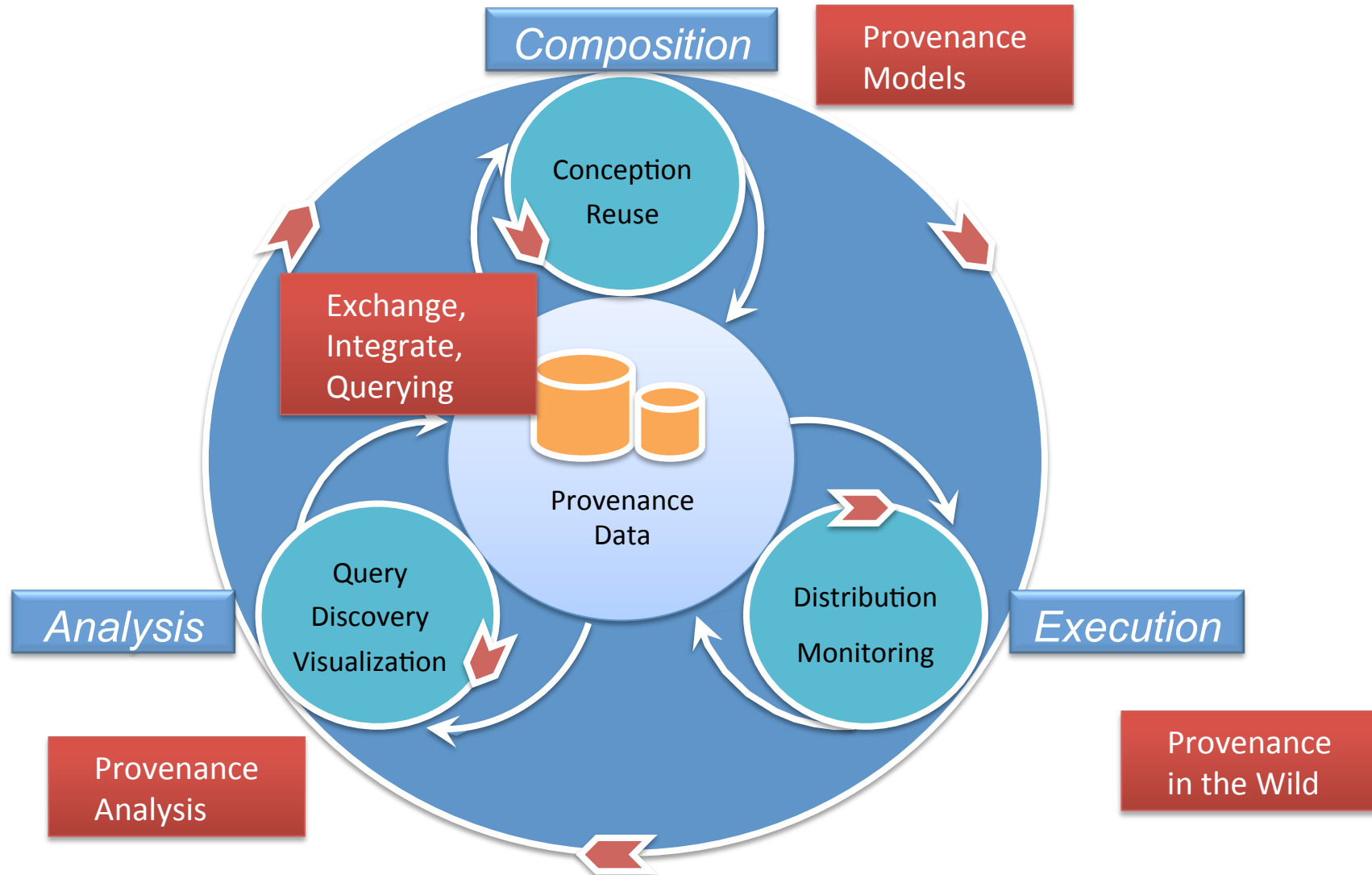
all these queries are related to the ability of reproducing and validating a scientific experiment
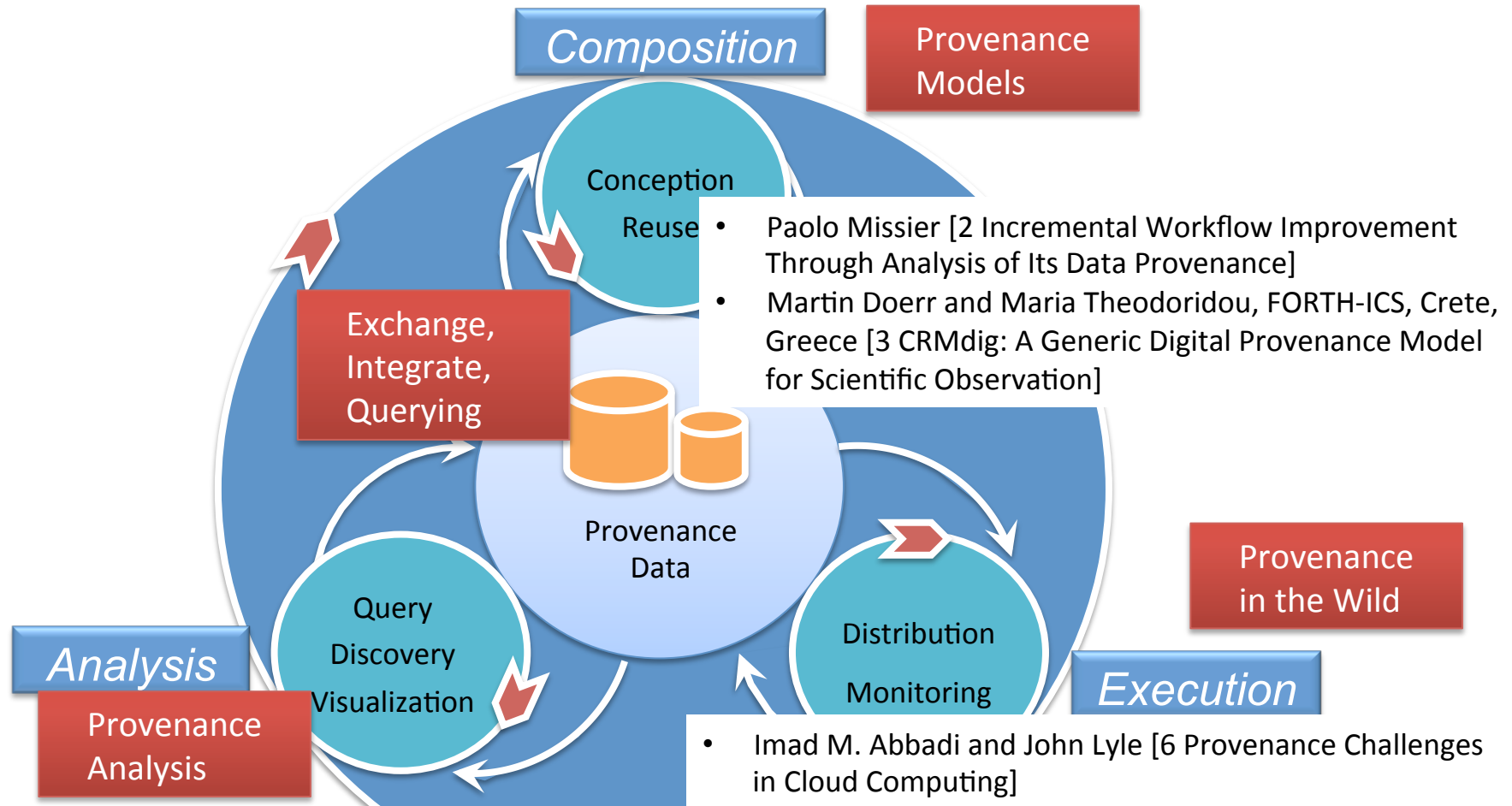
# Experiment Life Cycle*



*Mattoso et al, 2010 - Towards Supporting the Life Cycle of Large Scale Scientific Experiments. IJBPIM

# Experiment Life Cycle &TaPP Sessions

# Experiment Life Cycle &TaPP Papers

**Composition**

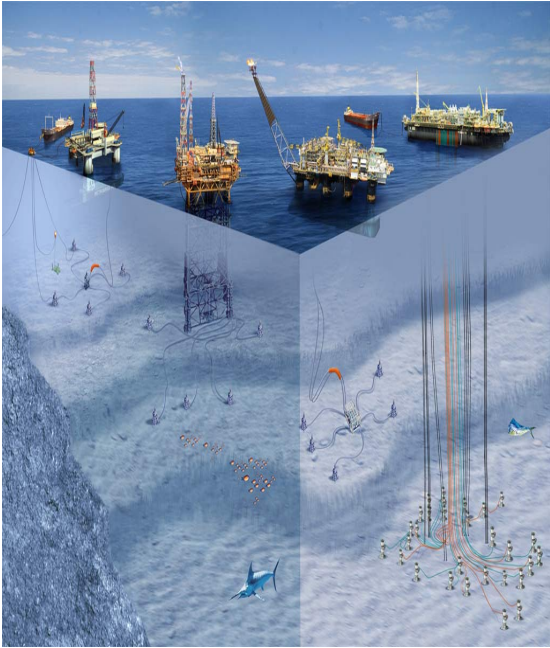Provenance Models

Conception Reuse

- Paolo Missier [2 Incremental Workflow Improvement Through Analysis of Its Data Provenance]
- Martin Doerr and Maria Theodoridou, FORTH-ICS, Crete, Greece [3 CRMdig: A Generic Digital Provenance Model for Scientific Observation]

Exchange, Integrate, Querying

Provenance Data

*Analysis*

Query Discovery Visualization

Provenance Analysis

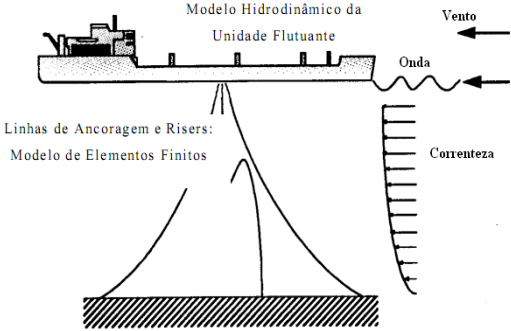Distribution Monitoring

Provenance in the Wild

*Execution*

- Imad M. Abbadi and John Lyle [6 Provenance Challenges in Cloud Computing]
- Peter Macko, Marc Chiarini, and Margo Seltzer [18 Collecting Provenance via the Xen Hypervisor]
- Elaine Angelino, Uri Braun, David A. Holland, Peter Macko, Daniel Margo, and Margo Seltzer [23 Provenance Integration Requires Reconciliation]

- Reng Zeng, Xudong He, Jiafei Li, Zheng Liu, W.M.P. van der Aalst [1 A Method to Build and Analyze Scientific Workflows from Provenance through Process Mining]

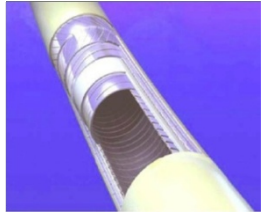# Risers' fatigue analysis in oil elevation from ultra-deep waters following a coupled analysis



Input Data to simulate Movements:
Waves, wind, currents, batimetryDados de onda vento, correnteza, bathymetry, etc. :
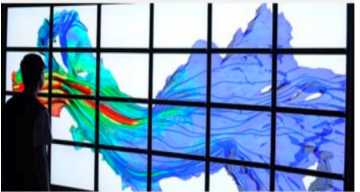
1. Coupled movement Analysis (TPN or Prosim)

Modelo Hidrodinâmico da Unidade Flutuante

Vento

Onda

Correnteza

Linhas de Ancoragem e Risers: Modelo de Elementos Finitos
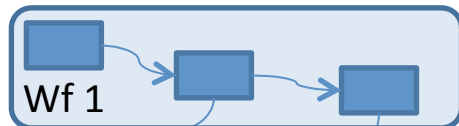
Generates large quantity of Data ...
(finite element meshes )

Estimate risers lifetime

3. Results are analyzed POSFAL
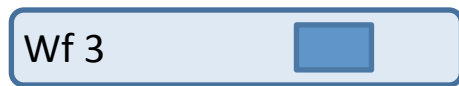
2. ... To do structural Analysis ofRisers (ANFLEX)

# Distributed Provenance



**Scientific experiment**

Wf 1

Wf 2

Offline analysis (vis cave)

Wf 3

Provenance Systems

Publish Experiment and Workflow

Analysis of Movements of Platform

Movements Filtering

Analysis of Risers' Structure

Analysis of Fatigue
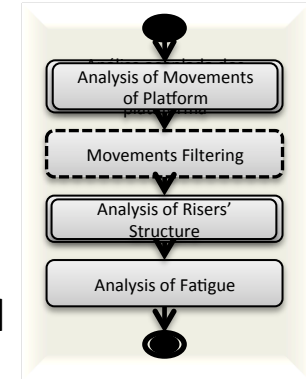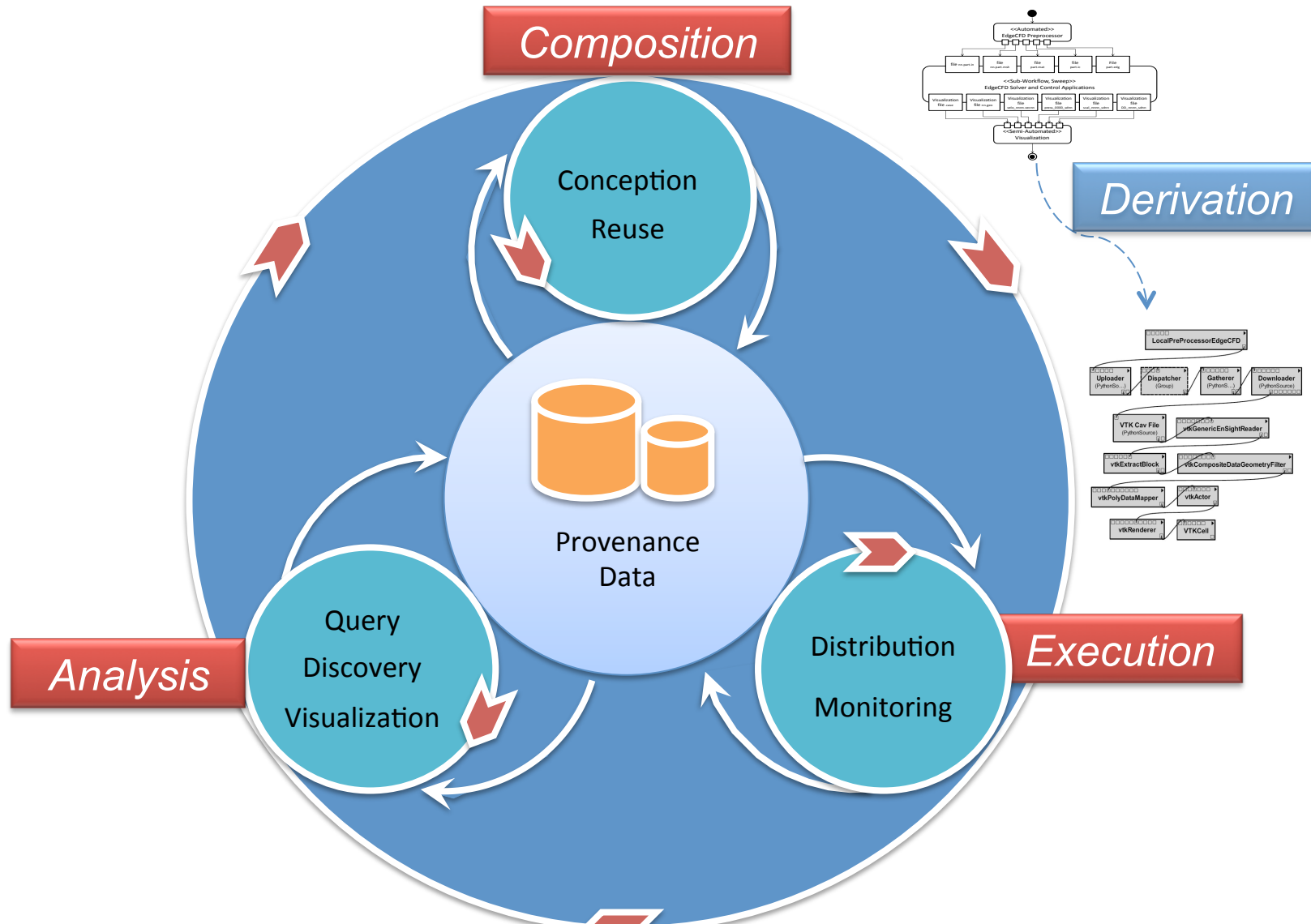
Sub-workflow parallel execution in HPC clusters, clouds

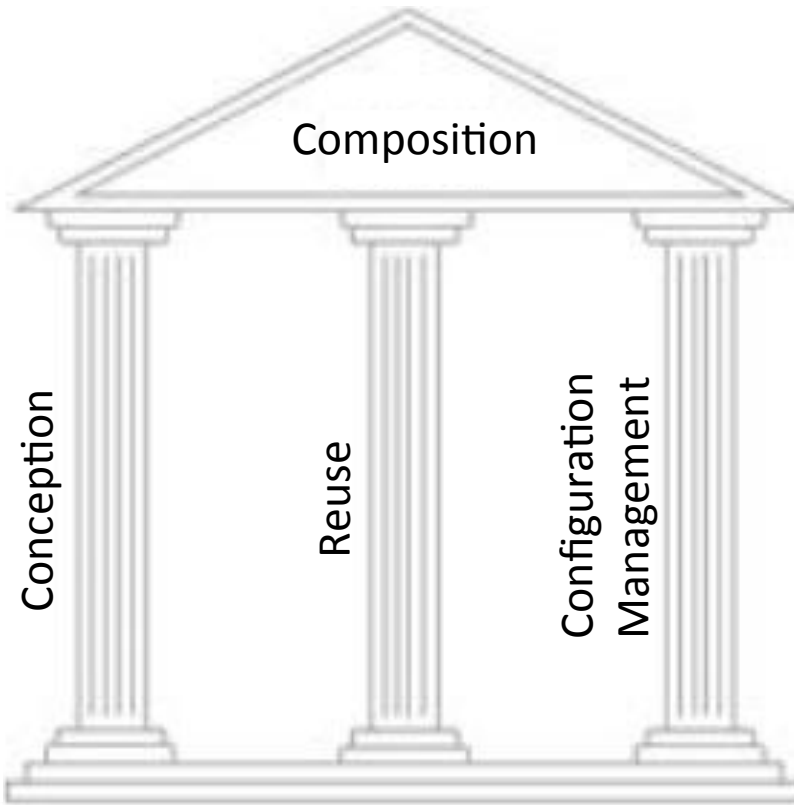Visualize and Share provenance data with others scientists
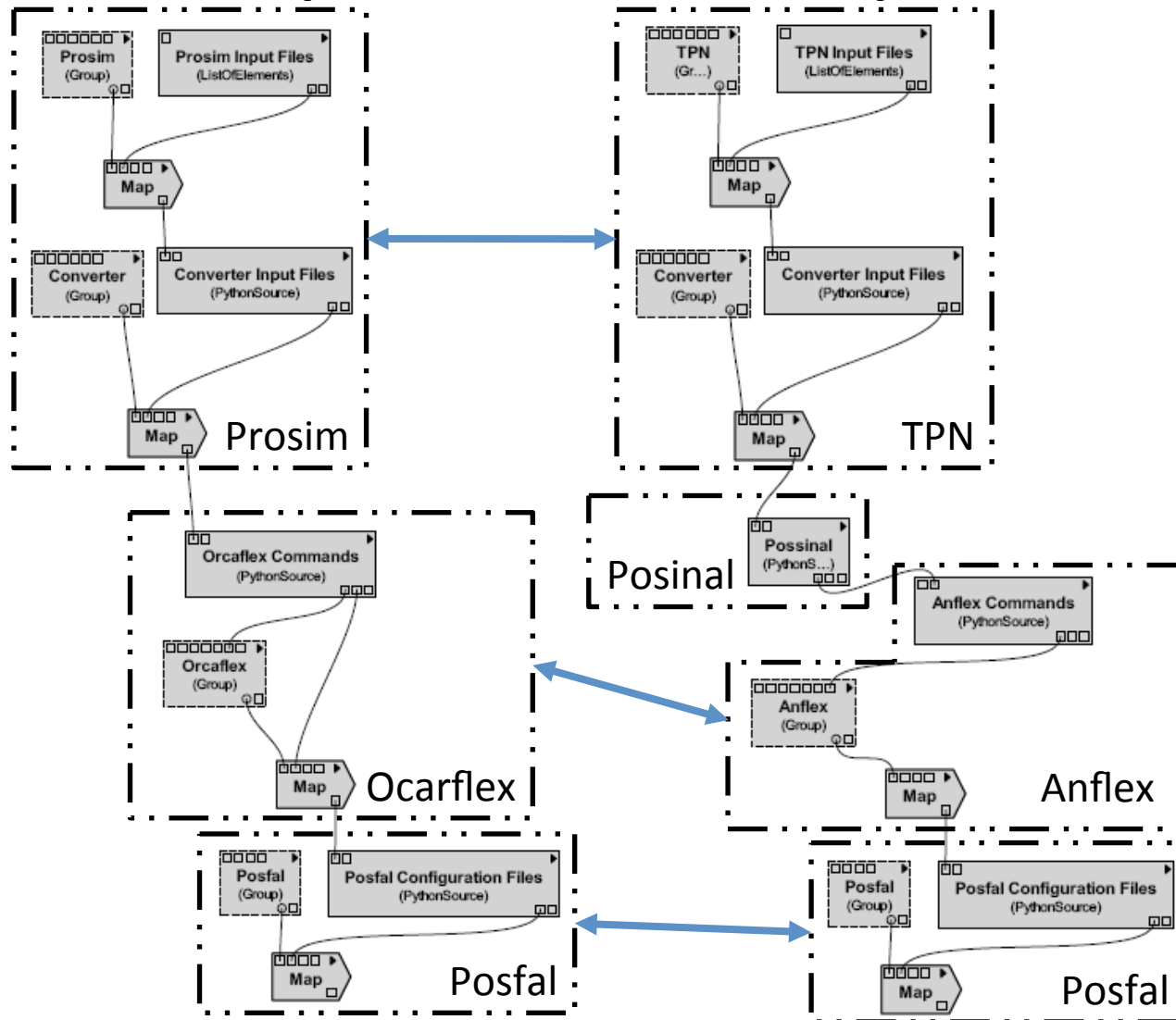
# Some aspects of Composition

# Pillars of composition

Provenance is orthogonal to those pillars and it is generated in each one of them

The composition should be supported in scientific experiments

Composition

Conception

Reuse
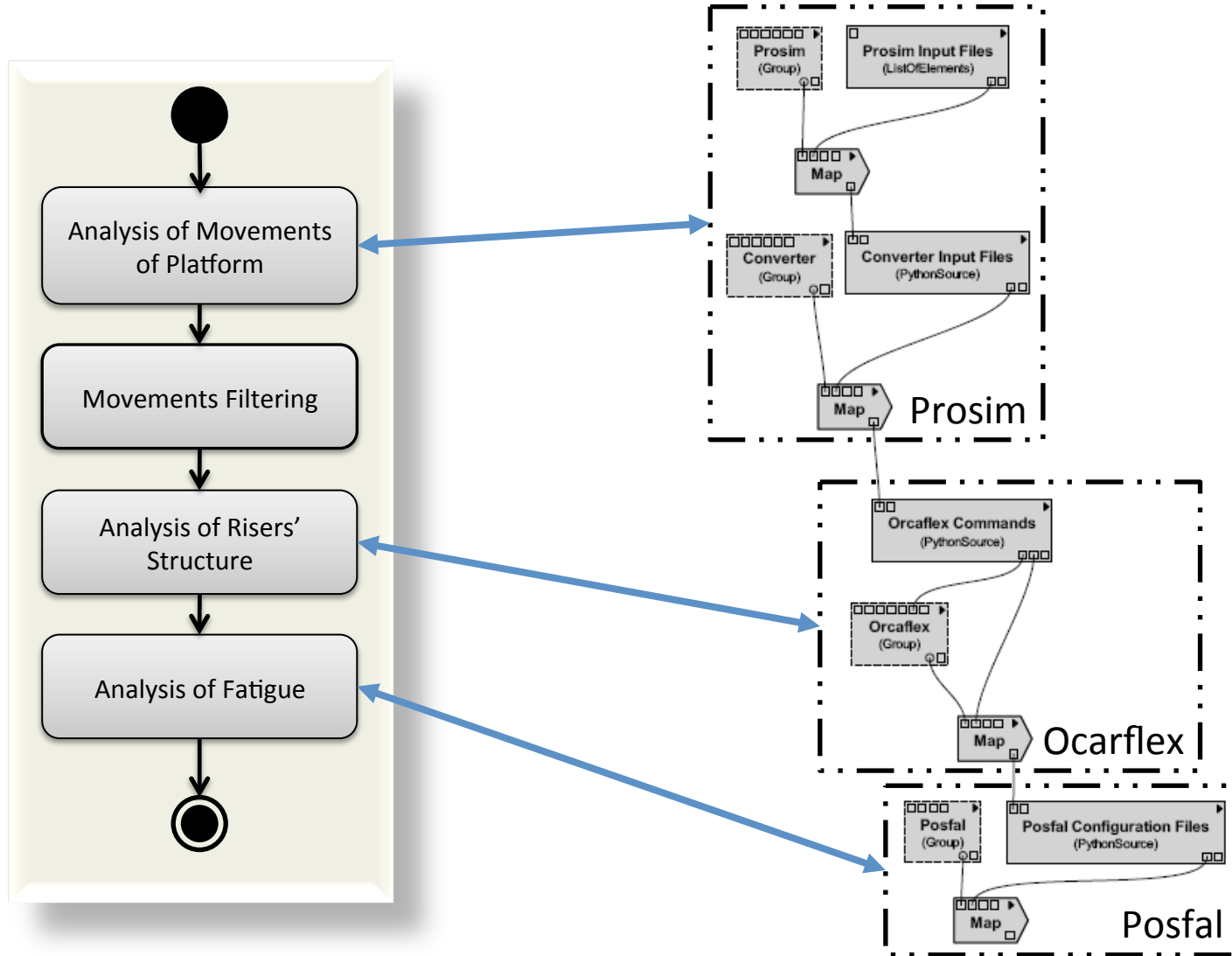
Configuration Management

Provenance

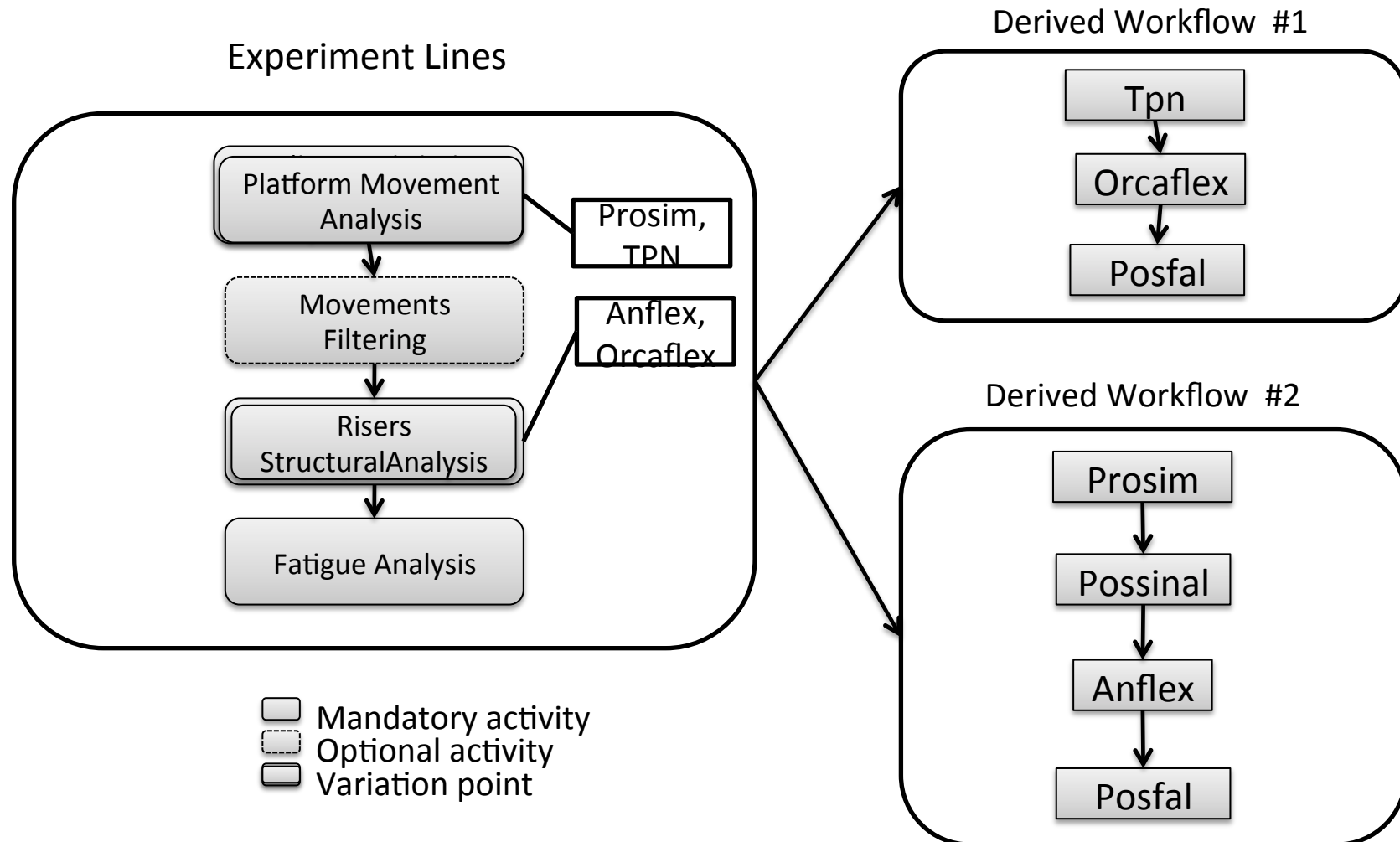# Concrete Workflows for
# Ultra Deep Water Oil Exploration



Workflow #1

Workflow #2

# Conceptual Workflows and Concrete Workflows Limitations

# Experiment Lines*- abstract workflow



Experiment Lines

Platform Movement Analysis

Prosim, TPN

Movements Filtering

Anflex, Orcaflex

Risers StructuralAnalysis

Fatigue Analysis

Mandatory activity
Optional activity
Variation point

Derived Workflow #1

Tpn
Orcaflex
Posfal

Derived Workflow #2

Prosim
Possinal
Anflex
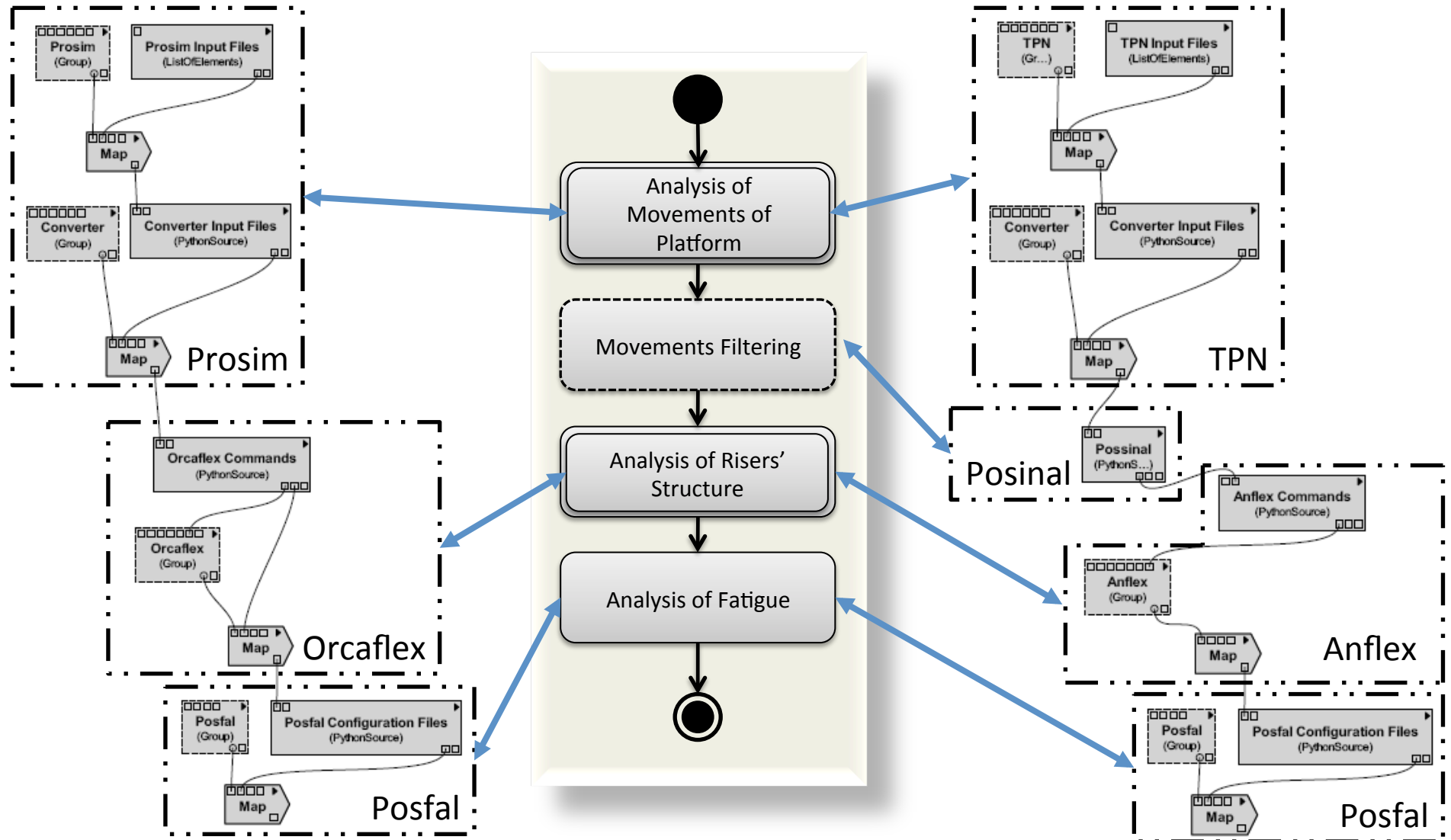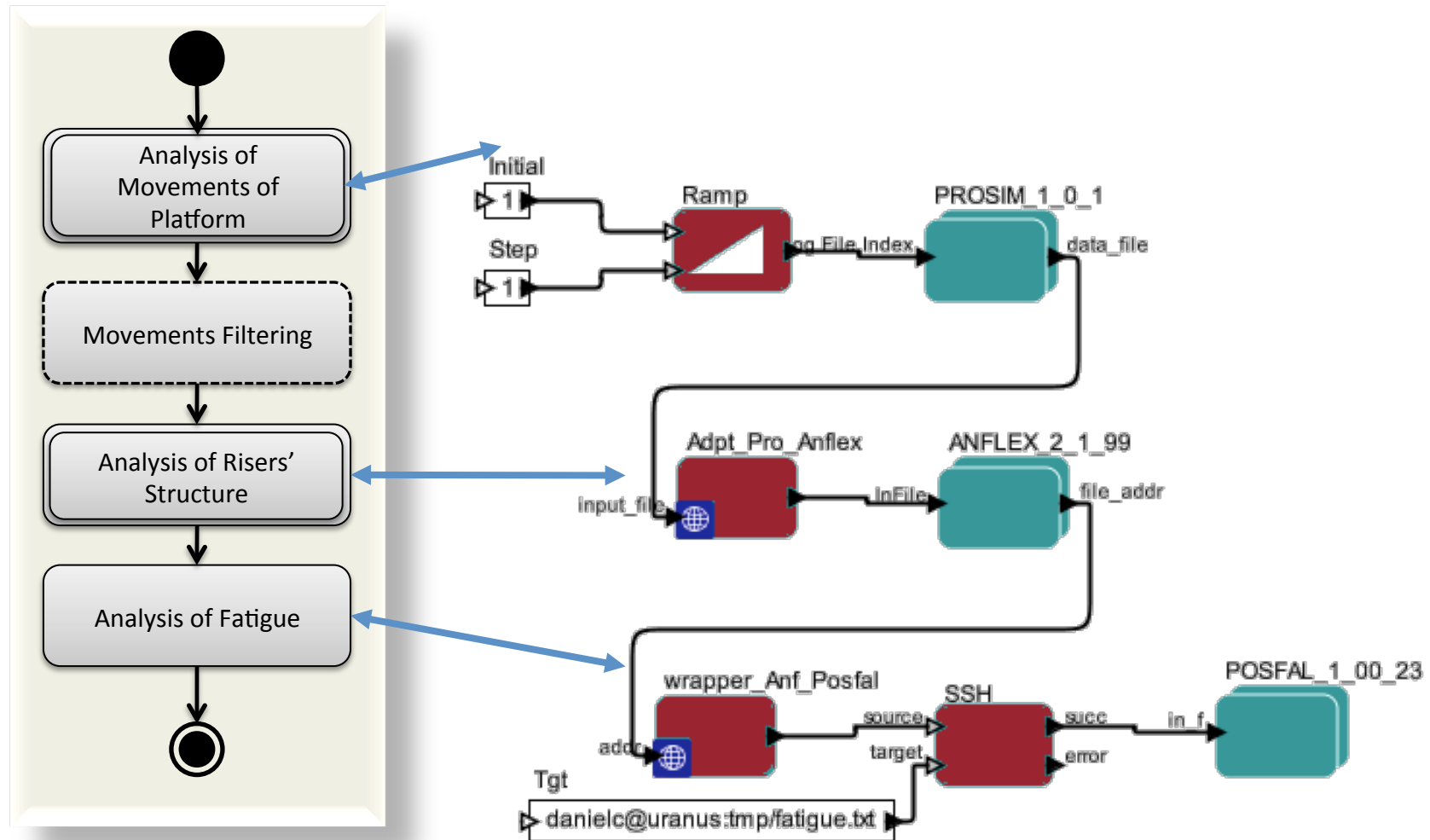Posfal

*Experiment Line: Software Reuse in Scientific Workflows. SSDBM 2009: 264-272
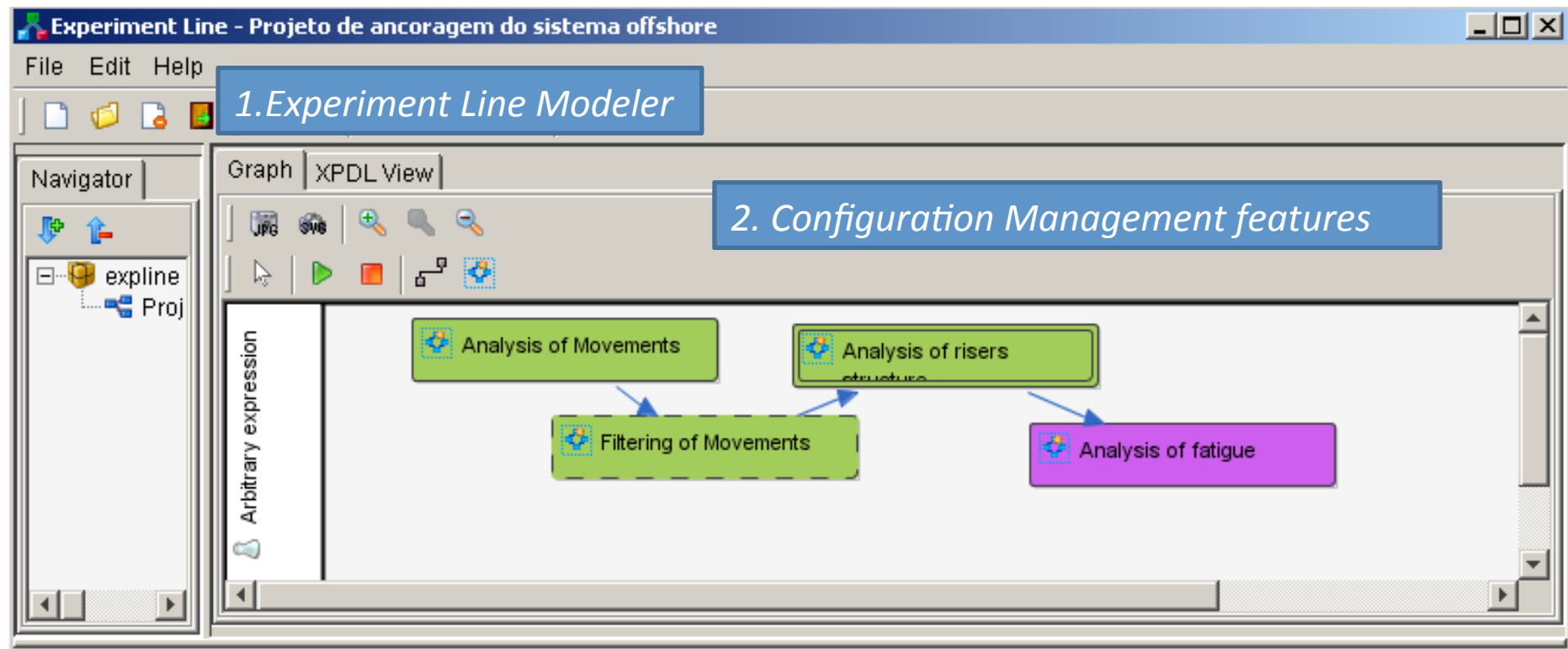
# Workflow Derivation - VisTrails

# Workflow Derivation – Kepler

# Derivation in GExpLine



**1. Experiment Line Modeler**

**2. Configuration Management features**

**3. Workflow Importing**

**4. Workflow derivation**

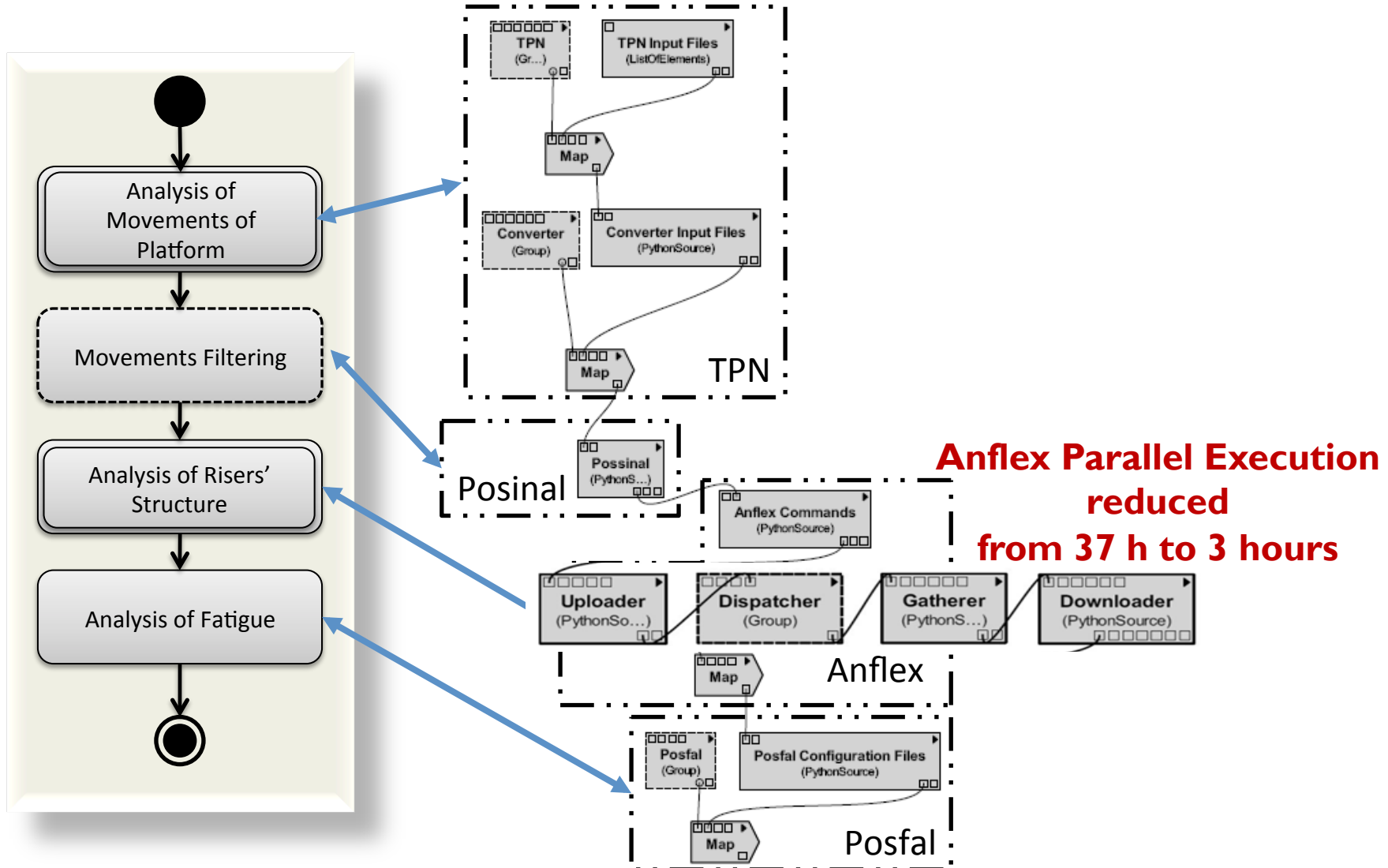**5. Prospective Provenance Querying Support**

# Derivation Process

Derive concrete workflows from a conceptual workflow

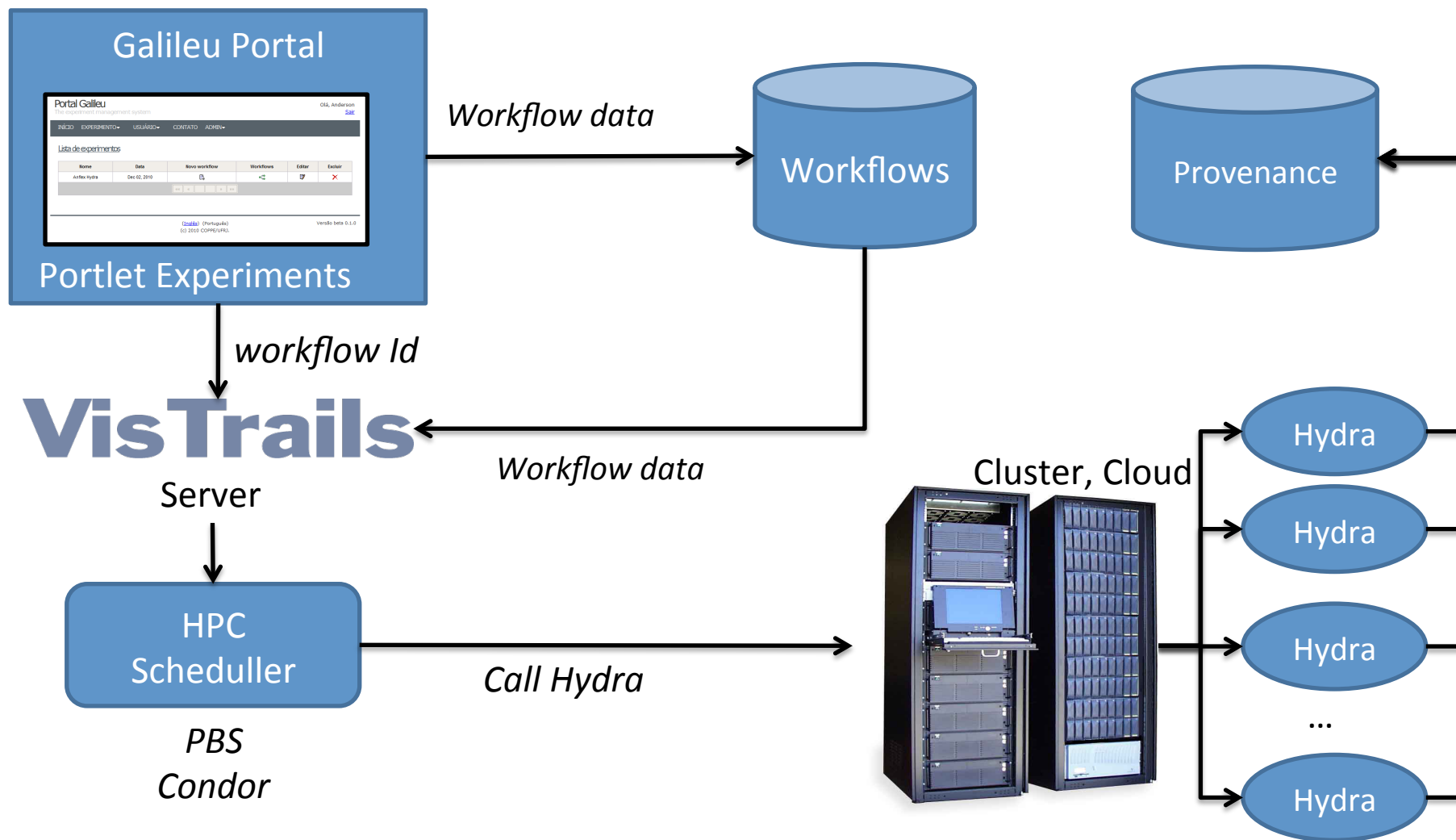Derivation information is an important provenance data

It relates all concrete workflows (trials) for a single experiment (conceptual)

# Workflow Derivation – VisTrails and HPC



**Anflex Parallel Execution reduced from 37 h to 3 hours**

# Workflow Execution

# Issues in distributed provenance

- Provenance integration (local SWfMS and HPC wf execution)

- Provenance gathering in distributed/ heterogeneous environments

- Controlling provenance from parallel execution in distributed environments

- Using provenance for steering activities in distributed environments

# Provenance can support analyzing scientific experiments

- Before execution:
  - What programs may be used? Is there any alternative methodology to explore?
  - Is there any dependency between activities? Which activities are mandatory?

- After execution:
  - What were the parameters used in the best result ?
  - What was the scientific workflow used to obtain such result?
  - Where are the output files generated by the distributed activity A using the parameters P?
  - How many times the activity A in version V was used in the experiment E?

all these queries are related to the ability of reproducing and validating a scientific experiment

# Provenance
## Exchange, Integration and Querying

Contributors:

- M. David Allen, Adriane Chapman, Barbara Blaustein, Len Seligman
  [5 Getting It Together: Enabling Multi-organization Provenance Exchange]

- Anderson Marinho, Marta Mattoso, Claudia Werner, Vanessa Braganholo and Leonardo Murta
  [33 Challenges in managing implicit and abstract provenance data: experiences with ProvManager]

- Luiz M. R. Gadelha Jr., Marta Mattoso, Michael Wilde, Ian Foster
  [26 Provenance Query Patterns for Many-Task Scientific Computing]

# Provenance Exchange, Integration and Querying

**Marta Mattoso**

**Federal University of Rio de Janeiro, Brazil**