

Advanced Data Management (CSCI 640/490)

Provenance

Dr. David Koop

Sharing Data

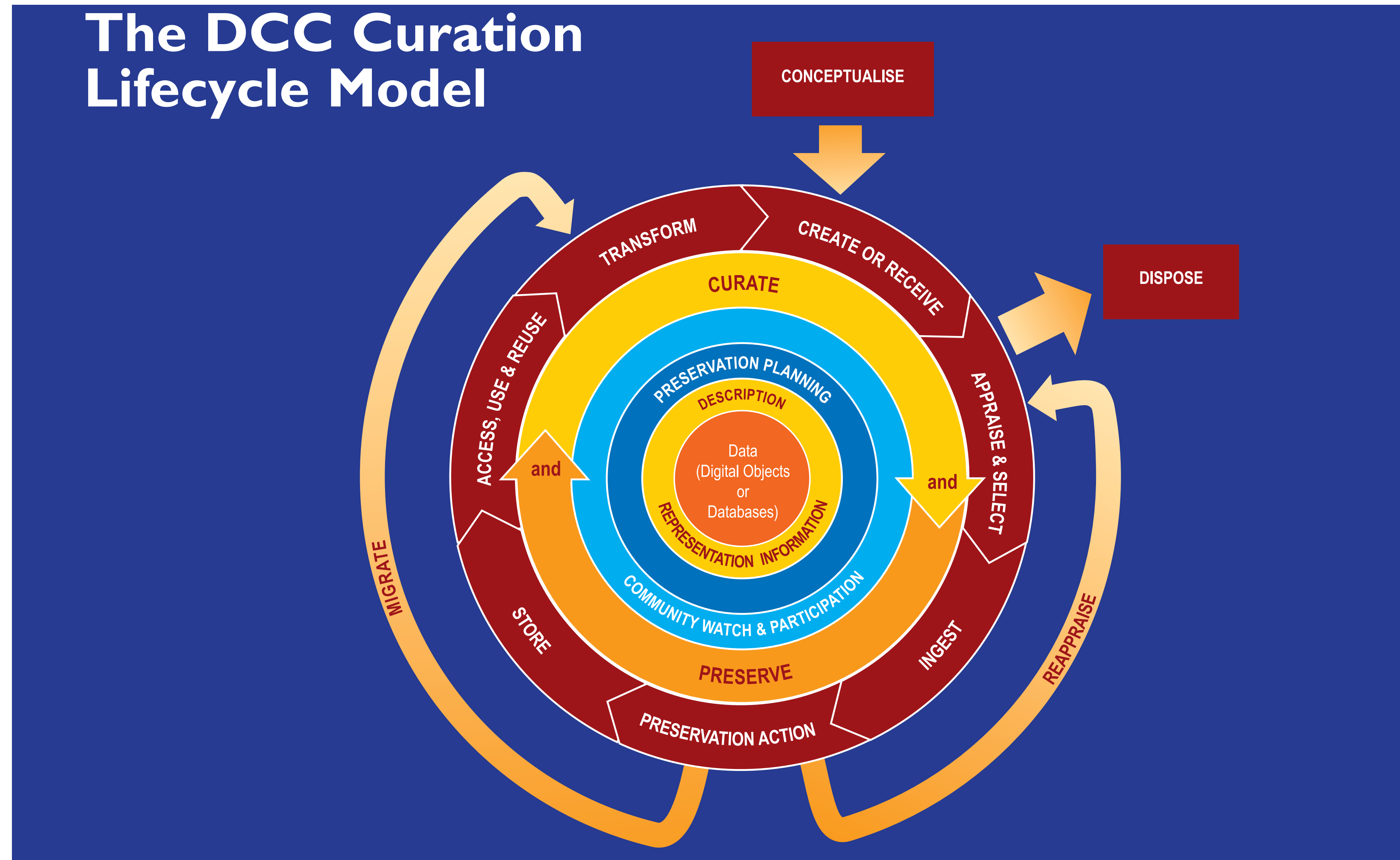
- Required/encouraged by universities, funding agencies, publishers
- "Publications are arguments made by authors, and **data are the evidence** used to support the arguments." [C. L. Borgman]
- Questions:
 - How is data maintained? Who is responsible?
 - What is the process for curating data?
 - How long should data be kept?
 - How should data collection and curation be acknowledged?

Research Data Infrastructure Stakeholders

- Research Funding Agencies
- Individual Scientists and Scholars
 - Data collection/analysis, managing teams/technology
- Academic Institutions
 - Academic Leadership: Regulations, Governance, Financial Management
 - Research Computing
 - University Libraries: Maintain knowledge resources, provide access, steward
 - Schools and Departments

[C. L. Borgman & P. E. Bourne]

Data Curation Lifecycle



[DCC]

Sequential Actions in Data Curation

- Conceptualize: Plan creation of data—capture method and storage options.
- Create or Receive: Create/receive data and make sure metadata exists
- Appraise and Select: Evaluate data and select for long-term curation and preservation
- Ingest: Transfer data to an archive, repository, data centre or other custodian
- Preservation Action: Data cleaning, validation (ensure that data remains authentic, reliable and usable)
- Store: Store the data in a secure manner adhering to relevant standards
- Access, Use and Reuse: Make sure is accessible to users and reusers
- Transform: Create new data from the original (migrate formats, subsets, etc.)

[DCC]

FAIR Principles

- **Findable:** Metadata and data should be easy to find for both humans and computers
- **Accessible:** Users need to know how data can be accessed, possibly including authentication and authorization
- **Interoperable:** Can be integrated with other data, and can interoperate with applications or workflows for analysis, storage, and processing
- **Reusable:** Optimize the reuse of data. Metadata and data should be well-described so they can be replicated and/or combined in different settings

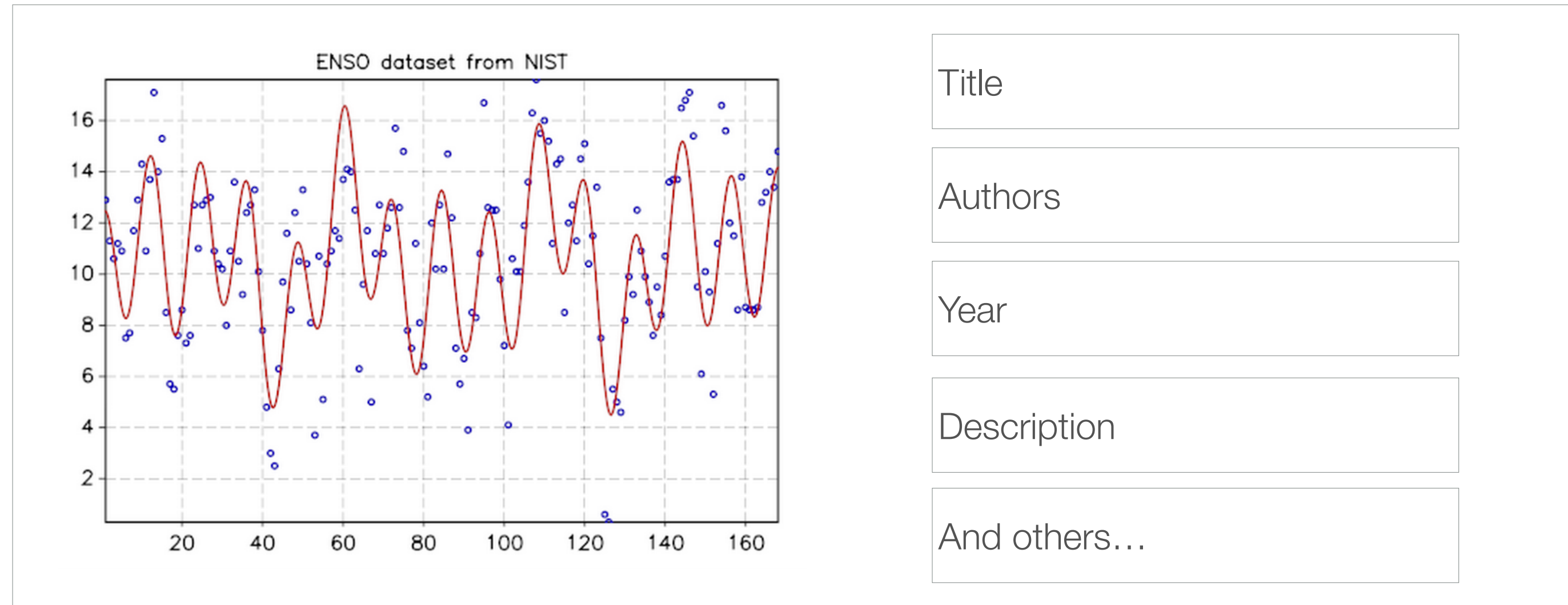
[\[GO FAIR\]](#)

Findable: DataCite Workflow

1. Take a dataset

2. Describe it

3. Assign a DOI

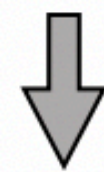


10.1234/exampledata

Proxy

Prefix

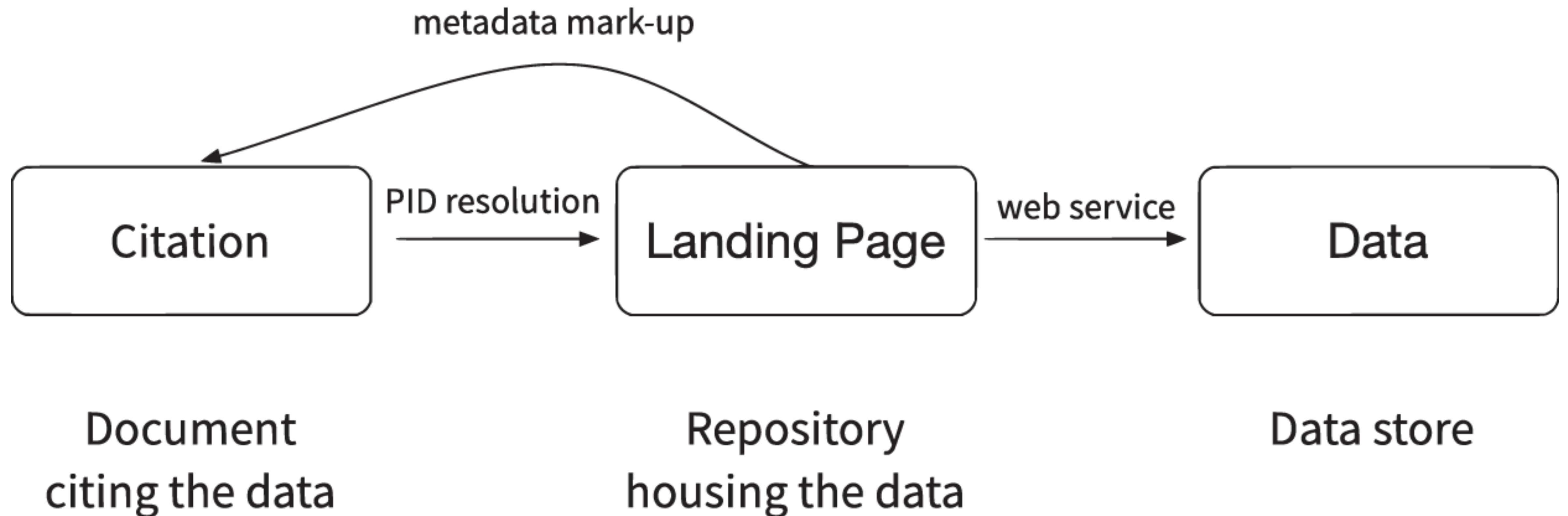
Suffix



<https://doi.org/10.5438/n138-z3mk>

[DataCite]

Accessible: DOI to Landing Page with Metadata



[M. Fenner et al., 2019]

Interoperable: Standard vocabularies

View as TableView as Grid

Sort byName

Recommended Records

Recommended

Associated Publication?

No PublicationHas Publication

Claimed?

No MaintainerHas Maintainer

Record Status

UncertainDeprecatedIn developmentReady

Standard Type

Terminology Artifact771

Model/Format405

Reporting Guideline163

Metric30

Identifier Schema15

Show More

Domains

Report141

Data Transformation134

Showing records 1 - 50 of 1384.

«12345678910111213141516171819202122232425262728»

Registry	Name	Abbreviation	Type	Subject	Domain	Taxonomy	Related Database	Related Standard	Related Policy	In Collection/Recommendation	Status
	ABA Adult Mouse Brain	ABA	Standard	Neuroscience	BrainGene ExpressionBrain Imaging	Mus musculus	NeuroMorpho.Org	None	None	None	R
	Access to Biological Collection Data	ABCD	Standard	BiodiversityBiologyLife Science	None	All	GBIFALA IPT - GBIF Australia RepositoryGBIF Spain IPT - GBIF Spain RepositoryCanadensys IPT - GBIF Canadensys RepositorySiB Colombia IPT - GBIF Colombia RepositoryPlus 1 more...	ABCDDNAABCDEF	None	TDWG Biodiversity Information Standards	R
	Access to Biological Collection Databases Extended for Geosciences	ABCDEF	Standard	Earth ScienceGeologyPaleontologySoil Science	None	All	GeoCAsE Data Portal	XMLABCD	None	None	R
	Access to Biological Collection Data DNA extension	ABCDDNA	Standard	BiodiversityBiologyLife Science	DNA Sequence DataExperiment MetadataSequenceDeoxyribonucleic AcidPolymerase Chain ReactionPlus 1 more...	All	GenBank	MOD-COABCD	None	TDWG Biodiversity Information Standards	Dev
	.ACE format	.ACE format	Standard	Life Science	DNA Sequence DataContigDeoxyribonucleic AcidGenome	All	None	None	None	None	R
	AdaLab-meta ontology	ADALAB-META	Standard	None	None	All	None	None	None	None	R
	AdaLab ontology	ADALAB	Standard	None	None	All	None	None	None	None	R
	Adverse Drug Reaction Markup Language	EU-ADR ML	Standard	None	Adverse ReactionElectronic Health Record	Homo sapiens	None	XML	None	None	U

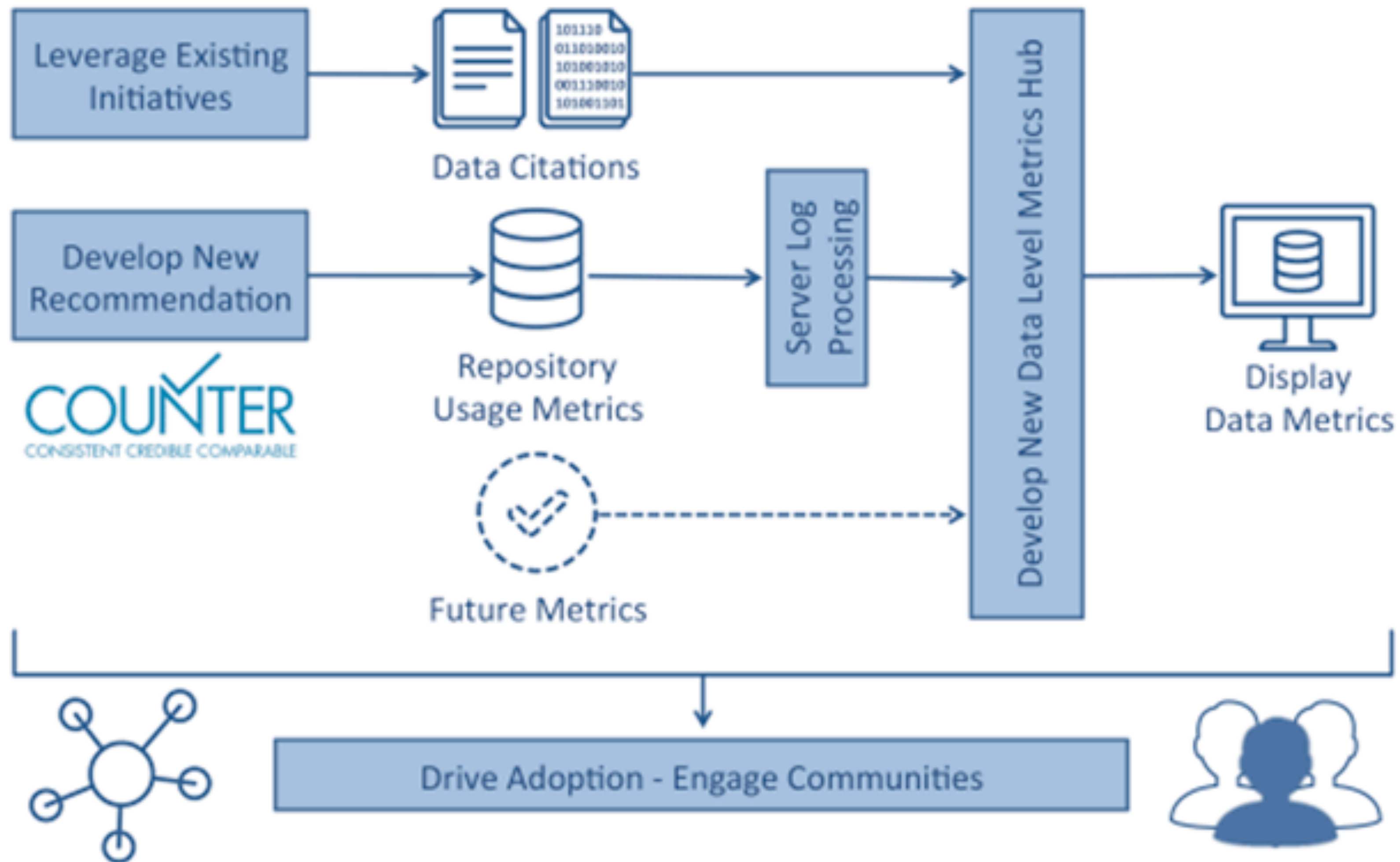
[fairsharing.org]

Reusable: Licensing

- Citation of a dataset is expected as a scholarly norm, not by law
- CC0:
 - "I hereby waive all copyright and related or neighboring rights together with all associated claims and causes of action with respect to this work to the extent possible under the law"
- CC BY: license, not a waiver as CC0
 - "You must give appropriate credit, provide a link to the license, and indicate if changes were made."
- Data Use Agreements (DUA): Used when data are restricted due to proprietary or privacy concerns.

[M. Crosas]

Reusable: Data Citation & Metrics



[H. Cousijn et al., 2019]

Assignment 5

- FAA & ADS-B aircraft data
- Spatial data processing, visualization, time series
- Due at the end of the semester

geopandas example

Provenance

What actually happened in a computational experiment?

Provenance in Art



Rembrandt van Rijn

Dutch, 1606 - 1669

Self-Portrait, 1659

oil on canvas

Andrew W. Mellon Collection

1937.1.72

Provenance

George, 3rd Duke of Montagu and 4th Earl of Cardigan [d. 1790], by 1767;[1] by inheritance to his daughter, Lady Elizabeth, wife of Henry, 3rd Duke of Buccleuch of Montagu House, London; John Charles, 7th Duke of Buccleuch; (P. & D. Colnaghi & Co., New York, 1928); (M. Knoedler & Co., New York); sold January 1929 to Andrew W. Mellon, Pittsburgh and Washington, D.C.; deeded 28 December 1934 to The A.W. Mellon Educational and Charitable Trust, Pittsburgh; gift 1937 to NGA.

[1] This early provenance is established by presence of a mezzotint after the portrait by R. Earlom (1743-1822), dated 1767. See John Charrington, *A Catalogue of the Mezzotints After, or Said to Be After, Rembrandt*, Cambridge, 1923, no. 49.

Associated Names

- Buccleuch, Henry, 3rd Duke of
- Buccleuch, John Charles, 7th Duke of
- Colnaghi & Co., Ltd., P. & D.
- Knoedler & Company, M.
- Mellon, Andrew W.
- Mellon Educational and Charitable Trust, The A.W.
- Montagu, and 4th Earl of Cardigan, George, 3rd Duke of

[National Gallery of Art]

Provenance in Art



Rembrandt van Rijn

Dutch, 1606 - 1669

Self-Portrait, 1659

oil on canvas

Andrew W. Mellon Collection

1937.1.72

Provenance

George, 3rd Duke of Montagu and 4th Earl of Cardigan [d. 1790], by 1767;[1] by inheritance to his daughter, Lady Elizabeth, wife of Henry, 3rd Duke of Buccleuch of Montagu House, London; John Charles, 7th Duke of Buccleuch; (P. & D. Colnaghi & Co., New York, 1928); (M. Knoedler & Co., New York); sold January 1929 to Andrew W. Mellon, Pittsburgh and Washington, D.C.; deeded 28 December 1934 to The A.W. Mellon Educational and Charitable Trust, Pittsburgh; gift 1937 to NGA.

[1] This early provenance is established by presence of a mezzotint after the portrait by R. Earlom (1743-1822), dated 1767. See John Charrington, *A Catalogue of the Mezzotints After, or Said to Be After, Rembrandt*, Cambridge, 1923, no. 49.

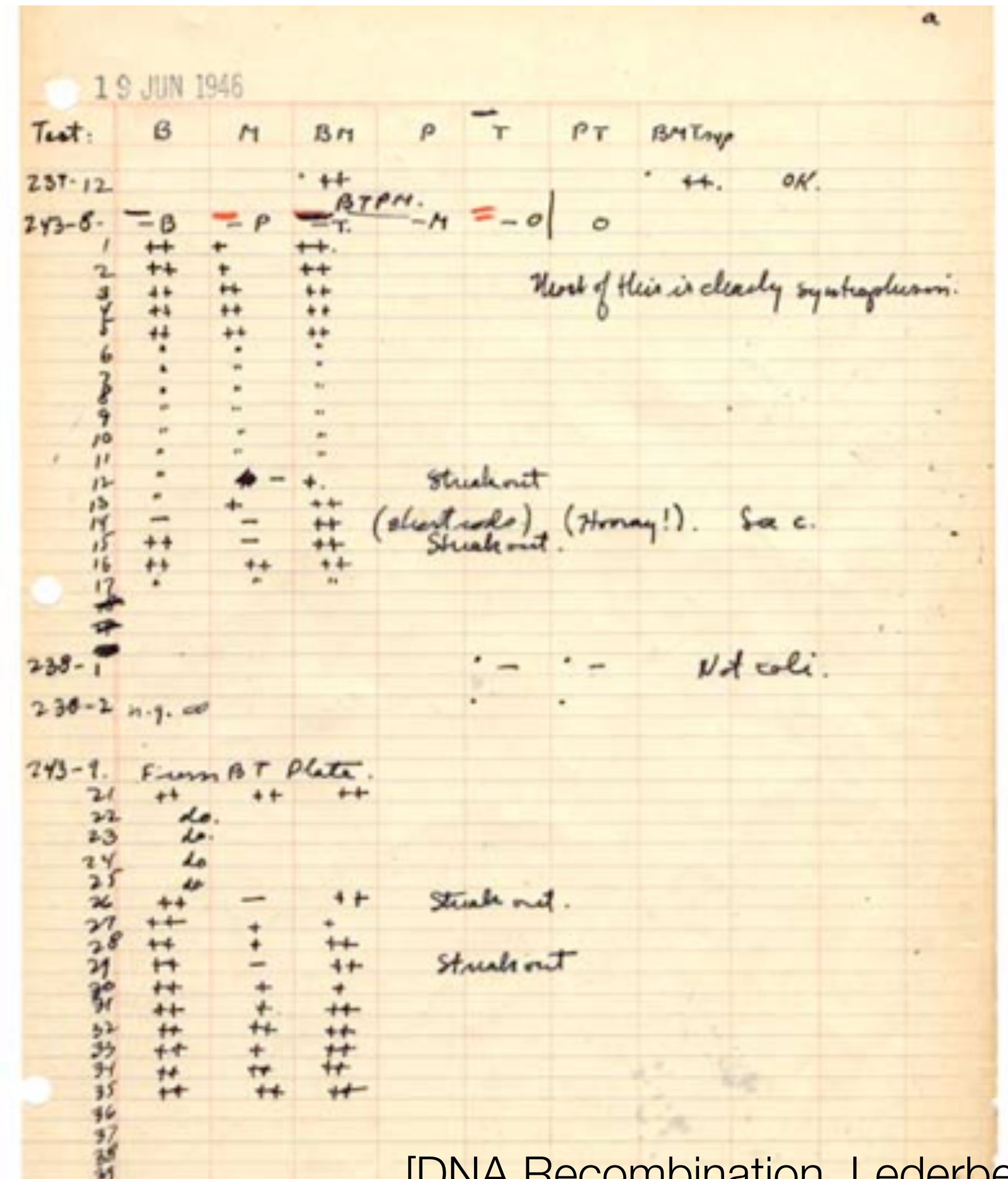
Associated Names

- Buccleuch, Henry, 3rd Duke of
- Buccleuch, John Charles, 7th Duke of
- Colnaghi & Co., Ltd., P. & D.
- Knoedler & Company, M.
- Mellon, Andrew W.
- Mellon Educational and Charitable Trust, The A.W.
- Montagu, and 4th Earl of Cardigan, George, 3rd Duke of

[National Gallery of Art]

Provenance in Science

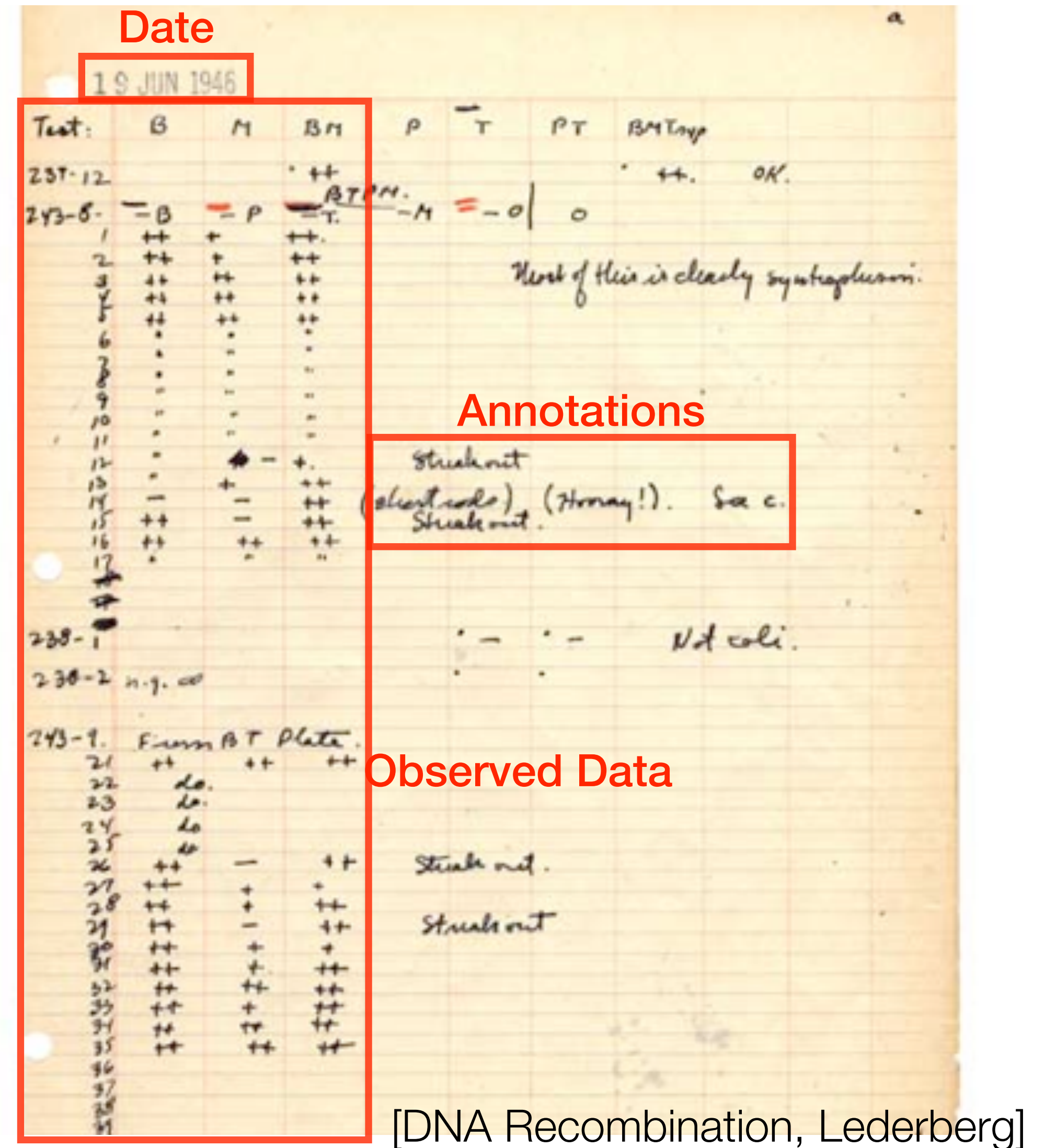
- Provenance: the lineage of data, a computation, or a visualization
- **Provenance is as (or more) important as the result!**
- Old solution:
 - Lab notebooks
- New problems:
 - Large volumes of data
 - Complex analyses
 - Writing notes doesn't scale



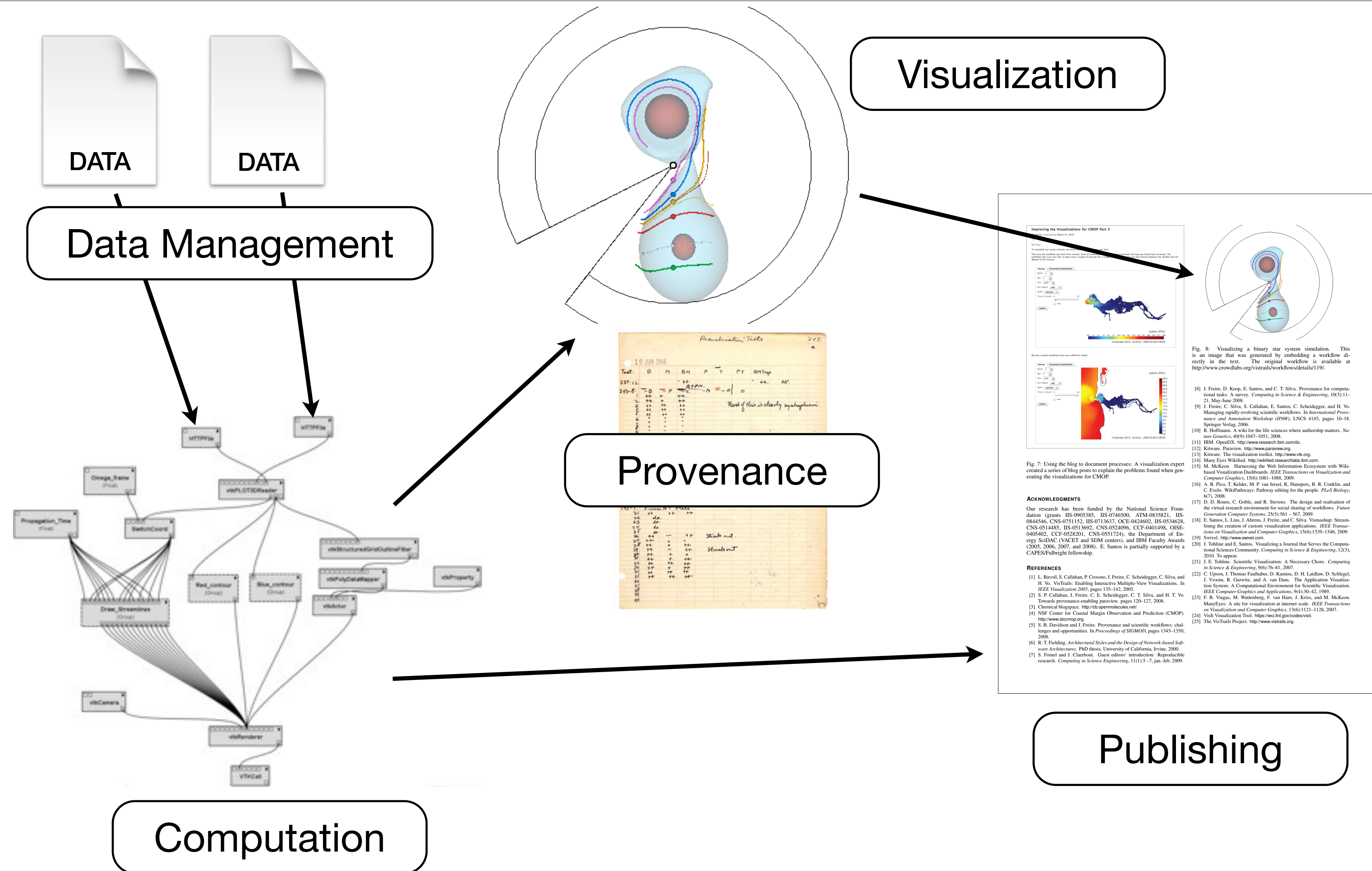
[DNA Recombination, Lederberg]

Provenance in Science

- Provenance: the lineage of data, a computation, or a visualization
- **Provenance is as (or more) important as the result!**
- Old solution:
 - Lab notebooks
- New problems:
 - Large volumes of data
 - Complex analyses
 - Writing notes doesn't scale



Provenance in Computational Science



Evolution of Publication

- Publish paper
- Publish code
- Publish computational experiments/tests
- Publish provenance (what actually happens during your runs)

Galois Conjugates of Topological Phases

M. H. Freedman,¹ J. Gukelberger,² M. B. Hastings,¹ S. Trebst,¹ M. Troyer,² and Z. Wang¹

¹Microsoft Research, Station Q, University of California, Santa Barbara, CA 93106, USA

²Theoretische Physik, ETH Zurich, 8093 Zurich, Switzerland

(Dated: July 6, 2011)

Galois conjugation relates unitary conformal field theories (CFTs) and topological quantum field theories (TQFTs) to their non-unitary counterparts. Here we investigate Galois conjugates of quantum double models, such as the Levin-Wen model. While these Galois conjugated Hamiltonians are typically non-Hermitian, we find that their ground state wave functions still obey a generalized version of the usual code property (local operators do not act on the ground state manifold) and hence enjoy a generalized topological protection. The key question addressed in this paper is whether such non-unitary topological phases can also appear as the ground states of Hermitian Hamiltonians. Specific attempts at constructing Hermitian Hamiltonians with these ground states lead to a loss of the code property and topological protection of the degenerate ground states. Beyond this we rigorously prove that no local change of basis (IV.5) can transform the ground states of the Galois conjugated doubled Fibonacci theory into the ground states of a topological model whose Hermitian Hamiltonian satisfies Lieb-Robinson bounds. These include all gapped local or quasi-local Hamiltonians. A similar statement holds for many other non-unitary TQFTs. One consequence is that the “Gaffnian” wave function cannot be the ground state of a gapped fractional quantum Hall state.

PACS numbers: 05.30.Pr, 73.43.-f

I. INTRODUCTION

Galois conjugation, by definition, replaces a root of a polynomial by another one with identical algebraic properties. For example, i and $-i$ are Galois conjugate (consider $z^2 + 1 = 0$) as are $\phi = \frac{1+\sqrt{5}}{2}$ and $-\frac{1}{\phi} = \frac{1-\sqrt{5}}{2}$ (consider $z^2 - z - 1 = 0$), as well as $\sqrt[3]{2}$, $\sqrt[3]{2}e^{2\pi i/3}$, and $\sqrt[3]{2}e^{-2\pi i/3}$ (consider $z^3 - 2 = 0$). In physics Galois conjugation can be used to convert non-unitary conformal field theories (CFTs) to unitary ones, and vice versa. One famous example is the non-unitary Yang-Lee CFT, which is Galois conjugate to the Fibonacci CFT $(G_2)_1$, the even (or integer-spin) subset of $\text{su}(2)_3$.

In statistical mechanics non-unitary conformal field theories have a venerable history.^{1,2} However, it has remained less clear if there exist physical situations in which non-unitary models can provide a useful description of the low energy physics of a quantum mechanical system – after all, Galois conjugation typically destroys the Hermitian property of the Hamiltonian. Some non-Hermitian Hamiltonians, which surprisingly have totally real spectrum, have been found to arise in the study of PT -invariant one-particle systems³ and in some Galois conjugate many-body systems⁴ and might be seen to open the door a crack to the physical use of such models. Another situation, which has recently attracted some interest, is the question whether non-unitary models can describe 1D edge states of certain 2D bulk states (the edge holographic for the bulk). In particular, there is currently a discussion on whether or not the “Gaffnian” wave function could be the ground state for a *gapped* fractional quantum Hall (FQH) state albeit with a non-unitary “Yang-Lee” CFT describing its edge.⁵⁻⁷ We conclude that this is not possible, further restricting the possible scope of non-unitary models in quantum mechanics.

We reach this conclusion quite indirectly. Our main thrust is the investigation of Galois conjugation in the simplest non-

Abelian Levin-Wen model.⁸ This model, which is also called “DFib”, is a topological quantum field theory (TQFT) whose states are string-nets on a surface labeled by either a trivial or “Fibonacci” anyon. From this starting point, we give a rigorous argument that the “Gaffnian” ground state cannot be locally conjugated to the ground state of any topological phase, within a Hermitian model satisfying Lieb-Robinson (LR) bounds⁹ (which includes but is not limited to gapped local and quasi-local Hamiltonians).

Lieb-Robinson bounds are a technical tool for local lattice models. In relativistically invariant field theories, the speed of light is a strict upper bound to the velocity of propagation. In lattice theories, the LR bounds provide a similar upper bound by a velocity called the LR velocity, but in contrast to the relativistic case there can be some exponentially small “leakage” outside the light-cone in the lattice case. The Lieb-Robinson bounds are a way of bounding the leakage outside the light-cone. The LR velocity is set by microscopic details of the Hamiltonian, such as the interaction strength and range. Combining the LR bounds with the spectral gap enables us to prove locality of various correlation and response functions. We will call a Hamiltonian a *Lieb-Robinson Hamiltonian* if it satisfies LR bounds.

We work primarily with a single example, but it should be clear that the concept of Galois conjugation can be widely applied to TQFTs. The essential idea is to retain the particle types and fusion rules of a unitary theory but when one comes to writing down the algebraic form of the F -matrices (also called $6j$ symbols), the entries are now Galois conjugated. A slight complication, which is actually an asset, is that writing an F -matrix requires a gauge choice and the most convenient choice may differ before and after Galois conjugation.

Our method is not restricted to Galois conjugated DFib^G and its factors Fib^G and $\overline{\text{Fib}}^G$, but can be generalized to infinitely many non-unitary TQFTs, showing that they will not arise as low energy models for a gapped 2D quantum mechan-

non-Hermitian DYL model

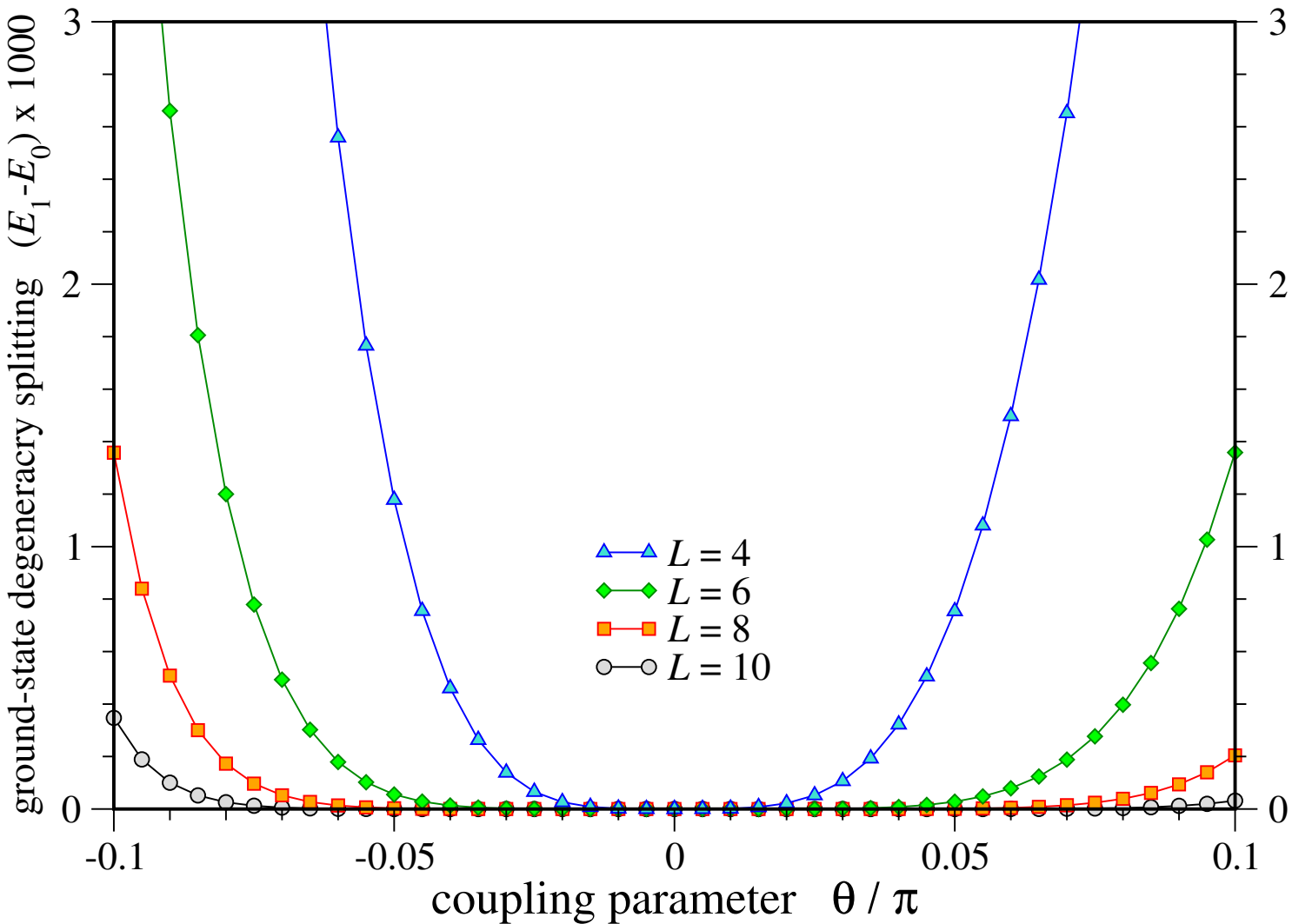


FIG. 6. (color online) Ground-state degeneracy splitting of the non-Hermitian doubled Yang-Lee model when perturbed by a string tension ($\theta \neq 0$).

[Freedman et al., 2012]

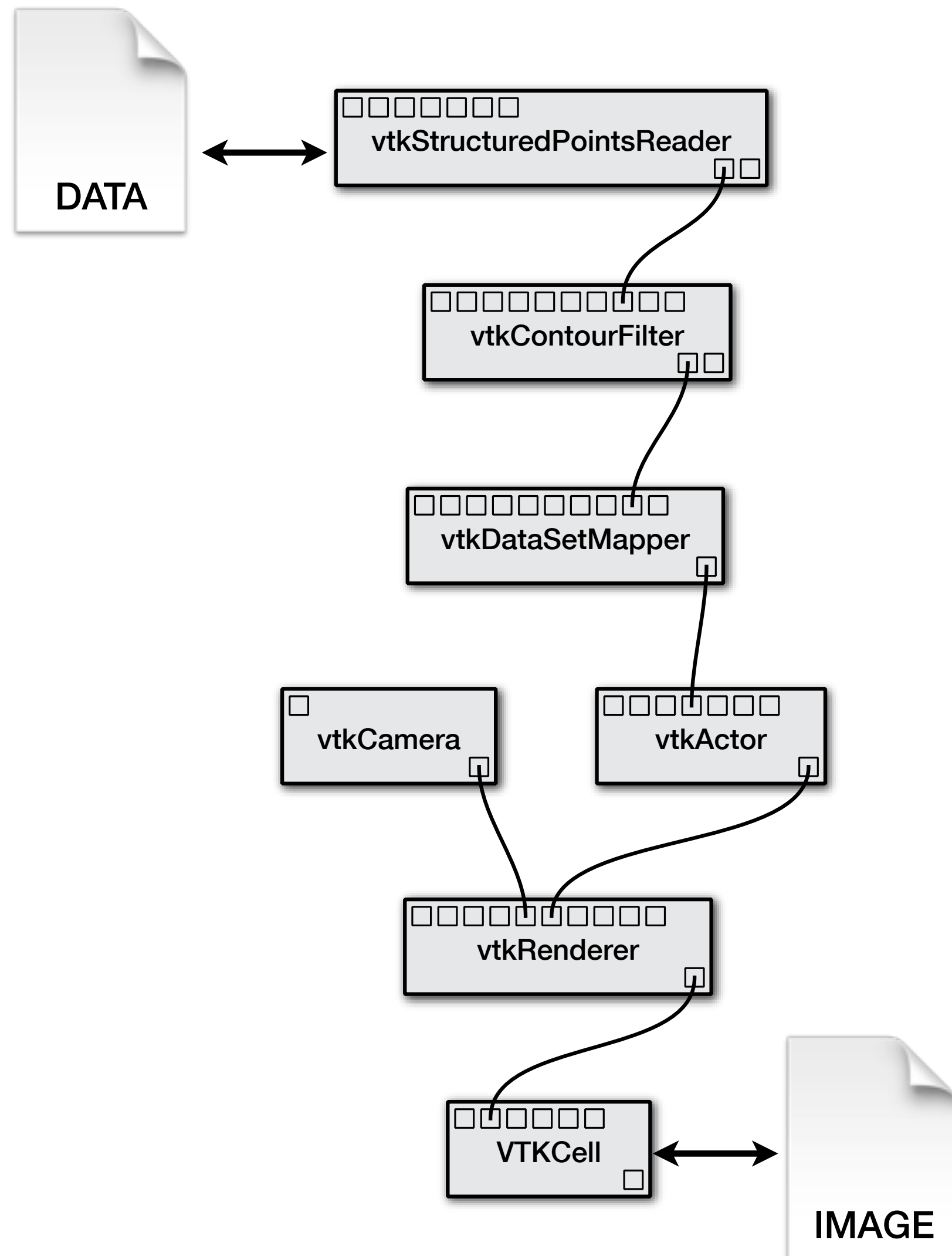
Benefits of Provenance-Rich Publications

- Produce more knowledge—not just text
- Allow scientists to stand on the shoulders of giants (and their own)
- Science can move faster!
- Higher-quality publications
- Authors will be more careful
- Many eyes to check results
- Describe more of the discovery process: people only describe successes, can we learn from mistakes?
- Expose users to different techniques and tools: expedite their training; and potentially reduce their time to insight

Provenance Definitions

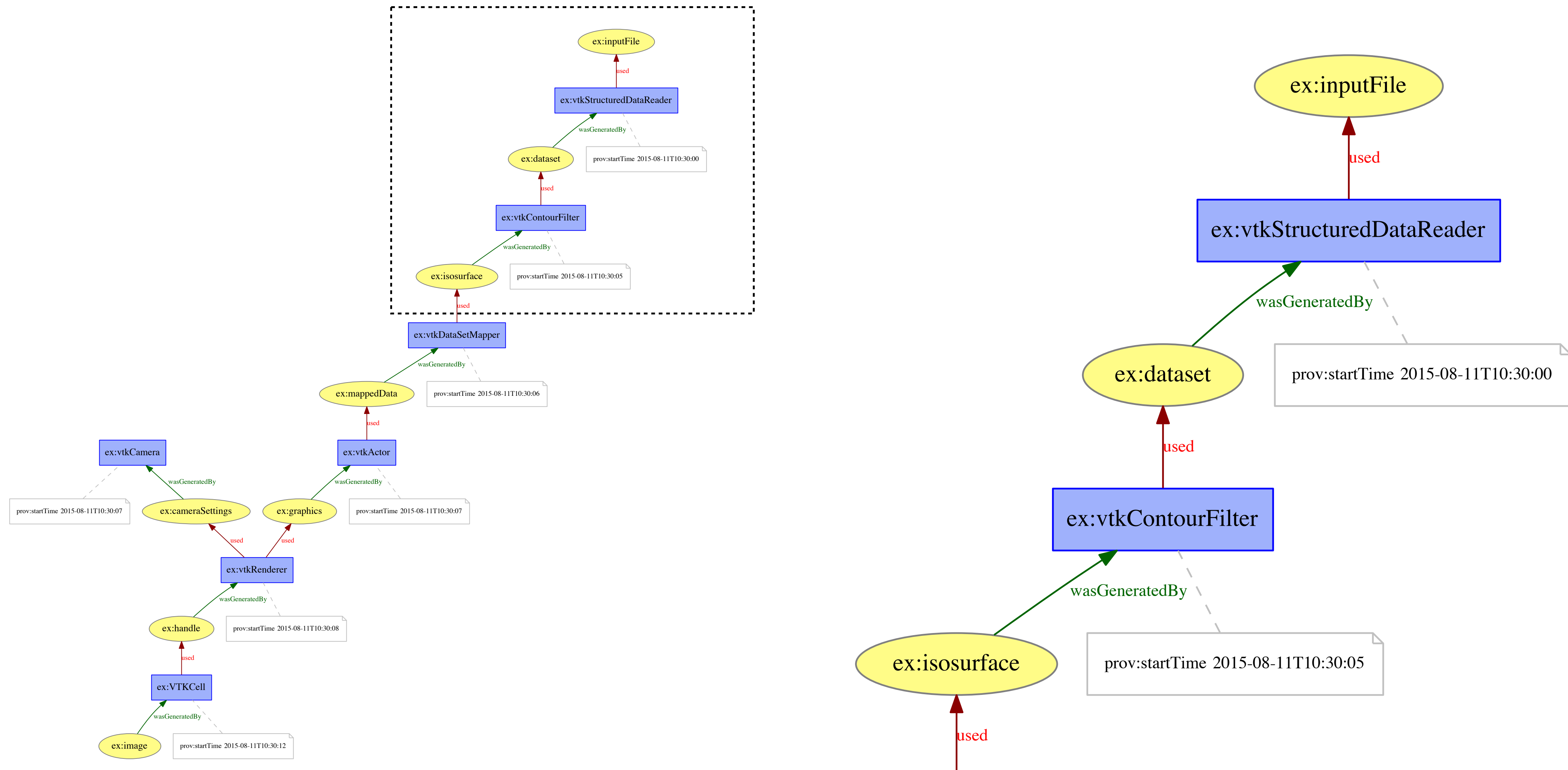
- Dictionary: "the source or origin of an object; its history and pedigree; a record of the ultimate derivation and passage of an item through its various owners."
- Focus on **causality**—the sequence of steps that detail how a result was generated and/or **derivation**—what data a result depended on
- Provenance itself is **data**, this list of steps along with metadata for each step: when it occurred, who initiated it, notes about it
- Can be used to preserve information about an experiment and to answer many questions

Workflows

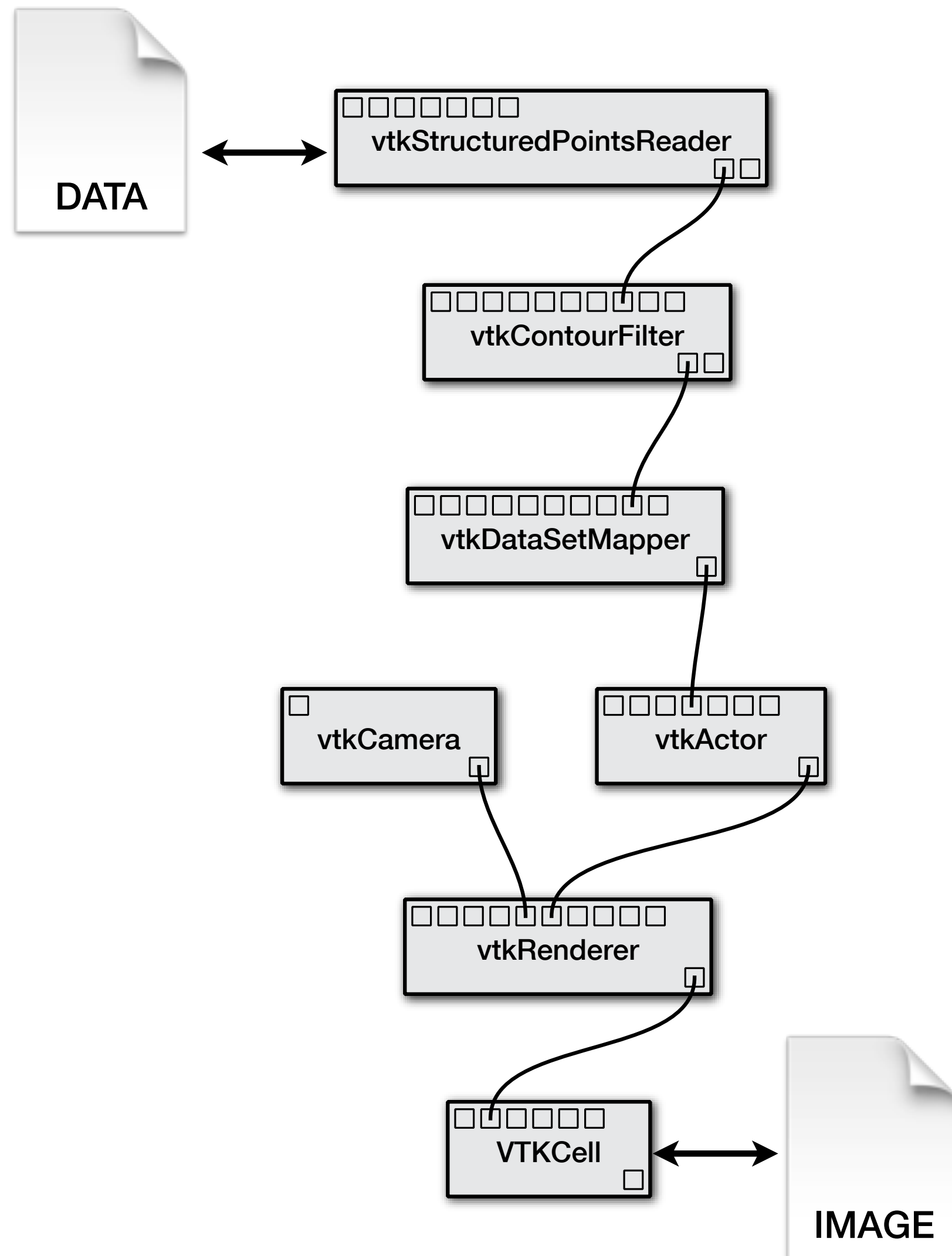


- Abstract computation
- Computational modules connected through input and output ports
- Data flows along the connections

Provenance Graph



Provenance Questions



- What process led to the output image?
- What input datasets contributed to the output image?
- What workflows create an isosurface with isovalue 57?
- Who create this data product?
- When was this data file created?
- Why was `vtkCamera` used?
- Why do two output images differ?

Questions about Provenance

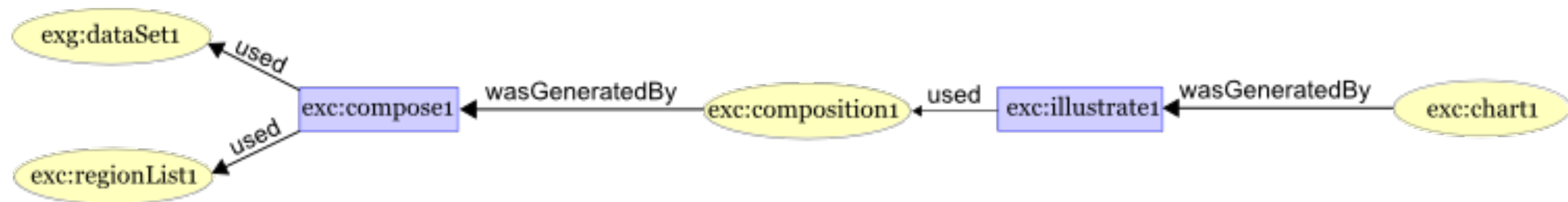
- How does one capture provenance?
- How does one manage provenance for later use?
- How do we answer questions about our provenance?
- How do we use provenance for good?

Provenance Management

- Provenance can be generated from tasks/programs/scripts/etc.
- Properties of provenance are related to the **computational model**
 - a specific application with a graphical interface
 - a script that automates the use of several command-line tools
 - a scientific workflow that combines several tools

Provenance & Causality

- Knowing what data/steps influenced other data/steps is important!
- Data dependencies: this output file depended on this input file
- Data-process dependencies: this output figure depended on these processes
- Causality can often be represented as a **graph** where connections represent dependencies



User-defined provenance

- Goal: capture lots of provenance automatically based on what steps are executed
- Problem: not everything can be captured automatically
- Annotations offer ability to keep notes about processes
- Users might also specify known causal links that cannot be automatically determined (e.g. a step depends on three system files that were not specified as inputs in the workflow)

Provenance Management

- What is needed to capture, store, and use provenance?
 1. Capture mechanism
 2. Model for representing provenance
 3. Tools to store, query, and analyze provenance

Provenance Capture Mechanisms

- **Workflow-based:** Since workflow execution is controlled, keep track of all the workflow modules, parameters, etc. as they are executed
- **Process-based:** Each process is required to write out its own provenance information (not centralized like workflow-based)
- **OS-based:** The OS or filesystem is modified so that any activity it does it monitored and the provenance subsystem organizes it
- Tradeoffs:
 - Workflow- and process-based have better abstraction
 - OS-based requires minimal user effort once installed and can capture "hidden dependencies"

Provenance Granularity

- How detailed should our provenance be?
 - **Coarse**: "This program ran with inputs x, y, z and produced outputs a, b, c"
 - **Fine**: "Input x was read into register 4, input y was read in register 5, add operation was performed using registers 4 and 5, ..."
- More queries are possible with fine-grained provenance, but...
 - Storage concerns
 - Performance concerns
- Abstraction can help here

Abstraction: Script, Workflow, Abstract Workflow

```
data = vtk.vtkStructuredPointsReader()
data.SetFileName("../examples/data/head.120.vtk")

contour = vtk.vtkContourFilter()
contour.SetInput(data.GetOutput())
contour.SetValue(0, 67)

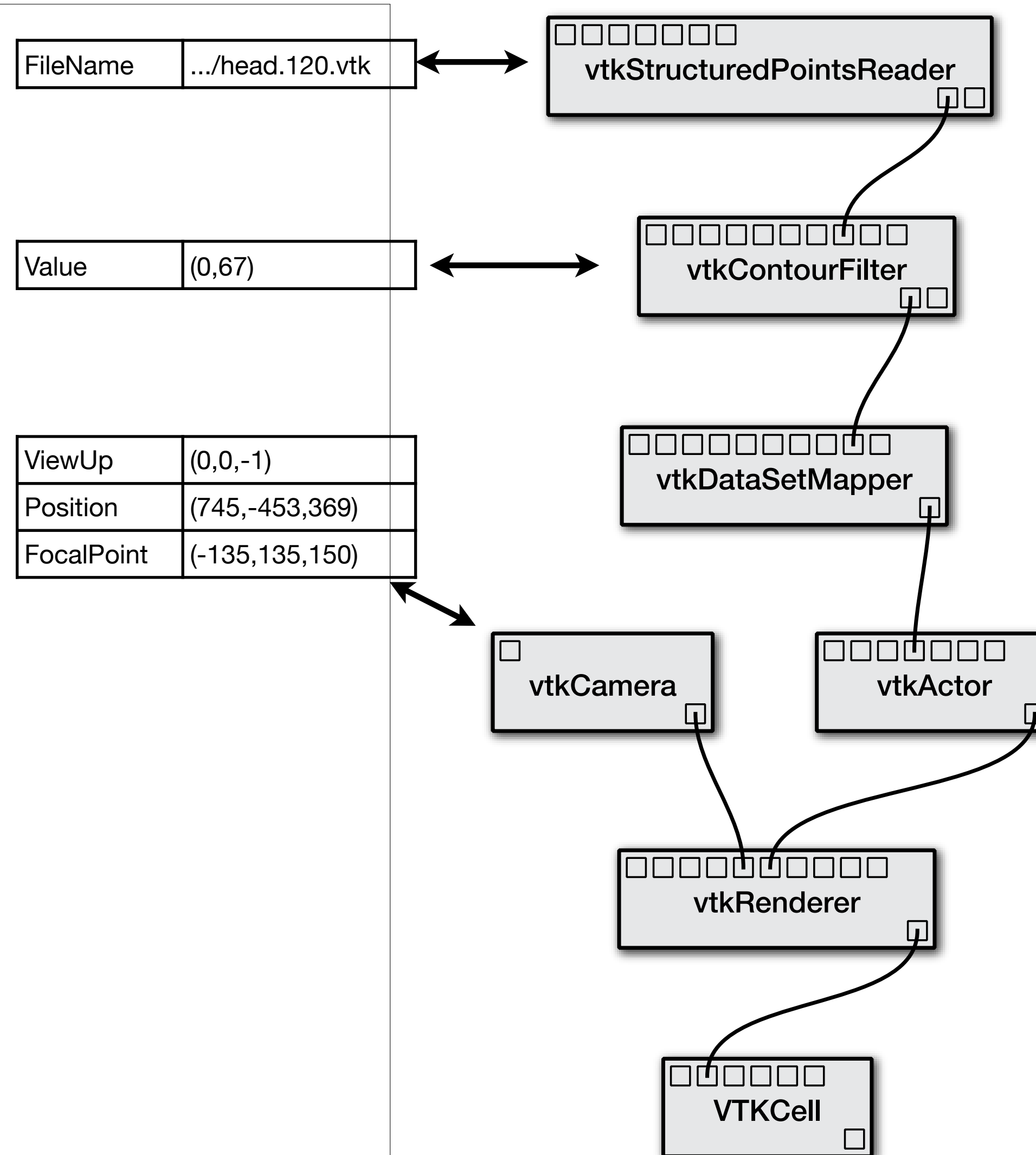
mapper = vtk.vtkPolyDataMapper()
mapper.SetInput(contour.GetOutput())
mapper.ScalarVisibilityOff()

actor = vtk.vtkActor()
actor.SetMapper(mapper)

cam = vtk.vtkCamera()
cam.SetViewUp(0,0,-1)
cam.SetPosition(745,-453,369)
cam.SetFocalPoint(135,135,150)
cam.ComputeViewPlaneNormal()

ren = vtk.vtkRenderer()
ren.AddActor(actor)
ren.SetActiveCamera(cam)
ren.ResetCamera()
renwin = vtk.vtkRenderWindow()
renwin.AddRenderer(ren)

style = vtk.vtkInteractorStyleTrackballCamera()
iren = vtk.vtkRenderWindowInteractor()
iren.SetRenderWindow(renwin)
iren.SetInteractorStyle(style)
iren.Initialize()
iren.Start()
```



Abstraction: Script, Workflow, Abstract Workflow

```

data = vtk.vtkStructuredPointsReader()
data.SetFileName("../examples/data/head.120.vtk")

contour = vtk.vtkContourFilter()
contour.SetInput(data.GetOutput())
contour.SetValue(0, 67)

mapper = vtk.vtkPolyDataMapper()
mapper.SetInput(contour.GetOutput())
mapper.ScalarVisibilityOff()

actor = vtk.vtkActor()
actor.SetMapper(mapper)

cam = vtk.vtkCamera()
cam.SetViewUp(0,0,-1)
cam.SetPosition(745,-453,369)
cam.SetFocalPoint(135,135,150)
cam.ComputeViewPlaneNormal()

ren = vtk.vtkRenderer()
ren.AddActor(actor)
ren.SetActiveCamera(cam)
ren.ResetCamera()
renwin = vtk.vtkRenderWindow()
renwin.AddRenderer(ren)

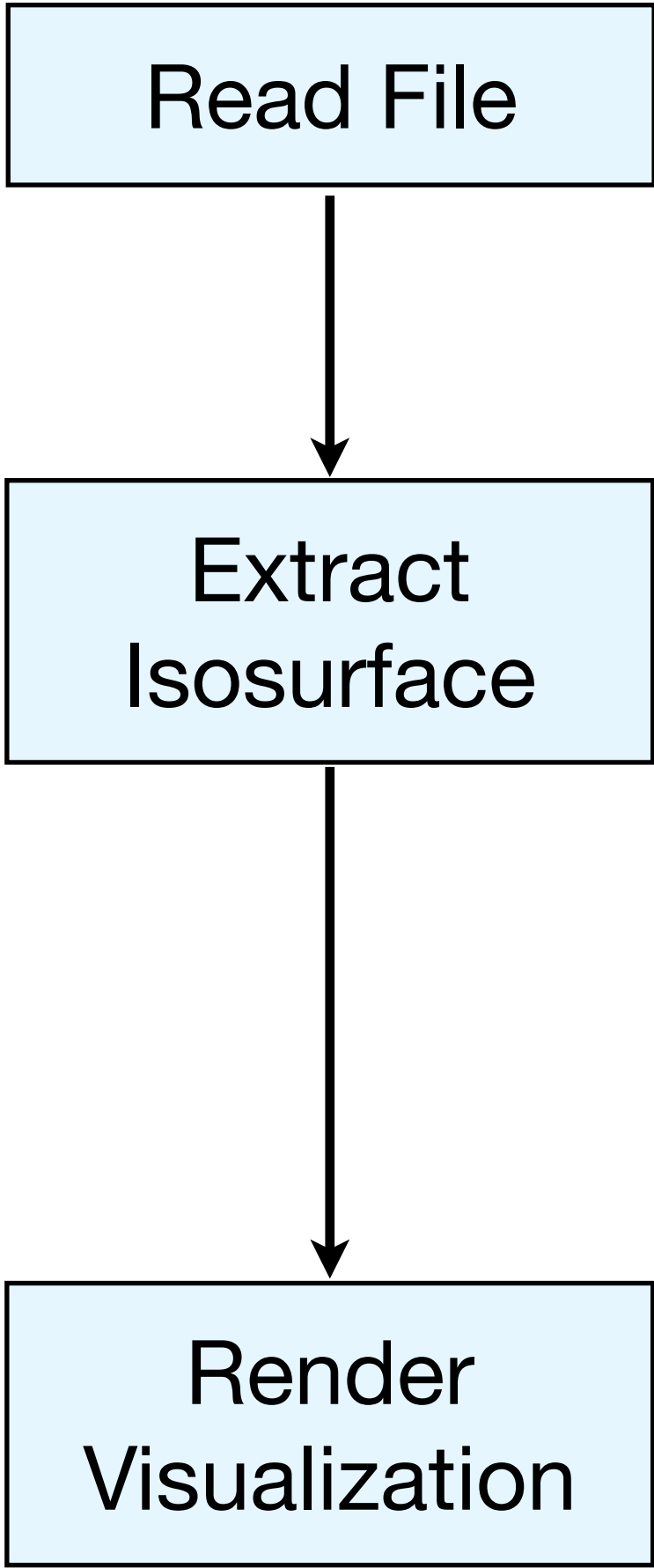
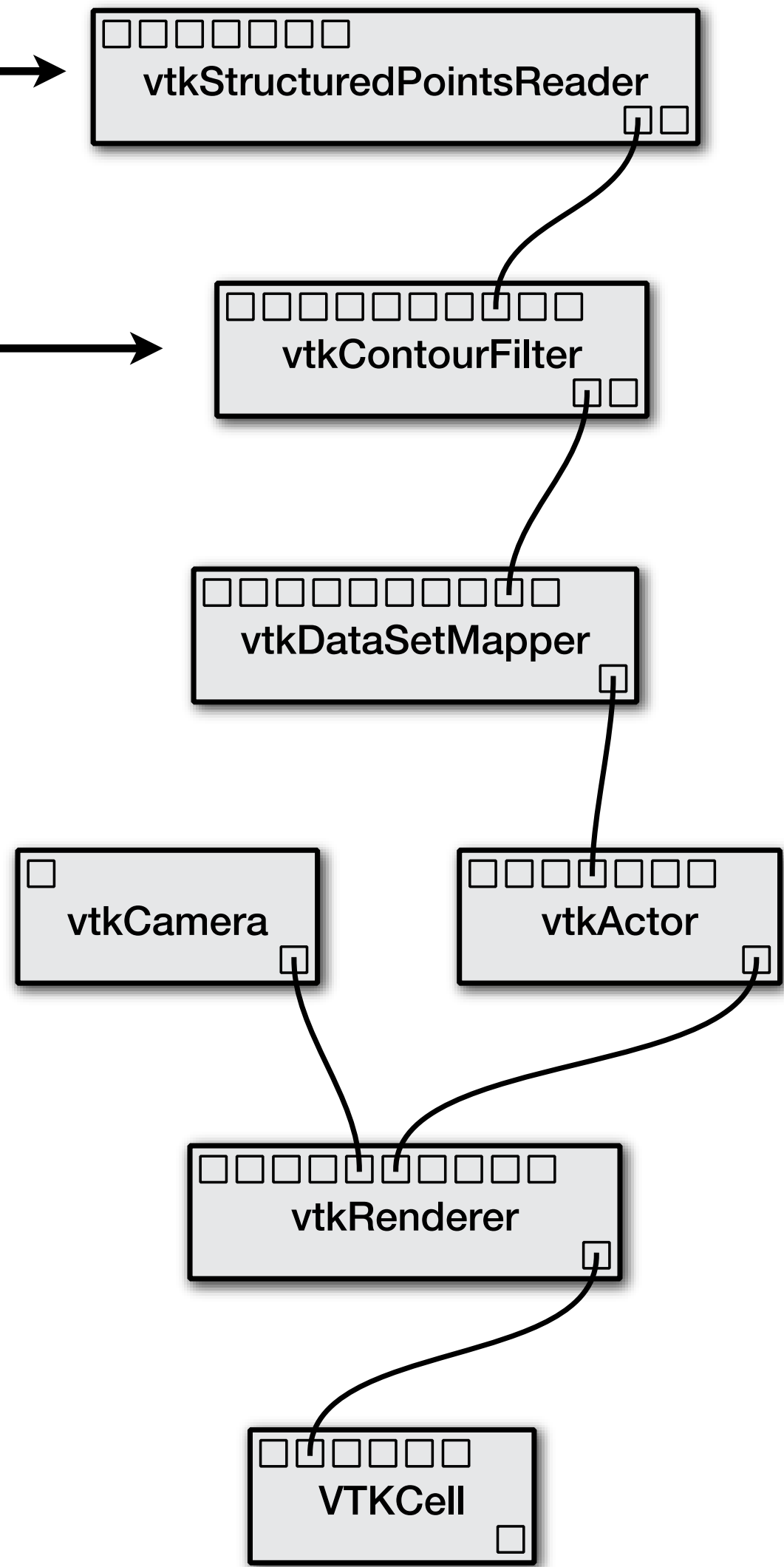
style = vtk.vtkInteractorStyleTrackballCamera()
iren = vtk.vtkRenderWindowInteractor()
iren.SetRenderWindow(renwin)
iren.SetInteractorStyle(style)
iren.Initialize()
iren.Start()

```

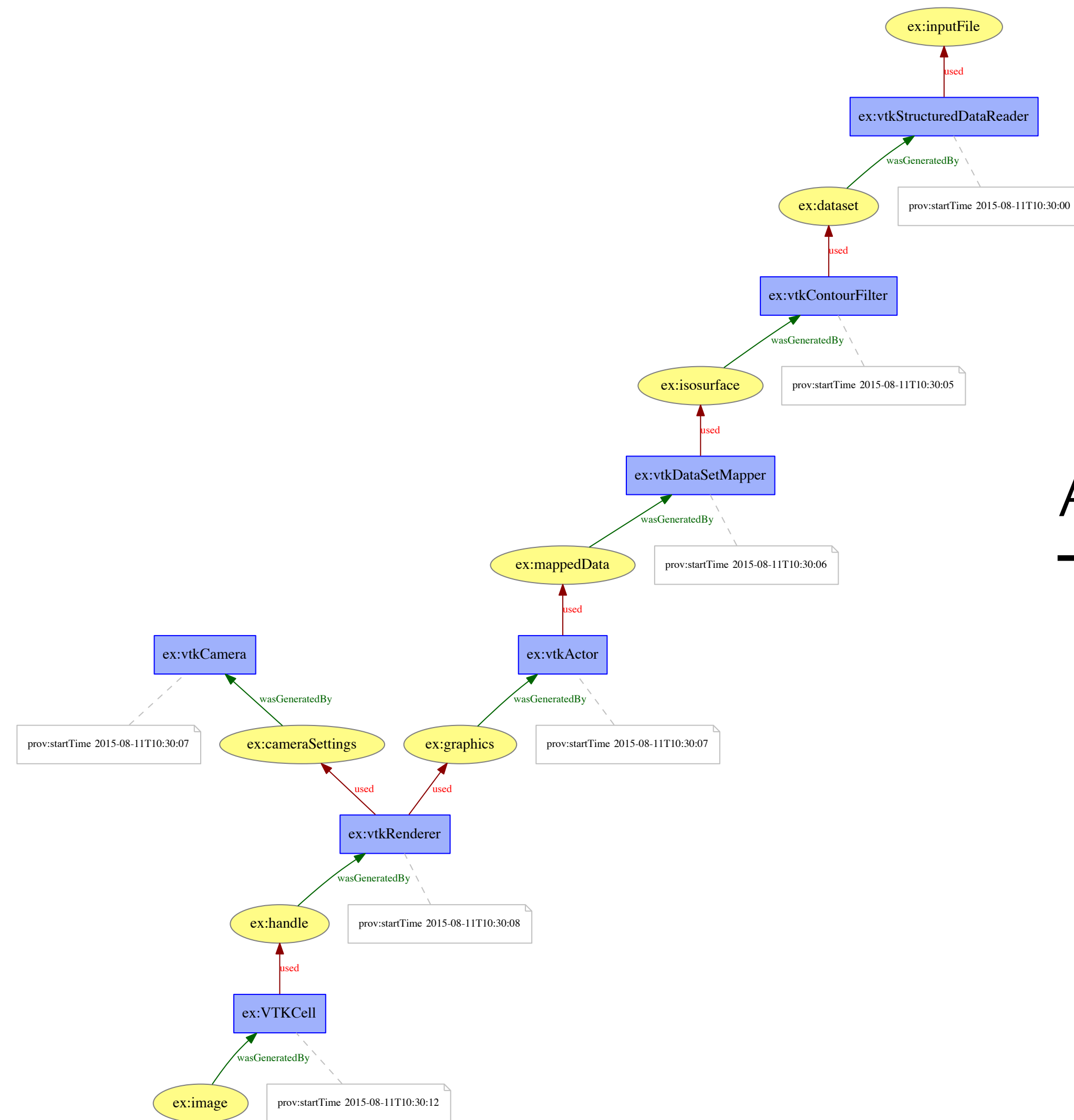
FileName	../head.120.vtk
----------	-----------------

Value	(0,67)
-------	--------

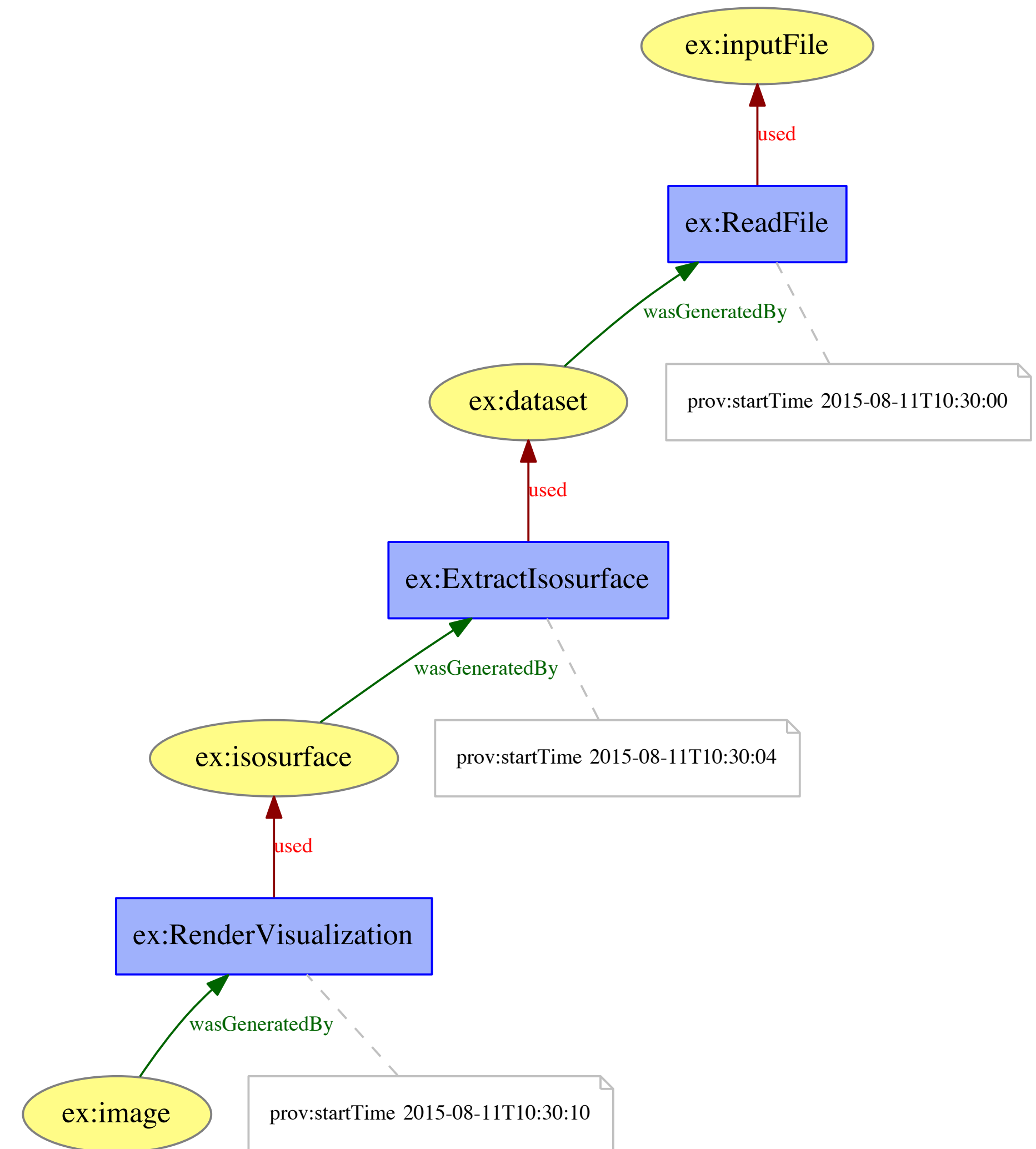
ViewUp	(0,0,-1)
Position	(745,-453,369)
FocalPoint	(-135,135,150)



Abstraction: Provenance Views



Abstract
→



Provenance Storage

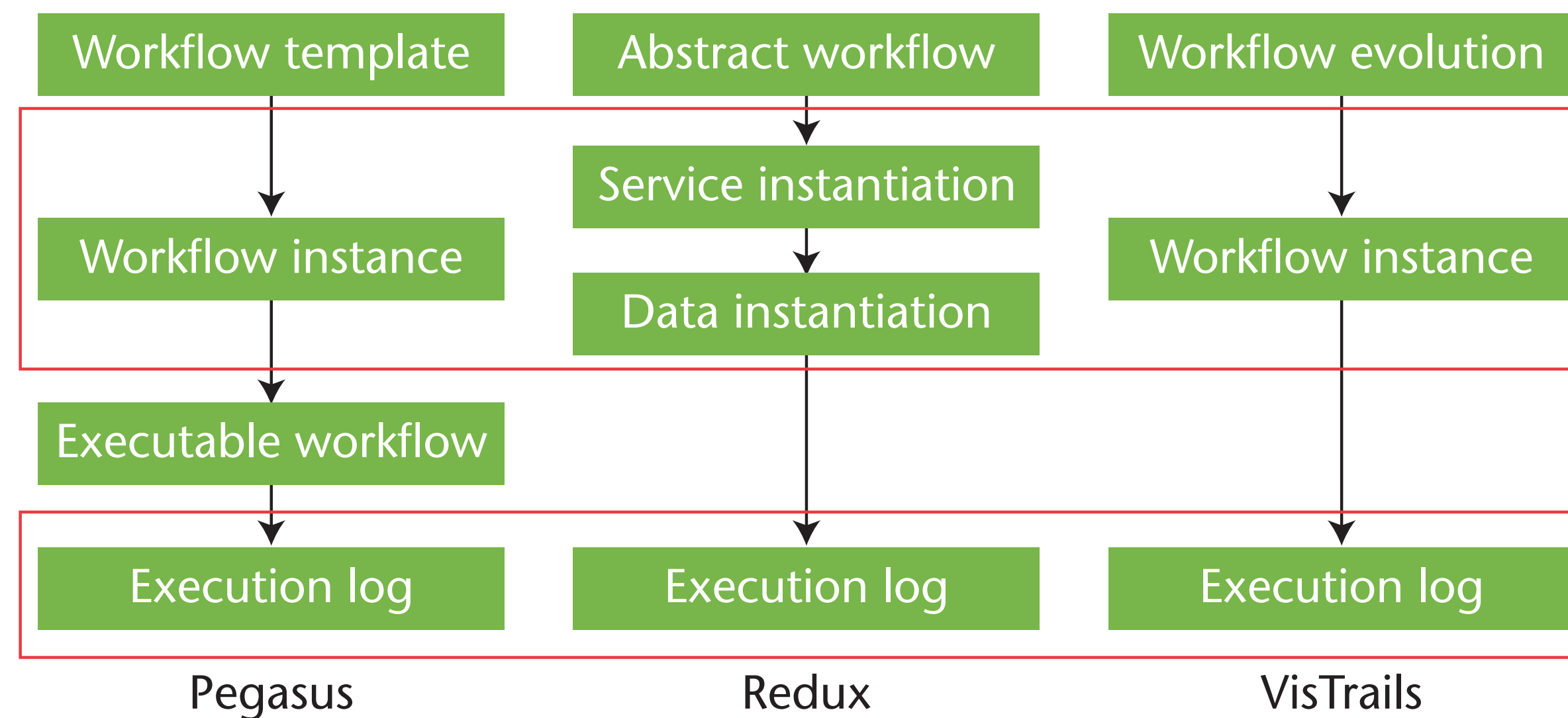
- Keeping provenance for each data item means lots of **repetition**
- Nested data storage also induces repetition
- Coarse provenance is naturally more compact, but how to decide what (not) to store?
- Repeated provenance is not uncommon:
 - Repeating the same computation with a different parameter
 - Creating a new computation that has a very similar structure to one that was run two weeks ago
- Provenance compression/factorization techniques (e.g. [Chapman et al., 2008], [Anand et al., 2009]) take advantage of that to reduce storage costs

Provenance Storage Formats

- Files, relational databases, XML databases, RDF (linked data)
- Log files are good for preserving data but can be bad to query or analyze
- Relational databases are great for column-specific queries but can be bad for dependency queries
- XML databases are more portable than relational databases but are usually less efficient for queries
- RDF triples are better for dependencies and integrating domain-specific knowledge but can be slower

Layered Provenance

- As with relational databases, want to normalize provenance to **minimize redundant information**
- Example: Don't store workflow specification each time that workflow is executed—store it once and reference it
- Also allow different layers for different aspects of provenance



[Freire et. al, 2008]

Provenance Models

- How provenance is represented (more abstract than the details of how it is actually stored)
- PROV (W3C Standard) has different storage backends for provenance but all of it conforms to the same model
- Model the objects involved and their relationships (e.g. activities, dependencies)
- Interoperability is a concern
 - Why? May use multiple tools/techniques to achieve a result, want to analyze the entire provenance chain

Prospective and Retrospective Provenance

- Prospective provenance is what was specified/intended
 - a workflow, script, list of steps
- Retrospective provenance is what actually happened
 - actual data, actual parameters, errors that occurred, timestamps, machine information
- **Do not need** prospective provenance to have retrospective provenance!
- Retrospective provenance is often the same type of information as prospective plus more
- Could have multiple retrospective provenance traces for one prospective provenance listing

Prospective and Retrospective Provenance

- **Example:** Baking a Cake
- Prospective Provenance (Recipe):
 1. Gather ingredients ($\frac{3}{4}$ cup butter, $\frac{3}{4}$ cocoa, $\frac{3}{4}$ cup flour, ...)
 2. Preheat oven to 350 degrees
 3. Grease cake pan
 4. Mix wet ingredients in large bowl
 5. Mix dry ingredients in a separate bowl
 6. Add dry mixture to wet mixture
 7. Pour batter into cake pan
 8. Put pan in the oven and bake for 30 minutes
 9. Take cake out of oven and let it cool



Prospective and Retrospective Provenance

- Retrospective Provenance (What actually happened)

1. Went to store to buy butter

↕ 2. Gathered ingredients (3/4 cup butter, 3/4 cocoa, **1 cup flour**, ...)

3. Greased cake pan

4. Preheated oven to 350 degrees

5. Mixed wet ingredients in large bowl

6. Mixed dry ingredients in a separate bowl

7. Added **wet** mixture to **dry** mixture

8. Poured batter into cake pan

9. Put pan in the oven and baked for **35 minutes**

10. Took cake out of oven and let it cool for **10 minutes**



Provenance Model History

- Community organized provenance challenges (2006-2009)
- First Provenance Challenge assessed capabilities of systems
- Second Provenance Challenge examined interoperability
- Led to development of Open Provenance Model (OPM), (2007)
 - Sought to establish interchange format for provenance
- Further work led to PROV W3C Recommendations (2013)
 - Some confusion from name changes from OPM to PROV even though concepts are similar
 - Focus is on **model** not formats

PROV: Three Key Classes



An **entity** is a physical, digital, conceptual, or other kind of thing with some fixed aspects; entities may be real or imaginary.



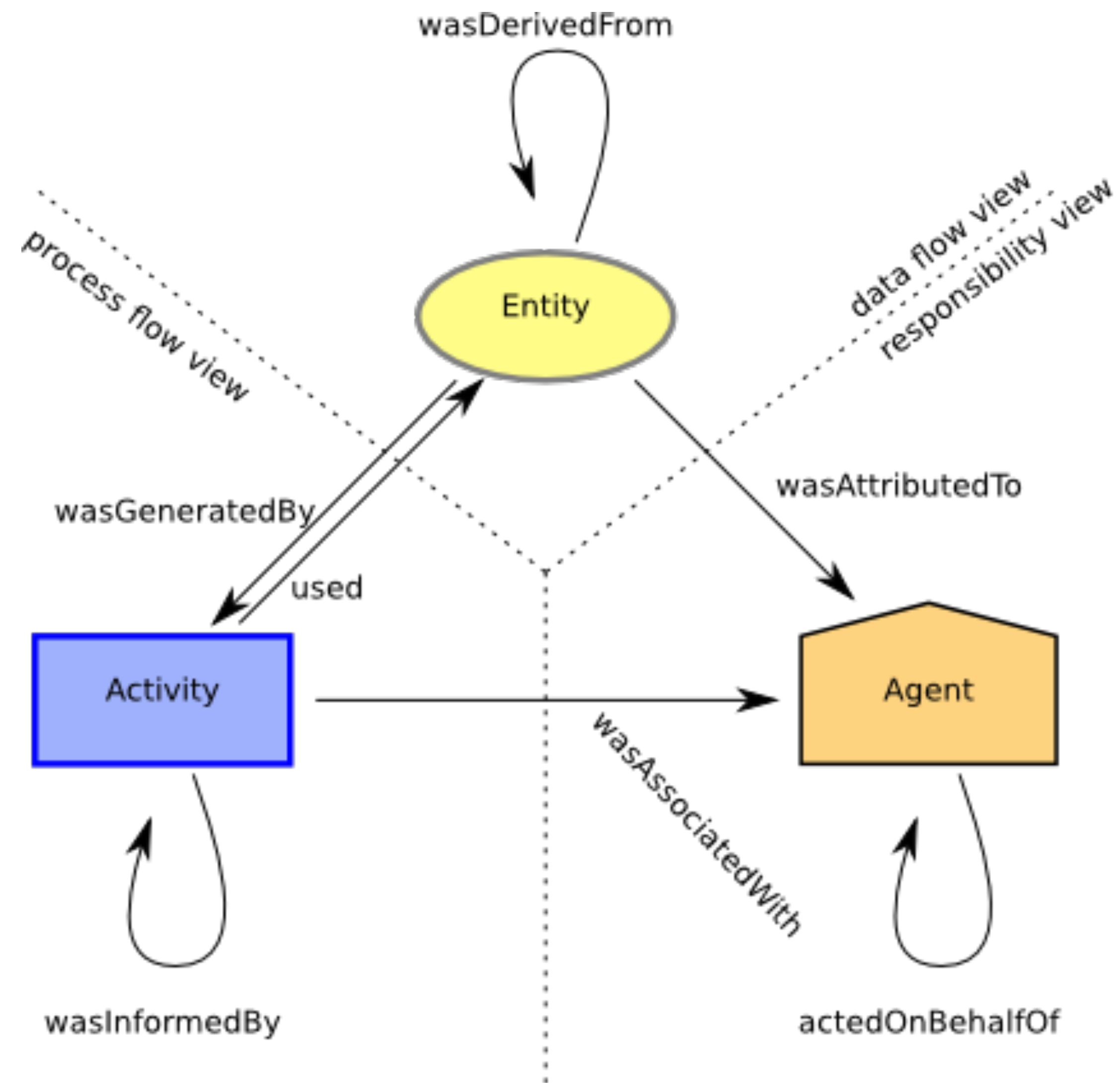
An **activity** is something that occurs over a period of time and acts upon or with entities; it may include consuming, processing, transforming, modifying, relocating, using, or generating entities.



An **agent** is something that bears some form of responsibility for an activity taking place, for the existence of an entity, or for another agent's activity.

[Moreau et al., 2014]

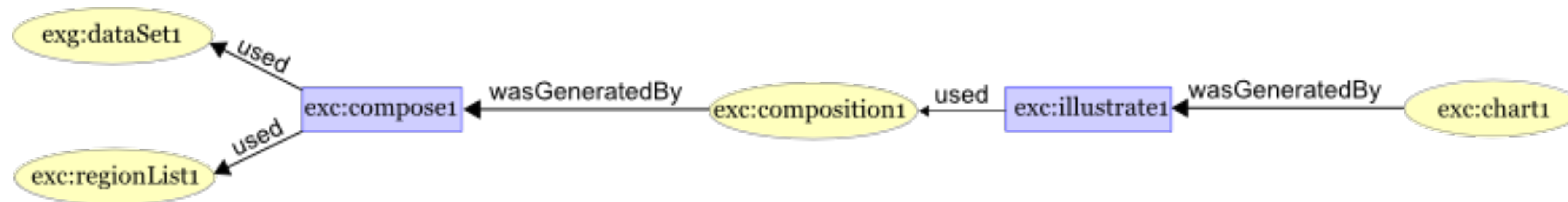
PROV: Three Views of Provenance



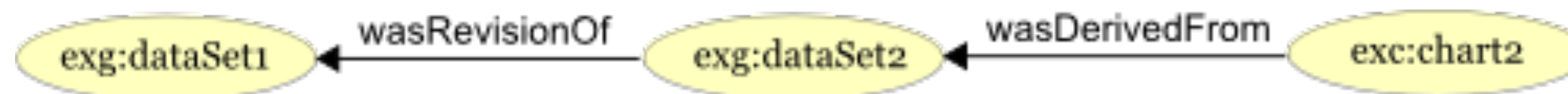
[Moreau et al., 2014]

PROV Edges: Derivation

- Derivation Edges:
 - wasGeneratedBy: entity \rightarrow activity
 - used: activity \rightarrow entity

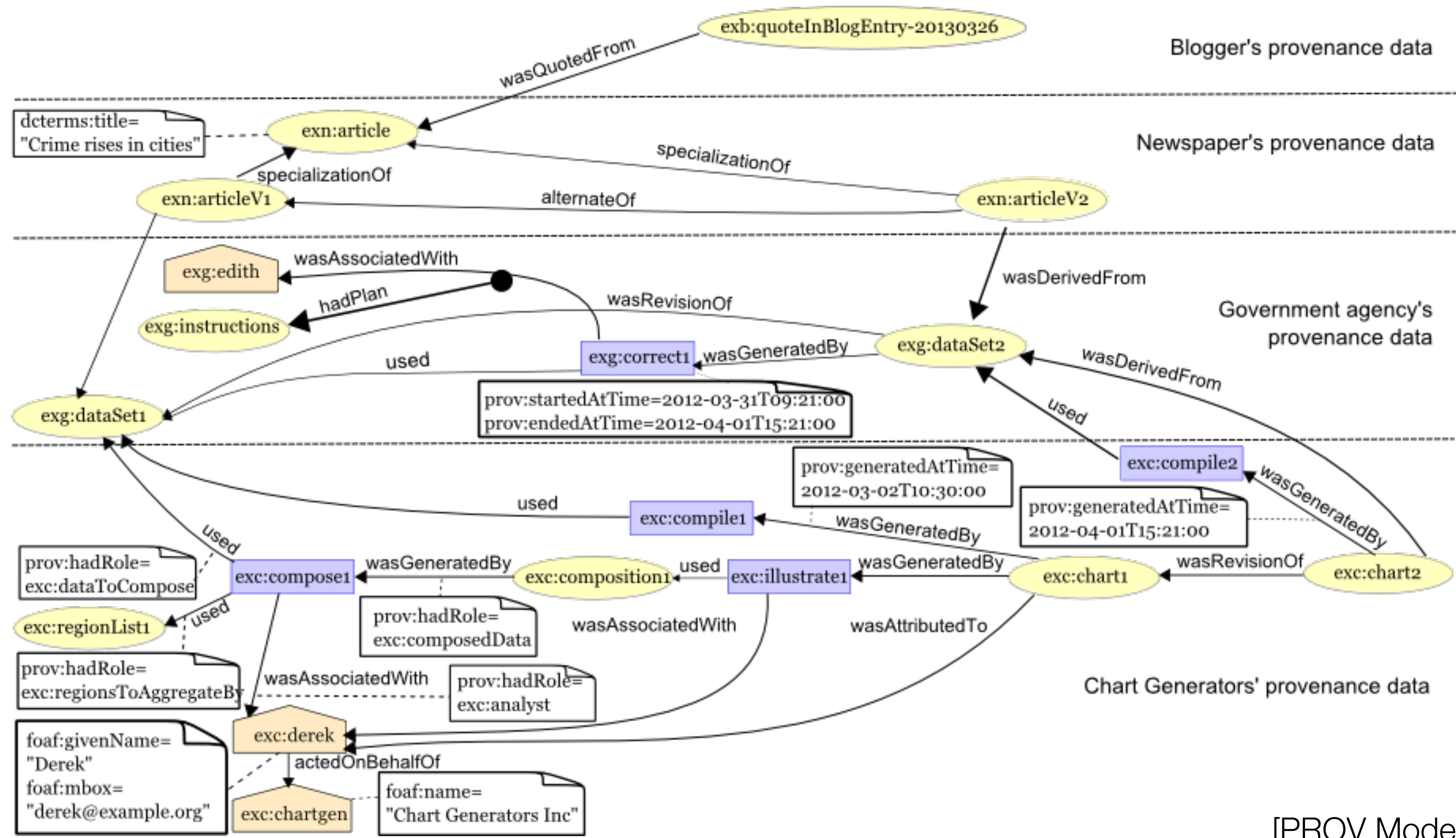


- wasDerivedFrom: entity \rightarrow entity



[PROV Model Primer, 2013]

PROV Example



[PROV Model Primer, 2013]

Querying Provenance

- Query methods are often tied to storage backend
- SQL, XQuery, Prolog, SPARQL, ...

REDUX

```
SELECT Execution.ExecutableWorkflowId, Execution.ExecutionId, Event.EventId, ExecutableActivity.ExecutableActivityId
from Execution, Execution_Event, Event, ExecutableWorkflow_ExecutableActivity, ExecutableActivity,
     ExecutableActivity_Property_Value, Value, EventType as ET
where Execution.ExecutionId=Execution_Event.ExecutionId
and Execution_Event.EventId=Event.EventId
and ExecutableActivity.ExecutableActivityId=ExecutableActivity_Property_Value.ExecutableActivityId
and ExecutableActivity_Property_Value.ValueId=Value.ValueId and Value.Value=Cast('-m 12' as binary)
and ((CONVERT(DECIMAL, Event.Timestamp)+0)%7)=0 and Execution_Event.ExecutableWorkflow_ExecutableActivityId=
     ExecutableWorkflow_ExecutableActivity.ExecutableWorkflow_ExecutableActivityId
and ExecutableWorkflow_ExecutableActivity.ExecutableWorkflowId=Execution.ExecutableWorkflowId
and ExecutableWorkflow_ExecutableActivity.ExecutableActivityId=ExecutableActivity.ExecutableActivityId
and Event.EventType=ET.EventType and ET.EventTypeName='Activity Start';
```

VisTrails

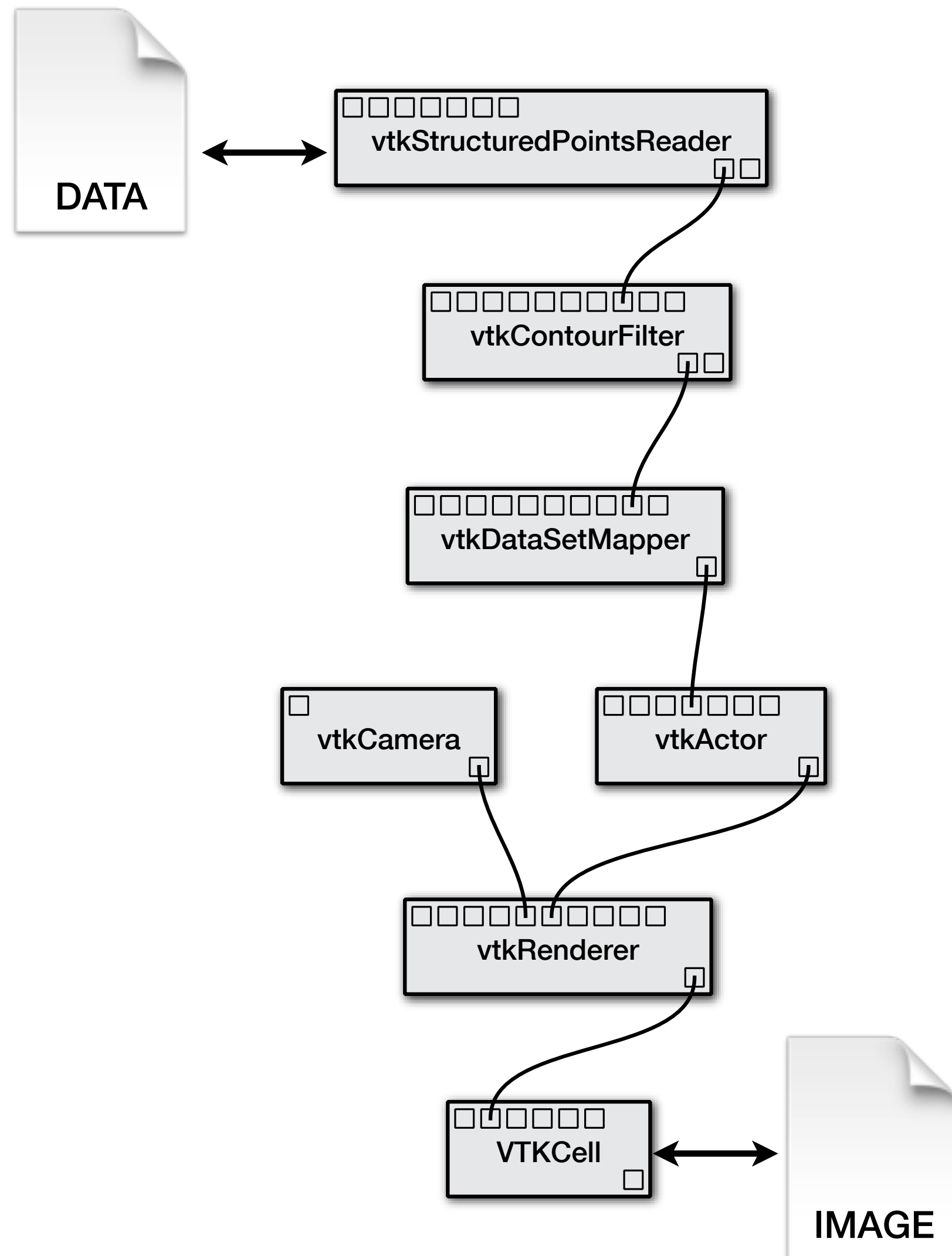
```
wf{*}: x where x.module='AlignWarp' and x.parameter('model')='12'
and (log{x}: y where y.dayOfWeek='Monday')
```

MyGrid

```
SELECT ?p
where (?p <http://www.mygrid.org.uk/provenance#startTime> ?time) and (?time > date)
using ns for <http://www.mygrid.org.uk/provenance#> xsd for <http://www.w3.org/2001/XMLSchema#>

SELECT ?p
where <urn:lsid:www.mygrid.org.uk:experimentinstance:HXQOVQA2ZI0>
(?p <http://www.mygrid.org.uk/provenance#runsProcess> ?processname .
?p <http://www.mygrid.org.uk/provenance#processInput> ?inputParameter .
?inputParameter <ont:model> <ontology:twelfthOrder>)
using ns for <http://www.mygrid.org.uk/provenance#> ont for <http://www.mygrid.org.uk/ontology#>
```

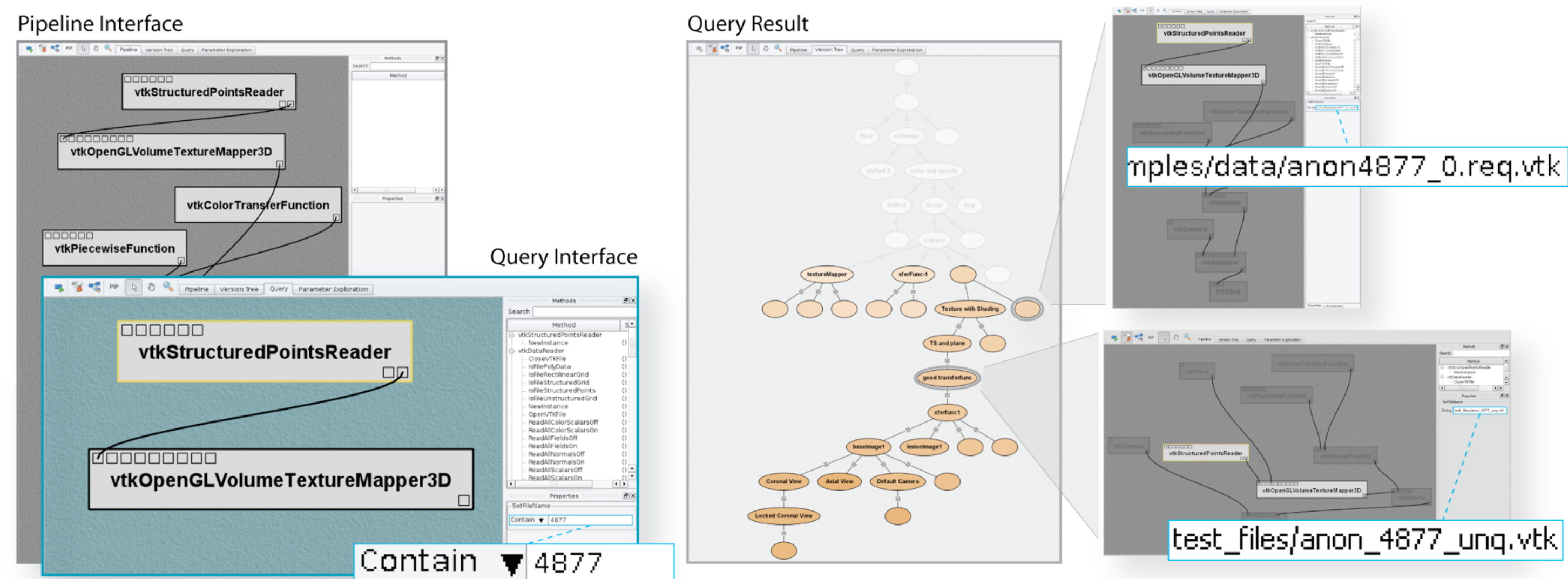
Querying Provenance



- *What process led to the output image?*
- *What input datasets contributed to the output image?*
- *What workflows include resampling and isosurfacing with isovalue 57?*
- Graph traversal or graph patterns
 - How do we write such queries?

Querying Provenance by Example

- Provenance is represented as graphs: hard to specify queries using text!
- Querying workflows by example [Scheidegger et al., TVCG 2007; Beerli et al., VLDB 2006; Beerli et al. VLDB 2007]
 - WYSIWYQ -- What You See Is What You Query
 - Interface to create workflow is same as to query



Stronger Links Between Provenance and Data

```
<workflow_exec id="1">
  <m_exec id="5"
    name="vtkStructuredData"
    package="edu.utah.sci.vistrails.vtk"
    version="5.6.0">
    <param id="2" name="SetFile"
      value="/MyData/0512s12.dat"/>
    </m_exec>
    <m_exec id="6"
      name="vtkContourFilter"
      package="edu.utah.sci.vistrails.vtk"
      version="5.6.0">
      <param id="3" name="SetValue"
        value="[1, 57]"/>
      <param id="4" name="ComputeScalarsOn"
        value="True"/>
      </m_exec>
      ...
      <m_exec id="11"
        name="FileSink"
        package="edu.utah.sci.vistrails.basic"
        version="1.5">
        <param id="15" name="path"
          value="/home/a/results/2312s12.dat"/>
        </m_exec>
```



FILE NOT FOUND

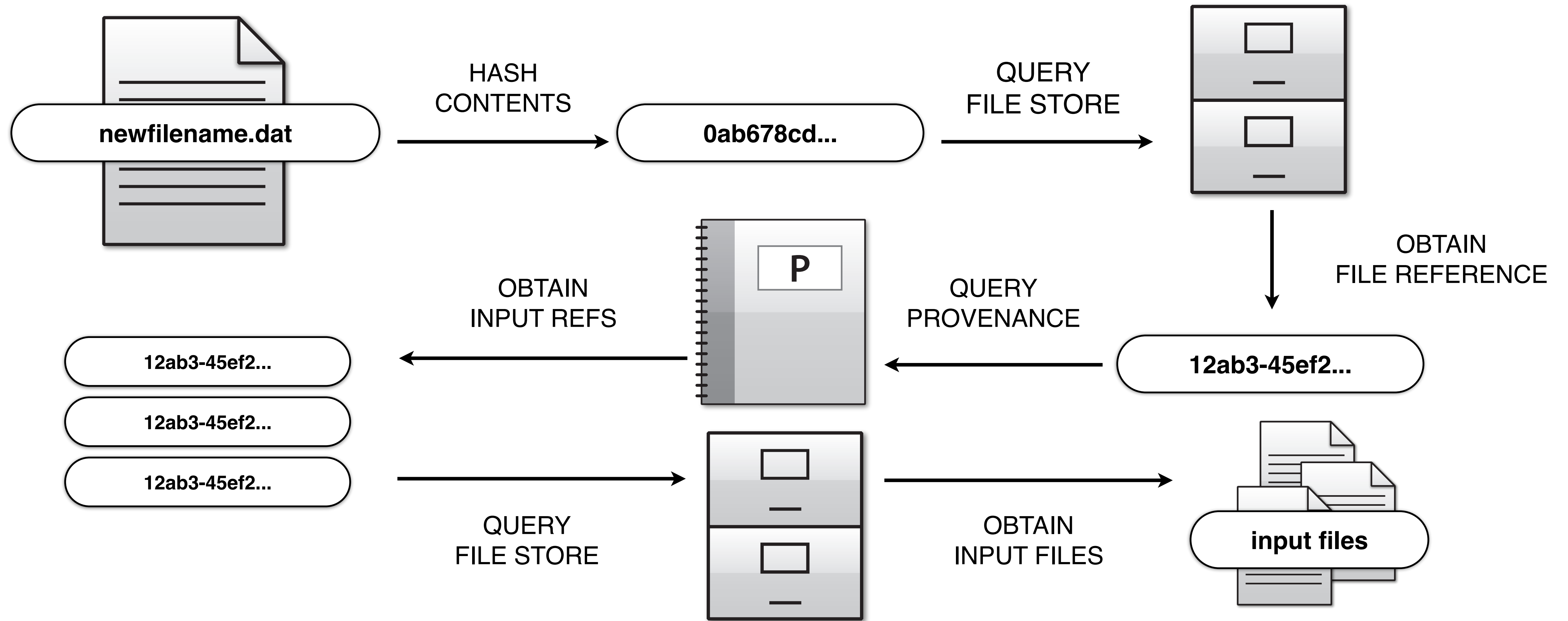


FILE NOT FOUND

- Filenames are often the mode of identification in data exploration
- We might also use URIs or access curated data stores
 - Always expected for exploratory tasks?
 - What happens if offline?
- Solution:
 - Managed store for data associated with computations
 - Improved data identification
 - Automatic versioning

[Koop et. al, 2010]

Provenance from Data



[Koop et. al, 2010]

Provenance-Enabled Systems

Table 1. Provenance-enabled systems.

System	Capture mechanism	Prospective provenance	Retrospective provenance	Workflow evolution
REDUX	Workflow-based	Relational	Relational	No
Swift	Workflow-based	SwiftScript	Relational	No
VisTrails	Workflow-based	XML and relational	Relational	Yes
Karma	Workflow- and process-based	Business Process Execution Language	XML	No
Kepler	Workflow-based	MoML	MoML variation	Under development
Taverna	Workflow-based	Scufl	RDF	Under development
Pegasus	Workflow-based	OWL	Relational	No
PASS	OS-based	N/A	Relational	No
ES3	OS-based	N/A	XML	No
PASOA/PreServ	Process-based	N/A	XML	No

[Freire et. al, 2008]

Provenance-Enabled Systems

Table 1. Provenanc			
System	Storage	Query support	Available as open source?
REDUX	Relational database management system (RDBMS)	SQL	No
Swift	RDBMS	SQL	Yes
VisTrails	RDBMS and files	Visual query by example, specialized language	Yes
Karma	RDBMS	Proprietary API	Yes
Kepler	Files; RDBMS planned	Under development	Yes
Taverna	RDBMS	SPARQL	Yes
Pegasus	RDBMS	SPARQL for metadata and workflow; SQL for execution log	Yes
PASS	Berkeley DB	nq (proprietary query tool)	No
ES3	XML database	XQuery	No
PASOA/PreServ	Filesystem, Berkeley DB	XQuery, Java query API	Yes

[Freire et. al, 2008]

Provenance-Enabled Systems

Table 1. Provenanc

System	Storage	Query support	Available as open source?
REDUX	Relational database management system (RDBMS)	SQL	No
Swift	RDBMS	SQL	Yes
VisTrails	RDBMS and files	Visual query by example, specialized language	Yes
Karma	RDBMS	Proprietary API	Yes
Kepler	Files; RDBMS planned	Under development	Yes
Taverna	RDBMS	SPARQL	Yes
Pegasus	RDBMS	SPARQL for metadata and workflow; SQL for execution log	Yes
PASS	Berkeley DB	nq (proprietary query tool)	No
ES3	XML database	XQuery	No
PASOA/PreServ	Filesystem, Berkeley DB	XQuery, Java query API	Yes



[Freire et. al, 2008]

Database Provenance

- Motivation: Data warehouses and curated databases
 - Lots of work
 - Provenance helps check correctness
 - Adds value to data by how it was obtained
- Three Types:
 - Why (Lineage): Associate each tuple t present in the output of a query with a set of tuples present in the input
 - How: Not just existence but routes from tuples to output (multiple contrib.'s)
 - Where: Location where data is copied from (may have choice of different tables)

[Cheney et al., 2007]

Why Provenance

Agencies

	name	based_in	phone
t_1 :	BayTours	San Francisco	415-1200
t_2 :	HarborCruz	Santa Cruz	831-3000

ExternalTours

	name	destination	type	price
t_3 :	BayTours	San Francisco	cable car	\$50
t_4 :	BayTours	Santa Cruz	bus	\$100
t_5 :	BayTours	Santa Cruz	boat	\$250
t_6 :	BayTours	Monterey	boat	\$400
t_7 :	HarborCruz	Monterey	boat	\$200
t_8 :	HarborCruz	Carmel	train	\$90

Q1:
SELECT a.name, a.phone
FROM Agencies a, ExternalTours e
WHERE a.name = e.name AND e.type='boat'

Result of Q_1 :

name	phone
BayTours	415-1200
HarborCruz	831-3000

- Lineage of (HarborCruz, 831-3000) :
{Agencies(t_2), ExternalTours(t_7)}
- Lineage of (BayTours, 415-1200):
{Agencies(t_1), ExternalTours(t_5, t_6)}
- This is not really precise because we don't need both t_5 and t_6 —only one is ok

[Cheney et al., 2007]

How Provenance

Agencies

	name	based_in	phone
t_1 :	BayTours	San Francisco	415-1200
t_2 :	HarborCruz	Santa Cruz	831-3000

ExternalTours

	name	destination	type	price
t_3 :	BayTours	San Francisco	cable car	\$50
t_4 :	BayTours	Santa Cruz	bus	\$100
t_5 :	BayTours	Santa Cruz	boat	\$250
t_6 :	BayTours	Monterey	boat	\$400
t_7 :	HarborCruz	Monterey	boat	\$200
t_8 :	HarborCruz	Carmel	train	\$90

```
Q2:
SELECT  e.destination, a.phone
FROM    Agencies a,
        (SELECT name,
                based_in AS destination
         FROM Agencies a
        UNION
         SELECT name, destination
         FROM ExternalTours ) e
WHERE   a.name = e.name
```

Result of Q₂:

destination	phone
San Francisco	415-1200
Santa Cruz	831-3000
Santa Cruz	415-1200
Monterey	415-1200
Monterey	831-3000
Carmel	831-3000

$t_1 \cdot (t_1 + t_3)$
 t_2^2
 $t_1 \cdot (t_4 + t_5)$
 $t_1 \cdot t_6$
 $t_1 \cdot t_7$
 $t_1 \cdot t_8$

- How provenance gives more detail about how the tuples provide witnesses to the result
- Prov of (San Francisco, 415-1200):
 $\{\{t_1\}, \{t_1, t_3\}\}$
- t_1 contributes **twice**
- Uses provenance semirings (the "polynomial" shown on the right)

[Cheney et al., 2007]



Where Provenance

Agencies

	name	based_in	phone
t_1 :	BayTours	San Francisco	415-1200
t_2 :	HarborCruz	Santa Cruz	831-3000

ExternalTours

	name	destination	type	price
t_3 :	BayTours	San Francisco	cable car	\$50
t_4 :	BayTours	Santa Cruz	bus	\$100
t_5 :	BayTours	Santa Cruz	boat	\$250
t_6 :	BayTours	Monterey	boat	\$400
t_7 :	HarborCruz	Monterey	boat	\$200
t_8 :	HarborCruz	Carmel	train	\$90

Q_1 :
SELECT
FROM
WHERE
 $a.name, a.phone$
Agencies a , ExternalTours e
 $a.name = e.name$
AND $e.type='boat'$

Q'_1 :
SELECT
FROM
WHERE
 $e.name, a.phone$
Agencies a , ExternalTours e
 $a.name = e.name$
AND $e.type='boat'$

Result of Q_1 :

name	phone
BayTours	415-1200
HarborCruz	831-3000

- Where provenance traces to specific locations, not the tuple values
- Q and Q' give the same result but the name comes from different places
- Prov of HarborCruz in second output: $(t_2, name)$
- Important in annotation-propagation

[Cheney et al., 2007]