Advanced Data Management (CSCI 640/490)

Provenance

Dr. David Koop





Sharing Data

- Required/encouraged by universities, funding agencies, publishers
- 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?

D. Koop, CSCI 640/490, Spring 2023

"Publications are arguments made by authors, and data are the evidence





2

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













Data Curation Lifecycle

The DCC Curation Lifecycle Model













Sequential Actions in Data Curation

- Create or Receive: Create/receive data and make sure metadata exists
- preservation
- 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

• Conceptualize: Plan creation of data—capture method and storage options. Appraise and Select: Evaluate data and select for long-term curation and

Ingest: Transfer data to an archive, repository, data centre or other custodian

Transform: Create new data from the original (migrate formats, subsets, etc.)





5

FAIR Principles

- 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 welldescribed so they can be replicated and/or combined in different settings

• Findable: Metadata and data should be easy to find for both humans and











Findable: DataCite Workflow













Accessible: DOI to Landing Page with Metadata



Document citing the data

D. Koop, CSCI 640/490, Spring 2023

Repository housing the data

Data store











Interoperable: Standard vocabularies

View as Table	View as Grid				« 1	2	3	4 5	6	7
ort by										
Name		\$	Registry	Name	Abbreviati	ion		Туре	Subjec	t
ecommended Records			க்	ABA Adult Mouse Brain	ABA			Standard	🕜 Neu	iroscieno
Recomm	ended		E	Access to Biological	ABCD			Standard	 Biod 	diversity
Associated Publication?				Collection Data					🖌 Life	Science
No Publication	Has Publication									
Claimed?		_								
No Maintainer	Has Maintainer									
Uncertain Deprecated	n development Re	eady		Access to Biological Collection Databases Extended for Geosciences	ABCDEFG			Standard	EartPale	th Scien eontolog
Standard Type				Access to Biological	ABCDDNA			Standard	Biod	diversity
Terminology Artifact		771		Collection Data					Life	Scienci
Model/Format		405								
Reporting Guideline		163								
Metric		30		.ACE format	.ACE forma	at		Standard		Science
Identifier Schema		15								
	Shov	v More	கீ	AdaLab-meta ontology	ADALAB-N	IETA		Standard	None	
Domains			Å	AdaLab ontology	ADALAB			Standard	None	
Report		141		Adverse Drug	EU-ADR M	IL		Standard	None	
		134		Language						

						nc5	~	/ \C/ V	ancea											
				Sho	owing r	records	1 - 50	of 138 4	ŀ.											
8 9	10) 11	12	13	14	15	16	17	18 1	9 20	21	22	23	24	25	26	27	28	»	
		Domain				Taxonoi	my	Relate	d Database		Related	Standard	i t	Related F	Policy	In Co	llection/	/Recom	mendatior	n Status
		BrainBrain II	Gene I Gene I	Expression		🖌 Mus i	musculus	Neuro	lorpho.Org		None		I	None		No	ne			R
Biology		None				All		GBIF ALA IP Reposi GBIF S Reposi Canade SiB Co Colomb Plus 1	T - GBIF Aus tory pain IPT - G tory ensys IPT - C ensys Repos lombia IPT - oia Reposito more	stralia BIF Spain GBIF sitory GBIF ry	ABCDDN	IA G	1	None		T	DWG Biodive	ersity Informa	ation Standards	R
GeologySoil Scier	nce	None				II AII		GeoCA	Se Data Poi	rtal	XML ABCD		I	None		No	ne			R
Biology		 DNA S Experin Sequence Deoxyr Polyme Plus 1 	equence Data ment Metadat nce ribonucleic Ac erase Chain R more	a id Reaction		I Ali		GenBa	nk		MOD-CC ABCD		1	None			ƊWG Biodive	ersity Informa	ation Standards	
		DNA SDeoxyr	equence Data ribonucleic Ac	a 🛷 Co id 🗣 G	enome	II All		None			None		I	None		No	ne			R
		None				🛷 All		None			None		I	None		Noi	ne			R
		None				🗣 All		None			None		I	None		Noi	ne			R
		AdversElectro	e Reaction	cord			o sapiens	None			XML		1	None		No	ne		ſ	fairs









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.









Reusable: Data Citation & Metrics



D. Koop, CSCI 640/490, Spring 2023





Northern Illinois University



<u>Assignment 5</u>

- Divvy Bikes Data
- Spatial, Graph, and Temporal Data Processing
- Use pandas, geopandas, neo4j, (modin for extra credit)

D. Koop, CSCI 640/490, Spring 2023

Processing odin for extra credit)





geopandas example





Provenance





What actually happened in a computational experiment?





Provenance in Art



D. Koop, CSCI 640/490, Spring 2023

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









Provenance in Art



D. Koop, CSCI 640/490, Spring 2023

oil on canvas 1937.1.72

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

Rembrandt van Rijn Dutch, 1606 - 1669 Self-Portrait, 1659 Andrew W. Mellon Collection

Provenance

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









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

Test:	в	м	BM	p	T	PT	BATON		
237-12			. ++				. ++.	OK.	
		- 0	BTI	PM.	= 0	-			
143-0-	-0	- P	T.	-1	-0	0			
2	++	+	++						
3	4+	++	++		N	cost of th	his is cli	inally sympto	John
¥	+4	++	**			0		11	
•	++	++	**						
1									
9	-	-							
10		-							
11		4		4	1 +				
12		-	1.	34	unnil	-	11-20		
N	-	-	++ (elevent	ale)	(Hon	au!).	Sac.	
15	++	-	++ (Sh	rale mit		1		
.16	++	++	++						
17		-							
-									
	S				1				-
230-1					-		N.	t cole.	
2.34-2									
	n.1								
		10.00							
243-1.	Fin	nBT	plate.						
21	++	++	++						
23	1								
24	Ĩ.								
25	40		1 States						
26	++	-	4+	Sti	ich ni				•
n	++-	+	*						
20		-	44	S.	usland	t			
30	++	+	+	-		-			
51	++-	+.	++-						
32	++	++	++						
32	++	+	#						
71	14	14	at						
46		14							
\$7									
35						000		otion	~





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

NIU



Provenance in Computational Science



D. Koop, CSCI 640/490, Spring 2023





18

Evolution of Publication

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

D. Koop, CSCI 640/490, Spring 2023

/tests happens during your runs)





inverse system size 1/L Provenance-Rich Publication

0.05

Galois Conjugates of Topological Phases

0.1

0.15

0.2

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. Othere we invest of the Galois con gates 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 where 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 nonunitary 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 $su(2)_3$.

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 lightcone. 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 $Fib^{\mathcal{G}}$ and $Fib^{\mathcal{G}}$, 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-

201 Jul S .str-el] mat. cond 267 $\hat{\mathbf{O}}$ 00 $\overline{}$ arXi

0.25

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 $(\theta \neq 0)$.







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 owners."
- generated and/or derivation what data a result depended on
- when it occurred, who initiated it, notes about it
- many questions

record of the ultimate derivation and passage of an item through its various

Focus on causality—the sequence of steps that detail how a result was

• Provenance itself is **data**, this list of steps along with metadata for each step:

Can be used to preserve information about an experiment and to answer







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)
                                                                          .../head.120.vtk
                                                               FileName
contour = vtk.vtkContourFilter()
contour.SetInput(data.GetOutput())
contour.SetValue(0, 67)
mapper = vtk.vtkPolyDataMapper()
mapper.SetInput(contour.GetOutput())
                                                                              (0, 67)
                                                                Value
mapper.ScalarVisibilityOff()
actor = vtk.vtkActor()
actor.SetMapper(mapper)
cam = vtk.vtkCamera()
cam.SetViewUp(0, 0, -1)
                                                                             (0,0,-1)
                                                               ViewUp
cam.SetPosition(745,-453,369)
                                                                          (745, -453, 369)
                                                               Position
cam.SetFocalPoint(135,135,150)
cam.ComputeViewPlaneNormal()
                                                                          (-135,135,150)
                                                              FocalPoint
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)
                                                                          .../head.120.vtk
                                                               FileName
contour = vtk.vtkContourFilter()
contour.SetInput(data.GetOutput())
contour.SetValue(0, 67)
mapper = vtk.vtkPolyDataMapper()
mapper.SetInput(contour.GetOutput())
                                                                              (0, 67)
                                                                Value
mapper.ScalarVisibilityOff()
actor = vtk.vtkActor()
actor.SetMapper(mapper)
cam = vtk.vtkCamera()
cam.SetViewUp(0, 0, -1)
                                                                             (0,0,-1)
                                                               ViewUp
cam.SetPosition(745,-453,369)
                                                                          (745, -453, 369)
                                                               Position
cam.SetFocalPoint(135,135,150)
cam.ComputeViewPlaneNormal()
                                                                          (-135,135,150)
                                                              FocalPoint
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: Provenance Views







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

- 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



D. Koop, CSCI 640/490, Spring 2023

• As with relational databases, want to normalize provenance to **minimize**











Provenance Models

- 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

How provenance is represented (more abstract than the details of how it is











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 (3/4 cup butter, 3/4 cocoa, 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

D. Koop, CSCI 640/490, Spring 2023









40

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.





PROV: Three Views of Provenance



D. Koop, CSCI 640/490, Spring 2023



[Moreau et al., 2014]

Northern Illinois University





PROV Edges: Derivation

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







PROV Example



D. Koop, CSCI 640/490, Spring 2023



Northern Illinois University



46

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.EventTypeId=ET.EventTypeId 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; Beeri et al., VLDB 2006; Beeri et al. VLDB 2007]
- - WYSIWYQ -- What You See Is What You Query
 - Interface to create workflow is same as to query



D. Koop, CSCI 640/490, Spring 2023







49

Stronger Links Between Provenance and Data



- 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







Provenance from Data









Provenance-Enabled Systems

Table 1 Provenance enabled system

Table I. Provenal	nce-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, 2
Koop, CSCI 640/49	0, Spring 2023			Northern Illinois University





Provenance-Enabled Systems

Table 1. Provenanc					
System	Storage	Query support	Available source?	as open	
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	







Provenance-Enabled Systems

Table 1. Provenanc			
System	Storage	Query support	Available as open source?
REDUX	Relational database management system (RDBMS)	SQL	No
Swift	RDBMS	SC	Yes
VisTrails	RDBMS and files	Visual query by example, spe <mark>cialized</mark>	Yes
Karma	RDBMS Jupyter	Proprietary API	Yes
Kepler	Files; RDBMS plann	Under development	Yes
Taverna	RDBMS		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, 2
Koop, CSCI 640/490,	Spring 2023		Northern Illinois University





Today: Two types of provenance

- Database Provenance
- Evolution Provenance







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)













Provenance in Databases

A. Amarilli





Why Provenance

Agencies

	0		
	name	based_in	phone
t_1 :	BayTours	San Francisco	415-1200
t_2 :	HarborCruz	Santa Cruz	831-3000

ExternalTours

name	destination	type	price
BayTours	San Francisco	cable car	\$50
BayTours	Santa Cruz	bus	\$100
BayTours	Santa Cruz	boat	\$250
BayTours	Monterey	boat	\$400
HarborCruz	Monterey	boat	\$200
HarborCruz	Carmel	train	\$90
	name BayTours BayTours BayTours BayTours HarborCruz HarborCruz	namedestinationBayToursSan FranciscoBayToursSanta CruzBayToursSanta CruzBayToursMontereyHarborCruzMontereyHarborCruzCarmel	namedestinationtypeBayToursSan Franciscocable carBayToursSanta CruzbusBayToursSanta CruzboatBayToursMontereyboatHarborCruzMontereyboatHarborCruzCarmeltrain

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(t2), ExternalTours(t7)}
- Lineage of (BayTours, 415-1200): {Agencies(t1), ExternalTours(t5,t6)}
- This is not really precise because we don't need both t5 and t6—only one is ok













How Provenance

Agencies

	U		
	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_2 :

SELECT	e. destination, a. phone	Result of Q_2 :		
FROM	Agencies a ,	destination	phone	
	(SELECT name,	San Francisco	415-1200	$t_1 \cdot (t_1 + t_3)$
	based_in AS destination	Santa Cruz	831-3000	t_{2}^{2}
	FROM Agencies a	Santa Cruz	415-1200	$t_1 \cdot (t_4 + t_5)$
	UNION	Monterey	415-1200	$t_1 \cdot t_6$
	SELECT name, destination	Monterey	831-3000	$t_1 \cdot t_7$
	FROM External Tours) e	Carmel	831-3000	$t_1 \cdot t_8$
WHERE	a.name = e.name			

- How provenance gives more detail about how the tuples provide witnesses to the result
- Prov of (San Francisco, 415-1200): $\{ \{ t1 \}, \{ t1, t3 \} \}$
- t1 contributes **twice**
- Uses provenance semirings (the
- "polynomial" shown on the right)
- $t_5)$









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

name, a.phone gencies a, ExternalTours ename = e.nameND 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: (t2, name)
- Important in annotation-propogation















D. Koop, CSCI 640/490, Spring 2023

Evolution Provenance







Data Exploration



D. Koop, CSCI 640/490, Spring 2023

[Modified from Van Wijk, Vis 2005]









Data Exploration



- Data analysis and visualization are iterative processes
- In exploratory tasks, change is the norm!

D. Koop, CSCI 640/490, Spring 2023

[Modified from Van Wijk, Vis 2005]









Exploration and Creativity Support

- Reasoning is key to the exploratory processes
- "Reflective reasoning requires the ability to store temporary results, to make inferences from stored knowledge, and to follow chains of reasoning backward and forward, sometimes backtracking when a promising line of thought proves to be unfruitful. ...the process is slow and laborious" -Donald A. Norman
- Need external aids—tools to facilitate this process - "Creativity support tools" — Ben Shneiderman
- Need aid from people—collaboration







Change-based Provenance: Photo Editing

User Actions



Undo/Redo History









Change-based Provenance: Photo Editing

• User Actions



Undo/Redo History









Version Trees

- Undo/redo stacks are linear!
- We lose history of exploration
- Old Solution: User saves files/state
- VisTrails Solution:
 - Automatically & transparently capture entire history as a tree
 - Users can tag or annotate each version
 - Users can go back to **any** version by selecting it in the tree











VisTrails



© 2011-2013 NYU-Poly. © 2006-2011 University of Utah. All Rights Reserved. J. Freire, C. Silva, E. Anderson, L. Bavoil, C. Brooks, J. Callahan, S. Callahan, T. Ellqvist, L. Carlo, D. Koop, L. Lins, P. Mates, D. Rees, E. Santos, C. Scheidegger, N. Smith, H. Vo











VisTrails

- Comprehensive provenance infrastructure for computational tasks
- Focus on exploratory tasks such as simulation, visualization, and data analysis
- Transparently tracks provenance of the discovery process from data acquisition to visualization
 - The trail followed as users generate and test hypotheses
 - Users can refer back to any point along this trail at any time
- Leverage provenance to streamline exploration
- Focus on usability—build tools for scientists






Workflow Evolution Provenance











Workflow Evolution Provenance



D. Koop, CSCI 640/490, Spring 2023

with labels
filtered

delete module "GMapCell"

delete module "CellLocation"

delete module "ProjectTable"

delete module "SelectFromTable"

. . .

. . .

add module "SelectFromTable"

add parameter "float_expr" to "SelectFromTable" with value "latitutde > 40.6"

delete parameter "float_expr" from "SelectFromTable"

add parameter "float_expr" to "SelectFromTable" with value "latitutde > 40.7"

delete parameter "float_expr" from "SelectFromTable"

add parameter "float_expr" to "SelectFromTable" with value "latitutde > 40.8"











Execution Provenance









Execution Provenance

```
<module id="12" name="vtkDataSetReader"
        start time="2010-02-19 11:01:05"
        end time="2010-02-19 11:01:07">
 <annotation key="hash"</pre>
            value="c54bea63cb7d912a43ce"/>
</module>
<module id="13" name="vtkContourFilter"
        start time="2010-02-19 11:01:07"
        end time="2010-02-19 11:01:08"/>
<module id="15" name="vtkDataSetMapper"
        start time="2010-02-19 11:01:09"
        end time="2010-02-19 11:01:12"/>
<module id="16" name="vtkActor"
        start time="2010-02-19 11:01:12"
        end time="2010-02-19 11:01:13"/>
<module id="17" name="vtkCamera"
        start time="2010-02-19 11:01:13"
        end time="2010-02-19 11:01:14"/>
<module id="18" name="vtkRenderer"
        start time="2010-02-19 11:01:14"
        end time="2010-02-19 11:01:14"/>
• • •
```









Capturing and querying fine-grained provenance of preprocessing pipelines in data science

A. Chapman, P. Missier, L. Lauro, R. Torlone



