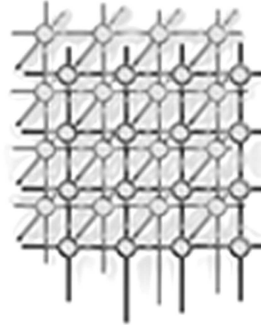


The First Provenance Challenge



Luc Moreau*, Bertram Ludäscher, Ilkay Altintas,
Roger S. Barga, Shawn Bowers, George Chin,
Shirley Cohen, Sarah Cohen-Boulakia, Ben Clifford,
Susan Davidson, Ewa Deelman, Luciano
Digiampietri, Ian Foster, Juliana Freire, James
Frew, Joe Futrelle, Tara Gibson, Yolanda Gil,
Carole Goble, Jennifer Golbeck, Paul Groth, David
A. Holland, Sheng Jiang, Jihie Kim, Ales Krenek,
Timothy McPhillips, Gaurang Mehta, Simon Miles,
Dominic Metzger, Steve Munroe, Jim Myers, Beth
Plale, Norbert Podhorszki, Varun Ratnakar, Karen
Schuchardt, Margo Seltzer, Yogesh L. Simmhan,
Peter Slaughter, Eric Stephan, Robert Stevens,
Daniele Turi, Mike Wilde, Jun Zhao, Yong Zhao

University of Southampton

SUMMARY

The first Provenance Challenge was a community activity aiming at understanding the expressiveness of provenance representations and capabilities of provenance systems. To this end, a Functional Magnetic Resonance Imaging workflow was defined, which participants had to either simulate or run in order to produce some provenance representation, from which a set of identified queries had to be implemented and executed. Seventeen teams responded to the challenge, and submitted their inputs. In this paper, we present the challenge workflow and queries, and summarise the participants contributions.

KEY WORDS: Provenance, representation, storing, recording, capabilities, queries

*Correspondence to: Electronics and Computer Science, University of Southampton Southampton, SO17 1BJ, U.K

†E-mail: contact author: l.moreau@ecs.soton.ac.uk



1. Introduction

The term *provenance* is commonly used in the context of art to denote the documented history or the chain of ownership of an art object. Provenance helps determine the authenticity and therefore the value of art objects. If the provenance of data produced by computer systems could be determined as it can for some works of art, then users, in their daily applications, would be able to interpret and judge the quality of data better [10]. In particular, the scientific and grid communities consider that, in order to support reproducibility, workflow management systems will be required to track and integrate provenance information as an integral product of the workflow [4]. Several surveys of provenance are available [11, 2, 7].

Against this background, the *International Provenance and Annotation Workshop* (IPAW'06), held on May 3-5, 2006 in Chicago, involved some 50 participants interested in the issues of data provenance, process documentation, data derivation, and data annotation [9, 1]. During a session on provenance standardisation, a consensus began to emerge, whereby the provenance research community needed to understand better the capabilities of the different systems, the representations they used for provenance, their similarities, their differences, and the rationale that motivated their designs. Hence, the first Provenance Challenge was born, and from the outset, the challenge was set up with an aim of being *informative* rather than *competitive*. In this editorial, we describe the challenge and provide a view on the contributions by the participating teams.

2. The Provenance Challenge

The challenge was defined by Simon Miles and Luc Moreau (U. of Southampton) and Mike Wilde and Ian Foster (U. of Chicago/Argonne Nat. Lab.) in May 2006; it was then reviewed by a larger group, including Juliana Freire (U. of Utah) and Jim Myers (NCSA), before a public review period by the IPAW'06 participants. It was published on June 19th and concluded by a two-day workshop, at GGF 18, in Washington, DC, on Sep. 13-14, 2006.

2.1. Instructions to Participants

The provenance challenge aimed to establish an understanding of the capabilities of available provenance-related systems and, in particular, examine the following.

- The representations that systems use to document details of processes that have occurred;
- The capabilities of each system in answering provenance-related queries;
- What each system considers to be within scope of the topic of provenance (regardless of whether the system can yet achieve all problems in that scope).

To help achieve the aims, a simple example workflow was defined to form the basis of the challenge. This workflow is inspired from a real experiment [12], in the area of Functional Magnetic Resonance Imaging (fMRI). Here, the term *workflow* [5] is used to denote a series of procedures being performed in a system, each taking some data as input and producing other data as output. The procedures and workflows in the challenge problem are not defined



in terms of any particular technology (e.g., EXE files, Web Services, in the case of procedures; BPEL, compiled executable, batch file, in the case of workflows). Instead, participants could adopt their technology of choice.

Our focus in this challenge was on provenance and not on running the experiment. Hence, to facilitate take-up, we allowed challenge participants to implement procedures as “dummies”, i.e., as fake procedures that make use of the input, output, and intermediate data we provided, take the right input and produce the right output, but do not execute real code. Alternatively, participants could execute the real workflow after installing the necessary libraries.

Different systems use different representations for provenance information. In order to explore the capabilities of these different representations, we also defined a set of core queries, and asked participants to show how they addressed those queries.

Challenge participants were invited to upload the following information to the Provenance Challenge TWiki [3], to then allow comparison.

- Representations of the workflow in their system.
- Representations of provenance for the example workflow.
- Representations of the result of the core queries.
- Contributions to a matrix of queries vs. systems, indicating for each whether: (1) the query can be answered by the system, (2) the system cannot answer the query now but considers it relevant, (3) the query is not considered relevant to the project.

Each participant was also invited to optionally contribute the following.

- Additional queries (beyond the core queries) that illustrate the scope of their system;
- Extensions to the example workflow that the participant feels illustrates the unique aspects of their system;
- Any categorisation of queries that the project considers to have practical value.

2.2. The Provenance Challenge FMRI workflow

The purpose of the challenge workflow is to create population-based “brain atlases” from the fMRI Data Center’s archive of high resolution anatomical data. It comprises procedures and data items flowing between them, respectively shown as ovals and rectangles in Figure 1. The workflow can be seen as having five stages, where each stage is depicted as a horizontal row of the same procedure in the figure. Note that the term “stage” is introduced only to help description of the workflow; we do not specify how “stages” should be realised in a concrete implementation.

Individual procedures employ the AIR (automated image registration) suite (bishopw.loni.ucla.edu/AIR5/index.html) to create an averaged brain from a collection of high resolution anatomical data, and the FSL suite (www.fmrib.ox.ac.uk/fsl) to create 2D images across each sliced dimension of the brain. In addition to the data items shown in the figure, there are other inputs to procedures (constant string options), details of which can be found on the Challenge TWiki [3].

The inputs to a workflow are a set [6] of new brain images (Anatomy Image 1 to 4) and a single reference brain image (Reference Image). All input images are 3D scans of a brain of

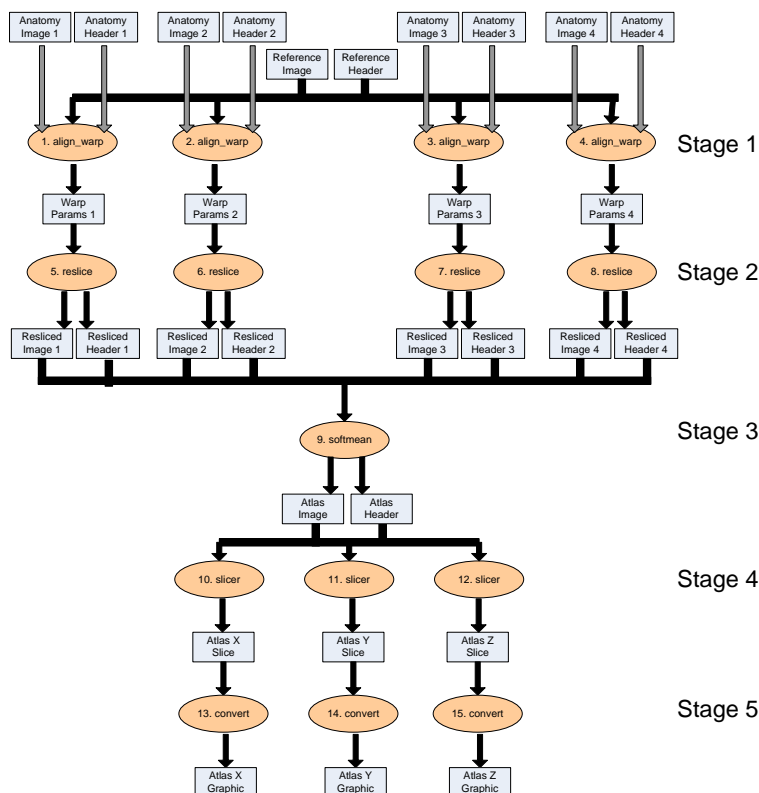


Figure 1. The Provenance Challenge Workflow

varying resolutions, so that different features are evident. For each image, there is the actual image and the metadata information for that image (Anatomy Header 1 to 4).

The stages of the workflow are as follows.

1. For each new brain image, `align_warp` compares the reference image to determine how the new image should be warped, i.e., the position and shape of the image adjusted, to match the reference brain. The output of each procedure in the stage is a *warp parameter set* defining the spatial transformation to be performed (Warp Params 1 to 4).
2. For each warp parameter set, the actual transformation of the image is done by `reslice`, which creates a new version of the original new brain image with the configuration defined in the warp parameter set. The output is a `reslice`'d image.
3. All the `reslice`'d images are averaged into one single image using `softmean`.



4. For each dimension (x, y and z), the averaged image is sliced, with the utility `slicer` , to give an atlas data set, i.e., a 2D atlas along a plane in that dimension, taken through the centre of the 3D image.
5. Each atlas data set is converted into a graphical atlas image using (the ImageMagick utility) `convert` .

2.3. Core Provenance Queries

In addition to the workflow, the challenge specified an initial set of provenance-related queries.

- Q1.** Find the process that led to Atlas X Graphic / everything that caused Atlas X Graphic to be as it is. This should tell us the new brain images from which the averaged atlas was generated, the warping performed etc.
- Q2.** Find the process that led to Atlas X Graphic, excluding everything prior to the averaging of images with `softmean` .
- Q3.** Find the Stage 3, 4 and 5 details of the process that led to Atlas X Graphic.
- Q4.** Find all invocations of procedure `align_warp` using a twelfth order nonlinear 1365 parameter model (see model menu describing possible values of parameter “-m 12” of `align_warp`) that ran on a Monday.
- Q5.** Find all Atlas Graphic images outputted from workflows where at least one of the input Anatomy Headers had an entry `global maximum=4095` . The contents of a header file can be extracted as text using the `scanheader AIR` utility.
- Q6.** Find all output averaged images of `softmean (average)` procedures, where the warped images taken as input were `align_warp` 'ed using a twelfth order nonlinear 1365 parameter model, i.e. “where `softmean` was preceded in the workflow, directly or indirectly, by an `align_warp` procedure with argument -m 12.”
- Q7.** A user has run the workflow twice, in the second instance replacing each procedures (`convert`) in the final stage with two procedures: `pgmtoppm` , then `pnmtjpeg` . Find the differences between the two workflow runs. The exact level of detail in the difference that is detected by a system is up to each participant.
- Q8.** A user has annotated some anatomy images with a key-value pair `center=UChicago` . Find the outputs of `align_warp` where the inputs are annotated with `center=UChicago` .
- Q9.** A user has annotated some atlas graphics with key-value pair where the key is `studyModality` . Find all the graphical atlas sets that have metadata annotation `studyModality` with values `speech` , `visual` or `audio` , and return all other annotations to these files.



3. An Analysis of Contributions to the Provenance Challenge

Following its publication, 17 teams responded to the challenge and submitted an entry to the challenge TWiki [3]. This special issue contains each participating team's contribution to the challenge. In this section, we introduce a classification of the different approaches to help the reader gain a better understanding of provenance systems and their differences. To contrast the different approaches, we have identified a set of criteria, which have been grouped according to two categorisations.

Categorisation 1 is concerned with the broad characteristics of provenance systems, such as the environment in which they are embedded and the technologies they use. Such systems are usually developed in the context of research projects that have specific foci: understanding research motivations is also useful to appreciate some design decisions. Given that the purpose of provenance systems is to build a computer-based representation of provenance that can be queried and reasoned over, Categorisation 2 groups criteria pertaining to such representations; these criteria allow the reader to extract some of the fundamental concepts underpinning representations, and therefore, capabilities of systems. All findings are summarised in Figure 2.

Categorisation 1: Characteristics of Provenance Systems

C1.1 Execution Environment In many cases (but not all), provenance systems are embedded in a specific execution environment. The most common environments are workflow systems and operating systems. When embedded in a single execution environment, provenance representation may become (though does not have to be) dependent on the execution technology. On the one hand, such approaches may offer opportunities for optimisation, which indeed were exploited by some teams. On the other hand, it makes representations technology specific, and brings difficulties if applications are composed of several execution environments.

C1.2 Execution Environment (for the challenge) When systems allow for multiple execution environments, we indicate which one was actually used for the challenge.

C1.3 Representation Technology Provenance is represented and stored using a range of technologies, including relational databases (RDBMS), semantic web technologies (RDF, OWL), and internal private formats. Several systems also expose provenance according to an XML view.

C1.4 Query Language Systems offer query interfaces that operate over the stored representation of provenance. In some systems, the supported language is standard, whereas for others, it is purpose-built.

C1.5 Research Emphasis Teams have different research objectives when investigating provenance concepts. Their research may focus variously on techniques for *executing* (E) workflows such as the one defined in the challenge; *recording* (R) a description of a process being executed; *storing* (S) descriptions of process in persistent storage; and/or *querying* (Q) stored descriptions, in a way that captures the user's interest.



	Redux	Mindswap	Kama	JP	myGrid	VisTrails	ES3	ZOOM	RWS	COMAD	PASS	SDG	NCSDJK	NCSCI	VDL	OPA	USC/ISI	
1. Characteristics of Provenance Systems																		
1.1 Execution Environment	Workflow system	Web	Workflow system	Workflow system and job submission	Workflow system	Workflow system	Operating system	Workflow system	Workflow system	Workflow system	Operating system	Workflow system	Visual Program. Env.	Workflow system	Workflow system	Technology	Workflow system	
(actual system)	WWF		XBaya and BPEL	EGEE gLite (inc Condor and Dagman and JDL)	VisTrails system and VisTrails shell	IDL and Bash	Technology independent	Kepler, Ptolemy	Kepler, Ptolemy	Kepler, Ptolemy	Linux	Kepler, Ptolemy	Cyber D2K	Cyber Integrator	VDS		Wings and Pegasus	
1.2 Execution Environment (for the challenge)				EGEE VOCE VO				SQL scripts			Shell script	RDF (external and internal)			Grid	Java	Grid	
1.3 Provenance Representation	RDBMS	OWL	XML View and RDBMS	key-value pairs in RDBMS	XML View and RDBMS	XML View	RDBMS	RDBMS	Internal	Internal and XML view	Internal	RDF (external and internal)	RDF	RDF	RDBMS	Internal and XML view	OWL and RDBMS	
1.4 Query Language	SQL	SPARQL	SQL	JRIS + JPRS + Perl	SQL	ES3 Queries	SQL + transitive closure	Internal Graph OL	Internal Graph OL	Internal Graph OL	Custom	Semantic Web DASH	ITOL	ITOL	X-XML querying	XQuery, POQuery + Java	SPARQL +	
1.5 Research Emphasis	E/R/S/Q	R/S/Q	E/R	R/S/Q	E/R/Q	R/S/Q	S/Q	E/R/Q	E/R/Q	E/R/Q	E/R/S	R/S/Q	Q	Q	E/R/S	R/Q/S	E/R/S/Q	
1.6 Challenge Implementation	Run	Run	Run	Run	Simulated	Run	Simulated	Partial	Partial	Partial	Run	Partial	Run	Run	Run	Partial	Run	
2. Properties of Provenance Representation																		
2.1 Provenance representation	yes	no	no	no	yes	yes	yes	yes	yes	yes	no	yes	no	no	yes	no	yes	
2.2 Data Derivation vs Causal Flow of Events	E	D	E	E	D	D	D	D	D	D	D	D	D	D	D	D	D	
2.3 Arbitrary annotations in scope/implemented	(+AS/+AI)	(+AS/+AI)	(-AS/-AI)	(+AS/+AI)	(+AS/+AI)	(-AS/-AI)	(-AS/-AI)	(-AS/-AI)	(+AS/+AI)	(+AS/+AI)	(+AS/+AI)	(+AS/+AI)	(+AS/+AI)	(+AS/+AI)	(+AS/+AI)	(-AS/+AI)	(+AS/+AI)	
2.4 Time supported/required	(+TS/+TR)	(+TS/+TR)	(+TS/+TR)	(+TS/+TR)	(+TS/+TR)	(+TS/+TR)	(+TS/+TR)	(+TS/+TR)	(+TS/+TR)	(-TS/-TR)	(+TS/+TR)	(+TS/+TR)	(+TS/+TR)	(+TS/+TR)	(+TS/+TR)	(+TS/+TR)	(+TS/+TR)	
2.5 Naming required (if yes, then what)	keys for ports and data	no	GUIDs for data, services, workflows	no	no	changes to workflow	no	no	no	no	no	no	no	no	no	no	logical file names and file domain metadata attributes	
2.6 Tracked data, and granularity	port level provenance contents	all data	Any GUID assignable data	I/O of any type of workflow	file or process log	file or process	uniform streams of tokens	collections of tokens	collections of tokens	collections of tokens	file or process	I/O of any type of content	any relationship content (contents stored)	any relationship content (contents stored)	file	anything	files and nested file collections reusable templates - workflow instances - execution details	
2.7 Abstraction mechanisms	layered provenance and model	grouping on files, processes, workflows	aspects	layered provenance model	script or job steps	user view of composite	user view of composite	user view of composite	user view of composite	user view of composite	user view of composite	user view of composite	user view of composite	user view of composite	user view of composite	user view of composite	user view of composite	user view of composite

Figure 2. Summary of Contributions



C1.6 Challenge Implementation Some teams executed the challenge workflow (run); others executed the challenge with fake image processing components (partial), making use of the data and intermediary results published with the challenge definition; finally, others fully simulated its execution (simulated).

Categorisation 2: Properties of Provenance Representation

At some level of abstraction, provenance captures a notion of a causal graph, explaining how a data product or event came to be produced in an execution. However, there are variations on this theme, as indicated by the following criteria.

C2.1 Includes Workflow Representation Some systems assume that an explicit representation of a workflow is part of the provenance representation, whereas others do not have such an assumption, and hence rely on other means to describe executions.

C2.2 Data Derivation vs. Causal Flow of Events Some systems describe derivation of data (e.g., conversion was applied to “.pgm” input to produce “.gif” output), whereas others document causal flow of events (e.g., writing of a file is followed by its opening for reading). Some are capable of characterising both data and event oriented views.

C2.3 Annotations Annotations entered by users may provide valuable information pertaining to data products or executions. While most systems were able support the challenge queries related to annotations, not all systems considered annotations to be in the scope of provenance. In the matrix +AS (resp. -AS) denotes that annotations are in scope of provenance (resp. not in scope), whereas +AI (resp. -AI) indicates annotations were implemented (respectively, not implemented) for the challenge queries.

C2.4 Time A representation of provenance does not have to include time, but it is perceived that it is practical for users to be able to refer to time. Therefore, most systems support a notion of time, so that users can refer to executions or data products according to the time they took place or were produced. However, this requirement brings the challenge of identifying which clock to use, given that distributed clocks may return different times. In the matrix, +TS (resp. -TS) indicates that time is supported for challenge queries (resp. not supported), whereas +TR (resp. -TR) time denotes that time is required (resp. not required) for capturing a correct representation of provenance.

C2.5 Naming In order to be able to identify data products, some systems require each product to be identified by a unique name, typically created during workflow execution; such a name can then be used to query about the provenance of data products. Other systems do not require names to be assigned, but see the identification of data items as a query in itself.

C2.6 Tracked data, Granularity Systems are capable of tracking the provenance of different kinds of data; some introduce restrictions on the granularity of data they can track the



provenance of. For instance, systems may or may not deal with collections, files, bytes, or bits.

C2.7 Abstraction mechanisms When processes or data products are complex, it is useful to describe them with different levels of abstractions, sometimes hiding details of execution or representation, and at other times providing them. Some provenance systems provide support for this, by introducing new concepts in their provenance representation.

4. Conclusions

The rest of this special issue consists of papers describing the different systems summarised in Figure 2. We judge that the provenance challenge was highly successful, as measured by the number of participating teams, the quality of their submissions, and the discussions that resulted during the Challenge workshop. A number of lessons were learned from the challenge.

- At times, provenance queries were considered ambiguous. In the future, it would be interesting to specify them better, more precisely and unambiguously, and to characterise the performance implications of the queries.
- While most participating teams could tackle all queries, it is unclear yet whether they all obtained the same or equivalent answers.
- The community lacks consistent and coherent terminology for provenance-related concepts [8]. A consistent terminology would help outsiders to easily grasp issues and compare systems.

Following discussions at the two-day Challenge workshop, the provenance research community has decided to organise a second provenance challenge to address some of these issues in a systematic manner.

REFERENCES

1. Raj Bose, Ian Foster, and Luc Moreau. Report on the international provenance and annotation workshop (ipaw06). *Sigmod Records*, September 2006.
2. Rajendra Bose and James Frew. Lineage retrieval for scientific data processing: A survey. *ACM Computing Surveys*, 37(1):1–28, March 2005.
3. <http://twiki.ipaw.info/bin/view/Challenge/FirstProvenanceChallenge>, June 2006.
4. Ewa Deelman and Yolanda Gil (Eds.). Workshop on the challenges of scientific workflows. Technical report, Information Sciences Institute, University of Southern California, May 2006.
5. Geoffrey C. Fox and Dennis Gannon. Special issue: Workflow in grid systems. *Concurrency and Computation: Practice and Experience*, 18(10), 2006.
6. Denise Head, Abraham Z. Snyder, Laura E. Girton, John C. Morris, and Randy L. Buckner. Frontal-hippocampal double dissociation between normal aging and alzheimer’s disease. *Cerebral Cortex*, (doi:10.1093/cercor/bhh174), September 2004. fMRI Data Center Accession Number: 2-2004-1168X.
7. Simon Miles, Paul Groth, Miguel Branco, and Luc Moreau. The requirements of recording and using provenance in e-science experiments. *Journal of Grid Computing*, 2006.
8. Luc Moreau. Usage of ‘provenance’: A Tower of Babel. Towards a concept map — Position paper for the Microsoft Life Cycle Seminar, Mountain View, July 10, 2006. Technical report, University of Southampton, June 2006.



-
9. Luc Moreau and Ian Foster, editors. *Provenance and Annotation of Data — International Provenance and Annotation Workshop, IPAW 2006*, volume 4145 of *Lecture Notes in Computer Science*. Springer-Verlag, May 2006.
 10. Luc Moreau, Paul Groth, Simon Miles, Javier Vazquez, John Ibbotson, Sheng Jiang, Steve Munroe, Omer Rana, Andreas Schreiber, Victor Tan, and Laszlo Varga. The Provenance of Electronic Data. *Communications of the ACM*, 2007.
 11. Y. Simmhan, B. Plale, and D. Gannon. A survey of data provenance in e-science. *SIGMOD Record*, 34, 2005.
 12. Yong Zhao, Jed Dobson, Ian Foster, Luc Moreau, and Michael Wilde. A Notation and System for Expressing and Executing Cleanly Typed Workflows on Messy Scientific Data. *Sigmod Record*, 34(3), September 2005.