# Provenance and DataONE: Facilitating Reproducible Science

Bertram Ludäscher, Chris Jones, Lauren Walker

GRADUATE SCHOOL OF LIBRARY AND INFORMATION SCIENCE
The iSchool at Illinois
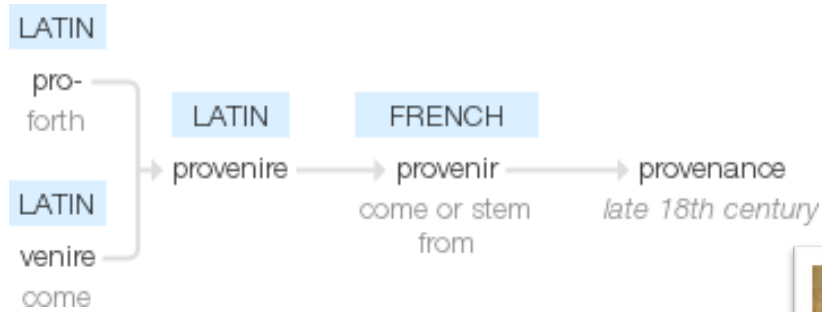
UCSB
UNIVERSITY OF CALIFORNIA
SANTA BARBARA

NCEAS

# Outline

1. **Overview** on Provenance (*Bertram*)
2. **Searching** and **Navigating** Provenance (*Lauren*)
3. **Further Details**, *"look behind the scenes"* (*Chris*)

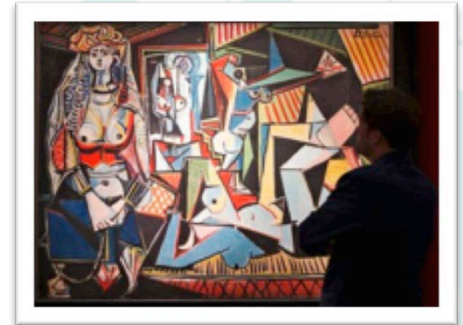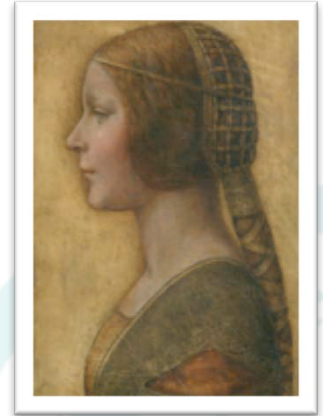**Acknowledgments & special thanks to**:
- NSF, DataONE CI Team, WG members (Phases I & II), others contributors (YW)

# **Provenance**

1. The place of **origin** or **earliest known history** of something
2. The **beginning** of something's existence
3. A **record of ownership** of a work of art or an antique, used as a guide to **authenticity** or **quality**

Related terms: lineage, genealogy, pedigree, …
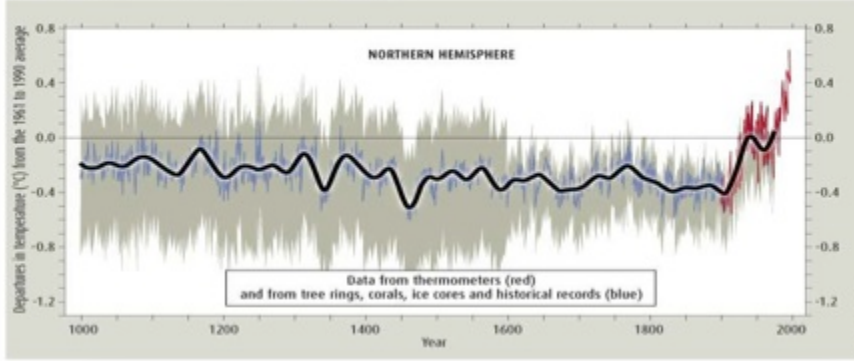
# Computational Provenance

- **Origin** and **processing history** of an artifact
- usually: data products, figures, …
- sometimes: workflow & script evolution …

- **Provenance sightings**:
- Data science, eScience, CI, Big Data, computational science, 4[th] paradigm …
- Bio(diversity)informatics, ecoinformatics, geoinformatics, ..
- Computer science, library & information science, …
- Scientific workflows & scripts …
- Databases, programming languages, …
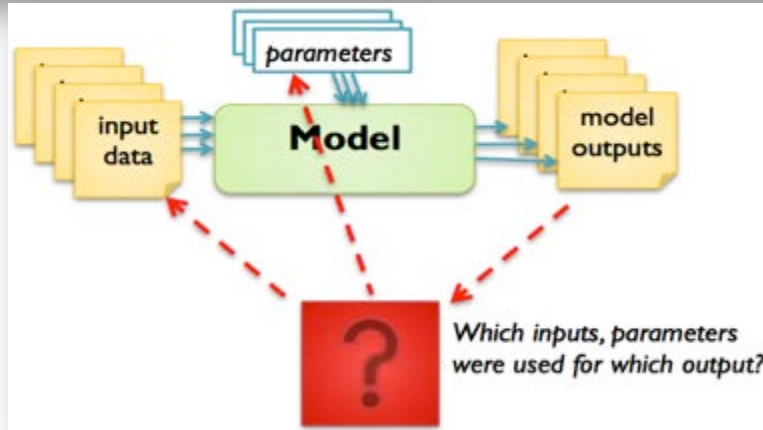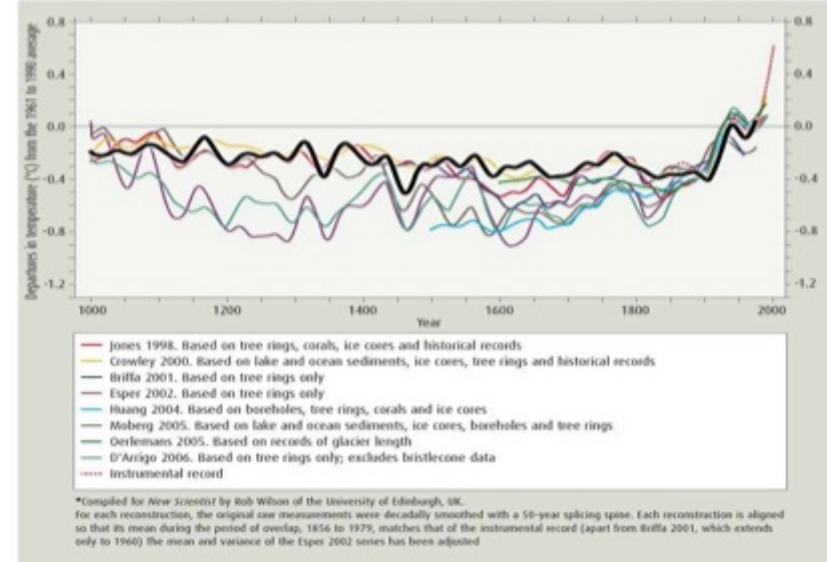- Privacy vs. provenance, …

# Reproducible Science



THE HOCKEY STICK: THE ORIGINAL AND LATER VERSIONS

The 2001 IPCC version: "Variations of the Earth's surface temperature over the past 1000 years"
The error bars (in grey) show the 95 per cent confidence range

NORTHERN HEMISPHERE

Data from thermometers (red)
and from tree rings, corals, ice cores and historical records (blue)



The IPCC version compared with some other northern hemisphere temperature reconstructions*

Jones 1998. Based on tree rings, corals, ice cores and historical records
Crowley 2000. Based on lake and ocean sediments, ice cores, tree rings and historical records
Briffa 2001. Based on tree rings only
Esper 2002. Based on tree rings only
Huang 2004. Based on boreholes, tree rings and ice cores
Moberg 2005. Based on lake and ocean sediments, ice cores, boreholes and tree rings
Oerlemans 2005. Based on records of glacier length
D'Arrigo 2006. Based on tree rings only; excludes bristlecone data
Instrumental record

*Compiled for New Scientist by Rob Wilson of the University of Edinburgh, UK.
For each reconstruction, the original raw measurements were decadally smoothed with a 50-year splicing spine. Each reconstruction is aligned so that its mean during the period of overlap, 1856 to 1979, matches that of the instrumental record (apart from Briffa 2001, which extends only to 1960) the mean and variance of the Esper 2002 series has been adjusted



Capturing **provenance** is crucial for
transparency, interpretation, debugging, …
=> *repeatable experiments,*
=> *reproducible science*

5

# Scientific Workflows: ASAP!

**A**utomation

wfs to **automate** computational aspects of science

**S**caling (exploit and optimize *machine* cycles)

wfs should make use of **parallel compute resources**
wfs should be able handle **large data**

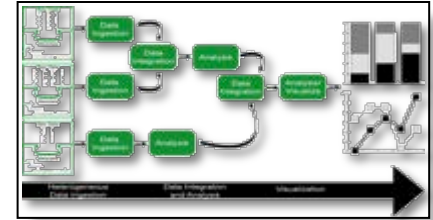**A**bstraction**, Evolution, Reuse** *(human cycles)*

wfs should be easy to **(re-)use, evolve, share**

**P**rovenance

wfs should capture **processing history**, **data lineage**
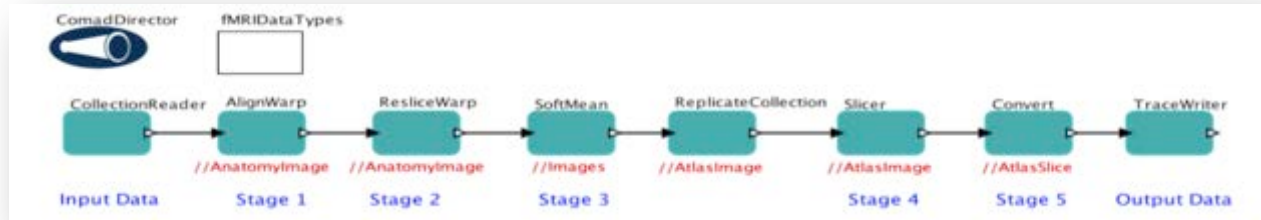=> traceable data- and wf-evolution
=> **Reproducible Science**

# Common Uses of Provenance in Science

- **Audit trail**: trace data generation and possible errors
- **Attribution**: determine ownership and responsibility for data and scientific results
- **Data quality**: from quality of input data, computations
- **Discovery**: enable searching of data, methodologies and experiments
- **Replication**: facilitate repeatable derivation of data to maintain currency

$\Rightarrow$ **Reproducible Science**

# Kinds of Provenance

- **Prospective** Provenance
- method/workflow description ("***workflow-land***")



- **Retrospective** Provenance
- runtime provenance tracking ("***trace-land***")

- **Q**: Which one is more important?

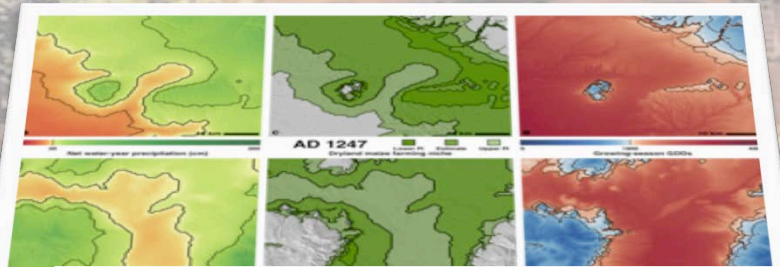# Prospective Provenance
## (A Data Curation Workflow: FilteredPush)

# SKOPE: Synthesized Knowledge Of Past Environments

Bocinsky, Kohler *et al.* study rain-fed maize of **Anasazi**

Four Corners; AD 600–1500. **Climate change** influenced **Mesa Verde Migrations**; late 13th century AD. Uses **network of tree-ring chronologies** to **reconstruct a spatio-temporal climate** field at a fairly high resolution (~800 m) from AD 1–2000. Algorithm estimates joint information in tree-rings and a climate signal to identify "best" tree-ring chronologies for climate reconstructing.

K. **Bocinsky**, T. **Kohler**, A 2000-year reconstruction of the rain-fed maize agricultural niche in the US Southwest. *Nature Communications*. doi:10.1038/ncomms6618



```
203    ## Gene Ontology Statistics are Calculated Here.
204
205    # Gene Ontology Categories that were shown to be relatively Higher (more expressed) in the Experimental Condition.
206    gostatshigher <- higheridrlinkedtogenes[1]
207    higherstatsfilename <- paste(outputDirectory, "/", runName, "_", conditions[1],"_GOStatsHigher_",mytestcond[1], "_v
208    write.table(gostatshigher,file=higherstatsfilename, row.names=FALSE, col.names=FALSE, quote=FALSE, sep="\t")
209    geneListHigherCHR <- gostatshigher$SYMBOL
210    geneListHigherLinkedtoEntrezIds <- select(hgu133plus2.db, keys= geneListHigherCHR, "ENTREZID", "SYMBOL")
211    GOstatsGenesH <- geneListHigherLinkedtoEntrezIds[,2]
212
213    x <- org.Hs.egACCNUM
214    mapped_genes <- mappedkeys(x)
215    xx <- as.list(x[mapped_genes])
216    geneUniverse <- (unique(names(xx)))
```
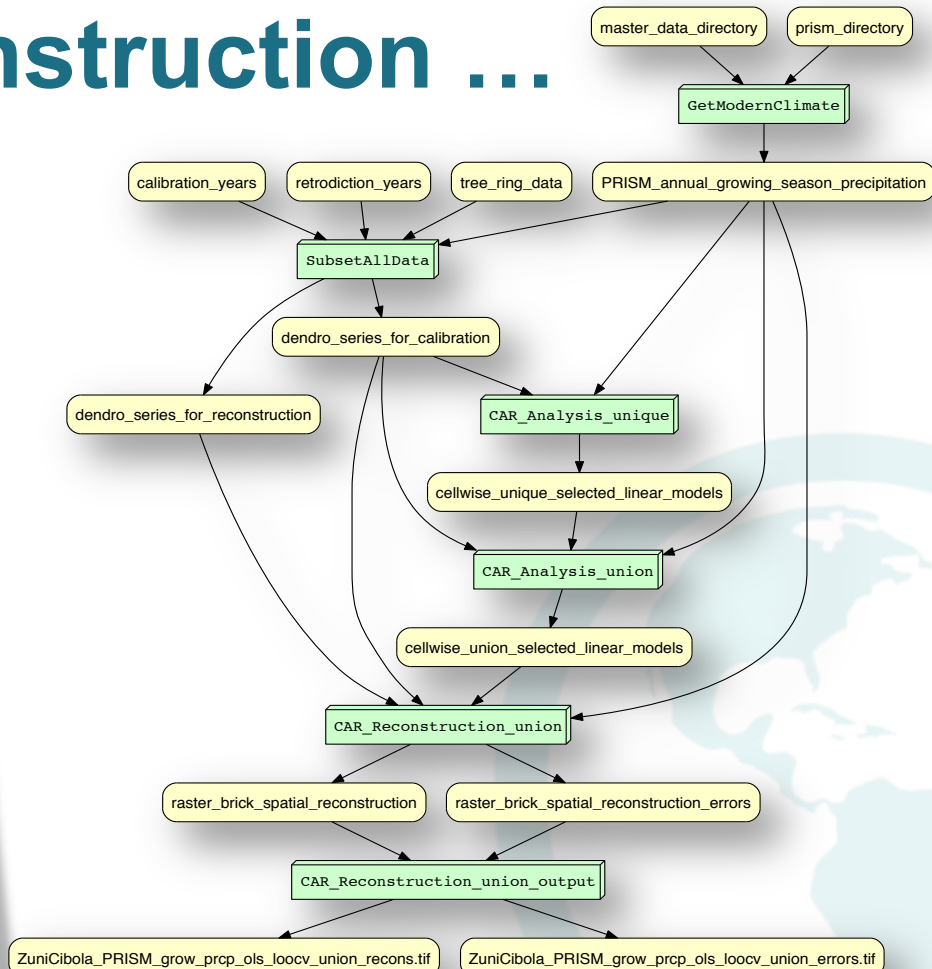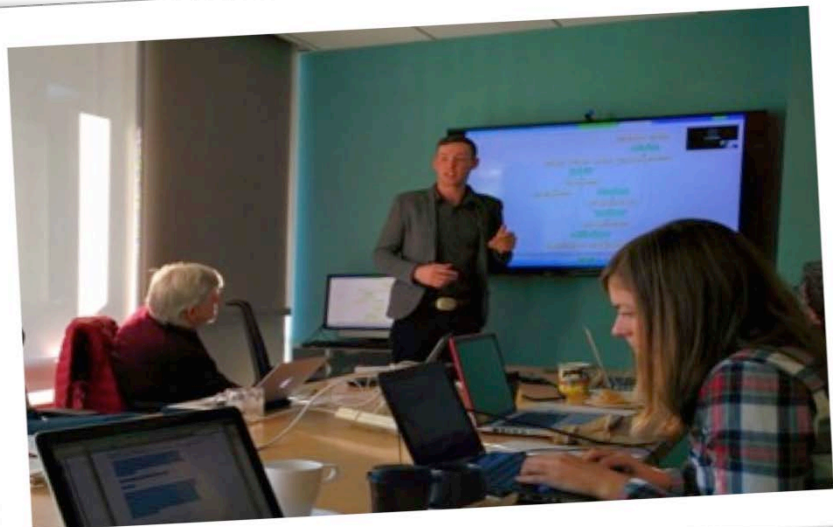
*… implemented as an R Script …*

# Paleoclimate Reconstruction ...

... explained using **YesWorkflow**

**https://github.com/yesworkflow-org/**

Kyle B., (computational) archeologist:

*"It took me about 20 minutes to comment. Less than an hour to learn and YW-annotate, all-told."*

# User Comments: YW Annotations

```
188    ## @begin GO_Analysis
189    # @in hgCutoff @as GO_stats_p_value_cutoff
190    # @in higheridrlinkedtogenes @as DEG_list_higher_in_test_condition
191    # @in loweridrlinkedtogenes @as DEG_list_lower_in_test_condition
192    # @out gostatshigher @as GO_stats_gene_list_higher_in_test_condition
193    # @out BP_SummH_File @as GO_stats_BP_higher_in_test_condition
194    # @out CC_SummH_File @as GO_stats_CC_higher_in_test_condition
195    # @out MF_SummH_File @as GO_stats_MF_higher_in_test_condition
196    # @out gostatslower @as GO_stats_gene_list_lower_in_test_condition
197    # @out BP_SummL_File @as GO_stats_BP_lower_in_test_condition
198    # @out CC_SummL_File @as GO_stats_CC_lower_in_test_condition
199    # @out MF_SummL_File @as GO_stats_MF_lower_in_test_condition
200
201    ######################### Begin GOStats Block #########################
202
203    ## Gene Ontology Statistics are Calculated Here.
204
205    # Gene Ontology Categories that were shown to be relatively Higher (more expressed) in the Experimental Condition.
206    gostatshigher <- higheridrlinkedtogenes[1]
207    higherstatsfilename <- paste(outputDirectory, "/", runName, "_", conditions[1],"_GOStatsHigher_",mytestcond[1], "_vs_",baseline,".
208    write.table(gostatshigher,file=higherstatsfilename, row.names=FALSE, col.names=FALSE, quote=FALSE, sep="\t")
209    geneListHigherCHR <- gostatshigher$SYMBOL
210    geneListHigherLinkedtoEntrezIds <- select(hgu133plus2.db, keys= geneListHigherCHR, "ENTREZID", "SYMBOL")
211    GOstatsGenesH <- geneListHigherLinkedtoEntrezIds[,2]
212
213    x <- org.Hs.egACCNUM
214    mapped_genes <- mappedkeys(x)
215    xx <- as.list(x[mapped_genes])
216    geneUniverse <- (unique(names(xx)))
```

**@begin** GO_Analysis

**@in** hgCutoff
**@in** …

**@out** BP_Summl_file
**@out** …

**@end** GO_Analysis

. . .

12

# Multi-Scale Synthesis and Terrestrial Model Intercomparison Project (MsTMIP)
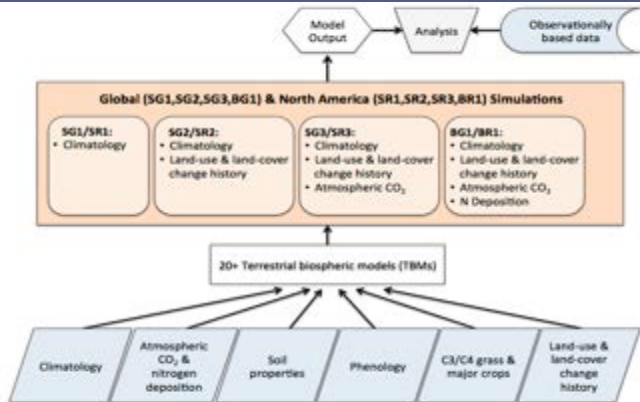


Fig. 1. Schematic of the Multi-Scale Synthesis and Terrestrial Model Intercomparison Project (MsTMIP) framework. Global simulations (SG1,SG2, SG3,BG1) are run at 0.5° by 0.5° resolution; North American simulations (SR1, SR2, SR3, BR1) are run at 0.25° by 0.25° resolution).



Fig. 3. Dendrogram showing general differences/similarities in how MsTMIP models formulate and parameterize (A) energy, (B) carbon, (C) vegetation, and (D) nitrogen process dynamics. Clusters are determined by Hamming distance. Models in the same "tree" share similar structural model characteristics. For example, models in the "tree" to the left (e.g., ISAM, CABLE-JPL, ORCHIDEE-JPL/LSCE) in (A) simulate ground heat flux and canopy heat storage, while models in the "tree" to the right (e.g., MC1, TEM6, VEGAS) do not. A majority of models separate live carbon into various pools (with exception of SiB-JPL), but they do so in various ways (e.g., left "tree" in (B)). Refer to the Supplement for the binary data used to create this diagram.

**Provenance**
- *Externally facing*
- *Internally facing*

**The North American Carbon Program Multi-Scale Synthesis and Terrestrial Model Intercomparison Project**

D. N. Huntzinger1, C. Schwalm2, A. M. Michalak3, K. Schaefer4,5, A. W. King6, Y. Wei6, A. Jacobson4,7, S. Liu6, R. B. Cook6, W. M. Post6, G. Berthier8, D. Hayes6, M. Huang9, A. Ito10, H. Lei11,12, C. Lu13, J. Mao6, C. H. Peng14,15, S. Peng8, B. Poulter8, D. Riccuito6, X. Shi6, H. Tian13, W. Wang16, N. Zeng17, F. Zhao17, and Q. Zhu15
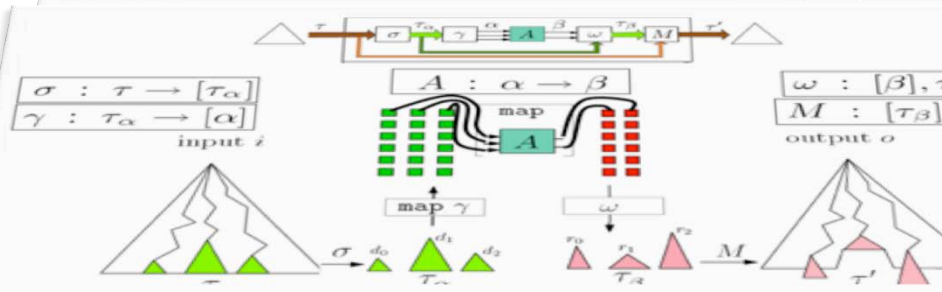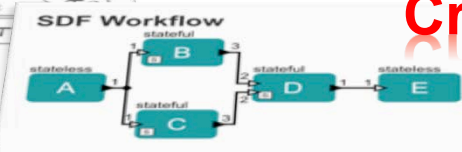
# Scientific Workflow Research



Modeling & Design

Provenance

Fault-Tolerance, Crash Recovery

Parallel Execution

# Wait, there is more …

- **Fine-grained** vs **coarse-grained** provenance
- **Black-box** vs **white-box** provenance

- **Standards**:
- OPM ➜ **PROV**
- D-OPM ➜ **ProvONE**

- **Database Community**:
- why-, where-, how-, why-not provenance
- links to causality
- … logical derivations, proofs, …

# Live Demonstration
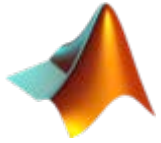
# Facilitate reproducible science

- Creating and managing provenance information

- Communicating script and model workflows

- Storing and sharing

- Using provenance information for search

# Creating and managing provenance information
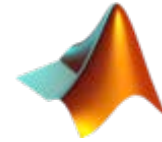
# Investigator tools (ITK)

Matlab DataONE Toolbox

Recordr R Library

Java YesWorkflow Tool

# Functions being added

| |
|---|
| record() |
| startRecord() |
| endRecord() |
| listRuns() |
| deleteRuns() |
| viewRun() |
| publish() |

| |
|---|
| set() |
| get() |
| saveConfig() |
| loadConfig() |
| listConfig() |

See: Run Manager API document

# Example: R programming

```
1   # Generate map of locations by type

2   library(recordr)

3   recordr <- new("Recordr")

4   pkg <- record(recordr, "./hcdbSites.R", "locations-by-type-png")
```

# R: managing script runs

```
>   listRuns(recordr)

Script          StartTime               EndTime                 PublishedTime   Tag                         RunID

hcdbSites.R     2015-05-07T18:53:09Z    2015-05-07T18:53:09Z    unpublished     locations-by-type-png
C85A ...

>   deleteRuns(recordr, "locations-by-type-png")

C85A188-B72E-49F1-AEF4-7BFC24DA186B


>   viewRun(recordr, "locations-by-type-png")

… details about the run listed here ...

>   publishRun(recordr, "locations-by-type-png")

C85A188-B72E-49F1-AEF4-7BFC24DA186B
```
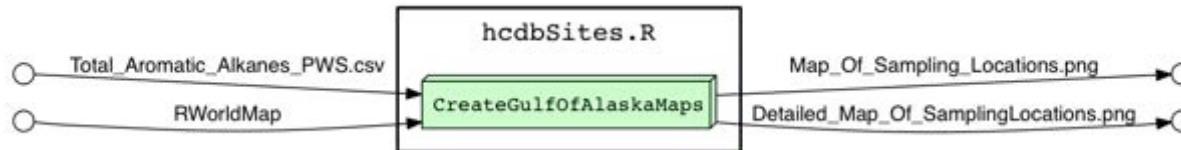
# Communicating
# workflow provenance

# YesWorkflow tool

```
1    # @begin CreateGulfOfAlaskaMaps

2    # @in hcdb @as Total_Aromatic_Alkanes_PWS.csv

3    # @in world @as RWorldMap

4    # @out map @as Map_Of_Sampling_Locations.png

5    # @out detailMap @as Detailed_Map_Of_SamplingLocations.png

     ... mapping code is here ...

25   # @end CreateGulfOfAlaskaMaps
```
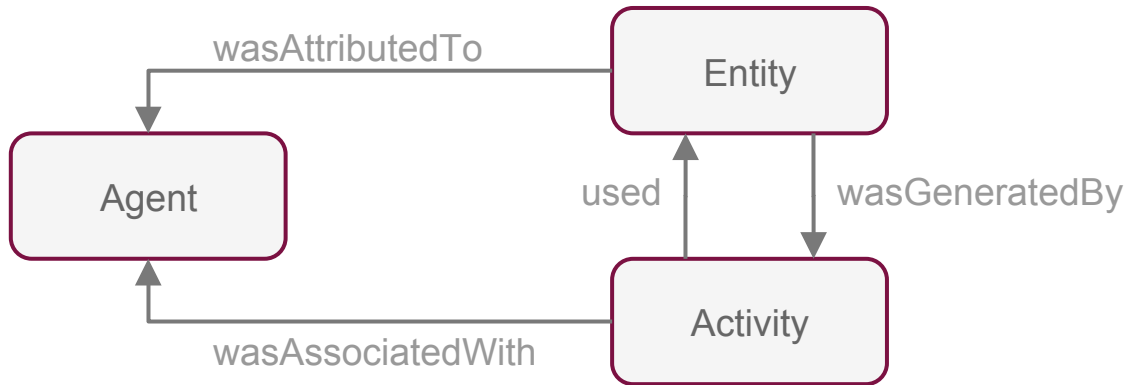
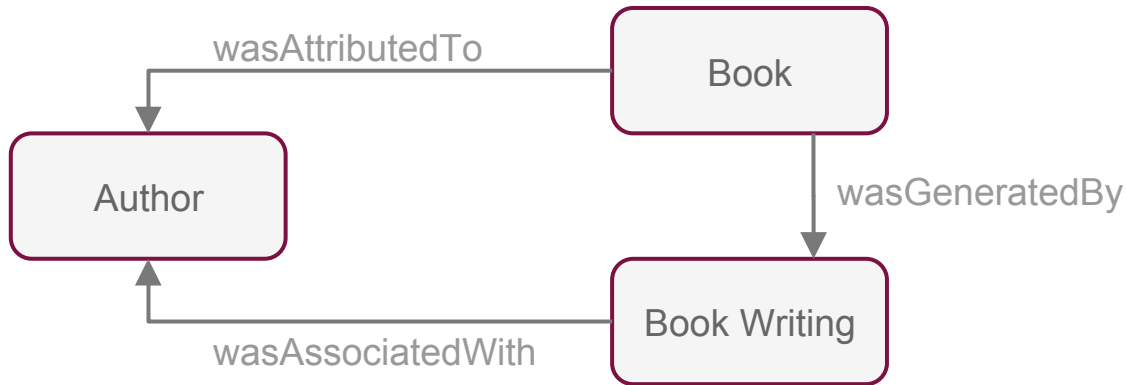# Storing and sharing provenance information

# Using a common model

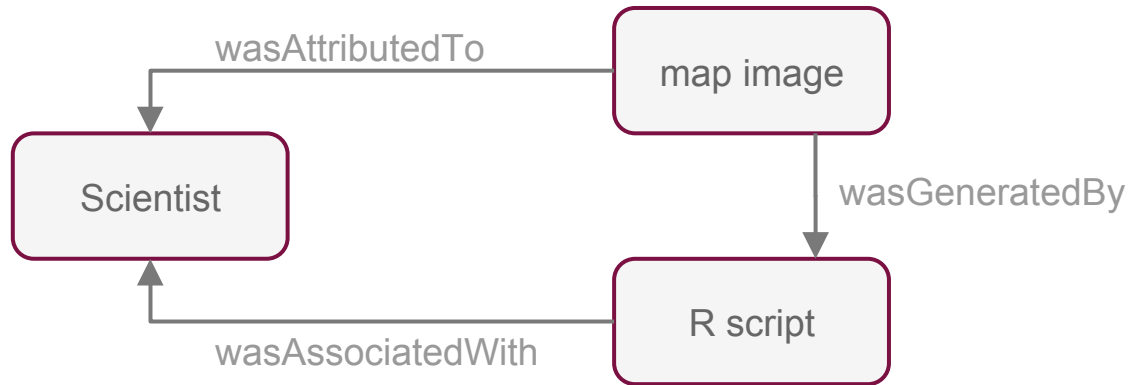W3C has published the 'PROV' family of recommendations



See

# Using a common model

## Example: Book writing activity

# **Using a common model**

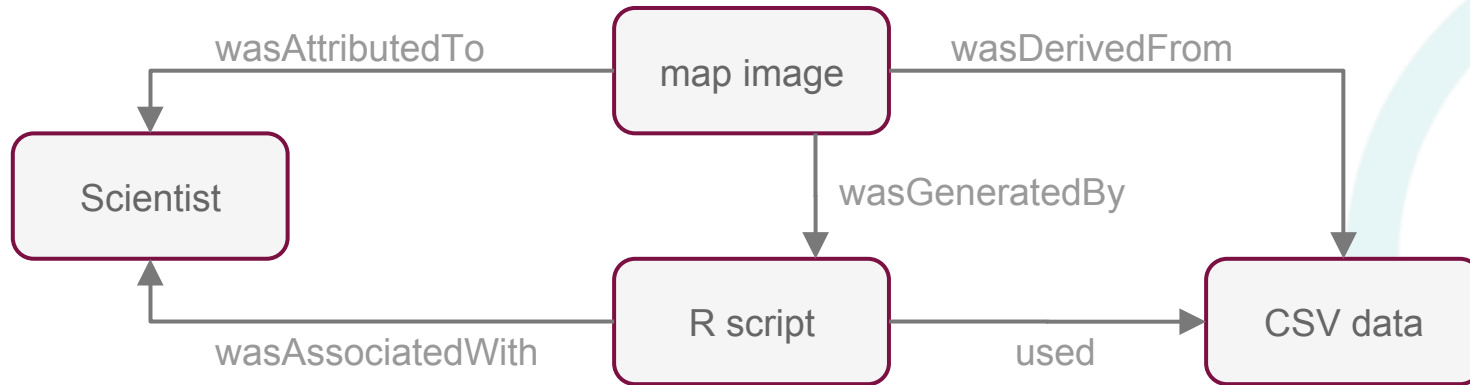## Example: Scientific workflow
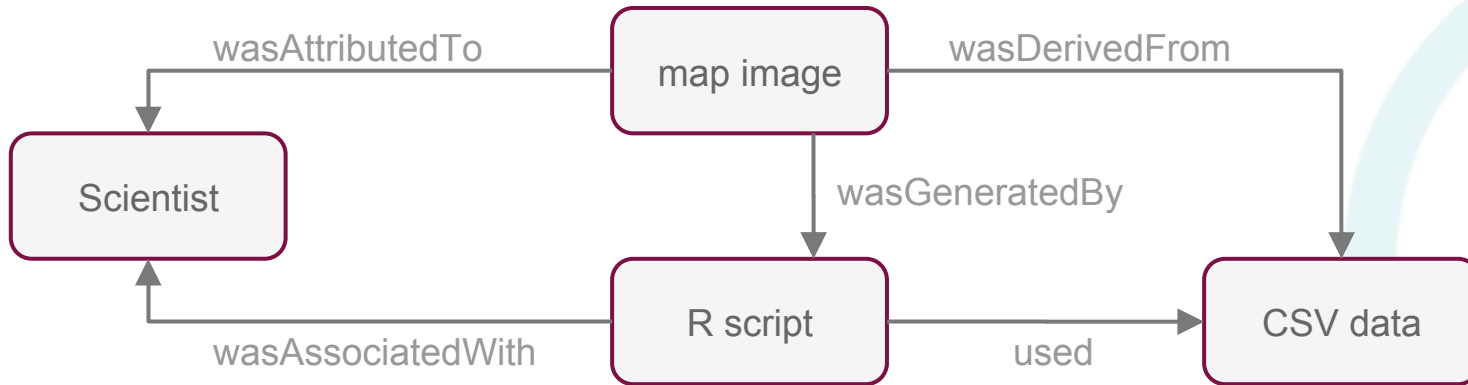
# **Using a common model**

## Example: Scientific workflow

# Using a common model

## Example: Scientific workflow



< "map image" wasDerivedFrom "CSV data" >

```
                      map image
Scientist   wasAttributedTo        wasDerivedFrom
              wasGeneratedBy
              R script      used       CSV data
Scientist   wasAssociatedWith
```
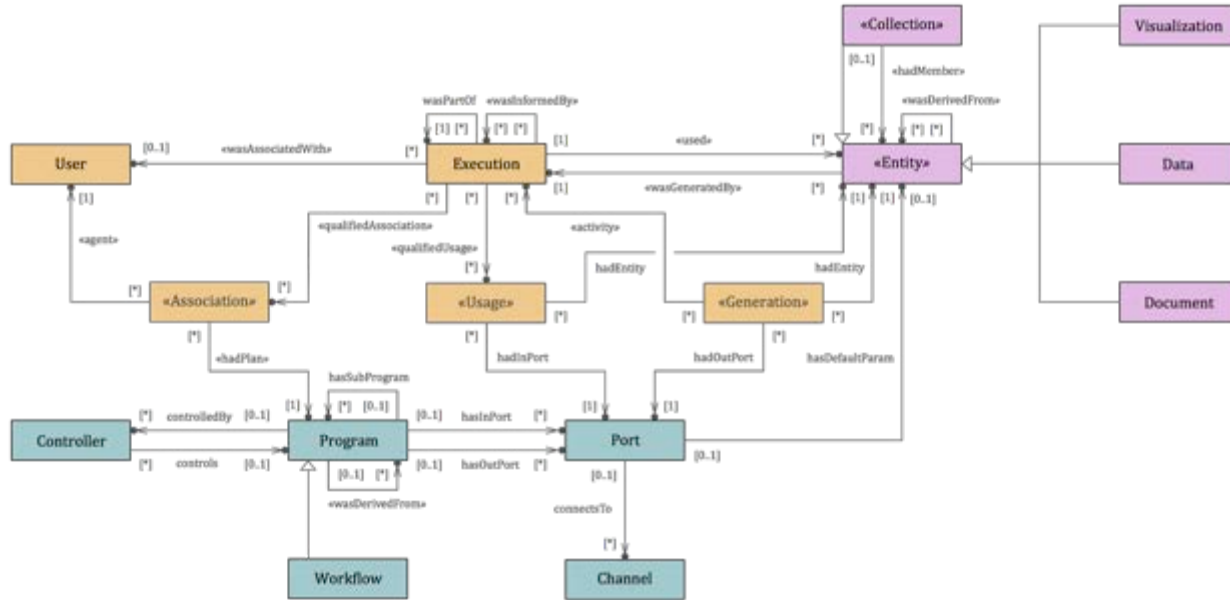
# Using a common model

- W3C PROV model : very generic, universal

- Tracking provenance in scientific workflows

  requires specialization of PROV

- The ProvONE model extends PROV to

  provide this

# ProvONE builds on W3C PROV

# Using provenance information for search

# **Facilitating search**

DataONE harvests provenance information and indexes it



```
ITK Client
   |
   | publish
   v
Member
Node   <---- harvest ---- Coordinating
                          Node
```