# Provenance, Workflows, and Reproducibility
# - *Cui Bono*?

**Bertram Ludäscher**

Director, Center for Informatics Research in Science & Scholarship (CIRSS)

**School of Information Sciences**

**University of Illinois, Urbana-Champaign**

`ludaesch@illinois.edu`

School of **Information Sciences**
The iSchool at Illinois

NCSA

# What's up on Day 1? *(evolving)*

- On DataONE (Bill M)
- On Provenance (yours truly)
- Status/overview of provenance tools (Matt J)
- Goals, agenda check (Dave V, Amber, Kyle B)
- *Break*
- Introductions come first (sort of)
- *Lunch*
- Tool demos & tutorials (hands-on!)
- *... Day 2 ...*

https://github.com/DataONEorg/provathon-2017

https://dataoneorg.slack.com  #prov-a-thon

# A **quick** *Tour de Provenance*

- **Provenance**
  - *... by whom, for whom, for what, how-to*
  - **Prospective** provenance (≈ scientific workflows)
  - **Retrospective** provenance (≈ runtime events, traces)
  - ➔ **Hybrid** provenance ...

- **Reproducibility** & **Transparency**
  - *... of what and for what?*

- **... projects**: DataONE, Whole Tale, SKOPE, ...

- **... and tools:** recordR, YesWorkflow, WT

- My provenance:

  CS/DB@{KIT,Freiburg} .. SDSC .. UC Davis ... {iSchool,NCSA,CS}@UIUC

# Provenance: The Million $$$ Question ..
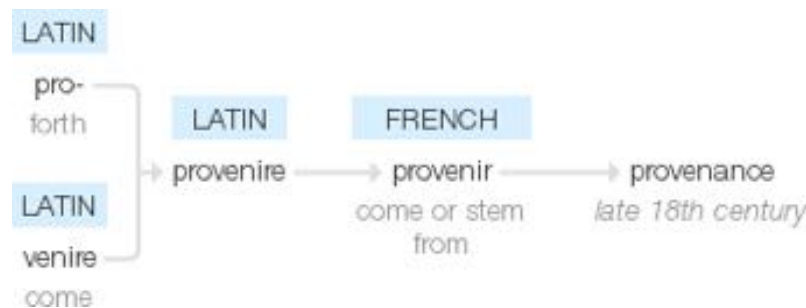




- One of these is has been sold for nearly $180 million.
- The other *could* be worth as much or more.
- Which is which?
- What is the difference?

# **Provenance** defined …

- **Oxford English Dictionary**
  - The **place of origin** or earliest **known history** of something:
    - *an orange rug of Iranian **provenance***
  - The **beginning** of something's existence; its **origin**:
    - *they try to understand the whole universe, its **provenance** and fate*
  - A **record of ownership** of a work of art or an antique, used as a guide to **authenticity** or **quality**:
    - *the manuscript has a distinguished **provenance***
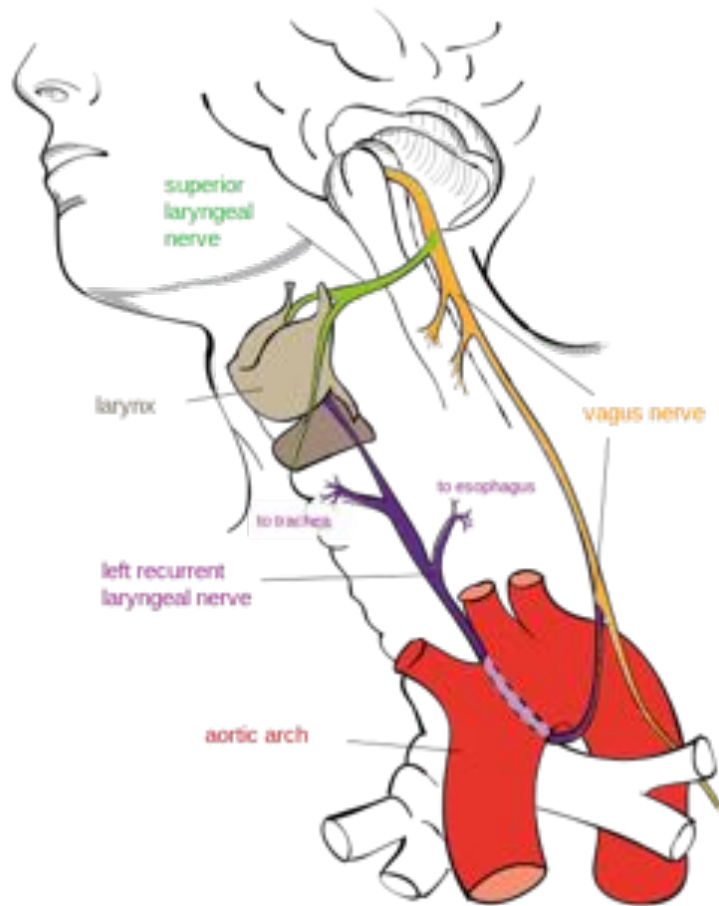- **Similar but different: *Provenience***

# Provenance in the Science Sciences



- Can you "see provenance" in this image?

- Grand Canyon's rock layers are a record of the early **geologic history** of North America. The ancestral puebloan granaries at Nankoweap Creek tell archaeologists about more recent **human history**. (By Drenaline, licensed under CC BY-SA 3.0)

# Science Science: Biology & Natural History
## *Provenance = **Understanding** what happened…*
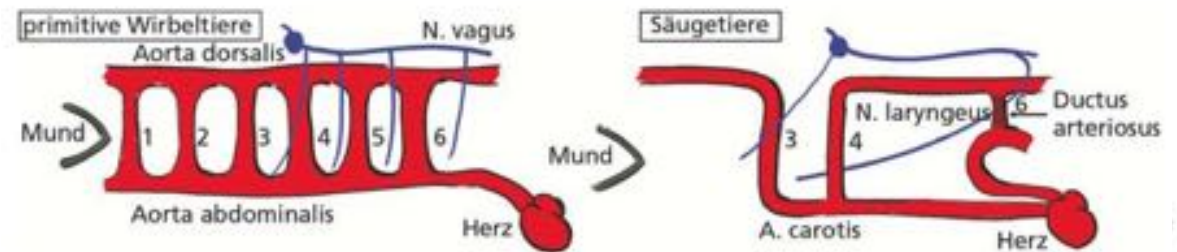


Author: Jkwchui (Based on drawing by Truth-seeker2004)

Zrzavý, Jan, David Storch, and Stanislav Mihulka. Evolution: Ein Lese-Lehrbuch. Springer-Verlag, 2009.

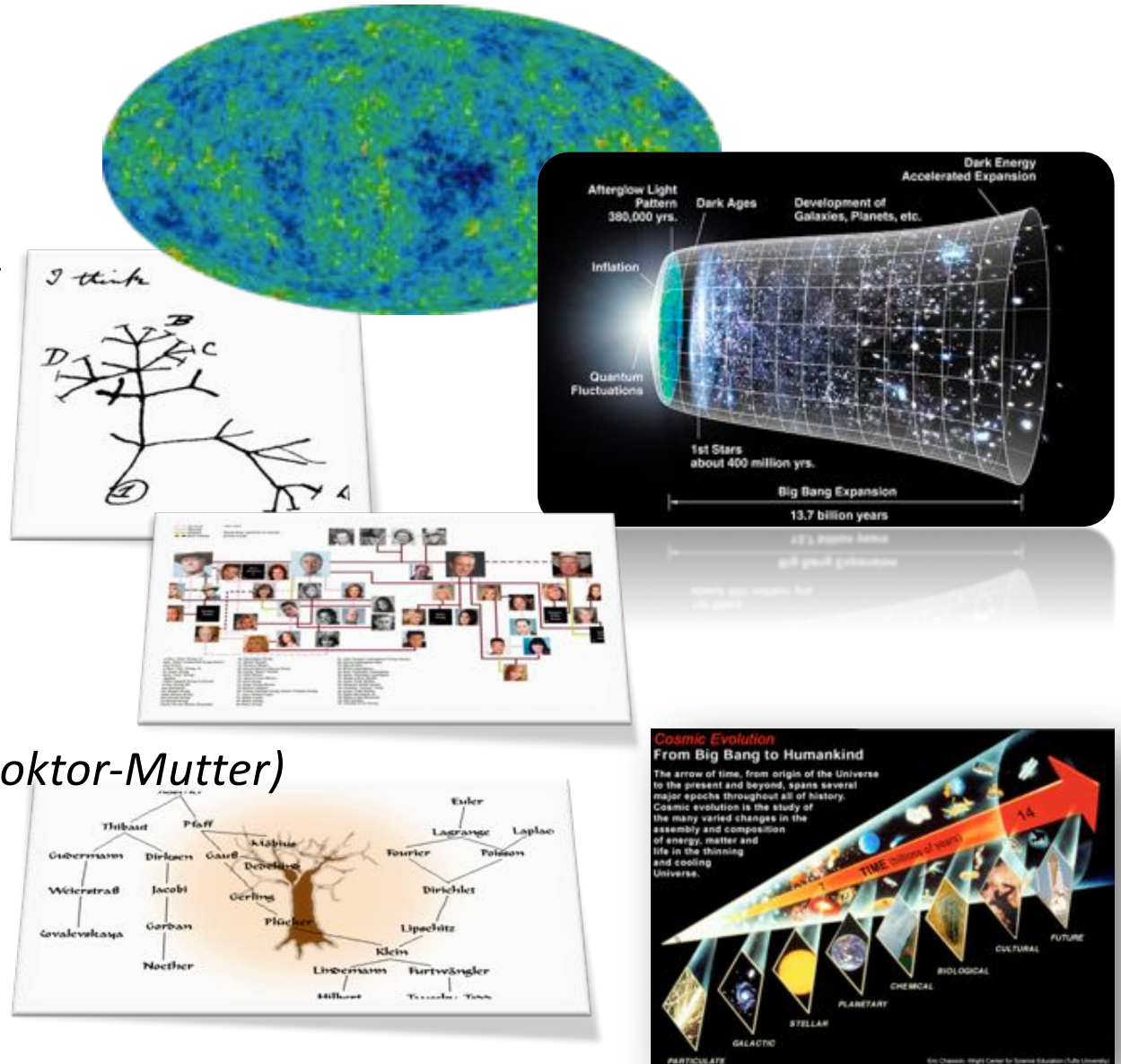**5.17** Suboptimale evolutio-näre Konstruktionslösung:

334    5 Adaptation



**5.16** Evolution der Schleife des rückläufigen Kehlkopfnervs (Nervus laryngeus recurrens) der Wirbeltiere. Dieser Nerv stellt den vierten Ast des Nervus vagus dar. Bei ursprünglichen Wassertieren sandte der Vagusnerv seine Äste zu den Kiemenarterien, die die Bauch- und die Rückenaorta verbanden. Während der Phylogenese der Wirbeltiere haben sich allerdings die Kiemenbögen und mit ihnen auch die Kiemenarterien verändert und das Herz wurde nach kaudal verschoben. Aus der sechsten Arterie wurde bei den Säugetieren der Ductus arteriosus; der vierte Ast des Vagus, der heute den Kehlkopf (Larynx) innerviert, liegt stets *hinter* der ehemaligen sechsten Arterie, also hinter dem Ductus arteriosus. Daher führt dieser Nerv vom Gehirn aus nach hinten, windet sich unter dem Ductus hindurch und kehrt nach vorne zurück, um den Larynx zu innervieren.

# Provenance in Science- Sciences Palooza

- What are those?
- *Cosmology*
- *Geology, Stratigraphy*
- *Phylogeny*
  - the *Tree of Life*
- *Genealogy*
  - your family: literally
- *Academic Pedigree*
  - *"Doktorvater" (oder Doktor-Mutter)*
- *Etymology*
- *Chain of custody*
  - of art(ifacts)
- Yes: all about **origins** and **history** …

# All Science is …

- … physics or stamp collecting
- … ~~physics or stamp collecting~~ **provenance**!

- **Provenance as …**
  - *evidence* (in science … computational science… )
  - *explanation* (in science … computational science .. )
- As we discuss, work with tools, ask …
  - … is this for provenance *recording*?
  - … is this for (method, …) *explaining*?
  - … *by* whom?
  - … *for* whom?

# Why provenance?

- … to *document* what happened
- … to *understand* what happened
- … to *explain* what happened
- … to *anticipate* what (might/will) happen?



Angela Bassa  (Follow)
Professional nerd by day, amateur nerd by night. Opinions all mine. http://angelabassa.com
Aug 28 · 8 min read

## Data Alone Isn't Ground Truth

… and why you should always carry a healthy dose of skepticism in your back pocket.

# ... never forget ...

- *"The government are very keen on amassing statistics. They collect them, add them, raise them to the $n^{th}$ power, take the cube root and prepare wonderful diagrams. But you must never forget that every one of these figures comes in the first instance from the village watchman, who just puts down what he damn pleases."*

SOME ECONOMIC FACTORS
IN MODERN LIFE

BY
SIR JOSIAH STAMP, G.B.E.
*Fellow of the British Academy*

LONDON
P. S. KING & SON, LTD.
ORCHARD HOUSE, WESTMINSTER
1912

# **Computational** Provenance defined …

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

- Different sub-communities:
  - Provenance in **(scientific) workflows** …
  - Provenance in **databases** …
  - … and of course there is more:
    - … programming languages, systems/security, …
    - … information science, archival science, diplomatics
    - … **science science**!

# Using **Provenance** for Transparency, Reproducibility

- What ***input data*** went into this study?

- What ***methods*** were used?

- … with what ***parameter*** settings, ***calibrations***, …?

- *Can we **trust** the data and methods?*



- **Provenance** (*lineage*): track **origin** and **processing history** of data ➔ trust, data quality ~ audit trail for attribution, credit

- **Discovery** of data, methodologies, experiments

# Climate Change: Whodunnit?

# Tracing the sources (data, code)

# Provenance today: Important but *hard*



Climate Change Impacts in the United States

U.S. National Climate Assessment
U.S. Global Change Research Program

"This report is the result of a **three-year** analytical effort by a team of **over 300 experts**, overseen by a broadly constituted Federal Advisory Committee of **60 members**. It was developed from information and analyses gathered in over 70 workshops and listening sessions held across the country."

➔many research projects, groups conduct R&D on provenance methods, tools, ...

Example: DataONE

# A scientific data federation: **DataONE**
# **D**ata **O**bservation **N**etwork for **E**arth



**Data⊙NE**Search
On Provenance

search.dataone.org

# Provenance in Action: Benefits & Impact



*A DataONE search (here: "grass") yields different packages with provenance*

# DataONE: Support for Provenance

*Yaxing's **script** with **inputs** & **output** products*

*Christopher's YesWorkflow model*

*Christopher's results can be traced back all the way to Yaxing's input*

*Christopher **using** Yaxing's outputs as inputs for his script*
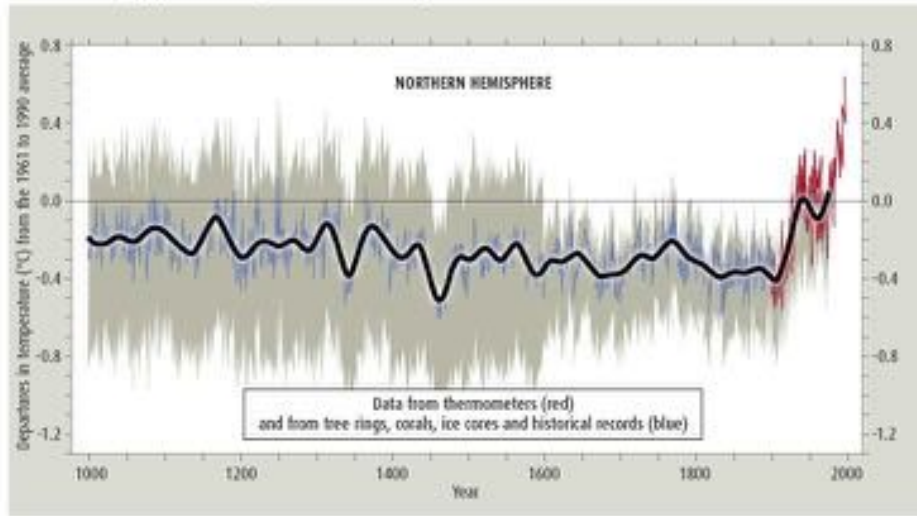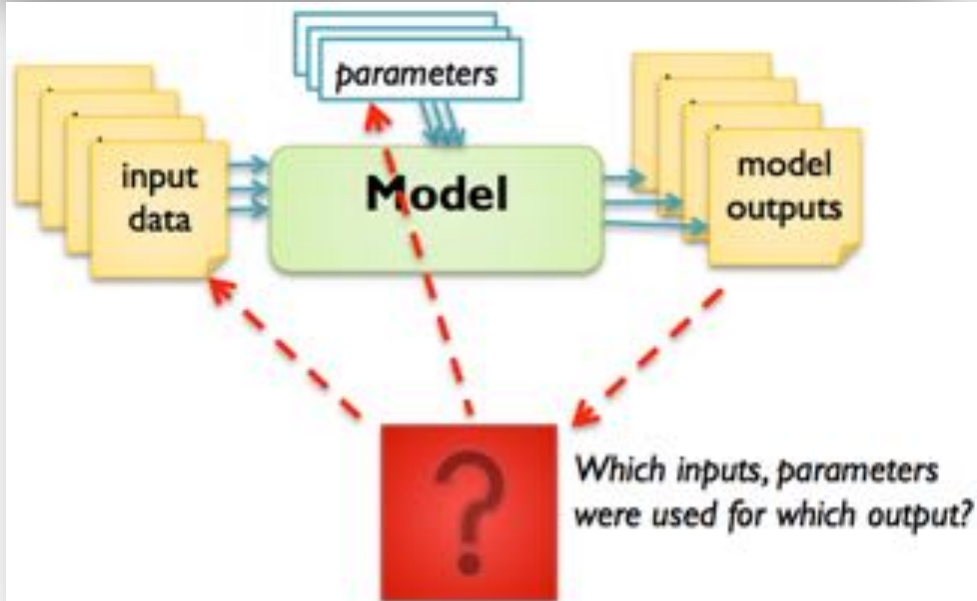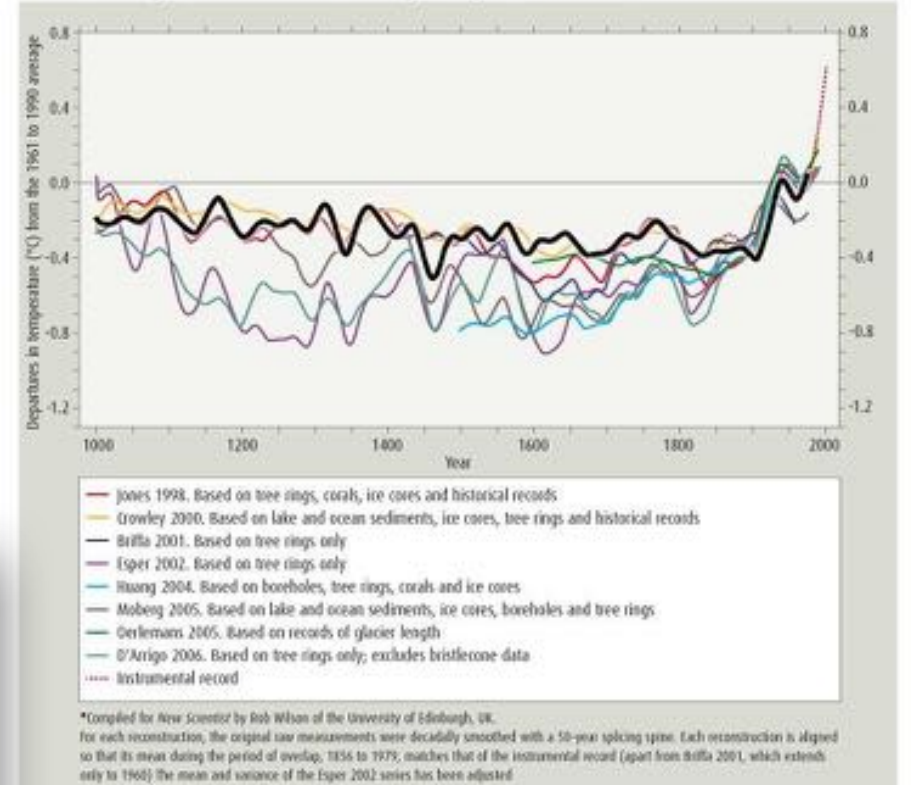
# REWIND: From Provenance to Reproducible Science …



THE HOCKEY STICK: THE ORIGINAL AND LATER VERSIONS

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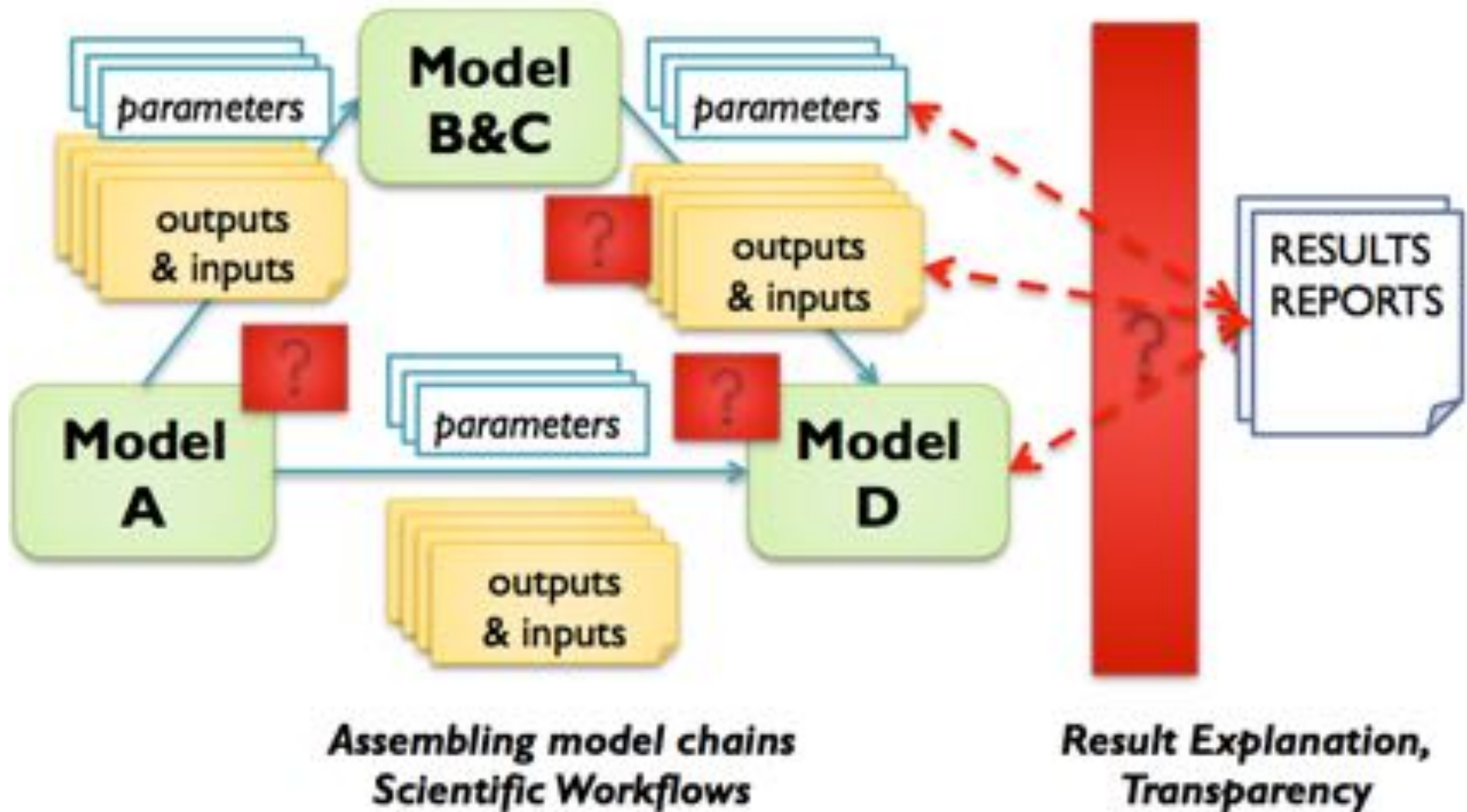The error bars (in grey) show the 95 per cent confidence range

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, corals 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.



Which inputs, parameters were used for which output?

Capturing **provenance** is crucial for transparency, interpretation, debugging, …
=> *repeatable experiments,*
=> *reproducible science*
=> *need workflow–system agnostic model*

# ... via scientific workflows (… and scripts)



Assembling model chains
Scientific Workflows

Result Explanation,
Transparency

# Tour Stop: **Scientific Workflows: ASAP**

- **A**utomation
  - wfs to **automate** computational aspects of science

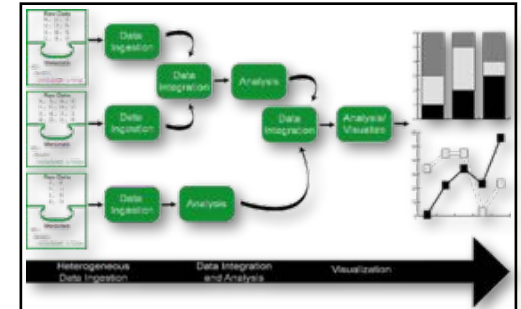- **S**caling (exploit and optimize *machine* cycles)
  - wfs should make use of **parallel compute resources**
  - wfs should be able handle **large data**

- **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
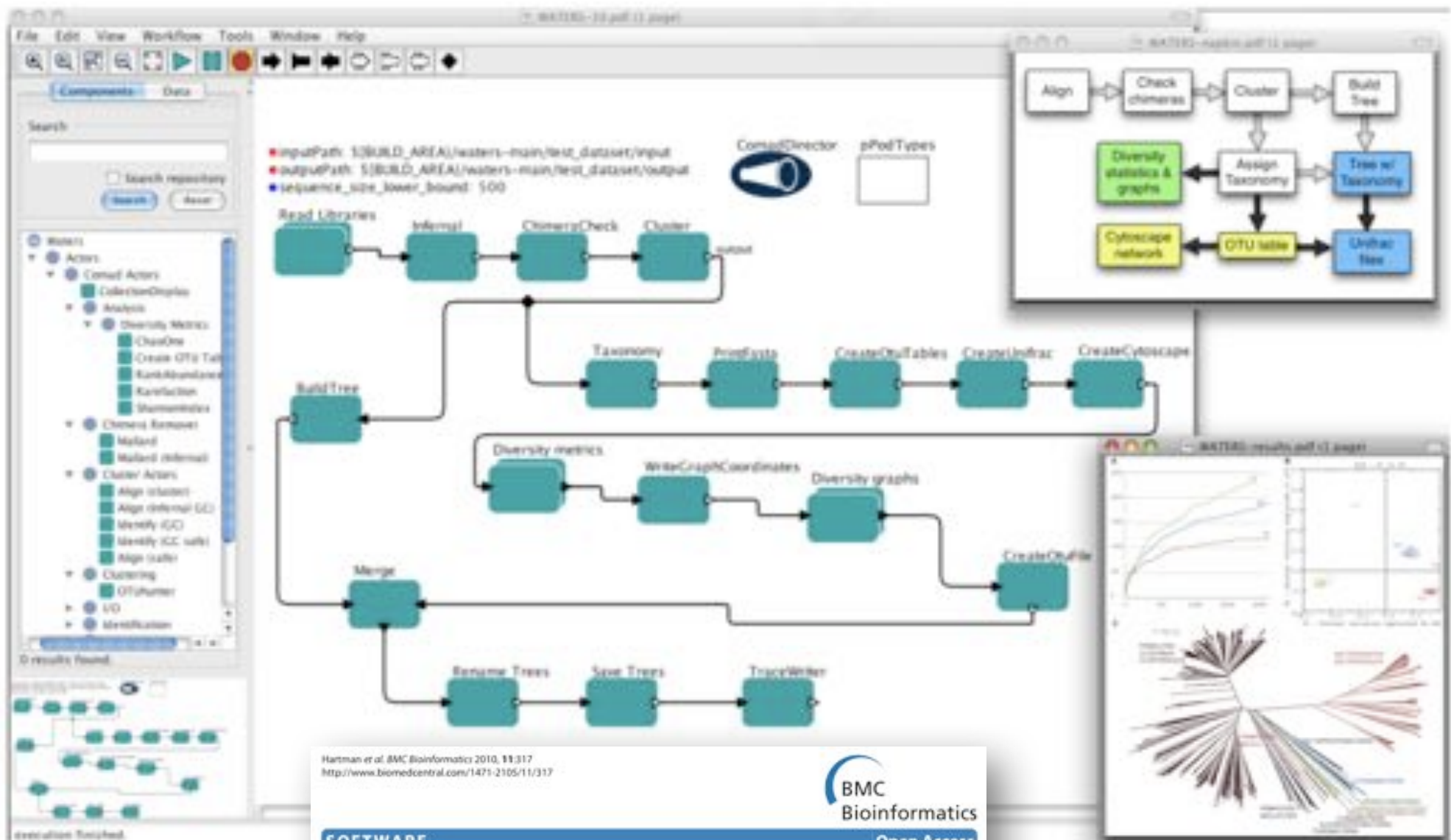  - ➜ **Reproducible Science**

Kepler

swift

Taverna

ASKALON

VisTrails

TRIANA

Trident
Workbench
*Es war einmal ...*

WINGS

pegasus

# Executable WATERS Workflow in Kepler

# Data Curation Workflows
## (Filtered-Push … Kepler … Kurator projects)



Human Curation      Google Cloud Service      Visualization In Google Map      Kepler Provenance Browser

# Provenance ⇔ Workflows

**Workflow Modeling & Design**
*(a.k.a. <span style="color:red">prospective</span> provenance "Workflow-land")*



**Runtime Provenance**
*(a.k.a. traces, logs, <span style="color:red">retrospective</span> provenance, "Trace-land")*



Past ⟵ ⟶ Future

On Provenance

# ProvONE: PROV for scientific workflows DataONE

*(Transfer station to any of several other "standard extensions")*

*"Trace-Land"* (**retrospective** provenance)



*"Data-Land"*

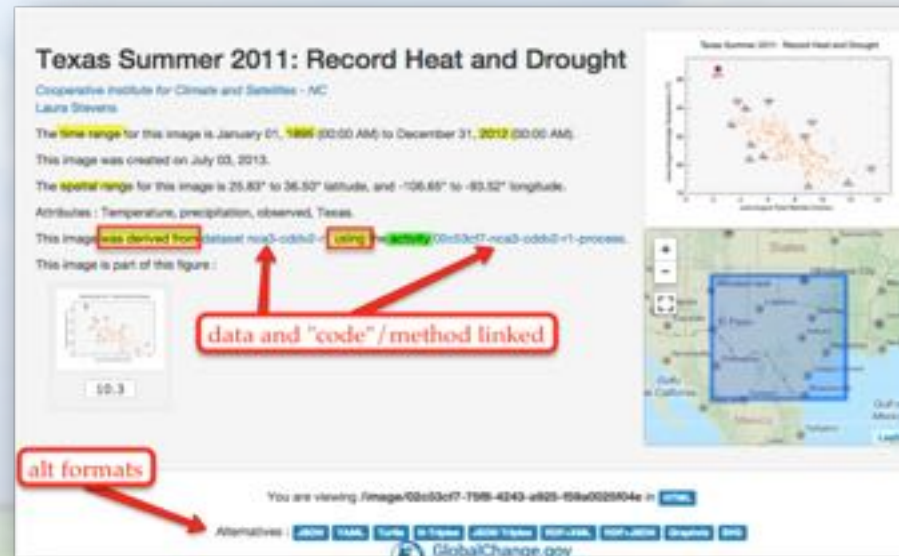*"Workflow-Land"* (**prospective** prov.)

**Also:** OPM-W (G & G et al), others …

# But … **how to prime the provenance pump??**
## Must support **"Provenance for Self"** !



*Provenance for Self?!*

*Provenance for Others*

✔ Provenance *capture* (Matlab, R, Python, … scientific workflow systems)

✔ Uploading, *sharing*, *linking* provenance through various provenance tools

✗ Tools for scientists to **exploit** (≠ *capture*, *share*, *link*) provenance **for their own day-to-day work**.

➔ **Prime the provenance pump** and **increase provenance generation**

➔ Scientists **accelerate their work via new, *active* uses of provenance**.

# From Workflows & Provenance to **Provenance for Script-based Workflows …**

- What workflow tools are (most) scientists using?
  - Workflow systems
  - … vs scripts (Python, R, MATLAB, …)

- What provenance tools are their?
  - Workflow system support
  - Tools for "workflow" scripts!?

# 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
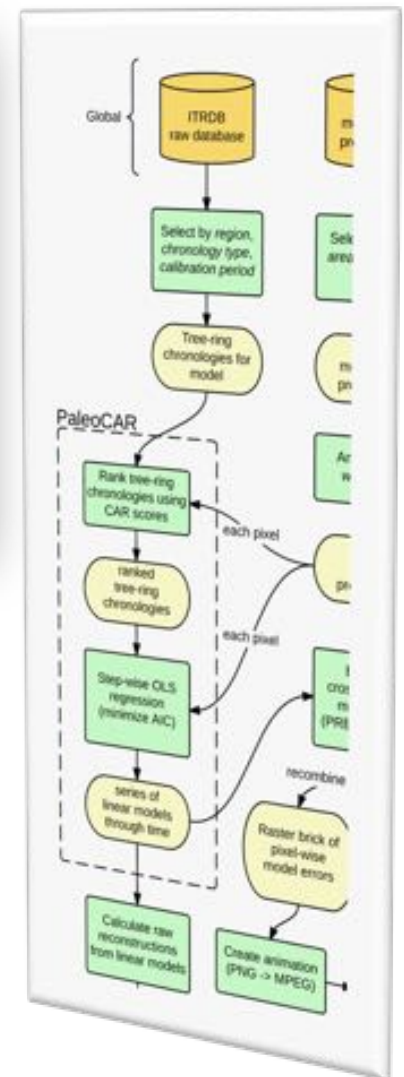
*… implemented as an R Script …*

# Provenance Support for Reproducible Science
# Example: Paleoclimate Reconstruction

Science paper (OA) uses:

- open source code:
  - R, PaleoCAR, …

- Is that all we need?

- What was the "workflow"?

- Is there prospective and/or retrospective provenance?

# YesWorkflow:
## *Yes, scripts are workflows, too!*



- **Script** *vs* Workflows/**ASAP**:
  - **A**utomation: **\*\*\*\*\***
  - **S**caling: **\*\***
  - **A**bstraction: **\***
  - **P**rovenance: **\*\***

# **YesWorkflow: Prospective** & Retrospective Provenance … (almost) for free!

```
# @BEGIN collect_data_set
# @PARAM cassette_id @PARAM accepted_sample @PARAM num_images @PARAM energies
# @OUT sample_id @OUT energy @OUT frame_number
# @OUT raw_image_path @AS raw_image
# @URI file:run/raw/{cassette_id}/{sample_id}/e{energy}/image_{frame_number}.raw
run_log.write("Collecting data set for sample {0}".format(accepted_sample))
sample_id = accepted_sample
for energy, frame_number, intensity, raw_image_path in collect_next_image(
                 cassette_id, sample_id, num_images, energies,
                 "run/raw/{cassette_id}/{sample_id}/e{energy}/image_{frame_number}.raw"):
    run_log.write("Collecting image {0}".format(raw_image_path))
# @END collect_data_set
```
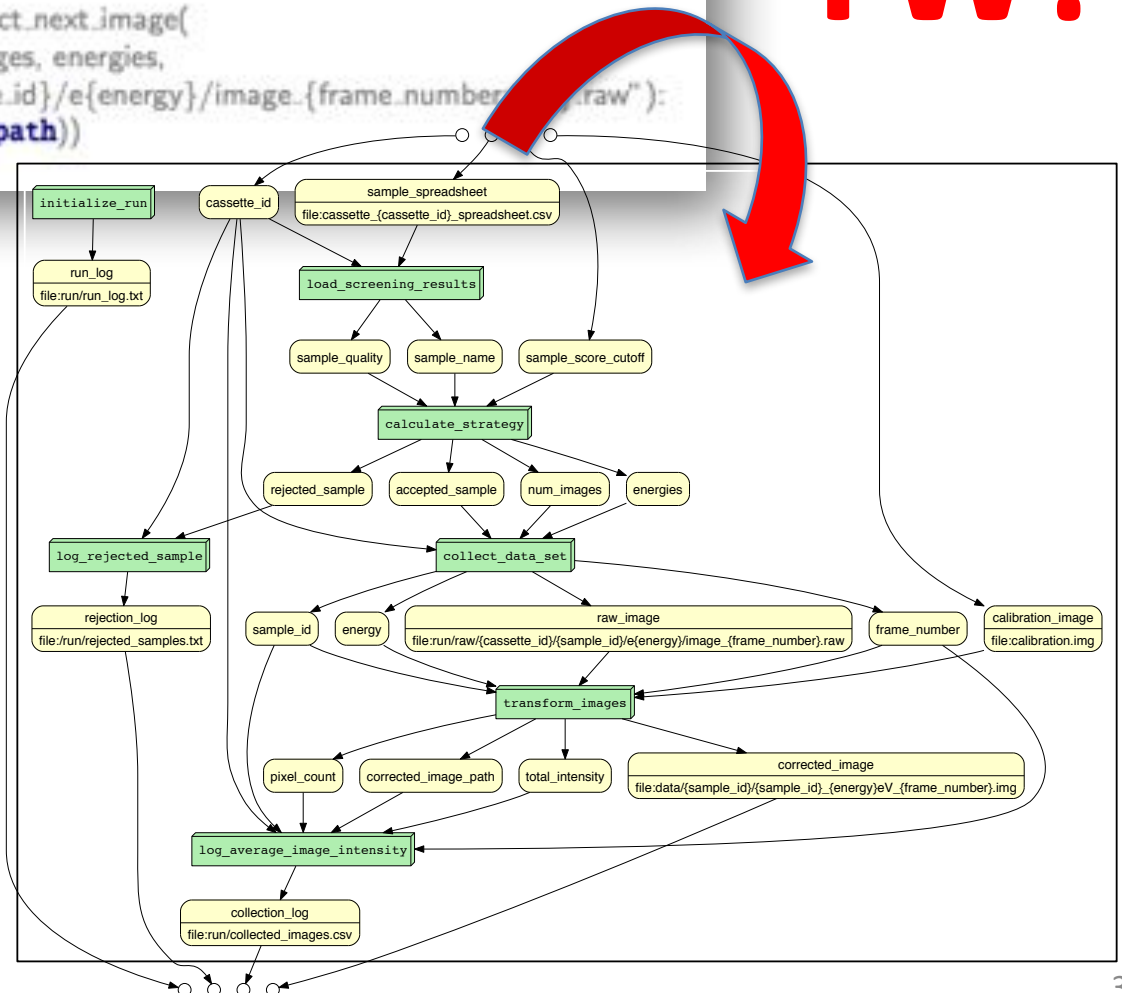
**YW!**

- YW annotations in the script (R, Python, Matlab) are used to **recreate the workflow view** from the script …
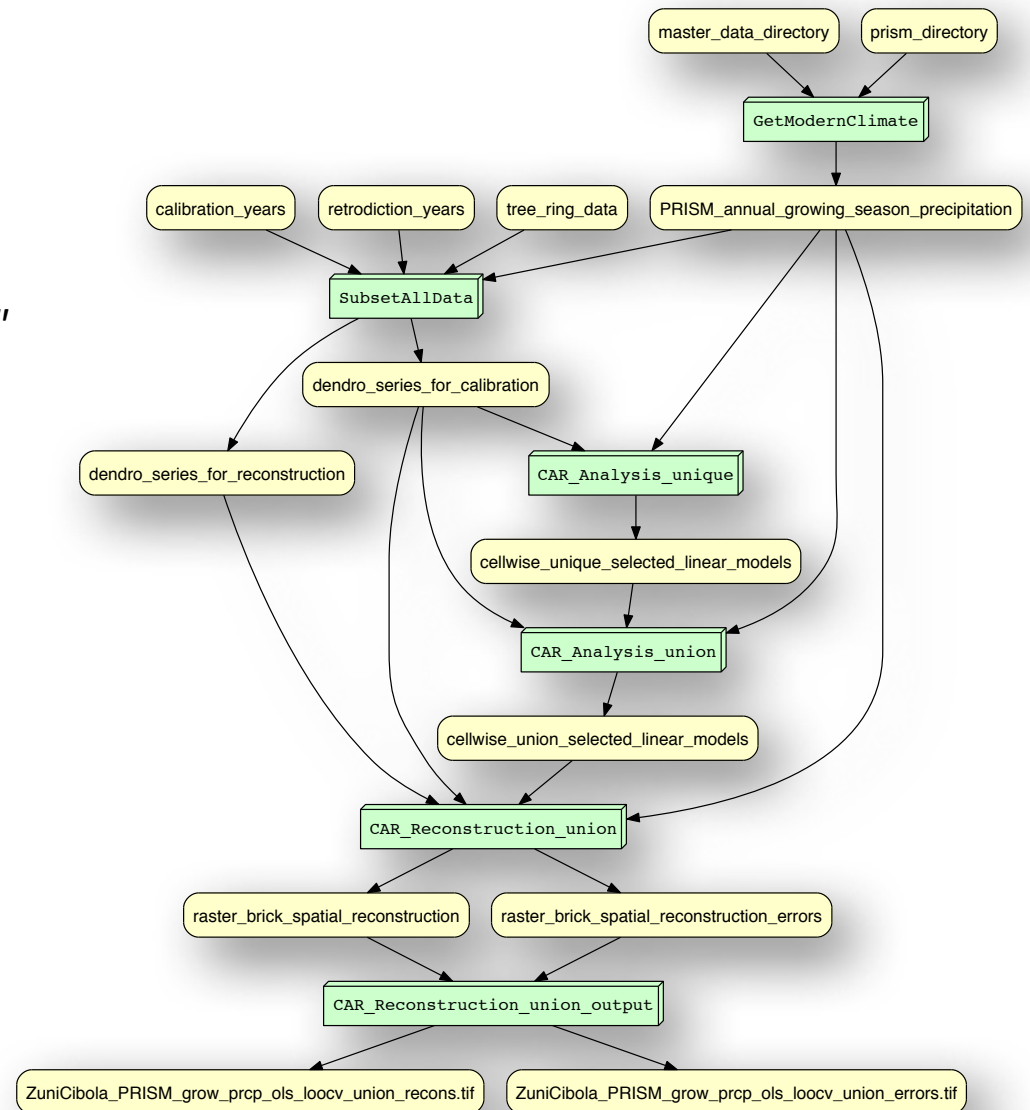
# Paleoclimate Reconstruction (openSKOPE.org)

- … explained using **YesWorkflow!**

Kyle B., (computational) archaeologist:

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



**SKOPE** + **Kurator**

+ **DataONE**
Data Observation Network for Earth

## => YesWorkflow.org

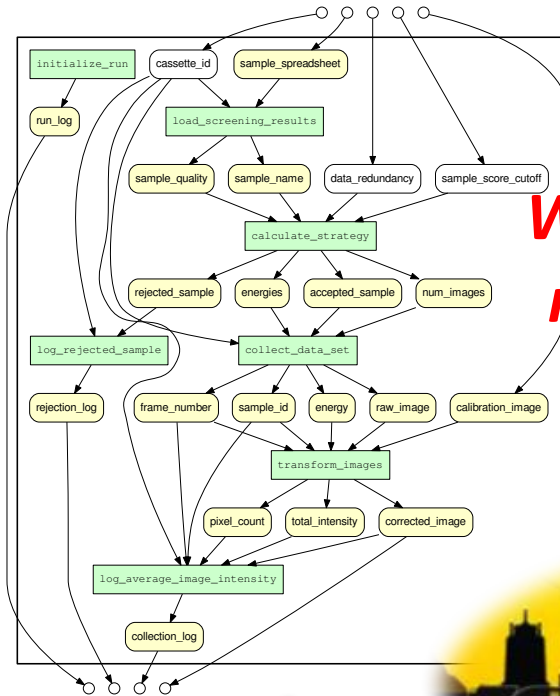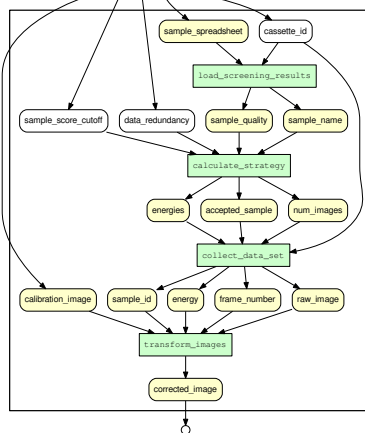# **noWorkflow**:
# **n**ot **o**nly
# **Workflow!**



- **Transparently** capture some/all provenance from Python script **runs**.

- Use filter **queries** to "zoom" into relevant parts ..

**YesWorkflow:**
**Conceptual** *workflow model*

**noWorkflow:**
**Python** *trace model*

*Would like to use YW model to query NW data!*

*lineage query*

*But how do we bridge this gap???*
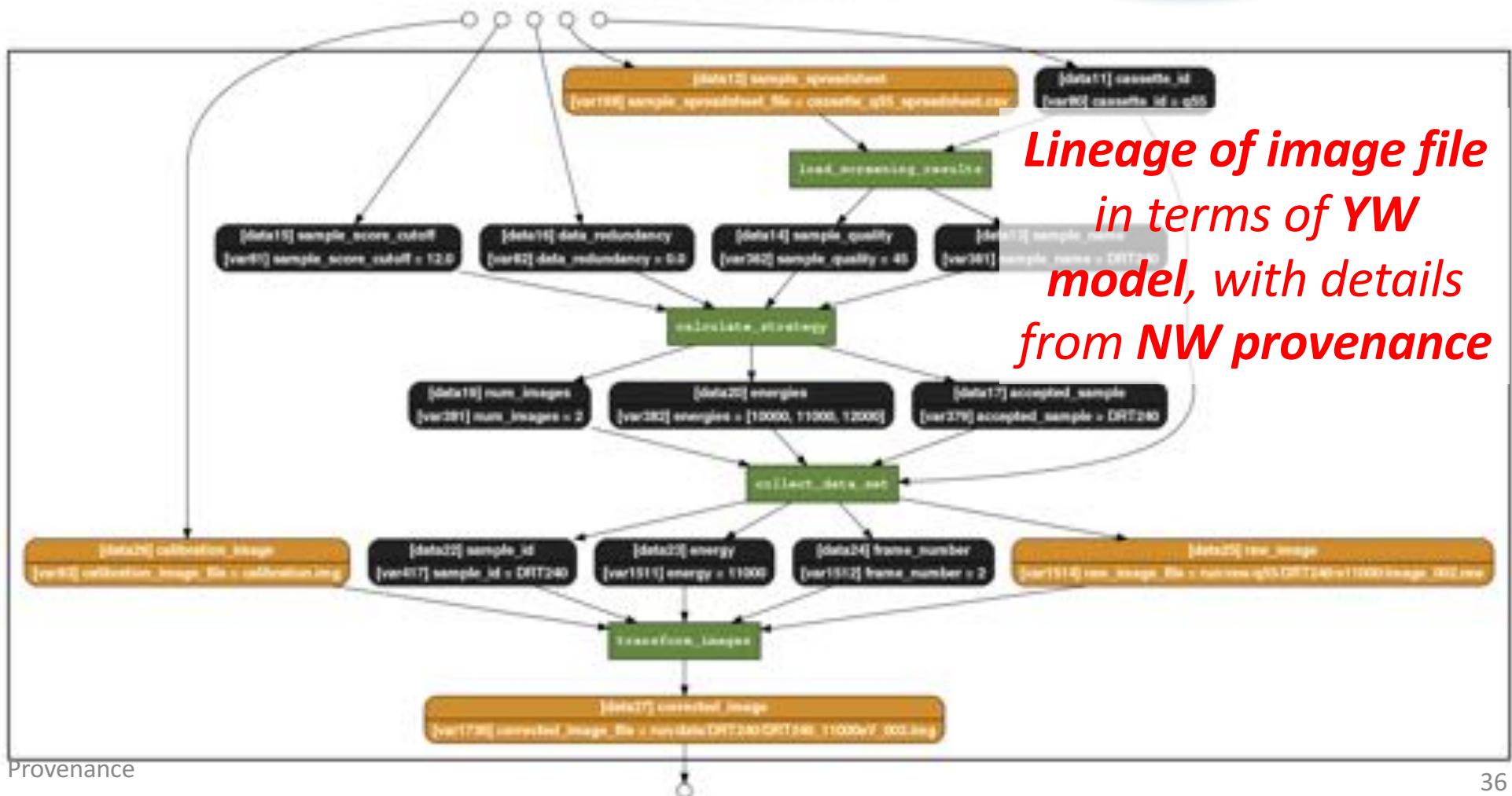
*lineage query*

# Habemus Pons!

## We've got the Bridge!

The bridge is the journey..

(The journey is the destination)



*Lineage of image file in terms of **YW model**, with details from **NW provenance***

# YW-IDCC'17 Demo Use Cases

| Domain | Use case | Programming language | Provenance methods |
|---|---|---|---|
| Climate science | C3C4 | MATLAB | YW + MATLAB RunManager |
| Astrophysics | LIGO | Python | YW + NW (code-level) |
| Protein crystal samples | Simulate data collection | Python | YW + NW (code-level) |
| Biodiversity data curation | kurator-SPNHC | Python | YW-recon + YW-logging |
| Social network analysis | Twitter | Python | YW + NW (file-level) |
| Oceanography | OHIBC Howe Sound (multi-run multi-script) | R | YW + R RunManager |






On Provenance

# Finer-grained Provenance: User Log Files!

# Reproducibility: (yesterday's discussion cont'd)

- What **questions** should we ask?
- What **queries** should we enable?

- **Cui bono?** (others, publishers, ... ?)
- **Provenance-for-Self** vs Provenance-for-others
- **Reproducibility-for-Self** vs Reproducibility-for-others

- For key terms, e.g., Carole's
  - ... rerun, repeat, replicate, reproduce, reuse, ...
- .. ask *"what information/insight do I gain from reproducing, repeating, replicating... ?"*
- *What is fixed and what does the study vary?*

=> **R**esearch **O**bjective, **M**ethod/Algorithm, **I**mplementation, **P**latform/Environment, **A**ctors/People, input **D**ata (params, raw data)

On Provenance

# Reproducibility Crisis *(reprised)*

- **Successful** reproducibility study:
    - **increases trust** in prior study ☺
    - ... but **no surprises** ☹

- **Failed** reproducibility study :
    - **decreases trust** (or *falsifies*) prior study ☹
    - ... but **surprising** failure yields **new info/knowledge** ☺

- Learning from failures!
    – Not really a new, revolutionary idea..
    – What is a positive vs negative result anyways?
    – *... fail early, fail often ...*

# PRIMAD *(what have you "primed"?)*

## 6.1.2 The PRIMAD Model

As a starting point, we defined a preliminary list of "variables" that could potentially be changed:

- (R) or (O) Research Objectives / Goals
- (M) Methods / Algorithms
- (I) Implementation / Code / Source-Code
- (P) Platform / Execution Environment / Context
- (A) Actors / Persons
- (D) Data (input data and parameter values)

This spells: OMIPAD. Rearranging the letters that we use to represent the several aspects that can be changed, it can be remembered as PRIMAD: (P)latform, (R)esearch Goal, (I)mplementation, (M)ethod, (A)ctor, (D)ata (both input and parameter data), which allows us to ask: What variables have you "primed" in your reproducibility study?

*Dagstuhl Seminar #16041 Report*

**Outputs = Exec(M,I,P,D) | RO, A**
- M = parsimony/bootstrap/..
- I = package XYZ
- P = MacOS ..
- D = (Params, Files)

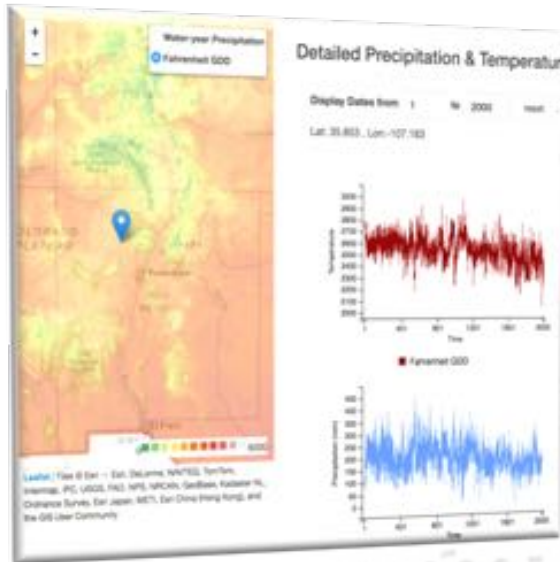On Provenance

# PRIMAD *(what have you "primed"?)*



130    16041 – Reproducibility of Data-Oriented Experiments in e-Science

| Label | Data | | Platform / Stack | Implementation | Method | Research Objective | Actor | Gain |
|---|---|---|---|---|---|---|---|---|
| | Parameters | Raw Data | | | | | | |
| Repeat | - | - | - | - | - | - | | Determinism |
| Param. Sweep | x | - | - | - | - | - | | Robustness / Sensitivity |
| Generalize | (x) | x | - | - | - | - | | Applicability across different settings |
| Port | - | - | x | - | - | - | | Portability across platforms, flexibility |
| Re-code | - | - | (x) | x | - | - | | Correctness of implementation, flexibility, adoption, efficiency |
| Validate | (x) | (x) | (x) | (x) | x | - | | Correctness of hypothesis, validation via different approach |
| Re-use | - | - | - | - | - | x | | Apply code in different settings, Re-purpose |
| Independent x (orthogonal) | | | | | | | x | Sufficiency of information, independent verification |

**Figure 1** PRIMAD Model: Categorizing the various types of reproducibility by varying the (P)latform, (R)esearch Objective, (I)mplementation, (M)ethod, (A)ctor and (D)ata, analyzing the gain they bring to computational experiments. x denotes the variable primed i.e. changed, (x) a variable that may need to be changed as a consequence, whereas – denotes no change.
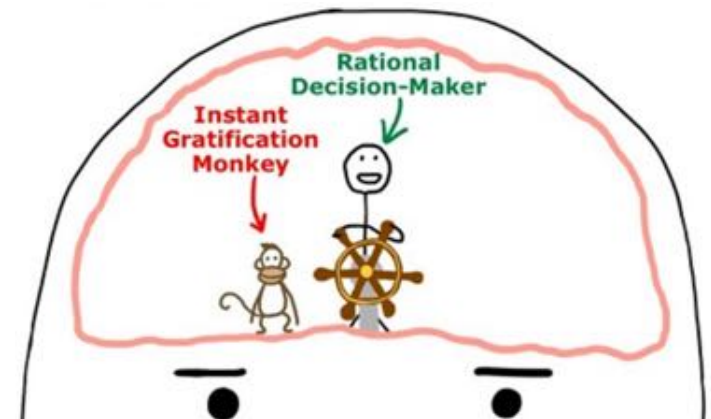
*Dagstuhl Seminar #16041 Report*

# Some related Projects ..





On Provenance

- **DataONE:** …
  - ➔ **Posters** by Linh Hoang & Hui Lyu, Xiaoliang Jiang
- **SKOPE**: system and tools to discover, access, analyze, visualize **paleoenvironmental data**
  - unprecedented ability to explore provenance (detailed, comprehensible record of computational derivation of results)
  - for researchers, tinkerers, and modelers
  - ➔ **WT /SKOPE poster** by Pratik Shrivastava
- **Whole Tale:**
  - **leverage & contribute to existing CI** to support the whole tale ("living paper"), from workflow run to scholarly publication
  - integrate tools & CI (DataONE, Globus, iRODS, NDS, …) to **simplify use and promote best practices.**
  - **driven by science WGs** (Archaeology/SKOPE, materials science, astro, bio ..)

# Preliminary Conclusions

- **Goal**: allow researchers (+*tinkers*, +*modelers*) to tell the whole tale of a science study, transparently and reproducibly.

- **Provenance** …
  - … is key to transparency, reproducibilty, comprehensibility
  - … comes in many (hybrid) forms (*workflow graphs*, *log files*, *trace events*, …)
  - … is metadata (=> "*a love note to the future*")
  - … should be *actionable today* (feed both, your IGM & RDM)

- **Provenance-for-Self** …
  - … asks: *how does provenance help me get my work done today?*
  - … is what provenance technologists and tool builders could/should do more of!

*Inside the mind of a master procrastinator (TED Talk by Tim Urban)*

# Truth or Consequences ... (ice breaker)

- Which of these are true/false?
- As a high-school student I ...
  - ... worked at a nuclear power plant
  - ... worked at a historic Roman bath
  - ... migrated code from FORTRAN to Pascal

- Later in life I ...
  - ... toured the backwaters from Cochin to Alleppey
  - ... toured the Isle of Skye on motorbike
  - ... toured Sri Lanka on motorbike
  - ... became a national master in chess
  - ... discovered game provenance, invented provenance games