

Theory & Practice of Data Cleaning

What is provenance and why is it important?

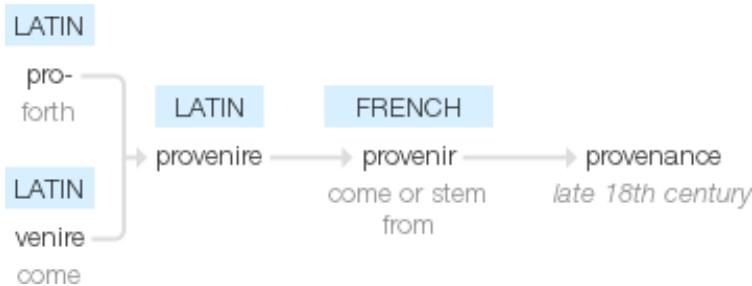
Example: The Fine Arts



- One of these 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*
- What is the origin (provenance!) of “provenance” ?



What is “provenance” (sensu W3C) ?

- *Provenance is a description of how things came to be, and how they came to be in the state they are in today (*)*

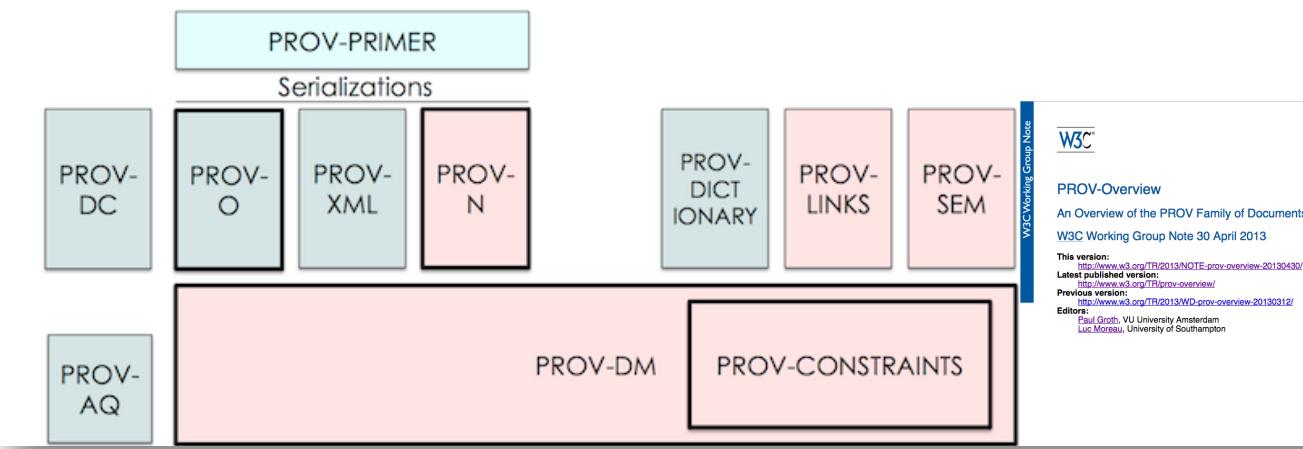


- *Provenance refers to the sources of information, including entities and processes, involved in producing or delivering an artifact (*)*
- *Provenance is a record that describes the people, institutions, entities, and activities, involved in producing, influencing, or delivering a piece of data or a thing in the world*

A screenshot of the PROV-DM: The PROV Data Model W3C Recommendation page. The page header includes the W3C logo and the title "PROV-DM: The PROV Data Model". Below the title, it says "W3C Recommendation 30 April 2013". The page contains several links to related documents and a list of contributors. The contributors listed include Khalid Belhajjame, Reza B'Far, James Cheney, Sam Coppens, Stephen Cresswell, Yolanda Gil, Paul Groth, Graham Klyne, Timothy Lebo, Jim McCusker, Simon Miles, James Myers, Satya Sahoo, and Curt Tilmes. The page also mentions Luc Moreau and Paolo Missier as editors.

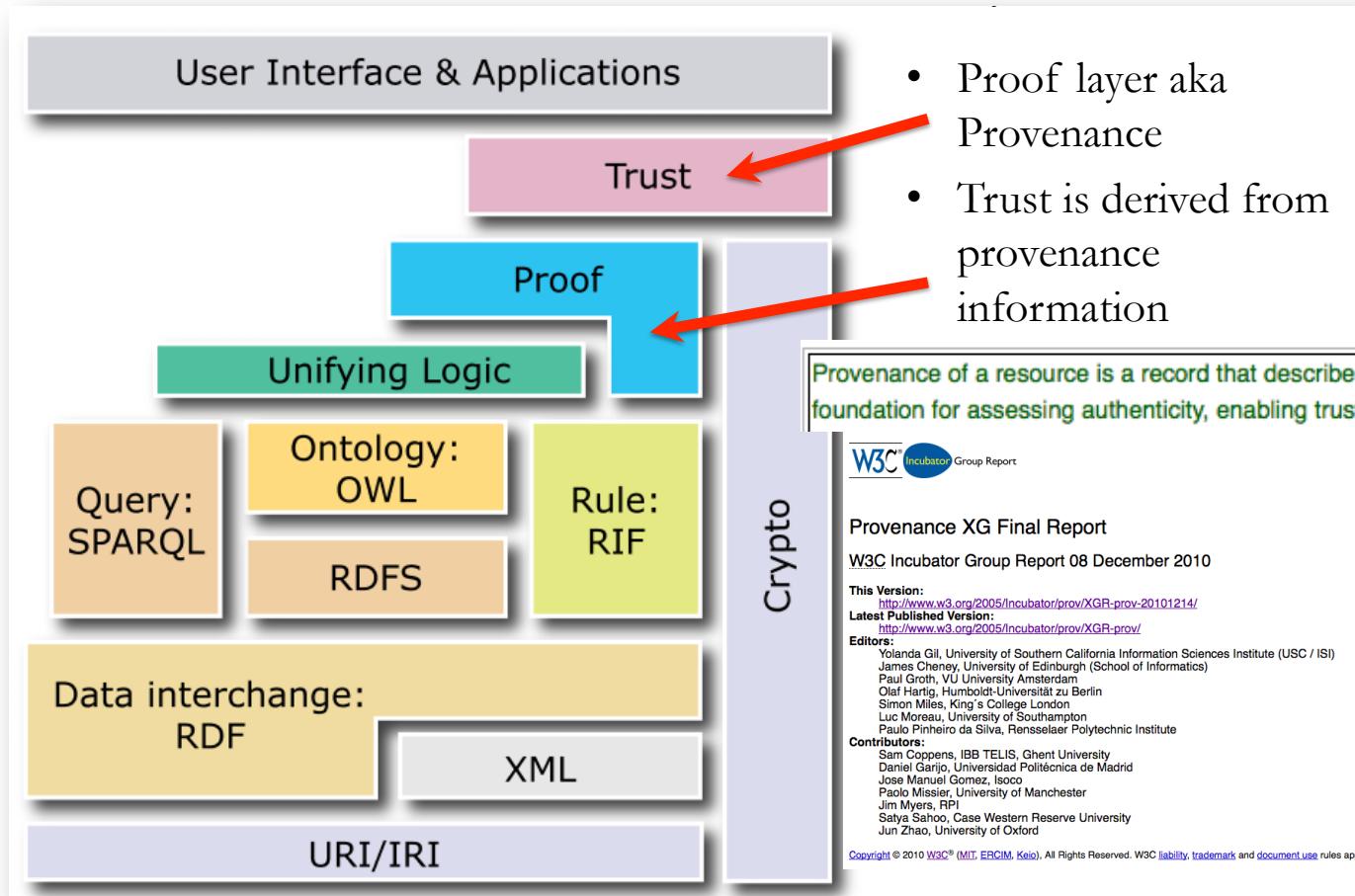
W3C PROV Family of Specifications: Modeling

- W3C Recommendations
 - PROV Data Model (PROV-DM)
 - PROV Ontology (PROV-O)
 - PROV-Constraints
 - PROV Notations (PROV-N)
- PROV Working Group Notes (selected)
 - PROV-Access and Querying (AQ)
 - PROV Dictionary
 - PROV XML
 - PROV and Dublin Core Mappings (PROV-DC)
 - PROV Semantics (using first-order logic) (PROV-SEM)



Provenance Analysis and RDF Query Processing,
Satya S. Sahoo, Praveen Rao, ISWC, October, 2015.

Provenance & Semantic Web Layer Cake



Provenance Analysis and RDF Query Processing,
Satya S. Sahoo, Praveen Rao, ISWC, October, 2015.

Back to the basics: Open Provenance Model (OPM) => W3C Prov

The Open Provenance Model aims to capture the causal dependencies between the artifacts, processes, and agents. Therefore, a provenance graph is defined as a directed graph, whose nodes are artifacts, processes and agents, and whose edges belong to one of the following categories depicted in Figure 1. An edge represents a causal dependency, between its source, denoting the effect, and its destination, denoting the cause.

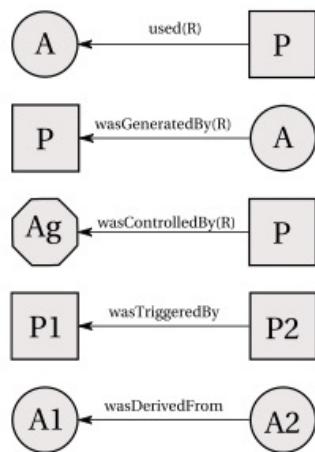


Figure 1: Edges in the Open Provenance Model: sources are effects, and destinations causes

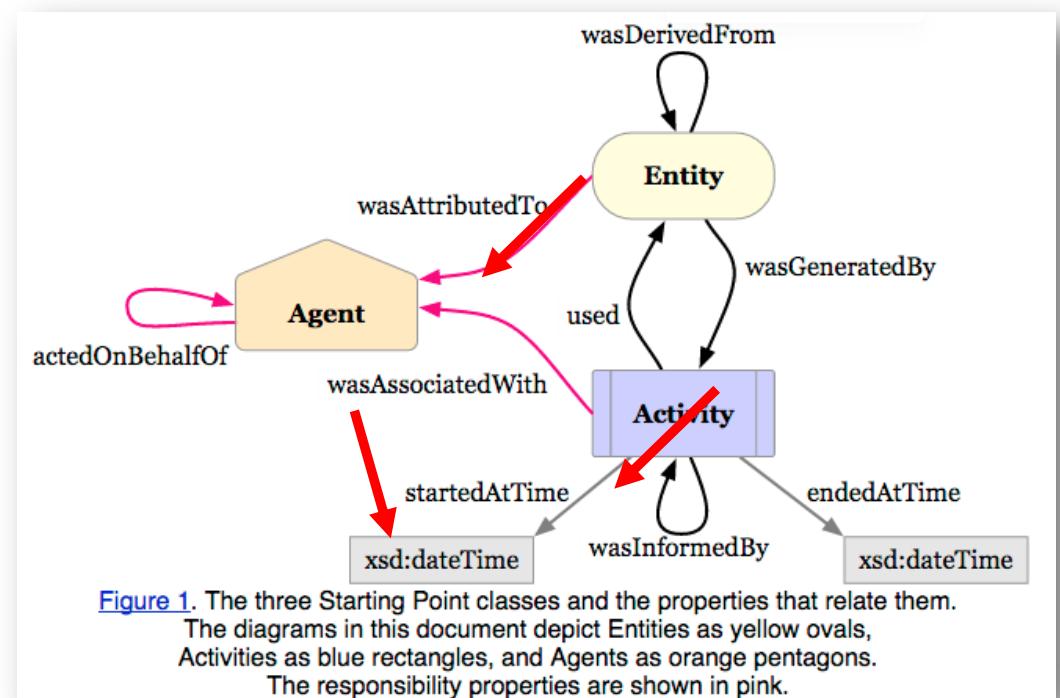


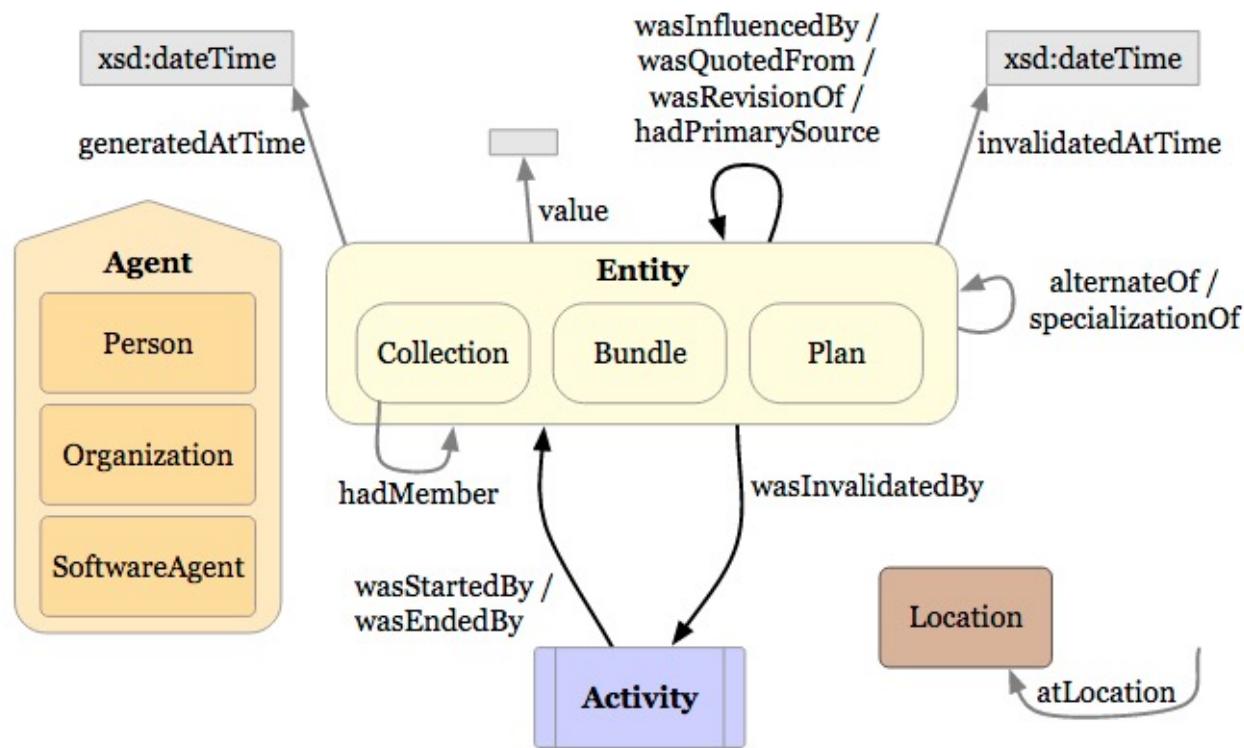
Figure 1. The three Starting Point classes and the properties that relate them.

The diagrams in this document depict Entities as yellow ovals, Activities as blue rectangles, and Agents as orange pentagons.

The responsibility properties are shown in pink.

Readings: Moreau, Luc, Ben Clifford, Juliana Freire, Joe Futrelle, Yolanda Gil, Paul Groth, Natalia Kwasnikowska et al. "The open provenance model core specification (v1. 1)." Future generation computer systems 27, no. 6 (2011): 743-756.

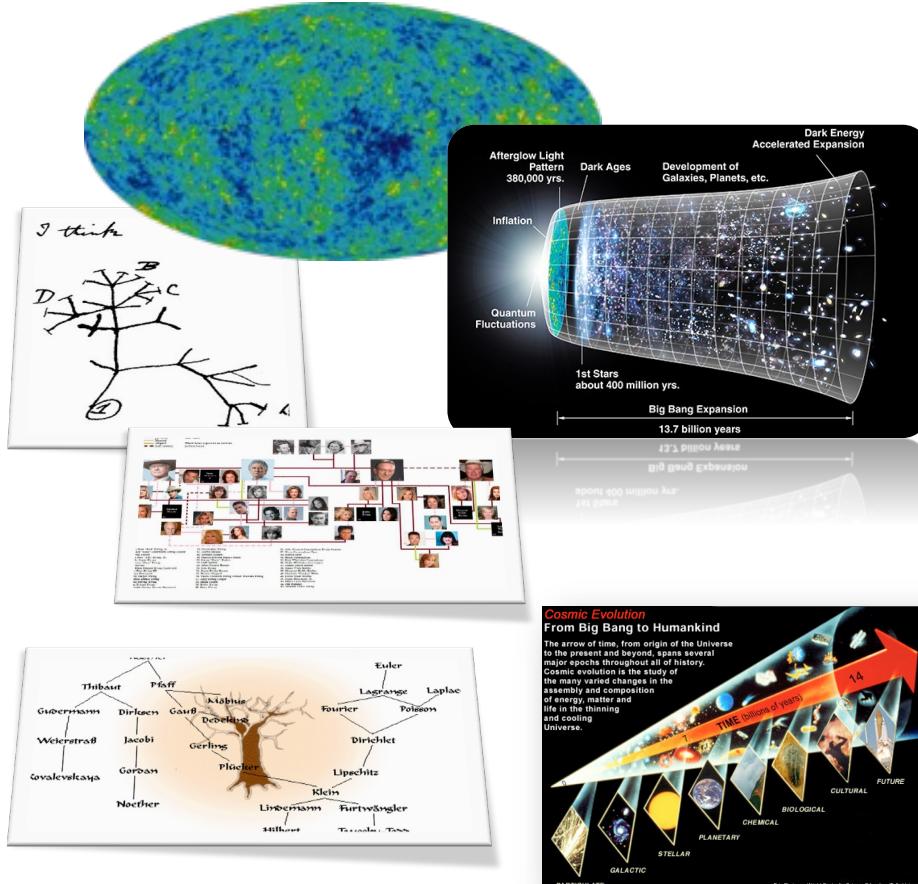
W3C Prov: some finer points



[Figure 3](#). The expanded terms build upon those in the [Starting Points section](#).

Provenance-in-Science Palooza

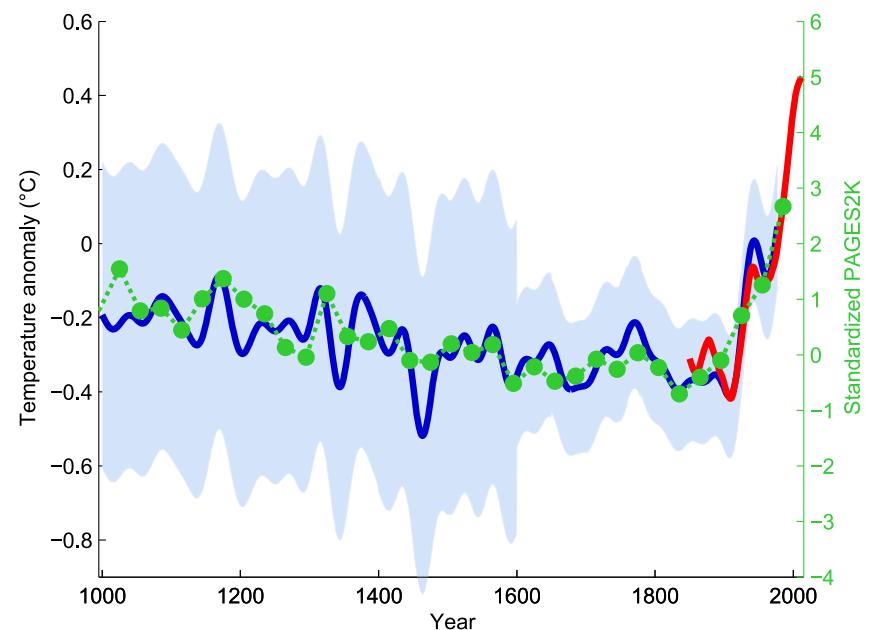
- What are those?
- Cosmology
- Geology, Stratigraphy
- Phylogeny
 - the *Tree of Life*
- Genealogy
 - your family: literally
- Academic Pedigree
 - “Doktorvater” (oder Doktor-Mutter)
- Etymology
- Chain of custody
 - of art(ifacts)
- Yes: all about origins and history



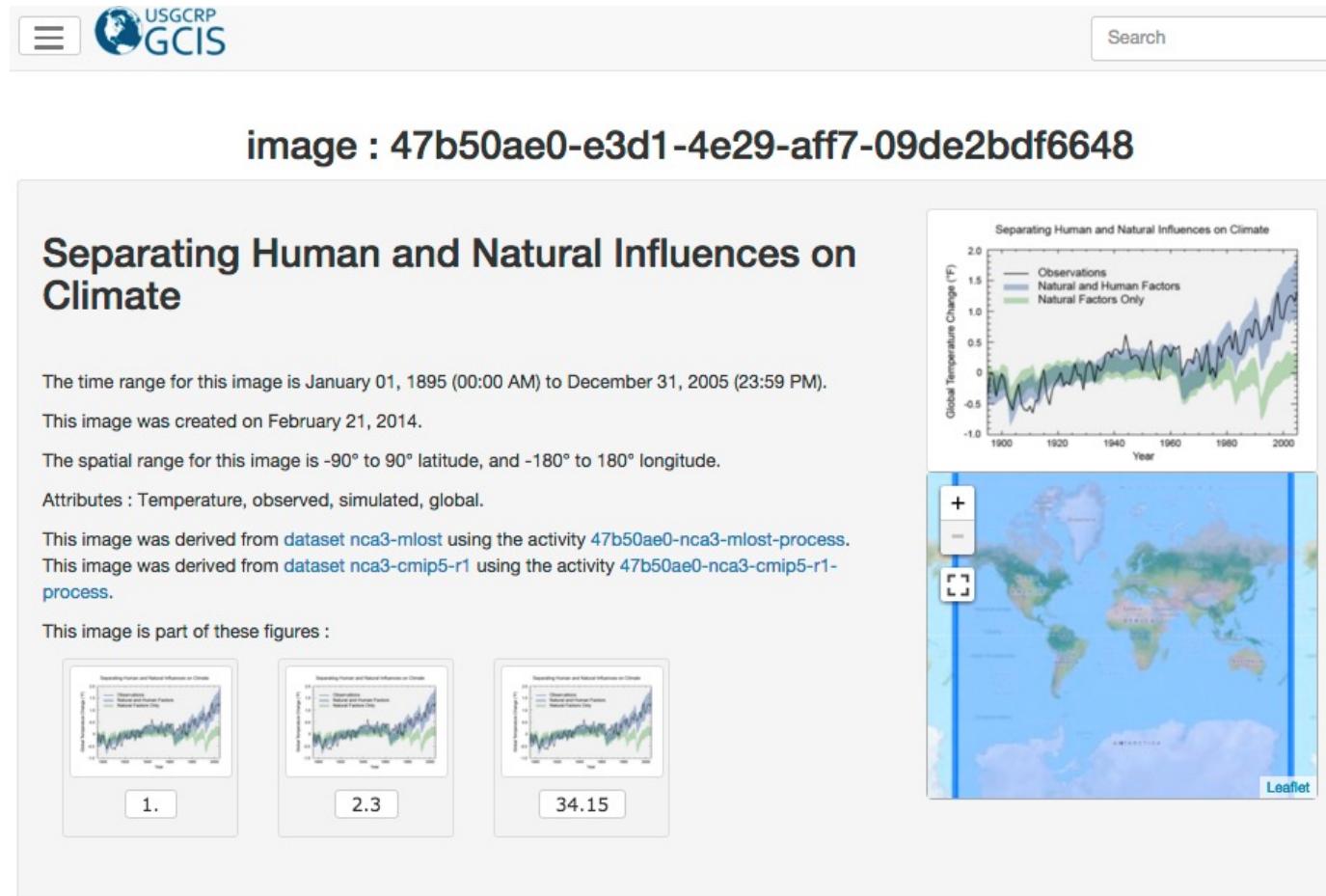
...

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?



You are viewing /image/47b50ae0-e3d1-4e29-aff7-09de2bdf6648 in [HTML](#)

Alternatives : [JSON](#) [YAML](#) [Turtle](#) [N-Triples](#) [JSON Triples](#) [RDF+XML](#) [RDF+JSON](#) [Graphviz](#) [SVG](#)

Tracing the sources (data, code)

Texas Summer 2011: Record Heat and Drought

Cooperative Institute for Climate and Satellites - NC
Laura Stevens

The time range for this image is January 01, 1895 (00:00 AM) to December 31, 2012 (00:00 AM).

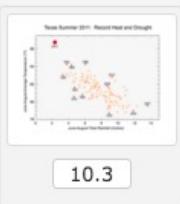
This image was created on July 03, 2013.

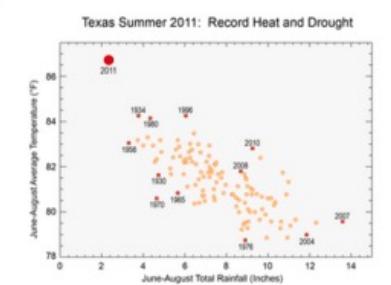
The spatial range for this image is 25.83° to 36.50° latitude, and -106.65° to -93.52° longitude.

Attributes : Temperature, precipitation, observed, Texas.

This image was derived from dataset nca3-cddv2-r using the activity 02c53cf7-nca3-cddv2-r1-process.

This image is part of this figure :

 10.3





data and "code" / method linked

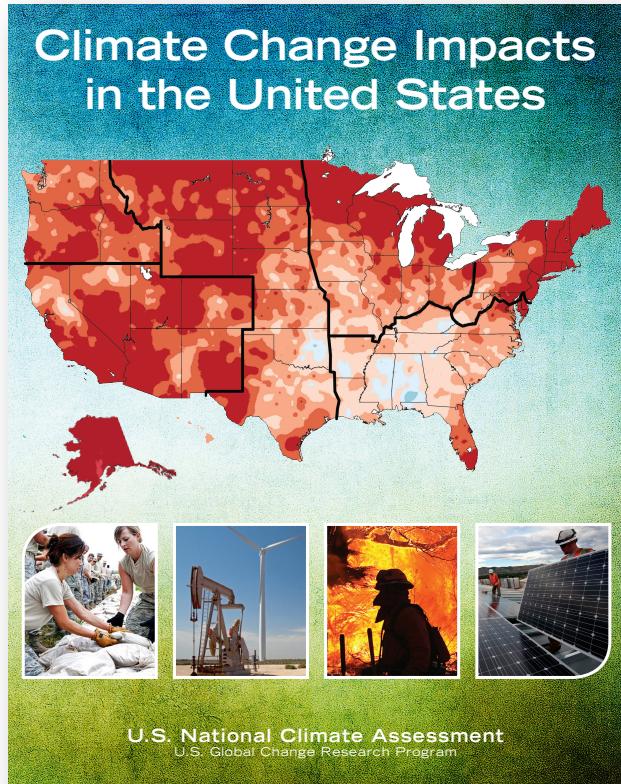
alt formats

You are viewing /image/02c53cf7-75f8-4243-a925-f59a0025f04e in [HTML](#)

Alternatives : [JSON](#) [YAML](#) [Turtle](#) [N-Triples](#) [JSON Triples](#) [RDF+XML](#) [RDF+JSON](#) [Graphviz](#) [SVG](#)

 GlobalChange.gov

Provenance today: Important but *hard*



"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 Data Observation Network for Earth

DataONE

About News Participate Resources Education Data

DATAONE SEARCH: Search Summary Jump to: DOI or ID Go Sign in or Sign up

Search ?
Search phrase

Datasets 1 to 25 of 238,334
1 2 3 ... 9,534 Next Sort by Most recent

Filter by:
Data attribute
Data files
Member Node
 Cornell Lab of Ornithology - e...
 Dryad Digital Repository
 EDAC Gstore Repository
 ESA Data Registry
Show 28 more

Creator
Year
Identifier
Taxon
Location

Cooney, Feargus, Vitikainen, Emma, Marshall, Harry, Smith, Robert, Cant, Michael, Goodey, Nicole, and van Royen, Wilmie. 2016. **Data from: Lack of aggression and apparent altruism towards intruders in a primitive termite.** Dryad Digital Repository. <http://dx.doi.org/10.5061/dryad.tb0c7?ver=2016-10-06T09:41:23.158-04:00>. [i](#)

Cooney, Feargus, Vitikainen, Emma, Marshall, Harry, Smith, Robert, Cant, Michael, Goodey, Nicole, and van Royen, Wilmie. 2016. **Data from: Lack of aggression and apparent altruism towards intruders in a primitive termite.** Dryad Digital Repository. <http://dx.doi.org/10.5061/dryad.tb0c7?ver=2016-10-06T09:43:10.046-04:00>. [i](#)

Greenway, Ryan, Drexler, Shannon, Arias-Rodriguez, Lenin, and Tobler, Michael. 2016. **Data from: Adaptive, but not condition-dependent, body shape differences contribute to assortative mating preferences during ecological speciation.** Dryad Digital Repository. <http://dx.doi.org/10.5061/dryad.2dn85?ver=2016-10-06T09:35:26.892-04:00>. [i](#)

Laidre, Kristin, Moon, Twila, Hauser, Donna, McGovern, Richard, Heide-Jørgensen, Mads Peter, Dietz, Rune, and Hudson, Benjamin. 2016. **Data from:**

Hide Map »

1727 171 7 1 NU 3 9 12

529 231 68 5 Canada 3 5 2 5 7 6

120 156 93 7 259 89 6 9 22 11

44 85 1621 507 267 ND 4902 MN 1418 64 NB PEI 63 11

+ 34 9080 OR 4369 ID 3424 SD 119 2748 NS 7285 251 11

- 39 75 418 38139 4337 NE 1767 2501 KY WV DE VA NC 4665 426 11

34 119 186 3808 6179 TX 511 LA 3540 747 172 5

83 160 19 52 56 Mexico 239 6803 841 74 6

225 19 11 29 65 86 110 266 Puerto Rico 1201 2

9 19 11 17 45 60 93 55 80 4

Satellite Terrain 17 20 29 23 113 51 Venezuela 107 3

Google 58 35 35 Map data ©2016 Google, INEGI 1000 km Terms of Use

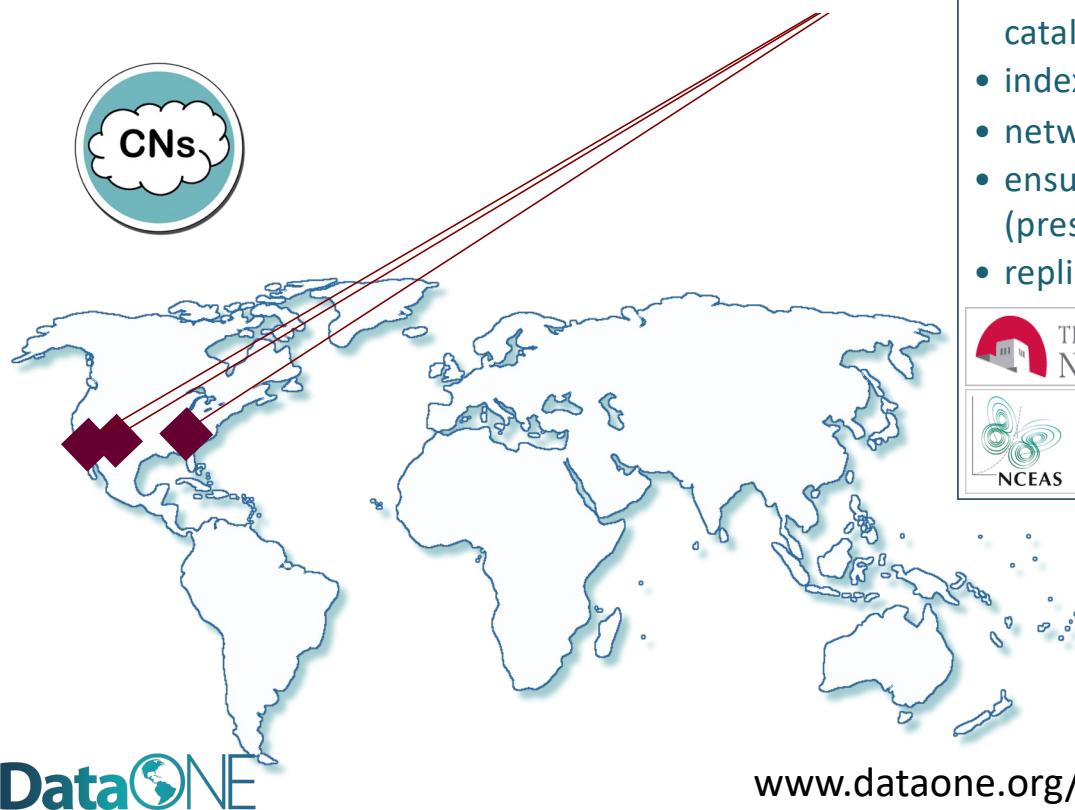
DataONE is a collaboration among many partner organizations, and is funded by the US National Science Foundation (NSF) under a Cooperative Agreement. Acknowledgment: This material is based upon work supported by the National Science Foundation under Grant Numbers 0830944 and 1430508. Disclaimer: Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.

DataONESearch

search.dataone.org

DataONE Cyberinfrastructure: Coordinating Nodes

Components for a flexible, scalable, sustainable network



DataONE

Coordinating Nodes

- retain complete metadata catalog
- indexing for search
- network-wide services
- ensure content availability (preservation)
- replication services



THE UNIVERSITY of
NEW MEXICO



OAK RIDGE
National Laboratory



NCEAS



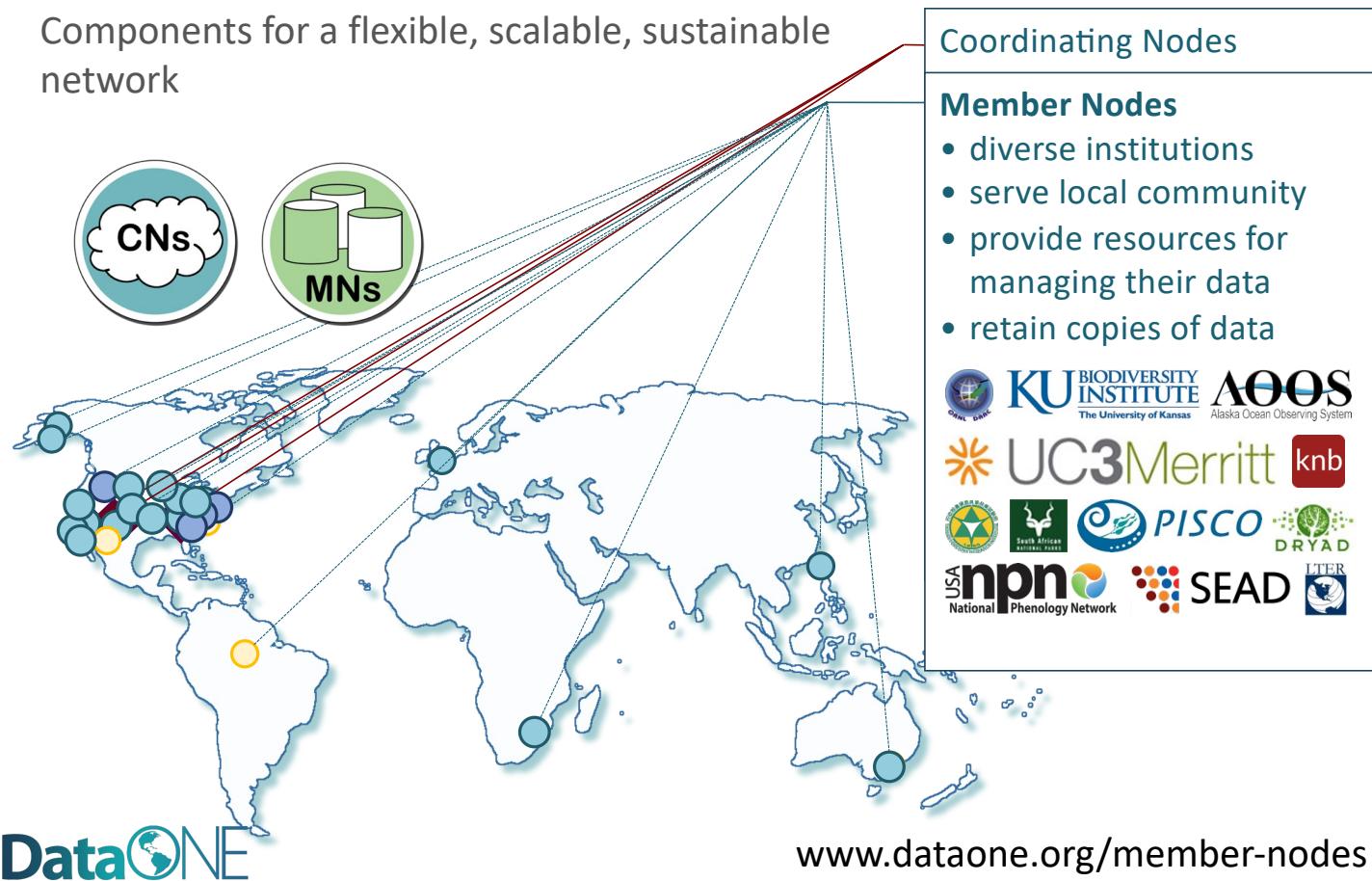
UCSB



www.dataone.org/coordinating-nodes

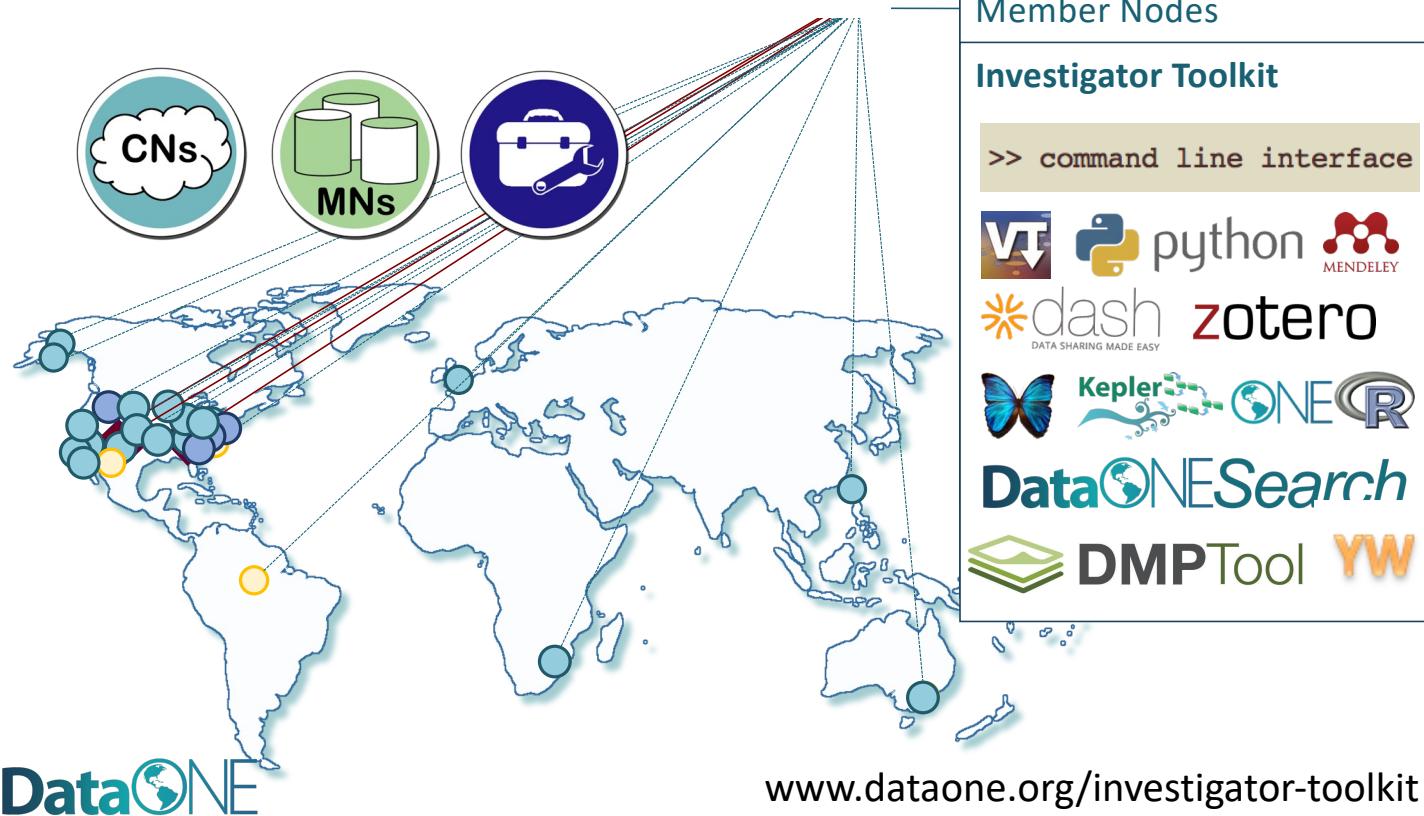
DataONE Cyberinfrastructure: Member Nodes

Components for a flexible, scalable, sustainable network



DataONE Cyberinfrastructure: Investigator Toolkit

Components for a flexible, scalable, sustainable
network



Provenance in Action: Benefits & Impact

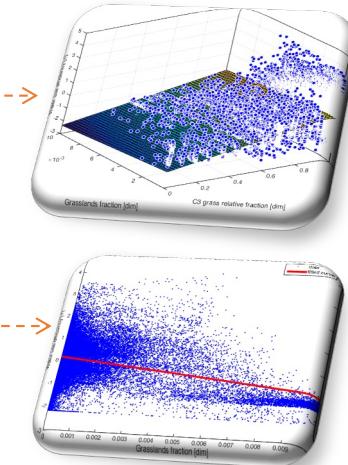
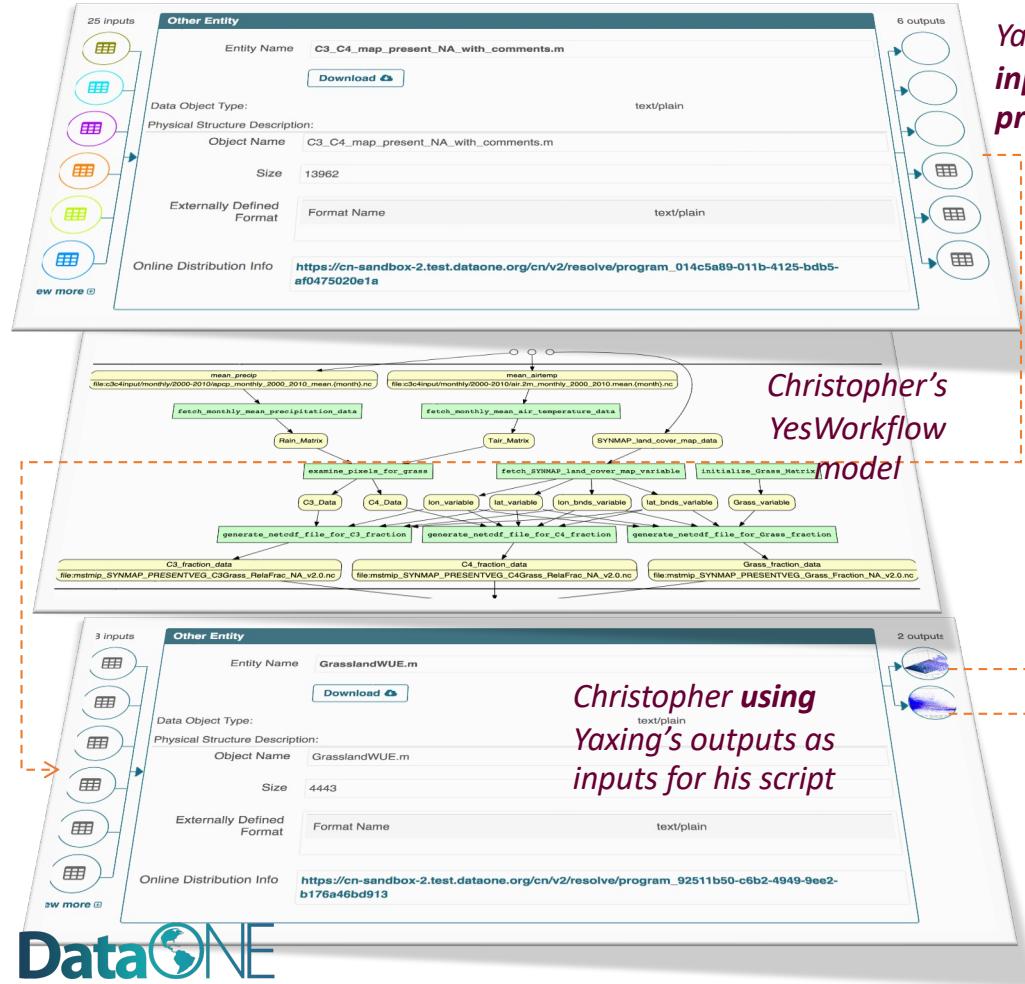
The screenshot shows a search interface for 'grass'. The search bar contains 'grass' and a magnifying glass icon. To the right, there's a 'Sort by' dropdown set to 'Most recent' with a downward arrow. Below the search bar is a navigation bar with pages 1 through 160, and a 'Next' button. The main area is titled 'My Search' and shows two search results:

- Christopher Schwalm. 2016. Grassland Water Use Efficiency (WUE) Analysis: Run of GrasslandWUE.m on 20160317T154050.** MN Demo 2. metadata_07277c1f-b2c2-467c-8aa2-792863524a21.xml.
This item has a circled 'p' icon with a dashed arrow pointing to the second item.
- Yaxing Wei. 2016. MsTMIP: C3 C4 soil map processing: Run of C3_C4_map_present_NA_with_comments.m on 20160311T181011.** MN Demo 2. metadata_e859d2dd-c5e6-4ec6-892f-1b00bb6f8f65.xml.
This item also has a circled 'p' icon with a dashed arrow pointing to the first item.

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



DataONE: Support for Provenance



Theory & Practice of Data Cleaning

YesWorkflow: Provenance for script-based workflows

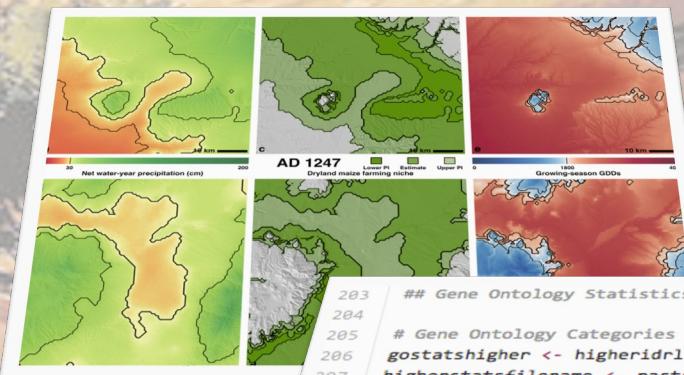
From Workflows & Provenance to Provenance for Script-based Workflows ...

- 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

```
## Gene Ontology Statistics are Calculated Here.  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215  
216 # Gene Ontology Categories that were shown to be relatively Higher (more expressed) in the Experimental Condition.  
gostatshigher <- higheridlinkedtogenes[1]  
higherstatsfilename <- paste(outputDirectory, "/", runName, "_", conditions[1], "_GOSTatsHigher_", mytestcond[1], ".v  
write.table(gostatshigher,file=higherstatsfilename, row.names=FALSE, col.names=FALSE, quote=FALSE, sep="\t")  
geneListHigherCHR <- gostatshigher$SYMBOL  
geneListHigherLinkedtoEntrezIds <- select(hgu133plus2.db, keys= geneListHigherCHR, "ENTREZID", "SYMBOL")  
G0statsGenesH <- geneListHigherLinkedtoEntrezIds[,2]  
  
x <- org.Hs.egACCNUM  
mapped_genes <- mappedkeys(x)  
xx <- as.list(x[mapped_genes])  
geneUniverse <- (unique(names(xx)))
```

... 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?

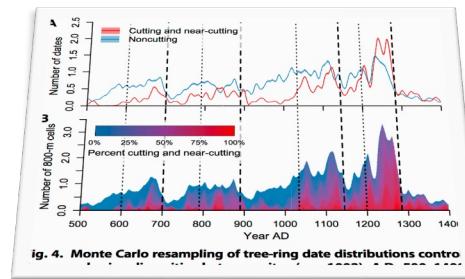
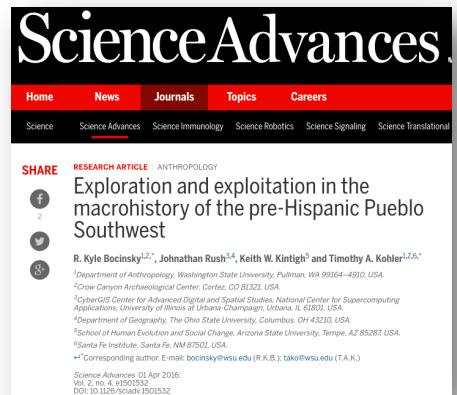
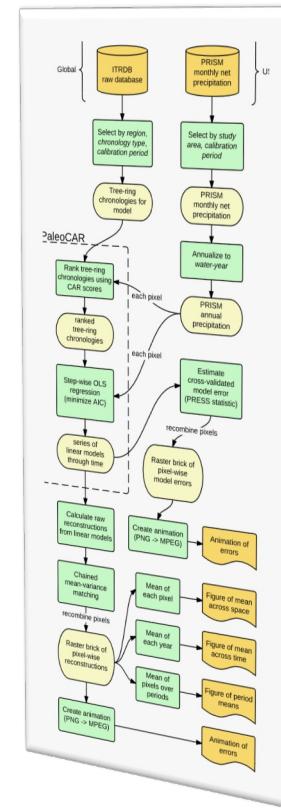
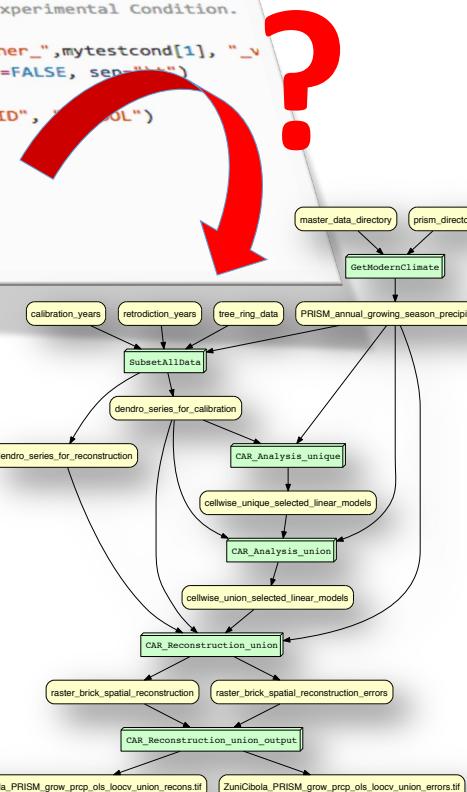


fig. 4. Monte Carlo resampling of tree-ring date distributions contro



YesWorkflow: Yes, scripts are workflows, too!

```
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 gosatshigher <- higheridrlinkedtogenes[1]  
207 higherstatsfilename <- paste(outputDirectory, "/", runName, "_", conditions[1], "_GOSTatsHigher_", mytestcond[1], ".v  
208 write.table(gosatshigher, file=higherstatsfilename, row.names=FALSE, col.names=FALSE, quote=FALSE, sep="\t")  
209 geneListHigherCHR <- gosatshigher$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)))
```



- Script vs Workflows/ASAP:

- Automation: ****
- Scaling: **
- Abstraction: *
- Provenance: **

YW annotations: Model your Workflow!

```
1 # @BEGIN collect_data_set
2 # @PARAM cassette_id @PARAM accepted_sample @PARAM num_images @PARAM energies
3 # @OUT sample_id @OUT energy @OUT frame_number
4 # @OUT raw_image_path @AS raw_image
5 # ... @URI file:run/raw/{cassette_id}/{sample_id}/e{energy}/image_{frame_number}.raw
6 run.log.write("Collecting data set for sample {0}".format(accepted_sample))
7 sample_id = accepted_sample
8 for energy, frame_number, intensity, raw_image_path in collect.next_image(
9         cassette_id, sample_id, num_images, energies,
10        "run/raw/{cassette_id}/{sample.id}/e{energy}/image_{frame_number:03d}.raw"):
11    run.log.write("Collecting image {0}".format(raw_image_path))
12 # @END collect_data_set
13
14 # @BEGIN transform_images
15 # @PARAM sample_id @PARAM energy @PARAM frame_number
16 # @IN raw_image_path @AS raw_image
17 # @IN calibration_image @URI file:calibration.img
18 # @OUT corrected_image @URI file:run/data/{sample_id}/{sample_id}-{energy}eV-{frame_number}.img
19 # @OUT corrected_image_path @OUT total_intensity @OUT pixel_count
20     corrected.image_path = "run/data/{0}/{0}_{1}eV_{2:03d}.img".format(sample_id, energy, frame_number)
21     (total_intensity, pixel_count) = transform.image(raw_image_path, corrected.image_path, "calibration.img")
22     run.log.write("Wrote transformed image {0}".format(corrected.image_path))
23 # @END transform_images
```

mark the code block

... and data inputs/outputs

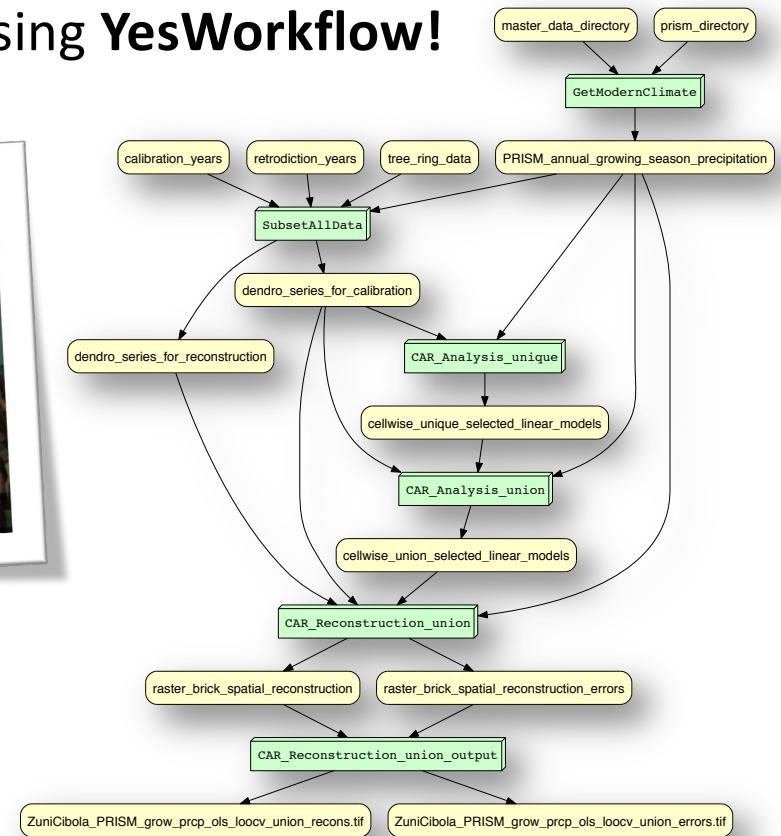
Figure 1. YW-annotated fragment of a Python script for data collection from protein crystal samples. YW-annotations @BEGIN and @END delimit code blocks; @IN and @OUT tags model relevant input and output data elements of a block; @PARAM identifies a block's parameters. @URI templates for raw images (line 5) and corrected images (line 18) link conceptual-level data elements such as raw_image with runtime resources (data files and their file paths). Executable script code is greyed out to emphasize YW-annotations. A program variable (raw_image_path) is highlighted in the code (lines 8, 11, 21): aliases (lines 4, 16) are used to link such program-level objects to the scientist's concepts (here: raw_image). Full example available from [MBL15].

Paleoclimate Reconstruction (EnviRecon.org) ...

... explained using **YesWorkflow!**

Kyle B.,
(computational)
archaeologist:

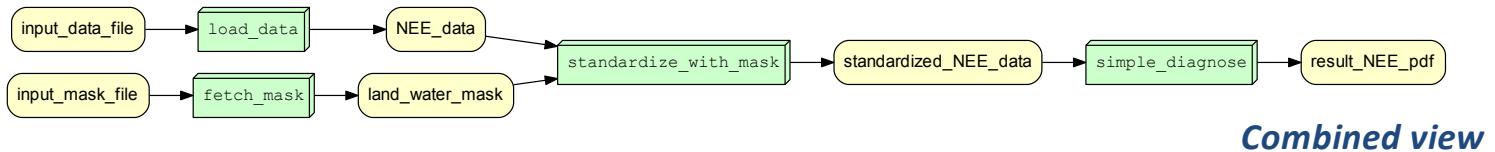
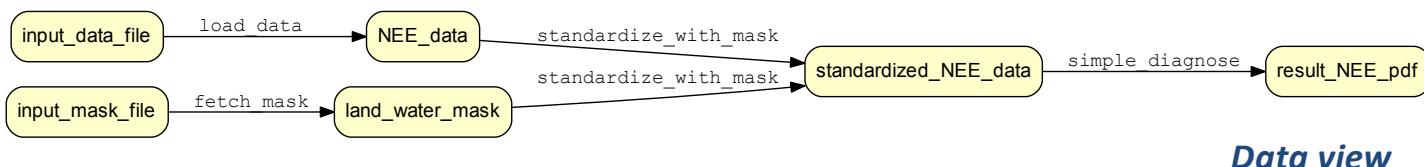
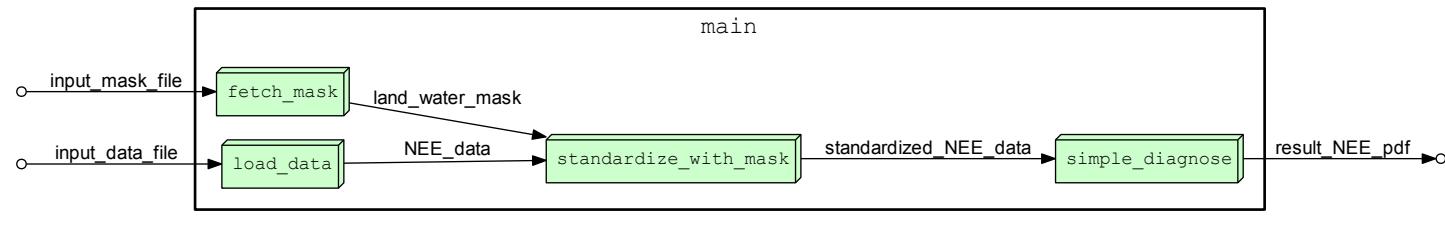
"It took me about 20 minutes to comment. Less than an hour to learn and YW-annotate, all-told."



SKOPE + **Kurator**
+ **DataONE**
Data Observation Network for Earth

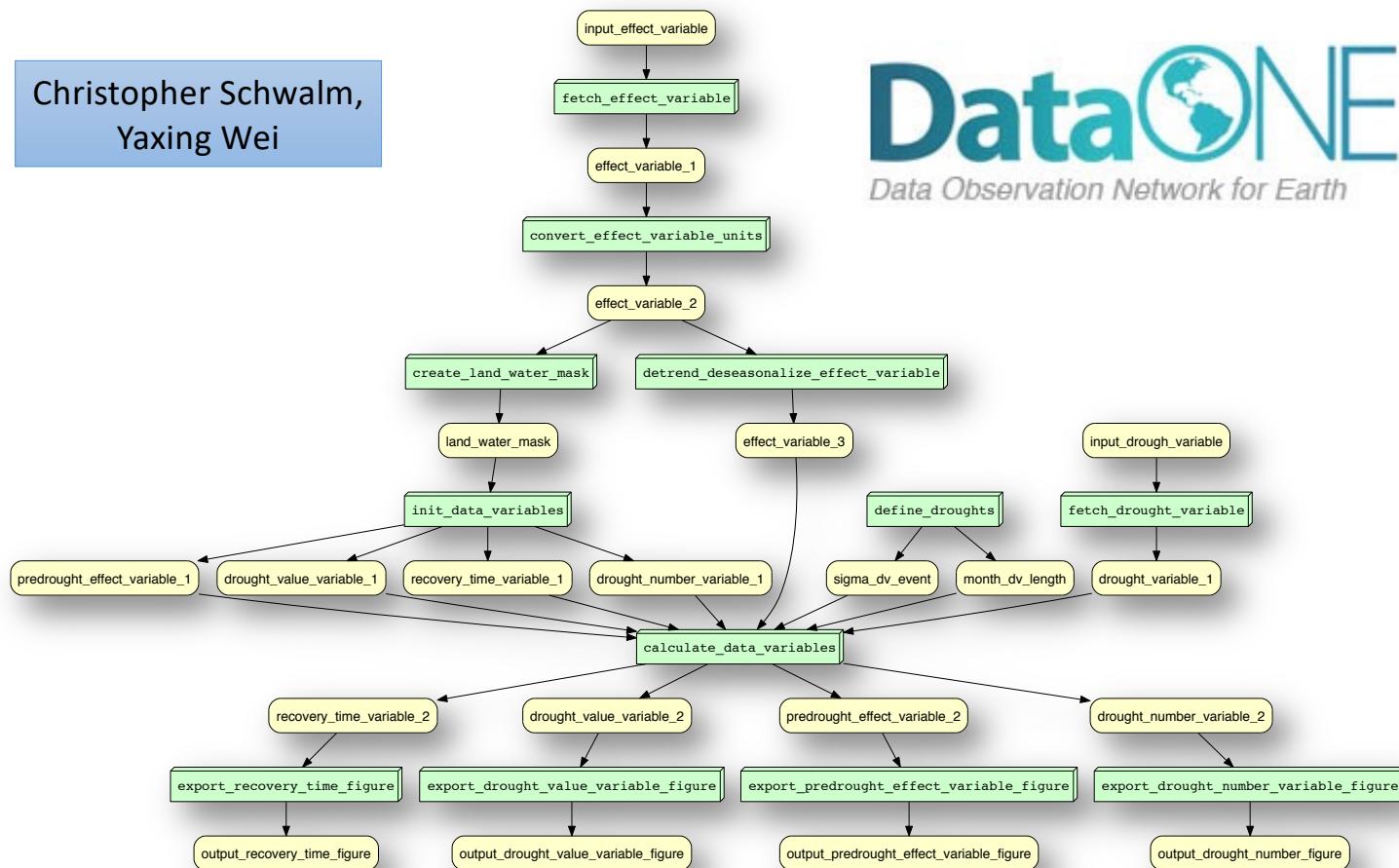
=> **YesWorkflow.org**

Get 3 views for the price of 1!

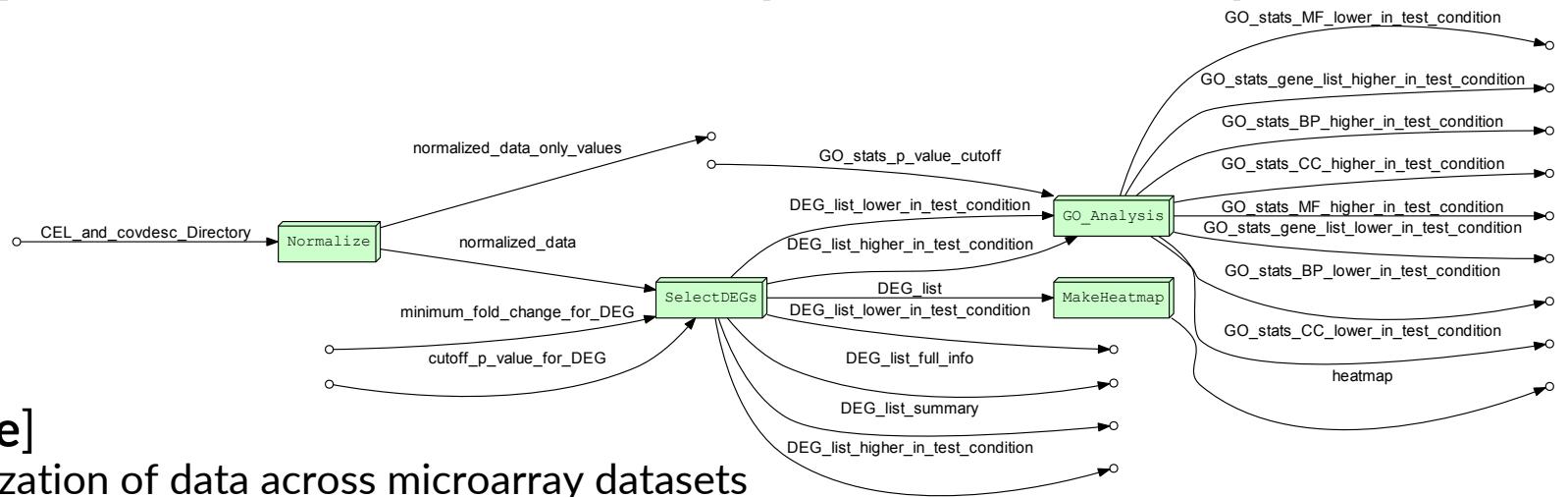


Multi-Scale Synthesis and Terrestrial Model Intercomparison Project (MsTMIP)

Christopher Schwalm,
Yaxing Wei



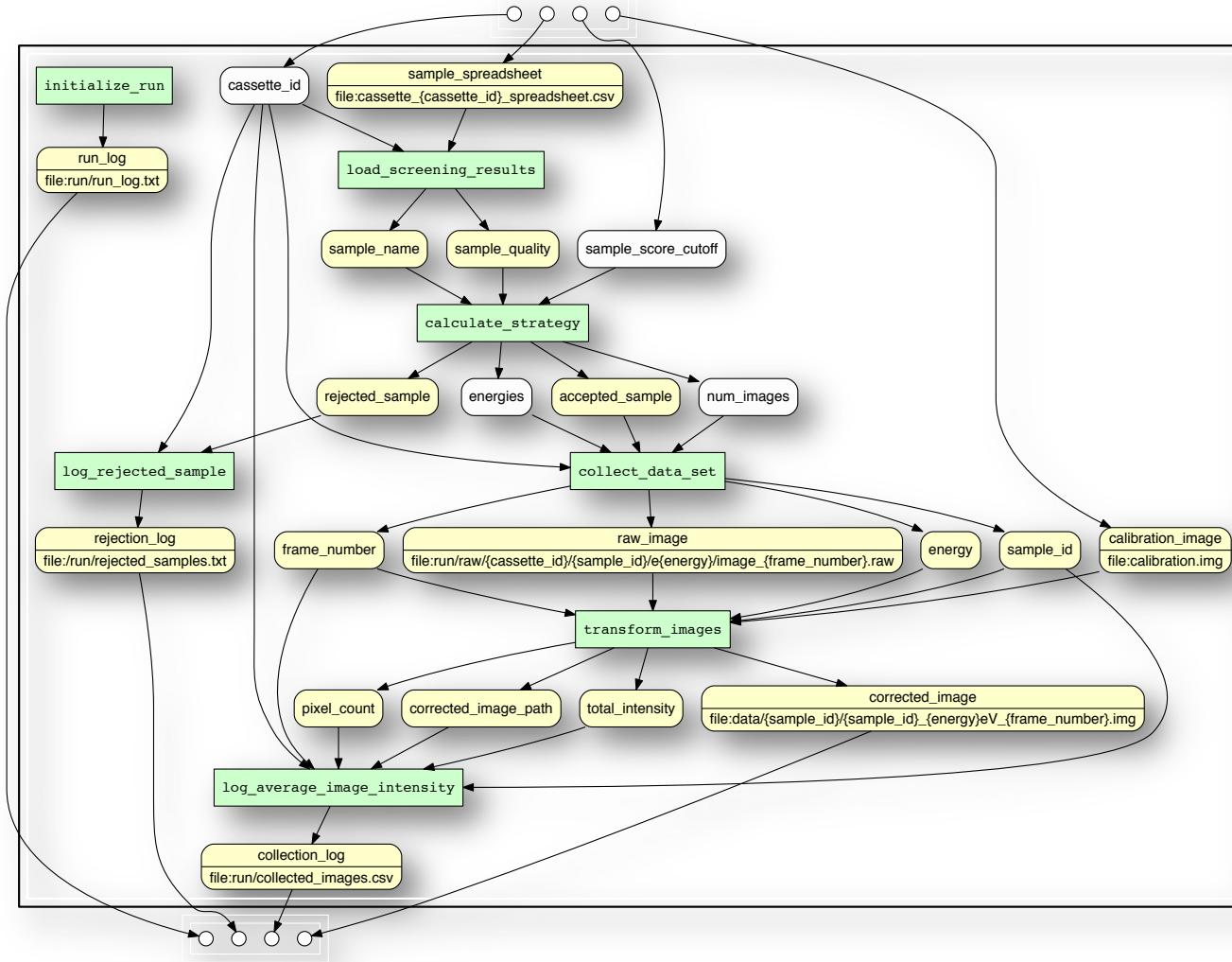
Gene Expression Microarray Data Analysis



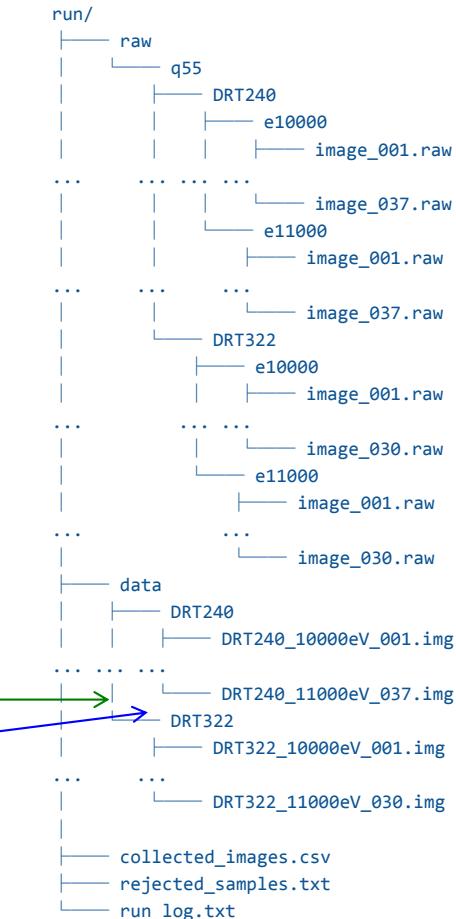
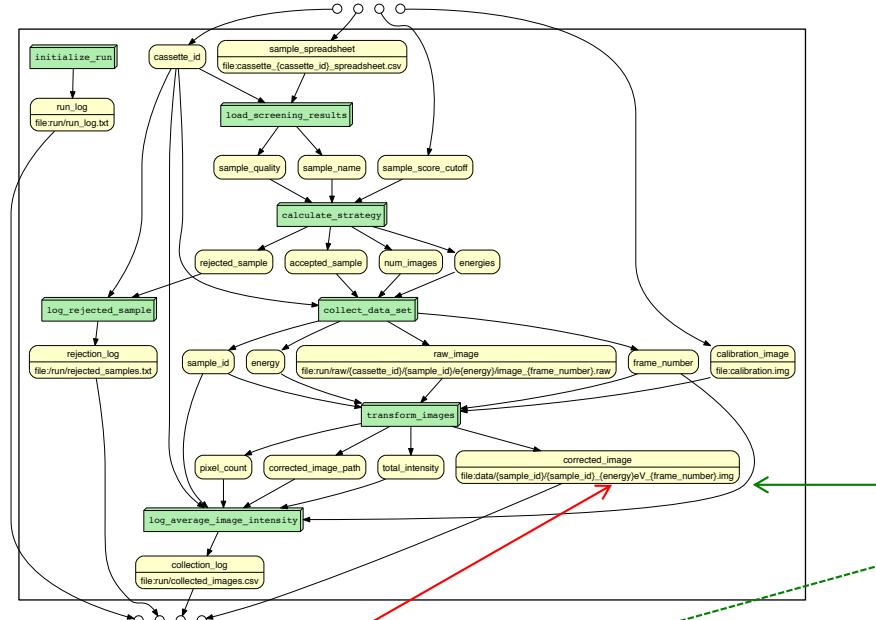
- **[Normalize]**
 - Normalization of data across microarray datasets
- **[SelectDEGs]**
 - Selection of differentially expressed genes between conditions
- **[GO Analysis]**
 - determination of gene ontology statistics for the resulting datasets
- **[MakeHeatmap]**
 - creation of a heatmap of the differentially expressed genes

Tyler Kolisnik, Mark Bieda

Data collection workflow (X-ray diffraction)



YW-RECON: Prospective & Retrospective Provenance ... (almost) for free!

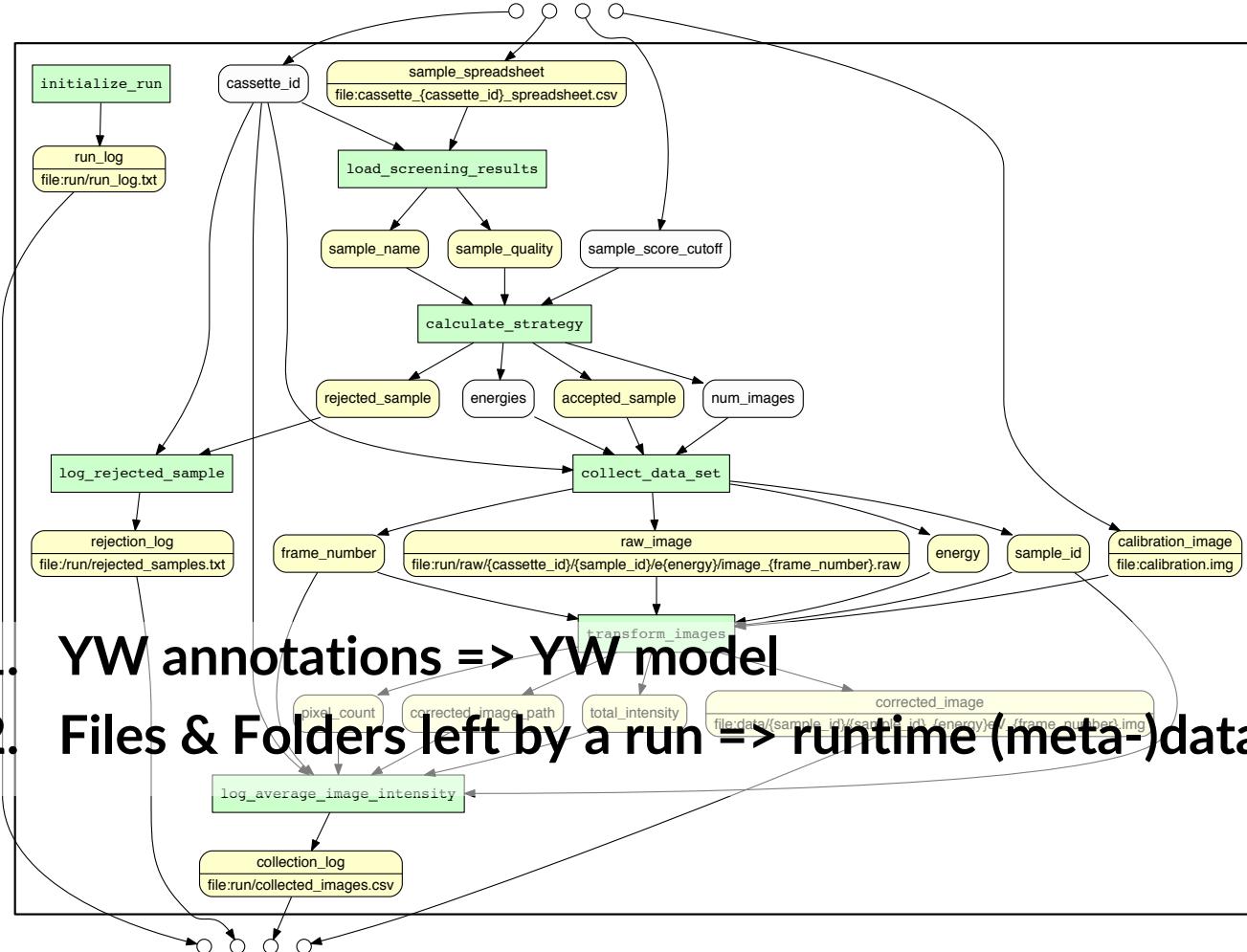


- **URI-templates** link conceptual entities to **runtime provenance** “left behind” by the script author ...
- ... facilitating provenance **reconstruction**

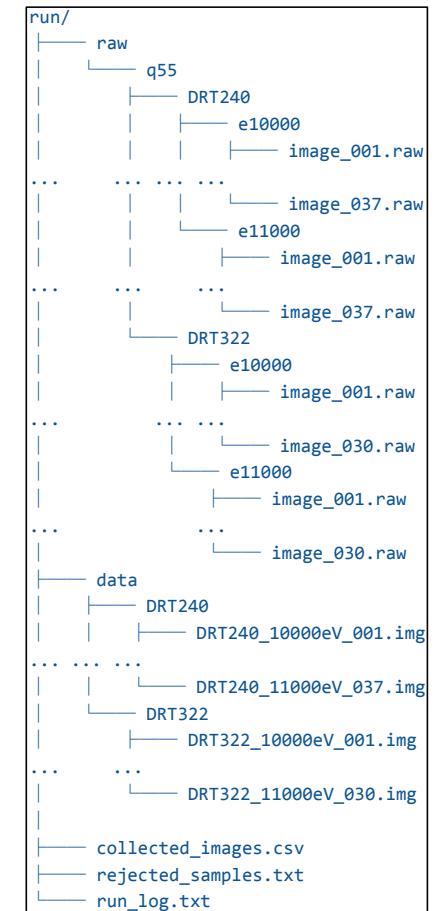
YW (*prospective*) and YW-Recon (*retrospective*) Provenance

- **1. YW: Annotate Script => YW Model**
 - Annotate @BEGIN..@END, @IN, @OUT
 - Visualize, share, be happy ☺
- **2. Run script**
 - Files are read and written
 - Folder- & Filenames have metadata
- **3. YW-Recon**
 - Use **@URI** tags that link YW Model \Leftrightarrow Persisted Data
 - Run URI-template queries
 - cf. “ls -R” & RegEx matching
- **4. YW-Query**
 - Answer the user’s provenance queries

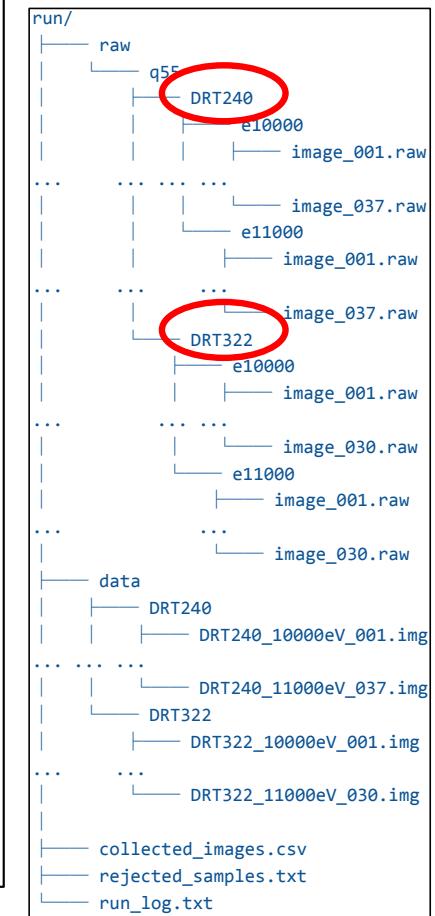
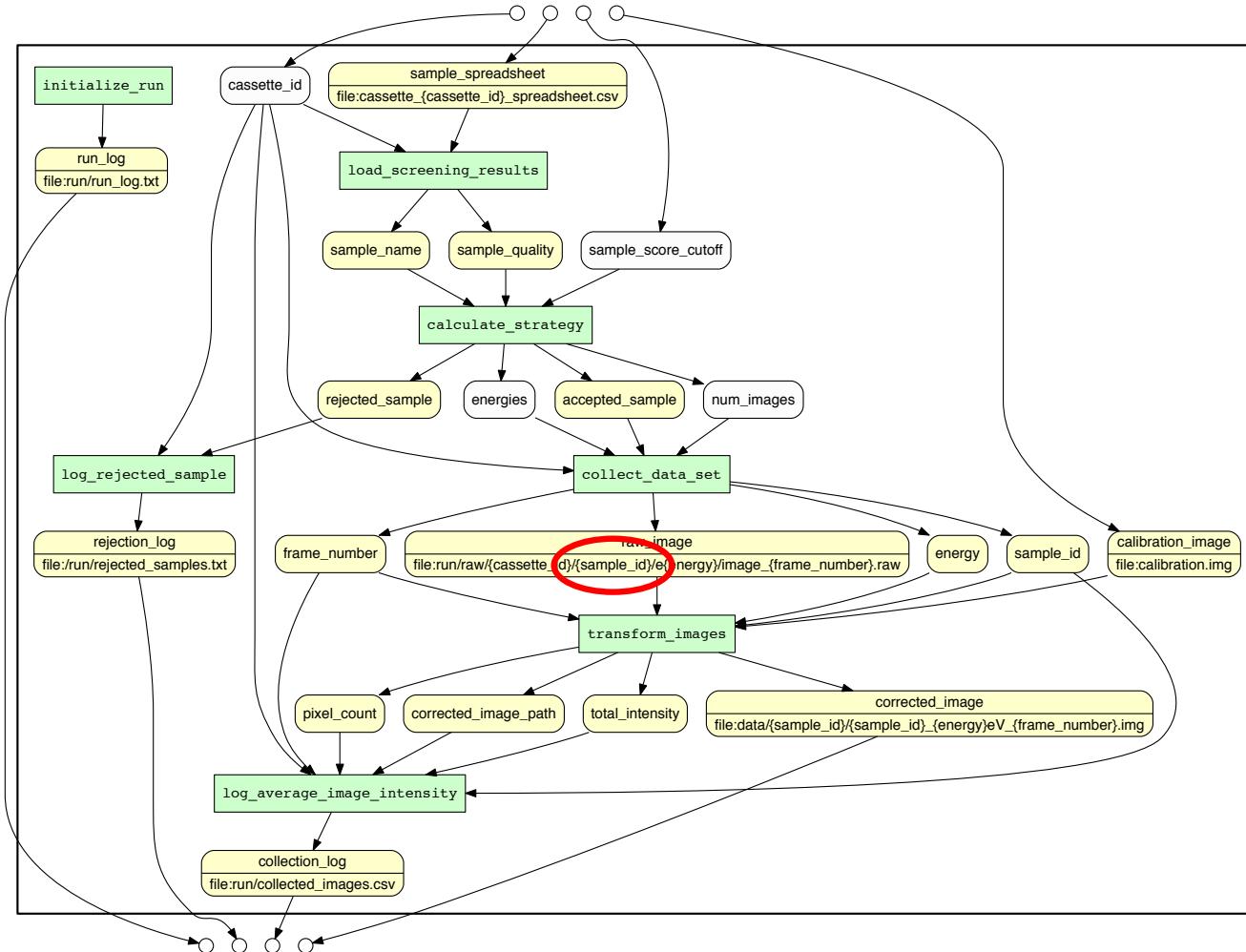
Data collection workflow: runtime data



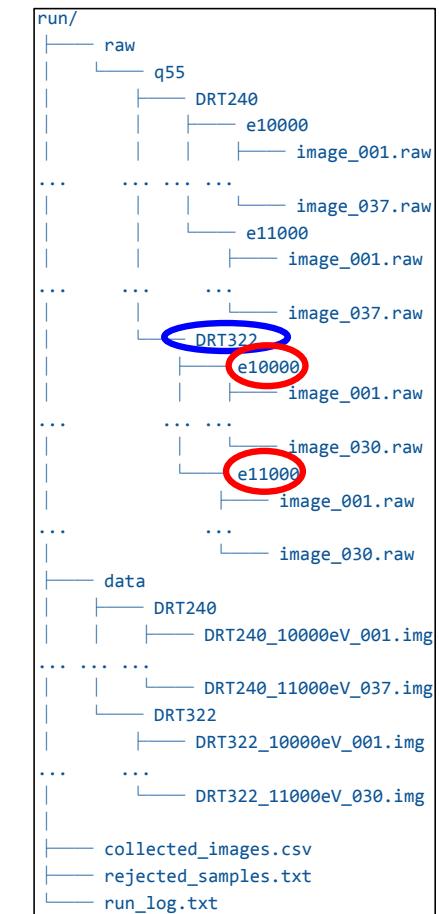
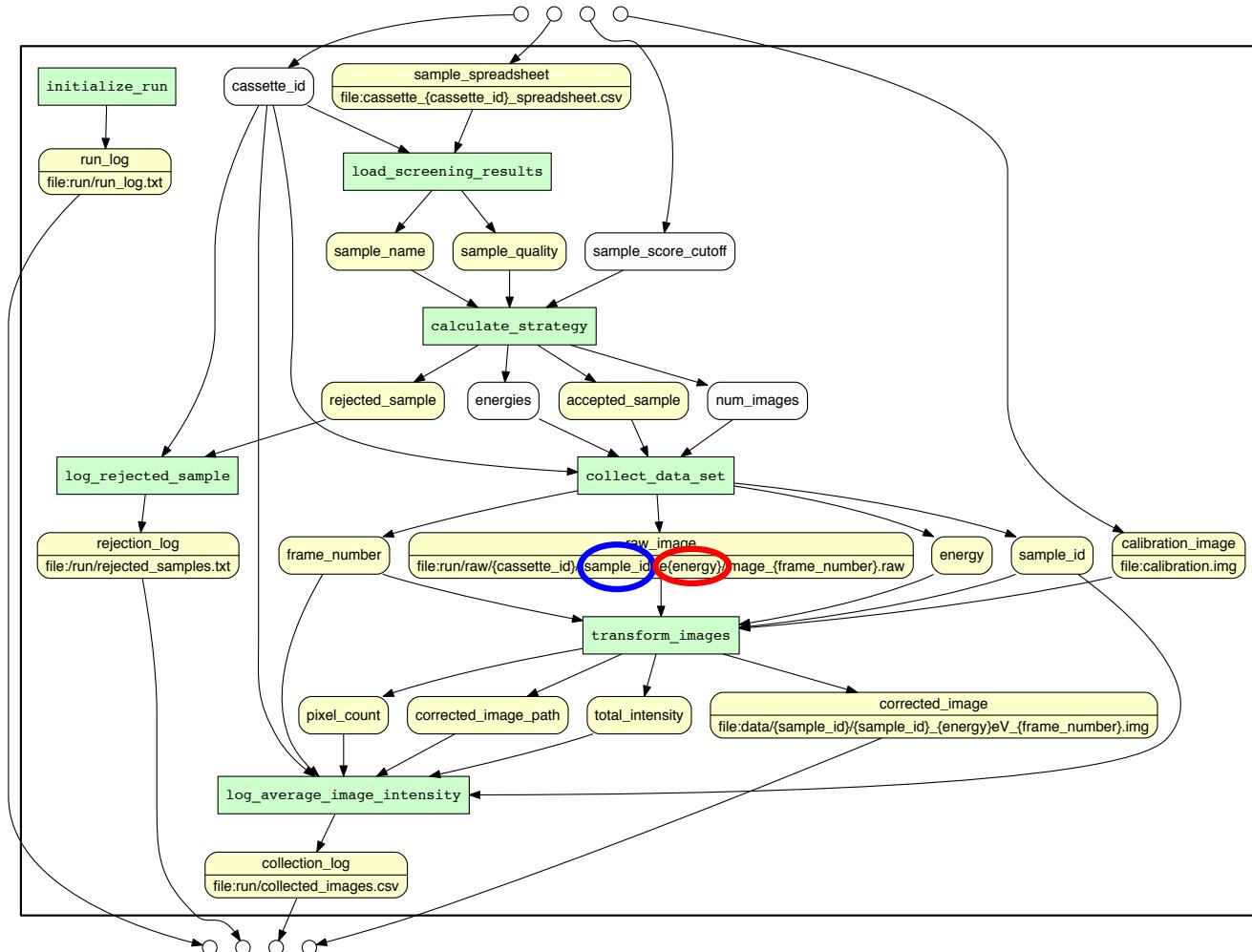
1. YW annotations => YW model
2. Files & Folders left by a run => runtime (meta-)data



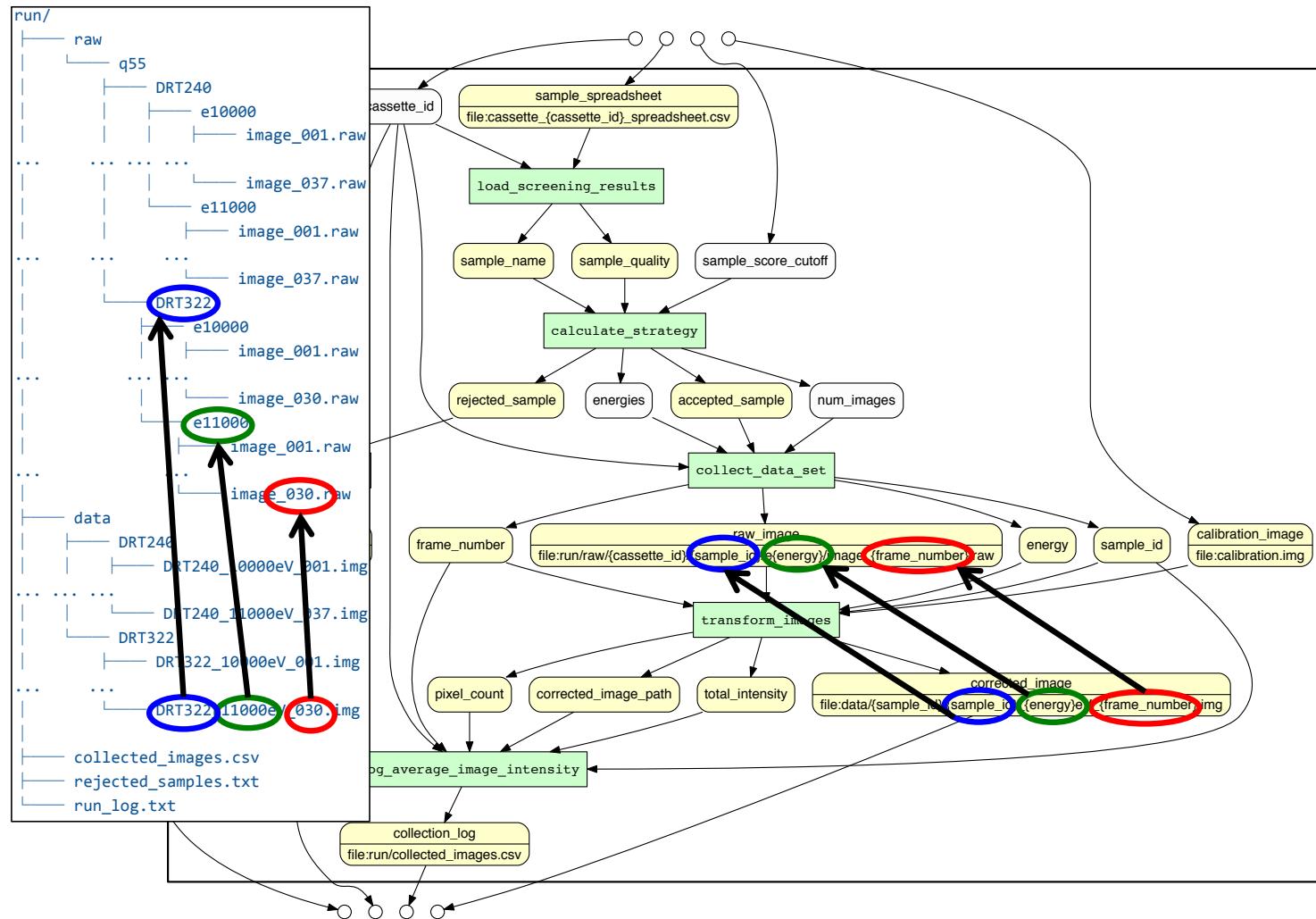
Q₁: What **samples** did the script run collect images from?



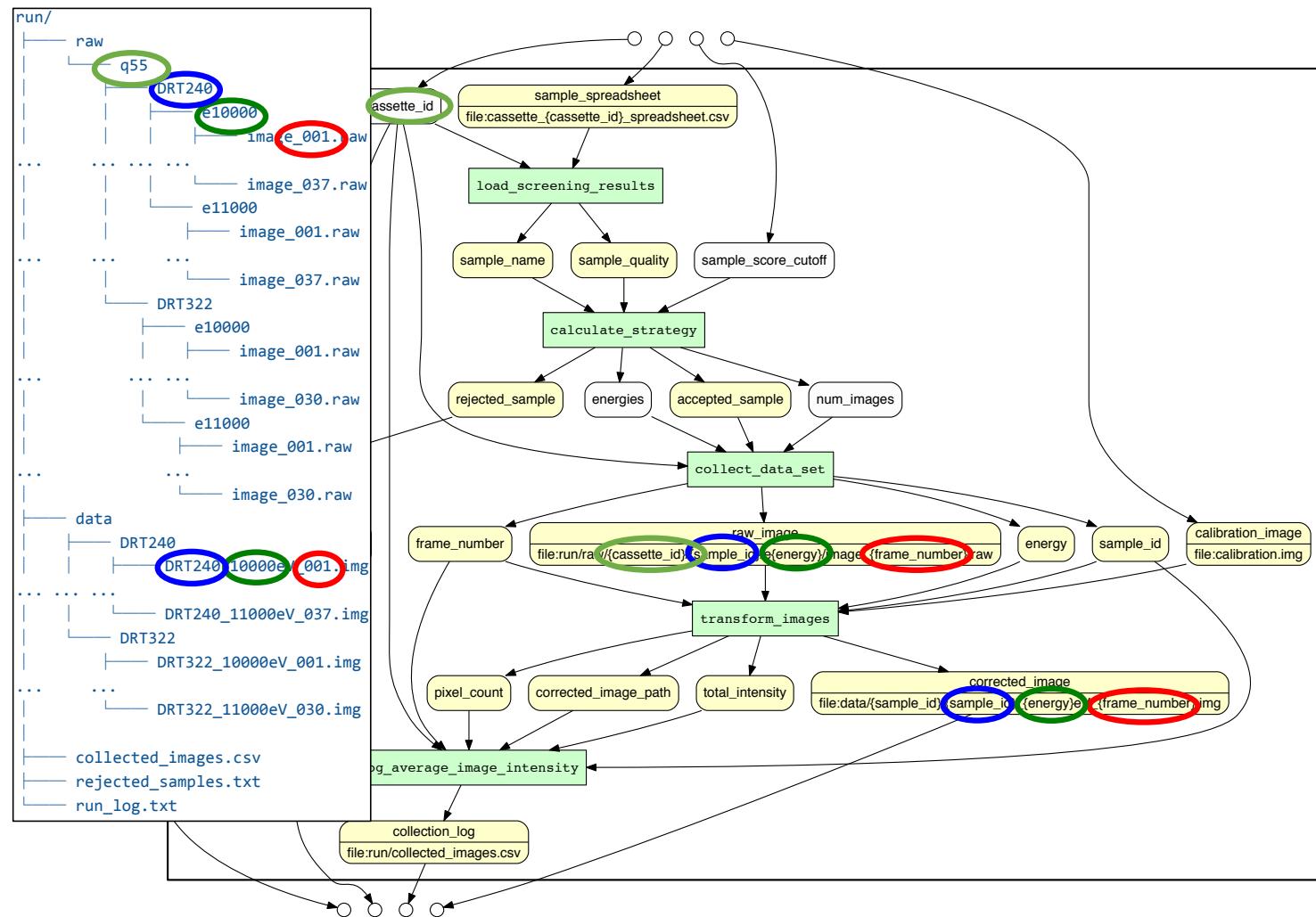
Q₂: What **energies** were used for image collection from sample **DRT322**?

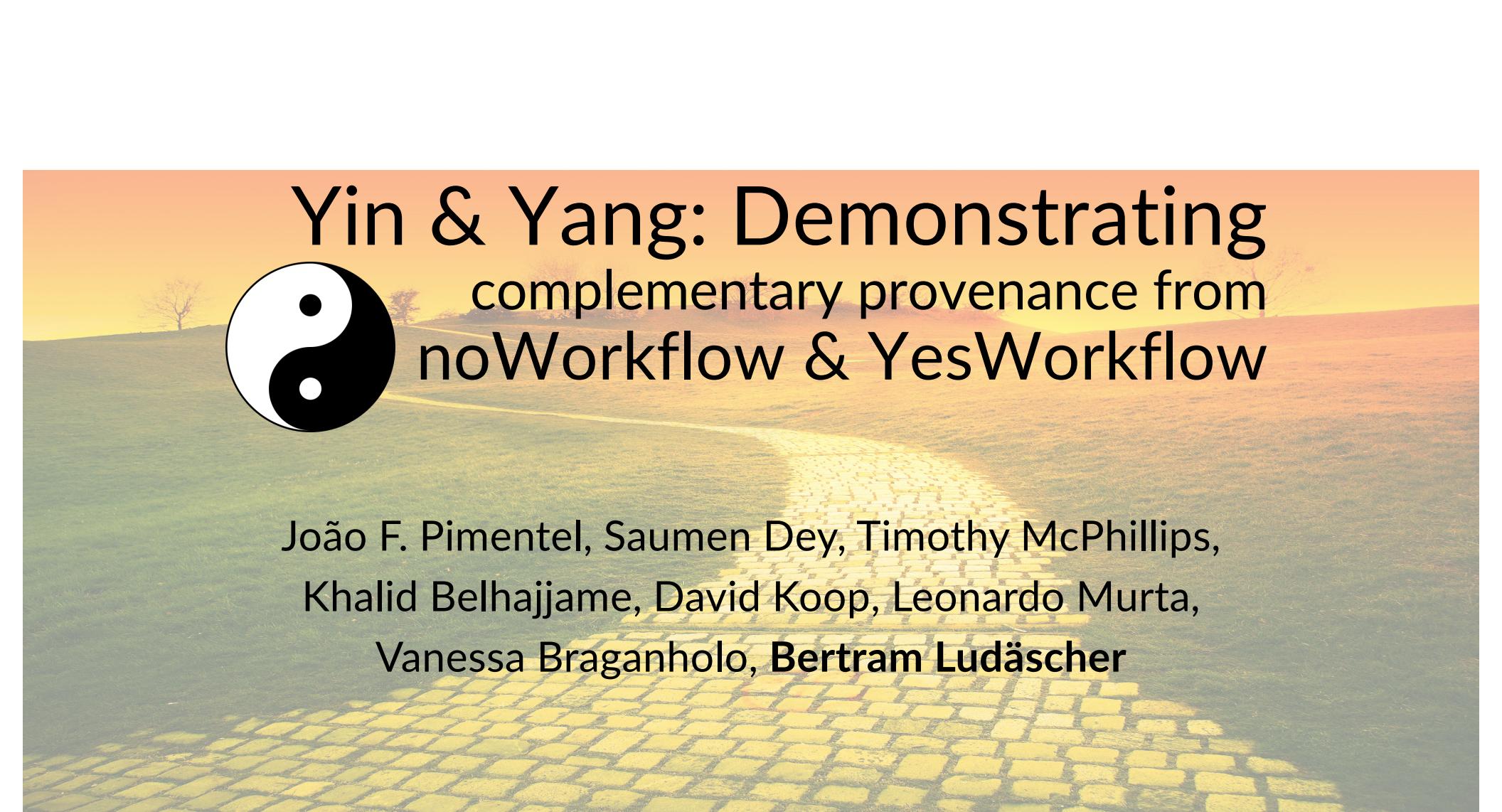


Q3: Where is the raw image of corrected image DRT322_11000ev_030.img?



Q₅: What cassette-id had the sample leading to DRT240_10000ev_001.img?





Yin & Yang: Demonstrating complementary provenance from noWorkflow & YesWorkflow



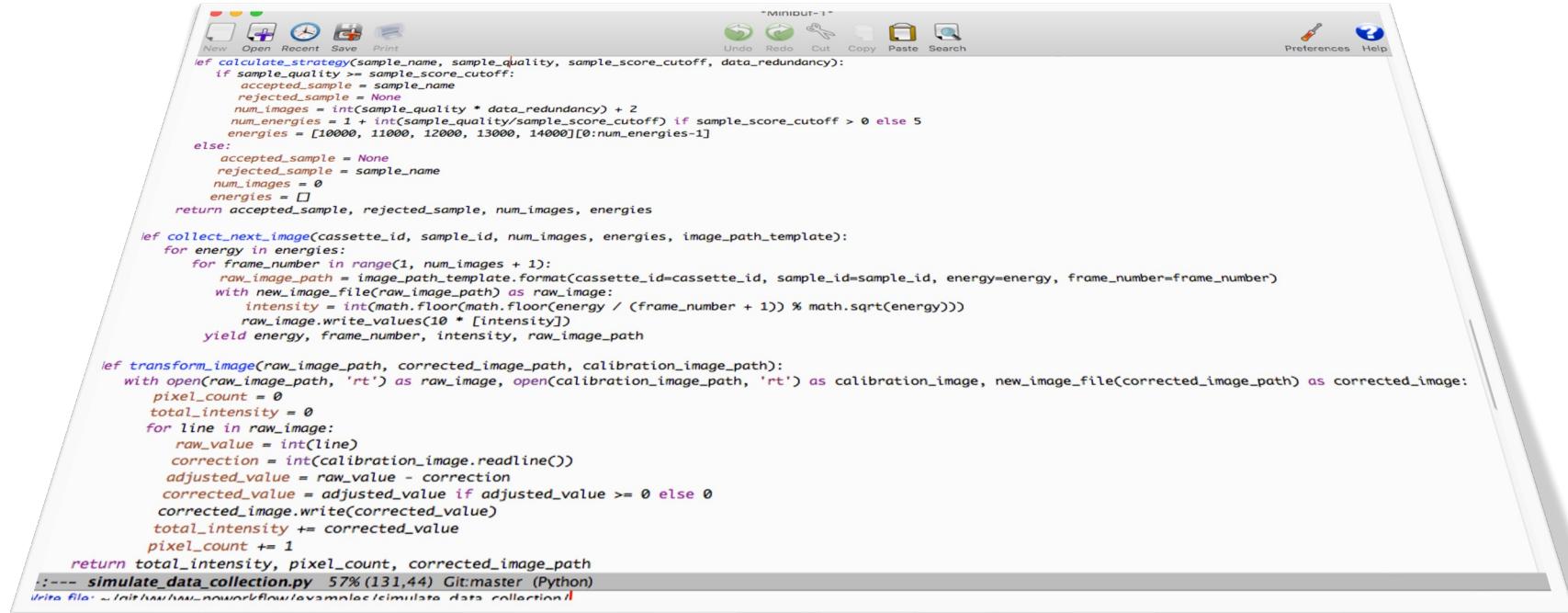
João F. Pimentel, Saumen Dey, Timothy McPhillips,
Khalid Belhajjame, David Koop, Leonardo Murta,
Vanessa Braganholo, Bertram Ludäscher



ILLINOIS
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN



Using Provenance from Script Runs



```
def calculate_strategy(sample_name, sample_quality, sample_score_cutoff, data_redundancy):
    if sample_quality >= sample_score_cutoff:
        accepted_sample = sample_name
        rejected_sample = None
        num_images = int(sample_quality * data_redundancy) + 2
        num_energies = 1 + int(sample_quality/sample_score_cutoff) if sample_score_cutoff > 0 else 5
        energies = [10000, 11000, 12000, 13000, 14000][0:num_energies-1]
    else:
        accepted_sample = None
        rejected_sample = sample_name
        num_images = 0
        energies = []
    return accepted_sample, rejected_sample, num_images, energies

def collect_next_image(cassette_id, sample_id, num_images, energies, image_path_template):
    for energy in energies:
        for frame_number in range(1, num_images + 1):
            raw_image_path = image_path_template.format(cassette_id=cassette_id, sample_id=sample_id, energy=energy, frame_number=frame_number)
            with new_image_file(raw_image_path) as raw_image:
                intensity = int(math.floor(math.floor(energy / (frame_number + 1)) % math.sqrt(energy)))
                raw_image.write_values(10 * [intensity])
            yield energy, frame_number, intensity, raw_image_path

def transform_image(raw_image_path, corrected_image_path, calibration_image_path):
    with open(raw_image_path, 'rt') as raw_image, open(calibration_image_path, 'rt') as calibration_image, new_image_file(corrected_image_path) as corrected_image:
        pixel_count = 0
        total_intensity = 0
        for line in raw_image:
            raw_value = int(line)
            correction = int(calibration_image.readline())
            adjusted_value = raw_value - correction
            corrected_value = adjusted_value if adjusted_value >= 0 else 0
            corrected_image.write(corrected_value)
            total_intensity += corrected_value
            pixel_count += 1
    return total_intensity, pixel_count, corrected_image_path
:--- simulate_data_collection.py 57% (131,44) Git:master (Python)
Uruta filo: ~ /init/han/hsu-nmrworkflows/exmaple/simulate_data_collection.d
```

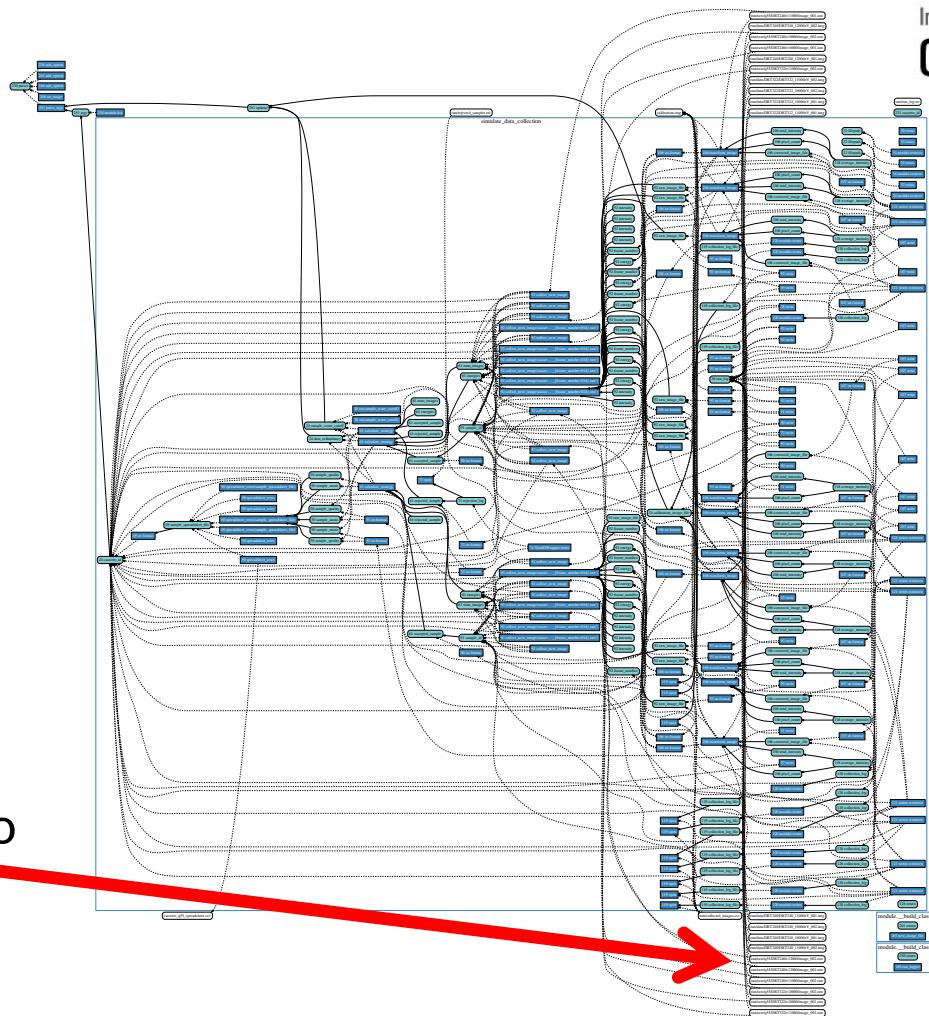
Example from the log-file:

2016-06-07 20:32:36 Wrote run/data/DRT240/**DRT240_11000eV_002.img**

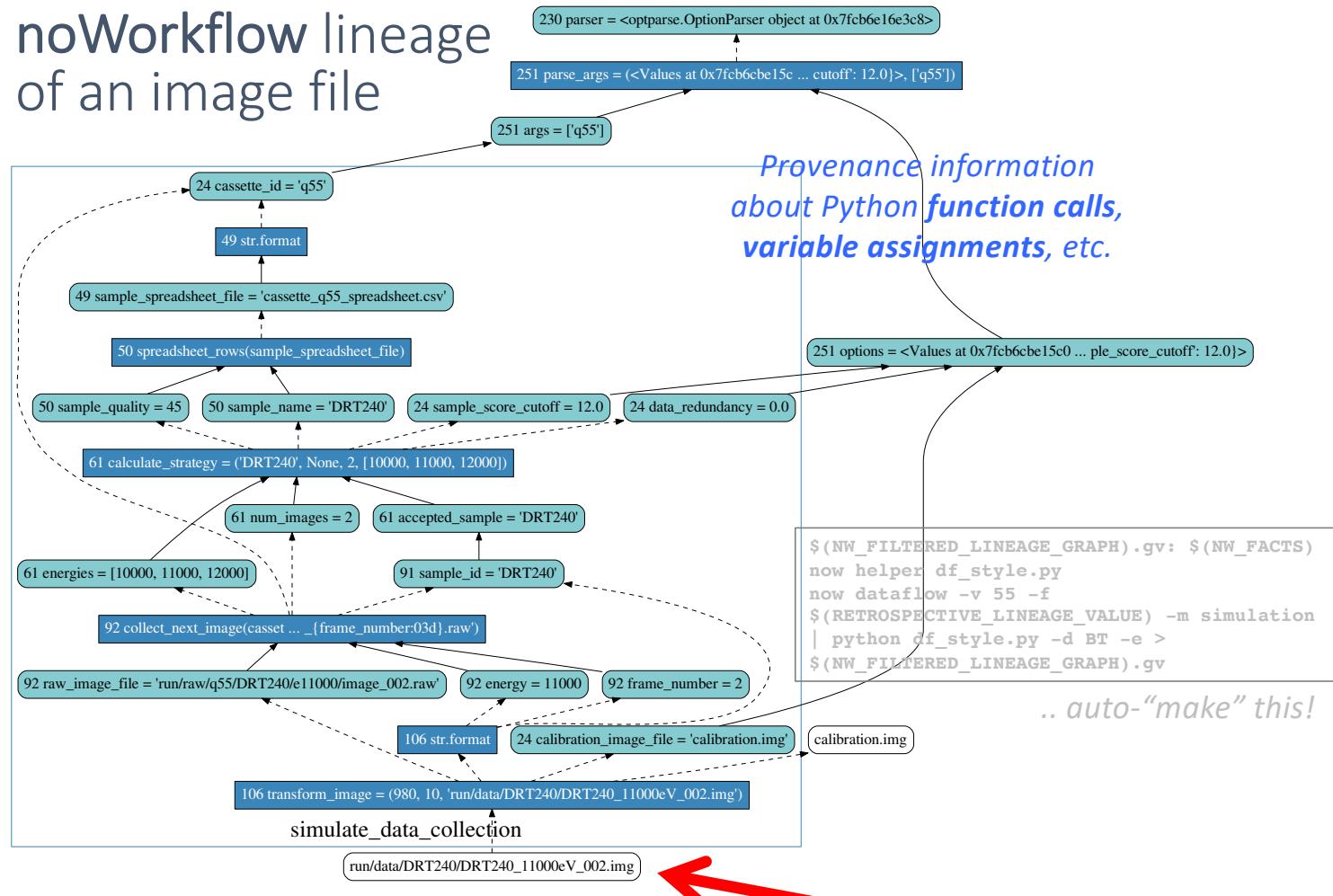
But how was that image derived?? ("Provenance for Self!")

noWorkflow: not only Workflow!

- Scripts have provenance, too!
- Transparently capture some/all provenance from Python script runs.
- Use filter queries to “zoom” into relevant parts ..



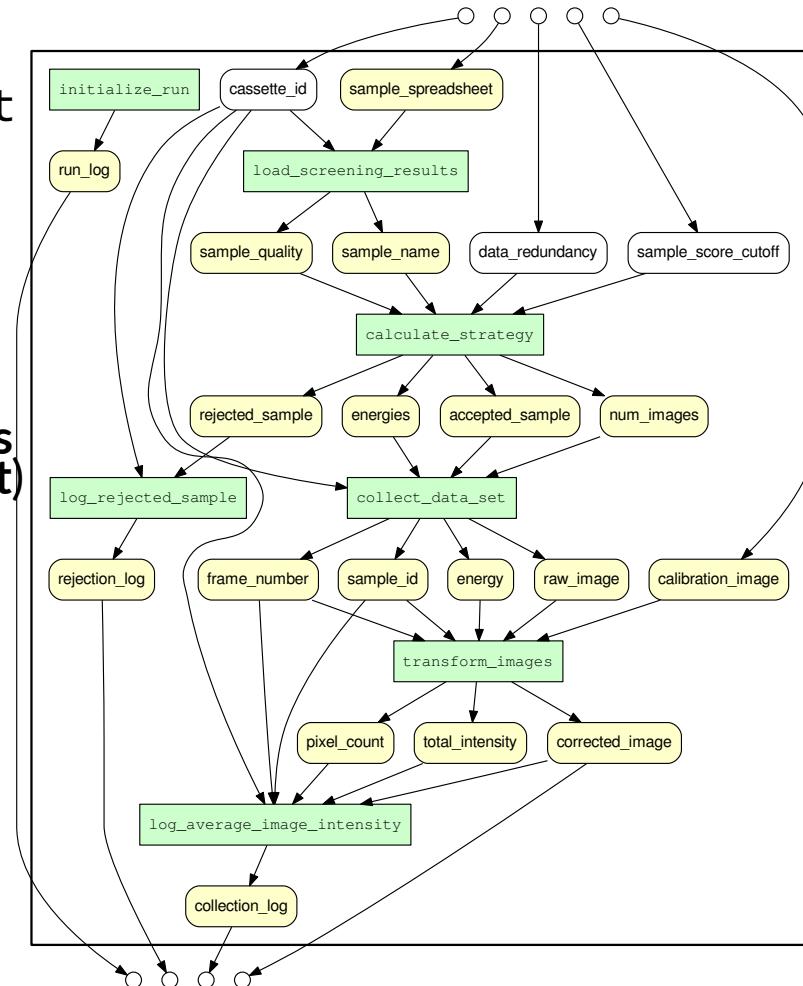
noWorkflow lineage of an image file



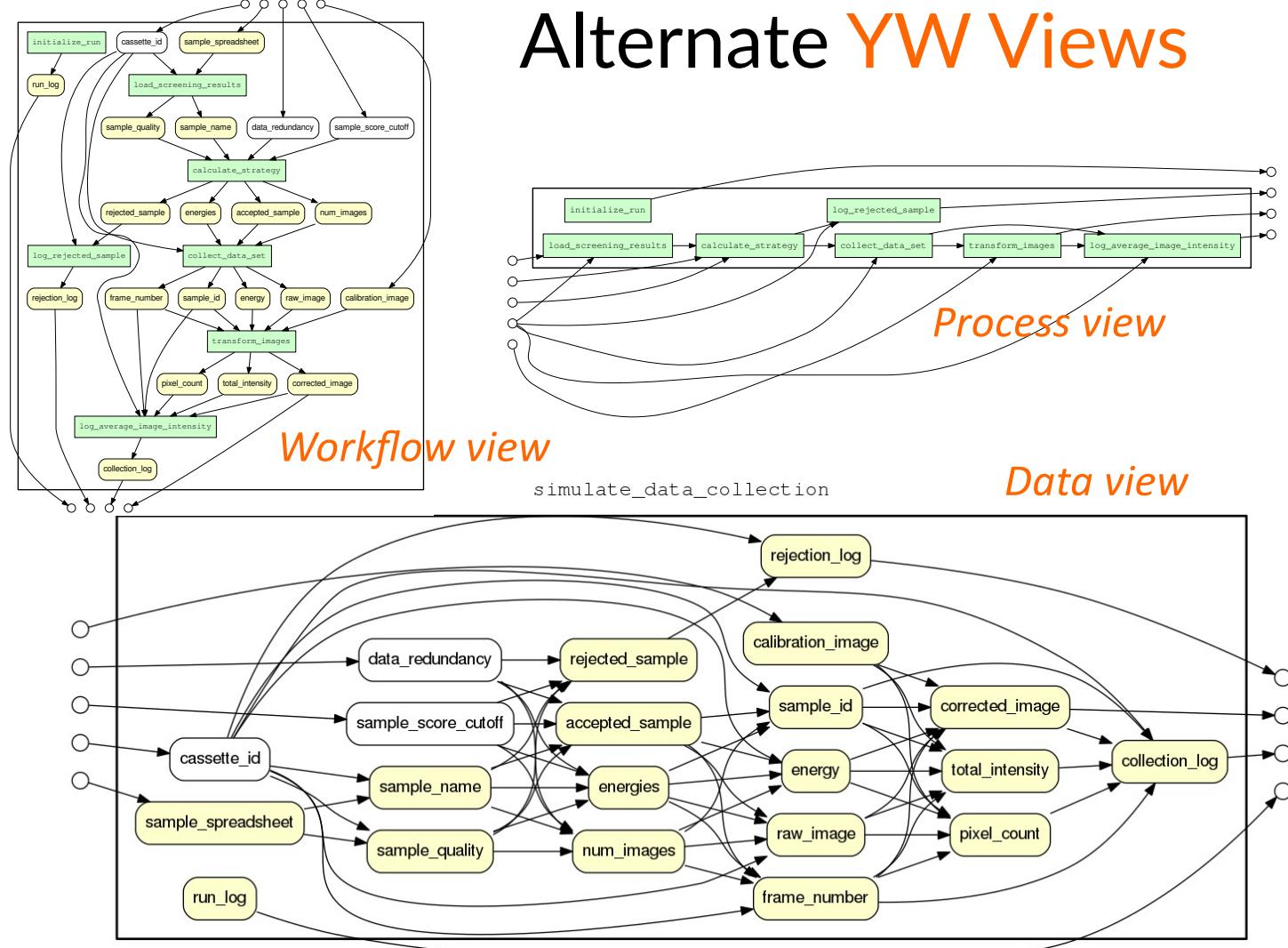
\$ now dataflow -f "run/data/DRT240/**DRT240_11000eV_002.img**"

YesWorkflow: Yes, scripts are Workflows, too!

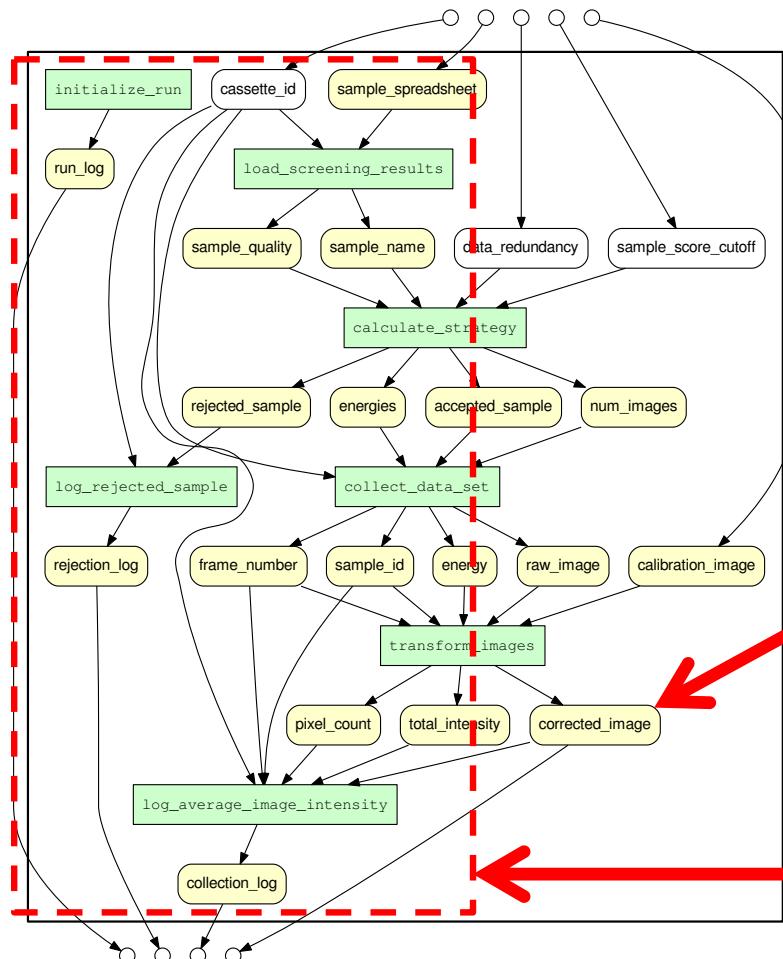
- Use **YW annotations**
@begin...@end, @in, @out
to reveal hidden
conceptual workflow
(prospective provenance)
- Script isn't changed:
 - annotations via **comments**
(=> language independent)
- For understanding and sharing the “big picture”
- **Query** and visualize!



Alternate YW Views



What is the lineage of “corrected_image”?

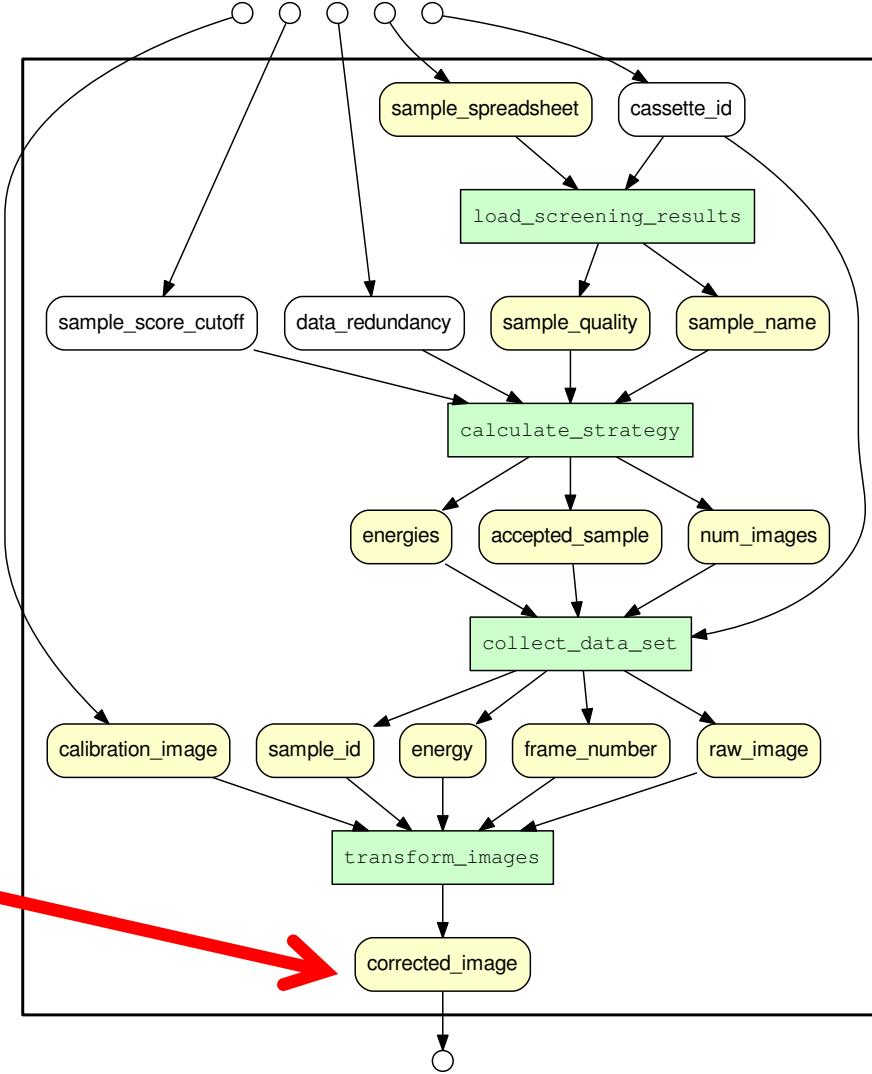


*From here on “upwards”:
What led (leads) to this?*

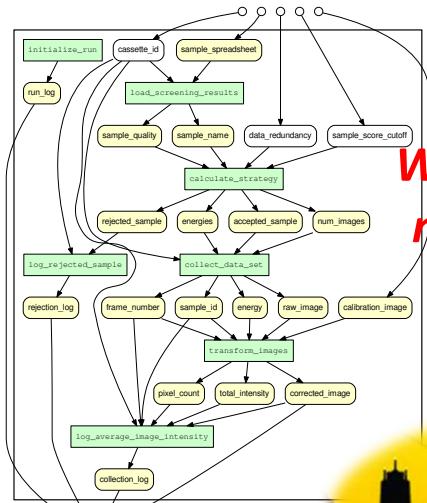
*.. and what is irrelevant
and should be pruned??*

Subgraph resulting from lineage query on YW workflow model

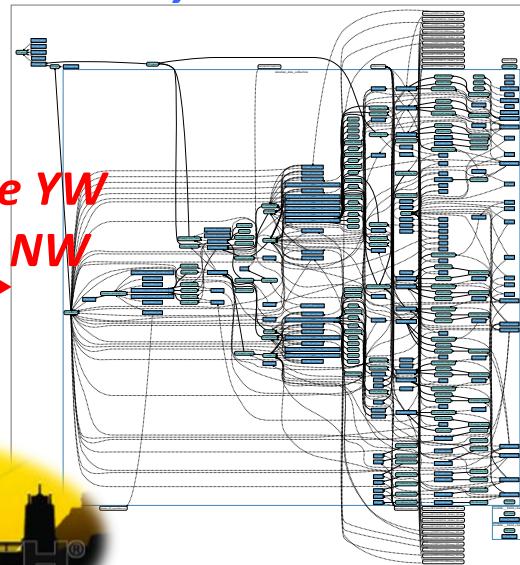
What is the lineage of corrected_image?



YesWorkflow:
Conceptual workflow model

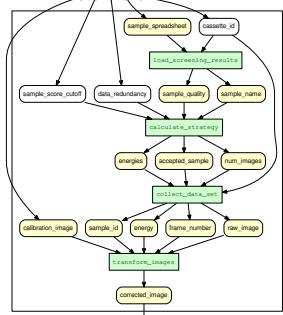


noWorkflow:
Python trace model



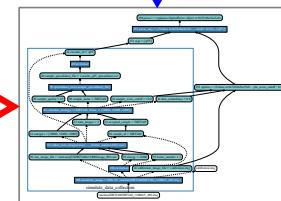
Would like to use YW
model to query NW
← - - - →
data!

lineage query



But how do we
bridge this gap???

lineage query



“Workflow-Land” (*prospective* provenance) \Leftrightarrow “Trace-Land” (*retrospective* provenance)

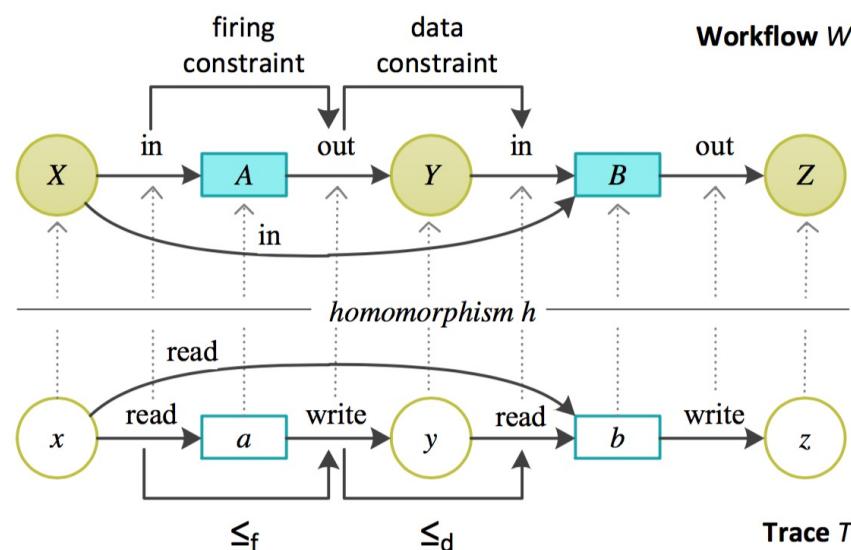
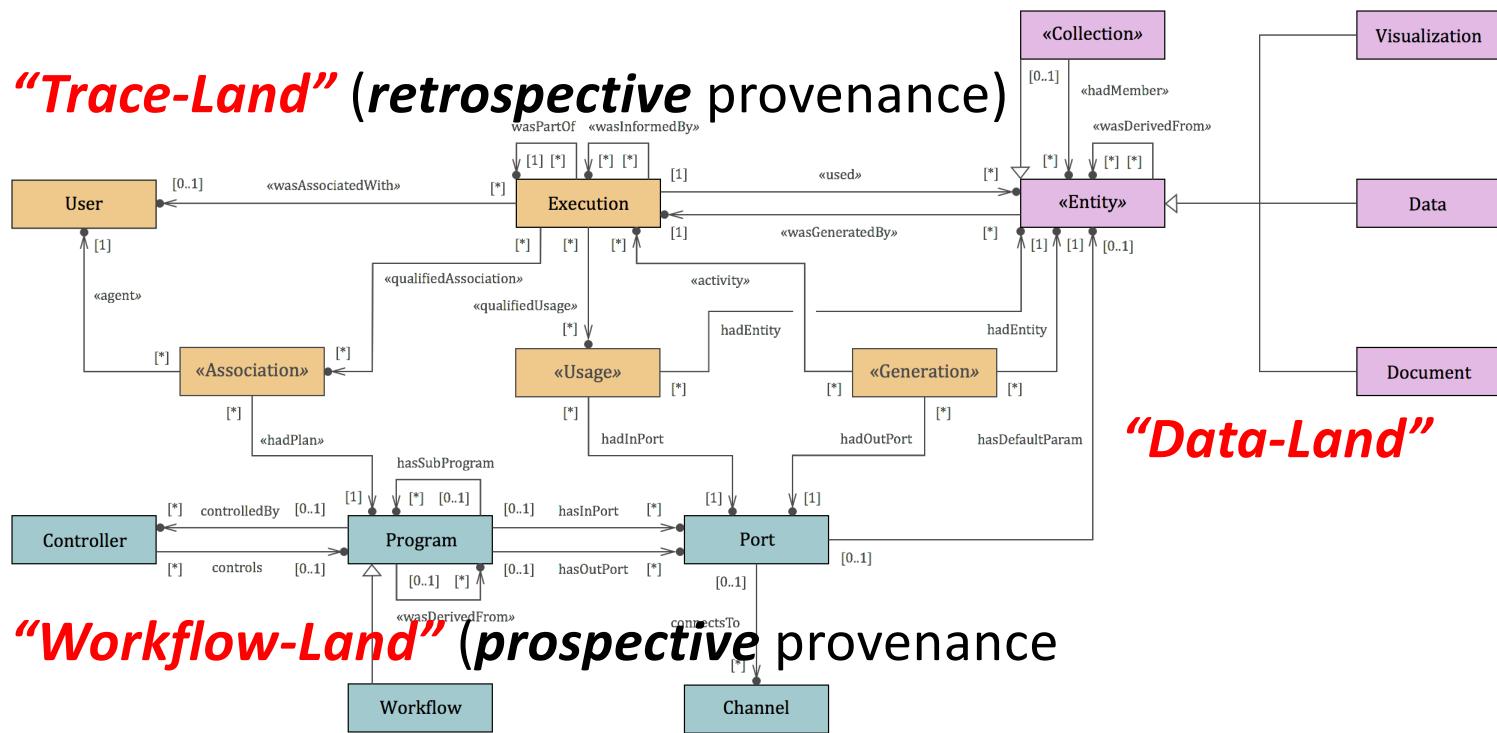


Figure 6: A homomorphism h from trace T to workflow W guarantees structural validity. Workflow-level constraints induce temporal constraints \leq_f and \leq_d on traces [DKBL12].

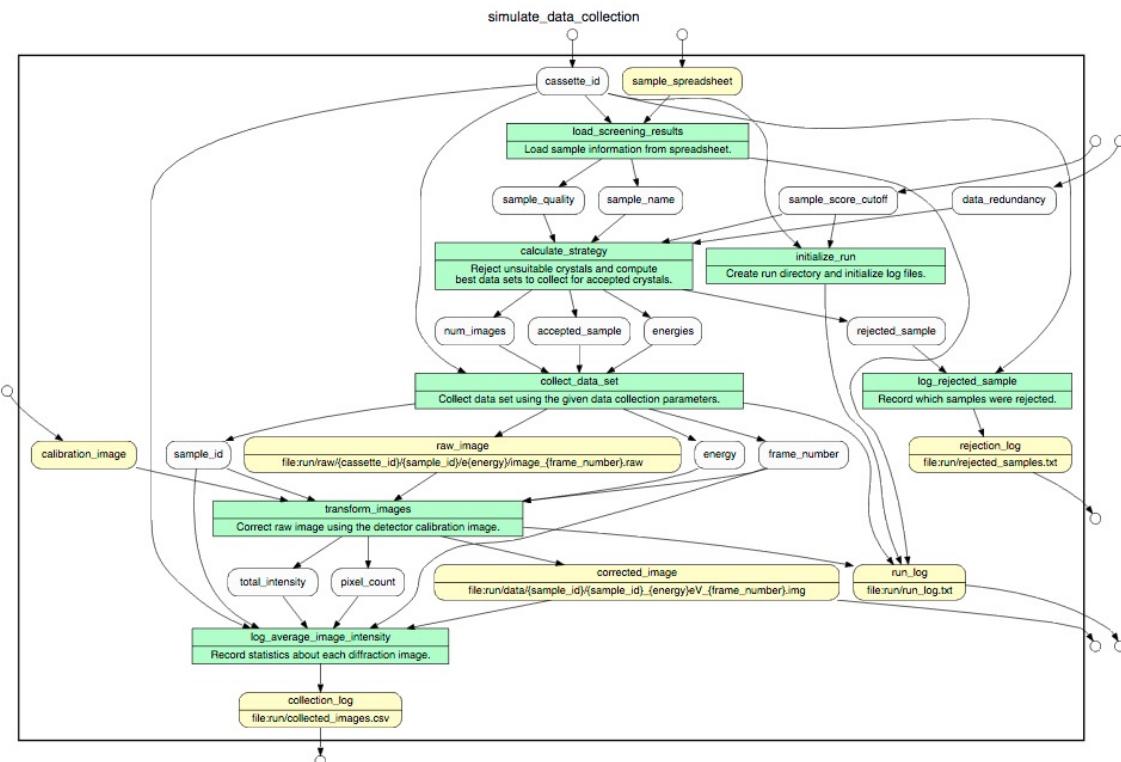
ProvONE: PROV for scientific workflows



<https://purl.dataone.org/provone-v1-dev>

Theory & Practice of Data Cleaning

YesWorkflow Demo



New: Container-based YW Web-app

The screenshot shows the GitHub repository page for `CIRSS/yw-web-app`. The repository is public and was generated from `repos-dev/repo-template`. The main interface includes a code editor with the `master` branch, a list of 21 commits by Timothy McPhillips, and sections for About, Releases, Packages, and Languages.

About
REPRO for easily running the service supporting the YesWorkflow web application.

Releases
No releases published [Create a new release](#)

Packages
No packages published [Publish your first package](#)

Languages

| Language | Percentage |
|------------|------------|
| Dockerfile | 66.2% |
| Makefile | 10.4% |
| Shell | 20.8% |
| Batchfile | 2.6% |

<https://github.com/CIRSS/yw-web-app>