

Retrospective Provenance **without** a Runtime Provenance Recorder

Timothy McPhillips

Shawn Bowers

Khalid Belhajjame

Bertram Ludäscher*

TaPP 2015
Edinburgh, Scotland
July 8-9, 2015



Overview / Exec Summary

- Scientific Workflows: ASAP!
- **Scripts are Scientific Workflows, too!**
 - .. so how can we help workflow script authors?
- **YesWorkflow (YW)**
 - **Prospective** provenance through YW-annotations:
 - script + annotations @begin, @end, @in, @out => workflow model
- **YW-Recon: Retrospective** provenance **without** a provenance recorder:
 - .. add @URI template-annotations
=> linking script persisted (meta-)data with workflow model

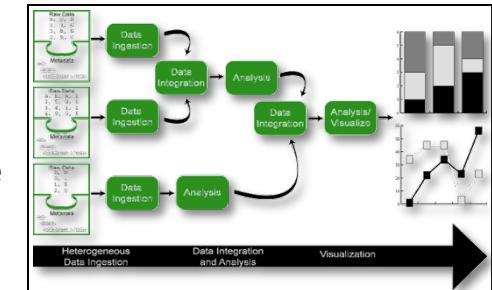
... => YW Retrospective Provenance Queries

- .. using YW workflow model to query YW-reconstructed runtime provenance!

Scientific Workflows: ASAP!

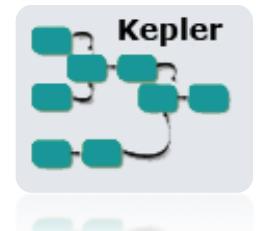
- **Automation**

- wfs to automate computational aspects of science



- **Scaling** (exploit and optimize *machine cycles*)

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



- **Abstraction, Evolution, Reuse** (*human cycles*)

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



- **Provenance**

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



Scientific Workflows ...

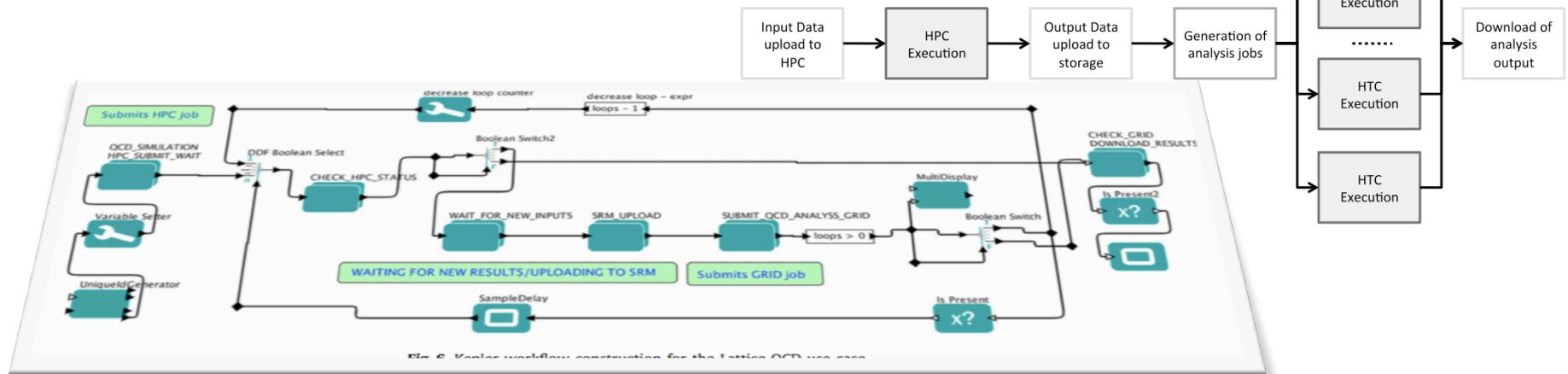


Fig. 6. Kepler workflow construction for the Lattice QCD use case.

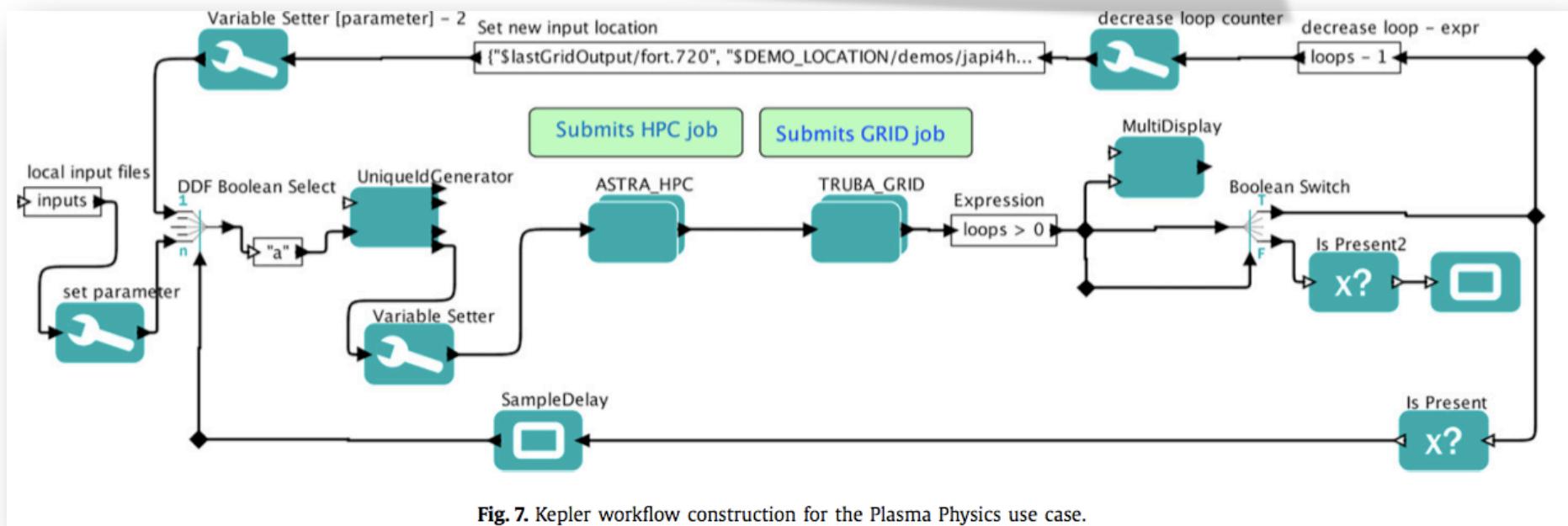


Fig. 7. Kepler workflow construction for the Plasma Physics use case.

... are a wonderful thing ...

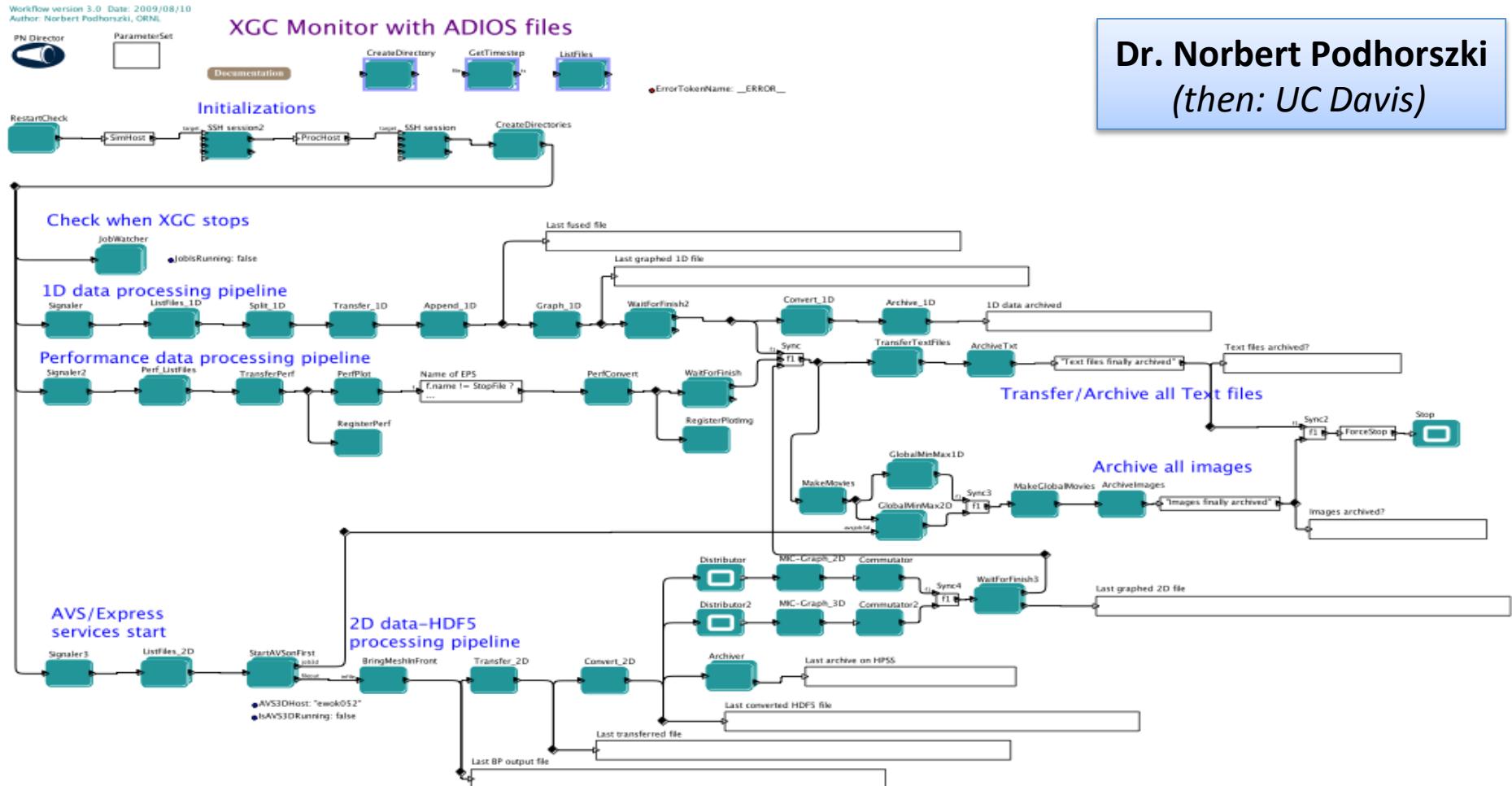
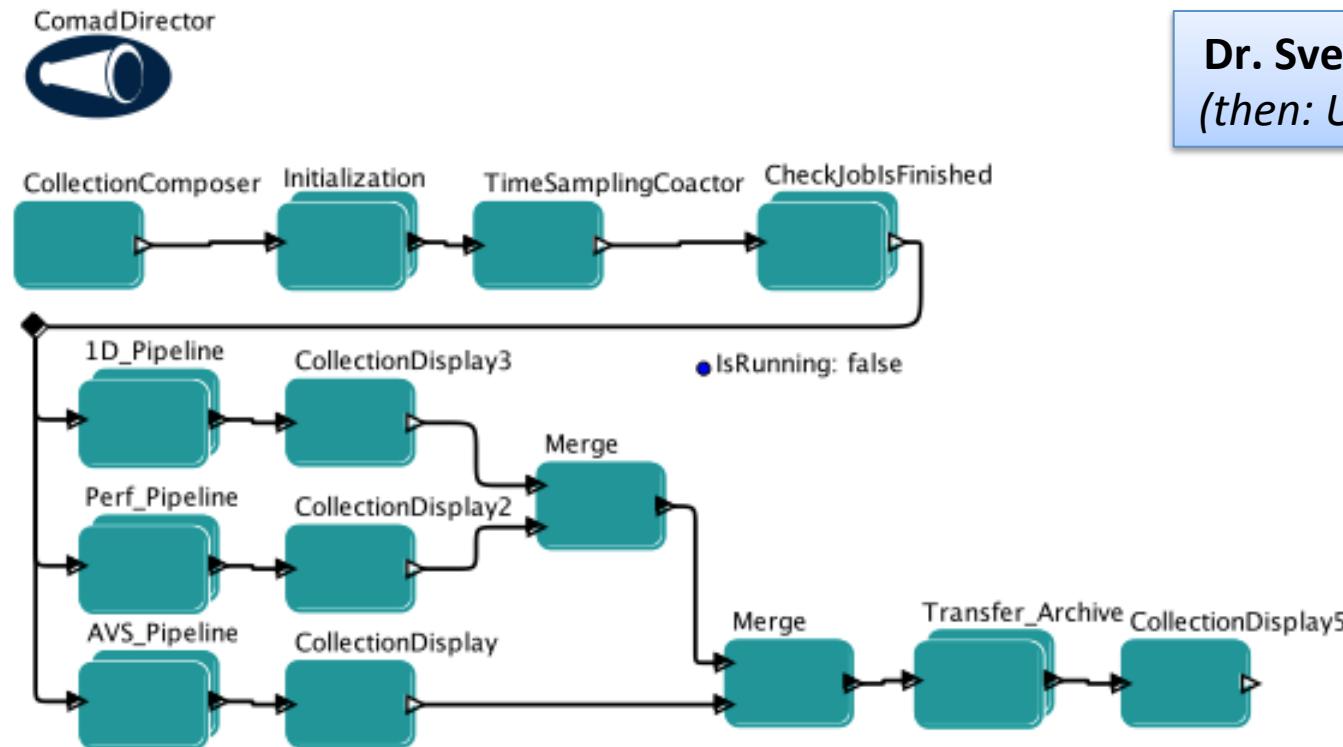


Figure 5.1: The original Monitoring Workflow. The top-level PN workflow shows all steps of the workflow. Additional parameters are contained in the parameter set in the top left corner. In the upper left corner, the workflow starts with initialization actors. Then, four pipelines that use a polling technique to observe a job execution follow: The fist pipeline just contains one actor that observes the job execution and adjusts the flag "JobIsRunning" according to the status of the job. The other three pipelines start with a sampling actor

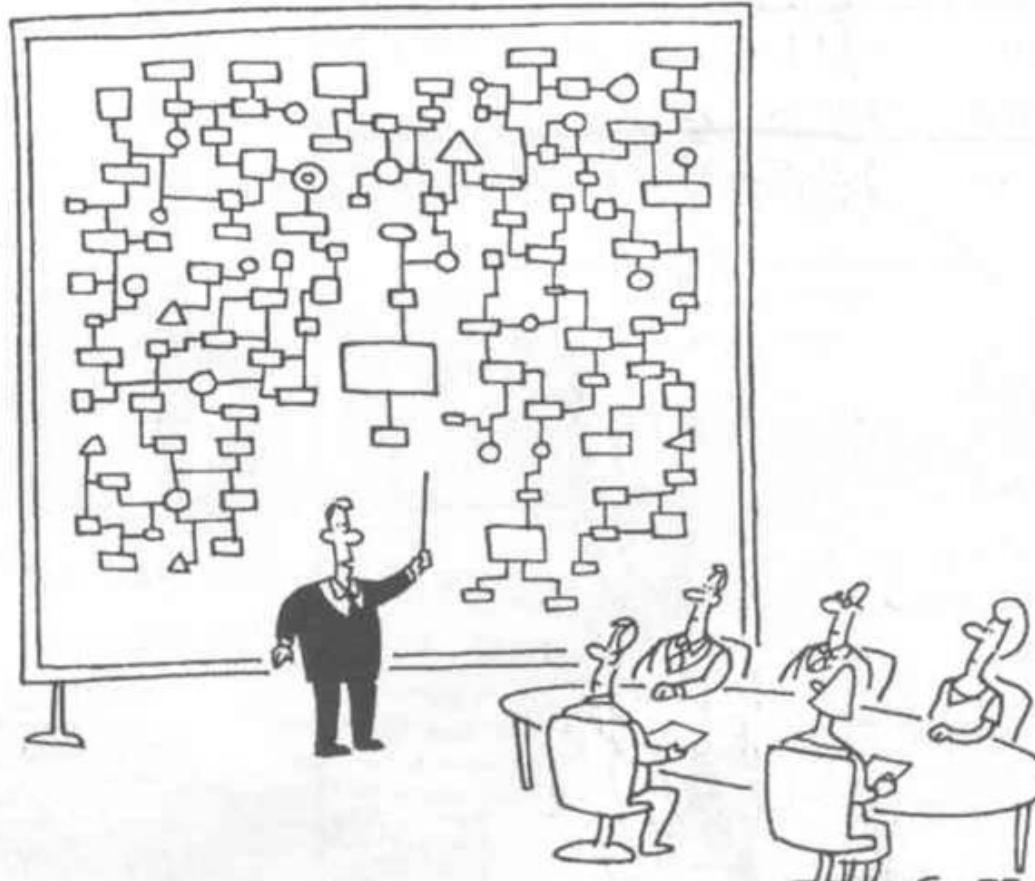
... after simplifying a bit (here: Kepler/COMAD)



Dr. Sven Köhler
(then: UC Davis)

Figure 5.2: Remodeled Monitoring Workflow using COMAD. Top-level model uses the COMAD director. Composites use SDF and DDF directors, but they do not need to handle stop tokens anymore.

I beg your pardon, I never promised you ..



*“Thanks to our **Graphical UI** your scientific workflows will be much easier to develop, understand and maintain!”*

Hmm... this was supposed to be **easier** than programming!

Meanwhile, on a nearby planet ...

Reproducible academic publications

This section contains academic papers that have been published in the peer-reviewed literature or pre-print sites such as the [ArXiv](#) that include one or more notebooks that enable (even if only partially) readers to reproduce the results of the publication. If you include a publication here, please link to the journal article as well as providing the nbviewer notebook link (and any other relevant resources associated with the paper).

- Automatic segmentation of odor maps in the mouse olfactory bulb using regularized non-negative matrix factorization, by J. Soelter et al. (Neuroimage 2014, Open Access). The [notebook](#) allows to reproduce most figures from the paper and provides a deeper look at the data. The [full code repository](#) is also available.
 - Multi-tiered genomic analysis of head and neck cancer ties TP53 mutation to 3p loss, by A. Gross et al. (*Nature Genetics* 2014). The [full collection of notebooks](#) to replicate the results.
 - powerlaw: a Python package for analysis of heavy-tailed distributions, by J. Alstott et al.. [Notebook of examples in manuscript](#), [ArXiv link](#) and [project repository](#).
 - Collaborative cloud-enabled tools allow rapid, reproducible biological insights, by B. Ragan-Kelley et al.. The [main notebook](#), the [full collection of related notebooks](#) and the [companion site](#) with the Amazon AMI information for reproducing the full paper.
 - A Reference-Free Algorithm for Computational Normalization of Shotgun Sequencing Data, by C.T. Brown et al.. [Full notebook](#), [ArXiv link](#) and [project repository](#).
 - The kinematics of the Local Group in a cosmological context by J.E. Forero Romero et al.. [The Full notebook](#) and also all the data in a [github repo](#).
 - Warming Ocean Threatens Sea Life, an article in *Scientific American* by Roberto de Almeida from [MarineBio](#). [The notebook](#) for its main plot. By [display\(i\)](#)



Data-driven journalism

- The Need for Openness in Data Journalism, by Eric Singer-Vine.
 - St. Louis County Segregation Analysis , analysis Area Is Even More Segregated Than You Probably Think, by Eric Singer-Vine.

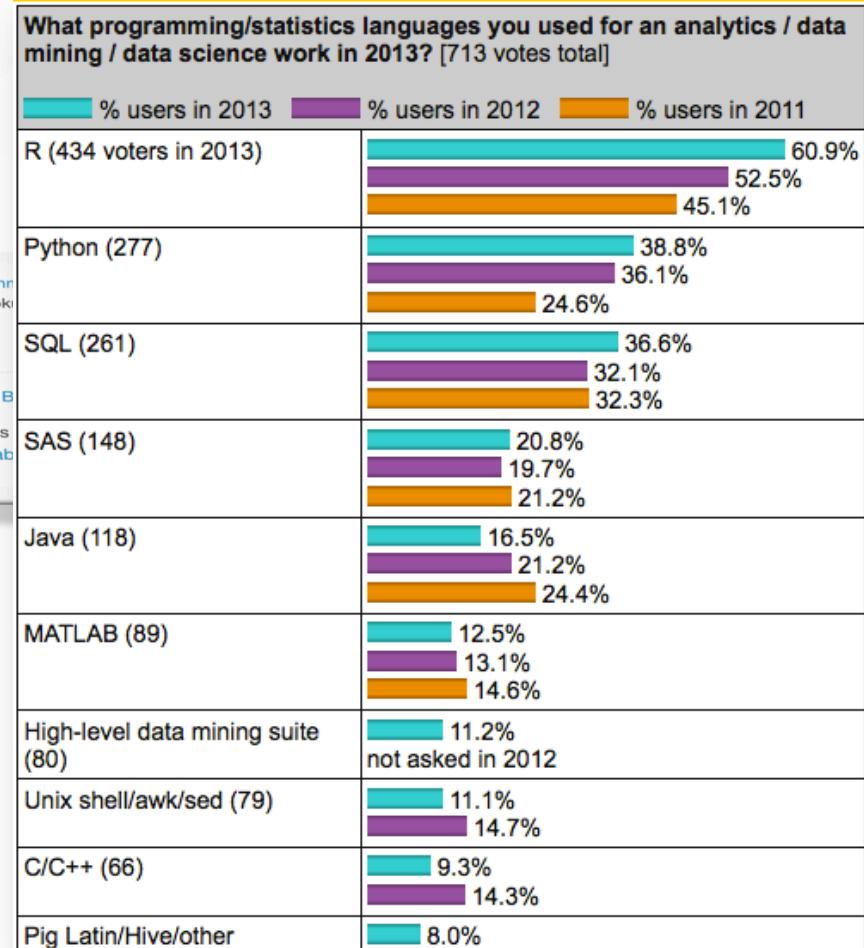
IP[y]: IPython Interactive Computing

```
In [3]: from IPython.display import SVG  
        SVG(filename='python-logo.svg')
```

Out[3]:



python™



SKOPE: Synthesized Knowledge Of Past Environments

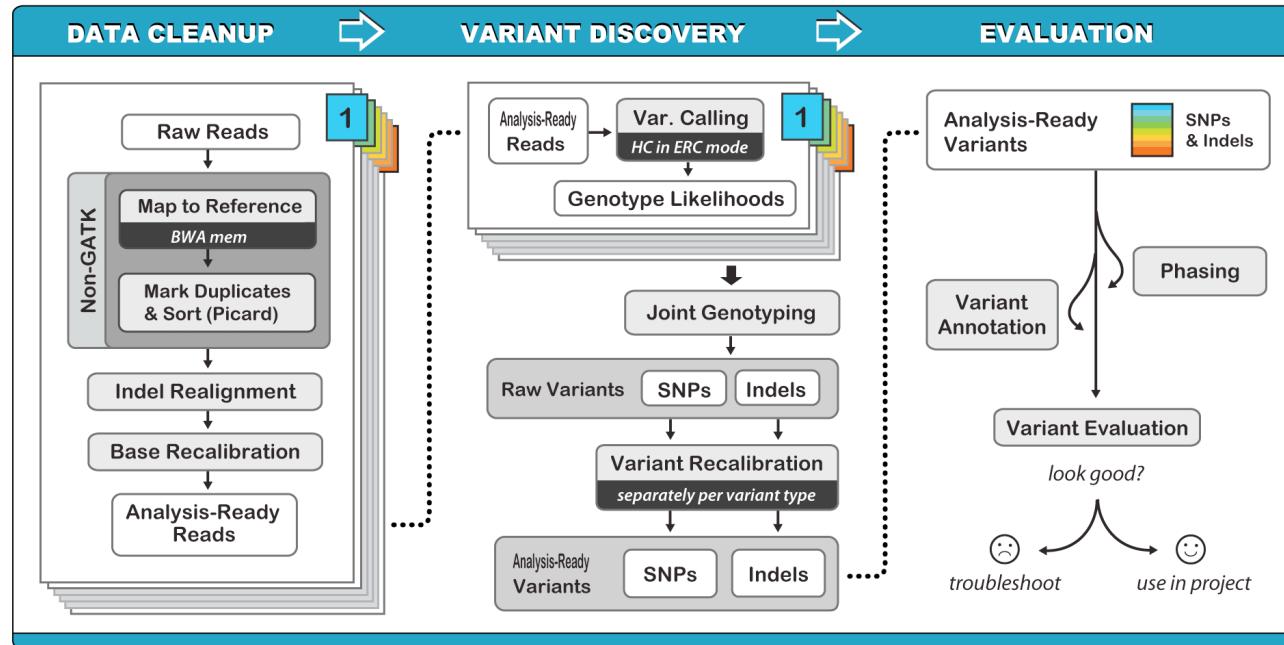
Bocinsky, Kohler *et al.* study rain-fed maize of Anasazi

- 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.



... HPCBio Workflows @ Illinois

Broad Institute: Recommended workflow for variant analysis



Quickly, say: `#!/bin/bash`

*Liudmila Mainzer,
Victor Jongeneel
HPC Bio @ Illinois*



YesWorkflow Provenance @ TaPP'15



*National Petascale
Computing Facility*

It's time to shift control ...



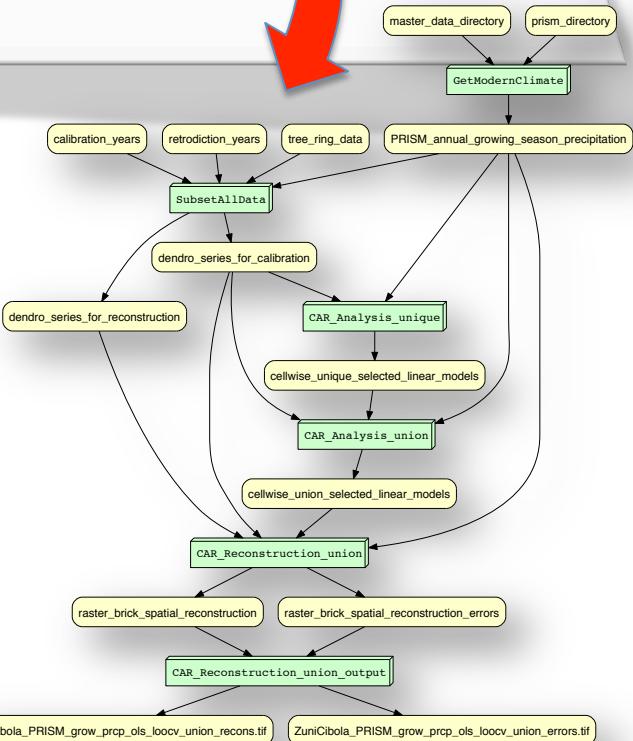
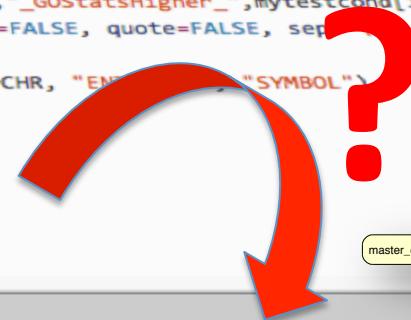
- ... back from being **consumers** of someone else's (= our) tools ..
 - “Just click here!”
- ... to **tool makers!**
 - Scientists who author workflows as scripts!
- **Go where the wild things (*users!*) are ...**
 - Yes, develop for “end users” ...
 - ... but don’t forget the **tool makers!**
- Can we do this **together?**



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

Enter: YesWorkflow! (yesworkflow.org)

- **YesWorkflow (YW)**
 - Grass-roots effort
 - ... meeting the scientists/users **where they R!**
 - R, Matlab, (i)Python, Jupyter, ...
 - Scripts + **simple** user annotations
- => **Reveal the workflow model/abstraction**
 - ... that underlies the (script) *implementation*
- => YW can give us more of **ASAP!**
 - First **YW**: ASAP (Abstraction)...
 - Then **YW-recon**: ASAP (reconstructing **runtime Provenance**)

YesWorkflow.org

YesWorkflow 

Grass-roots effort bringing workflow modeling to scripting languages.

<http://yesworkflow.org/wiki>

[Filters](#) [+ New repository](#)

yw-gui Python ★ 0 ⚡ 0
Graphical interface for rendering graphs produced by YesWorkflow
Updated 8 days ago

yw-tapp-15-recon TeX ★ 0 ⚡ 1
LaTeX sources and examples for the TAPP'15 article on YW provenance reconstruction
Updated 10 days ago

yw-idcc-15 TeX ★ 0 ⚡ 1
LaTeX sources and examples for the IDCC'15 and IJDC articles.
Updated 10 days ago

yw-prototypes Java ★ 6 ⚡ 1
Early implementations of YesWorkflow.
Updated 10 days ago

People 10 >



[Invite someone](#)

Teams 3 >

Owners
3 members · 4 repositories

Developers
10 members · 2 repositories

Paper authors
8 members · 2 repositories

[Create new team](#)

Related Work, other Approaches

... to bring workflow/provenance benefits to scripts:

- **Runtime Provenance Recorders:**
 - use (R, Python, ...) **libraries** and/or **code instrumentation** to capture **runtime observables**
 - file read/write, function calls, program variables & state, ...
 - **noWorkflow** system
 - [Murta-Braganholo-Chirigati-Koop-Freire-IPAW14]
 - exploit Python profiling library to capture runtime provenance
- => helps with "S" and "P"
- **OS-level capture of (system) provenance**
 - Some talks at TaPP !?

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

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}/{1}{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

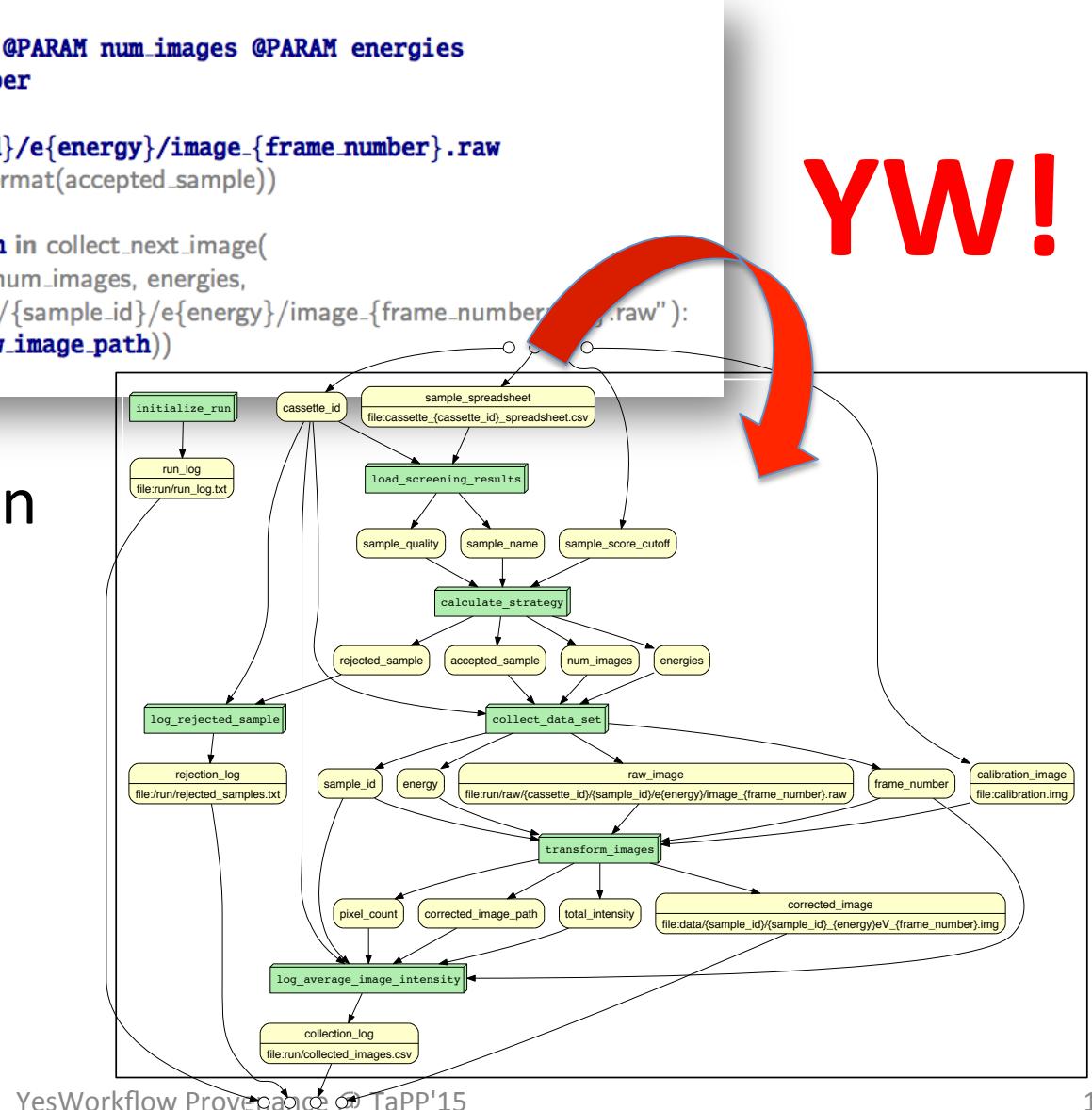
... and data inputs/outputs

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].

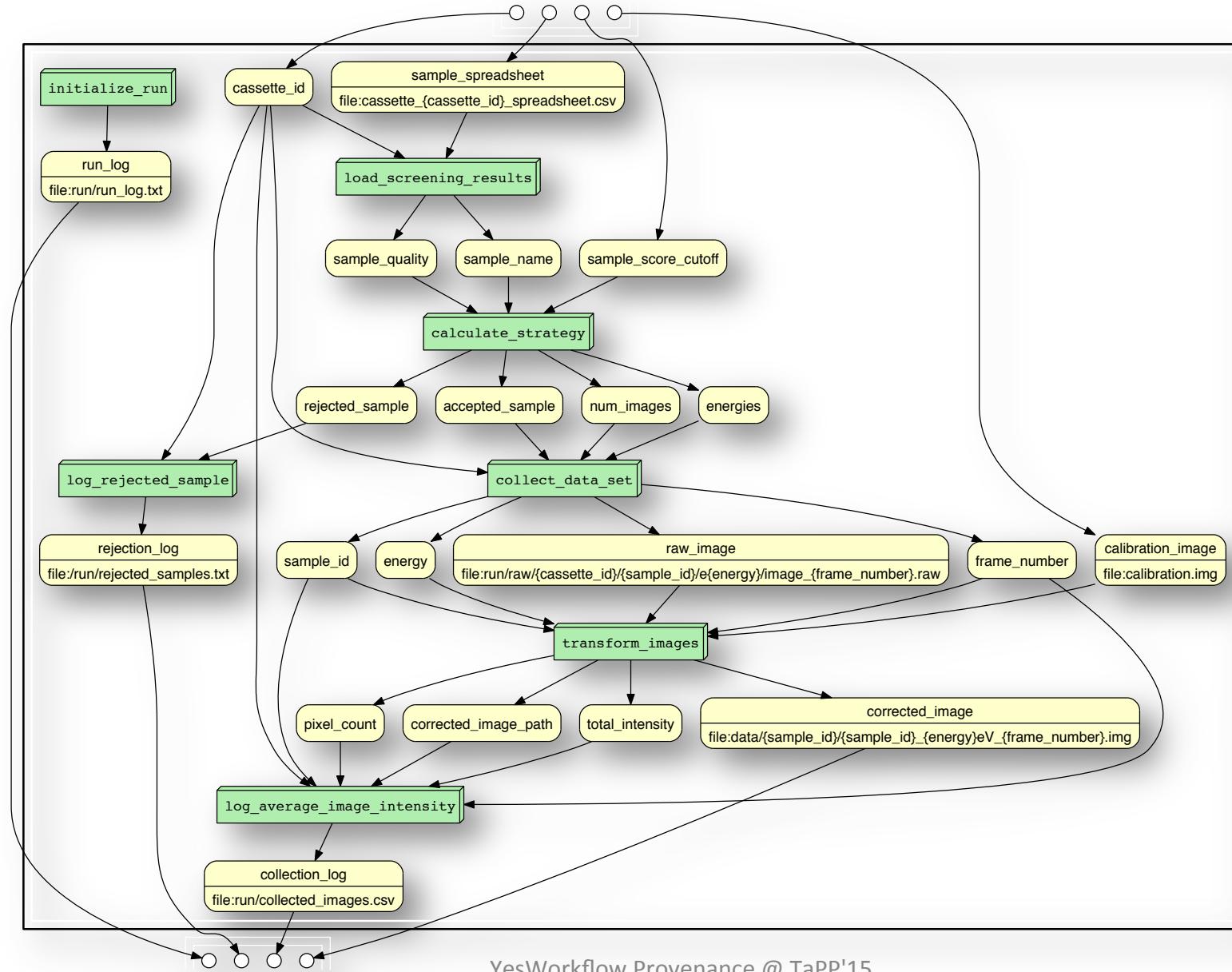
YesWorkflow: Prospective & Retrospective Provenance ... (almost) for free!

```
# @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
```

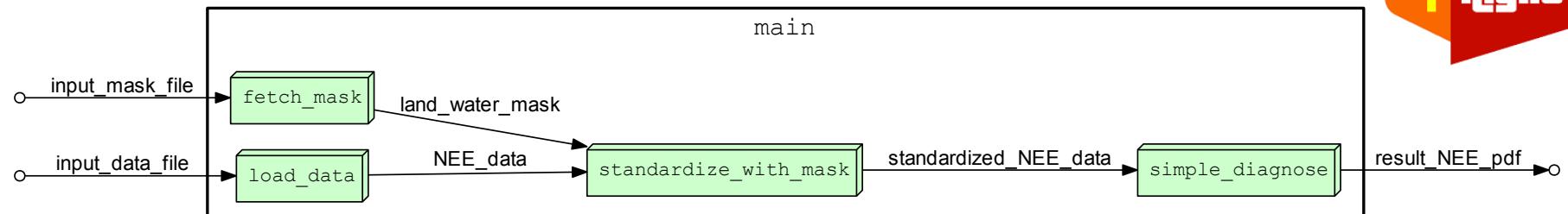
- YW annotations in the script (R, Python, Matlab) are used to recreate the workflow view from the script ...



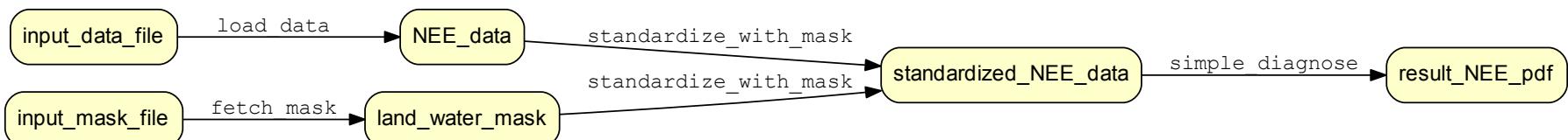
Voila! The Workflow revealed!



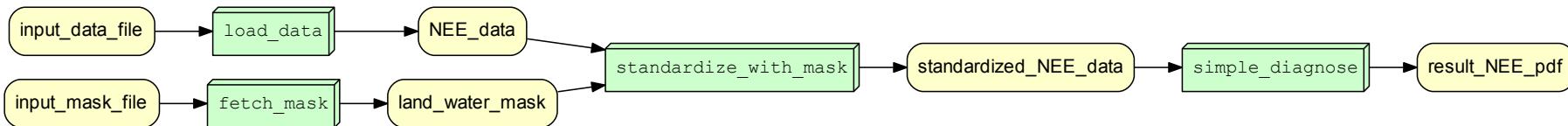
Get 3 views for the price of 1!



Process view



Data view



Combined view

Paleoclimate Reconstruction (EnviRecon.org)

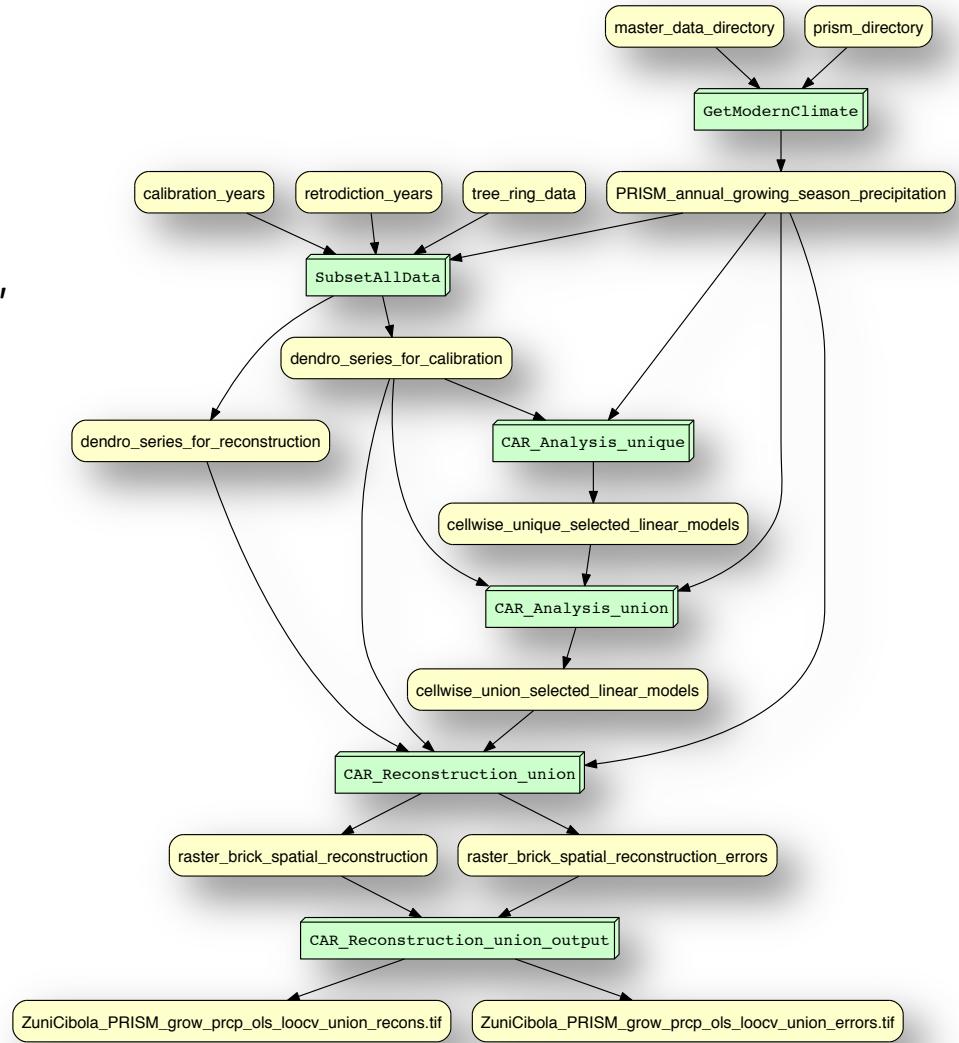
- ... explained using YesWorkflow!

Kyle B., (computational) archaeologist:

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

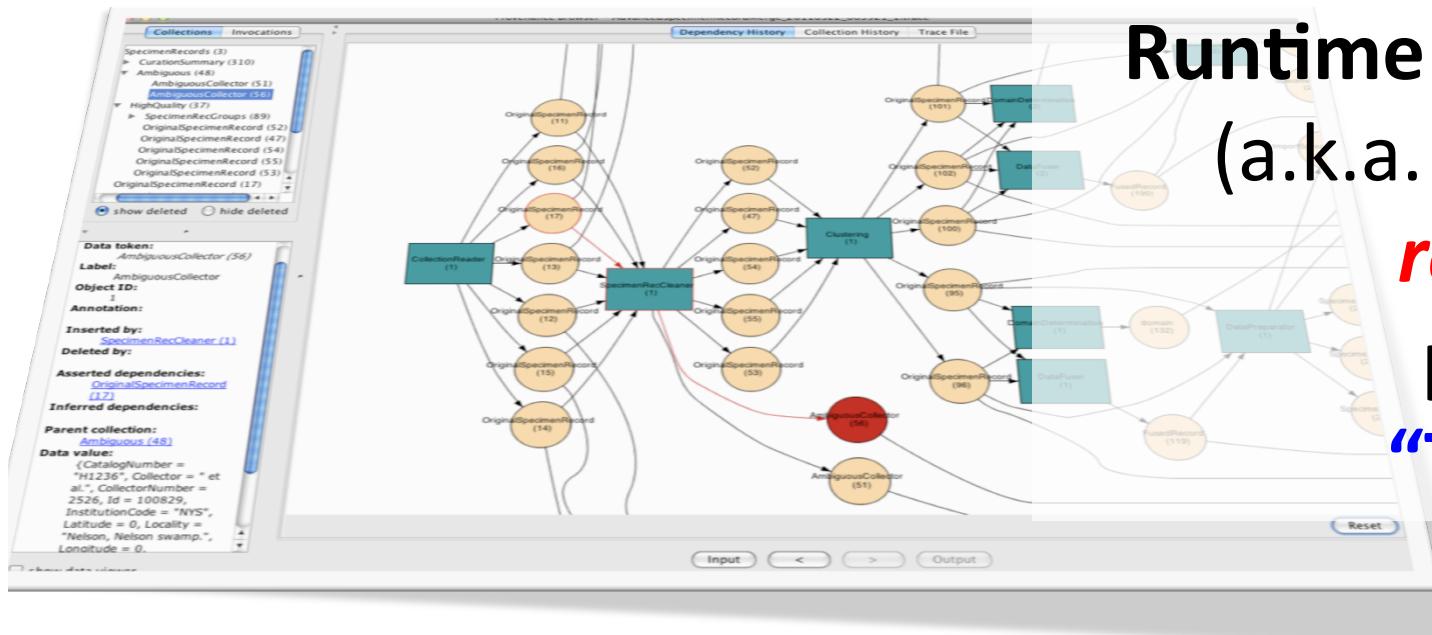
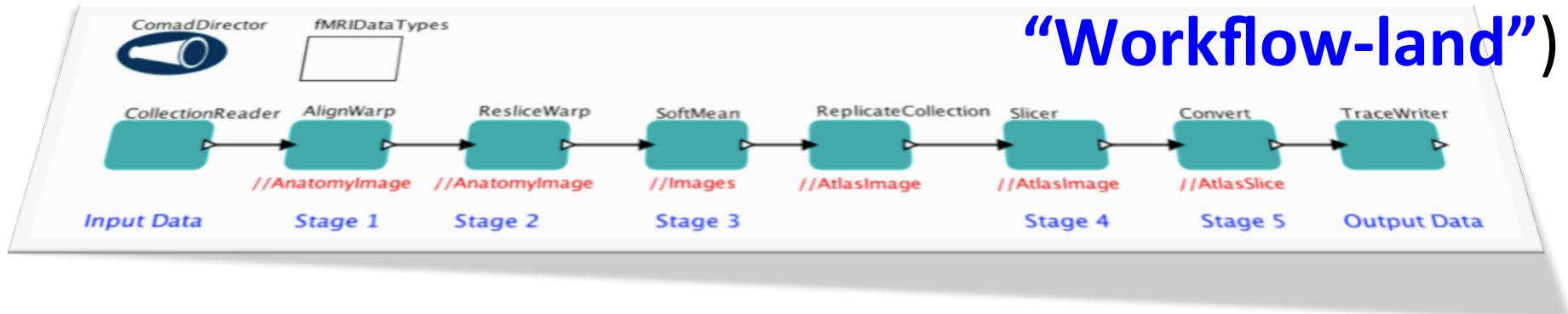


SKOPE + **Kurator**
+ **DataONE**
Data Observation Network for Earth
=> **YesWorkflow.org**



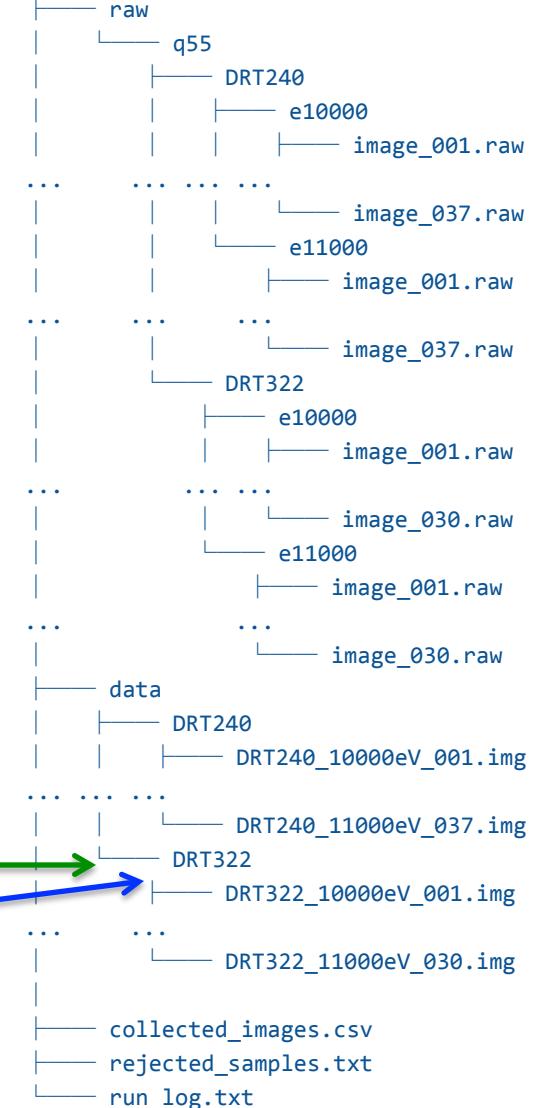
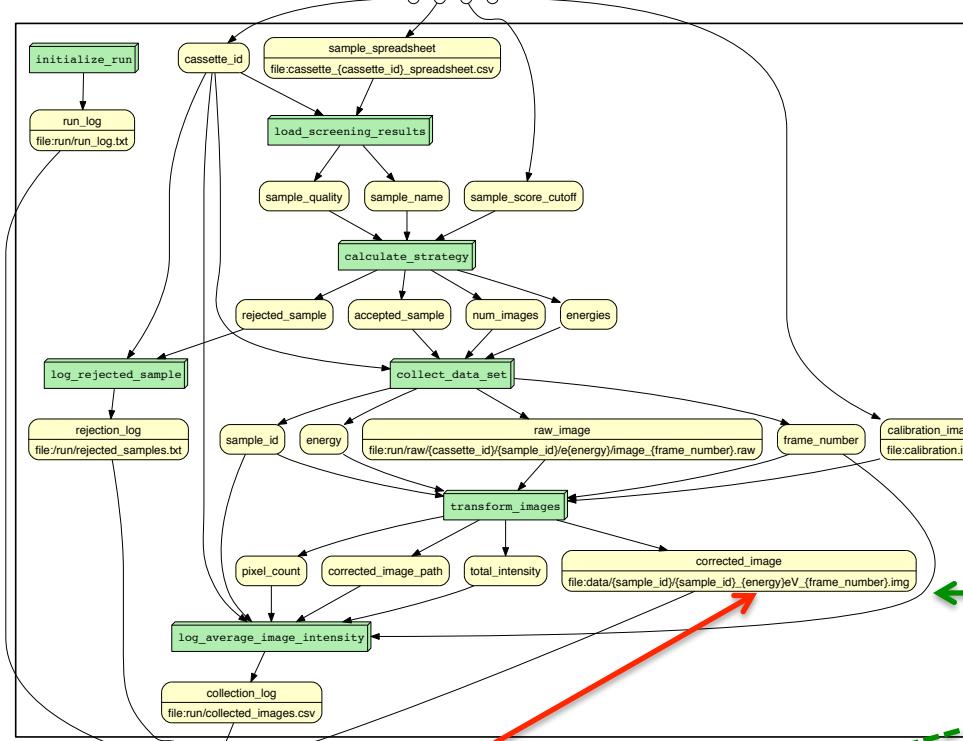
Provenance Lands

Workflow Modeling & Design (a.k.a. *prospective* provenance “Workflow-land”)



Runtime Provenance (a.k.a. traces, logs, *retrospective* provenance, “Trace-land”)

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

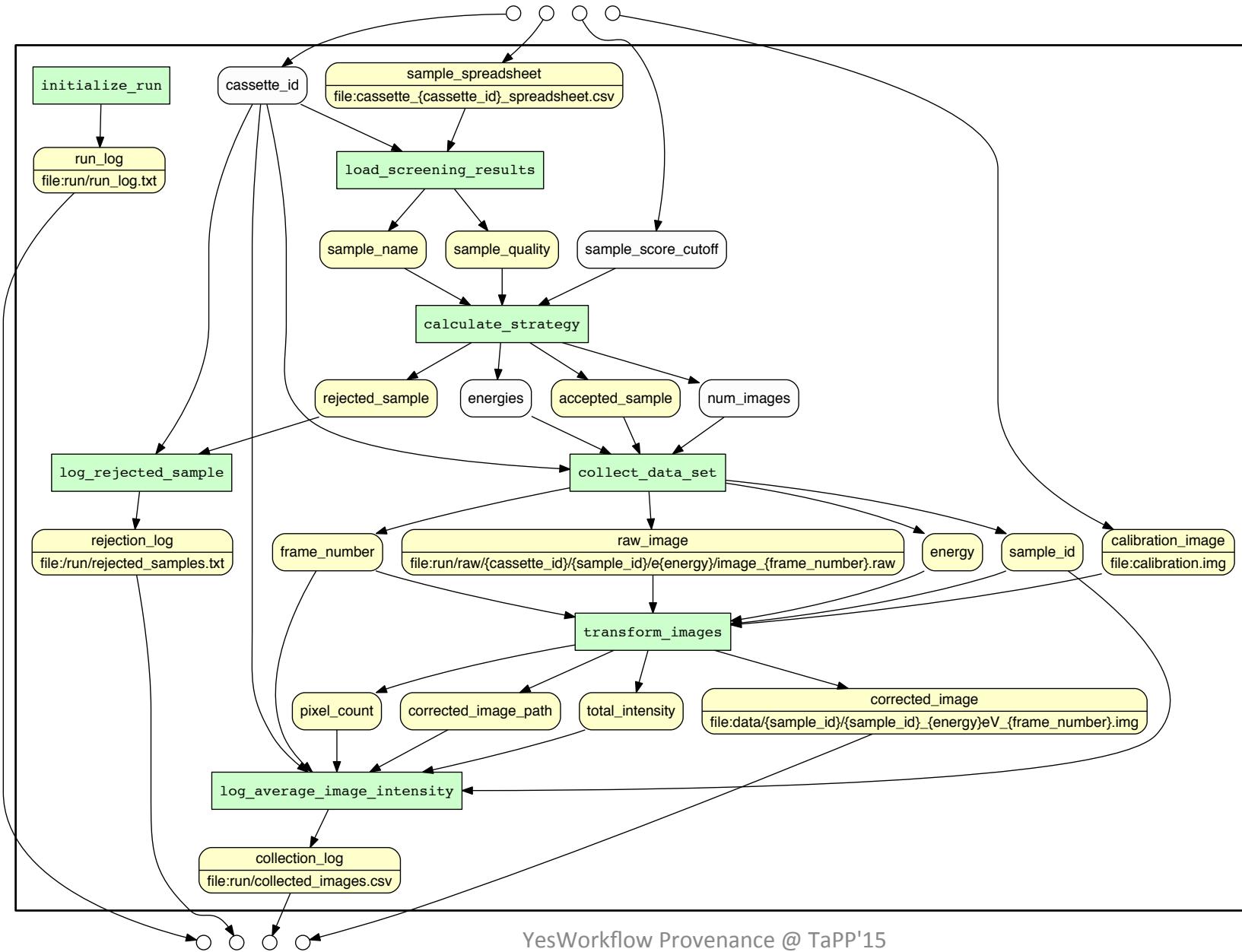


- **URI-templates link** conceptual entities to **runtime provenance** “left behind” by the script author ...
- ... facilitating provenance reconstruction

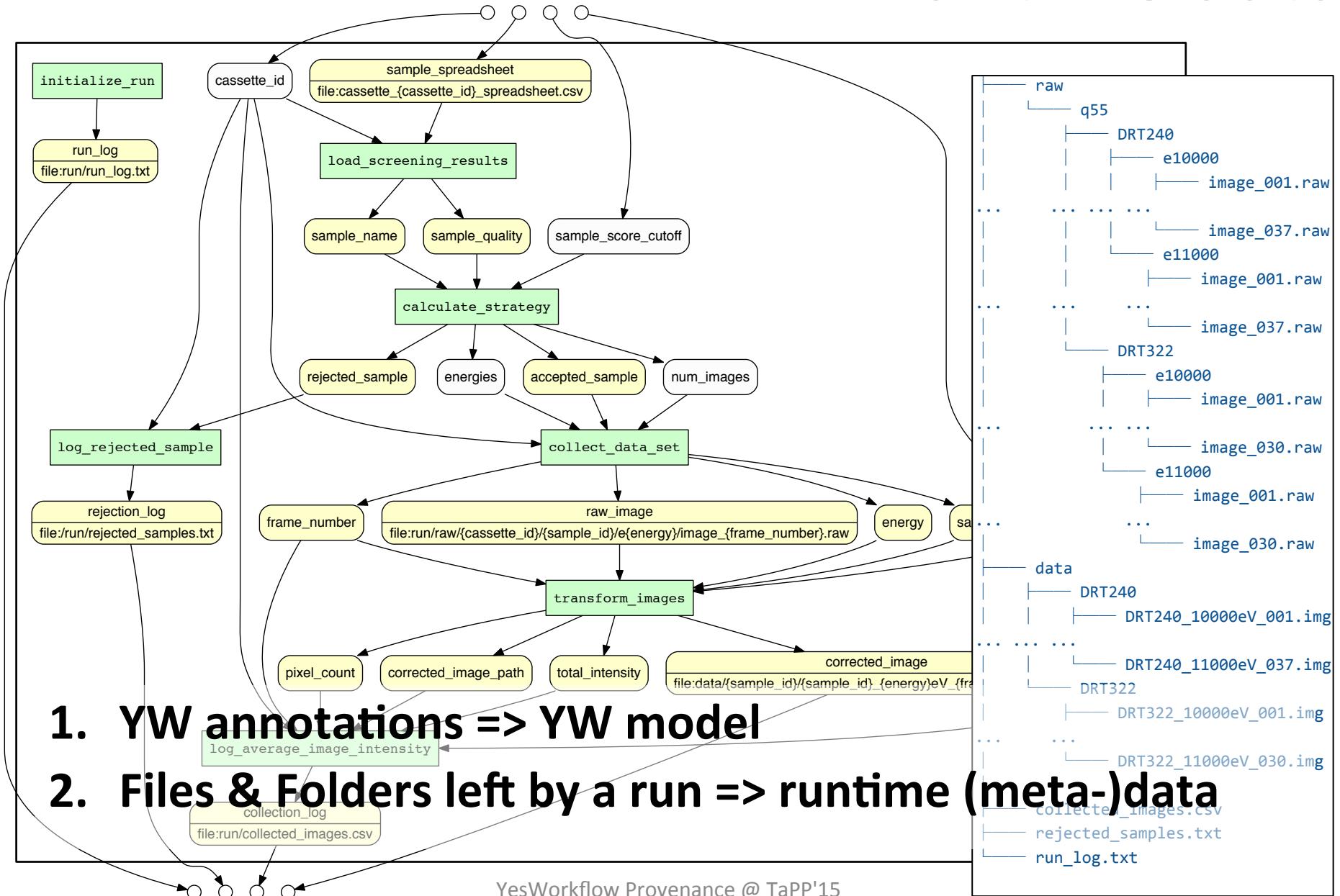
YW (*prospective*) and YW-Recon (*retrospective*) Provenance

- **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 ⇔ Persisted Data**
 - Run URI-template queries
 - cf. “ls -R” & RegEx matching
- **4. YW-Query**
 - Answer the user’s provenance queries

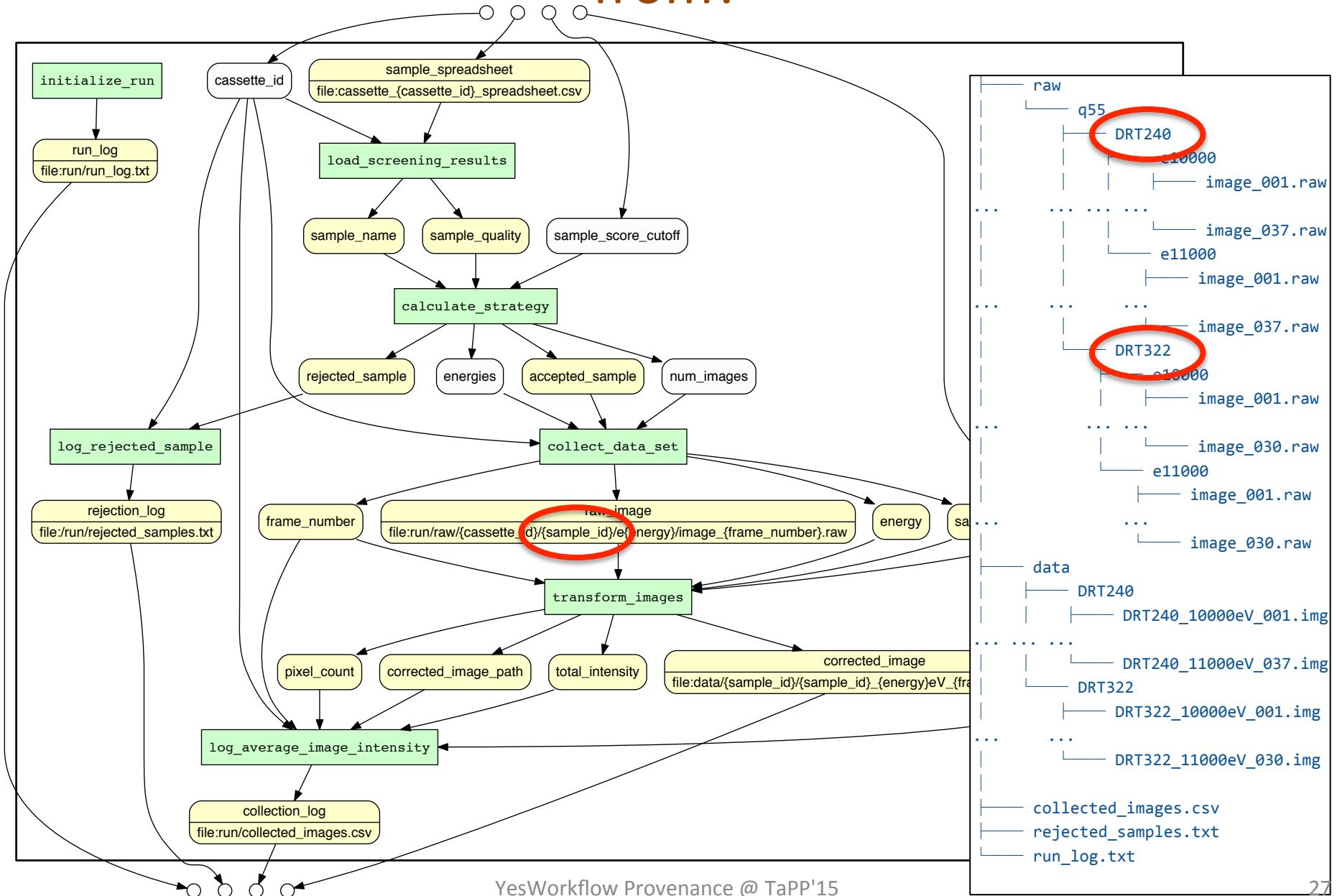
Data collection workflow (X-ray diffraction)



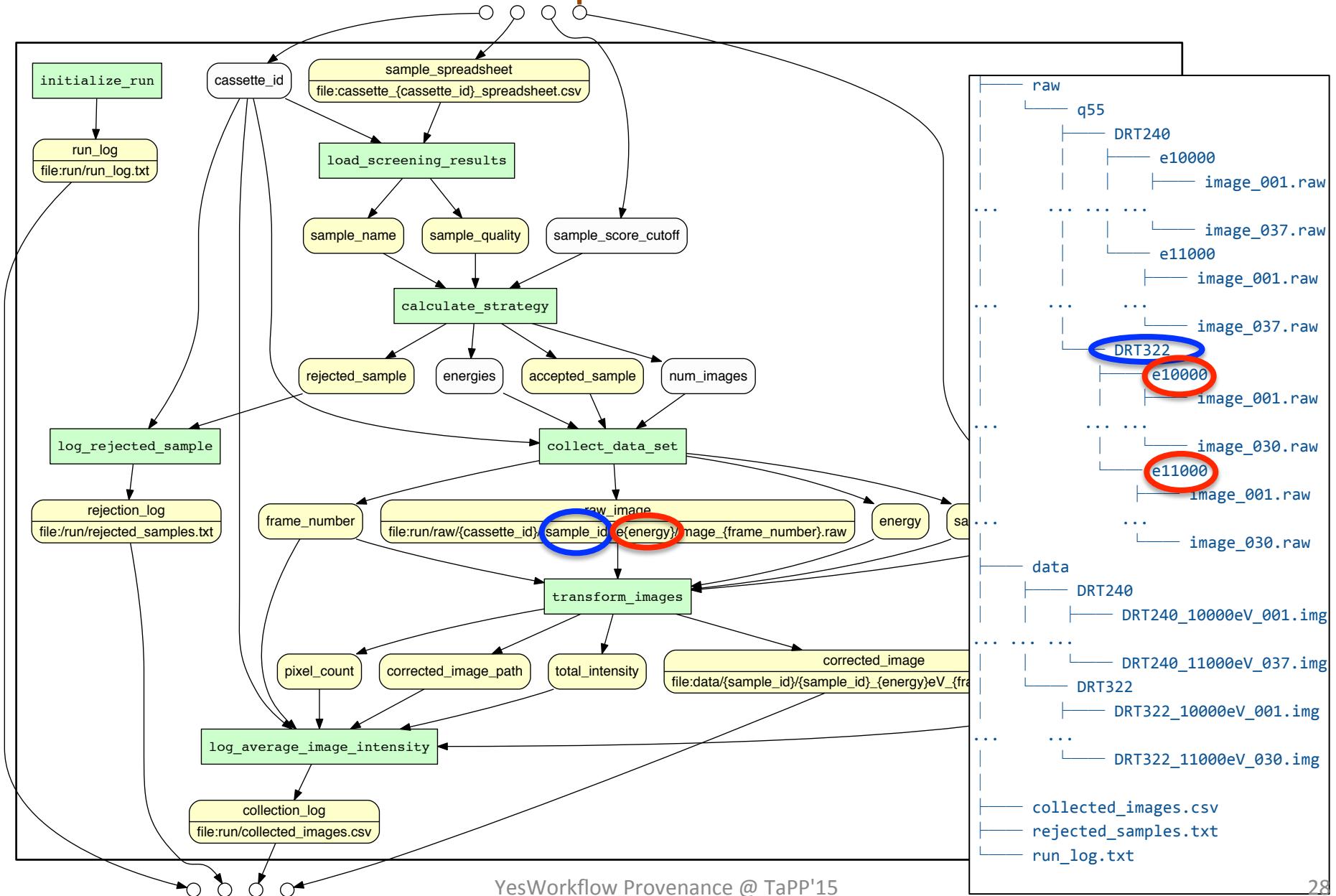
Data collection workflow: runtime data



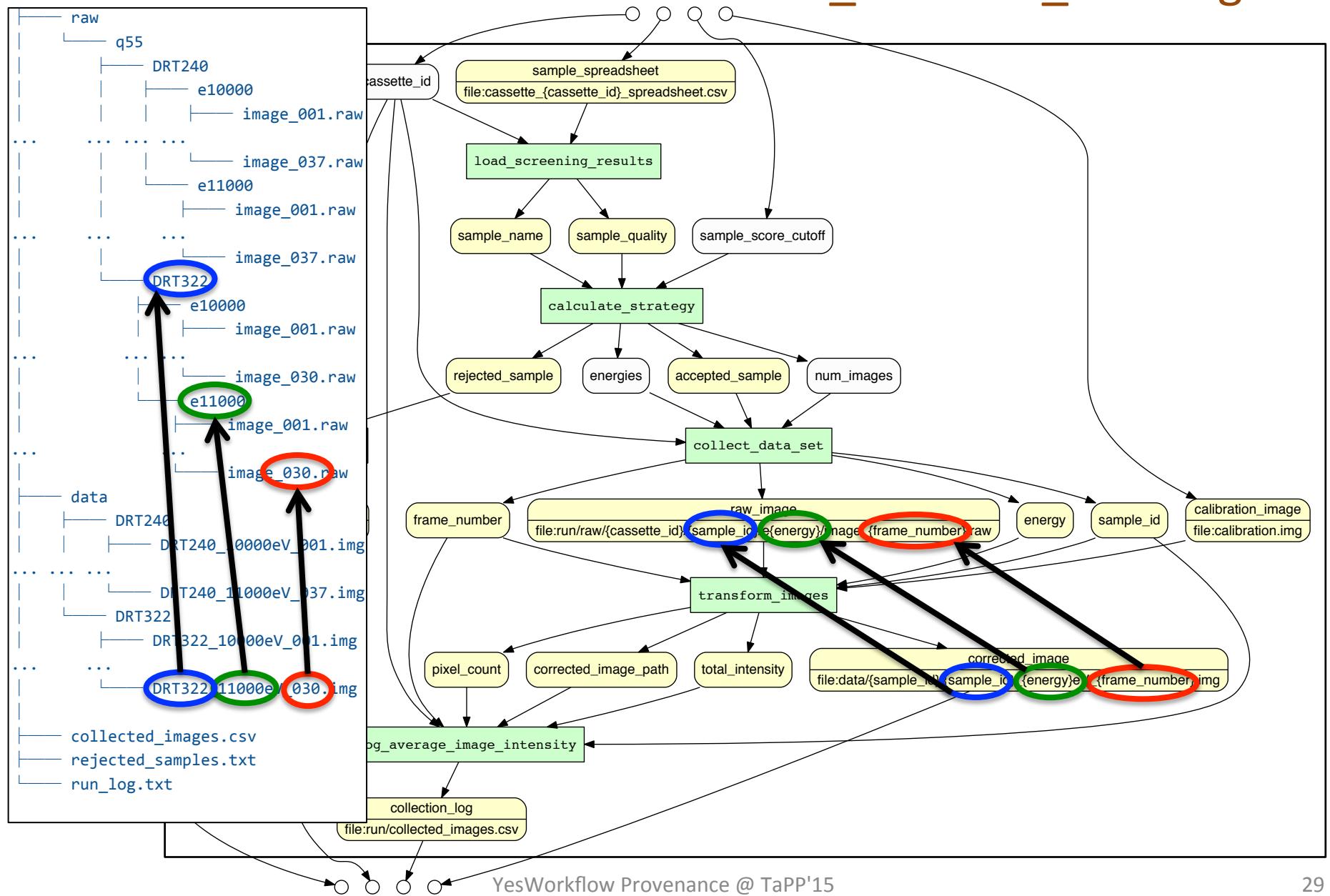
Q₁: What **samples** did the script run collect images from?



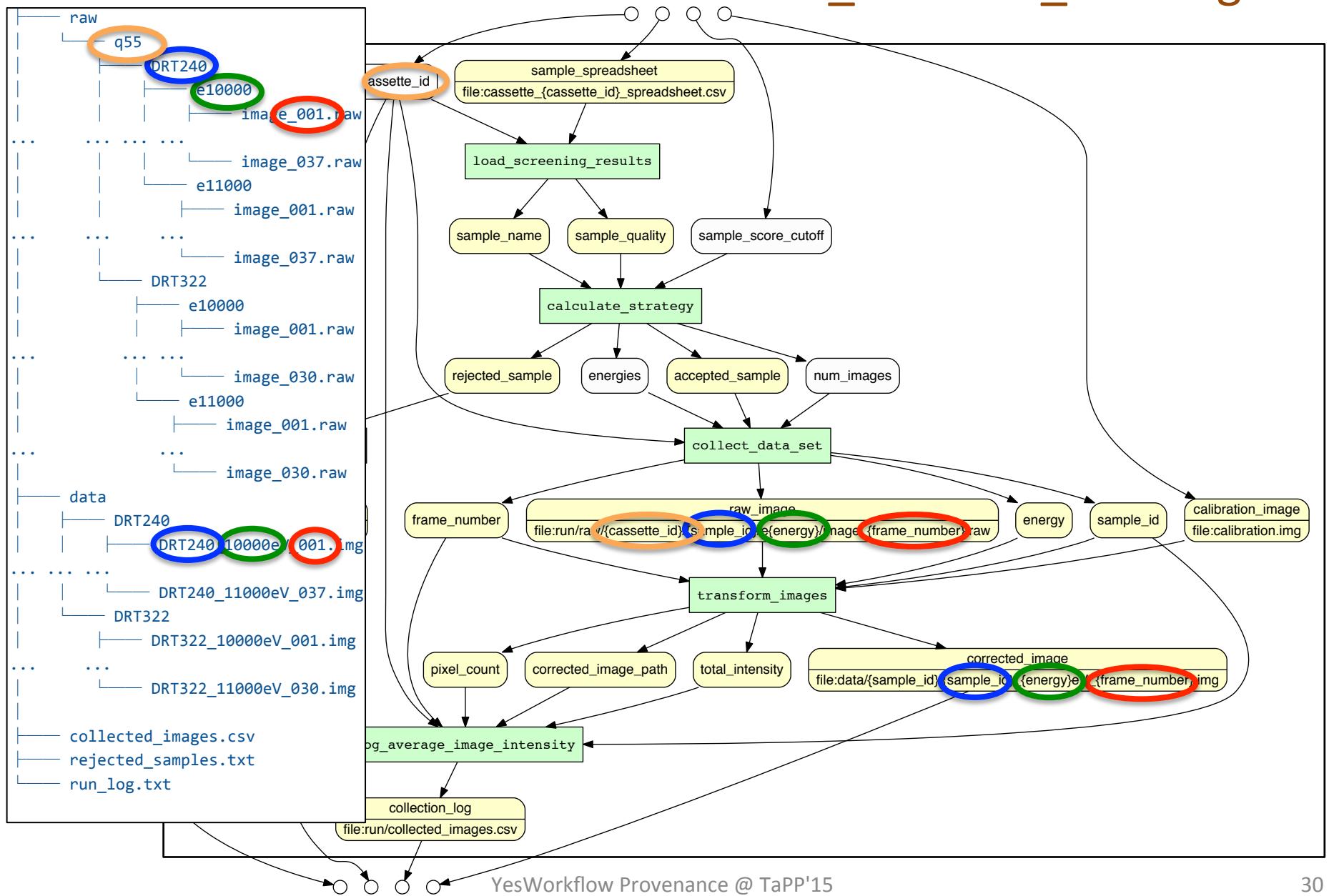
Q₂: What **energies** were used for image collection from sample DRT322?



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



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



Querying Provenance

(Q₂) *What energies were used during collection of images from sample DRT322?* The scientist's solution is to look at the contents of the run/raw/q55/DRT322 directory. The names of subdirectories indicate the values of the energies.

General solution using YW: Assume that 'what energies' and 'from sample DRT322' mean 'what values of energy as seen by collect_data_set when sample_id equals DRT322 as seen by collect_data_set'. Then the solution is to look for all persisted outputs of collect_data_set that include both energy and sample_id in the expanded URI template for the output. Extract the value of energy from each such path for which sample_id equals DRT322 and return the set of unique values. Q₂ can be expressed as the Datalog query:

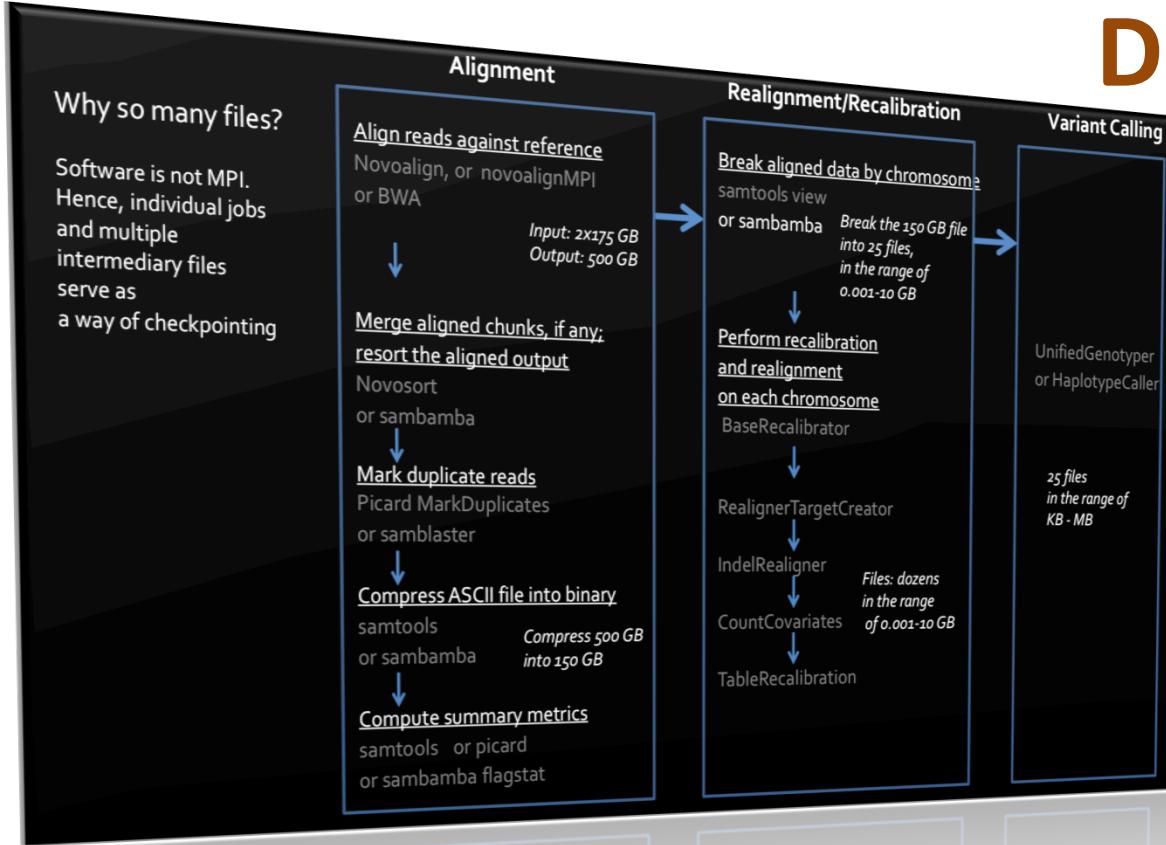
```
energies_used(EnergyValue) :-  
    resource_metadata(collect_data_set, raw_image, ResourceId,  
                      sample_id, "DRT322"),  
    resource_metadata(collect_data_set, raw_image, ResourceId,  
                      energy, EnergyValue).
```

For our example, the above query will yield the answers 10000 and 11000.

Base Relation	Description
program(<i>id, name, begin_annot_id, end_annot_id</i>)	Identifier, name, and annotation identifiers for begin and end of workflow program blocks
port(<i>id, type, name, annot_id</i>)	Identifier, type (<i>in, out, or param</i>), name, and annotation identifier of program ports
port_alias(<i>port_id, alias_name</i>)	Alias names optionally given to ports (specified via @AS annotations)
has_in_port(<i>program_id, port_id</i>)	Input ports of program blocks
has_out_port(<i>program_id, port_id</i>)	Output ports of program blocks
channel(<i>id, binding</i>)	Assignment of port names (script variables) or aliases to channels
port_connects_to_channel(<i>port_id, channel_id</i>)	Assignment of ports to channels
port_uri(<i>port_id, uri_template</i>)	URI template optionally assigned to port
uri_variable(<i>id, name, port_id</i>)	Identifier, name, and associated port of URI metadata variables
resource(<i>id, uri</i>)	Identifier and expanded URI (file path) of resources created by a run of the script
resource_channel(<i>resource_id, channel_id</i>)	Resources that were input to or output by a channel during a run of the script

Table 1 Base relations generated by the YawWorkflow provenance reconstruction process

Data & Workflow Management



1. Large **total data footprint**
2. Large **number of files**
3. Large number of **simultaneous but independent non-MPI computations**
4. **Keeping track** of what was done to the data: large amount of **metadata**
5. Workflow bottlenecks: fans & merges; more fans

Source: L Mainzer, V Jongeneel (IGB & NCSA)

Taking YW for a spin ...

- “To document **on-the fly**, specifically for a given workflow configuration invoked:
 - **do not insert annotations into code,**
 - **but rather have code print annotations into a special log during execution,**
 - **then parse that log!**”
- Liudmila Mainzer

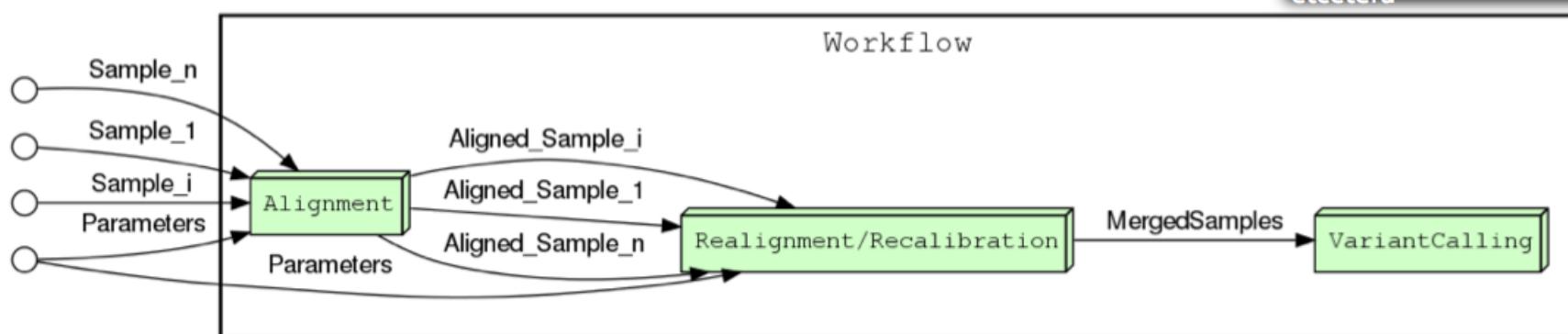
```
#!/bin/bash

# @begin Workflow
# @in x @as InputFileFolder
# @in y @as Parameters
# @in x1 @as Sample_1
# @in x1 @as Sample_i
# @in x1 @as Sample_n

# get input x somehow
# get input y somehow

# @begin Alignment
# @in y @as Parameters
# @in x1 @as Sample_1
# @in xi @as Sample_i
# @in xn @as Sample_n

etcetera
```



Conclusions & the Road ahead ...

- **YW: Go where the users are!**
 - ... they already capture provenance through metadata!
- **Beware your level of provenance abstraction**
 - Let the user provide a workflow model easily!
- **YW-Recon:**
 - ... finishing support for **retrospective provenance** *without* using a runtime provenance recorder!
 - Key insight: scientists already leave provenance “bread crumbs” behind! (it’s not an accident!)
 - Exploit this **annotations**: **URI**-templates
 - Extend YW to exploit log files
- **YW-GUI:**
 - Exploring & experimenting with a UI
 - Let script **author** immediately see effect of YW-annotations
 - Let scrip (**re-**)**user** explore the code and data provenance
- **YW-Query:**
 - ? Datalog, RPQ, ... for querying provenance ...