

Provenance as a requirement for large-scale complex astronomical instruments

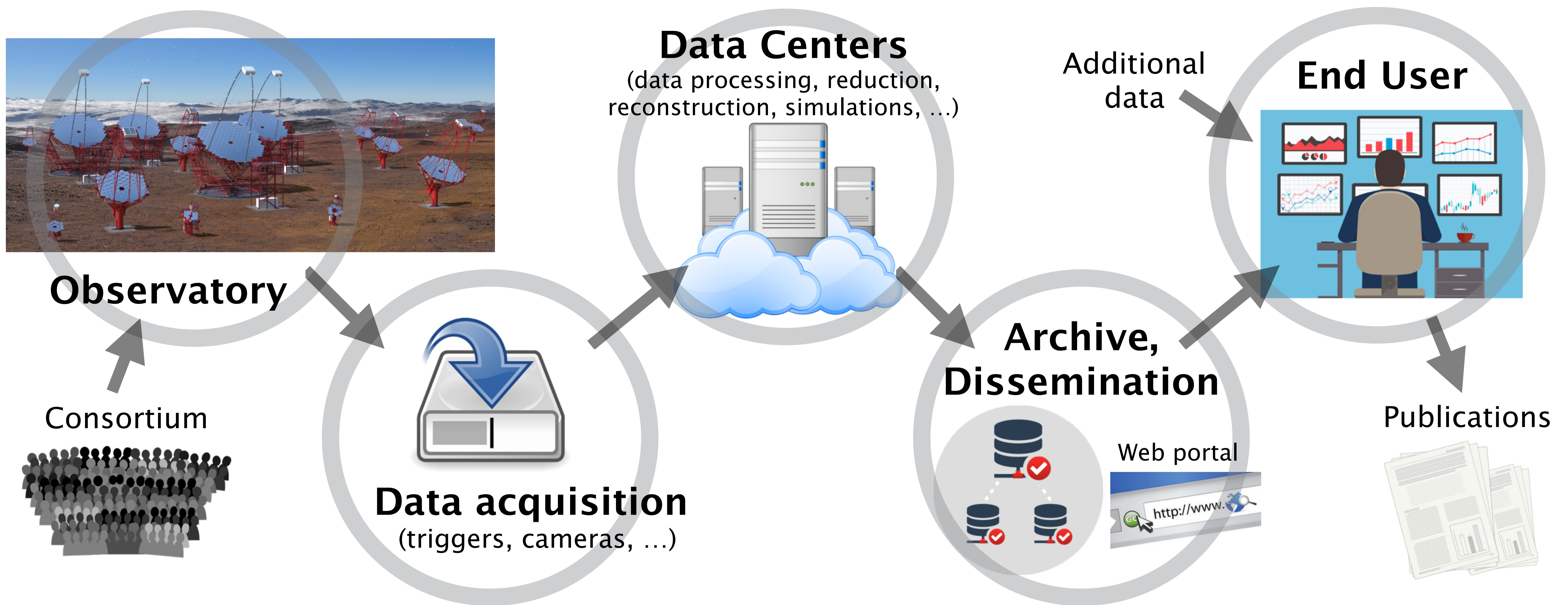


Mathieu Servillat,¹ Catherine Boisson,¹ Julien Lefaucheur,¹ Karl Kosack,² Michèle Sanguillon,³ Mireille Louys,^{4,5} François Bonnarel⁴

¹ LUTH, ² CEA Saclay, ³ LUPM, ⁴ CDS, ⁵ ICube

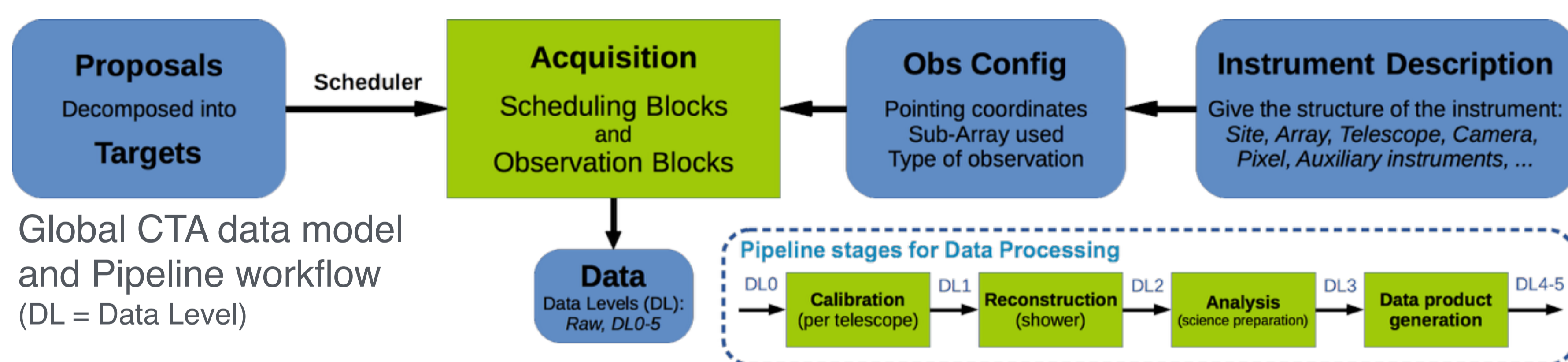
Context: State of the art observations are now performed by large-scale complex astronomical instruments. A **consortium** of specialists is generally responsible for the development and the operation of large **observatories**, as it is the case for example for the Cherenkov Telescope Array (CTA). The path of the data production from **acquisition** to **dissemination**, through e.g. **data centers**, **archives** and **web portals**, can be extremely obscure to the **end user**.

Provenance: to assess the **usefulness** and the **quality** of the data for their own scientific work, end users need a flowchart explaining the large number of steps and complexity involved in the data preparation. This can be done by collecting **provenance information** at each step of the data preparation. We followed the **IVOA Provenance data model** (see Poster 129) to develop solutions for CTA.



1/ How to collect provenance information during CTA data production?

- Include the relevant **metadata** in a structured CTA data model

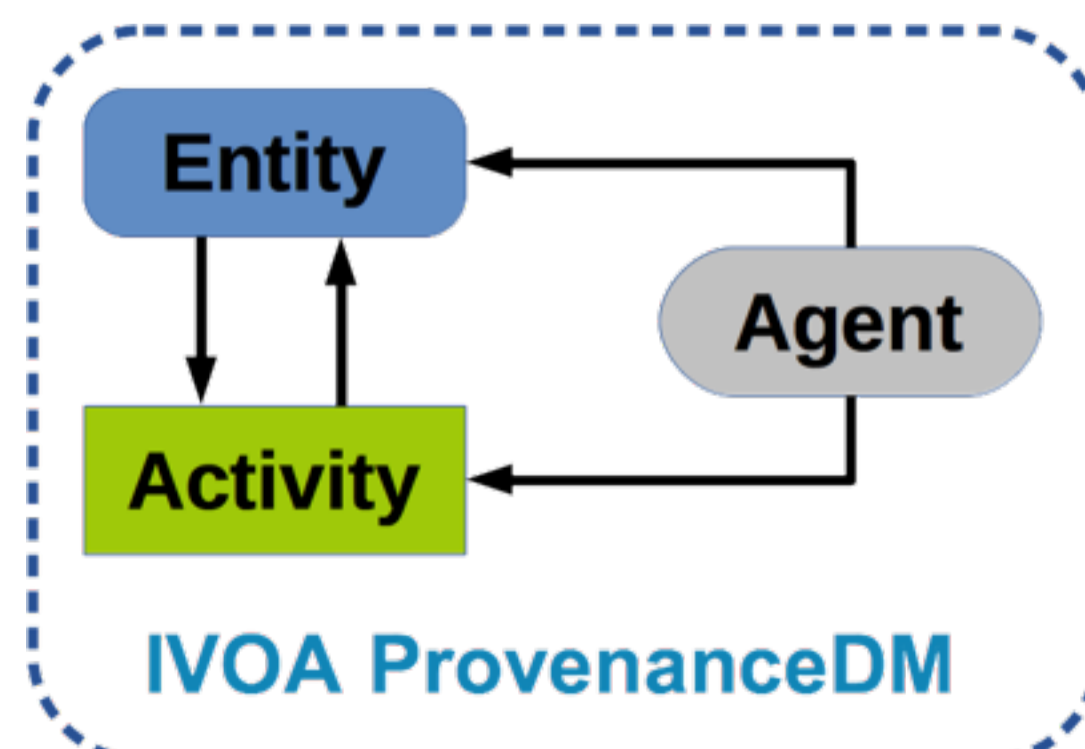


- Follow the **IVOA Provenance data model** for the generated **Data**

- **Collect** provenance information at each step of the data processing:

- ➔ Use **unique identifiers** for entities, activities and agents
- ➔ **Describe** each activity
- ➔ Keep a list of all **used** and **generated** entities during the execution of an activity

- A **Provenance Python class** has been developed for the CTA Pipeline framework **ctape**



```
from ctape.core import Provenance

prov = Provenance()
# prov a singleton, so this gives you the same provenance class

prov.start_activity("some_activity")

... # do things
prov.add_input_file("test.txt")
prov.add_output_file("out.txt")

prov.start_activity("some_sub_activity")

# do more things
prov.add_output_file("out2.txt")

prov.finish_activity() # finish some_activity
prov.finish_activity() # finish some_sub_activity
```

<https://github.com/cta-observatory/ctape>

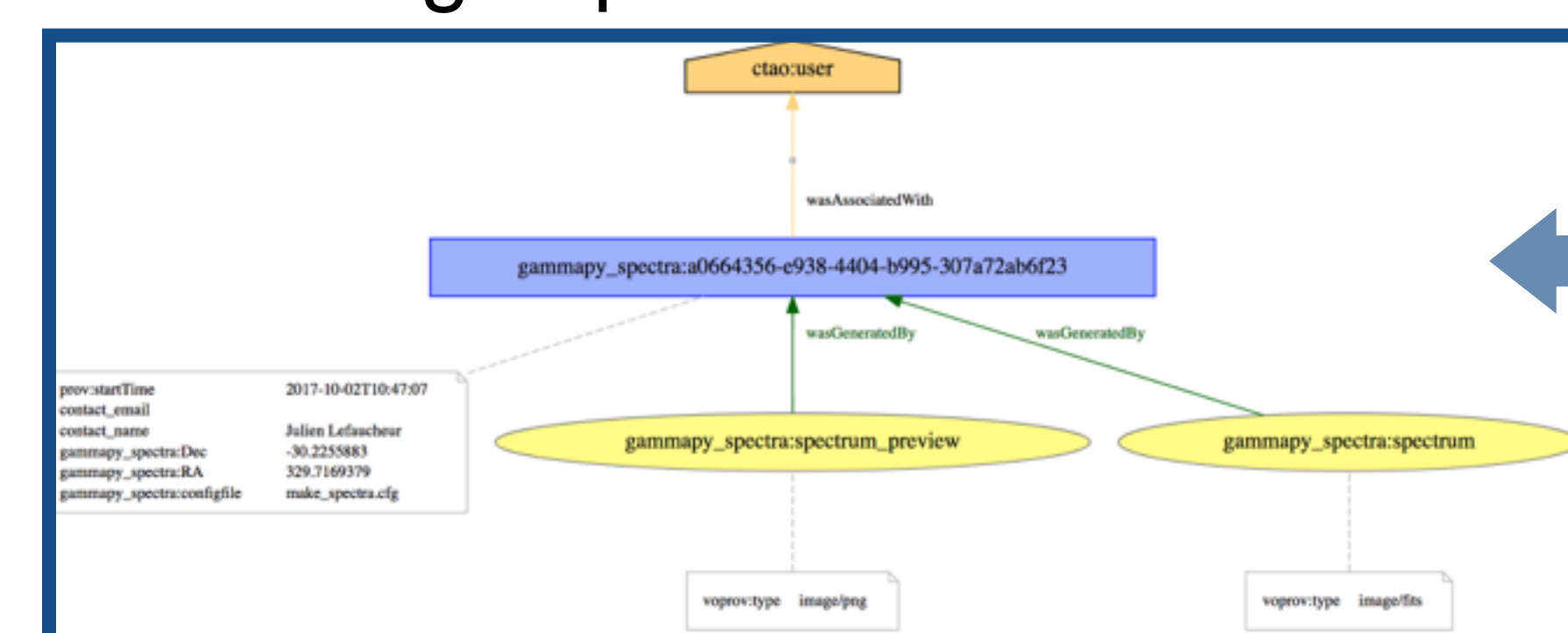
2/ How to store and expose the provenance information in a standard format?

We developed a **job control system** that stores provenance information following the IVOA UWS pattern and Provenance data model. The following features have been implemented:

- ➔ **Edit and fill** Activity Descriptions
- ➔ Run jobs **asynchronously** on a **work cluster**
- ➔ **Generate** and return Provenance files after job completion

Type	Start Time	Destruction Time	Phase	Control
gammamy_spectra	2017-10-02 10:47:07	2017-11-01 10:47:05	COMPLETED	Results
gammamy_spectra		2017-11-01 10:47:03	PENDING	
gammamy_spectra	2017-09-29 15:07:52	2017-10-29 15:07:51	COMPLETED	
gammamy_spectra	2017-09-29 14:55:10	2017-10-29 14:55:09	ABORTED	
gammamy_spectra	2017-09-29 14:21:20	2017-10-29 14:21:19	COMPLETED	

Tracking of provenance information



<https://github.com/mservillat/OPUS>

