

# Capture of provenance information : example of logprov for gammapy

Mathieu Servillat & J. Enrique Ruiz



Laboratoire Univers et Théories



# Objectives and requirements

- Capture provenance transparently
  - For Python tools (in particular data analysis on a laptop)
  - Make provenance info compatible between different software (ctapipe, gammipy, ...)
  - minimum effort from the developers, reuse existing tools
- Development drivers
  - Use **logging** system
    - known by developers, already integrated (base Python package)
    - log **dictionaries** to structure the information
  - **Non-intrusive** addition of code
    - use **decorators** for classes and methods
  - Fill with useful information
    - **descriptions** in a `definition.yaml` file (activity description, parameters, usage, generation...)
    - log parameters, arguments...

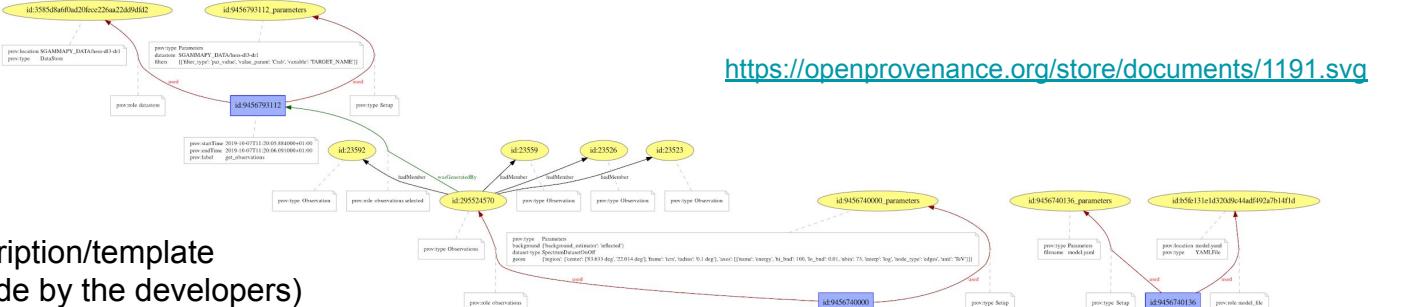
# Development within `gammapy`

- `logprov`, now an independent package:
  - <https://github.com/mservillat/logprov>
- `gammapy`: a community-developed, open-source Python package for gamma-ray astronomy:
  - <https://docs.gammapy.org/>
- **gammapy High Level Interface (HLI)**
  - Analysis class with methods
    - `get_observations()`
    - `set_model()`
    - `run_fit()`
    - ...
  - Fixes the <https://github.com/mservillat/logprov>: 1 gammapy HLI method = 1 activity
- Démo...

# Provenance in gammappy

1/ definition.yaml file for description/template  
(already **integrated** to the code by the developers)

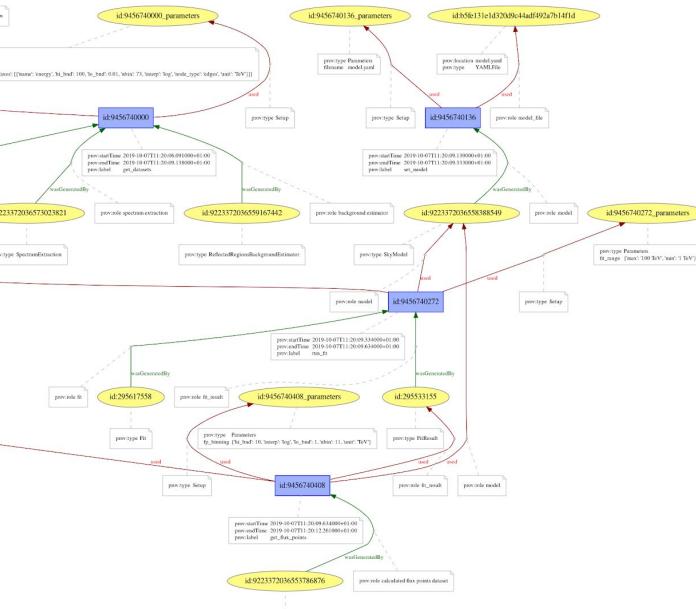
```
activities:
  get_observations:
    description:
      "Fetch observations from the data store according to criteria defined in the configuration"
    parameters:
      - name: datastore
        description: "DataStore path as string"
        value: settings.observations.datastore
      - name: filters
        description: "Filter criteria to select observations"
        value: settings.observations.filters
    usage:
      - role: datastore
        description: "DataStore object file"
        entityType: DataStore
        location: settings.observations.datastore
    generation:
      - role: observations selected
        description: "Observations selected"
        entityType: Observations
        value: observations
        has_members:
          entitytype: Observation
          list: observations.list
          id: obs_id
          namespace: ""
  get_datasets:
    description: "Produce reduced datasets"
    parameters:
```



2/ entries automatically stored in the log file

```
INFO:gammappy.utils.provenance.provenance:_PROV_2019-10-07T11:20:05
     .884436.PROV_{'activity_id': 9456793112, 'activity_name': 'get_observations', 'startTime':
'2019-10-07T11:20:05.884419'}
INFO:gammappy.utils.provenance.provenance:_PROV_2019-10-07T11:20:06
     .091102.PROV_{'activity_id': 9456793112, 'parameters': {'datastore':
'$GAMMAPY_DATA/hess-d13-dr1', 'filters': {'filter_type': 'par_value', 'value_param':
'Crab', 'variable': 'TARGET_NAME'}}}
INFO:gammappy.utils.provenance.provenance:_PROV_2019-10-07T11:20:06
     .091413.PROV_{'activity_id': 9456793112, 'used_role': 'datastore', 'used_id':
'35858a6f0ad20fec226aa22dd9df0d2', 'entity_type': 'DataStore', 'entity_location':
'$GAMMAPY_DATA/hess-d13-dr1'}
INFO:gammappy.utils.provenance.provenance:_PROV_2019-10-07T11:20:06
     .091474.PROV_{'activity_id': 9456793112, 'generated_role': 'observations selected',
'generated_id': '295524570', 'entity_type': 'Observations'}
INFO:gammappy.utils.provenance.provenance:_PROV_2019-10-07T11:20:06
     .091527.PROV_{'entity_id': '295524570', 'member_id': '23592', 'member_type': 'Observation'}
INFO:gammappy.utils.provenance.provenance:_PROV_2019-10-07T11:20:06
     .091571.PROV_{'entity_id': '295524570', 'member_id': '23523', 'member_type': 'Observation'}
INFO:gammappy.utils.provenance.provenance:_PROV_2019-10-07T11:20:06
     .091613.PROV_{'entity_id': '295524570', 'member_id': '23526', 'member_type': 'Observation'}
INFO:gammappy.utils.provenance.provenance:_PROV_2019-10-07T11:20:06
     .091653.PROV_{'entity_id': '295524570', 'member_id': '23559', 'member_type': 'Observation'}
INFO:gammappy.utils.provenance.provenance:_PROV_2019-10-07T11:20:06
     .091691.PROV_{'activity_id': 9456793112, 'endTime': '2019-10-07T11:20:06.091068'}
```

<https://openprovenance.org/store/documents/1191.svg>



3/ export to W3C PROV or search in provenance records

# Code: set up the capture

1/ YAML file following the IVOA  
Provenance DM description part

```
activities:  
  get_observations:  
    description:  
      "Fetch observations from the data store according  
      to criteria defined in the configuration"  
    parameters:  
      - name: datastore  
        description: "DataStore path as string"  
        value: settings.observations.datastore  
      - name: filters  
        description: "Filter criteria to select observations"  
        value: settings.observations.filters  
    usage:  
      - role: datastore  
        description: "DataStore object file"  
        entityType: DataStore  
        location: settings.observations.datastore  
    generation:  
      - role: observations selected  
        description: "Observations selected"
```

2/ class decorator (or function decorator @trace)

```
prov_capture = ProvCapture(definitions=definition, config=provconfig)
```

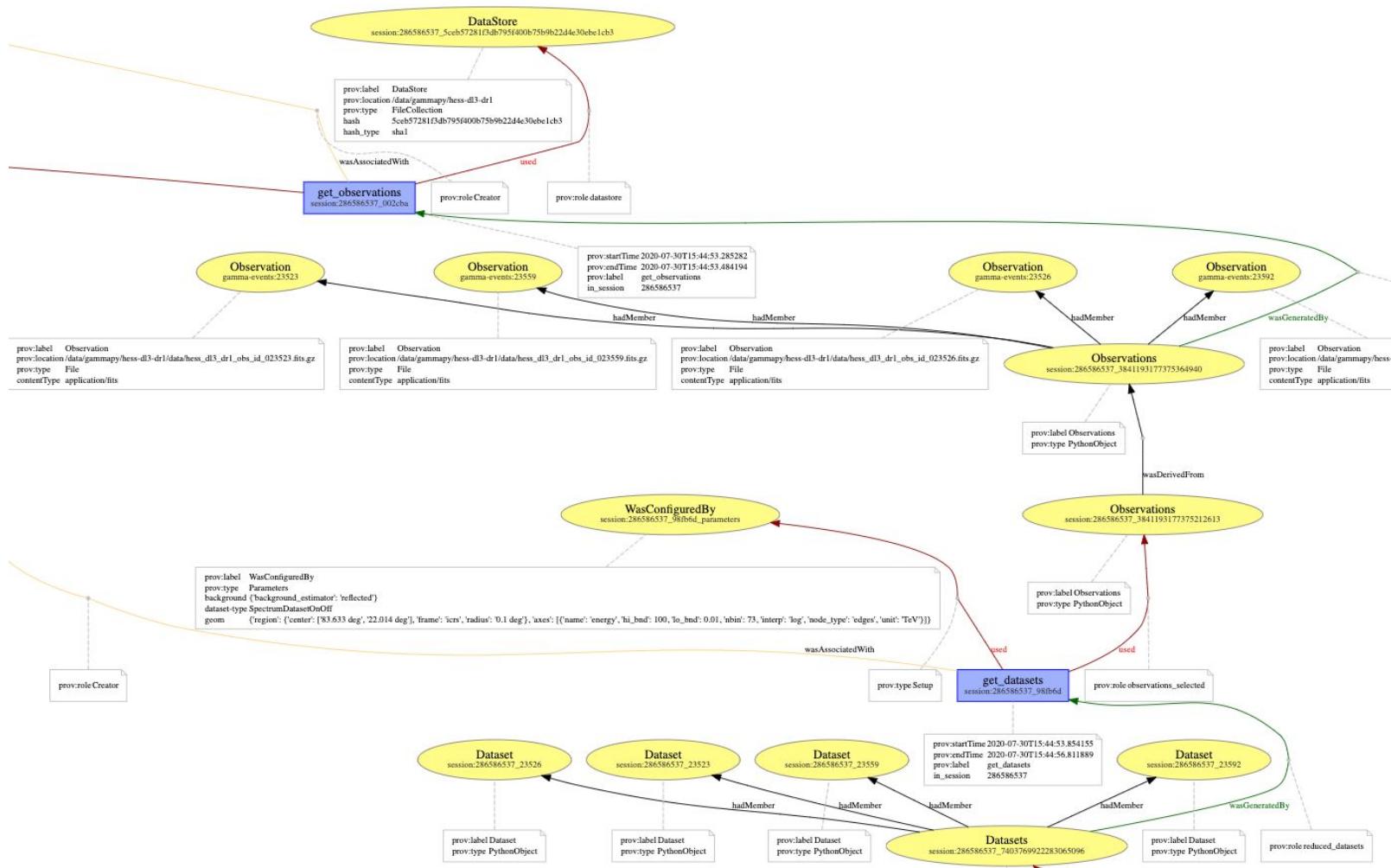
```
@prov_capture.trace_methods
```

```
class Analysis:
```

```
    """Config-driven high-level analysis interface following gammipy structure
```

# Code: what the user runs

```
# config = AnalysisConfig.from_template("3d")
config = AnalysisConfig.from_template("1d")
config.settings["general"]["logging"]["level"] = "INFO"
os.environ['GAMMAPY_DATA'] = '/data/gammapy'
# config.settings["general"]["logging"]["filename"] = logname
# config.settings["general"]["logging"]["filemode"] = 'w'
start = datetime.datetime.now().isoformat()
analysis = Analysis(config)
analysis.get_observations()
analysis.observations = copy.deepcopy(analysis.observations)    ← this transformation is not captured
analysis.get_datasets()
with open("/data/gammapy/model.yaml") as f:
    model = yaml.safe_load(f)
analysis.set_model(model=model)
# analysis.set_model(filename="/data/gammapy/model.yaml")
# analysis.datasets["stacked"].background_model.tilt.frozen = False
analysis.run_fit()
# analysis.datasets = copy.deepcopy(analysis.datasets)    ← model set with a filename or a dictionary
analysis.get_flux_points()
analysis.flux_points.write("flux_points.fits")
```



# Implementation details

- Identifiers and unicity...
  - the identifier is recomputed from the entity each time it is used
  - id(), hash(), str(), file hash (md5, sha1...)
- Parameter values directly attached to the activity
- Check for entity modifications before usage
  - indicate a derivation if entity has changed (but no more information about what happened)
- session identifier
  - when the Analysis class is instanciated
  - place to keep the system information, the global configuration, ...
- limits
  - only wrapped methods are traced
  - need definition for usage/generation

# OPUS + gammapy

- 1 OPUS job
  - Runs several gammapy functions
  - Stores result entities
- Internal provenance
  - Store objects
  - **Link to OPUS job**
    - Sub-activities ?
    - W3C PROV Bundle ?
  - **link to result**
    - Stored in OPUS archive
    - Derivation ?
    - Copy ?

