

Provenance tracking and inspection in Istchain and Istosa

Jose Enrique Ruiz
LST General Meeting – 19/11/2021

What is provenance ?

CTA General Requirement

A-USER-0110 The CTA Observatory must ensure that data processing is traceable and reproducible

Specific Requirements

- Keep the traceability of the data products
- Perform data quality and reliability checks
- Reproduce and reprocess the data
- Debug the pipelines

- Provide provenance information to the final user
- Possibility to filter out results based on provenance analysis



Terminology

Configuration provenance

- software dependencies and execution environment

Data provenance

- usually gives some context and info on last processing step on the dataset
- mostly recorded in FITS HISTORY headers
- also metadata describing the dataset

Execution provenance

- commands and parameters/values with timestamps
- mostly recorded in log files

Full provenance

- provides a graph view of the whole process
- chain of commands datasets parameters/values

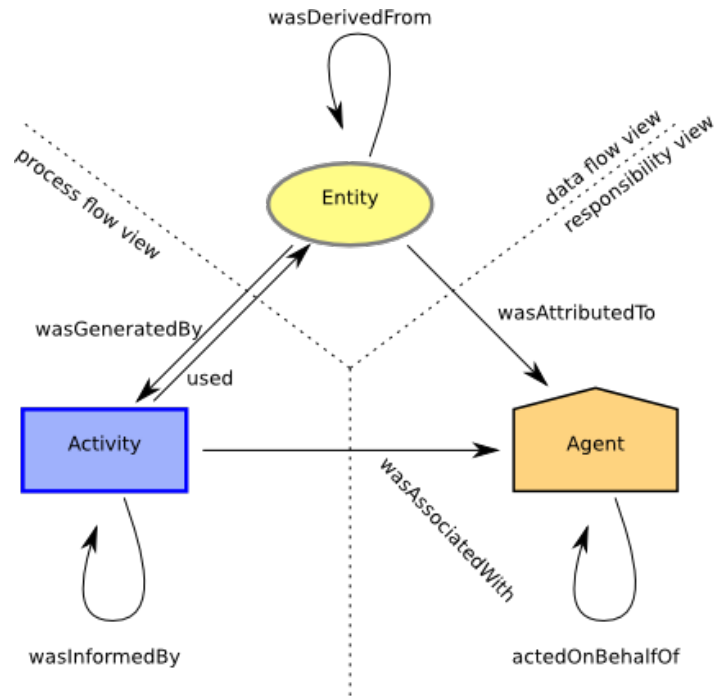
Standards in Astronomy

International Virtual Observatory Alliance

IVOA Documents

IVOA Provenance Data Model
Version 1.0

IVOA Recommendation 11 April 2020



On-going efforts within IVOA to develop standards for astronomical archives to expose and access provenance information

ProvSAP

Provenance Simple Access Protocol

ProvTAP

Provenance Tabular Access Protocol

voprov

Python package to work with IVOA Provenance Model objects

logprov

Python package to capture provenance following IVOA Provenance Model

work-in-progress

<https://www.ivoa.net/documents/ProvenanceDM>

Capturing provenance



On top

- data products collection have been already produced
- extract only information present in data and save it in a structured way elsewhere

Inside

- capture and save provenance at the moment of data processing

Granularity and level of details

- which steps, processes, datasets, ancillary files and parameters
- execution environment and dependencies

Format

- model and syntax that accounts for describing the pipeline components and relationships
- serializing provenance info as attached metadata, structured logs, graphs, relational DB,..

Software considered



lstosa

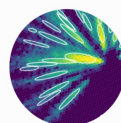
Python **pipeline** for the **On-site Analysis of low-level** data observations from the LST1 curated and developed by GAE-UCM group.

<https://github.com/gae-ucm/lstosa>

lstchain

Python **library** for the **processing of low-level** data observations from the LST1 curated and developed by the LST Collaboration.

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



ctapipe

Python **framework** for prototyping the **low-level data processing** algorithms for the Cherenkov Telescope Array.

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



- software dependencies and execution environment

- commands and parameters/values with timestamps
- mostly recorded in log files

Provenance in Tools

- provides a very simple API to **code provenance capture features in ctapipe Tools**
- provides serialization in JSON format and as Python dictionaries
- captures execution environment and profiling by default
- allows capture of processes, input and output datasets
- compliant with IVOA Provenance model

Cons

- tracks provenance execution **within a single Tool**
- needs the use of ctapipe Tools

Up to now it **does not allow chaining** provenance products of different Tools but allows concatenation of a provenance records for the same Tool into a single log file

```
from ctapipe.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
```

Data provenance

- usually gives some context and info on last processing step on the dataset
- mostly recorded in FITS HISTORY headers
- also metadata describing the dataset

Store configuration execution metadata in HDF5 files #747

Merged [ropezcoto](#) merged 18 commits into [master](#) from [metadata](#) 16 days ago

Conversation 22

Commits 18

Checks 8

Files changed 10



Write

Preview

This PR stores the config execution of scripts as metadata attributes of the HDF5 files and tables produced in the R0->calibration->DL1->DL1AB->DL2 workflow.

The scripts involved in this workflow are:

- `/scripts/lstchain_data_create_time_calibration_file.py`
- `/tools/lstchain_create_calibration_file.py`
- `/scripts/lstchain_data_r0_to_dl1.py`
- `/scripts/lstchain_dllab.py`
- `/scripts/lstchain_dl1_to_dl2.py`
- `/scripts/lstchain_add_source_dependent_parameters.py`

The global metadata stored at the root level of the HDF5 files, as well as into each of the tables/containers produced, account for the versions of the software used in the execution of the script that produces de file i.e. `CTAPIPE_VERSION`, `LSTCHAIN_VERSION`, `CTAPIPE_IO_LST`

The configuration parameters used for each script/tool are stored as an attribute `config` of the tables/containers produced.

Istchain

Data provenance

- usually gives some context and info on last processing step on the dataset
- mostly recorded in FITS HISTORY headers
- also metadata describing the dataset

These metadata may be easily retrieved with `pytables` library and `vitables` software, as described in the python snippets and screenshots below.

```
with tables.open_file(filename) as file:  
    ...:     print(file.root._v_attrs["LSTCHAIN_VERSION"])
```

0.7.4.post16+git0deb44e

```
with tables.open_file(filename) as file:  
    ...:     config = yaml.safe_load(file.root.d1.event.telescope.parameters.LST_LSTCam.attrs["config"])  
    print(config["tailcut"])
```

```
{'picture_thresh': 6,  
  'boundary_thresh': 3,  
  'keep_isolated_pixels': False,  
  'min_number_picture_neighbors': 2,  
  'use_only_main_island': False,  
  'delta_time': 2}
```

In the case where, during one of the workflow steps, a specific table is directly retrieved from one previous existing file and copied into the new produced HDF5 file, then the software versions and the config execution metadata are kept untouched.

Istchain

Data provenance

- usually gives some context and info on last processing step on the dataset
- mostly recorded in FITS HISTORY headers
- also metadata describing the dataset

ViTables 3.0.0

Árbol de bases de datos

calibcharge.hdf5

Propiedades de la tabla

General Atributos de Sistema

Atributos de usuario: 21

Nombre	Valor	Tipo de dato
CONTACT	LST Consortium	string
CTAPIPE_IO_LST_VERSI...	0.11.1	string
CTAPIPE_VERSION	0.11.0	string
LSTCHAIN_VERSION	0.7.6.dev315+geb990341	string
SOURCE_FILENAMES	['LST-1.1.Run05979.000...	python
charge_mean_DESC	np array of pedestal ave...	string
charge_median_DESC	np array of the pedestal ...	string
charge_median_outliers...	Boolean np array of the ...	string
charge_std_DESC	np array of the pedestal ...	string
charge_std_outliers_DE...	Boolean np array of the ...	string
config	{'EventSource': {'input...	string
n_events_DESC	Number of events used f...	string

ndow_width": 12, "apply_integration_correction": false}, "attach_subarray": "None")

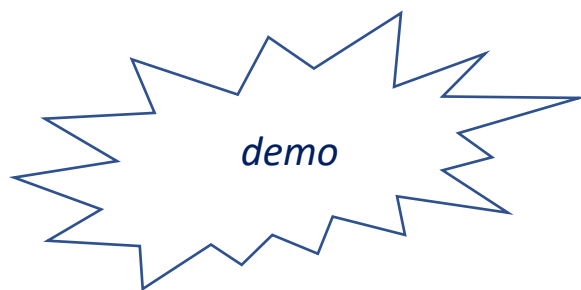
Añadir Borrar Qué es esto

Cancel OK

ViTables 3.0.0
Copyright (c) 2008-2017 Vicent Mas.
All rights reserved.
Muestra el diálogo de propiedades para el nodo seleccionado calibcharge.hdf5->/

```
CalibrationHDF5Writer:
  output_file: "calibcharge.hdf5"
  config_file: "/Users/jer/git/cta-observatory/cta-lstchain/lstchain/data/onsite_camera_calibration_param.json"
  one_event: true
  events_to_skip: 1000
  calibration_product: "LSTCalibrationCalculator"
  log_level: "DEBUG"
  log_config: {}
  EventSource:
    input_url: "/Users/jer/DATA/LST/R0/20210902/LST-1.1.Run05979.0000.fits.fz"
  LSTEventSource:
    LSTR0Corrections:
      drs4_pedestal_path: "/Users/jer/DATA/LST/calibration/20210902/v0.7.3/drs4_pedestal.Run05978.0000.fits"
      drs4_time_calibration_path: "/Users/jer/DATA/LST/calibration/20210902/v0.7.3/time_calibration.Run05979.0000.hdf5"
      select_gain: false
    EventTimeCalculator:
      run_summary_path: "/Users/jer/DATA/LST/monitoring/RunSummary/RunSummary_20210902.ecsv"
    allowed_tels:
      0: 1
  LSTCalibrationCalculator:
    squared_excess_noise_factor: 1.222
    flatfield_product: "FlasherFlatFieldCalculator"
    pedestal_product: "PedestalIntegrator"
  PedestalIntegrator:
    sample_duration: 100000
    tel_id: 1
    charge_median_cut_outliers:
      0: -10
      1: 10
    charge_std_cut_outliers:
      0: -10
      1: 10
    charge_product: "FixedWindowSum"
  FlasherFlatFieldCalculator:
    sample_duration: 100000
    tel_id: 1
    charge_product: "LocalPeakWindowSum"
    charge_median_cut_outliers:
      0: -0.5
      1: 0.5
    charge_std_cut_outliers:
      0: -10
      1: 10
    time_cut_outliers:
      0: 2
      1: 38
  LocalPeakWindowSum:
    window_shift: 5
    window_width: 12
    apply_integration_correction: false
  FixedWindowSum:
```





- provides a graph view of the whole process
- chain of commands datasets parameters/values

Provenance capture

structured logging

How?

- Using standard Python logging mechanism and a **provenance model defined in a YAML file**
- Non-intrusive implementation with **function/class decoration** in existing code
- Info within decorated functions can be extracted and mapped into W3C/IVOA Prov syntax
- Python **logging configuration** is set in an independent configuration file

Which info?

- Used and produced datasets, input parameters, config files and processes

What do we get?

- **Post-processed** log files as merged/filtered logs with info in W3C/IVOA Prov syntax
- W3C provenance **JSON** files and **PDF** graphs as final provenance products

- provides a graph view of the whole process
- chain of commands datasets parameters/values

Detailed considerations

- Captured **subrun-wise** and post-processed to **run-wise**
- Post-processing/merging yields provenance products at **different levels of granularity**
 - calibration
 - r0_to_dl1
 - dl1_to_dl2
 - all processes
- Execution **environment** is captured and stored as a session property
- Info *hidden* in **configuration files**
 - compared with hash-content algorithm and **copied** for reproducibility purposes
- Montecarlo simulated training datasets are not copied but **referenced**
- **Dry execution** mechanism allows provenance capture and merging avoiding data processing
- **Session properties**: observation date, run number, Istosa config file, Istchain version

run wise at different processing levels

```

activities:
  r0_to_dl1:
    description:
      "Create DL1 files for an observation run and subrun"
    parameters:
      - name: ObservationRun
        description: "Observation run number"
        value: ObservationRun
      - name: ObservationSubRun
        description: "Observation subrun number"
        value: ObservationSubRun
      - name: CalibrationRun
        description: "Calibration run number"
        value: CalibrationRun
      - name: PedestalRun
        description: "Pedestal run number"
        value: PedestalRun
      - name: ProdID
        description: "Production ID"
        value: ProdID
    usage:
      - role: "Observation subrun"
        description: "Observation subrun used"
        entityName: R0SubrunDataset
        value: R0SubrunDataset
        # filepath: /fefs/aswg/data/real/R0/20200218/LST1.1Run02006.0001.fits.fz
      - role: "Pedestal file"
        description: "Pedestal file used"
        entityName: PedestalFile
        value: PedestalFile
        # filepath: /fefs/aswg/data/real/calibration/20200218/v00/drs4_pedestal.Run02005.0000.fits
      - role: "Coefficients calibration file"
        description: "Coefficients calibration file"
        entityName: CoefficientsCalibrationFile
        value: CoefficientsCalibrationFile
        # filepath: /fefs/aswg/data/real/calibration/20200218/v00/calibration.Run02006.0000.hdf5
      - role: "Time calibration file"
        description: "Time calibration file"
        entityName: TimeCalibrationFile
        value: TimeCalibrationFile
        # filepath: /fefs/aswg/data/real/calibration/20191124/v00/time_calibration.Run1625.0000.hdf5
      - role: "Pointing file"
        description: "Pointing filename for DL1"
        entityName: PointingFile
        value: PointingFile
        # filepath: /fefs/home/lapp/DrivePositioning/drive_log_20_02_18.txt
        # description: "Pointing configuration file"

```

```

version: 1
formatters:
  simple:
    format: '%(levelname)s %(name)s %(message)s'
    #format: '%(asctime)s.%(msecs)03d%(message)s'
    datefmt: '%Y-%m-%dT%H:%M:%S'
handlers:
  provHandler:
    class: logging.handlers.WatchedFileHandler
    level: INFO
    formatter: simple
    filename: prov.log
loggers:
  provLogger:
    level: INFO
    handlers: [provHandler]
    propagate: False
disable_existing_loggers: False
PREFIX: __PROV__
HASH_METHOD: md5
HASH_BUFFER: path
capture: True

# Conda environment for provenance package
# conda env update -f environment.yaml

channels:
  - conda-forge

dependencies:
  - python
  - pyyaml
  - prov
  - pydot
  - pydotplus
  # dev dependencies
  - pytest
  - pytest-cov
  - black
  - isort

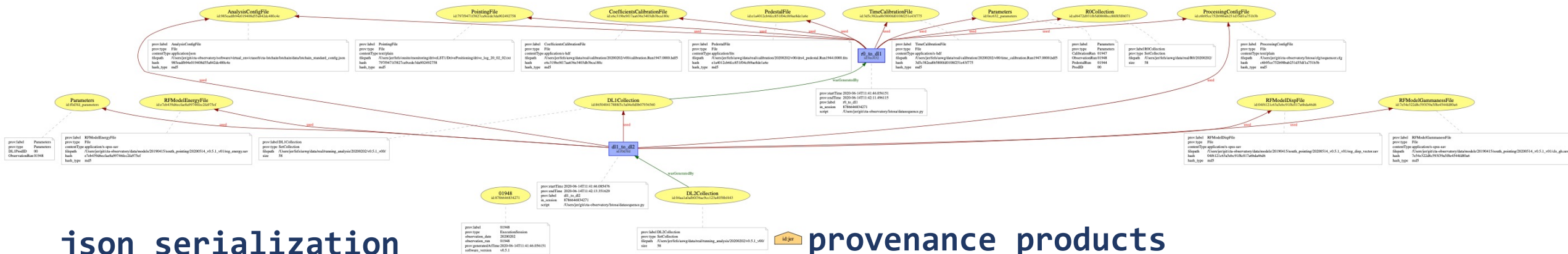
```

prov.log

```
INFO provLogger __PROV_2020-05-18T14:18:30.445713_PROV__{'session_id': 8739478486569, 'name': '01618', 'startTime': '2020-05-18T14:18:06.362321', 'system': {'executable': '/fefs/aswg/software/virtual_env/anaconda3/envs/osa/bin/python', 'platform': {'architecture_bits': '64bit', 'architecture_linkage': '', 'machine': 'x86_64', 'processor': 'x86_64', 'node': 'cp15', 'version': '#1 SMP Thu Nov 8 23:39:32 UTC 2018', 'system': 'Linux', 'release': '3.10.0-957.el7.x86_64', 'libcver': "(('glibc', '2.10')", 'num_cpus': 32, 'boot_time': '2020-03-24T03:48:43'}, 'python': {'version_string': '3.7.6 | packaged by conda-forge | (default, Mar 23 2020, 23:03:20) \n[GCC 7.3.0]', 'version': '3.7.6', 'compiler': 'GCC 7.3.0', 'implementation': 'CPython'}, 'environment': {'CONDA_DEFAULT_ENV': 'osa', 'CONDA_PREFIX': '/fefs/aswg/software/virtual_env/anaconda3/envs/osa', 'CONDA_PYTHON_EXE': '/fefs/aswg/software/virtual_env/anaconda3/bin/python', 'CONDA_EXE': '/fefs/aswg/software/virtual_env/anaconda3/bin/conda', 'CONDA_PROMPT_MODIFIER': '(osa)', 'CONDA_SHLVL': '2', 'PATH': '/local/home/lstanalyzer/usr/bin:/local/home/lstanalyzer/.local/bin:/fefs/aswg/software/virtual_env/anaconda3/envs/osa/bin:/fefs/aswg/software/virtual_env/anaconda3/condabin:/usr/lib64/qt-3.3/bin:/usr/local/bin:/usr/bin:/usr/sbin:/opt/ibutils/bin:/local/home/lstanalyzer/.local/bin:/local/home/lstanalyzer/bin', 'LD_LIBRARY_PATH': '/local/home/lstanalyzer/usr/lib:', 'DYLD_LIBRARY_PATH': None, 'USER': 'lstanalyzer', 'HOME': '/local/home/lstanalyzer', 'SHELL': '/bin/bash'}, 'arguments': ['/fefs/aswg/lstosa/datasequence.py', '-c', 'cfg/sequencer_Nov2019_dragontime_v03.cfg', '-d', '2019_11_23', '--prod_id', 'v0.5.1_v03', '/fefs/aswg/data/real/calibration/20191123/v03/calibration.Run1614.0000.hdf5', '/fefs/aswg/data/real/calibration/20191123/v03/drs4_pedestal.Run1611.0000.fits', '/fefs/aswg/data/real/calibration/20191123/v03/time_calibration.Run1614.0000.hdf5', '/fefs/aswg/scripts-osa/corrected_drive_logs_Nov19/drive_log_19_11_23.txt', '0', '0', '0', '0', '--stderr=sequence_LST1_01618_2432982.err', '--stdout=sequence_LST1_01618_2432982.out', '01618.0008', 'LST1'], 'start_time_utc': '2020-05-18T14:18:30.445684'}, 'software_version': 'v0.5.1', 'observation_date': '20191123', 'observation_run': '01618', 'session_tag': 'r0_to_dl1:01618'}
INFO provLogger __PROV_2020-05-18T14:18:30.447360_PROV__{'activity_id': '621ca2', 'name': 'r0_to_dl1', 'startTime': '2020-05-18T14:18:06.362321', 'in_session': 8739478486569, 'agent_name': 'lstanalyzer', 'script': '/fefs/aswg/lstosa/datasequence.py', 'session_tag': 'r0_to_dl1:01618'}
INFO provLogger __PROV_2020-05-18T14:18:30.447499_PROV__{'activity_id': '621ca2', 'parameters': {'ObservationRun': '01618', 'ObservationSubRun': '0008', 'CalibrationRun': '01614', 'PedestalRun': '01611', 'ProdID': '03'}, 'session_tag': 'r0_to_dl1:01618'}
INFO provLogger __PROV_2020-05-18T14:18:30.447628_PROV__{'entity_id': '446f45dd1c878559585395eedce5bc7a', 'name': 'R0SubrunDataset', 'filepath': '/fefs/aswg/data/real/R0/20191123/LST-1.1.Run01618.0008.fits.fz', 'hash': '446f45dd1c878559585395eedce5bc7a', 'hash_type': 'md5', 'type': 'File', 'contentType': 'application/fits', 'session_tag': 'r0_to_dl1:01618'}
INFO provLogger __PROV_2020-05-18T14:18:30.447752_PROV__{'activity_id': '621ca2', 'used_id': '446f45dd1c878559585395eedce5bc7a', 'used_role': 'Observation subrun', 'session_tag': 'r0_to_dl1:01618'}
INFO provLogger __PROV_2020-05-18T14:18:30.447859_PROV__{'entity_id': '7404bb00748454d63badfe247c774a13', 'name': 'PedestalFile', 'filepath': '/fefs/aswg/data/real/calibration/20191123/v03/drs4_pedestal.Run1611.0000.fits', 'hash': '7404bb00748454d63badfe247c774a13', 'hash_type': 'md5', 'type': 'File', 'contentType': 'application/fits', 'session_tag': 'r0_to_dl1:01618'}
INFO provLogger __PROV_2020-05-18T14:18:30.447966_PROV__{'activity_id': '621ca2', 'used_id': '7404bb00748454d63badfe247c774a13', 'used_role': 'Pedestal file', 'session_tag': 'r0_to_dl1:01618'}
INFO provLogger __PROV_2020-05-18T14:18:30.448065_PROV__{'entity_id': 'd8077e3bbdbc371f688ae65b0972e212', 'name': 'CoefficientsCalibrationFile', 'filepath': '/fefs/aswg/data/real/calibration/20191123/v03/calibration.Run1614.0000.hdf5', 'hash': 'd8077e3bbdbc371f688ae65b0972e212', 'hash_type': 'md5', 'type': 'File', 'contentType': 'application/x-hdf', 'session_tag': 'r0_to_dl1:01618'}
INFO provLogger __PROV_2020-05-18T14:18:30.447752_PROV__{'session_id': 8773094569769, 'name': '01618', 'startTime': '2020-05-18T14:18:06.362323', 'system': {'executable': '/fefs/aswg/software/virtual_env/anaconda3/envs/osa/bin/python', 'platform': {'architecture_bits': '64bit', 'architecture_linkage': '', 'machine': 'x86_64', 'processor': 'x86_64', 'node': 'cp15', 'version': '#1 SMP Thu Nov 8 23:39:32 UTC 2018', 'system': 'Linux', 'release': '3.10.0-957.el7.x86_64', 'libcver': "(('glibc', '2.10')", 'num_cpus': 32, 'boot_time': '2020-03-24T03:48:43'}, 'python': {'version_string': '3.7.6 | packaged by conda-forge | (default, Mar 23 2020, 23:03:20) \n[GCC 7.3.0]', 'version': '3.7.6', 'compiler': 'GCC 7.3.0', 'implementation': 'CPython'}, 'environment': {'CONDA_DEFAULT_ENV': 'osa', 'CONDA_PREFIX': '/fefs/aswg/software/virtual_env/anaconda3/envs/osa', 'CONDA_PYTHON_EXE': '/fefs/aswg/software/virtual_env/anaconda3/bin/python', 'CONDA_EXE': '/fefs/aswg/software/virtual_env/anaconda3/bin/conda', 'CONDA_PROMPT_MODIFIER': '(osa)', 'CONDA_SHLVL': '2', 'PATH': '/local/home/lstanalyzer/usr/bin:/local/home/lstanalyzer/.local/bin:/fefs/aswg/software/virtual_env/anaconda3/envs/osa/bin:/fefs/aswg/software/virtual_env/anaconda3/condabin:/usr/lib64/qt-3.3/bin:/usr/local/bin:/usr/bin:/usr/sbin:/opt/ibutils/bin:/local/home/lstanalyzer/.local/bin:/local/home/lstanalyzer/bin', 'LD_LIBRARY_PATH': '/local/home/lstanalyzer/usr/lib:', 'DYLD_LIBRARY_PATH': None, 'USER': 'lstanalyzer', 'HOME': '/local/home/lstanalyzer', 'SHELL': '/bin/bash'}, 'arguments': ['/fefs/aswg/lstosa/datasequence.py', '-c', 'cfg/sequencer_Nov2019_dragontime_v03.cfg', '-d', '2019_11_23', '--prod_id', 'v0.5.1_v03', '/fefs/aswg/data/real/calibration/20191123/v03/calibration.Run1614.0000.hdf5', '/fefs/aswg/data/real/calibration/20191123/v03/drs4_pedestal.Run1611.0000.fits', '/fefs/aswg/data/real/calibration/20191123/v03/time_calibration.Run1614.0000.hdf5', '/fefs/aswg/scripts-osa/corrected_drive_logs_Nov19/drive_log_19_11_23.txt', '0', '0', '0', '0', '--stderr=sequence_LST1_01618_2432982.err', '--stdout=sequence_LST1_01618_2432982.out', '01618.0006', 'LST1'], 'start_time_utc': '2020-05-18T14:18:30.447727'}, 'software_version': 'v0.5.1', 'observation_date': '20191123', 'observation_run': '01618', 'session_tag': 'r0_to_dl1:01618'}
INFO provLogger __PROV_2020-05-18T14:18:30.448191_PROV__{'activity_id': '621ca2', 'used_id': 'd8077e3bbdbc371f688ae65b0972e212', 'used_role': 'Coefficients calibration file', 'session_tag': 'r0_to_dl1:01618'}
INFO provLogger __PROV_2020-05-18T14:18:30.448243_PROV__{'activity_id': '7e8065', 'name': 'r0_to_dl1', 'startTime': '2020-05-18T14:18:06.362323', 'in_session': 8773094569769, 'agent_name': 'lstanalyzer', 'script': '/fefs/aswg/lstosa/datasequence.py', 'session_tag': 'r0_to_dl1:01618'}
INFO provLogger __PROV_2020-05-18T14:18:30.448296_PROV__{'entity_id': 'fcf52d425d9033504c154f25986978c2', 'name': 'TimeCalibrationFile', 'filepath': '/fefs/aswg/data/real/calibration/20191123/v03/time_calibration.Run1614.0000.hdf5', 'hash': 'fcf52d425d9033504c154f25986978c2', 'hash_type': 'md5', 'type': 'File', 'contentType': 'application/x-hdf', 'session_tag': 'r0_to_dl1:01618'}
```



data processing provenance graph



json serialization

```

{
  "prefix": {
    "id": "id:",
    "default": "param:"
  },
  "entity": {
    "id:8756102484710": {
      "prov:label": "r0_to_d11",
      "prov:type": "ExecutionSession",
      "prov:generatedAtTime": "2020-05-08T11:51:47.051025",
      "script": "/home/daniel.morcuende/lstosa/datasequence.py",
      "software_version": "v0.4.5.post566+gite395ef2",
      "observation_date": "20200118",
      "observation_run": "01832"
    },
    "id:4ble22_parameters": {
      "ObservationRun": "01832",
      "CalibrationRun": "01831",
      "PedestalRun": "01830",
      "ProdID": "02",
      "prov:type": "Parameters",
      "prov:label": "Parameters"
    },
    "id:6d23620c4d92ef92b642a0116908fad1": {
      "prov:label": "PedestalFile",
      "prov:type": "File",
      "filepath": "/fefs/aswg/workspace/daniel.morcuende/data/real/calibration/20200118/v02/drs4_pedestal.Run1830.0000.fits",
      "hash": "6d23620c4d92ef92b642a0116908fad1",
      "hash_type": "md5",
      "contentType": "application/fits"
    },
    "id:58f0d6b71f12a4e8a73eaad09345d64c": {
      "prov:label": "CoefficientsCalibrationFile",
      "prov:type": "File",
      "filepath": "/fefs/aswg/workspace/daniel.morcuende/data/real/calibration/20200118/v02/calibration.Run1831.0000.hdf",
      "hash": "58f0d6b71f12a4e8a73eaad09345d64c",
      "hash_type": "md5",
      "contentType": "application/x-hdf"
    },
    "id:a2856fc4fffcc1c9add2227bac130fb": {
      "prov:label": "TimeCalibrationFile",

```


provenance products

- DL1
 - 20200218
 - v0.4.3_v00
 - v0.5.0_v00
 - log
 - calibration.Run2006.0000.hdf5
 - drive_log_20_02_18.txt
 - drs4_pedestal.Run2005.0000.fits
 - lstchain_standard_config.json
 - r0_to_d11_02007_prov.json
 - r0_to_d11_02007_prov.log
 - r0_to_d11_02007_prov.pdf
 - r0_to_d11_02008_prov.json
 - r0_to_d11_02008_prov.log
 - r0_to_d11_02008_prov.pdf
 - r0_to_d11_02009_prov.json
 - r0_to_d11_02009_prov.log
 - r0_to_d11_02009_prov.pdf
 - sequencer.cfg
 - time_calibration.Run2006.0000.hdf5
- DL2
 - 20200218
 - v0.4.3_v00
 - v0.5.0_v00
 - log
 - dl1_to_dl2_02007_prov.json
 - dl1_to_dl2_02007_prov.log
 - dl1_to_dl2_02007_prov.pdf
 - dl1_to_dl2_02008_prov.json
 - dl1_to_dl2_02008_prov.log
 - dl1_to_dl2_02008_prov.pdf
 - dl1_to_dl2_02009_prov.json
 - dl1_to_dl2_02009_prov.log
 - dl1_to_dl2_02009_prov.pdf
 - lstchain_standard_config.json
 - r0_to_dl2_02007_prov.json
 - r0_to_dl2_02007_prov.log
 - r0_to_dl2_02007_prov.pdf
 - r0_to_dl2_02008_prov.json
 - r0_to_dl2_02008_prov.log
 - r0_to_dl2_02008_prov.pdf
 - r0_to_dl2_02009_prov.json
 - r0_to_dl2_02009_prov.log
 - r0_to_dl2_02009_prov.pdf
 - sequencer.cfg



- provides a graph view of the whole process
- chain of commands datasets parameters/values

Improvements on provenance tracking #34

 OpenBultako wants to merge 14 commits into `main` from `prov`  Conversation 0 Commits 14 Checks 1 Files changed 8

Bultako commented 1 hour ago

Collaborator



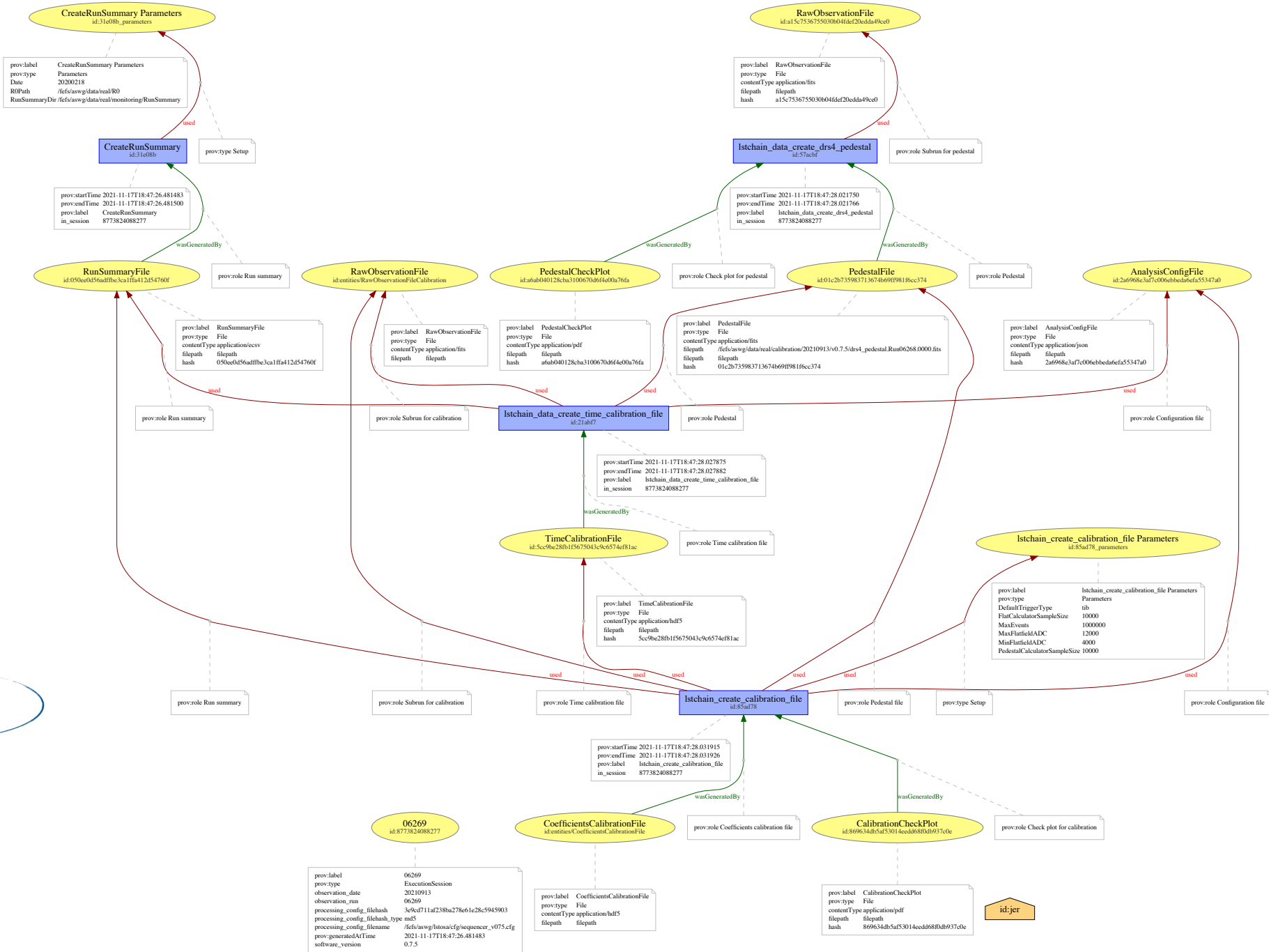
The content of this PR is the following:

- fixes some issues derived from changes in the data processing steps that have not been translated into the model up to now.
- adds calibration processes to the simulate processing script
- tracks and saves python packages dependencies of the conda environment
- adds `lstchain_d11ab` process into the data sequence provenance
- adds `d11_datacheck` process into the data sequence provenance
- adds missing DL2 merged file to the provenance
- adds Istosa config file used (i.e. `sequencer.cfg`) as a session parameter

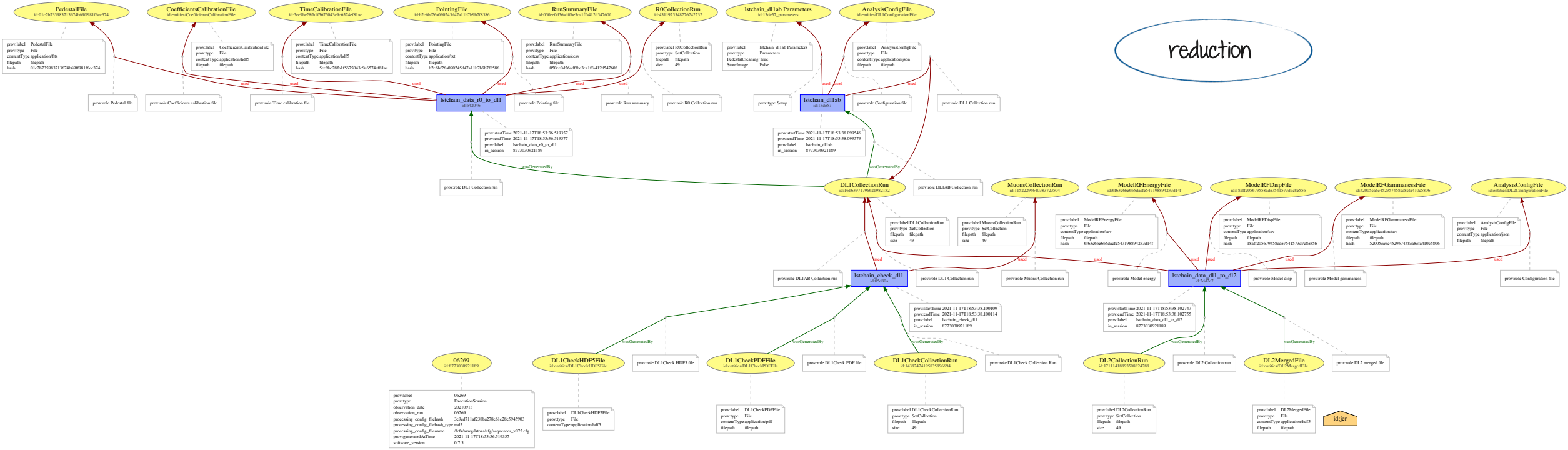
The final graph produced is as below.

[provDLshort.pdf](#)

calibration



id:jer

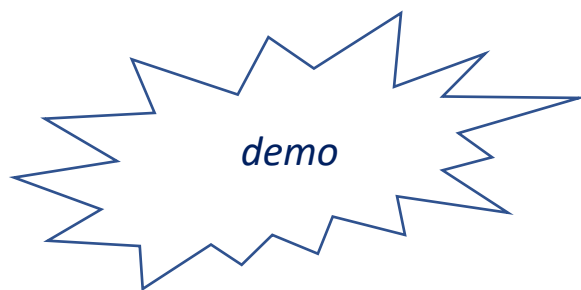


reduction

```

prov-label 06269
prov-type ExecutionSession
observation_date 20210913
observation_run 06269
processing_config_file_hash 1c9d711d238ba27861e2b5945903
processing_config_file_hash_type md5
processing_config_filename /cds/www/bdow/cfg/sequenceer_v075.cfg
prov-generatedAtTime 2021-11-17T18:53:36.519357
software_version 0.7.5
  
```

adger



Lessons learnt and next steps

- Continuous update and addition of captured info – **flexible model and implementation**
- **Configuration** may be stored as metadata attached to datasets (*fits headers/hdf5 attributes*)
- **Capturing relationships** among activities and entities (*datasets*) needs a provenance model
- **Independent capture for the different dependent software** packages is possible
Istosa requires Istchain requires ctapipe
- Massive logging subrun-wise makes a very poor performance because of **huge usage of disk**
- **Post-processing** of prov info is needed to produce run-wise info and different **levels of granularity**
- **Structured logging** in text files may be a solution for storage of small sessions provenance
- Considering storing provenance in a RDBMS or *better* in a noSQL **database** (i.e. *mongo + json*)
- Development of a **provenance query mechanism** for detailed analysis and inspection is needed
*...considering using url query params to produce on-the-fly **SVG graphs with access-links in a browser***

