

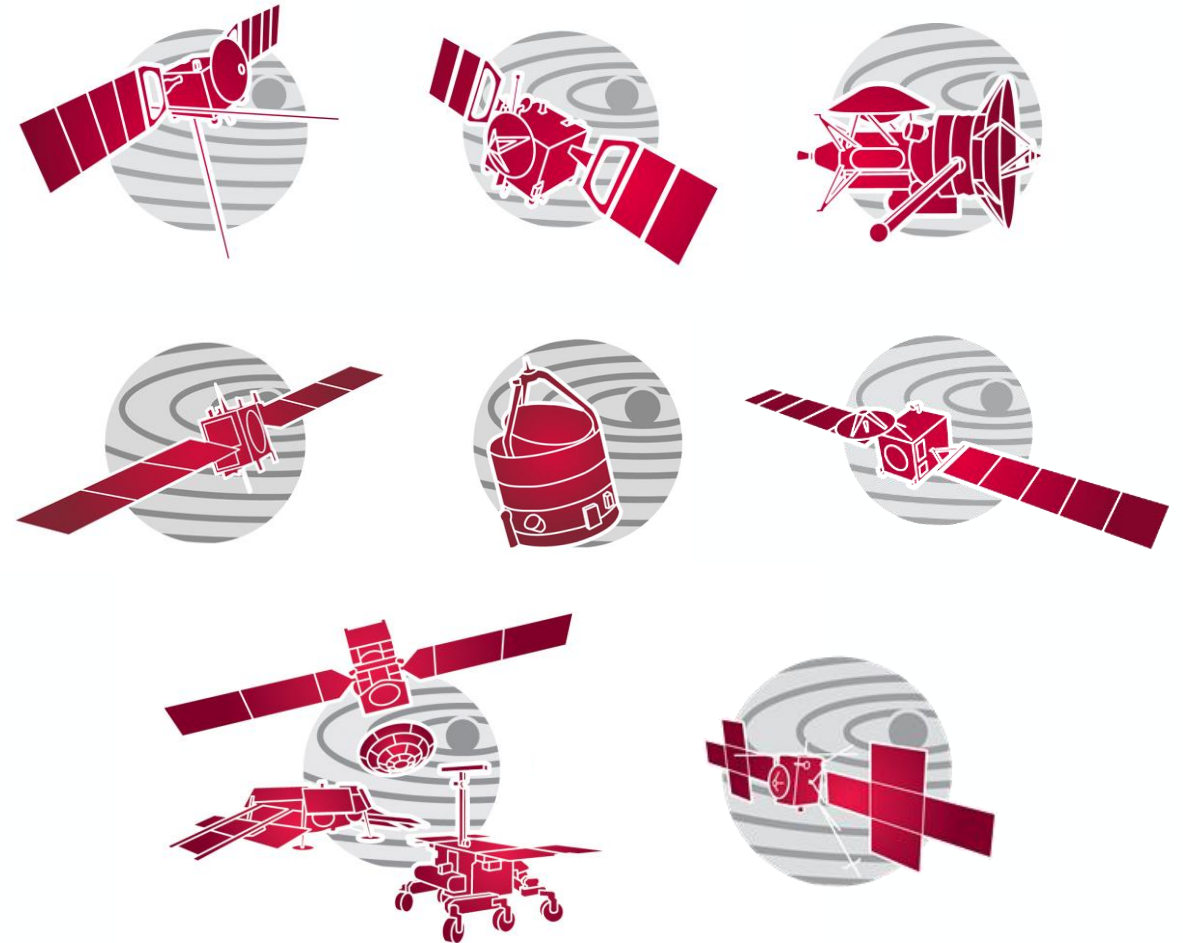
DataLabs for the Planetary Science Archive

Mark S. Bentley
on behalf of the PSA team

06/10/2022

Reminder: What is the PSA

- The Planetary Science Archive (PSA) is:
 - a multi-mission science archive
 - for ESA's planetary missions
 - orbiters, landers, rovers, etc.
- Long term archive
 - data should be usable in 50 years
- But also used operationally
 - PDS4 missions deliver ~daily
 - data private until reviewed and released
 - (data volumes have a mix of public/private data)
- Adopts the NASA PDS standard
 - missions delivering PDS3 and 4
 - these “wrap” many standard formats
 - extensive meta-data based on an information model



The scope for using DataLabs is clearly enormous, and we won't know what the community will do with it until we open it up and provide the relevant tools and APIs. However, the first items likely to be of use are:

- Data Tutorials
 - PSA contains data from many different instrument types, stored in many different ways
 - User feedback is that **finding** data is easy, but working with it is hard
 - One way to help is to provide tutorials in Jupyter Labs or similar, e.g.
 - replicate the analysis from a paper
 - repeat the calibration from raw to calibrated, or derived
 - DataLabs works well because we can provide a working environment
 - and users can copy/paste and develop their own more complex workflow from here
- Work with arbitrary meta-data
 - PSA only indexes a small fractions of the tens of thousands of meta-data attributes
 - DataLabs allows databases to be built/used without download thousands of data product
 - note that this may only be a temporary use case until we have all meta-data available

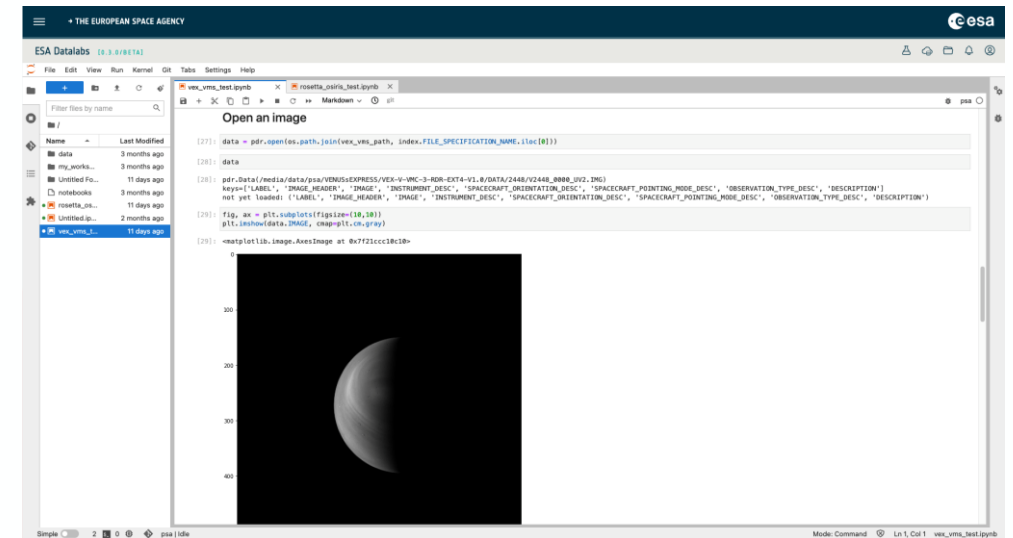
What do we need from DataLabs?

- In the first instance, the main low hanging fruit needs only Jupyter Labs
 - plus an appropriate python environment
 - fortunately the superset of useful packages is fairly constrained
- Note two extremely useful open source packages for reading actual data:
 - PDS3: <https://github.com/MillionConcepts/pdr>
 - PDS4: https://github.com/Small-Bodies-Node/pds4_tools
- The examples shown today use the generic ESDC Jupyter Labs DataLab
 - with a “psa” conda environment loaded with the relevant packages
- Moving forward we need to engage the community to see what else they may need

Approaches to DataLabs integration

The PSA is following this approach for integrating data with DataLabs:

- Internally moved data from legacy missions to two volumes (legacy + Rosetta) and mount them
 - ongoing missions with a mix of public/private data will be handled separately
 - discussion needed – internally we are setting up FUSE for this purpose
- A new TAP column will be added with the path to each data product
 - needs a stable mount point in a PSA DataLab
 - allows users to query and load data dynamically
- We need a PSA DataLab
 - Jupyter Labs + common python tools
 - pdr, pds4_tools etc
 - currently using generic lab for testing
- Eventually may link UI to DataLabs
 - e.g. “send to DataLabs” option



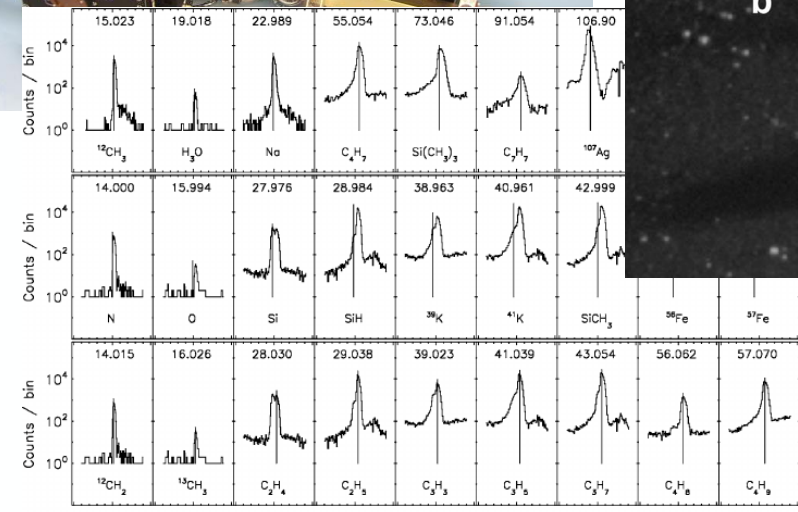
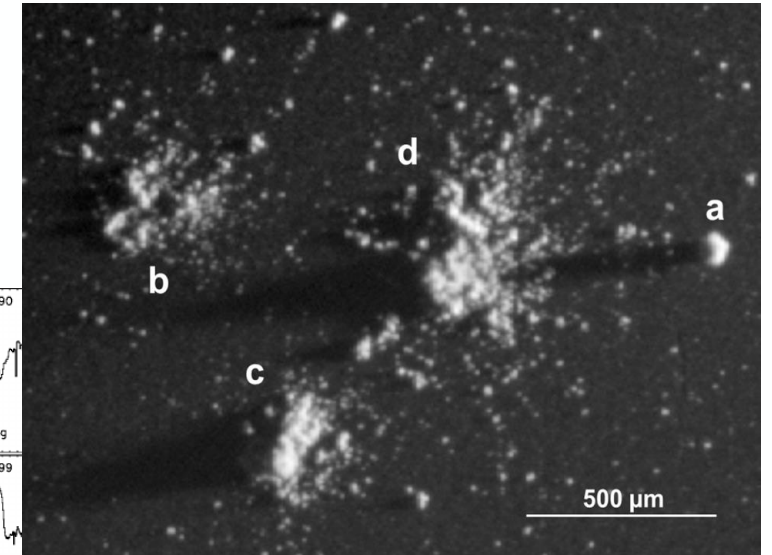
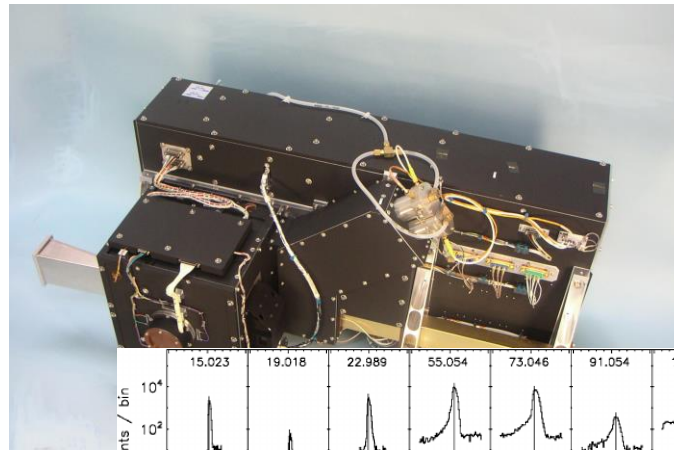
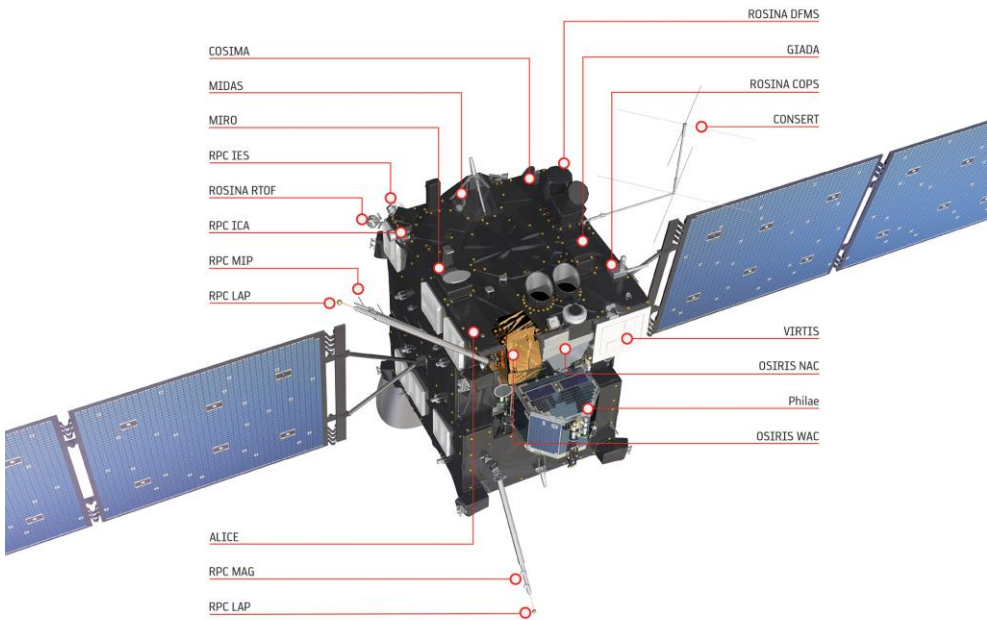
Coming soon – PDS registry integration

- PSA plans to ingest all PDS4 meta-data into the PDS registry
 - an AWS-hosted OpenSearch instance
 - structured using the PDS information model
- This will enable use of the PDS API on PSA data
 - including full meta-data search/access which is a highly requested feature
- DataLabs/Jupyter notebooks may be updated to use the API
 - as an alternative to TAP, allowing queries of arbitrary meta-data
- TBD if we need a local copy of the registry in-house, sync'd with the PDS one
 - for local searches from PSA UI and/or DataLabs
 - and to map the datasource to a local path rather than https URL
 - or we do that translation on the fly



Examples: Rosetta/COSIMA

- This example was prepared by Maddalena Buggati as part of her internship
 - and updated slightly for this workshop
 - it uses data from the COSIMA microscope and dust mass spectrometer on the Rosetta orbiter



(a)

→ THE EUROPEAN SPACE AGENCY

ESA Datalabs [0.3.0/BETA]

TUTORIALCOSIMA.IPYNB

Tutorial COSIMA

- Instrument description:
- Data description:
- Tutorial description:
- 0. Import libraries ...
- 1. Locate the data
- 2. Substrate (target) history
- 3. Microscope images
- 4. Dust grain table
 - 4a) Grain identification using the grain list table (GR_.TAB)
 - 4b) Grain identification using a different method
 - 4c) Analysis of one dust grain dimension
- 5. Mass spectra
- 6. Identified peaks
- 7. Dust grains per day (exposure)

Image before - 7 August

1st exposure - 17 August

2nd exposure - 24 August

4. Dust grain table

This section deals with the identification of the dust grains on the COSISCOPE image. A list of grain positions is stored in the `GR_.TAB` table where the position is given in substrate coordinates. These coordinates are in microns and refers to the x and y position in microns starting from the bottom left corner of the substrate (0-10000 microns in X and Y direction). To better understand the coordinates system and how the left bottom corner position is constrained please refer to <https://archives.esac.esa.int/psa/ftp/INTERNATIONAL-ROSETTA-MISSION/COSIMA/RO-C-COSIMA-3-V6.0/CATALOG/DATASET.CAT> (almost at the end of the PDF are explained different coordinate systems).

COSISCOPE onboard COSIMA automatically calculates the origin and rotation offsets due to mechanical target holder grasping inaccuracy of the substrate and store these information in the file header. To convert the substrate coordinates in pixel coordinates on the COSISCOPE image it is then sufficient to access the file header.

Unfortunately the `pdr` library is not able to access fits header and then we must use another library (`astropy.io.fits` library).

In the header it's stored the position of the bottom left corner of the substrate under this name:

Simple 0 3 | psal | Idle Mode: Command Ln 1, Col 1 TutorialCosima.ipynb

Examples: OSIRIS image matches

- Developed as part of the Zooniverse [Rosetta Zoo](#) project
 - Citizen science project to identify and classify changes on comet 67P using OSIRIS images
 - originally downloaded 20k meta-data label files from PSA
 - with DataLabs we don't have to do this!
- rosetta_zoo_metadata.ipynb (the real use of DataLabs is here!)
 - scrapes and cleans meta-data from matching data products and adds to pandas DataFrame
 - adds URLs to browse and data products
 - dumps to CSV
- rosetta_zoo_orientation.ipynb (not shown)
 - uses SPICE to add a north vector to each image to aid orientation
 - dumps to CSV
- rosetta_zoo_pairs (the final product)
 - defines a “matching function” to identify images of similar areas of comet 67P
 - displays candidate pairs (with optional image rotation)

Examples: OSIRIS image matches



→ THE EUROPEAN SPACE AGENCY

ESA Datalabs [0.3.0/BETA]

File Edit View Run Kernel Git Tabs Settings Help

ROSETTA_ZOO_METADATA.IPYNB

rosetta_zoo_pairs.ipynb x rosetta_zoo_metadata.ipynb

```
[4]: datasets = glob.glob('/media/data/rosetta/RO-C-OSINAC-3-*V3.0')
```

Now we are going to read the meta-data of the labels for each NAC image in the above datasets:

```
[12]: meta = {}
for dset in datasets:
    dset_name = os.path.split(dset)[-1]
    meta[dset_name] = Database(files='N*.LBL', directory=os.path.join(dset, 'DATA/FIT'), config_file=config_file, recursive=False)
```

Ultimately we want to combine each of these datasets into a master set of image meta-data:

```
[22]: df_list = [meta[dset].get_table('osiris_nac_zoo') for dset in meta.keys()]
zoo_meta = pd.concat(df_list)
```

Now we can take a look at the database we have produced:

```
[27]: zoo_meta.head()
```

```
[27]:
```

	filename	dataset	prod_id	mission_id	start_time	stop_time	instr_id	surf_int_x	surf_int_y	surf_int_z	distance	sc_altitude	phase_angle	sc_position_x	sc_positi
0	/media/data/rosetta/RO-C-OSINAC-3-ESC1-67PCHUR...	RO-C-OSINAC-3-ESC1-67PCHURYUMOV-M10-V3.0	N20141122T044453801ID30F22.FIT	ROSETTA	2014-11-22 04:46:05.229	2014-11-22 04:46:05.427	OSINAC	-1.214	-1.229	-0.802	29.498	29.41187	94.60917	19.961	
1	/media/data/rosetta/RO-C-OSINAC-3-ESC1-67PCHUR...	RO-C-OSINAC-3-ESC1-67PCHURYUMOV-M10-V3.0	N20141122T044501581ID30F23.FIT	ROSETTA	2014-11-22 04:46:13.009	2014-11-22 04:46:13.302	OSINAC	-1.216	-1.228	-0.803	29.497	29.41182	94.60928	19.96	3
2	/media/data/rosetta/RO-C-OSINAC-3-ESC1-67PCHUR...	RO-C-OSINAC-3-ESC1-67PCHURYUMOV-M10-V3.0	N20141122T044509456ID30F24.FIT	ROSETTA	2014-11-22 04:46:20.884	2014-11-22 04:46:21.171	OSINAC	-1.218	-1.227	-0.803	29.497	29.41177	94.60939	19.96	3
3	/media/data/rosetta/RO-C-OSINAC-3-ESC1-67PCHUR...	RO-C-OSINAC-3-ESC1-67PCHURYUMOV-M10-V3.0	N20141122T044517756ID30F27.FIT	ROSETTA	2014-11-22 04:46:29.184	2014-11-22 04:46:30.084	OSINAC	-1.219	-1.226	-0.804	29.495	29.41172	94.6095	19.959	3

Simple 0 4 psa | Idle Mode: Command Ln 1, Col 1 rosetta_zoo_metadata.ipynb



Examples: OSIRIS image matches

→ THE EUROPEAN SPACE AGENCY

ESA Datalabs [0.3.0/BETA]

File Edit View Run Kernel Git Tabs Settings Help

ROSETTA_ZOO_PAIRS.IPYNB

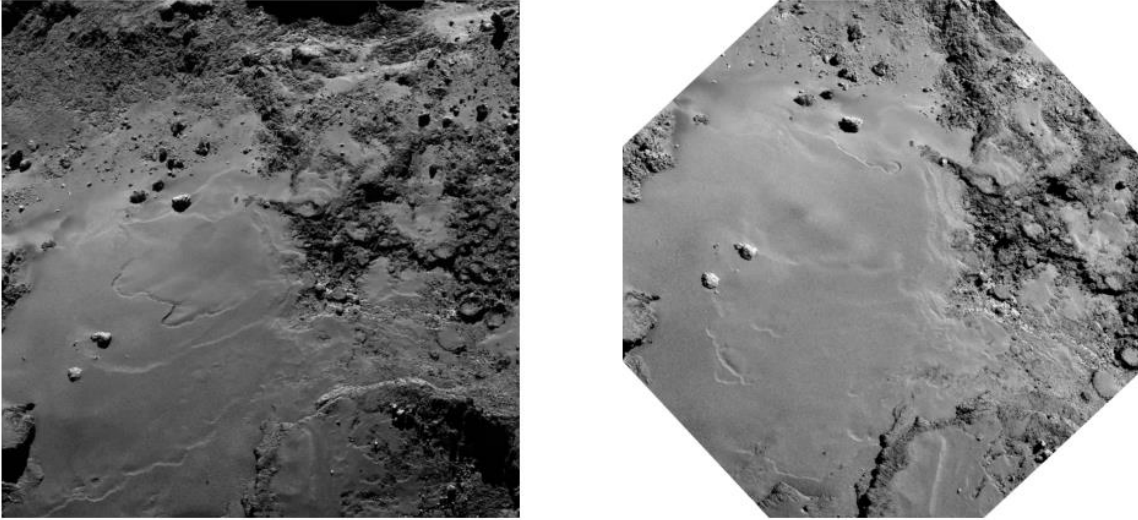
- Introduction
- Load the data
- Find some image pairs
- Find *all* image pairs

rosetta_zoo_pairs.ipynb x rosetta_zoo_metadata.ipynb rosetta_zoo_orientation.ipynb

So we have >3k matches according to the criteria in the best_match function. Many of these will not be useful and so human inspection will be needed to quickly discard some. Now we can loop through this list and accept or discard each match. Here we show a few examples:

```
[25]: idx = 50
      show_pair(pre, post, match_pre[idx], match_post[idx], rotate=True)
```

N20141122T071302582ID30F23 N20160613T060353637ID30F22



```
[38]: idx = 400
      show_pair(pre, post, match_pre[idx], match_post[idx], rotate=True)
```

N20141203T122934611ID30F51 N20160622T1527388 [Table of Contents](#)

Simple 0 5 psa | Idle Mode: Command Ln 1, Col 1 rosetta_zoo_pairs.ipynb

- PSA DataLab
 - request DataLab (once static mountpoints can be defined with a Lab)
 - update TAP with datalabs_path column
 - investigate how to provide authenticated access to private data
 - we are currently implementing an SFTP service which offers this
- Data tutorials
 - tidy up Rosetta notebooks and ask instrument teams to review, if possible
 - create new notebooks for other missions and instruments
 - BepiColombo – data tutorials are a WP in ongoing EXPRO contracts with the teams
 - other missions – developed in house, by interns, or archive scientists?

- Tutorials
 - how do we validate these, especially when instrument teams are no longer around?
 - reproducing results of a paper would be one way, reproducing calibration/processing steps would be another
 - how do we maintain the notebooks/tutorials
 - in the face of library updates, API changes, new data deliveries etc.
 - we can freeze our environment, but we want to stay more-or-less current
- Jupyter @ DataLabs in general
 - what is the typical user workflow?
 - often – develop in notebook, refactor into a module, re-write notebook to use module
 - do we offer any Git integration or similar? or rely on GitHub and co?
 - can we expose running notebooks externally so IDEs can hook into them?
 - i.e. connect to the running kernel (current URIs are localhost and probably behind a proxy)
 - can users manage Lab extensions themselves?
 - can users update packages in the base environment themselves?
 - in general how often do we update the Jupyter base environments