

University of Groningen

The Euclid Archive Processing and Data Distribution Systems

Williams, O. R.; Begeman, K.; Boxhoorn, D.; Droge, B.; Nutma, T. A.; Tsyganov, A.; Valentijn, E. A.; Vriend, W.-J.; Dabin, C.

Published in:
Astronomical Data Analysis Software and Systems XXIX

IMPORTANT NOTE: You are advised to consult the publisher's version (publisher's PDF) if you wish to cite from it. Please check the document version below.

Document Version
Publisher's PDF, also known as Version of record

Publication date:
2020

[Link to publication in University of Groningen/UMCG research database](#)

Citation for published version (APA):

Williams, O. R., Begeman, K., Boxhoorn, D., Droge, B., Nutma, T. A., Tsyganov, A., Valentijn, E. A., Vriend, W.-J., & Dabin, C. (2020). The Euclid Archive Processing and Data Distribution Systems: A Distributed Infrastructure for Euclid and Associated Data. In R. Pizzo, E. R. Deul, J-D. Mol, J. de Plaa, & H. Verkouter (Eds.), *Astronomical Data Analysis Software and Systems XXIX: Proceedings of the 29th annual conference on Astronomical Data Analysis Software and Systems (ADASS XXIX)*. (pp. 291-294). (Astronomical Data Analysis Software and Systems XXIX. ASP Conference Series; Vol. 527). Astronomical Society of the Pacific. <http://adsabs.harvard.edu/abs/2020ASPC..527..291W>

Copyright

Other than for strictly personal use, it is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license (like Creative Commons).

The publication may also be distributed here under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license. More information can be found on the University of Groningen website: <https://www.rug.nl/library/open-access/self-archiving-pure/taverne-amendment>.

Take-down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

Downloaded from the University of Groningen/UMCG research database (Pure): <http://www.rug.nl/research/portal>. For technical reasons the number of authors shown on this cover page is limited to 10 maximum.

The Euclid Archive Processing and Data Distribution Systems: a Distributed Infrastructure for Euclid and Associated Data

O. R. Williams,¹ K. Begeman,² A. N. Belikov,² D. Boxhoorn,² B. Droge,¹
T. A. Nutma,² A. Tsyganov,¹ E. A. Valentijn,² W-J. Vriend,² and C. Dabin³

¹*Centre for Information technology, University of Groningen, Groningen;*
o.r.williams@rug.nl

²*Kapteyn Institute, University of Groningen, The Netherlands*

³*CNES, Toulouse, France*

Abstract. The Euclid Archive System is an ambitious information system, which sits at the heart of the Euclid Science Ground Segment. It is a joint development between the Euclid Consortium and the ESAC Science Data Centre. It encompasses both Euclid data and the large volume of associated ground based data (e.g. KiDS, DES and LSST).

The Euclid Science Ground Segment consists of the Euclid Science Operations Centre and ten national Science Data Centres. The large data volumes demand that data transfer is minimized and that the processing is taken to the data. This is supported by the Euclid Archive Data Processing System and the Euclid Archive Distributed Data System. The Data Processing System consists of a central metadata repository, which contains the information necessary to process any data item and full data lineage of any data product created. The Distributed Data System provides a cloud solution with a node at each of the national Science Data Centres, which controls data storage and transfer. It supports a large number of storage types, including POSIX, iRODS, gridftp and Xrootd. No limitations are placed on the storage implemented at an individual SDC. Further more, the user of the system needs no knowledge of where data is located. Jobs will be started at the most appropriate locations, or data transferred as necessary.

1. Introduction

The Euclid mission will be a milestone in the understanding of the geometry of the Universe (Laureijs et al. 2011). The Euclid Science Ground Segment (SGS) and the Euclid Archive System (EAS) have three main challenges during the data processing: firstly, the unprecedented accuracy which must be achieved in order to meet the scientific goals; secondly, the heavy dependence on the processing of ground-based data which will form the bulk of the stored data volume; thirdly the large scale reprocessing which will be required to meet the accuracy requirements (Pasian et al. 2014). In total Euclid may produce more than 26 PB per year of data during the later phases of the mission (Williams et al. 2014).

The EAS must provide the Euclid SGS with a distributed scientific information system, able to handle hundreds of PBs of data, together with tools to help in the assessment of the data quality for each produced item.

2. EAS and SGS

The EAS is responsible both for the support of data processing during the mission and for the delivery of the science-ready data to the astronomical community. The responsibility for delivery of science-ready data lies with the EAS Science Archive System which is described elsewhere (Nieto et al. 2019). The responsibility for data processing support, described in this paper, lies within two components of the EAS: the Data Processing System (EAS-DPS) and the Distributed Storage System (EAS-DSS).

The EAS-DPS stores metadata related to the data being processed, the orchestration of the processing, the quality assessment of the data products and the preparation of data releases. The EAS-DSS stores the data files themselves, from raw frames to calibrated images and spectra. Figure 1 shows an overview of the three components of the EAS and their principal interactions.

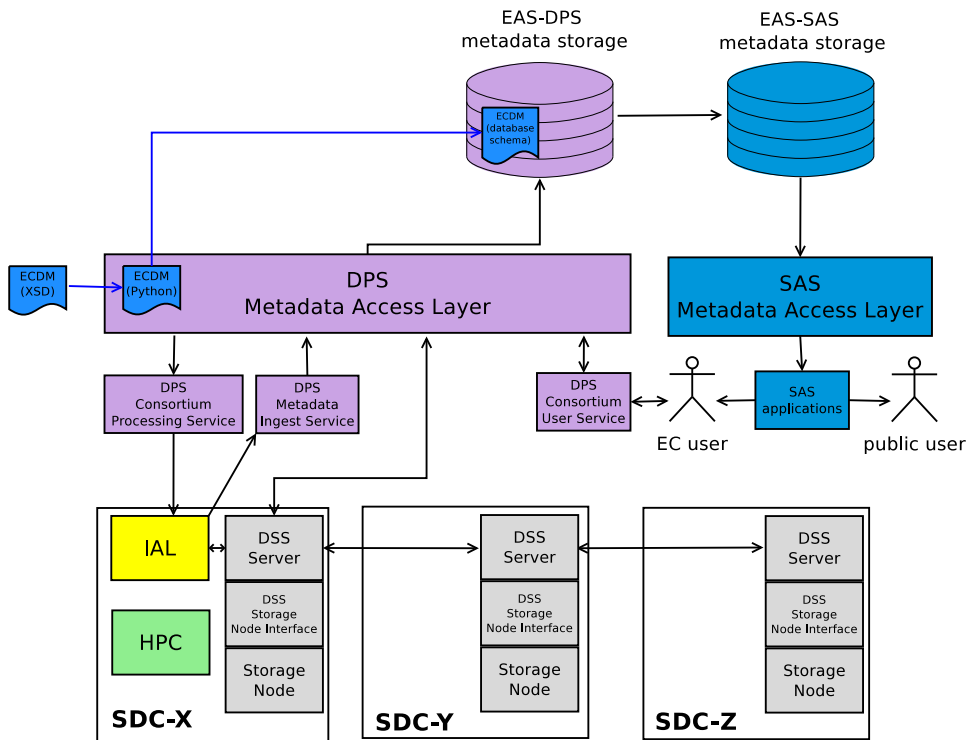


Figure 1. An Overview of the three components of the EAS

To implement the scientific requirement for traceability of the data, all operations during the data processing must be reflected in the metadata. Moreover, all data products necessary for the next step in the processing must be ingested into the EAS. The EAS is thus not merely an archive for the storage of data, but rather an information system which can provide an astronomer with a detailed overview of the status of the data which has been processed or is being processed. This information includes full backward and forward lineage for each data item. Such lineage is crucial for the quality assessment of the data and also to prevent unnecessary reprocessing.

Euclid processes data in a distributed environment which consists of ten national Science Data Centers and the Science Operation Center. The EAS-DPS and EAS-DSS must allow access to data from this distributed environment and guarantee data distribution according to the needs of the Euclid processing plan.

The designs of the EAS-DPS and EAS-DSS draw on lessons learned from earlier archives for OmegaCAM and the LOFAR Long Term Archive (Begeman et al. 2013).

3. EAS-DPS Design

The binding between the different processing steps in the Euclid SGS is defined by the Euclid Common Data Model (ECDM). The ECDM describes not only the input and output of each pipeline, but also contains the processing information for the Euclid SGS. The ECDM is based on the XML Schema Definition Language and forms an object-oriented data model. The EAS-DPS takes each stable release of ECDM and implements it by creating first Python stubs from each definition in ECDM and then generating a DDL schema to be created in the metadata database of the EAS-DPS (Williams et al. 2019).

The Python stubs form the core of the the Metadata Access Layer (MAL) of EAS-DPS. MAL can hide the complexity of the metadata database implementation from the user and allows the user to interact with objects formed according to ECDM instead of the selection of rows in tables. A Docker container has been developed to allow users to easily install the MAL on their own PCs and to use it within Jupyter notebooks.

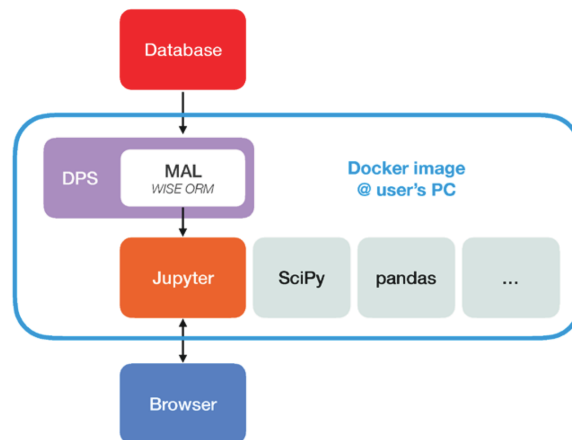


Figure 2. A Docker container allows users to install easily the MAL on their own workstations and to use it either with exiting Jupyter notebooks or their own.

The EAS-DPS provides a number of services for other SGS components and users. Firstly, the Consortium User Service allows users to browse the metadata database using a web browser. Secondly, the Consortium Processing Service is used by other SGS components to retrieve metadata from the EAS-DPS. Finally, the Metadata Ingest Service supports the transformation of XML to Python objects and their commitment to the metadata database.

4. EAS-DSS Design

The EAS-DSS serves both the EAS-DPS and the EAS-SAS as a common, distributed file storage solution. The EAS-DSS is a data storage grid with a single https-based user interface. The DSS servers act as an interface to the non-homogeneous data storage solutions deployed in the SDCs. Currently a DSS server can be deployed on top of a local POSIX filesystem, an iRODS server, an sftp server, a Grid storage element, an Astro-WISE dataserer, Openstack and Xrootd. At least one DSS server is installed at each SDC and stores the data files processed or created at this SDC.

5. EAS Status and Euclid Data Processing Challenges

An EAS Prototype, based on Oracle, was developed in 2013 and 2014. This prototype formed the basis of the first version of the EAS itself, released in 2015 (Belikov & co authors 2016). In 2017 a systematic study of RDBMS systems for the EAS-DPS was performed. It was decided that continuing with ORACLE gave the best performance.

The development schedule of the Euclid SGS is organized around large-scale integration and data processing challenges. These data processing challenges follow an incremental and iterative test-driven path and consolidate the interfaces, through concrete utilization of the common data mode (Dubath et al. 2017). The initial challenges conducted in 2016 allowed the testing of preliminary versions of the EAS-DPS and EAS-DSS. Design improvements were made, particularly to the EAS-DSS in order to improve the performance when transferring small files. Files are transferred now in groups, which decreases the overhead involved in authentications and database updates.

The EAS-DSS and EAS-DPS now meet the minimum performance requirements, although development to further improve the system continues. The latest versions will support the new generation of Euclid data processing challenges in 2020, as the SGS moves closer to processing realistic volumes.

References

- Begeman, K., Belikov, A. N., & co authors. 2013, *Experimental Astronomy*, 35, 1. 1208.0447
- Belikov, A., & co authors 2016, in *Proceedings of 2016 conference on Big Data from Space (BiDS'16)*, edited by P. Soille, & P. Marchetti (Publications Office of the European Union), vol. JRC100655, 212
- Dubath, P., Apostolakis, N., & co authors 2017, in *Astroinformatics*, edited by M. Brescia, S. G. Djorgovski, E. D. Feigelson, G. Longo, & S. Casulli, vol. 325 of IAU Symposium, 73. 1701.08158
- Laureijs, R., Amiaux, J., & co authors 2011, *ArXiv e-prints*. 1110.3193
- Nieto, S., de Teodoro, P., & co authors 2019, in *Astronomical Society of the Pacific Conference Series*, edited by P. J. Teuben, M. W. Pound, B. A. Thomas, & E. M. Warner, vol. 523 of *Astronomical Society of the Pacific Conference Series*, 437
- Pasian, F., Hoar, J., & co authors 2014, in *Astronomical Data Analysis Software and Systems XXIII*, edited by N. Manset, & P. Forshay, vol. 485 of *Astronomical Society of the Pacific Conference Series*, 505
- Williams, O., Begeman, K., & co authors 2019, in *Astronomical Society of the Pacific Conference Series*, edited by M. Molinaro, K. Shorridge, & F. Pasian, vol. 521 of *Astronomical Society of the Pacific Conference Series*, 120
- Williams, O., Belikov, A., & Koppenhoefer, J. 2014, in *Proc. of NETSPACE Workshop*, edited by O. Sykioti, & I. Dalgis (11), 491