

University of Groningen

The Euclid Archive System: A Datacentric Approach to Big Data

Belikov, Andrey; Williams, Owen; Altieri, Bruno; Boxhoorn, Danny; Buenadicha, Guillermo ; Droge, Bob; McFarland, John; Nieto, Sara; Salgado, Jesus; de Teodoro, Pilar

Published in:

Proceedings of 2016 conference on Big Data from Space (BiDS'16)

DOI:

[10.2788/854791](https://doi.org/10.2788/854791)

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

Early version, also known as pre-print

Publication date:

2016

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

Citation for published version (APA):

Belikov, A., Williams, O., Altieri, B., Boxhoorn, D., Buenadicha, G., Droge, B., McFarland, J., Nieto, S., Salgado, J., de Teodoro, P., Tsyganov, A., & Valentijn, E. (2016). The Euclid Archive System: A Datacentric Approach to Big Data. In P. Soille, & P. G. Marchetti (Eds.), *Proceedings of 2016 conference on Big Data from Space (BiDS'16)* (Vol. JRC100655, pp. 212-215). Publications Office of the European Union . <https://doi.org/10.2788/854791>

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 SYSTEM:

A DATA-CENTRIC APPROACH TO BIG DATA

A.N.Belikov⁽¹⁾, O.R.Williams⁽²⁾, B.Altieri⁽³⁾, D.Boxhoorn⁽¹⁾, G.Buenadicha⁽³⁾, B. Droege⁽²⁾, J.McFarland⁽¹⁾, S.Nieto⁽³⁾, J. Salgado⁽³⁾, P. de Teodoro⁽³⁾, A.Tsyganov⁽²⁾, E.A. Valentijn⁽¹⁾

(1) Kapteyn Astronomical Institute, University of Groningen, Groningen, The Netherlands

(2) Donald Smits Centre for Information Technology, University of Groningen, Groningen, The Netherlands

(3) European Space Astronomy Center, European Space Agency, Spain

ABSTRACT

The Euclid Archive System (EAS) is in the core of the Euclid Science Ground Segment (SGS). It supports the processing and storage of Euclid data from the input raw frames to the science-ready images and catalogs. We review the architectural design of the system, implementation progress and main challenges in the building of the EAS.

***Index Terms*— Euclid, Data Storage, Information System**

1. INTRODUCTION

Euclid is an ESA M2 mission and a milestone in the understanding of the geometry of the Universe [1]. Euclid faces two main challenges from the point of view of the data processing. Firstly, the unprecedented accuracy which must be achieved in order to meet the scientific goals. Secondly, the mission will depend heavily on the processing and reprocessing of the ground-based data which will form the bulk of the data volume to store [2]. In total Euclid will produce up to 26 PB per year of observations [3].

2. EAS AND SGS

The Euclid Science Ground Segment is a distributed data processing and data storage system which is responsible for the delivery of the science-ready data to ESA [3]. The SGS is formed by 9 national Science Data Centres and the EuclidScience Operations Centre. Each data centre provides resources for processing, along with expertise in the coding

of Euclid pipelines. The Euclid Archive System (EAS) is a core element of the SGS implementing data-centric approach [4] to the data processing.

In the SGS and the EAS the data is a combination of data files (images, spectra, catalogs) and metadata. The metadata includes the full data lineage, which is necessary to achieve the goals of the Euclid mission. The EAS must provide to SGS users and subsystems of SGS the ability to trace any bit of information produced in the SGS and must provide quality information on the produced images, spectra and catalogs. According to our experience in previous missions (Astro-WISE [4] and LOFAR Long-Term Archive [5]) the data volume of metadata will not exceed 5% of the total data volume. The clear difference between the bulk of the data stored in files and extensive metadata describing data products in files allows to move most of data mining operations to the relational database without access to the data files.

The EAS is a joint development between ESA and Euclid Consortium led by the SDC of the Netherlands and the ESDC (ESAC Science Data Centre) .

It should be stressed that the EAS is not an archive in the conventional sense. Instead, the EAS is a distributed scientific information system which ensures that any operation with the data is registered and can be traced back to the user and pipeline. EAS ensures that SGS can access and process hundreds of Petabytes of data.

3. EAS DESIGN

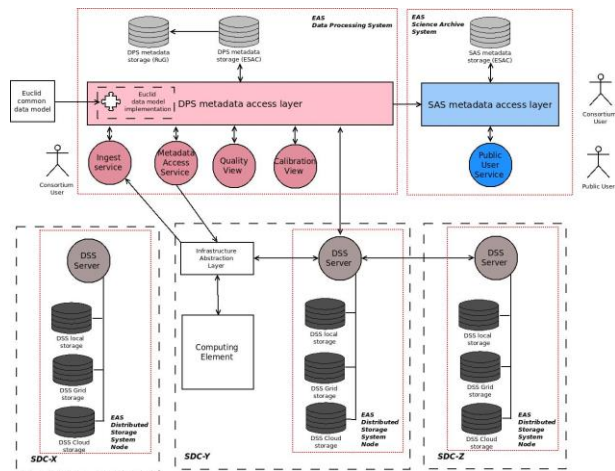


Figure 1: Overview of EAS and its place in Euclid SGS

The EAS design is based on requirements, formed by the ESA and the Euclid community, on the data processing in distributed environment. To allow the handling of numerous and often contradictory requirements on the EAS the whole system was divided on 3 independent parts:

- The Data Processing System (DPS) for the interactions within the SGS
- The Science Archive System (SAS) for the delivery of data releases, their long-term preservation and support of Euclid scientific use-cases,
- The distributed Storage System (DSS) for the storage of data files.

Figure 1 shows overview of the EAS design with primary users and their interaction with DPS, SAS and DSS.

The primary task of the DPS is to server as an information system for the production of the data release. To achieve this goal DPS keeps data lineage for each data object tracing any operation in the system to allow data reproduction and reprocessing. DPS allows to trace each frame, catalog or source catalog back to the original raw data, processing parameters used and pipeline which produced it. As well it allows to see the current status of the processing in SGS and quality of the raw and processed data.

The DSS serves both the DPS and SAS as a common distributed file storage solution. In this solution we utilize a

distributed approach for the data files storage and a centralized approach to the storage of the metadata. The DSS consists of DSS servers which are interfacing the non-homogeneous data storage solutions provided by the national SDCs. At least one DSS server is installed at each SDC and stores the data files processed or created during running of a pipeline in this SDC. The design of the processing plan for each pipeline allows to minimize data file transfer between SDCs. To ensure zero loss of the data at least 2 copies of each data file is created in DSS. Each copy is registered in the metadata storage.

The EAS implements the Euclid Common Data Model (ECDM) which describes both scientific data (data products generated by pipelines) and processing and operational metadata (processing and data distribution orders, location of the file at DSS, processing plan). The ECMD is implemented "as it is" inside the DPS and transformed to Science Exploration Data Model for the implementation in SAS. The interface between SAS and DPS will allow to SAS to retrieve the science metadata from the DPS for the Euclid data release and put it in long-term storage in SAS. The same ECMD is used to bind pipeline data flows.

The SAS is part of the EAS and it's aimed to support the needs for scientific data exploration for Euclid Consortium and wider astronomical community. The SAS, as part of the Euclid Archive, will face the challenge of Big Data, as it will store a huge and increasing amount of scientific metadata, and catalogues of 10 billion of galaxies. This will provide the worldwide astronomical community with an extremely large source of targets for future missions. Under these premises, the SAS will face the challenge to guarantee the preservation and public access of the data stored to the scientific community.

The SAS is being built at the ESAC Science Data Centre (ESDC), that is in charge of the development and operations of the scientific archives for the ESA Astronomy, Planetary and Heliosphysics missions. The SAS is focused on the needs from the scientific community. In this context, the SAS will provide access to the most valuable scientific metadata coming from the EAS-DPS through a set of public releases. According to the policy of public releases defined by the Euclid Consortium, the plan is to deliver 3 data releases to the scientific community every 2 years after the nominal start of the mission [6].

The design of the EAS is fully hardware independent and partially software independent. Metadata storage is based on

a RDBMS, currently Oracle 11g, but it will be possible to switch to other RDBMS in the future. Metadata storage in the DPS and the SAS are completely independent and can be based on different solutions. DSS server can use local or NFS-mounted filesystem, Grid SE, sftp server, iRODS, Astro-WISE dataserer and can be extended for the existing or upcoming storage cases. The ability to extend DSS storage practically to any storage solution was introduced in the design to cope with the possible changes during the lifetime of the project and to support the SGS for at least 8 years after the launch.

The design of the SAS follows the latest generation of archives being developed by the ESDC, taking full advantage of the existing knowledge, expertise and code. The SAS builds on top of the latest ESDC's common Archives Building System Infrastructure (ABSI), which defines the common components to the latest ESA Science Archives (i.e. Gaia) within a three-tier modular architecture (client, server and data layer).

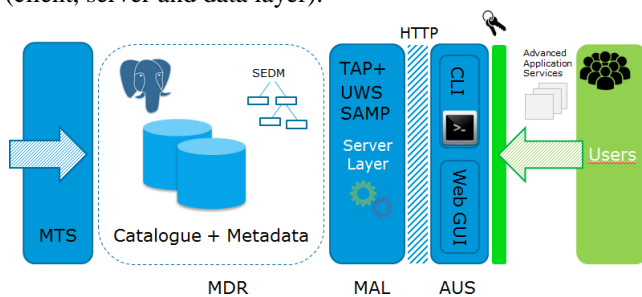


Figure 2: SAS Architecture Diagram

4. EAS INTERFACES

There are a number of interfaces and services which are implemented for EAS. EAS DPS will have to support interaction with SDC for the data processing (through Infrastructure Abstract Layer- IAL). IAL queries EAS DPS for commands for the data distribution and data processing and ingest metadata for newly created data products. The basic metadata browsing service allows to browse content of the metadata storage (according to the visibility scope for each user) and retrieve metadata for the data product in XML or other format.

DPS will be supplied with at least 2 other services – for the manual inspection of the quality of data products and for the retrieval and update of calibration frames.

The SAS will provide two ways of access through the Archive User Services (AUS), a web-based portal and a command line interface for programming access. The web-portal is based on the Google Web Toolkit technology [7]

and it has been designed to be easy to use and to provide a friendly interaction.

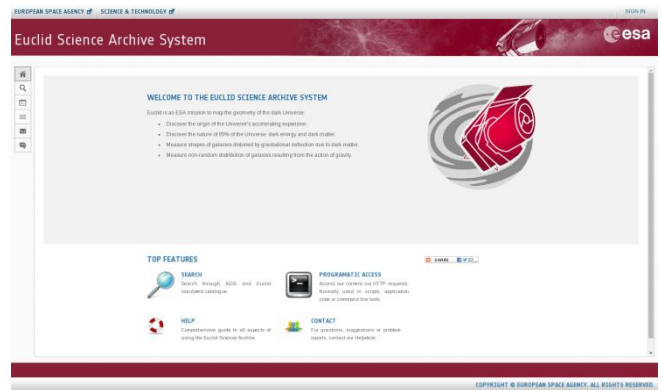


Figure 3: SAS Web User Interface

On the server side, the components are based on Java and it will integrate standard VO protocols to manage requests from the users.

At database level, the metadata repository will support the Science Exploitation Data Model (SEDM), that describes the metadata and catalogues oriented to scientific exploration. At the time of writing, the RDBMS is PostgreSQL join with pgSphere and Q3C modules that provides spherical data types, functions and operators to PostgreSQL.

Finally, the Metadata Transfer Service (MTS) will ingest the science metadata of the data release compliant with the SEDM from the ECDM at DPS, according to the policy of public releases defined by the Euclid Consortium.

The scientific requirements on the SAS mainly covers three different areas: parametric search for metadata and catalogues, data retrieval and visualization of images, spectra, etc. Regarding the maps visualization, it will be based on the technology developed for the ESA Sky [8], that allows the exploration of the astronomical resources using a useful and intuitive web interface. In addition, it will provide a set of tools to allow on-line research.

For the data retrieval, SAS stores the links pointing to the data products as part of its metadata. That links are managed by the DSS servers at the SDCs to deliver the data to the end users.

VO protocols play an important role within the architecture of the SAS, but also supporting added value tools integrated into the archive. Table Access Protocol (TAP+) [9] that provides efficient parametric search, has been developed for the Gaia Archive (Reference to Gaia BIDS Contribution), is also part of the SAS infrastructure and can be accessed from the Web interface.

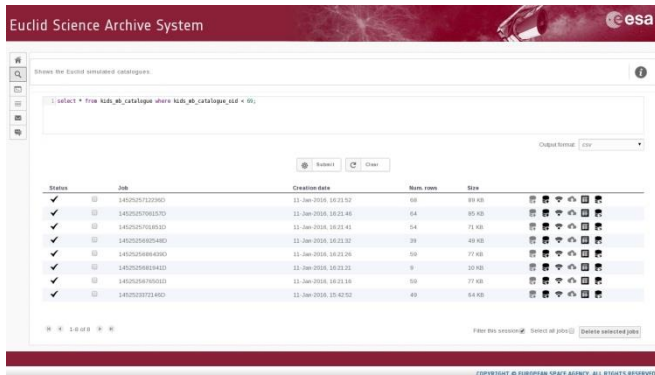


Figure 4: TAP Web Interface in SAS

VOSpace [10], the IVOA protocol for distributed data storage, is also integrated as part of the SAS infrastructure as an added value tool for the archive. It provides a storage abstraction layer and sharing capabilities transparent for the user.

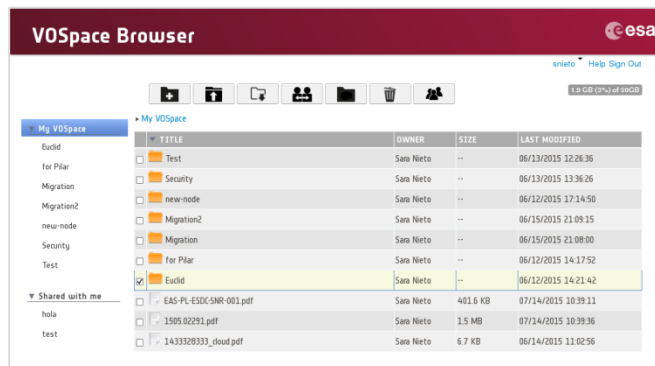


Figure 5: VOSpace Web User Interface

Other VO protocols like SAMP, will provide interoperability with astronomical analysis applications like Topcat among others.

To conclude, SAS will provide the tools and VO interfaces to enable the Software-to-Data Paradigm and "bring the software to the data" for the Euclid science.

4. EAS STATUS

The EAS Prototype was developed and tested in 2013 and 2014,. In 2015 the prototype formed the basis of the first version of EAS itself. First interfaces for the DPS and the DSS were released and tested during IT challenges

organized by Euclid Consortium producing and storing simulated Euclid data. We have successfully tested massive metadata ingestion for the data objects with extensive data lineage (KIDS DR1) [11].

In 2015 EAS team tested 2-copies configuration of EAS DPS metadata storage and metadata transfer between EAS DPS metadata storage mirrors in ESAC and Groningen.

5. FUTURE DEVELOPMENT

In the following years prior to the start of the mission the EAS will go through a number of crucial steps in the development including the final selection of an RDBMS. We plan as well to create a system which will support a dynamic data model: accommodating changes in ECDM without re-creation of the metadata storage scheme and migrating the data between different versions of the data model. Current DPS metadata storage will be separated on DPS metadata storage proper and DSS metadata storage.

6. REFERENCES

- [1] R. Laureijs et al., Euclid Definition Study Report, arxiv1110.3193, 2011
- [2] O.R.Williams et al., Data transmission, handling and dissemination of Euclid, Proceedings of NETSPACE 2014, 11, 2014
- [3] F. Pasian et al., Development Plans for the Euclid Science Ground Segment, Proceedings of ADASS XXIII, ASP Conf. Ser. 485, 505, 2014
- [4] K.Begeman et al., The Astro-WISE datacentric information system, Experimental Astronomy 35, 1, 2013
- [5] K.Begeman et al., LOFAR Information System, Future Generation Computer Systems 27, 319, 2011
- [6] Maurice P. et al., Euclid : Big Data From Dark Space, Big Data Spain, Oct. 2015
- [7] Dewsbury, R., Google Web Toolkit Applications, ISBN 978-0-321-50196-7, Prentice Hall, (2007).
- [8] Salgado, J., ESA Sky: A Tool to Discover Astronomical Data from ESA Missions, Big Data From Space 2016.
- [9] Dowler P. et at., IVOA Recommendation: Table Access Protocol Version 1.0, 2011arXiv1110.0497D
- [10] Graham M. et al., IVO Working Draft : VOSpace Version 2.1, WD-VOSpace-2.1-20150601
- [11] J.T.A. De Jong et al., Kilo Degree Survey, Experimental Astronomy 35, 24, 2013