

University of Groningen

## The Euclid Science Ground Segment Distributed Infrastructure: System Integration and Challenges

Frailis, Marco; Belikov, Andrey; Benson, Kevin; Bonchi, Andrea; Dabin, Christophe; Ealet, Anne; Fumana, Marco; Grenet, Catherine; Holliman, Mark; Maggio, Gianmarco

*Published in:*

**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:*

2019

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

*Citation for published version (APA):*

Frailis, M., Belikov, A., Benson, K., Bonchi, A., Dabin, C., Ealet, A., Fumana, M., Grenet, C., Holliman, M., Maggio, G., Maino, D., McCracken, H. J., Melchior, M., Piemonte, A., Polenta, G., Poncet, M., Scala, P. L., Serrano, S., & Williams, O. R. (2019). The Euclid Science Ground Segment Distributed Infrastructure: System Integration and Challenges. In M. Molinaro, K. Shortridge, & F. Pasian (Eds.), . (pp. 612-615). (ASP Conference series; Vol. 521). Astronomical Society of the Pacific.  
<https://ui.adsabs.harvard.edu/abs/2019ASPC..521..612F>

### 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 Science Ground Segment Distributed Infrastructure: System Integration and Challenges**

Marco Frailis,<sup>1</sup> Andrey Belikov,<sup>2</sup> Kevin Benson,<sup>3</sup> Andrea Bonchi,<sup>4,14</sup>  
Christophe Dabin,<sup>5</sup> Anne Ealet,<sup>6</sup> Marco Fumana,<sup>7</sup> Catherine Grenet,<sup>8</sup>  
Mark Holliman,<sup>9</sup> Gianmarco Maggio,<sup>1</sup> Davide Maino,<sup>10</sup>  
Henry J. McCracken,<sup>8</sup> Martin Melchior,<sup>11</sup> Antonello Piemonte,<sup>12</sup>  
Gianluca Polenta,<sup>4,14</sup> Maurice Poncet,<sup>5</sup> Paolo Luigi Scala,<sup>7</sup>  
Santiago Serrano,<sup>13</sup> and Owen Rees Williams<sup>2</sup>

<sup>1</sup>*INAF - Osservatorio Astronomico di Trieste, Trieste, Italy;*  
*frailis@oats.inaf.it*

<sup>2</sup>*University of Groningen, Groningen, Netherlands*

<sup>3</sup>*Mullard Space Science Laboratory, Dorking, United Kingdom*

<sup>4</sup>*Agenzia Spaziale Italiana Science Data Center, Roma, Italy*

<sup>5</sup>*CNES Toulouse Centre National d'Etudes Spatiales, Toulouse, France*

<sup>6</sup>*LAM - Laboratoire d'Astrophysique de Marseille, Marseille, France*

<sup>7</sup>*INAF - IASF Milano, Milan, Italy*

<sup>8</sup>*CNRS - Institut d'Astrophysique de Paris, Paris, France*

<sup>9</sup>*University of Edinburgh, Edinburgh, United Kingdom*

<sup>10</sup>*University of Milano, Milan, Italy*

<sup>11</sup>*University of Applied Science of Northwestern Switzerland, Windisch,  
Switzerland*

<sup>12</sup>*Max Planck Institute for Extraterrestrial Physics, Garching, Germany*

<sup>13</sup>*Institut d'Estudis Espacials de Catalunya, Barcelona, Spain*

<sup>14</sup>*INAF - Osservatorio Astronomico di Roma, Monte Porzio Catone, Rome,  
Italy*

**Abstract.** The Science Ground Segment (SGS) of the Euclid mission provides distributed and redundant data storage and processing, federating nine Science Data Centres (SDCs) and a Science Operations Centre. The SGS reference architecture is based on loosely coupled systems and services, broadly organized into a common infrastructure of transverse software components and the scientific data Processing Functions. The SGS common infrastructure includes: 1) the Euclid Archive System (EAS), a central metadata repository which inventories, indexes and localizes the huge amount of distributed data; 2) a Distributed Storage System of EAS, providing a unified view of the SDCs storage systems and supporting several transfer protocols; 3) an Infrastructure Abstraction Layer, isolating the scientific data processing software from the underlying IT infrastructure and providing a common, lightweight workflow manage-

ment system; 4) a Common Orchestration System, performing a balanced distribution of data and processing among the SDCs. Virtualization is another key element of the SGS infrastructure. We present the status of the Euclid SGS software infrastructure, the prototypes developed and the continuous system integration and testing performed through the Euclid “SGS Challenges”.

## 1. The Euclid Science Ground Segment

The Euclid SGS is composed of the ESA Science Operation Centre (SOC) and the Euclid Consortium SGS. The SOC is in charge of the mission planning, of the first consistency and quality checks and of the production of quick-look-quality data for public distribution. The Euclid Consortium SGS, composed of eight European SDCs and one USA SDC, is in charge of the instrument-related processing, production of science data products, simulations, ingestion of external data and, in general, of designing, developing, integrating and operating the scientific data processing.

The Euclid scientific data processing levels are decomposed into eleven Processing Functions (PFs), which are the highest-level break-down of the complete processing (Pasian et al. 2015). They are developed by distributed teams, with the constraint that each PF pipeline, except for external data processing, should run in any SDC.

The Euclid mission will generate a large amount of data: heavy simulations will be needed before flight operations, and several re-processing steps - from raw data up to the science products - will multiply the data volume by dozens. In addition, a large amount of external data will be gathered from ground-based observations. The current estimate of the total data volume at the end of the mission is about 90 Pbytes.

## 2. The SGS common infrastructure

The large amount of data generated during the Euclid mission, together with the multiple SDC organizations involved in the Euclid Consortium, leads to a distributed processing of the Euclid data, achieved through the following architecture key principles:

- Each SGS pipeline should run on any SDC
- Each SDC is both a processing and a storage “node”
- Separation of metadata (inventory) from data (storage)
- For lower processing levels, adoption of a “map/reduce” approach, where a pipeline is designed to process an independent quantum of locally available data

Following these principles, the SGS System Engineering team has defined a common infrastructure of loosely coupled software components for the development and running of the Euclid processing functions: 1) a **Common Data Model** (CDM), providing a central repository where all SGS components, interfaces and data structures are formalized in the XSD language; 2) the **EAS Data Processing System** (DPS), implementing the CDM metadata repository in a (relational) database management system and providing an abstract metadata access layer independent of the database implementation (Williams et al. 2019); 3) the **EAS Distributed Storage System** (DSS), supporting the copy, retrieval and movement of data files between SDCs and compatible with different storage solutions, including a GRID storage, an SFTP server or an iRODS

data system; 4) a **Common Orchestration System** (COORS) for the preparation, dispatch and monitoring of pipeline processing orders and performing a load balancing depending on the Euclid scanning strategy and the resources available in each SDC; 5) a **Monitoring and Control** system (M&C), tracking availability of computing resources and the Euclid data processing performance and progress; 6) the **Infrastructure Abstraction Layer** (IAL), which makes processing functions software isolated from the Euclid Archive, the SDCs computing infrastructure and the data storage and transfer services.

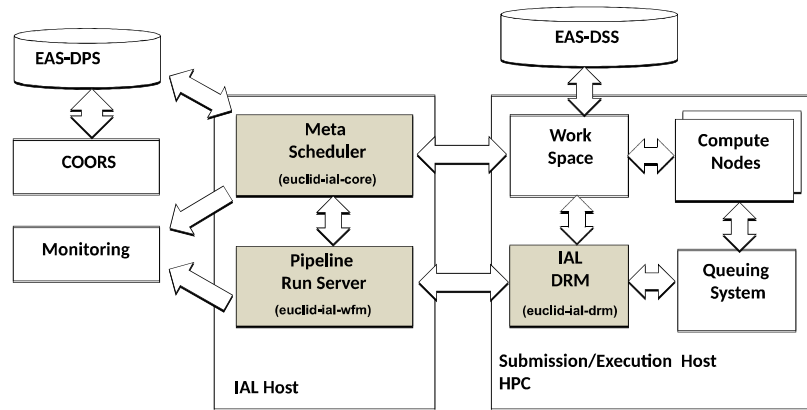


Figure 1. The architecture of the SGS common infrastructure

The IAL has three main components: a Meta Scheduler system, retrieving pipeline processing orders from the EAS and assuring that all needed pipeline inputs become locally available; a WorkFlow Manager (WFM), configuring and submitting jobs according to pipeline workflow definitions; a Distributed Resource Manager (DRM), providing an abstraction for the different batch queuing systems used by SDCs.

The IAL WFM translates pipeline specifications written in Python into a data flow graph. Then, it evaluates in which order to execute each task of the workflow and which of them can run in parallel. As soon as the last workflow task completes, the IAL retrieves the outputs and ingests them automatically in the EAS, together with their metadata.

In order to run each pipeline in any SDC and have reproducible results, another key element of the SGS infrastructure is the adoption of a reference Operating System with a set of reference software libraries and versions. Such a reference environment, called EuclidVM, is shared by all SDCs and deployed as a virtual node of the computing infrastructure using either virtualization technologies (e.g. KVM) or Linux Containers (e.g. Docker). The EuclidVM is kept as a lightweight system with the adoption of the CernVM-FS, a read-only network file system shared by all SDCs to distribute and update the SGS reference environment and processing functions (Poncet et al. 2019).

### 3. The SGS challenges

The assembly of such a complex system as the SGS infrastructure - developed by several distributed teams - requires a continuous and incremental integration and delivery process. This process has been implemented by the SGS as a set of planned Infrastruc-

ture and Scientific Challenges, taking into account the incremental development and maturity of the involved software components. The Infrastructure Challenges are dedicated to the integration and validation of the software components that are not “science dependent” and to their deployment in each SDC environment. Each Infrastructure Challenge is followed by a Scientific Challenge, focused on the integration and validation of the SGS processing functions.

Currently, Infrastructure Challenge 6 has been successfully performed, testing most of the SGS common infrastructure (EAS-DPS, EAS-DSS, IAL, M&C, CernVM-FS, etc.). It has been followed by Scientific Challenge 2, which has involved the integration of four processing functions: **SIM**, which produces all the simulated data necessary to validate the data processing stages, and to calibrate observational or method biases; **VIS** and **NIR**, producing fully calibrated photometric exposures respectively from the visible imaging and the near-infrared imaging of Euclid (da Silva et al. 2019); and **SIR**, producing fully calibrated 1D spectra extracted from the NISP spectroscopic exposures. In particular, SIM has produced VIS, NISP-P and NISP-S raw dither exposures, introducing the requested features and instrumental effects. Then VIS, NIR and SIR proto-pipelines have been successfully deployed and run in the SDCs computing infrastructure, processing the simulated data and ingesting the obtained output products in the EAS archive.

#### 4. Conclusions

The Euclid Science Ground Segment poses new challenges in terms of data volume and processing on geographically distributed Science Data Centres. The integration and validation of working prototypes is performed by periodic IT and Scientific end-to-end tests (SGS challenges) in order to assess both the architectural choices and their implementation. Currently, a first validation of the SGS common infrastructure, the production of simulations and the lower levels of the SGS scientific data processing has been performed. Next SGS challenges will involve the subsequent processing levels, from the production of source catalogs containing consistent photometric and spectroscopic measurements up to the high-level science products.

**Acknowledgments.** The authors acknowledge the Euclid Consortium, the European Space Agency and the support of a number of agencies and institutes that have supported the development of Euclid. A detailed complete list is available on the Euclid web site (<http://www.euclid-ec.org>).

#### References

- da Silva, R., et al. 2019, in ADASS XXVI, edited by Molinaro, M. and Shortridge, K. and Pasian, F. (San Francisco: ASP), vol. 521 of ASP Conf. Ser., 311
- Pasian, F., et al. 2015, in ADASS XXIV, edited by A. R. Taylor, & E. Rosolowsky, vol. 495 of ASP Conf. Ser., 207
- Poncet, M., Le Boulc’h, Q., & Holliman, M. 2019, in ADASS XXVI, edited by M. Molinaro, K. Shortridge, & F. Pasian (San Francisco: ASP), vol. 521 of ASP Conf. Ser., 588
- Williams, O. R., et al. 2019, in ADASS XXVI, edited by M. Molinaro, K. Shortridge, & F. Pasian (San Francisco: ASP), vol. 521 of ASP Conf. Ser., 120