



ASTERICS - H2020 - 653477

Data Format Survey

ASTERICS GA DELIVERABLE: D3.3

Document identifier:	ASTERICS-D3.3-draft-V2.0.pdf
Date:	26-04-2016
Work Package:	OBELICS – WP3
Lead Partner:	UCM
Document Status:	INCOMPLETE DRAFT
Dissemination level:	INTERNAL
Document Link:	https://www.asterics2020.eu/dokuwiki/lib/exe/fetch.php?media=intra:wp3:task3.2:ASTERICS-D3.3-draft-V2.0.pdf

Abstract

ASTERICS will benefit ESFRI projects and other related major research infrastructures, including ESFRI-precursor experiments. This survey seeks synergies between these experiments with the aim of fostering the use of open data standards that would enable interoperability and software re-use.

I. COPYRIGHT NOTICE

Copyright © Members of the ASTERICS Collaboration, 2015. See www.asterics2020.eu for details of the ASTERICS project and the collaboration. ASTERICS (Astronomy ESFRI & Research Infrastructure Cluster) is a project funded by the European Commission as a Research and Innovation Actions (RIA) within the H2020 Framework Programme. ASTERICS began in May 2015 and will run for 4 years.

This work is licensed under the Creative Commons Attribution-Noncommercial 3.0 License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc/3.0/> or send a letter to Creative Commons, 171 Second Street, Suite 300, San Francisco, California, 94105, and USA. The work must be attributed by attaching the following reference to the copied elements: “Copyright © Members of the ASTERICS Collaboration, 2015. See www.asterics2020.eu for details of the ASTERICS project and the collaboration”. Using this document in a way and/or for purposes not foreseen in the license, requires the prior written permission of the copyright holders. The information contained in this document represents the views of the copyright holders as of the date such views are published.

II. DELIVERY SLIP

	Name	Partner/WP	Date
Author(s)	Jaime Rosado and José Luis Contreras	UCM/WP3	21/04/2016
Amendments	Tammo Jan Dijkema, Tamas Gal, Dominique Boutigny	ASTRON/WP3 FAU/WP3 CNRS/WP3	26/04/2016
Reviewed by			
Approved by			

III. DOCUMENT LOG

Issue	Date	Comment	Author/Partner
v0.5	01/04/2016	Incomplete draft	J. Rosado/UCM
v1.2	21/04/2016	Complete draft	J. Rosado/UCM
V2.0	26/04/2016	Revised draft	J. Rosado/UCM

IV. APPLICATION AREA

This document is a formal deliverable for the GA of the project, applicable to all members of the ASTERICS project, beneficiaries and third parties, as well as its collaborating projects.

V. TERMINOLOGY

AMS	Alpha Magnetic Spectrometer
ANTARES	Astronomy with a Neutrino Telescope and Abyss environmental RESearch
ASTERICS	Astronomy ESFRI & Research Infrastructure Cluster
CASA	Common Astronomy Software Applications
CTA	Cherenkov Telescope Array
CTDS	Casacore Table Data System
DADI	Data Access, Discovery and Interoperability
E-ELT	European Extremely Large Telescope
EGO	European Gravitational Observatory
ESFRI	European Strategy Forum on Research Infrastructures
ESO	European Southern Observatory
e-VLBI	Electronic Very-Long-Baseline Interferometry
EVN	The European VLBI Network
FITS	Flexible Image Transport System
GW	Gravitational wave
HAWC	High-Altitude Water Cherenkov Observatory
HDF	Hierarchical Data Format
H.E.S.S.	The High Energy Stereoscopic System
IACT	Imaging Atmospheric Cherenkov Telescope
IGWD	Interferometric Gravitational Wave Detector
JIVE	Joint Institute for VLBI in Europe
KM3NeT	Cubic Kilometre Neutrino Telescope
LIGO	Laser Interferometer Gravitational Wave Observatory

LOFAR	The Low Frequency Array
LSST	The Large Synoptic Survey Telescope
MAGIC	Major Atmospheric Gamma-Ray Imaging Cherenkov
OBELICS	OBservatory E-environments Llnked by common ChallengeS
SKA	The Square Kilometre Array
UCM	Universidad Complutense de Madrid
VO	Virtual Observatory

A complete project glossary is provided at the following page:

<https://www.asterics2020.eu/glossary>

VI. EXECUTIVE SUMMARY

ASTERICS is intended to give services to research infrastructures identified in the ESFRI Roadmap (CTA, KM3NeT, SKA and E-ELT) as well as other major international projects, including precursor experiments, in the area astronomy, astrophysics and astroparticle physics. ASTERICS looks towards enabling interoperability between them, encouraging cross-fertilisation and developing joint multi-wavelength/multi-messenger capabilities. To do that, one of the primary targets is to seek commonalities in data handling, data storage, etc. between those projects. This initiative also goes in the direction of the global movement for open science data.

Within ASTERICS, the work package OBELICS aims to enable interoperability and software re-use for the data generation, integration and analysis of the ASTERICS ESFRI and pathfinder facilities. Therefore, one of its main priorities is to establish open standards and software libraries for multi-wavelength and multi-messenger data. This document is a deliverable from the work package OBELICS that has the objective of seeking synergies between experiments that could enable the use of appropriate common data formats. For this purpose, experiments have been grouped into three different categories according to the type of data that they produce: image-based, event-based and signal-based experiments.

As a result of this survey, two possibilities are identified, mainly for event-based experiments: defining standards for low-level data and pushing for common high-level data formats. For the second one, the UCM-ASTERICS group has proposed extending a recent initiative for common data formats within the gamma-ray community to encompass neutrino and cosmic-ray observatories as well.

Table of contents

I. COPYRIGHT NOTICE	1
II. DELIVERY SLIP.....	1
III. DOCUMENT LOG	1
IV. APPLICATION AREA	2
V. TERMINOLOGY	2
VI. EXECUTIVE SUMMARY	3
Table of contents.....	4
1. Introduction	5
2. ESFRI projects and related pathfinders	5
3. Types of data produced by experiments	6
4. Image-based experiments	7
5. Event-based experiments	8
6. Signal-based experiments.....	10
7. Possible synergies	11
8. Conclusions	12
References.....	13

1. Introduction

ASTERICS will benefit ESFRI projects and other related major research infrastructures, including ESFRI-precursor experiments. This survey seeks synergies between these experiments with the aim of fostering the use of open data standards that would enable interoperability and software re-use.

The projects included in this survey are identified in [section 2](#). In [section 3](#), the experiments are grouped in three categories according to the type of data that they produce. [Sections 4 to 6](#) describe the data formats that are in use or envisaged for the experiments belonging to each group. The possible synergies between experiments are discussed in [section 7](#). Conclusions are given in [section 8](#).

2. ESFRI projects and related pathfinders

This data format survey encompasses 13 projects identified as ASTERICS stakeholders. They include ESFRI projects, their pathfinders and other linked projects, most of them being ASTERICS-OBELICS partners. The table shows the field of application of each experiment and the type of data produced as categorized in [section 3](#).

Project	ESFRI	Field	Type of data
CTA	Yes	Cherenkov observatories for gamma-ray astronomy	Events
H.E.S.S.	Pathfinder for CTA	Cherenkov telescope array for gamma-ray astronomy	Events
MAGIC	Pathfinder for CTA	Cherenkov telescope array for gamma-ray astronomy	Events
KM3NeT	Yes	Neutrino telescope	Events
IceCube	Pathfinder for KM3NeT	Neutrino telescope	Events
ANTARES	Pathfinder for KM3NeT	Neutrino telescope	Events
E-ELT	Yes	Ground-based optical/near-infrared telescope	Images
LSST	No	Optical telescope	Images
EUCLID	No	Satellite mission to map the dark Universe	Images
SKA	Yes	Radio telescope arrays	Signals

e-EVN	Pathfinder for SKA	e-VLBI network for radio astronomy	Signals
LOFAR	Pathfinder for SKA	Radio interferometric array	Signals
Advanced LIGO	No	Gravitational wave detectors	Signals
Advanced Virgo	No	Gravitational wave detector	Signals

Table 1: Projects included in this data format survey. Experiments are classified according to the type of data that they produce.

3. Types of data produced by experiments

In this survey, experiments are classified according to the type of data that they produce. Although final science products are similar for all the astro(particle)physics experiments (e.g., skymaps, catalogues, spectra), raw and processed data below the science data level have fundamental differences depending on the experimental technique. In this way, imaging atmospheric Cherenkov telescopes (IACT) for gamma-ray astronomy differ fundamentally from visible-light astronomical ones in the fact the formers detect particles one by one, the so-called “events”, whereas optical telescopes collect light to form images. Moreover, both detection techniques are very different to those employed by radio interferometers or by gravitational wave detectors, which record signals in the time (or frequency) domain to be subsequently processed. On this basis, three groups of experiments are identified: image-based, event-based and signal-based experiments.

The above types of data have different requirements for formatting and processing. This also applies to a significant part of the metadata needed to find the stored data and to link them to calibration and other auxiliary instrument data. Synergies between experiments producing similar type of data are more natural and easier to identify. Therefore, software reuse and common standards are expected to be attained to much more extent in experiments belonging to the same group. In the next three sections, these groups of experiments are reviewed and their corresponding data formats (both in use and envisaged) are discussed.

It should be noted, however, that these groups are not completely bounded and often overlap each other. This leaves open the possibility to look for common solutions for a wider range of experiments.

reconstruction and event filtering is done in near real time so that data can be transferred over a satellite link to the data centres located in the Northern Hemisphere [17, 26].

Data from neutrino experiments are organized hierarchically and have several levels of complexity in a similar way to IACT ones. In general, data produced by event-based observatories share these features whichever the type of particle that they detect.

The raw data rates of the above experiments are listed in [table 2](#). Succinct information on the formats used for low and mid level data is also given in the table. Custom data formats are commonly used for low-level data, since they depend on the particularities of each experiment. The analysis software of most current event-based experiments is developed in the ROOT framework [27], which was originally designed for particle physics data analysis in CERN projects. Hence, ROOT-based data formats are used in the processing chain of H.E.S.S. [28], MAGIC [29] and ANTARES [24]. The archival and data exchange system of IceCube creates a metadata XML file for each data file following a specification called “DIF Plus” [17] compliant with the format required by the Office of Polar Programs for all experiments funded by the National Science Foundation. Then, data are processed in the data centers using a custom analysis framework.

Observatory	Raw data rate	Low and mid level data format
CTA	10 GB/s	To be determined
H.E.S.S.	46 MB/s	Based on ROOT
MAGIC	100 MB/s	Raw data stored in binary form and translated to ROOT format at preprocessing level
KM3NeT	9 GB/s (phase 2)	To be determined
IceCube	12 MB/s (total) 1.2 MB/s (satellite)	Custom format and analysis software
ANTARES	0.3 – 10 GB/s	Based on ROOT

Table 2: Data rates and formats used by the event-based experiments included in this report.

Following the path defined by the Fermi experiment [30], efforts are being made within the gamma-ray community to define unified standards for higher level data ready for distribution [31, 32]. The proposed standards are based on the widespread FITS file format, which is not only used for images, but also for a large variety of scientific data. In particular, the CTA data level 3 (DL3) is identified as the lowest high-level data to be delivered to basic users. It will consist in FITS files containing lists of well reconstructed gamma-like events along with associated instrumental response characterizations and any necessary technical data for science analysis [23]. DL3 is the first data level being nearly independent of the particularities of the detector and thus the most suitable for science analysis.

In addition, current IACT observatories are integrating themselves into the VO framework and already making available their final science products and some legacy observatory data (e.g., light curves, spectra and source catalogues) in formats compliant with VO tools [33, 34]. The final goal is to allow scientists to do multi-messenger astronomical studies by combining data from different facilities in a single analysis. In fact, CTA is intended to work as a “standard”, fully VO compliant astronomical observatory.

Neutrino observatories are also making their high-level data available. ANTARES has made public an ASCII table containing the data set for its 2007-2010 search for cosmic neutrino point sources [35]. The IceCube collaboration has also built a web page from which several data sets associated to journal papers can be downloaded [36], where event lists and most auxiliary data (e.g., background distributions) are provided in ASCII tables, whereas some other data are in HDF5 format [37]. Remarkably, IceCube also provides lower level data (i.e., calibrated data without reconstructions) in an open format based on python along with appropriate analysis tools [38]. Both ANTARES and IceCube (and later KM3NeT) participate in multi-messenger networks by sending alert messages in VOEvent format, as will be discussed in [section 7](#).

6. Signal-based experiments

In radio interferometric arrays, signals from all the individual telescopes are brought together and processed by the so-called correlator, which combines the signals to form an image of the observed radio source. Among the projects listed in [table 1](#), SKA [39], EVN [40] and LOFAR [41] belong to this category.

The European VLBI Network (EVN) comprises 21 very large and sensitive reflector telescopes (dishes) spread throughout Europe and beyond. In addition to “traditional” VLBI sessions in which the data are first recorded at the telescopes on tapes or discs and then physically delivered to the processor, the EVN consortium is conducting a development programme called e-EVN, which aims at creating an e-VLBI network where data is transferred and processed in real time. All data correlated at the EVN Data Processor at JIVE (raw FITS format data, processed images, calibration tables, etc.) are placed into the EVN archive and made public once the 12-month proprietary period has expired [40].

The LOFAR array, on the other hand, consists of about 7000 simple omni-directional antennas organised in stations containing local computing resources to perform beam-forming [42]. All these beam-formed data (about 19 GB/s for the entire array) are sent via a high-speed fibre network to the central processing facility, where they are pre-processed on-line. In interferometric imaging mode, pre-processed data are stored in the MeasurementSet format based on the Casacore Table Data System (CTDS) [43], which is part of the CASA package for radio-astronomical data processing [44]. On the other hand, the HDF5 format is used for pulsar processing. Then, several off-line processing steps are performed to obtain

either images or pulsar data products in PRESTO format [45]. The final scientific data are transferred to the LOFAR long-term archive for cataloguing and distribution to the community [46]. This archive is expected to grow by up to 5 PB per year.

The above two projects are recognised as pathfinders for SKA, which will operate both an array of dishes and an array of antennas grouped into stations. The SKA project is designed in two phases and, when both are complete, the observatory will consist of many thousands of connected radio telescopes. For the first phase, the summed data rate is estimated to be about 3 TB/s and the expected archive data volume is 100 PB per year [47]. The SKA project will profit from the lessons learned on data management in existing radio interferometric arrays; although the file formats to be used to store data are still to be defined. Radio interferometric arrays are also integrating into the VO network and SKA data will be delivered in formats compliant with VO standards (see, e.g., [48]).

The two Advanced LIGO detectors in EEUU [49] and the Advanced Virgo detector in Italy [50] included in [table 1](#) are the first ones of a network of very sensitive interferometric detectors of gravitational waves (GW) working together. These detectors also produce signal-based data, but they are very different to those generated by radio interferometric arrays. The data produced by these experiments are stored in “frames” with thousands of channels, where the GW strain channel only represents a small fraction of data and all the other channels are used for auxiliary instrumental and environmental monitoring. Each interferometer has a data rate of a few tens of MB/s [51, 52].

A standard data format, called IGWD Frame format, was already established by a LIGO-Virgo agreement in 1997 [53]. Both projects will make data publicly available in IGWD Frame format as well as in other standard formats easier to use for most open data users (e.g., HDF5). In addition, they plan to participate in the VOEvent network by communicating real-time alerts to follow-up observers.

7. Possible synergies

One of the main goals of ASTERICS is to foster common standards and software re-use between ESFRI and other relevant research infrastructures. This is more attainable for the highest-level data, since they do not rely on the detection technique. Prominently, the VO framework allows the wide scientific community to access astronomical data from different facilities in a single transparent system. Originally formed by optical observatories, the VO network now joins a growing and diverse list of experiments with the purpose of enabling multi-wavelength and multi-messenger science. Especially, science-ready data products from gamma-ray observatories are already being made accessible from VO tools.

Other astroparticle observatories are less integrated into the VO network, although they are developing systems to communicate real-time alerts to follow-up observers using the

VOEvent protocol within the VO framework. In particular, the Astrophysical Multi-messenger Observatory Network (AMON) seeks to perform a real-time correlation analysis of the high-energy signals across all known astronomical messengers, i.e., photons, neutrinos, cosmic rays and gravitational waves [54].

For gamma-ray observatories, the first high level data not depending on the particularities of the experiment contains lists of selected gamma-like reconstructed events as well as the instrument response functions and other relevant data for science analysis. In CTA, this data level is known as DL3 and basic users will have access to it in a similar way to how Fermi data is distributed [30]. A new initiative has been created to describe data formats that are in use for this data level in the gamma-ray community and to define a unified standard based on FITS to be adopted by current and future experiments, chiefly IACT observatories [31].

It can easily be seen that the above unified format for gamma-ray astronomy might also be used for DL3-equivalent data produced by other event-based astroparticle detectors, as neutrino and cosmic-ray observatories. Note that event lists basically contain information on the energy and direction of the particle, and therefore, it would only be required to add information about the nature of the particle. The effective area (or alternatively the exposure) and any other relevant high-level technical data should be linked to time periods short enough to neglect instrument response variations, regardless of the observation mode (e.g., source pointing runs or wide-aperture continuous observation). The UCM-ASTERICS group has brought this proposal to a recent meeting to work on this DL3 data format [55] and it is presently under study.

For lower level data, synergies between experiments are most likely limited to those belonging to the same category. For example, standard data formats and analysis software exist in radio interferometry. Interestingly, the GW community has already established a common data format and interoperability between the several detectors of the GW network is ensured. Most current event-based experiments use ROOT-based data formats and analysis tools, but no standards have been defined yet. Nevertheless, it could be conceived some standards for the hierarchical structure of data produced by event-based experiments. In this way, operational elements like “telescope” or “optical module” can be seen as instances of hierarchical levels of a more generic class structure. As an alternative to existing ROOT-based formats, HDF5 would be suitable for this kind of data and available HDF5 tools (e.g., viewers) would be of great help.

8. Conclusions

More and more experiments are making their science-ready data products publicly available in formats compliant with the VO standards. On the other hand, although big steps have been made toward multi-messenger astronomy, progresses are still to be achieved to enhance the integration of astroparticle observatories into the VO network. Within

ASTERICS, the work package DADI is contributing to the development and implementation of the VO framework by organising technology forums and training events [56].

A unified standard based on the FITS file format is currently being developed for the so-called DL3 data in the gamma-ray community. The UCM-ASTERICS group has proposed to extend this standard to other event-based experiments (i.e., neutrino and cosmic-ray observatories) and potentially to gravitational wave detectors too. This proposal will be explored trying to involve staff of these experiments. The data distribution in open formats would bring returns to the experiments, e.g., through a larger impact of its results, externally developed software and cross-check analyses.

The use of common formats for low-level data is less clear, since this data depends on the particularities of each experiment. Nevertheless, possible synergies between experiments producing similar type of data are found. In particular, low and mid level data generated by all event-based experiments have a hierarchical structure and therefore hierarchical formats as HDF5 are suitable for them. The UCM-ASTERICS group plans to propose and test a HDF5 format for raw data of CTA.

References

- [1] <https://www.eso.org/sci/facilities/eelt/>
- [2] <http://www.lsst.org/>
- [3] <http://www.euclid-ec.org/>
- [4] M. Jurić et al., The LSST Data Management System. Available at <http://arxiv.org/abs/1512.07914>
- [5] E-ELT Programme – Observatory Top Level Requirements. Available at <https://www.eso.org/sci/facilities/eelt/>
- [6] Euclid Definition Study Report. Available at <http://arxiv.org/abs/1110.3193>
- [7] M. Jurić et al., The LSST Data Products Definition Document. Available at <https://docushare.lsstcorp.org/docushare/dsweb/Get/LSE-163/>
- [8] http://fits.gsfc.nasa.gov/fits_home.html
- [9] W.D. Pence et al., Definition of the Flexible Image Transport System (FITS), version 3.0, A&A 524 (2010) A42.
- [10] <http://www.ivoa.net/>

- [11] VOTable Format Definition. Available at <http://www.ivoa.net/Documents/latest/VOT.html>
- [12] VOEvent: Sky Event Reporting Metada. Available at <http://www.ivoa.net/documents/VOEvent/>
- [13] <https://www.cta-observatory.org/>
- [14] <https://www.mpi-hd.mpg.de/hfm/HESS/>
- [15] <https://magic.mpp.mpg.de/>
- [16] <http://www.km3net.org/home.php>
- [17] <https://icecube.wisc.edu/>
- [18] <http://antares.in2p3.fr/>
- [19] www.hawc-observatory.org/
- [20] <http://fermi.gsfc.nasa.gov/>
- [21] <http://www.ams02.org/>
- [22] <https://www.auger.org/>
- [23] G. Lamanna et al., Cherenkov Telescope Data Management, PoS (ICRC2015) 947. Available at: <http://arxiv.org/abs/1509.01012>
- [24] J.A. Aguilar et al., The data acquisition system for the ANTARES neutrino telescope, NIM A 570 (2007) 107.
- [25] Letter of Intent for KM3NeT2.0. Available at: <http://arxiv.org/abs/1601.07459>
- [26] R. Abbasi et al., The IceCube data acquisition system: Signal capture, digitization, and timestamping, NIM A 601 (2009) 294.
- [27] <https://root.cern.ch/>
- [28] A. Balzer et al., The H.E.S.S. central data acquisition system, Astropart. Phys. 54 (2015) 67.
- [29] J. Aleksić, Optimized Dark Matter Searches in Deep Observations of Segue 1 with MAGIC, PhD Thesis (2013).
- [30] <http://fermi.gsfc.nasa.gov/ssc/data/>

- [31] <http://gamma-astro-data-formats.readthedocs.org/en/latest/>
- [32] <https://www-zeuthen.desy.de/multi-messenger/GammaRayData/>
- [33] <http://magic.pic.es/> (see section “Public area”).
- [34] S. Derrière et al., Using the Virtual Observatory: multi-instrument, multi-wavelength study of high-energy sources. Journées de la SF2A 2014. Available at <http://arxiv.org/abs/1411.6514>
- [35] <http://antares.in2p3.fr/publicdata.html>
- [36] <https://icecube.wisc.edu/science/data>
- [37] <https://www.hdfgroup.org/HDF5/>
- [38] <http://icecube.umd.edu/PublicData/>
- [39] <https://www.skatelescope.org/>
- [40] <http://www.evlbi.org/>
- [41] <http://www.lofar.org/>
- [42] M.P. van Harlem et al., LOFAR: The LOw-Frequency ARray, A&A 556 (2013) A2.
- [43] G.N.J. van Diepen, Casacore Table Data System and its use in the MeasurementSet, A&C 12 (2015) 174.
- [44] <https://casa.nrao.edu/>
- [45] <https://prestodb.io/>
- [46] <http://www.astron.nl/radio-observatory/astronomers/technical-information/lofar-technical-information>
- [47] SKA1 System Baseline v2. Available at <https://www.skatelescope.org/key-documents/>
- [48] <http://amiga.iaa.es/>
- [49] <http://www.ligo.org/index.php>
- [50] <http://www.virgo-gw.eu/>
- [51] LIGO Data Management Plant. Available at <https://dcc.ligo.org/LIGO-M1000066/public>

[52] F. Acernese et al., Advanced Virgo: a second-generation interferometric gravitational wave detector, *Class. Quantum Grav.* 32 (2015) 024001.

[53] Specification of a common data frame format for interferometric gravitational wave detectors. Available at <http://dcc.ligo.org/cgi-bin/DocDB/RetrieveFile?docid=329>

[54] <http://amon.gravity.psu.edu/index.shtml>

[55] https://github.com/open-gamma-ray-astro/2016-04_IACT_DL3_Meeting/blob/master/README.md

[56] <https://www.asterics2020.eu/dokuwiki/doku.php?id=open:wp4:start>