

## Archive Prototypes in the Cherenkov Telescope Array : INAF Contributions

### Stefano Gallozzi<sup>1</sup>, Fabrizio Lucarelli<sup>2</sup> and Angelo Antonelli<sup>1,2</sup> on behalf of the Whole CTA Rome Archive Team

<sup>1</sup>INAF, Astrophysical Observatory of Rome <sup>2</sup>SSDC, ASI- Space Science Data Center



Stefano Gallozzi

### **The ASTRI Data Management Team**



I.Bi.S.Co.



Astri/CTA Data Challeng

Matteo Perri

F. G. Saturni

INAF **ISTITUTO NAZIONALE DI ASTROFISICA** NATIONAL INSTITUTE FOR ASTROPHYSICS

**[**]

Antonio Stamerra Vincenzo Testa



Teche.it

Martina Cardillo



Antonino D'Aì



## How much BIG are DATA?

Without data compression and assuming 165 operational nights/yr:

ASTRI/Prot.	$\rightarrow$	~0.8 TB/night
	$\rightarrow$	~0.3 PB/year
Mini-Array	$\rightarrow$	~3 TB/night
	$\rightarrow$	~6.1 TB/night <u>A.R.</u>
	$\rightarrow$	~1.0 PB/year <u>A.R.</u>

... and for CTA ?

A pessimistic scenario can involve more than 100PB/year !

 $\rightarrow$  triggering systems?

 $\rightarrow$  data reduction on site?

or...

ightarrow challenging with big data



Once (period 2012-2016) the Archive's principal requirement was: CTA Archive system must store, manage, preserve and provide easy access to a such huge amount of data for a long time.

Stefano Gallozzi



#### Data Archives

Data archives are central to astronomy today, and their importance continues to grow. The science impact of these archives is large and increasing rapidly. Papers based on archival data from the Hubble Space Telescope now outnumber those based on new observations in any year and include some of the highest-impact science from the HST, as shown in Figures 5.6 and 5.7. Data from the 2 Micron All

Stefano Gallozzi

### Archive is not a... ...simple "repository"!!!



Stefano Gallozzi

#### Archive role is "central" for...



### ...an Astronomical Observatory

In the scientific data lifecycle of any **OBSERVATORY** the role of the Archive is central.

The major aim of a Scientific Archive is to guarantee **data preservation and access** information for the **Long Term** and for **all data science products**.

The **archived information** must be also usable by **different user categories** (*data consumers*) who are separate in time, space and background from the *data producers*.

Archive **MUST** be accessible well beyond the end of the operational life of the observatory

Stefano Gallozzi





Stefano Gallozzi





Stefano Gallozzi



**CTA Collaboration & Community** participated to the INDIGO-Data Cloud H2020 Project AS "**Use Case**" for the **INDIGO infrastructure**.

The aim of our commitment was the very fruitful multi-disciplinar collaboration with INDIGO Communities in order to include the **BigData challenges** coming from the CTA Archive as an INTERNAL INDIGO Use CASE / Case Study  $\rightarrow$  to be investigated with a distributed approach  $\leftarrow$ **BIG-DATA Archive** direct access to CTA cloud **Grid Providers** CTA Storage still needed? Scientist / End-users Broke ONECATA **ONEDAT** i2 ONEDATA LHC Grid Scheduler (+VOMS) CTA end USER INDIGO solution effort INDIGO - DataCloud In the Distributed Federation of Storage OneDATA solutions are ready for CTA A&A

Stefano Gallozzi





 $\begin{array}{ll} \underline{Prototype} \rightarrow & \text{Onedata's REST API's as well as oneclient command line tool for} \\ & \text{mounting virtual Onedata filesystem on the local machine} \end{array}$ 

Stefano Gallozzi



### Archive Prototype testbed ... (REAL+SIMUL ASTRI&CHEC DATA)



Currently working the Archive Prototype Solution using: → the ASTRI camera real data



• → the INAF-PRIN ASTRI CTA Data Challenge (AC-DC) for mini-array based simulation



### $\rightarrow$ the CHEC-M Camera real data

since in last weeks CHEC-M Camera was hosted to the ASTRI prototype design and The ASTRI Archive is going to manage and store CHEC real data. Three CTA SST: GATE (left), ASTRI (centre), SST-1M (right).



Stefano Gallozzi



Using the end-to-end ASTRI Archive System as data feeder



Stefano Gallozzi



### Archive Prototype testbed



• **CTA clients & End Users** able to access to the CTA federated storage cloud through One Data interfaces (clients) + Users A&A.

**OneClient:** a command-line based application for **accessing** and **managing** user spaces (mounted in local FS) via **virtual file system**.



Stefano Gallozzi



### Archive Prototype testbed



Stefano Gallozzi



### Archive Prototype testbed



Stefano Gallozzi

## **ASTRI Archive System**



Stefano Gallozzi

## **Physical & Logical Archives**



Stefano Gallozzi

## **Knowledge Discovery DBs**

### ...testing databases technology using ASTRI data-model...



Stefano Gallozzi

### **ASTRI Gateway**



nance to view observability of the following list of Tarmet Fields in 2019/05/31 at S.L.N. Catania (Italy)

#### Stefano Gallozzi

### **ASTRI User Access #1**



Stefano Gallozzi



#### Stefano Gallozzi



## **ASTRI Hardware**

- Switch T.O.R.  $1 \rightarrow 10$ Gb/s
- Other OneData nodes???
- OneData Provider#3 (@SSDC)
- OneData Provider#2 (@LNF)
- **OneData Provider#1 (@OAR-MPC)**
- Redundant ASTRI services
- ASTRI Gateway (gitLab & redmine)
- **Mirror ASTRI Gateway**
- Switch KVM
- Other Computing???
- DBs Service and File Catalogs <u>& Pipeline</u> (devel & Runtime)
- STORAGE (expanding ~1PB) Other Storage???
- <u>UPS</u> downstream and stabilized by the Institute's UPS



### ASTRI-miniArray environment is ready for both solutions even centralized data centers on site (incremental archive)

Stefano Gallozzi

### **QUESTIONS?**



Stefano Gallozzi

### **BACK-UP SLIDES**



S.Gallozzi 14<sup>th</sup>-18<sup>th</sup> May 2018 – C.T.A. General Meeting @ Paris (FRA)

# **OneData Overview**

**OneData** system **virtualizes** storage systems provided by storage resource providers **distributed** globally.

The most important concepts of the platform are:

- **Spaces** distributed virtual volumes, where users can organize their data
- **Providers** entities who support spaces with actual storage resources
- **Zones** federations of providers which enable creation of closed or interconnected communities.



# CTA OneZone

- OneZone is the gateway for users to the OneData system. It is responsible for connecting to the authentication and authorization infrastructure.
- It allows users to:
- ✓ create user spaces
- ✓ generate space support tokens, that can be used to support user spaces with storage from a dedicated storage provider
- monitor availability of storage providers that support user spaces
- ✓ see the geographical distribution of storage providers
- choose storage provider for spaces



## CTA OneZone



# CTA OneProvider(s)



# CTA OneProvider(s)

0)	ONEJATA					¥.	₽	s⊊® 8	•	4	1		<u> </u>	🚺 admin +
Ô	CTADATASPACE													
Data	CTADATASPACE	$\checkmark$	FILES										SIZE	MODIFICATION
Shared	Root directory		astri_000_41_001_00001_R_000004_000_1	002.lv0									916.88 KB	2017-01-13 12:01
$\bigcirc$			astri_000_41_001_00001_R_000005_000_1	002.lv0									916.88 KB	2017-01-16 11:01
Groups				File distribution		×	l							
کی ا				Distribution of file astri_000_41_001_0	blocks among providers for file 00001_R_000005_000_1002.lv0 Eile blocks									
				SSDC	0	916.88 KB								
Providers				LNF	0	916.88 KB								
				МРС	0	916.88 KB								
					Close									

# **CTA OneClient**

- OneClient is a command-line based application for accessing and managing user spaces via virtual file system.
- User spaces are **mounted** in the local file system tree (i.e. in a Grid Storage-Element FS as well).



## Metadata

Metadata in OneData are organized into 3 levels:

- ✓ Filesystem attributes basic metadata related to file system operations such as file size, creation and modification timestamps, POSIX access rights, etc.,
- Extended attributes these attributes enable assigning custom key-value pairs.
- User metadata this level provides most flexibility and OneData itself does not assume any schema related with these metadata. For each resource, user can assign a separate document in one of supported metadata formats (currently JSON and RDF).

The filesystem and extended level attributes are accessible via **REST-API** and **CDMI** or directly through queries to the embedded database.

## Metadata

0)	ONEDATA					s_0 (0	1	<u> </u>	4	Û S	🚺 admin -	
ß	CTADATASPACE											
Data	СТАДАТАЅРАСЕ	$\sim$	FILES							SIZE	MODIFICATION	Í
Shared	Root directory		astri_000_41_001_00001_	R_000004_000_1002.lv0					(4)	916.88 KB	2017-01-13 12:01	l
Spaces			BASIC JSON	RDF								l
品			DATATYPE	0000	×							l
Groups			DATA_LEVEL	lv0	×							l
			MODES_ID	R	×							l
Tokens			OBSERV_ID	00001	×							l
8			ORIGIN_ID	41	×							l
roviders			PACKET_TYPE	1002	$\times$							l
			PATH	/CTADATASPACE/astri_000_41_001_00001_R_000004_000_1002.lv0	×							l
			PROGRAM_ID	001	$\times$							l
			PROP_ID	00000000000001	×							l
			RUNS_ID	000004	$\times$							l
			SEQUENCE_NUM	000	$\times$							l
			TSTART	430580855	×							l
			тятор	430580965	×							l
			Attribute	Value								l
			Save all changes	Discard changes Remove metadata								
			astri_000_41_001_00001_	R_000005_000_1002.lv0						916.88 KB	2017-01-16 11:01	

## Metadata



# References

- CTA web page: http://www.cta-observatory.org/
- ASTRI web page: http://www.brera.inaf.it/astri/
- YouTube demo: https://youtu.be/UhOWnJluIgE
- INDIGO Data Cloud: https://www.indigo-datacloud.eu
- OneData documentation: https://onedata.org/docs/index.html
- OneData @ docker hub: https://hub.docker.com/u/onedata/



## SSDC as server of CTA data products

- The ASI-SSDC (Space Science Data Center):
  - wide experience as MWL data center, both for low-level data products (AGILE data center, Fermi-LAT/SWIFT/... data mirror center) and high-level data, data products and catalogs.
  - Data and data products integrated in a fully MWL environment (MMIA: Multi-Mission Interactive Archive).
  - Possibility to perform cross-catalog searches betweeen resident and external catalogs.
  - Powerful tools to extract SED of sources and modelization.
  - VHE catalog products from literature already integrated in the TeGeV Catalogue.