



INAF Italian Astronomical Archives facility

17/06/2019 - Cristina Knapic on behalf of IA2 team



Italian Center for Astronomical Archives IA2 Centro Italiano Archivi Astronomici



Home Services -Projects -Software Additional Info -IA2 Group ABOUT US IA2 (Italian center for Astronomical Archive) is an Italian Astrophysical research e-infrastructure project that aims at co-ordinating different national initiatives to improve the quality of astrophysical

data services. It aims at co-ordinating these developments and facilitating access to this data for research purposes. The IA2 is supported by INAF since 2005. IA2's main goals consist in data archiving systems and safety, including data hosting and data curation and preservation, data and metadata distribution over geographical sites, access services including publication within the VO scenario. IA2 provides also services and tool to the community, like data sharing (owncloud), project management (redmine), software collaboration (git-lab) and has available a workflow manager (Yabi) for computational needs.

contact us

MAIN ACTIVITIES

TELESCOPE ARCHIVES & SIMULATIONS MAIN TNG Asiado RADIO MAIN Exoclimates INTRIGOSS LBT OLD HOSTED ASTI INES Byurakan INES BaSTI IBIS-A other archives @ INAF **OTHER SERVICES** uabi EUROUD yabi ownCloud VO initiatives GitLab redmine PROJECTS expand [+]

https://ia2.inaf.it/

collapse [-]

collapse [-]





IA2 proposal preparation support



Generation of Observing programs links to registered/known users via administration tool:

CSV format

Program name	User name and surname	E-mail	Write into Postgres (1=true,	0=false)	User address	Schedule	time	Priority	Call	Support e-mail	
Select CSV file				s	elect range: From	yyyy-MM-dd		Ħ	То ууу	y-MM-dd	Q Search
Choose File No	file chosen			Date	Operation		Portal	llser	Information		Success
Upload COV				2017-10-25	GROUP_CREATION		TEST	Fake Admin	Group g1		~
Opload CSV				2017-10-25	MEMBERSHIP_CREATION		TEST	Fake Admin	User ID RAP:1 User descrip Mario Rossi Group g1	tion (IA2) - mario.rossi@ia2.inaf.it	*
E-mail	Program			2017-10-25	CONFIRMATION_EMAIL_SE	ENDING	TEST	Fake Admin	E-mail zorba@oats	.inaf.it	~
zorba@oats.inaf.it	test2			2017-10-25	MEMBERSHIP_CREATION		TEST	Fake Admin	User ID RAP:3 User descrip Paola Gialli Group g1	tion (Google) - paola.gialli@ia2.inaf.it	*
Operations				2017-10-25	CONFIRMATION_EMAIL_SE	ENDING	TEST	Fake Admin	E-mail zorba@oats	.inaf.it	~
Operation		Sta	tus					× 1 2 3 »			
Creation of group	test2							Total results: 14			
Creation of user m	embership on Grouper			г	Do not fo	raat ta			TED	voureolf	
Sending of confirmation e-mail to Zorba@oats.inaf.it Mail sender is in testing mode Confirm Reject			- - (wheneve allow you observat	to ea	sibl sily rep	le- in a y acco paratio	adva ess on a	ance, it will your nd data!		

https://sso.ia2.inaf.it/home/





Raw Data: the scientific and technical content



Observations are composed by several components:

- the schedule
- the scientific exposures
- the image/spectra description (metadata)
- the calibration exposures
- the engineering data
- the night shifts

All those components mixed together contribute to the correct data analysis and interpretation.

Having all the components available and interoperable thru different systems allow **data reuse**, preservation and curation.

• Continuous interaction with Data Providers to evolve the data models

- Telescope staff, instrumentation specialists, support astronomers, technical staff etc..
- continuous interaction with standards experts
 - IVOA, RDA, FAIR,

See and remember the FAIR principles!





From Observation to data products

- store data for long term preservation;
- store data in a repository for temporary store;
- store data in <u>on line</u> archive;
- serve data in VO compliant manner;
- allow operations on data in a user space;
- allow computation on data;
- allow for interactive/collaborative tools operations on data;
- DOI.....





- Hardware IA2 @ TS:
 - **800 TB**
 - 300 used + 400 free TB on line
 - backup : 100 TB for VMs
 - T950 HPE LTO-8 of 1.25 PB expandible to 12.5 PB (coming next month)
- Hardware @ other sites:
 - IRA: 60 TB on new machine (x testing)
 - SRT : 60 TB on new machine (buffer to Ts)
 - Hardware owned by others:
 - IRA : 40 TB Radio Distributed Archive
 - SRT : 1 TB (pulsar testing machine)
 - Serra La Nave : 500 GB on site
 - LBT : 12 TB upgraded 1TB /y Full LBT Archive
 - Asiago : 500 GB on site
- Bandwidth: 10Gb/s GARR





From data providers to science



Data Providers: Telescopes, Projects (Simulations or Surveys), Users

Types of data:

- <u>Structured data</u>: well organized data following defined standards and always containing all the descriptors to guarantee the FAIR principles;
 - Easy to handle by a automatic routine, can be ingested into a searchable archive;
- <u>Partially structured data</u>: data customized to fulfill specific needs but not always coherent with the previously produced one. Sometimes it follows a self define standard and can be seen as an evolving data model dataset;
 - Require frequent updates to the data model and can be difficulty be published in a unique catalogue;
- <u>Shared data</u>: fully incoherent data composed usually of a mix of data models and formats. Totally customised to the current utilization and defined by human approach (no standardization). It usually contains a mixture of raw, reduced data and lists of parameters, comments, wiki pages etc;
 - require a collaborative tool to support data sharing and cannot be published by an archiving service. It can be part of the DOI publication but in any case presents difficulties in data searching.



Structured Data: what is a data collection?



- Data collection is the **systematic approach** to gathering and measuring information from a variety of sources to get a **complete and accurate** picture of an area of interest in a **standardized and established manner** that enables the collector to **answer or test** hypothesis and evaluate outcomes of the particular collection. The goal for all data collection is to capture quality evidence that allows analysis to lead to the formulation of convincing and credible answers to the questions that have been posed.
- Accurate data collection is essential to maintaining the integrity of research, guaranteeing findability, accessibility, interoperability and reproducibility (FAIR principles);
- Impact of faulty data:
 - Inability to answer research questions accurately;
 - Inability to repeat and validate the study.



What is a data collection?



- A collection (used as a noun) is the topmost container for grouping related documents, data models, and datasets.
 - Data model: is an abstract model that organizes elements of data and standardizes how they relate to one another and to properties of the real world entities;
 - **Dataset**: is a collection of data. Most commonly a data set corresponds to the contents of a single database table, or a single statistical data matrix, where every column of the table represents a particular variable, and each row corresponds to a given member of the data set in question. The data set lists values for each of the variables, such as height and weight of an object, for each member of the data set. Each value is known as a datum. The data set may comprise data for one or more members, corresponding to the number of rows.

AND ASTROPISICA



Data models: ObsCore and CAOM



Data Models to Archives



C

Help

C New search

nadir.ira.inaf.it:8080/radio/faces/search.xhtml;isessionid=CDFFAA5328909FE06DA2C64A4023EF02?dswid=-8495

Q 1

Login

🖹 👻 Your files 💽 - Currently not logged in



IA2 support to Science



- **Store and preserve** astronomical data (observed or simulated);
- **Support data providers** in correctly set up Archives:
 - from raw data to calibrated one from Telescopes;
 - simulated data (exoclimates, intrigoss, cosmological)
- **Publish public data** through the VO services;

• Support Astronomers in data retrieval via web interfaces:

- search on public data without login;
- after login, a list of proprietary datasets are presented and filters can be applied;
- using filters to find data;
- single file direct download;
- user space for bunch of files;
- possibility to download VOTables of queried data;
- possibility to download a CSV file of queried data;
- possibility to download a URL list of files to download them using external tools like WGET;
- name resolver to find coordinates of objects;
- connection via SAMP HUB to link VO clients like Aladin or Topcat;



The acquired know how from data provider side made evident the collaboration between domain experts and scientific community is essential to correctly identify data models and to define the specific content of the exposed services. In this scenario, Virtual Observatory, FAIR policies, Open Data, RDA play a fundamental role.

User support for data retrieval is essential too!!



- ARC and AENEAS experience: for the future 2020 era telescopes a scientific support to extract the correct information from the huge amount of data is ESSENTIAL;
- CADC experience: user space available to observing program collaborators where to manage and handle data with available or provided pipelines is highly desirable in order also to interoperate and cross correlate different datasets and scientific products;
- DOI: papers and datasets!! After paper publication, datasets should be public!



http://archives.ia2.inaf.it/tng/faces/help.xhtml

 Ownload •

 • Create tar from selected
 • VOTable (all query results - 3 rows)
 • URL list (all query results - 3 rows) .txt
 • Construction of the selected of

2009-09-30

2009-09-30

KQEB0005.FTS.gz

KQEB0006.FTS.gz

Using external tools (e.g. wget)



IA2 SAMP





©0	Filename	P
	AF187351.fits.gz 🔊	S
	AF187352.fits.gz 🔊	S
	AF187353.fits.gz 🔊	S
	AF187354.fits.gz 🔊	S
* ▲	Your downloads 🗿 -	. (

If your SAMP hub is running the icon is green and you can register on it to see the other applications.

When you are registered to the hub you can select the target application.

If you are unregistered or you haven't select any application your messages will be sent to all active applications.

You can:

- Send a VOTable by SAMP (MType table.load.votable)
- Send an image by SAMP (MType image.load.fits)

When a SAMP hub is detected, a "signal" icon will appear near search result files and generated VOTable files.





Login to access tools

ON PSTROFISICA





WEB: Services, portals e VO compliancy

collaborate with

TWiki



Indi ∞









WEB: Services, portals and VO compliancy

Science ready data: the TNG example

INAF

- → reduced data using pipeline DRS published on TNG portal following observing program based policy (from 1 to 3 years);
- → all HarpsN and GIANOB observers have the possibility to reprocess data using in an interactive way the installed pipelines (DRS and Gofio). Pipelines are updated on YABI platform by IA2 team.

Yabi : platform for the management of data reduction pipelines.



Advantages: -

Disvantages:

- → software up-to-dated to the last pipeline release;
- → private pipelines, potentially public data;
- → intuitive and interactive pipeline usage;
- → users do not need computing power.
- → local data management (currently no reingestion of reduced data for privacy reasons)
- → currently not integrated with RAP



- IA2 manages data in a distributed manner on 3 continents!!
- GAPS experience bring important know-how. Pipelines and workflow management systems will be the must of 2020 Era Telescopes;
- IA2 allow state-of-the-art authentication and authorization mechanisms. Same results will be applied in pilot project for SKA and AENEAS.





Connecting Research and Researchers



N PSTROFISIC

INAF

The same person, multiple identities, multiple user ID



One user, multiple identities, the same user ID

Grouper



- secure http:	s.//source.main.ogrouper/grouper/or/app/or/zmain.index/ope	Ardton=01725t ¥
INTERNET	Sear	ch Q
r	Logged in as Cristina Knapic (eduGAI	N+Google+Linke · Log out
Create new group 🝷	Home > Root > OATS > LBT	
uick links _	LBT	Edit folder
ly groups	Large Binocular Telescope	More actions
ly folders	More ~	
ly favorites		
ly services	Folder contents Privileges More -	
ly activity		
liscellaneous	Filter fo Folder, group, or attribute nan Apply filter	
ite UI	Reset	
	Nama –	
rowse folders 🛛 🖓		
🔁 Root		
- 🔁 OATS	* ADMIN	
	e brothberg_01	
E CBT	Stothberg_02	
ADMIN ADMIN	<pre>eveillet_01</pre>	
<pre>brothberg_01</pre>	<pre>tellet_02</pre>	
<pre>weillet_01</pre>	<pre> cveillet_03 </pre>	
<pre>weillet_02</pre>	<pre> cveillet_04 </pre>	
<pre>ecveillet_03</pre>	<pre># dthompson_01</pre>	
<pre> cveillet_04 </pre>	dthompson_02	
dthompson_01	dthompson_03	
and dipompeop (1)		

You can give access to a set of data adding a member to your group (corresponding to a program in this case).

More than one application could query on the group management system and allow different operation and privileges to the members, so you can easily collaborate.

To run a specific temporary project, a user space is a very nice place where to share files and run applications that cannot be operated locally, but needs more resources than our laptop could offer. It is in development.....

AND F ROPACT	Addi	ng a use	r to your	group	·
INTERNET.		Logged in as Soni	Search Q a Zorba (eduGAIN+X.509+Google) · Log out	https://sso.ia2.inaf.	it/help-grouper/
+ Create new group 🔹	Home > Root > OATS > TNG_GROUPS				
Quick links —		S	Edit folder More actions▼		
My folders My favorites	Folder contents Privileges	More -			
My activity Miscellaneous	Filter for: Folder, group, or attribute name	∞ ≤ Ia2-sso:8080/group ← → C □ Ia2-sso:8080/grouper/gr	ouperUl/app/UiV2Main.index?operation=	:UIV2Group.vlewGroup&groupId=a24091cc3a014	a6bb5ef74d33942; 🛧 💩 🥝 🕅 :
Admin UI Lite UI	Name -	INTERNET		Logg	Search Q
Browse folders	A16TAC_10 A16TAC_11 A16TAC_12 A16TAC_13 A16TAC_14 A16TAC_15	Create new group Quick links My groups My folders My favorites	Home > Root > OATS > TNG_GROUPS	CO.CEP	+ Add members
	Co-I	My services My activity Miscellaneous Admin UI Lite UI Browse folders	More ~ Members More ~ The following table lists all e	entities which are members of this gru	bup.
Collaborator	Support staff Collaborator	 ➡ Coot ➡ OATS ➡ ➡ etc 	Filter for: All members Remove selected members	▼ Member name	Apply filter Reset
27			Entity name ▼	Membership	Actions 🔻
Post-Doc			Show: 50 V	5	Showing 1-1 of 1 · First Prev Next Last
		© Institute of Higher Education			27



Under development: VOSpace for storage and computation



VOSpace implementation compatible with CADC implementation, possible share of authorizations (thanks to S. Bertocco, G.Taffoni, S. Gaudet, P. Douler, B. Major experiments within EgiEngage)

Two levels of computation:

- 1. user approach to interactive pipeline with no HPC/HTC and small data volumes
 - a. RAP + Yabi;
 - b. Containers on IA2 infrastructure;
- 2. user approach using bash processing with HPC/HTC needs
 - a. Containers or VMs on Chipp;
 - b. GCloud.

Processing close to Data + POCs





Your data



- Remember your reduced data set is composed by several components:
 - \circ raw and reduced datasets (see above)
 - \circ Pipeline
 - reduction parameters (filters applied..etc..) and intermediate products
 - \circ tables / plots

DOI for Data Sets — supported by ICT What about creating catalogues of DOIs using additional metadata (data descriptors)? It allows datasets exploitation and burst the number of citations

See details in R. Smareglia talk

DOI (Digital Object Identifier):

DOI is a character string used to uniquely identify an object such as an electronic document. The **DOI** for a document is **permanent**, whereas its location and other meta-data may change. Referring to an on-line document by its DOI provides more stable linking than simply referring to it by its URL, because if its URL changes, the publisher need only update the meta-data for the DOI to link to the new URL.





The ideal plan



- Improve the services offered (link to proposal submission);
- Increase the number of services offered (possibility to automatic re-processing of data with last version pipeline, see ARI-L project; basic Pipelines for all the instruments);
- provide data curation before publication;
- Increase the number of supported data models if not compatible with CAOM in all the chain (archives, automatic data reduction software, VO);
- Grant access to all the offered services via SSO;
- Adopt a VOSpace compliant service to allow user space utilization (storage + computation);
- data/catalogues publication increasing DOI metadata;
- Improve the user experience in order to stimulate the Astronomical Data life cycle.



Conclusions



- Target of Data Centers:
 - Support science offering Archival services, computing power and maintenance of both
- Target of the Astronomers:
 - Use of data center services
 - Have their own User Space where to make science!
 - share data, papers, ideas, everything!!
 - Reporting problems
 - Suggest upgrade/necessities

Thanks for your attention!