

COEOSC FAIR-EASE Building Interoperable Earth Science & Environmental Services

Integrating a community

Earth sciences use cases

Samuel Keuchkerian - (CNRS)

Marie Jossé - Data Terra (CNRS)



FAIR-EASE has received funding from the European Union's Horizon Europe Framework Programme (HORIZON) - under grant agreement No. 101058785.

FAIR-EASE

Building an interdomain digital architecture for distributed and integrated use of environmental data



FAIR-EASE Data Discovery and Access Interdisciplinary Service FAIR-EASE Virtual environments

0 0 0

FAIR-EASE Earth sciences use cases

5 pilots for an earth system model



Coastal Water Dynamics: focuses on the coastal marine environment near river estuaries, where important processes

take place.



Earth Critical Zone: monitors land and soil degradation.



Volcano Space Observatory: monitors global volcanic activity, allowing the focus on any major volcanic eruption worldwide

Ocean Bio-Geochemical Observations:

addresses fundamental scientific

questions regarding the health of marine

ecosystems (e.g. ocean acidification, ...)

and needs for ocean resource

management.

Marine Omics Observation:

analyses of spatial- and time-comparable marine microbial metagenomics data sets for the exploration of biodiversity and its correlations with environmental quality

FAIR-EASE FAIR-EASE datalake infrastructure



9. 9. B

FAIR-EASE Pilots on D4science VLabs

1

Marie Jossé 👻

Go to -

🚳 🍃 🖉 🖂 🔍

Marine Omics Observations . Members G JupyterLab About Marine Omics Observations MarineOmicsObservat . ORecent 7 Name Owner Last modified Starting from an ongoing effort undertaken by the EMBRC infrastructure, with the establishment of the European Marine Omics Biodiversity Observation Network (EMO-BON), this pilot focuses on the challenge to set up a web-based VRE to provide products and services orientated to non-specialist researchers interested in omics approaches to study marine biodiversity. Today, data FM 25 Sep 14:53 24 EMO-BON includes several marine stations that will sample for genomic microbial marine biodiversity, essential ocean variables (EOVs), and essential biological variables (EBVs). Current limitations and needs: Data & Repositories: Analytical services: FAIREASE impacts: parquet files FM 20 Sep 10:59 24 See less Shared attachments AT 11 Mar 10:56 24 Other options .. retrieve dat ... FM 21 Oct 16:07 24 run alpha di ... FM 11 Dec 17:35 24 Previous Next Show 5 entries 1 to 5 of 12 items Go to shared workspace **VRE Managers and Groups Wiew Managers**

FAIR-EASE FAIR-EASE FAIR-EASE FAIR-EASE Pilots on Galaxy Europe



We are working in close collaboration with the Galaxy Training Network (GTN) to develop training materials of data analyses based on Galaxy. If you

Settings

Galaxy Training Network

An easy way to learn how to use Galaxy and improve your skills on various domains for instance a set of tutorials are available on FAIR management

🕀 The Workflow Run RO... 🗧 Galaxy | Europe 🐨 conda-force | commu... 🗧 GCC CoFest 2024 - Cro... 🦉 Galaxy | Configured b... 🗮 Galaxy 🗃 Home - Galaxy Comm... M Settings - Q Search Tutorials Galaxy Training! Learning Pathways @ Help *

Tutorials

Welcome to Galaxy Training!

Collection of tutorials developed and maintained by the worldwide Galaxy community

Galaxy for Scientists We have separated the tutorials into several categories based on field and technology. We are exploring other ways to organise the tutorials going forward! Start Here Introduction to Galaxy Analyses 13 Using Galaxy and Managing your Data 22 Not sure where to start? Try the NGS Basics Learning Path!

Scientific Fields

	7110100			
Торіс	Tutorials			
Climate	12			
Computational chemistry	9			

Quickstart Learning Pathways Galaxy for SysAdmins Galaxy for Galaxy for Teachers Developers



Upcoming Events Check out upcoming events around the Galaxy!

January 28, 2025 Galaxy at SURF Research Cloud workshop

February 4 - 6, 2025 Code & Collaborate: The FAIRytale of Software Development

A catalog of tutorials

- Pathways on a dedicated topic
- Classes, courses, webinars, and other

interactive events

A community for the community https://training.galaxyproject.org/t raining-material/



Galaxy Training Network

FAIR-EASE contribution (for earth sciences)

Earth sciences discovery tutorials

Thematic tutorials on ocean, land, atmosphere, and biosphere

Development tutorials end pathways to build a subdomain and a community with Galaxy

C Ocean C This tutotrial aims at familiarzing you with Earth Science and discovering the earth data available on Galaxy. The target audience is not a scientist but anyone interested in learning about Earth system. 🍇 Land 😹 Atmosphere Agenda R Biodiversity 🎭 In this tutorial, we will cover: Climate @ Learning Pathways @ Help • Conclusion Galaxy Training 1. C Ocean C Extra information 1. The Argo program Cetting your hands-on earth data 2. Copernicus Data Space Ecosystem Questions 3. EMODnet Chemistry Author(s) Marie Josse Feedback 4. Copernicus Marine Data Store Resciencers 🚱 🚯 Citing this Tutorial 2. 🍓 Land 🌌 1. Copernicus Data Space Ecosystem 2. QGIS (Geographical Information System) 3. Atmosphere 1. Copernicus Data Space Ecosystem 2. Climate Data Store 4. 💂 Biodiversity 🏇 💫 Galaxy Training! 🗋 Climate 🕼 Learning Pathways 🕐 Help 🔹 🏟 Settings 👻 Q. Search Tutoria 1. Marine biodiversity 2. Land biodiversity by 5. Conclusion 6. Extra information Ccean's variables study Author(s) Marie Josse Reviewers 🤗 🏟 🙆 Overview (?) Questions: © 0 · How to process extract ocean's variables? · How to use ODV collections? · How to create climatological estimates? Objectives: Tool development for a nice & · Deals with ODV collection with data originating from Emodnet chemistry shiny subdomain · Visualise ocean variables to study climate changes Discover Galaxy's communities and learn how to create your subdomain and enrich it by writing, testing and submiting your tools on Galaxy. This learning pathway will guide you through all the steps required to build a tool for Galaxy with Planemo for batch tools and how write an interactive tool subdomain community tool development 3-day course dev



Bérénice Batut Anika Erxleben-Eggenhofer Ph.D., Researcher Dr. rer. nat., Researcher bebatut@informatik.uni-freiburg.de erxleben@informatik.uni-freiburg.de O bebatut O erxleben +49(0) 761 - 203 54126 +49(0) 761 - 203 54130 Q Build.: 079, Room: -1006 A http://research.bebatut.fr 🗩 🎔 🙆 in 💿 🞖 R⁶ [m] in 💿

Helena Rasche B.Sc. Biochem., Technician

Freiburg Team

hxr@informatik.uni-freiburg.de O hexylena Galaxy Administrator



Paul Zierep Dr. rer. nat., Researcher

zierep@informatik.uni-freiburg.de O paulzierep +49(0) 761 - 203 54130



0 0 0

Services for EOSC in the proposal for a node

(and the EOSC EU node)

French national digital infrastructures

- Renater (Geant);
- France Grille (EGI), & Mesonet > GENCI / EuroHPC: IDRIS, CINES, TGCC;
- National & regional labelled data centres and meso-centres.

COCOSC FAIR-EASE

EOSC Federation

- D-T data and services accessible through the EOSC EU node;
- Services interoperability with thematic cluster nodes and related national nodes.

Core services

- · Distributed data storage and management;
- Large data transfer (files, objects);
- User spaces (interactive notebooks, virtual machine, container images);
- HPC/Cloud computing services
- Federated AAI.





Integrated Earth System Observation - Data

Terra



MAIN Develop a global system for accessing and processing observation data (satellite, in situ), value-added products and services to observe, understand and predict in an integrated manner the functioning and evolution of the Earth system.



Partners

French scientific organizations and universities





A multidisciplinary approach because it calls on work in several areas of Earth System sciences

An inclusive project that goes beyond the scientific circle with an approach also oriented towards field actors and participatory data

OURFacilitating the cross-referencing of observations and the modeling of Earth System dataSERVICESThe IR Data Terra offers services around Earth system observation data. The objective is to provide interoperable and

interdisciplinary services at all levels.



wednesday session



COEOSC FAIR-EASE Building Interoperable Earth Science & Environmental Services

Data storage and handling

Samuel Keuchkerian - (CNRS/FAIR-EASE)

Marie Jossé - Data Terra (CNRS)



FAIR-EASE has received funding from the European Union's Horizon Europe Framework Programme (HORIZON) - under grant agreement No. 101058785.







Carning Pathways

⑦ Help ▼

Search Tutorials

Q

Settings

Adding file-sources to Galaxy



This tutorial is not in its final state. The content may change a lot in the next months. Because of this status, it is also not listed in the topic pages.



Overview

② Questions:

· How to set up an S3 bucket

Objectives:

Add your S3 bucket on Galaxy

Time estimation: 15 minutes

- C Supporting Materials:
- The Published: Jan 20, 2025

📩 Last modification: Jan 20, 2025

Find the information you need

Add the S3 bucket

This tutorial demonstrates how to implement an S3 bucket as a Galaxy file-source within Galaxy. We will add here the public Argo data Amazon S3 bucket. Argo is an international program that observes the interior of the ocean with a fleet of profiling floats drifting in the deep ocean currents (https://argo.ucsd.edu). It started 20 years ago and is a dataset of 5 billion in situ ocean observations from 18.000 profiling floats (4.000 active). The Argo GDAC dataset is a collection of 18.000 NetCDF files. It is a major asset for ocean and climate science and a contributor to IOCCP reports.



	Find the information you need	
OSC	Add the S3 bucket	Agenda
030 117	Conclusion	In this tutorial, we will cover:
	Frequently Asked	1. Find the information you need
	Questions	2. Add the S3 bucket
	Feedback	3. Conclusion
	Citing this Tutorial	

Find the information you need



Hands-on: Find an S3 bucket

Go on Amazon Sustainability Data Initiative.

There you can visit the catalog of data, and by searching for Argo you can directly get to the Argo registry.

On this last page you'll find all the information you'll need to add the S3 bucket to Galaxy

O A ≈ https://registry.opendata.aws/argo-gdac-marinedata/

M - Agenda 🕀 The Workflow Run RO... 🕏 Galaxy | Europe 🐄 conda-forge | commu... 🗧 GCC CoFest 2024 - Cro... 🦉 Galaxy | Configured b... 💆 Galaxy | Configured b...

Registry of Open Data on AWS

Argo marine floats data and metadata from Global Data Assembly Centre (Argo GDAC)

pology memory climate paracenter olgitis assets geometricity geophysics g

Description

Argo is an international program to observe the interior of the ocean with a fleet of profiling floats drifting in the deep ocean currents (https://argo.ucd.edu). Argo GDAC is a dataset of 5 billion in situ ocean observations from 18.000 profiling floats (4.000 active) which started 20 years ago. Argo GDAC dataset is a collection of 18.000 NNeCDF flies. It is a major asset for ocean and climate science, a contributor to IOCCP reports.

Update Frequency

Data is updated daily.

License

Open data, there are no restrictions on the use of this data. https:// creativecommons.org/licenses/by/4.0/

Documentation

http://www.argodatamgt.org/Documentation

Managed By

Euro-Argo

See all datasets managed by Euro-Argo.

Contact

Resources on AWS

Description Argo GDAC data and metadata

Resource type S3 Bucket

Amazon Resource Name (ARN) arn:aws:s3:::argo-gdac-sandbox

AWS Region eu-west-3

AWS CLI Access (No AWS account required) aws s3 ls --no-sign-request s3://argo-gdac-sandbox/ Explore

aws

Browse Bucket

you need

Add the S3 bucket

Frequently Asked Questions

Feedback

Conclusion

meosc

Citing this Tutorial





Hands-on: Add on Galaxy

- If not already done clone the Galaxy Europe Infrastructure-playbook repo
- Create a branch on your fork
- Go to the file file_sources_conf.yml.j2 in templates/galaxy/config/

There you can edit the file and add your S3 bucket by adding a Argo specific section, like in the following:

- type: s3fs

label: Argo marine floats data and metadata from Global Data Assembly Centre (Argo GDAC) id: argo-gdac-sandbox

doc: Argo is an international program to observe the interior of the ocean with a fleet of profiling f bucket: argo-gdac-sandbox

anon: true



Finally, commit your changes and write a nice message for the admin when you open your Pull Request.



. . . .

0 0 0

Help - Community - About - Projects - @jxtx

🖓 Edit

Search

← Back to News

Bring Your Own Storage with Galaxy

Connect Galaxy with external storage resources easily

September 20, 2024

Introduction

The EuroScienceGateway Project is streamlining the way that users *Bring their Own Storage (BYOS)* to Galaxy. This post covers a specific case in which a user has access to storage resources in the EGI Federated Cloud and wants to connect it to the EU Galaxy instance. However, similar steps can be followed to connect a different Galaxy instance with storage capacity using Azure Blob Storage, Amazon Web Services S3 Storage, and Google Cloud Storage.

In order to connect Galaxy with external cloud storage, the user needs to click on the User drop-down menu, select Preferences and then Manage Your Storage Locations. This will display a list of the cloud storage services already connected to Galaxy. The first time it will be empty and the user can connect to a new cloud storage service clicking on Create. The following options are given to the user:



Azure Blob Storage Amazon Web Services S3 Storage Select storage location template to create new storage location with. These templates are configured by your Galaxy administrator.

Any S3 Compatible Storage

Google Cloud Storage

Generally speaking the user needs the following information to connect Galaxy with cloud storage:

1. Name of the bucket

2. Access key

3. Secret key

Getting S3 storage from EGI

Users with access to cloud resources in the EGI Federated Cloud can benefit from the available TOSCA template in Infrastructure Manager (or IM) for the automated deployment of MinIO.

Below are the steps that a user needs to follow to make use of computing resources in the EGI Federated Cloud:

- Create an EGI Check-in account.
 Enroll in a Virtual Organization (VO). You need to wait for approval before moving on.
 Once you are a member of a VO, configure credentials in Infrastructure Manager.
 In Infrastructure Manager select Deploy a VM and click Configure.
 Select MinIO and click Add.
- 6. Configure deployment details for MinIO and click Submit.

A prerequisite for the deployment of MinIO with IM is the use of some sort of Dynamic DNS service. EGI also offers a Dynamic DNS service for free for users with an EGI Check-in account (see a tutorial for further help). Using a Dynamic DNS service the user needs to register two DNS hostnames: one DNS hostname for the MinIO Console, and another one for the MinIO API endpoint.

IM will deploy a Virtual Machine with Docker and MinIO as a containerised service. Therefore, the first step for the user is to select Deploy a VM in IM:





Examind / Iceberg/WebODV effective from this month end.

More precise elements: Cloud IaaS : Openstack (community version) https://www.openstack.org/ total : 48 hypervisors, 800 CPU oores, RAM 7.6 TB,

Distributed storage/ S3 : Ceph (community version) https://ceph.io/ total : 600 TB SSD (for vm openstack), 5.5 PB HDD (S3) (replicated 3 times so /3)

FAIR-EASE specifically allocated resources: 44 CPUs, 152 GB RAM, 128 GB SSD, 100 TB via API S3

Cluster kubernetes FE : kubernetes v1.31, déployed withTerraform et Talos Linux (https://www.talos.dev/)

3 vm for control plane

3 vm for workers executed on client services instances (JupyterHub/JupyterLab)

UCA TESTBED INFRASTRUCTURE





. . . .

0 0 0

Apache Iceberg an open table format for huge analytic datasets (including metadata)

Iceberg adds tables to compute engines including Spark, Trino, PrestoDB, Flink, Hive and Impala using a high-performance table format that works just like a SQL table. https://iceberg.apache.org/docs/1.5.2/

- Iceberg avoids unpleasant surprises. Schema evolution works and won't inadvertently un-delete data. Users **don't need to know about partitioning to get fast queries.**
- Iceberg is used in production where a <u>single table can contain tens of petabytes</u> of data and even these huge tables can be read without a distributed SQL engine.
- Open standard designed and developed to be an open community standard with a specification to ensure compatibility across languages and implementations. <u>Apache Iceberg is open source</u>, and is developed at the Apache Software Foundation.

Summary

Iceberg excels in decoupling storage from compute, managing evolving schemas, and enabling scalable, FAIR-compliant data workflows, making it ideal for research infrastructures.

For FAIR, open data infrastructures like FAIR-EASE, Iceberg offers flexibility and adaptability while ensuring long-term usability and performance at scale.

Iceberg to ERDDAP data publication



OPEOSC FAIR-EASE Integrating New Data Storage and Access Paradigms in pilot's practices

JECT STO	127.0.0.1:90	1) fyrowser/warehouae/glodap%2Fglodap, v2, 2023%2Fdata%2F +* Object Browser	Q Start typing to filte	er objects in th	e bucket		01	A ☆ 0 0 ¢ 6 % … (@ 0 €				
Object Browser Access Keys Documentation		warehouse Created on. Thu, Jul 11 2024 10:07:15 (dMT+1) Ac warehouse / glodap, / glod	xess: PRIVATE 147.3 MiB - 31	l Objects Las	at Modified		Rewind 9	Refrech C. Upbod L C. Create new path uit Size				
		00003-4-2ff596d6-09b8-487c-8390-5bc325	45b35-0-00001.parquet	Too	iay, 10:12			51.5 MiB		Iceberg to E	RDDAP data publication	
ilstrator		00009-10-2ff596d6-09b8-487c-8390-5bc32	b45b35-0-00001.parquet	Too	iay, 10:12			25.9 MiB		-		
Buckets	💭 jupy	/ter Glodap_iceberg_ Last Checkpoint: 8	hours ago									
Dollalaa	File Edit	View Run Kernel Settings Help						Trusted	UCA 53 Object Store		Trino Clients SQL interfaces	ERDDAP.cosc-faircase
1	a + 3	< 🗋 🗎 🕨 🖩 🗘 🏎 Code 🗸						JupyterLab 🖸 🐞 Python 3 (ipykernel) 🔿 Format SQL	- https://s3.wesocentre.uca.fr			
1	[6]:	N%sql SELECT COUNT(*) as cnt FROM glodap.glodap_v2_2023 cnt							Arristic Sector	Catalog	Trino (query Engine)	
		1402829								\sim		
/	[7]:	%%sql SELECT doi,latitude,longitude,depth,MAD	(temperature) FROM g	lodap.glo	odap_v2_2023	3 GROUP BY doi,lat	itude,long:	itude,depth ORDER BY MAX(temperature) DESC;		+ pySpork		Bana In Cell 199 In State
1	[7]:		doi	latitude	longitude	depth max(tem	erature)			(Query Bryine)		
_		https://doi.	rg/10.25921/16y6-9e29	-18.0	-71.92	6.0	34.501	•				
		https://doi.	rg/10.25921/16y6-9e29	-22.0	-70.75	17.0	32.975					
		https://doi	org/10.25921/0sta-y820	16.0	41.5833	0.0	32.68					
		https://doi	org/10.25921/0sta-y820	16.0	41.5833	10.0	32.63					
		https://doi	org/10.25921/0sta-y820	16.0	41.5833	20.0	32.5					
		https://doi	org/10.25921/0sta-y820	14.7	42.1667	0.0	32.4					
		https://doi	org/10.25921/0sta-y820	14.7	42.1667	10.0	32.39					
		https://doi	org/10.25921/0sta-y820	14.7	42.1667	20.0	32.35					
		https://doi.org/10	.3334/cdiac/otg.ndp080	24.569	57.222	4.0	32.171		"Open Dat	a Lake House Format - Apa	ache Iceberg - Quick start" Dar	mian Smvth:
	1	https://doi.org/10	.3334/cdiac/otg.ndp080	26.411	56.57	4.0	31.912	-	Spon But		the second se	
		https://doi.org/10.3334/cdiac/otc	nacifica 49ho19930807	-2.0	175.0	0.0	31.8					

COROSC FAIR-EASE EXAMIND

The Examind provides all the features you need to build a Geographic, Hydrographic, Oceanographic and Meteorologic geospatial data processing infrastructure. Designed for **interoperability**, all the products in the suite comply with **OGC standards**, and integrate OGC Web services

Examind community is a open-source platform / server.

- It manages multiple data formats (Netcdf, geotiff, etc.), clouds native data formats (COGs, GIMI), multiple OGC (Open Geospatial Consortium) standards (WCS, WMS, WPS, OGC API, etc.)
- It has a number of ready-to-use processes for various uses in the geospatial world.
- It offers several ways of managing and structuring data. The server can connect to an existing data source, hold the data locally, or generate new data via different processes (via WPS, OGC API Process, openEO); all via different paradigms, such as data cube structuring.

Examind supports several data storage options local storage, FTP, S3 (AWS / Minio), via HTTP/HTTPS, or storage from a WMS / WMTS service, etc.

The service can be deployed via a docker image and a docker-compose file on kubernetes infrastructures.



Summary

Iceberg excels in decoupling storage from compute, managing evolving schemas, and enabling scalable, FAIR-compliant data workflows, making it ideal for research infrastructures.

For FAIR, open data infrastructures like FAIR-EASE, Iceberg offers flexibility and adaptability while ensuring long-term usability and performance at scale.

WEDNESDAY AFTERNOON

COREC FAIR-EASE FAIR-EASE in the context of the European green deal



0 0 0

. . .

0 0



0 0 0

FAIREASE in the context of the european green deal

European green deal topics

- Designing a set of deeply transformative policies
- Mainstreaming sustainability in all EU policie

- 2.1.7.
- Preserving and restoring ecosystems and biodiversity
- 2.2.2.
- Greening national budgets and sending the right price signals
- 2.2.3.
- Mobilising research and fostering innovation

0 0 0

. . .

- 2.1.7.
- Preserving and restoring ecosystems and biodiversity
- 2.2.2.
- Greening national budgets and sending the right price signals
- 2.2.3.
- Mobilising research and fostering innovation

https://www.aeris-data.fr/leruption-volcanique-s ur-les-iles-tonga-vue-par-les-satellites-geostation naires/

0 0 0

. . .

