# Integrating a community

*Earth sciences use cases*
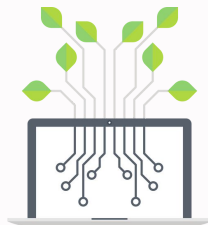
**Samuel Keuchkerian - (CNRS)**

**Marie Jossé - Data Terra (CNRS)**

# 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

# 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.
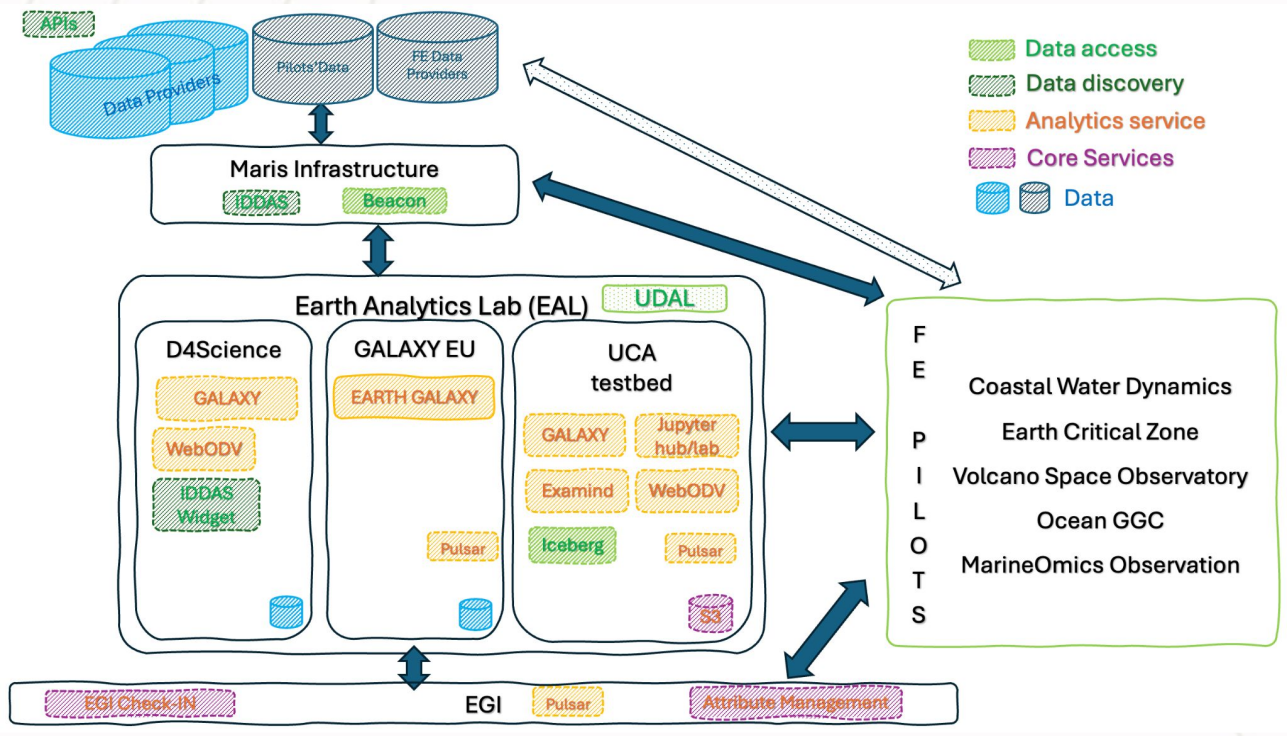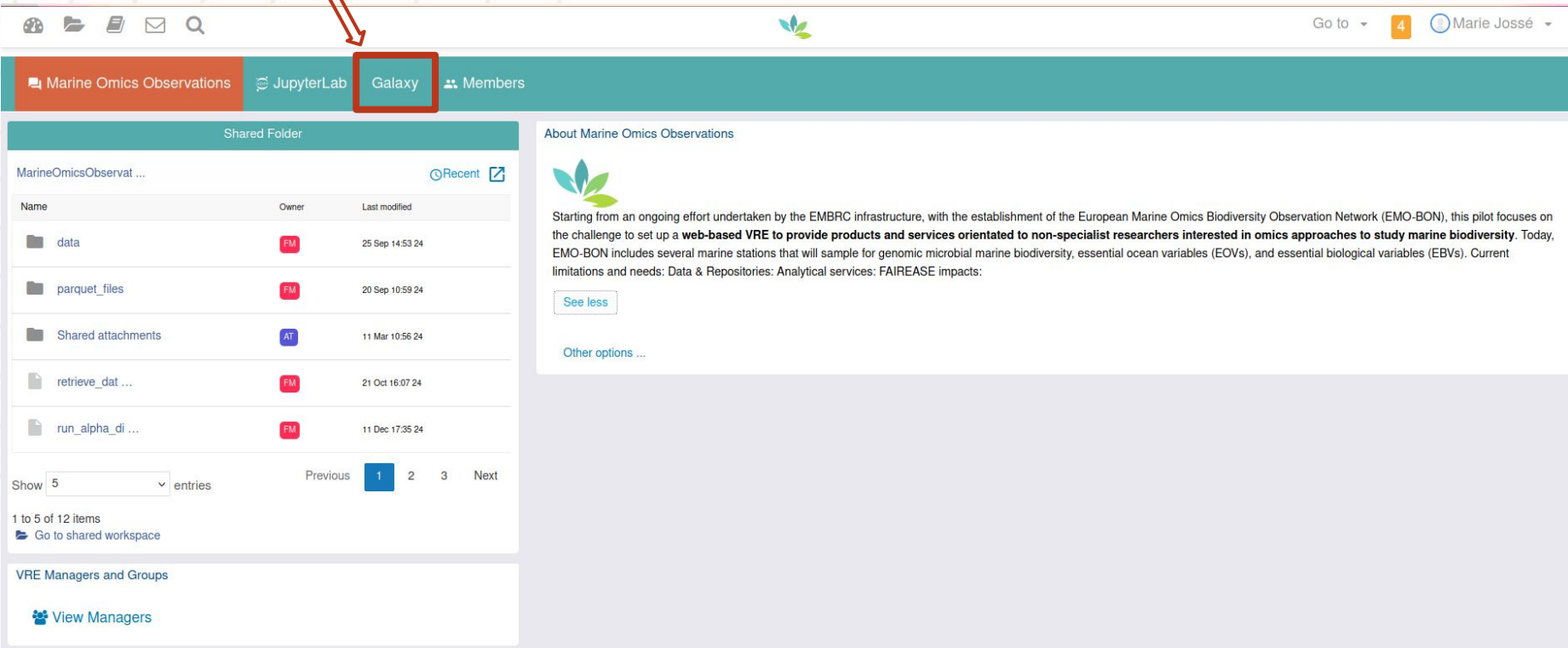
**Marine Omics Observation:** analyses of spatial- and time-comparable marine microbial metagenomics data sets for the exploration of biodiversity and its correlations with environmental quality

# FAIR-EASE datalake infrastructure

# FAIR-EASE Pilots on D4science VLabs

# FAIR-EASE Pilots on Galaxy Europe

# 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

• 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/training-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**

# Importance of the team in the sustainability of Galaxy

**A vibrant team of :**

**Administrators**

**Developpers**

**Trainers**

**BUT with a strong focus on life sciences**



Freiburg Team

**Björn Grüning**
Dr. rer. nat., Researcher, Head of the Team
✉ gruening@informatik.uni-freiburg.de
🔗 bgruening
📞 +49(0) 761 - 203 54130
🔍 Build.: 079, Room: -1006
💬 [m] 🐘 in 🔵 🌐 R⁶

**Bérénice Batut**
Ph.D., Researcher
✉ bebatut@informatik.uni-freiburg.de
🔗 bebatut
📞 +49(0) 761 - 203 54126
🏠 http://research.bebatut.fr
💬 🐦 🐘 in 🔵 🌐 R⁶

**Anika Erxleben-Eggenhofer**
Dr. rer. nat., Researcher
✉ erxleben@informatik.uni-freiburg.de
🔗 erxleben
📞 +49(0) 761 - 203 54130
🔍 Build.: 079, Room: -1006
[m] in 🔵

**Wolfgang Maier**
Dr. phil. nat., Researcher
✉ maierw@informatik.uni-freiburg.de
🔗 wm75
📞 +49(0) 761 - 203 54126

**Helena Rasche**
B.Sc. Biochem., Technician
✉ hxr@informatik.uni-freiburg.de
🔗 hexylena
🔵 Galaxy Administrator

**Paul Zierep**
Dr. rer. nat., Researcher
✉ zierep@informatik.uni-freiburg.de
🔗 paulzierep
📞 +49(0) 761 - 203 54130

**Need the same team for earth sciences**

# Earth community



Galaxy pyramid of needs

# Services for EOSC in the proposal for a node

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

**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.



**Interoperability and common standards**

**Code & Software artefacts repositories**
*(e.g. gitlab/github, dockerhub/quay.io, Galaxy toolshed)*

publish

**DATA TERRA**

**IAM & Security**

**Analytics Services**

**Tools**

**Data**

**Resources Catalogue and Registry Services**
*(Shared among all User Spaces)*
- **Analytic Platforms Assets** *(e.g. notebook, workflow)*
- **Analytic Platforms** *(e.g. JupyterLab, Galaxy, [Geo]viewer)*
- *Galaxy*
- **Appstore** *(e.g. QGis desktop, eo-dag)*
- **Dataspace** *(e.g. satellite, in-situ, model, biodiversity)*

**User Space** *(Bring your own ...)*
- **Analytic Platforms Assets** *(e.g. notebook, environment)*
- **Interactive Apps** *(e.g. shiny, dash)*
- **Any tools** *(e.g. data tools, algorithms)*
- **Work Data**

**Learn & Collaborate**

Read Only / Read Write

**Storage, Computing resources & Network** *(e.g. S3/iRods, L3VPN on Renater, CPU/GPU)*

**Provenance, Metrics & Accounting** *(e.g. RO-Crate)*

**Servicedesk**

Gaia data timeline allows to co-develop with EOSC EU node and third party (EOSC Federation)

Data Terra services designed to support EOSC EU Node and third party core services

**GAIA Data**

Development phase — June 2021 — May 2027 — Exploitation phase — May 2029

**eosc FAIR-EASE**

Sept 2022 — Aug.2025

**EOSC Federation (EOSC EU node and other nodes)**

- Identity Management
- Monitoring and Accounting
- Service Management System
- Resources Catalogues and Registry Services
- Application Workflow Management
- User Space

# Integrated Earth System Observation - Data Terra
## a Research Infrastructure to access, process and combine data

**MAIN MISSION**

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.

€42m (2020)

+1000 products & services

+15,000 users

100,000 TB (2022/2023)

## Partners

French scientific organizations and universities

## 5 thematic hubs

- Atmosphere — AERIS
- Ocean — ODATIS
- Solid earth — ForM@Ter
- Land Surfaces — Theia
- Biodiversity — PNDB (Pôle National de Données de Biodiversité)

## Cross-cutting device

Satellite images — DATA TERRA DINAMIS

- 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

**OUR SERVICES**

## Facilitating the cross-referencing of observations and the modeling of Earth System data

The IR Data Terra offers services around Earth system observation data. The objective is to provide interoperable and interdisciplinary services at all levels.

Connection with producers of data

1 Data Access

2 Daily production of data

3 On-demand analysis and processing

4 User support services

Software sharing, Analysis platform, Model evaluation

DATA TERRA

# wednesday session

# Integrating New Data Storage and Access Paradigms in pilot's practices

**"S3 Testbed Tests" Antoine Mahul, David Sarramia, Damian Smyth**

# Galaxy deals with remote compute resources and remote files system

# UCA TESTBED INFRASTRUCTURE

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)

Ressources allouée pour FAIR-EASE actuelles: 44 vCPUs, 152 GB RAM, 128 GB SSD, 100 TB via API S3

Cluster kubernetes FE : kubernetes v1.31, déployéed withTerraform et Talos Linux (https://www.talos.dev/)
 3 vm for control plane
 3 vm for workers qui vont executés les instances services client (JupyterHub/JupyterLab)

# 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**.

- Schema evolution supports add, drop, update, or rename, and has no side-effects
- Hidden partitioning prevents user mistakes that cause silently incorrect results or extremely slow queries
- Partition layout evolution can update the layout of a table as data volume or query patterns change
- Time travel enables reproducible queries that use exactly the same table snapshot, or lets users easily examine changes
- Version rollback allows users to quickly correct problems by resetting tables to a good state

**Reliability and performance**: Iceberg **was built for huge tables**. Iceberg is used in production where a single table can contain tens of petabytes of data and even these huge tables can be read without a distributed SQL engine.

- Scan planning is fast -- a distributed SQL engine isn't needed to read a table or find files
- Advanced filtering -- data files are pruned with partition and column-level stats, using table metadata
Iceberg was designed to solve correctness problems in eventually-consistent cloud object stores.
- Works with any cloud store and reduces NN congestion when in HDFS, by avoiding listing and renames
- Serializable isolation -- table changes are atomic and readers never see partial or uncommitted changes
- Multiple concurrent writers use optimistic concurrency and will retry to ensure that compatible updates succeed, even when writes conflict

Open standard designed and developed to be an open community standard with a specification to ensure compatibility across languages and implementations. Apache Iceberg is open source, and is developed at the Apache Software Foundation.

# Integrating New Data Storage and Access Paradigms in pilot's practices

**"Open Data Lake House Format - Apache Iceberg - Quick start" Damian Smyth:**

Joint meeting with D4T2 on New Data Storage and Access Paradigms on July 12th

https://fair-ease.atlassian.net/wiki/spaces/FAIREASE/pages/400818177/2024-07-12+Meeting

# EXAMIND

The Examind software suite, developed by Geomatys, provides all the features you need t**o 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 the open-source platform / server of the Examind ecosystem developed by Geomatys. This map server has a wide range of functions available, 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.), and has a number of ready-to-use processes for various uses in the geospatial world. Examind also 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. It operates standalone and is based on the Apache SIS geospatial processing library.

Apache Iceberg is an open table format designed to manage large-scale analytic datasets reliably and efficiently, making it ideal for democratized data infrastructures like those used in research and Earth Sciences. It provides versioning, schema evolution, and ACID compliance, enabling robust and transparent data management across diverse platforms. For infrastructures such as the ones you've described, Iceberg can ensure data access, support for distributed query engines, and adaptability to regional and cloud-based resources, empowering collaborative and FAIR data-driven science.

## 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.

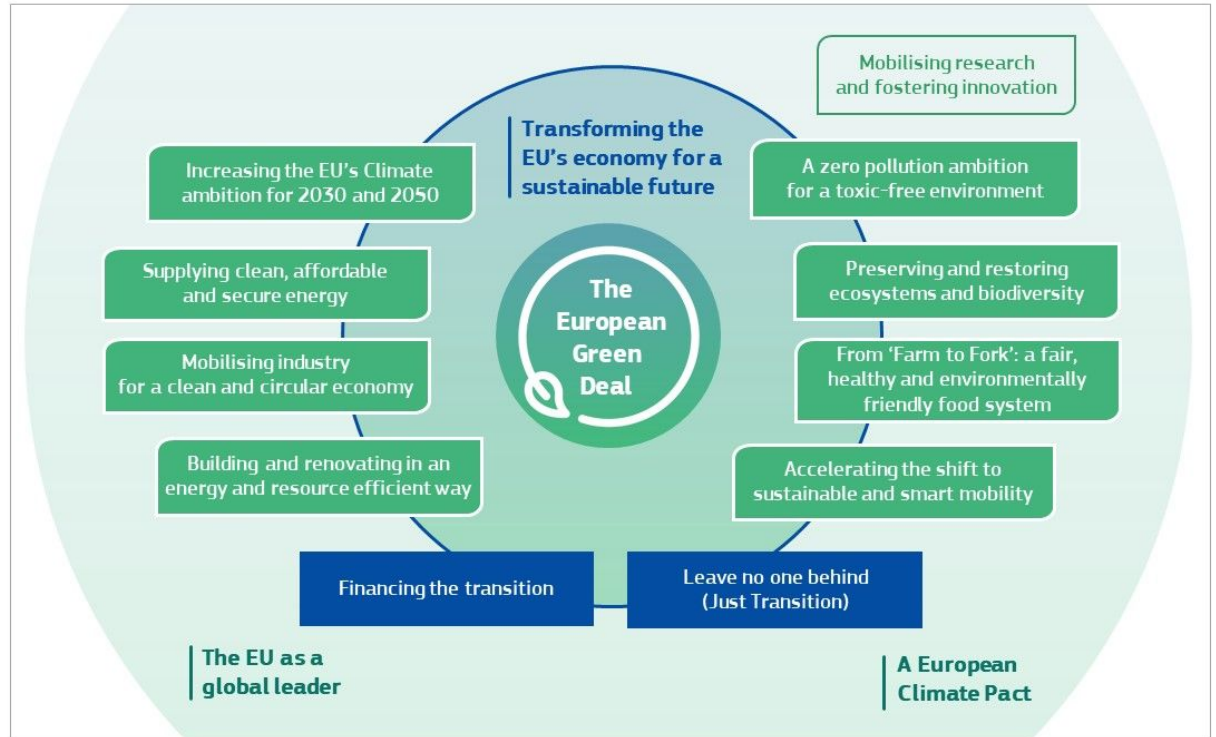| Feature | Apache Iceberg |
|---|---|
| **Primary Use Case** | Data lake management with strong ACID compliance |
| **\*ACID Transactions** | Fully supported, scalable for large datasets |
| **Schema Evolution** | Rich support for evolving schemas without rewriting data |
| **Partition Handling** | Hidden partitioning; avoids explicit user-defined partitions |
| **Integration** | Compatible with Spark, Trino, Presto, Flink, Hive, etc. |
| **Streaming Support** | Experimental, gaining traction |
| **Community** | Open, strong Apache ecosystem |
| **Adoption** | Used by Netflix, Apple, LinkedIn |

\* In computer science, **ACID (atomicity, consistency, isolation, durability)** is a set of properties of database transactions intended to guarantee data validity despite errors, power failures, and other mishaps

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

blablabla

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*

- *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-sur-les-iles-tonga-vue-par-les-satellites-geostationnaires/

- *2.1.7.*
- *Preserving and restoring ecosystems and biodi...*
- *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/le...ur-les-iles-tonga-vue-par-les-satellites-geostationnaires/



© Thomas Haessig