IBERGRID 2024
28-30 OCT
UNIVERSITY OF PORTO
better software for better science
13TH IBERIAN GRID CONFERENCE

# In situ processing of medical imaging data
## (An Open Challenge Experience)

Andrei S. Alic
Ignacio Blanquer
Pau Lozano
Damià Segrelles-Quilis

Institute of Instrumentation for Molecular Imaging
Universitat Politècnica de València

# Accelerating the lab to market transition of AI tools for cancer management: CHAIMELEON Project



Breast

Lung

Prostate

Colo / rectal

Images + Related
clinical data (e-form)

A Cloud-based cancer imaging repository as an online resource for the AI community working on the development of cancer management solutions

Not just a data warehouse…
- Incorporating all necessary functionalities to allow AI experimentation on the cloud (without downloading the data).
- Powered with automation tools.
- Interoperable with other existing initiatives.

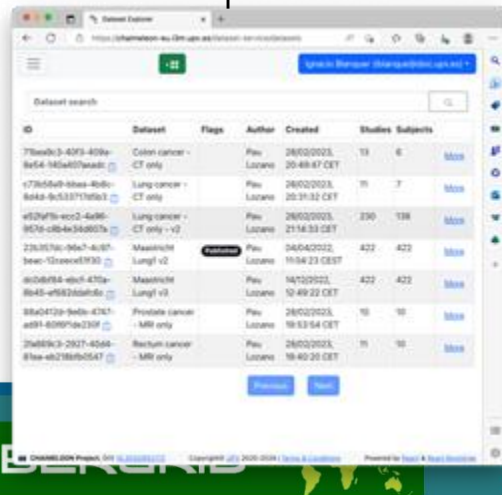# Use Cases and Design Principles

- Use Cases
  - Create, publish and explore a Dataset.
  - Access to the Images and clinical data from a Dataset through an interactive application.
  - Submit a processing job on the dataset to the infrastructure.
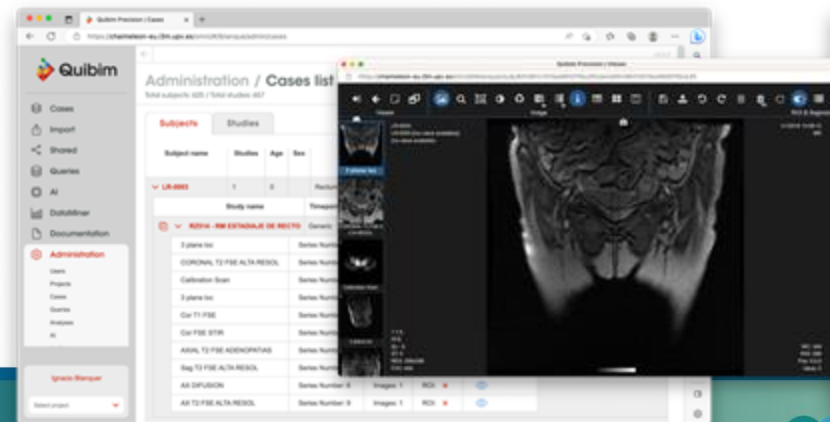  - Publish a processing application on the Marketplace.

- Design principles
  - The data is anonymised and uploaded to a central repository.
  - Data cannot be downloaded.
  - Data is organized into datasets (a coherent set of annotated image studies and the associated clinical data that have a persistent identifier - a FAIR citable research object).
  - Published Datasets have their metadata publicly accessible.
  - Data can be processed "in situ" using the tools available in the platform on the cloud resources of the platform.
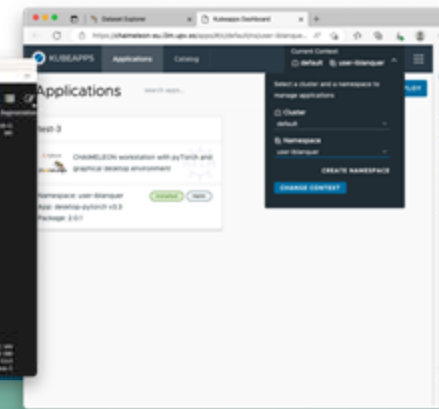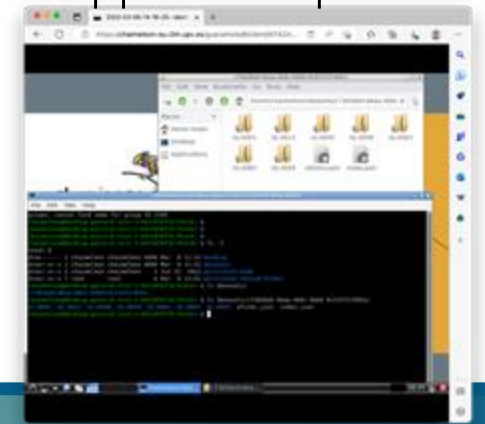
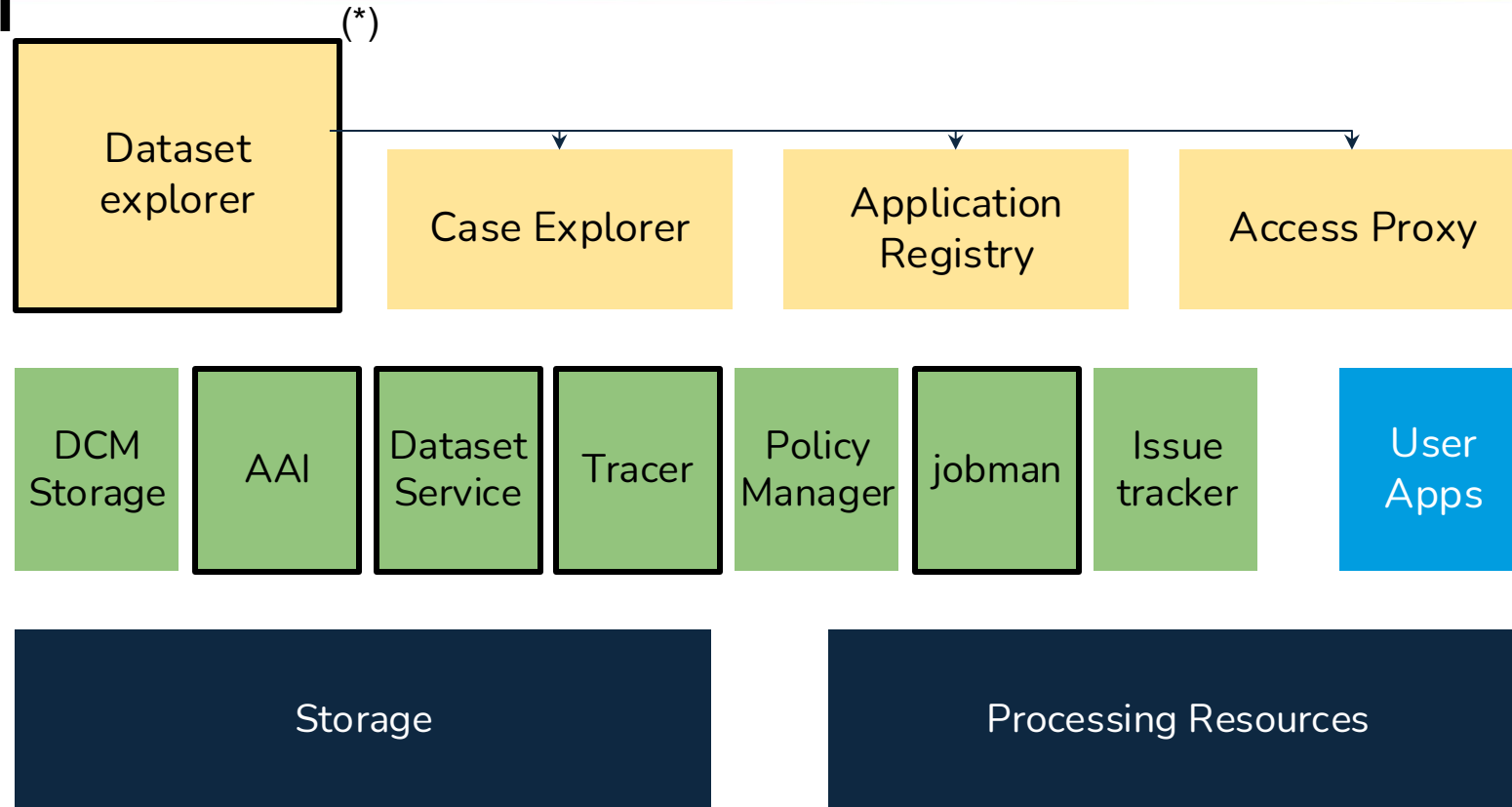Dataset Explorer          Case Explorer          Application Dashboard     App Desktop Access

# Central Platform

- All services managed as Kubernetes deployments and IaC recipes.

- **Secure**. All services are secured and access to data is restricted to the boundaries of the platform.

- **Auditable**. Access to data is registered in a Blockchain.

- **FAIR compliant**. It provides Findability, Accessibility, Interoperability and Reusability for the datasets.

- **High-Performance**. Integrates GPUs.

- **Reliable**. It uses Kubernetes deployments.

- **Reproducible**. It uses IaC, Open Source technologies, containers and cloud backends.

(*)

| Dataset explorer | Case Explorer | Application Registry | Access Proxy |
|---|---|---|---|

| DCM Storage | AAI | Dataset Service | Tracer | Policy Manager | jobman | Issue tracker | User Apps |
|---|---|---|---|---|---|---|---|

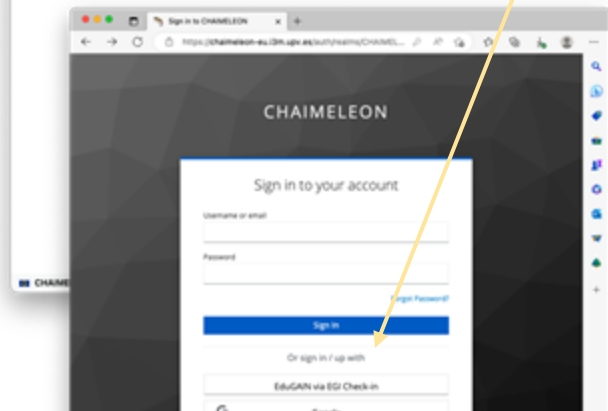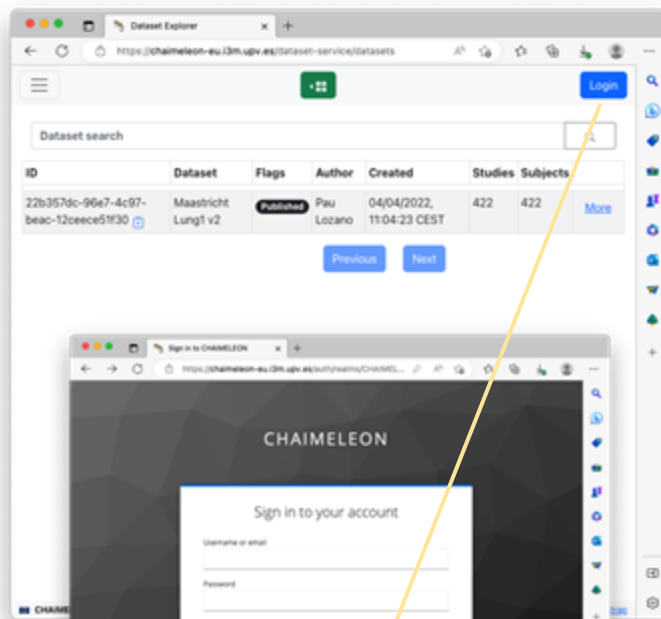| Storage | Processing Resources |
|---|---|

A Kubernetes elastic platform on top of a Cloud-based platform with 128 cores, 1,5 TB RAM, 3 V100 and 4 A30 GPUs (192 GB of GPU RAM)

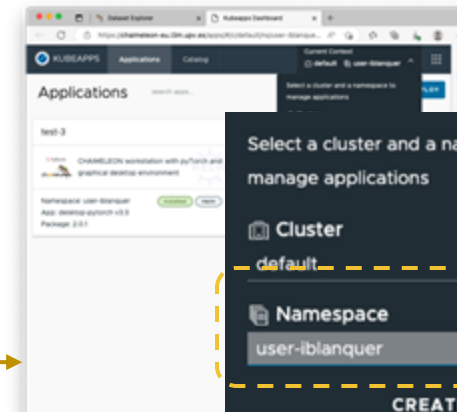*(*) Fully developed in CHAIMELEON*

# Authentication and Authorisation

https://chaimeleon-eu.i3m.upv.es/auth/realms/CHAIMELEON/account/#/



KubeAuthoriser
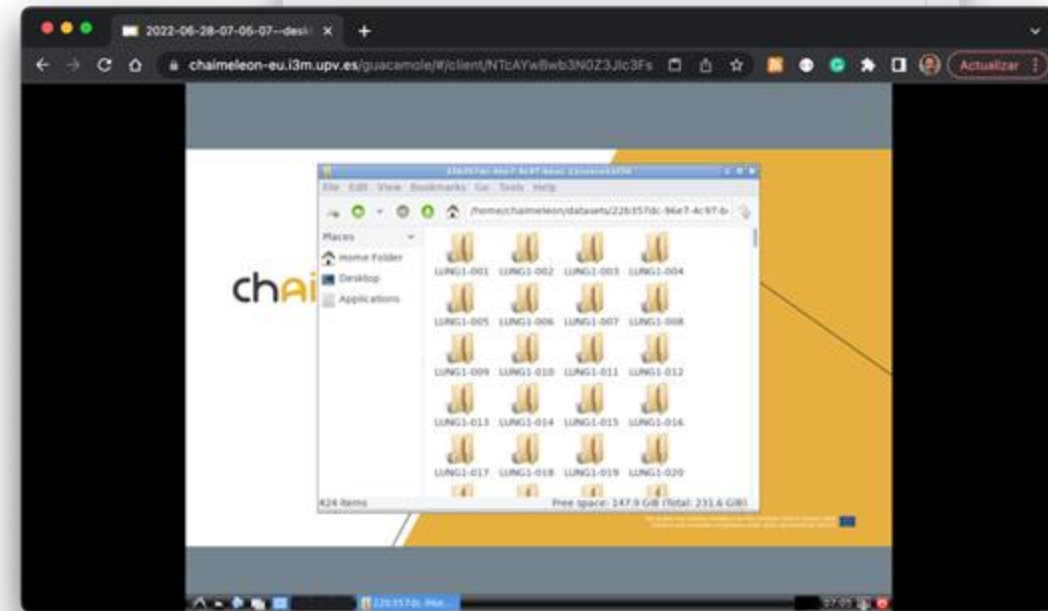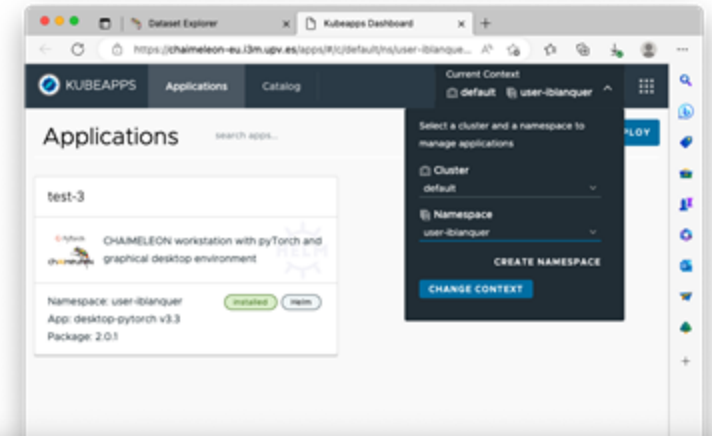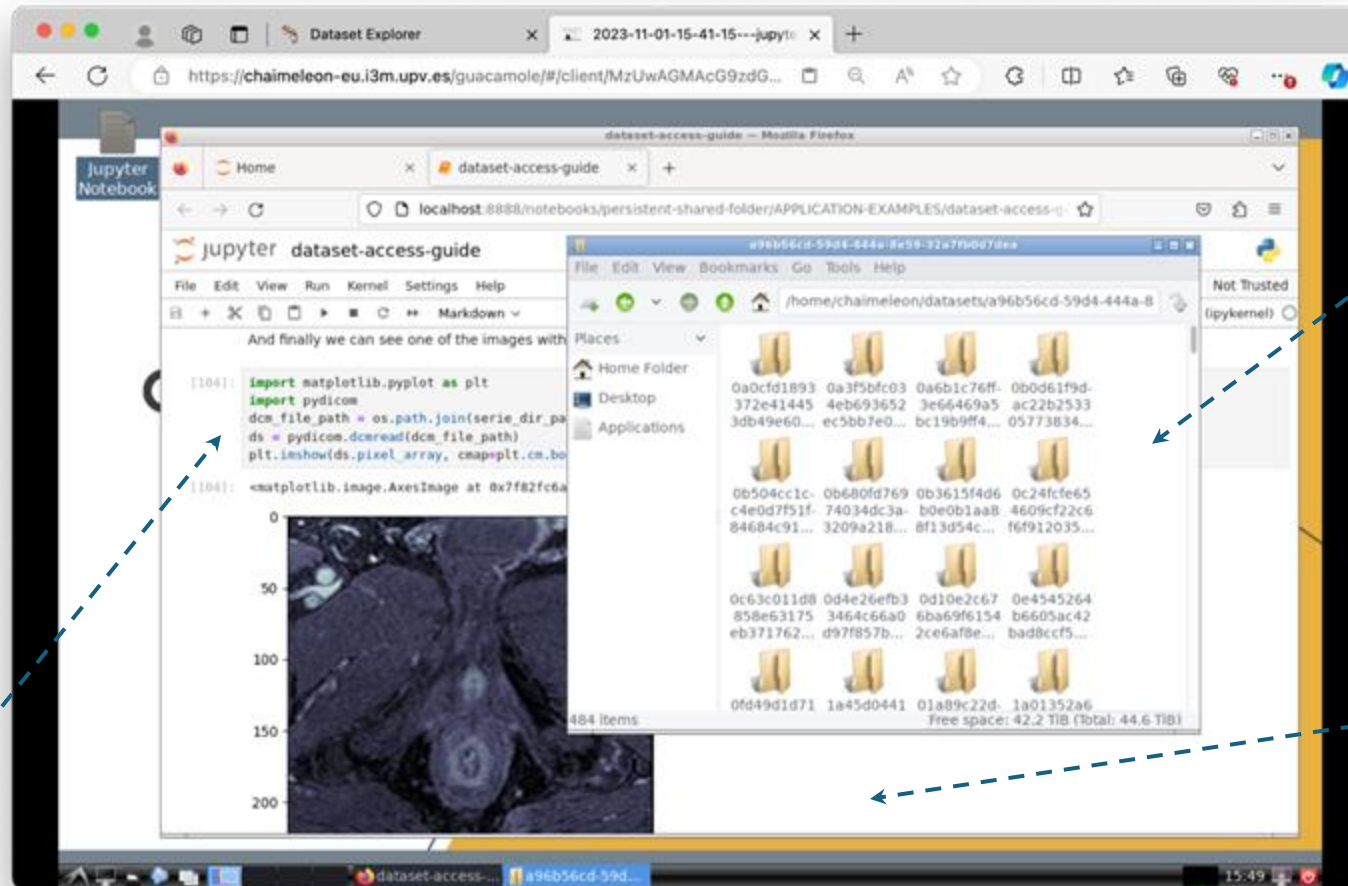
# Processing Model

- Processing is performed in situ by means of interactive and batch applications.

- Interactive applications
    - Accessed through a VNC proxy, to avoid the risk of downloading data.
    - Coded as Helm charts registered in KubeApps.

- Batch applications
    - Run unattended from the interactive virtual environments through a command line.
    - Run seamlessly in the Kubernetes infrastructure, accessing the data in the same way as the interactive applications.
    - Use a collection of verified Docker containers from an internal repository.

# Virtual Research Environments



A GUI with an ubuntu container in a network-restricted environment

Data Analytic SW libraries

Mounts the studies of the datasets as a POSIX volume.

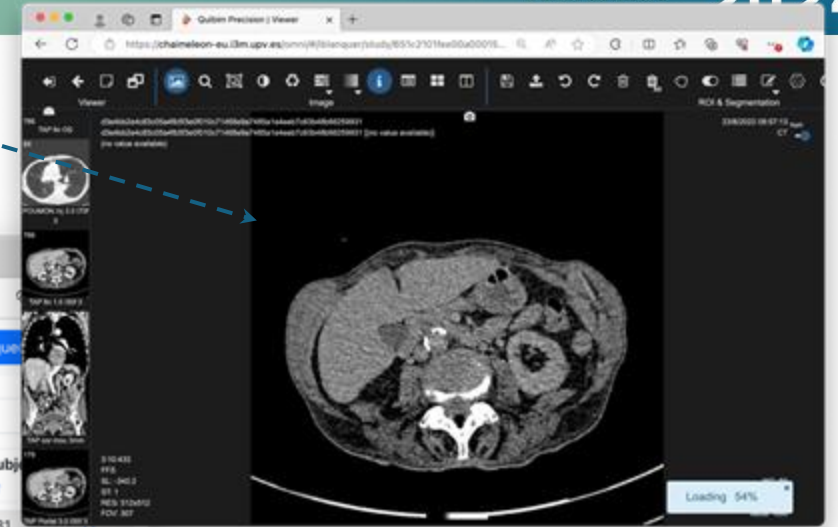Link to a batch queue system with GPUs and powerful resources.

# Dataset explorer



Explore images and DICOM metadata.

Links to the Application Catalogue, Case Explorer and Application access proxy

List of the available collections, with PIDs and basic metadata

Explore additional metadata and access logs.

# FAIR Compliance and Traceability

- A Dataset is a FAIR citable research object
  - Datasets have a metadata that contains aggregated information, which could be made public (just metadata).
  - Published Datasets have their metadata publicly accessible.

- Traceability is provided through a Web service
  - Logs users' actions (create / update / use datasets) in the CHAIMELEON repository (the traces)
  - Stores traces in private blockchain(s) (support for BigchainDB & Hyperledger Besu)
  - Anti-tampering, redundancy, and distribution of the complete set of traces
  - It does not store repository's users/patients private/sensitive information
  - API @ https://app.swaggerhub.com/apis/UPV-CHAIMELEON/Tra...

# Logging workloads

- Along with the usage of the datasets (through the Interactive Applications), access through the batch processing tools is also registered.

- Detailed information is persisted on a Blockchain
  - Dataset creation (User and Type)
  - Dataset Access (user, environment, start and end)
  - Job (Execution time, End Status, Job Type, Command)

# Testing the platform – Open Challenges



https://chaimeleon.eu/open-challenges/

- Five clinical questions
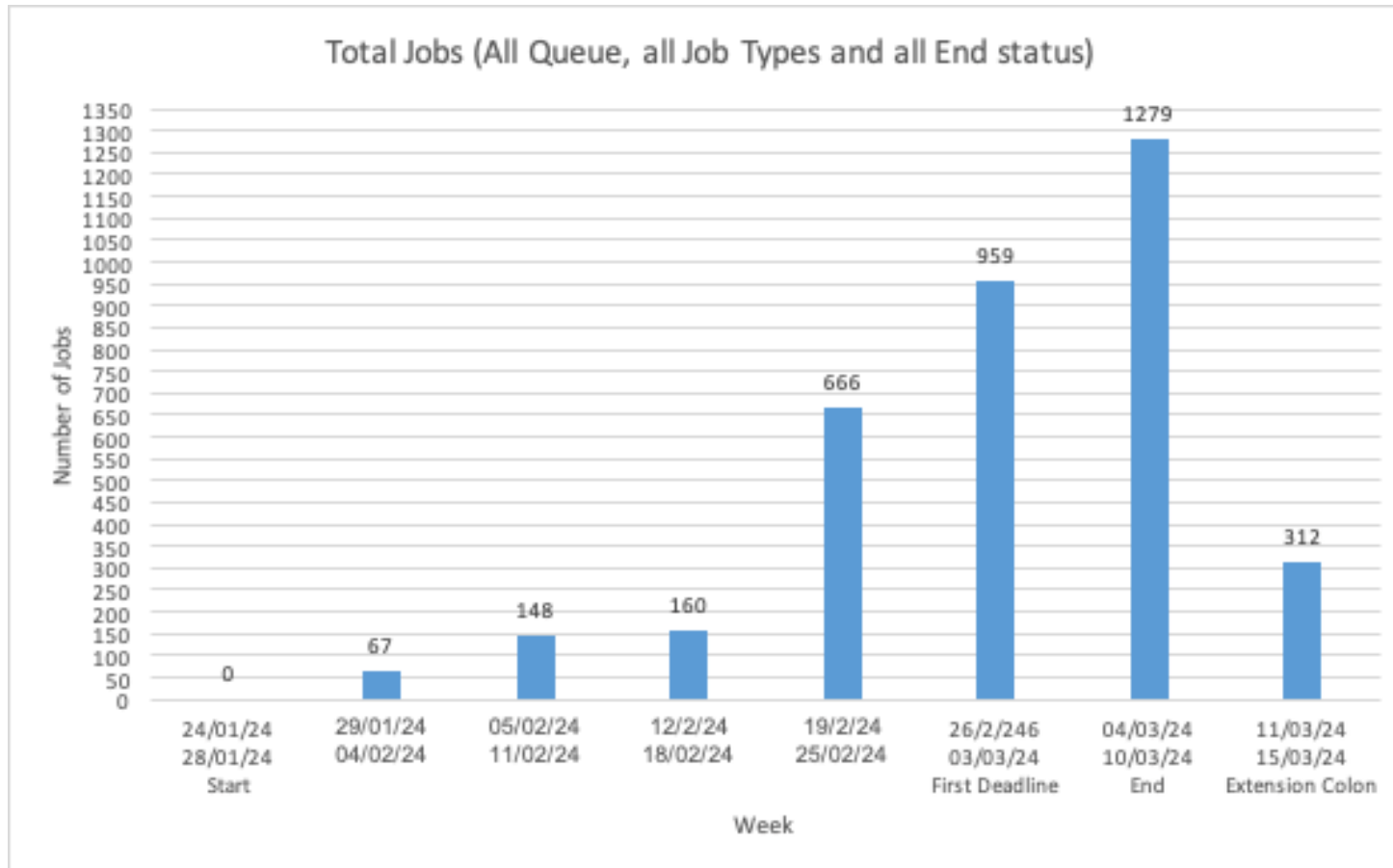  - Related to Lung, Prostate, Colon, Breast and Rectal Cancer

- Two-stage championship
  - Classification phase
    - 161 participants sent the organizers the T&C signed (240 got interest).
    - Provided with synthetic data for Prostate and Lung cancer.
  - Championship pase
    - 30 Qualified, 20 active participated.
    - In-situ Access.
    - Access provided for the full 5 datasets.

| VMs | Description | Jobs total | GPU | VRAM | Cores Total | RAM total | Disk for docker images | Disk for jobs containers | Total ephemeral disk |
|---|---|---|---|---|---|---|---|---|---|
| 7 | (2022.XLarge)<br>for the master node, core services | - | - | - | 56 | 224 GB | 280 GB | 280 GB | 560 GB |
| 3 | (2022.XLarge)<br>for the CEPH master and storage nodes | - | - | - | 24 | 96 GB | 120 GB | 120 GB | 240 GB |
| 5 | (2023.Chaimeleon.XLarge)<br>for desktops | 40 | - | - | 40 | 320 Gb | 400 GB | 300 GB<br>(60 each job) | 700 GB<br>(140 each VM) |
| 5 | (2023.Chaimeleon.XLarge-V100)<br>for large-gpu jobs | 5 | 5 x<br>V100 32GB | 5x32 GB | 40 | 320 GB | 160 GB | 300 GB<br>(60 each job) | 700 GB<br>(140 each VM) |
| 2 | (2023.Chaimeleon.XXL128-A30)<br>for small-gpu jobs | 8 | 2 x<br>A30 24GB | 8x6 GB | 32 | 256 GB | 160 GB | 360 GB<br>(45 each) | 520 GB |
| 2 | (2023.Chaimeleon.XXL128-A30)<br>for medium-gpus jobs | 4 | 2 x<br>A30 24GB | 4x12 GB | 32 | 256 GB | 160 GB | 360 GB<br>(60* each job) | 520 GB |
| 1 | (2023.CHAIMELEON.Ceph)<br>for Ceph persistent home | - | - | - | 8 | 128 GB | 80 GB | - | 80 GB |
| **25** | | **57** | | | **232** | **1.600 GB** | **1.360 GB** | **1.720 GB** | **3.320 GB** |

# Total VCPU Capacity: 284.801,93 VM/hours for the Open Challenge period

# Open Challenge workload



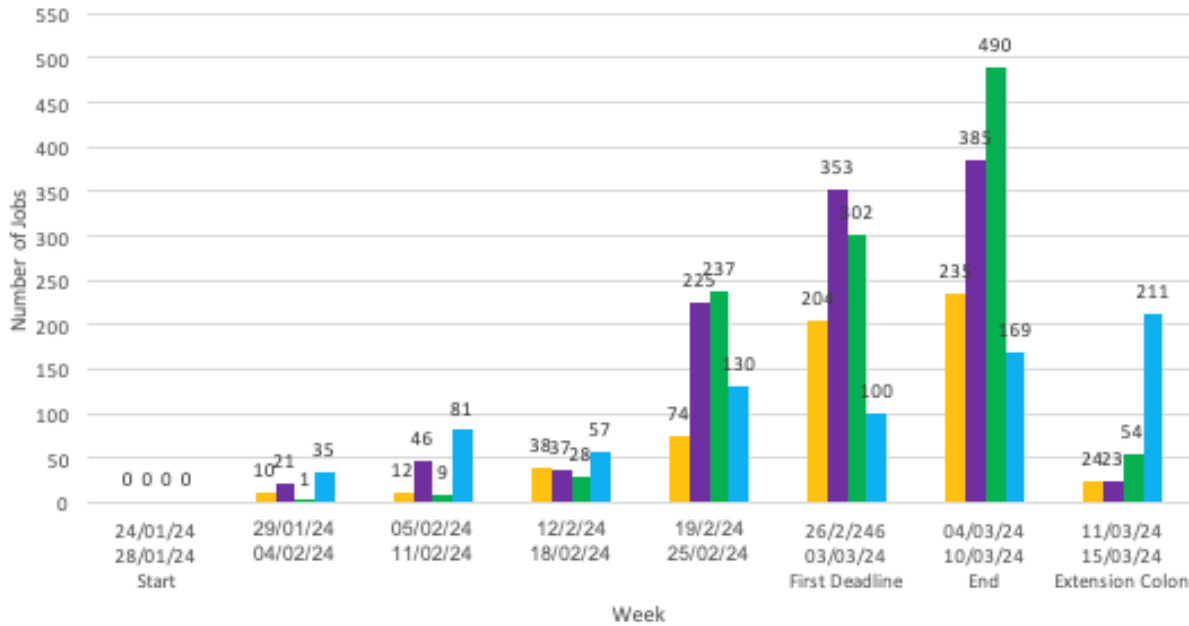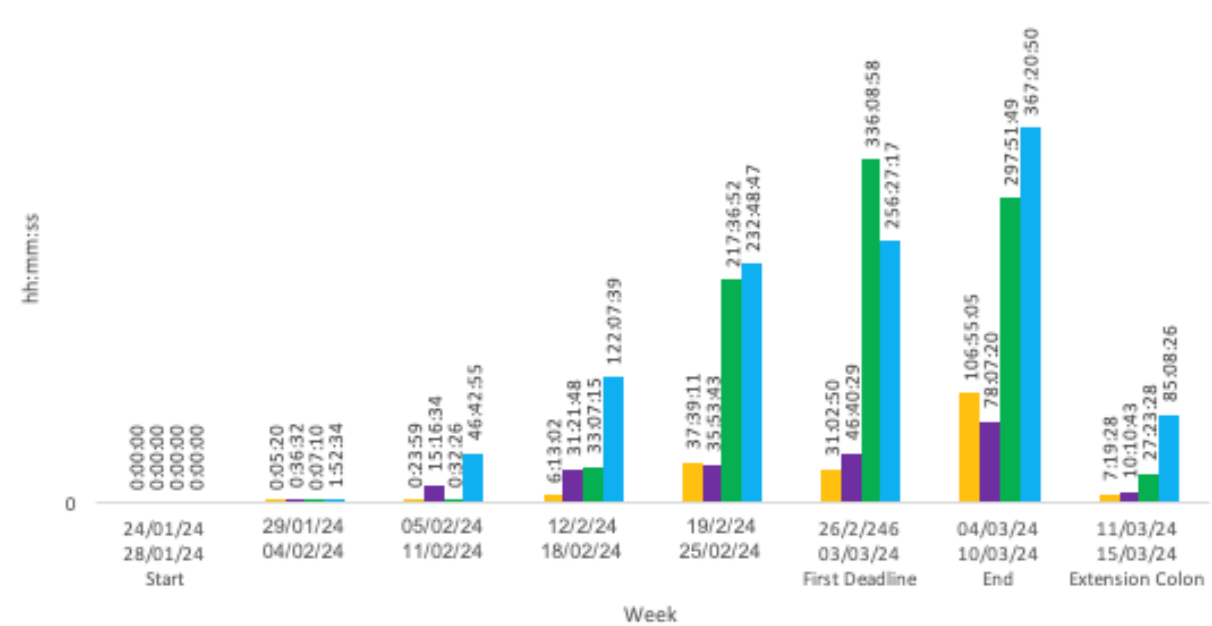Total Jobs (All Queue, all Job Types and all End status)

- **A total of 3.591 Batch Jobs** have been submitted during the championship phase. These corresponds to more than **2.433 hours.**

- The increase in weeks of batch jobs may be explained by the fact that users have gained experience with the platform and the closeness to the championship deadline.

- Significant increase after weeks 4 and 5.

# Open Challenge workload

- GPUs resources are the ones with the highest number of submissions and running time.
- "large-gpu" has a lowest number of jobs but the highest computational time.
  - This can be explained by the fact that final training has been performed in the on "large-gpu" meanwhile preliminary adjustments have been done on "small-gpu" or "medium-gpu".

# Conclusions

**Strengths**
- The amount of resources has been over-dimensioned
  - We expected 40 users but 29 were selected and 20 were active.
  - The queueing time has been negligible.
  - We estimate that with the current computing resources, at least 60 users could have been supported.
- The platform has very few technical incidences
  - Just 1 resource had network issues which were solved without losing the running workload.
  - The platform support worked well.
- The number of jobs submitted was reasonably high (over 3.500 jobs, mostly on three weeks)

**Weaknesses**
- The learning curve is steep
  - Most of the complaints were addressed in the documentation (including those related to the data format)
  - We need to consider a proactive learning process and a better documentation, as well as a mandatory training period.
- The number of failed jobs is considerably high
  - Many of them, according to the information in the forum, were due to a wrong understanding on how the storage system works.