



**CS<sup>3</sup>  
MESH<sup>4</sup>  
EOSC**

**Connecting European Data**



# Federated collaborative workflows for Jupyter

Diogo Castro (CERN), Marcin Sieprawski (Software Mind)

[diogo.castro@cern.ch](mailto:diogo.castro@cern.ch), [marcin.sieprawski@softwaremind.com](mailto:marcin.sieprawski@softwaremind.com)



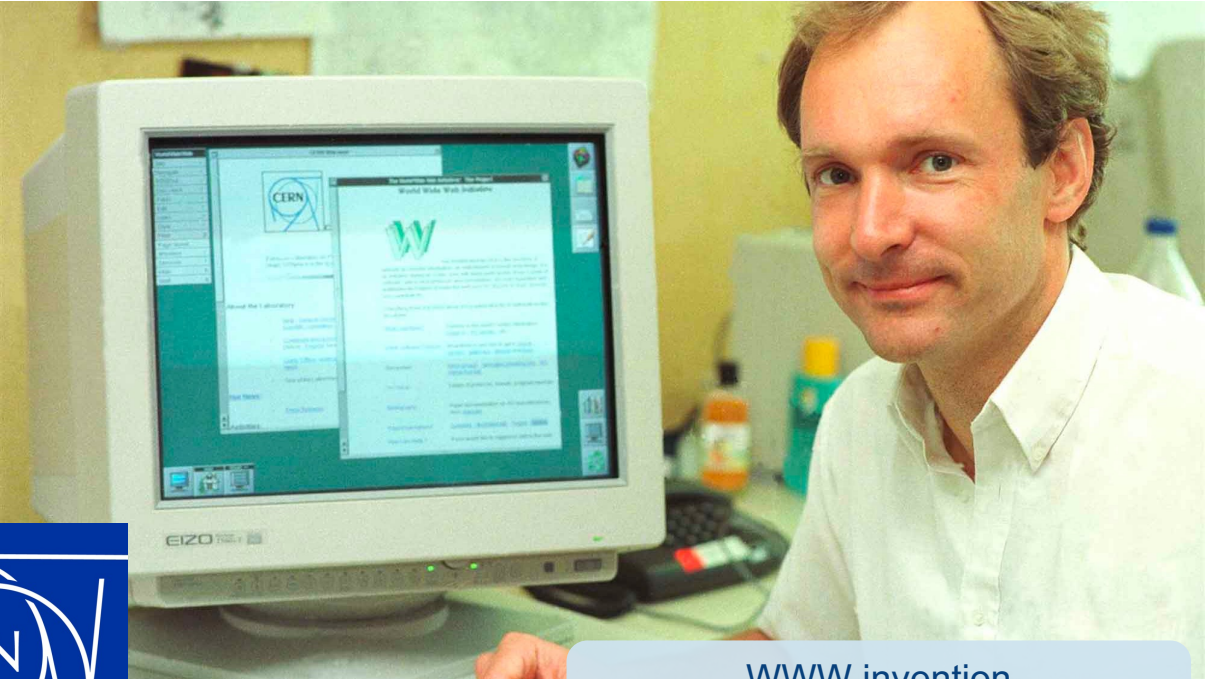
CS3MESH4EOSC has received funding from the European Union's Horizon 2020 Research and Innovation programme under **Grant Agreement No. 863353**.

Jupytercon 2023  
10/05/2023

# Introduction

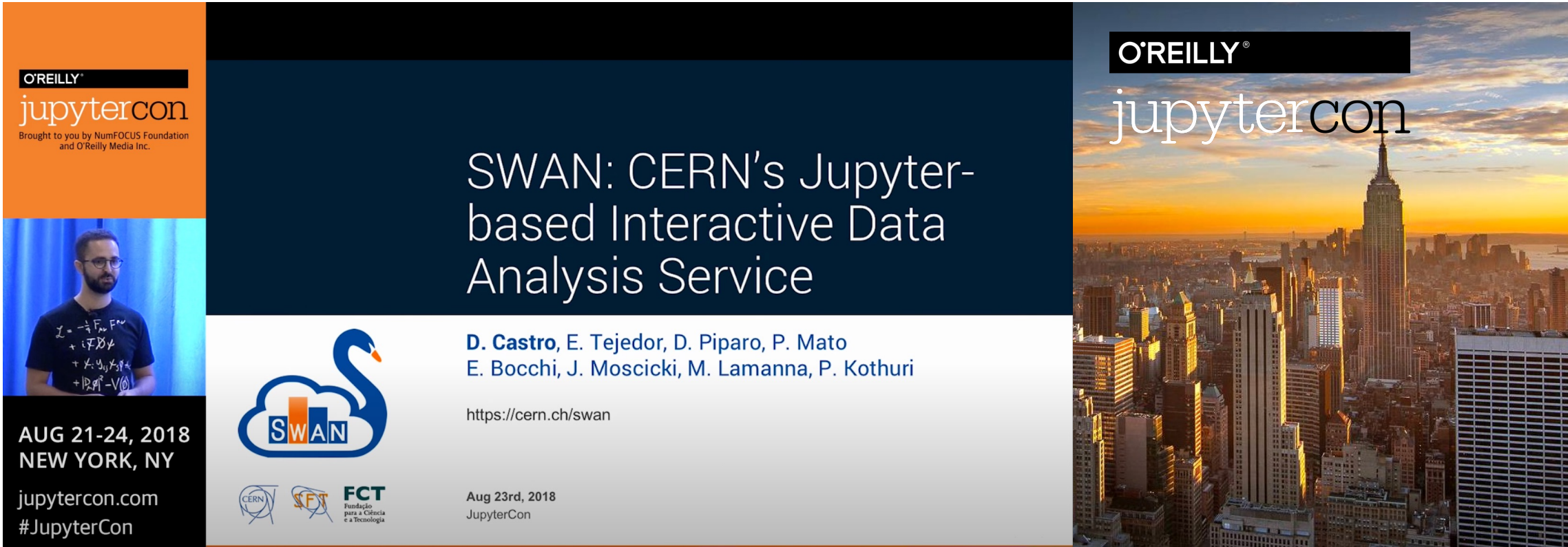


The Large Hadron Collider



WWW invention





**O'REILLY®**  
jupytercon  
Brought to you by NumFOCUS Foundation and O'Reilly Media Inc.




**SWAN: CERN's Jupyter-based Interactive Data Analysis Service**

**D. Castro, E. Tejedor, D. Piparo, P. Mato  
E. Bocchi, J. Moscicki, M. Lamanna, P. Kothuri**

<https://cern.ch/swan>


**AUG 21-24, 2018  
NEW YORK, NY**

[jupytercon.com](http://jupytercon.com)  
#JupyterCon

  
  **FCT**  
Fundação para a Ciência e a Tecnologia

Aug 23rd, 2018  
JupyterCon

**O'REILLY®**  
jupytercon



- # <https://conferences.oreilly.com/jupyter/jup-ny/public/schedule/detail/68359.html>
- # [https://www.youtube.com/watch?v=TDp\\_XlgtpDA](https://www.youtube.com/watch?v=TDp_XlgtpDA)

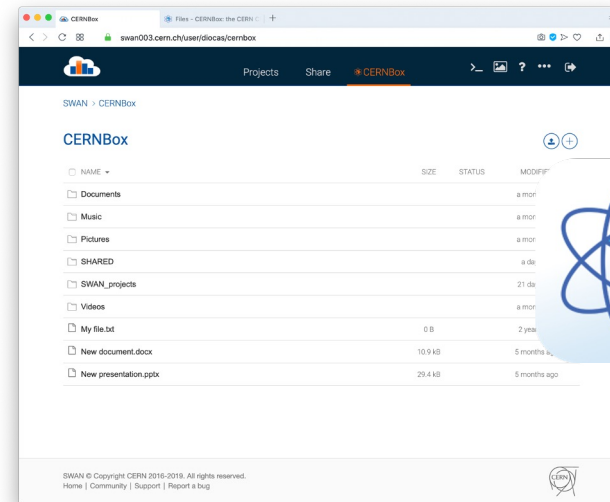


# All the data our users need for their analysis

- # CERNBox as home directory
- # Experiment repositories, projects, open data, ...

# Sync&Share

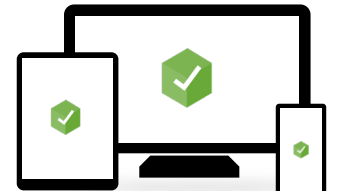
- # Files synced across devices and the Cloud
- # Simple collaborative analysis

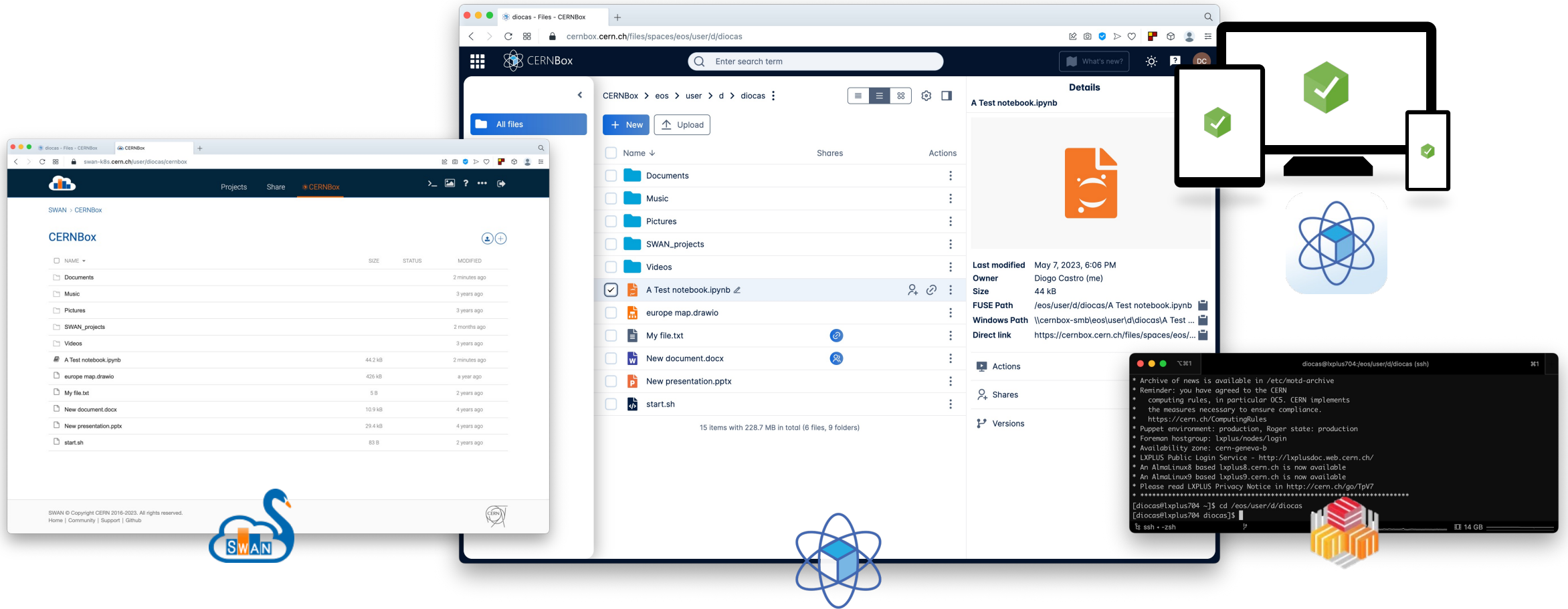


share



sync





The image displays three overlapping browser windows showing the CERNBox interface, illustrating a consistent user experience across different devices and contexts.

**Top Window (File List):** Shows the CERNBox interface with a search bar and a file list. The selected file is "A Test notebook.ipynb".

Name	Shares	Actions
Documents		
Music		
Pictures		
SWAN_projects		
Videos		
<input checked="" type="checkbox"/> A Test notebook.ipynb		
europe map.drawio		
My file.txt		
New document.docx		
New presentation.pptx		
start.sh		

15 items with 228.7 MB in total (6 files, 9 folders)

**Middle Window (Details):** Shows the details for "A Test notebook.ipynb".

**Details:**

- Last modified:** May 7, 2023, 6:06 PM
- Owner:** Diogo Castro (me)
- Size:** 44 kB
- FUSE Path:** /eos/user/d/diocas/A Test notebook.ipynb
- Windows Path:** \\cernbox-smb\eos\user\d\diocas\A Test ...
- Direct link:** https://cernbox.cern.ch/files/spaces/eos/...

**Bottom Window (Terminal):** Shows a terminal session with system boot logs.

```

diocas@lxplus704:/eos/user/d/diocas (ssh)
* Archive of news is available in /etc/motd-archive
* Reminder: you have agreed to the CERN
* computing rules, in particular OCS, CERN implements
* the measures necessary to ensure compliance.
* https://cern.ch/ComputingRules
* Puppet environment: production, Roger state: production
* Foreman hostgroup: lxplus/nodes/login
* Availability zone: cern-geneva-b
* LXPLUS Public Login Service - http://lxplusdoc.web.cern.ch/
* An AlmaLinux8 based lxplus8.cern.ch is now available
* An AlmaLinux9 based lxplus9.cern.ch is now available
* Please read LXPLUS Privacy Notice in http://cern.ch/go/IpW7
* .....
[diocas@lxplus704 ~]$ cd /eos/user/d/diocas
[diocas@lxplus704 diocas]$
ssh - -zsh
  
```

**Right Diagram:** Illustrates the concept of "The same view everywhere" by showing a laptop, a tablet, and a smartphone, all connected to a central server icon, representing a federated collaborative workflow for Jupyter.

### # SWAN acts as a client to other resources

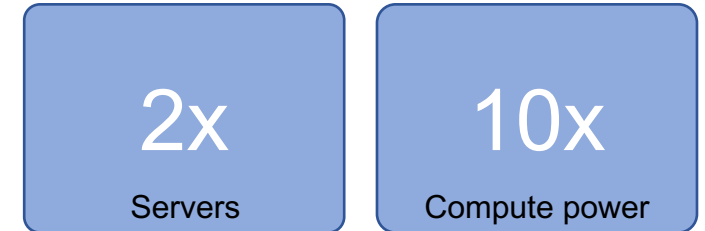
- # Allows support for both *single-node* and *distributed* analysis
- # Run lightweight analysis interactively, offload heavy computations
- # Storage as a shared layer across multiple services





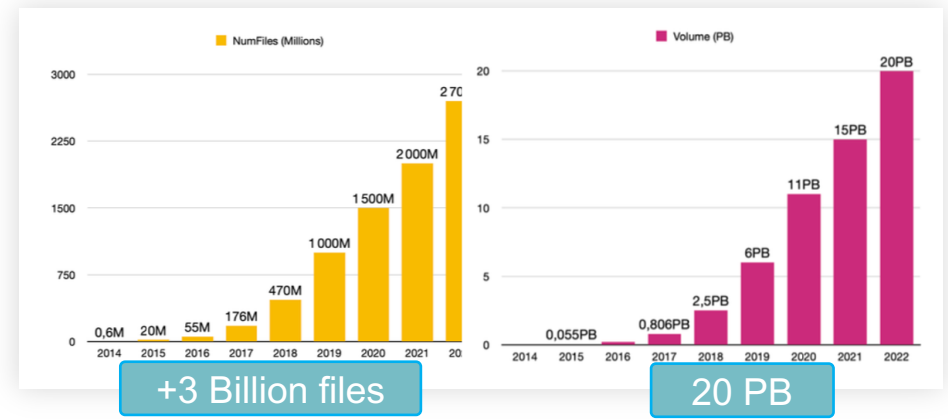
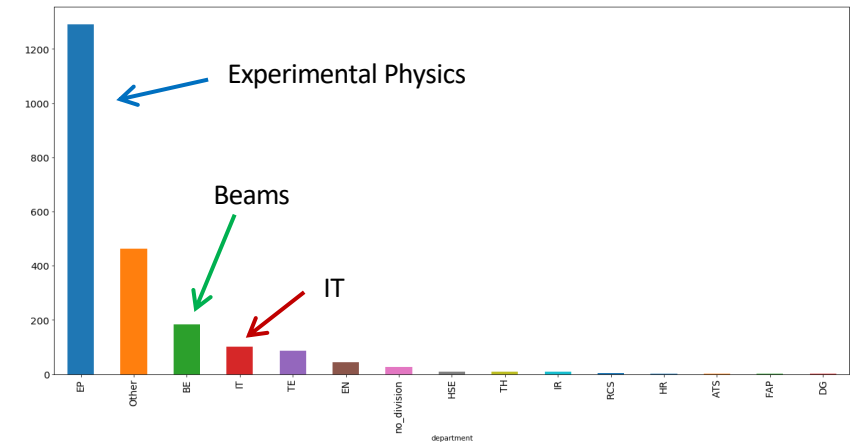
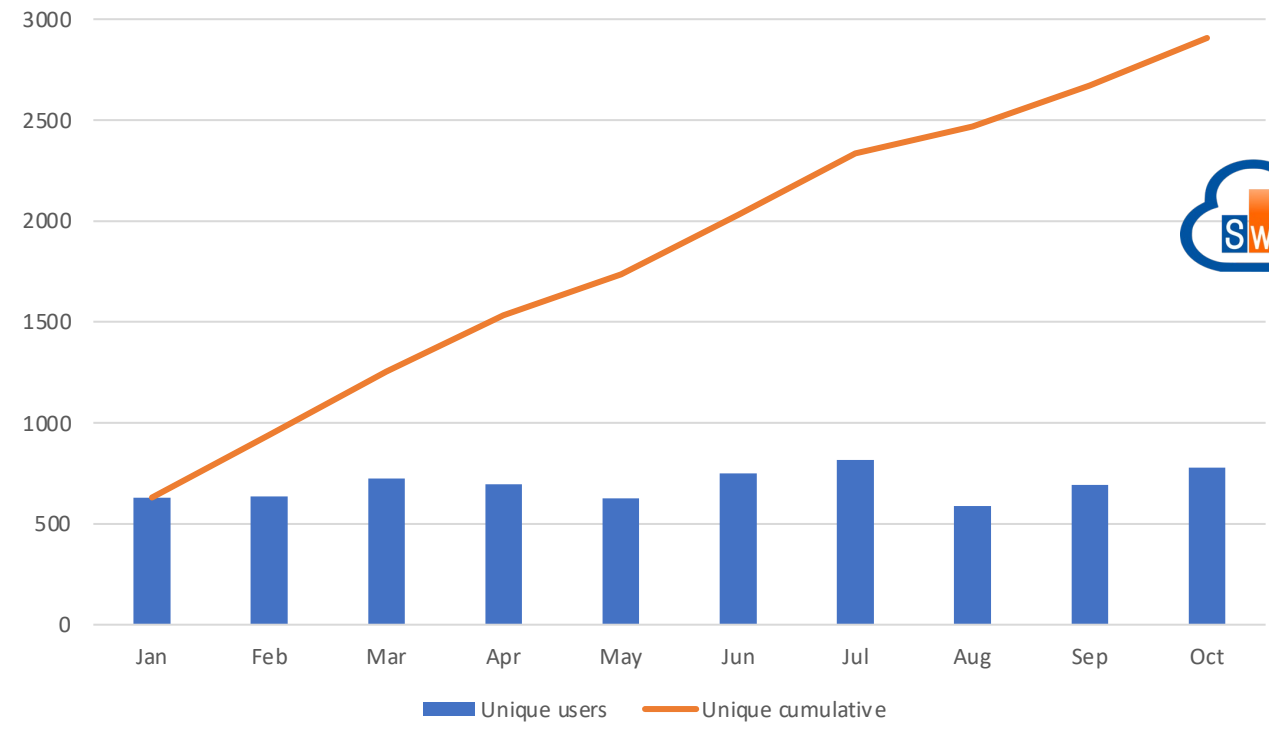
- # A major upgrade to LHC is coming
  - # Called Hi Luminosity (HL-LHC), planned for 2029
  - # 30x more data throughput
- # Major upgrade to analysis infrastructure is also required
- # HL-LHC needs are pushing us to build modern “Analysis Facilities”
  - # R&D effort across the whole physics community around the world (WLCG – worldwide LHC computing grid)

### Computing



### Storage





+3 Billion files

20 PB

- # Ability to perform fast research iterations on large datasets
- # Ability to convert interactive to batch-schedulable workloads
- # Ability to scale outside of the facility on occasion
- # Ability to reproducibly instantiate desired software stack
- # Ability to collaborate in a multi-organisational team on a single resource
- # ...

“In general, JupyterLab is the main choice when it comes to interactive analysis.”

CS3MESH4EOSC

### Commercial cloud providers



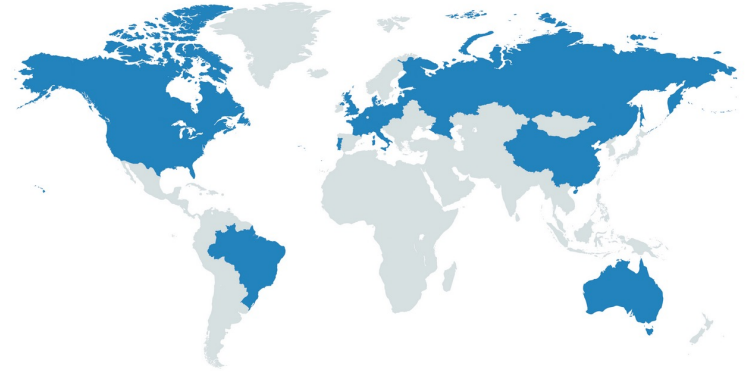
### Open source (self deployed) alternatives



\* Enterprise File Sync and Share

## Status quo

- # Many EFSS nodes, providing valuable **services** to the community
  - # **Mostly Sync and Share**, but not limited to that
- # User environments, higher level **applications**
  - # (e.g. editors, Jupyter data analysis...)
- # Basic **file sharing** possible

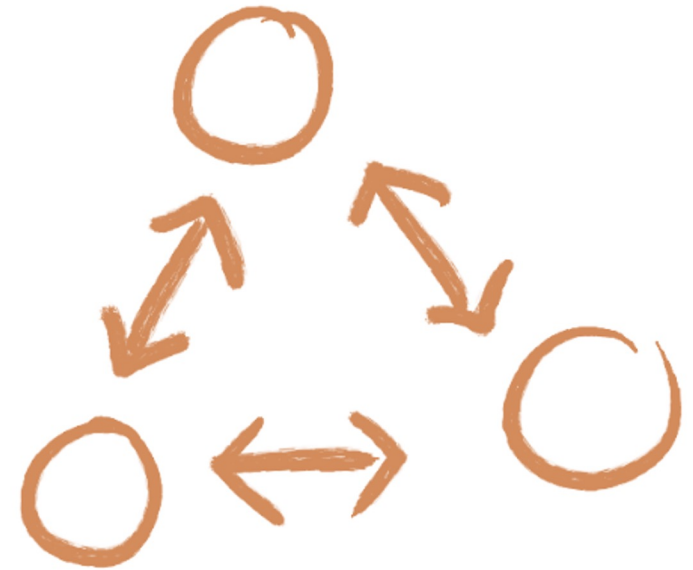


### CS3 Site Reports:

- # ~16 PB of data
- # > 20 nodes
- # > 400.000 users
- # > 3.5 billion files/dirs

## But...

- # Researchers remain **isolated on data islands** because these services aren't interconnected
- # **No common, ratified API**
  - # Hard to share add-ons between sites
  - # Hard to get traction with eScience community
- # **Suboptimal knowledge transfer** back to commercial and business environments.
  - # We can't make a joined-up front this way



## # 3-year Project

- # Ending in June 2023
- # Led by CERN

## # Objectives

- # Delivering a **Global Collaboration Service** for researchers, educators, data curators, analysts...
- # Providing an **interoperable platform** to easily share and deploy applications and software components
- # Leveraging the potential of the **CS3 Community** and expanding it

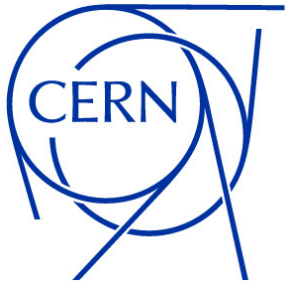
Jupyter is one of the main applications in the ecosystem

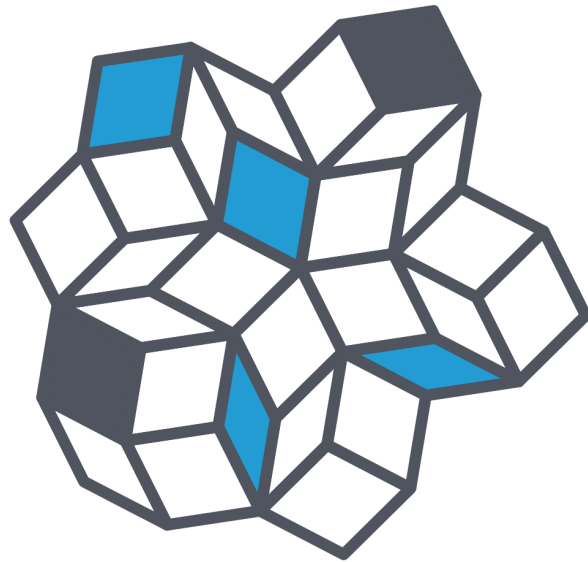


[www.cs3mesh4eosc.eu](http://www.cs3mesh4eosc.eu)





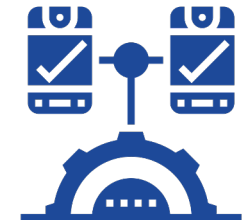
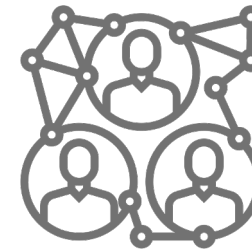




# Science Mesh

[www.sciencemesh.io](http://www.sciencemesh.io)

- # Decentralized **Mesh of EFSS nodes**
  - # Trusted federation of nodes
- # Based on **Open Standards** and **Open Source Software**
- # **Federated** environment where researchers can collaborate
  - # Usually on the the same UX as if they were local users
- # **Application Platform** for distributed collaborative tools



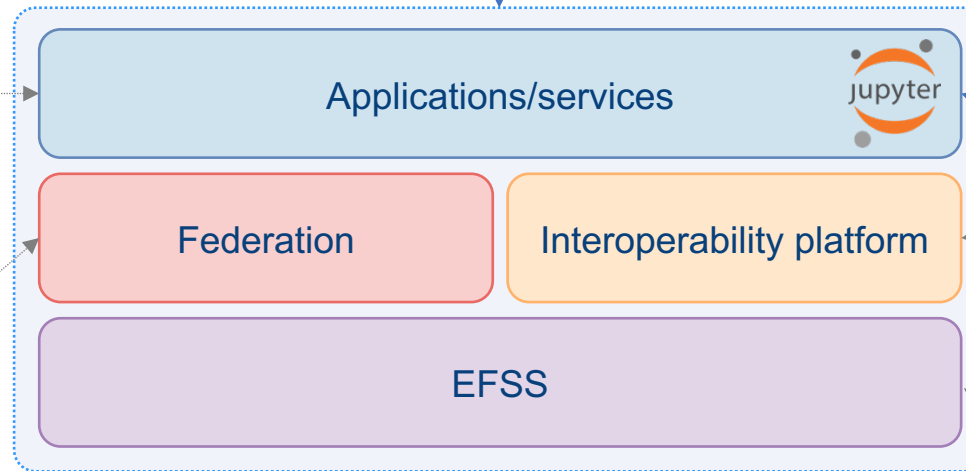


## Science Mesh



New domain-specific applications developed in the community

Build upon existing infrastructure and long-term efforts  
*eduGAIN, ...*



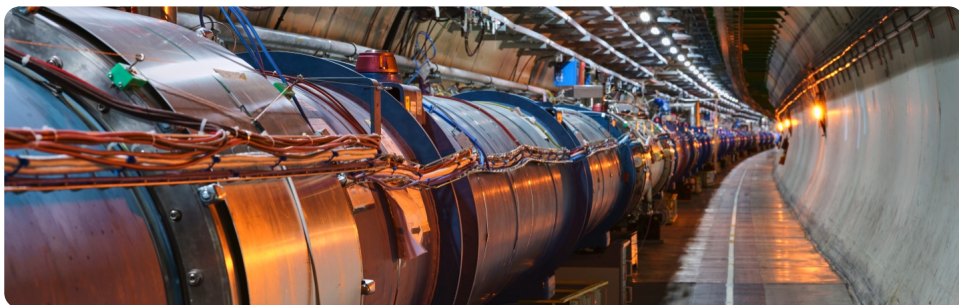
Some applications might need a connector  
*ScienceMesh plugin for Jupyter*

Lightweight add-on  
Easy to deploy and install new functionality  
*OCM, CS3APIs, Reva*

Connect to already deployed and commercially supported products  
*ownCloud, Nextcloud, seafile, ...*

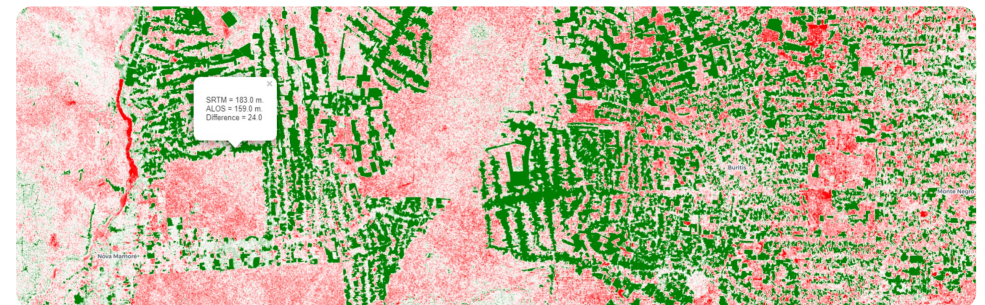
## # Physics

- # High Energy Physics analysis
- # Collaborative Analysis
- # LHC Accelerator logging and monitoring
- # Statistical and operational studies
- # Open data access
- # Education and outreach



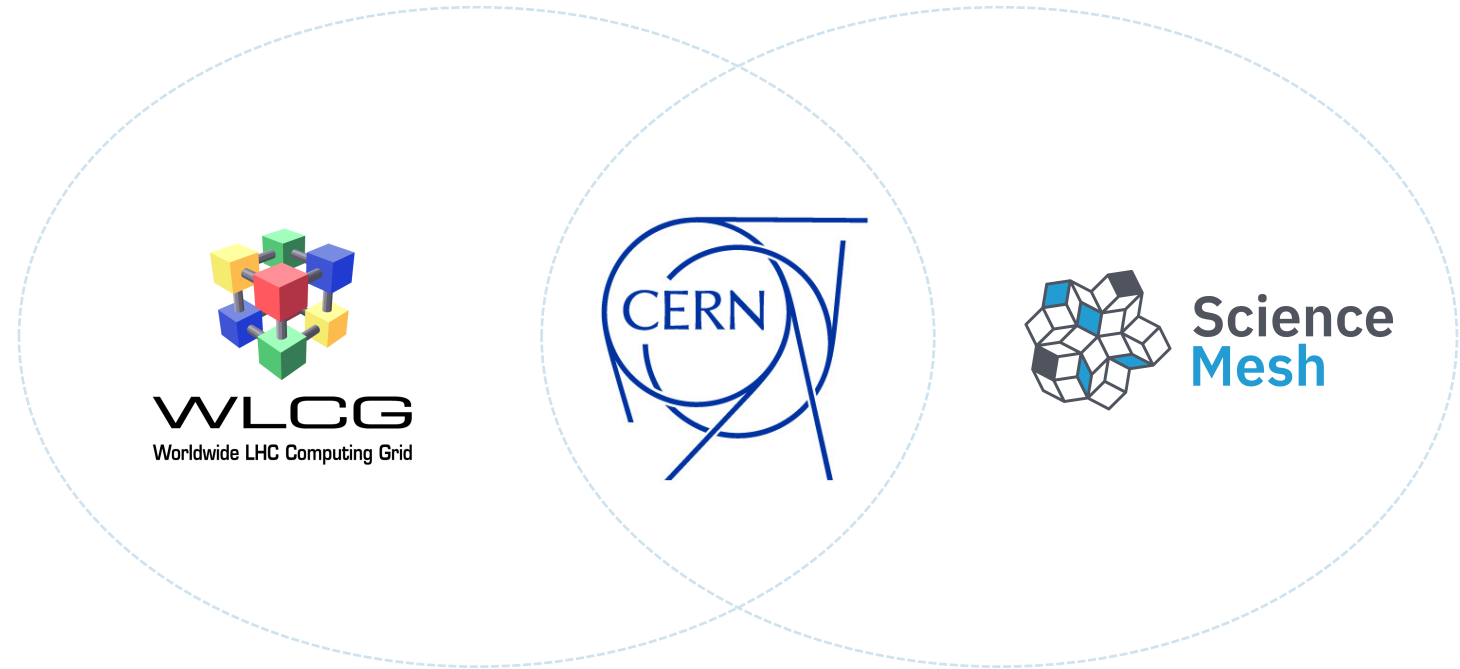
## # Earth Observation

- # Copernicus Earth Observation
- # Geo Visualization and Data Exploration
- # Interactive dashboards
- # Sustainable natural resources and water management (decision-making support); Land degradation (monitoring and assessment); ...



## # CERN sits at the intersection of these two communities

- # It has been trying to put it to work for the physics community
- # But this technology also has interest for long tail science or even education use cases outside of it



# ScienceMesh Plugin for Jupyter

- # In the beginning, notebooks could not be open in parallel
  - # Conflicts would happen, especially on shared filesystems
- # Now they can, and their data structures are synchronized
  - # This looks awesome!
  - # But optimal usage requires sharing the same Jupyter server and kernel (?)

### File Changed

"A Test notebook.ipynb" has changed on disk since the last time it was opened or saved. Do you want to overwrite the file on disk with the version open here, or load the version on disk (revert)?

Cancel

Revert

Overwrite



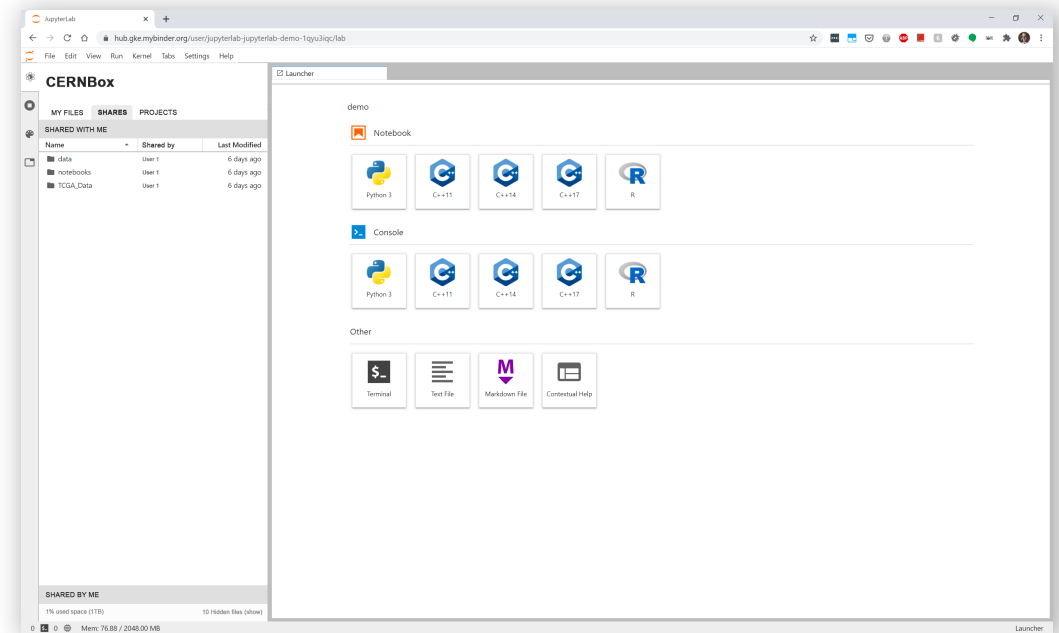
- # A shared filesystem might mean access from different Jupyter servers
  - # Or even other applications altogether
  - # The concurrent editing does not work fully
- # Collaboration requires coordination
  - # Which might not always be easy, especially if we don't know who is editing on the other side...
- # Sharing the same server + kernel is risky
  - # Full access to another user's account, storage and permissions on many resources
- # We're not aware of use cases that would benefit from true concurrent editing

We propose a complementary model better suited for large scale distributed environments

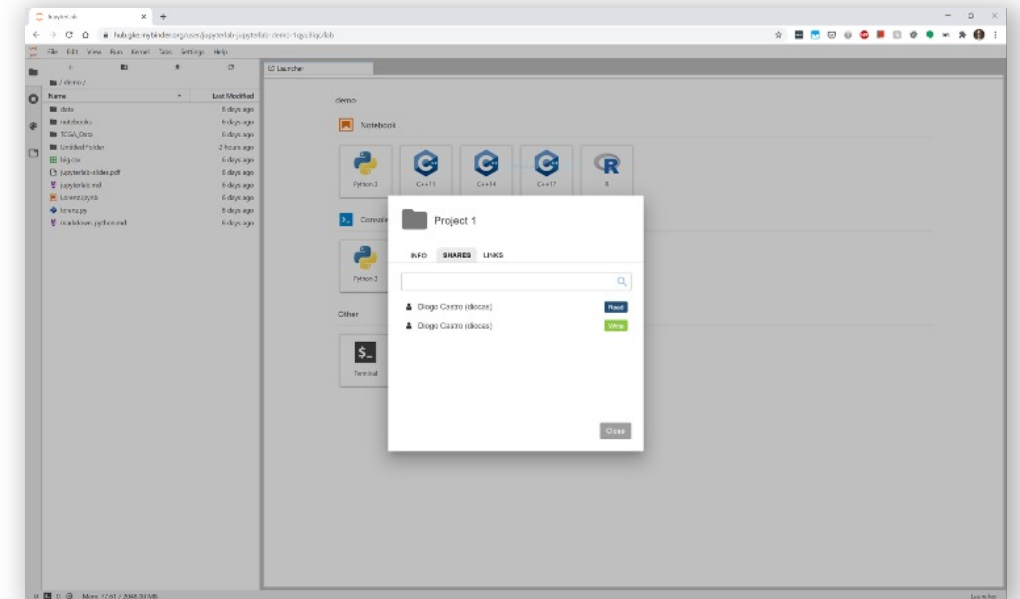
ScienceMesh Plugin for Jupyter

- # Connect to EFSS (and the Mesh)
  - # Using the IOP and CS3APIs
- # Generic JupyterLab extension
  - # Easy to install and configure
  - # Platform/Infrastructure independent

[github.com/sciencemesh/cs3api4lab](https://github.com/sciencemesh/cs3api4lab)



- # Same view as EFSS inside Jupyter
  - # Access files, different mounts, shares, versions, etc.
- # Sharing functionality
  - # Share with users or public links
  - # Same permissions everywhere
- # **Parallel access to notebooks**
  - # As alternative to concurrent editing
  - # Opening the same notebook without creating conflicts (both locally or remote)
  - # Execution environment independence

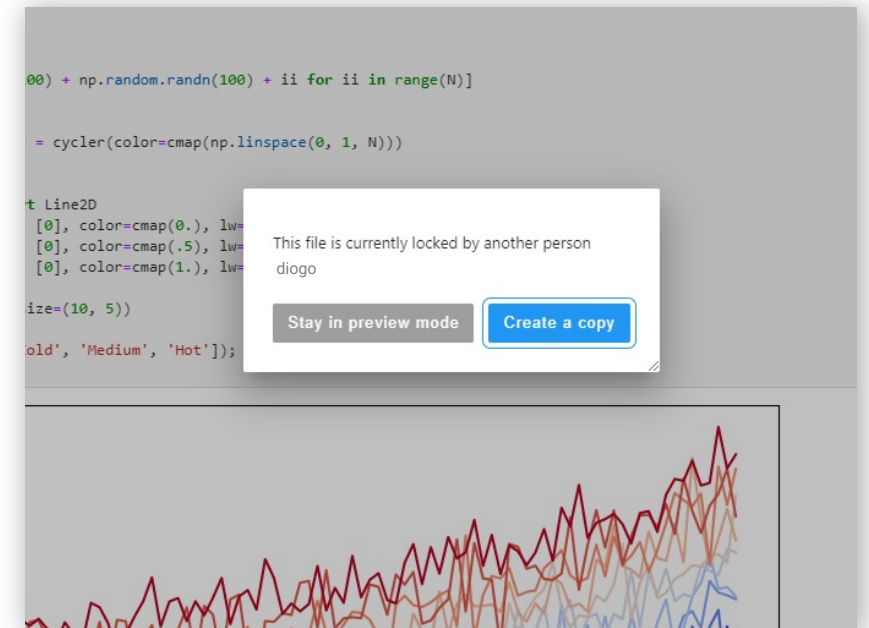


## # Uses locks to prevent other users from editing

- # Independent of application/server/etc opening it
  - # (they need to support it, but storage backend can enforce it)
- # Keeps track on the user that is holding the lock

## # A different approach

- # The first to open locks, the following users can see it in Read Only (or create a copy...)
- # Different but - actually - **complementary** to concurrent editing!



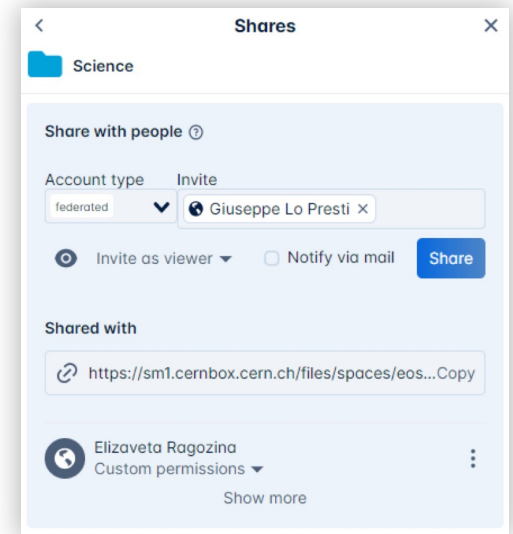
# <https://indico.cern.ch/event/1210538/contributions/5317088/>

## # Clone + merge

- # Visual diff of what changed between original and the copies (nbdime)
- # Allow merge from UI (automated if possible)

## # Enable collaboration across Mesh sites

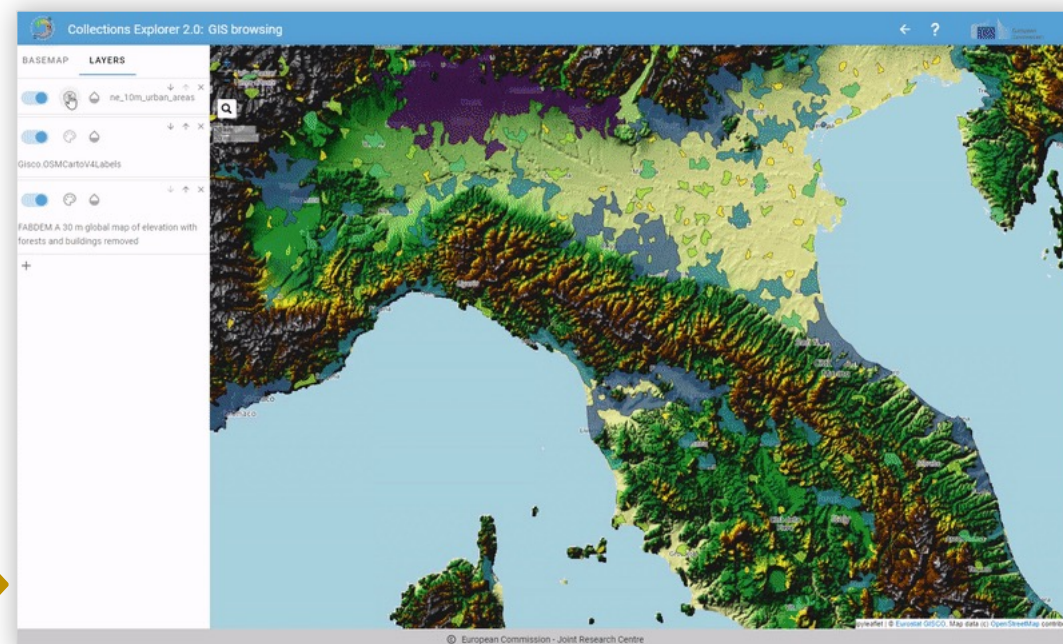
- # Finish the implementation for collaboration between different institutions (federated workflows)
- # Might also allow sharing of other types of resources
  - # But doing it with external users might be a security risk
  - # So we're just focusing on remote access to files



## # VOilà Simplification library

# Facilitate full exploitation of ipyvuetify/vuetify.js components with less code

- # Consistent usage of widgets variants/colors/ themes
- # Fullscreen apps, responsiveness, multipage, layered popups, etc

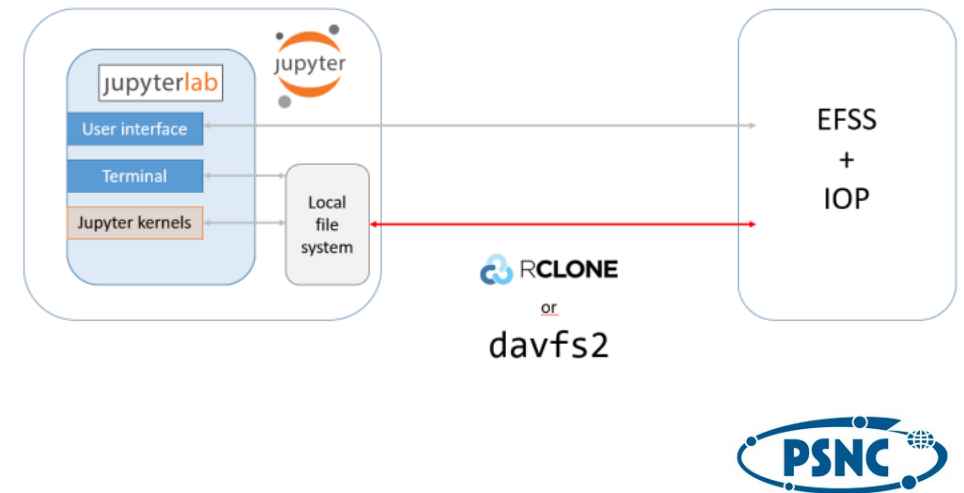


[vois.readthedocs.io](https://vois.readthedocs.io)



# <https://indico.cern.ch/event/1210538/contributions/5207924/>

- # VaaS - Voilà as a Service
- # File system integration (backend)
  - # For data access from kernels
  - # Tested for compatibility with various EFSS systems
  - # Local filesystem access
    - # Sync with native EFSS clients
    - # Sync with Rclone
  - # Online access
    - # FUSE mount with davfs2
    - # FUSE mount with native clients (i.e EOS client)





**CS<sup>3</sup>  
MESH<sup>4</sup>  
EOSC**

**Connecting European Data**

**Thank you!**  
Discover more on...

 [cs3mesh4eosc.eu](https://cs3mesh4eosc.eu)

 [company/cs3mesh4eosc](https://www.linkedin.com/company/cs3mesh4eosc)

 [@cs3mesh4eosc](https://twitter.com/cs3mesh4eosc)



CS3MESH4EOSC has received funding from the European Union's Horizon 2020 Research and Innovation programme under **Grant Agreement No. 863353**.