# Evolving SWAN through simplification

**Diogo Castro**

On behalf of the SWAN team

https://cern.ch/swan

CERN

# A reminder on SWAN

# Integrating (CERN) services



UI/Core

Software

Analysis platforms

Storage

Compute

Infrastructure

- Service for Web-based Analysis
  - Created in **2016**
  - Used by **~300** people daily
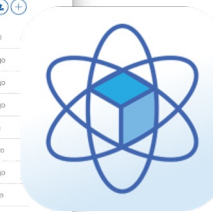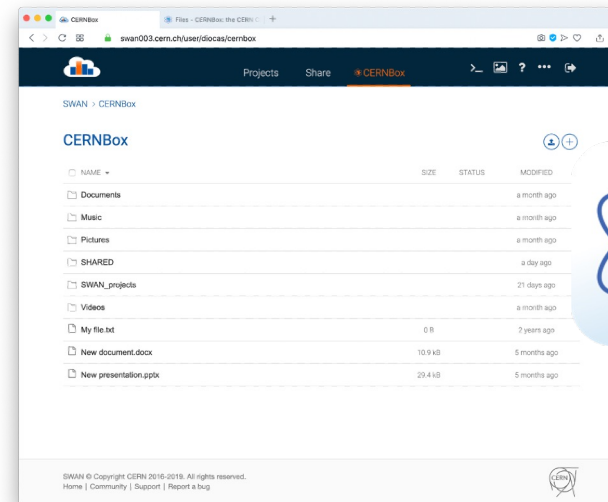
# Storage

> All the data our users need for their analysis
  - CERNBox as home directory
  - Experiment repositories, projects, open data, …
  - (EOS Fuse client)

> Sync&Share
  - Files synced across devices and the Cloud
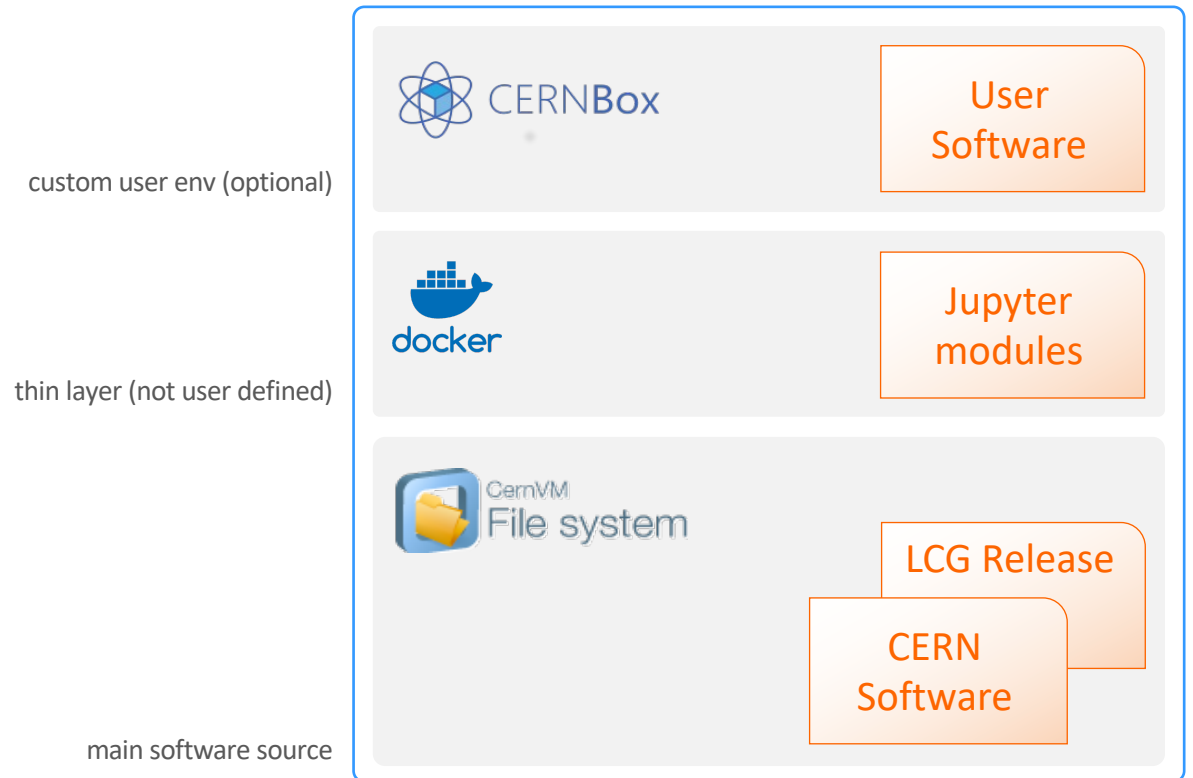  - Simple collaborative analysis

> Software distributed through CVMFS
  - "LCG Releases" - pack a series of compatible packages
  - Reduced Docker Images size
  - Lazy fetching of software

> Possibility to install libraries in user cloud storage
  - Good way to use custom/not mainstream packages
  - Configurable environment

CERNBox — User Software

custom user env (optional)

docker — Jupyter modules

thin layer (not user defined)

CernVM File system — LCG Release / CERN Software

main software source

# Latest updates

# Project priorities in 2023

**1**
- Conclude migration to Kubernetes
- Ensure scalability

**2**
- Conclude migration to Jupyterlab

**3**
- Migration to Alma 9 / simplification of current docker images
- Update to latest versions of upstream

**4**
- Conclude integration of more CERN services

**5**
- New ways to manage software
- Binder

> Migration campaign finished on March 5th, 2024
  - All physical nodes have been removed from the service
  - Improved operations and a single source of truth for metrics
  - Better UX for users (aligned set of features, single point of entry)

> Lessons learned on the way
  - Some components' maturity had to stabilize
  - Need to gain operational expertise
  - Physical machines were used as a fallback to disruptive updates

> Further improvements (in progress)
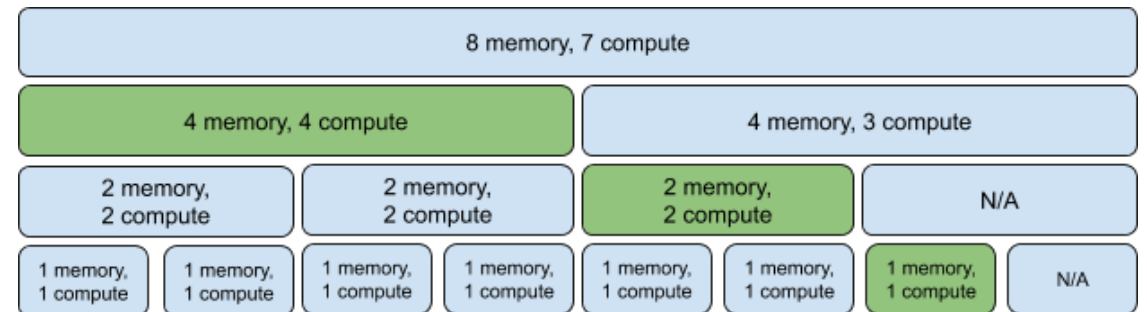  - Blue/Green deployment for cluster (disruptive) updates
  - DevOps automated tools

# GPU access

> SWAN allows to attaching a GPU to a user session
>   - Feature of the new SWAN k8s deployment, now available to anyone
>   - 18 GPUS (Tesla T4)

> New project: CERN IT-wide resources sharing
>   - Sharing all GPUs across different services in a pool
>   - Easy scalability in case of need (e.g., for tutorials)

# GPU partitioning

> A new NVIDIA GPU operator deployed
  - Supports newer GPU cards, including partitionable GPUs, i.e. Multi-Instance GPUs ([MIG](#) is viable with A100s, that are scarce at CERN)

> With GPU partitioning, users have exclusive access to one GPU fragment
  - No interferences

> Partitioning GPUs allows for resource sharing and better resource utilization
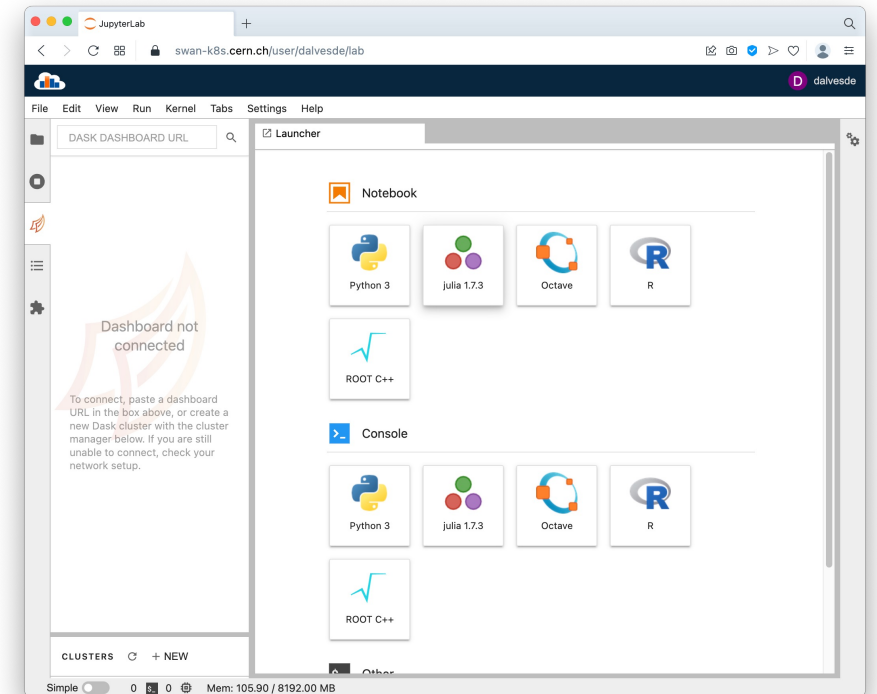  - E.g. assign one fragment per participant of a tutorial

# Migration to JupyterLab

> Deployed Jupyterlab v4
  - Extensions migrated to the new version

> Available as beta UI
  - Collection of user feedback underway
  - Users can use the old UI in parallel

> Deeper Sync&Share integration
  - Ongoing integration with CERNBox using the CS3 APIs Jupyterlab extension (CS3Mesh project)
  - Full sharing and collaborative capabilities
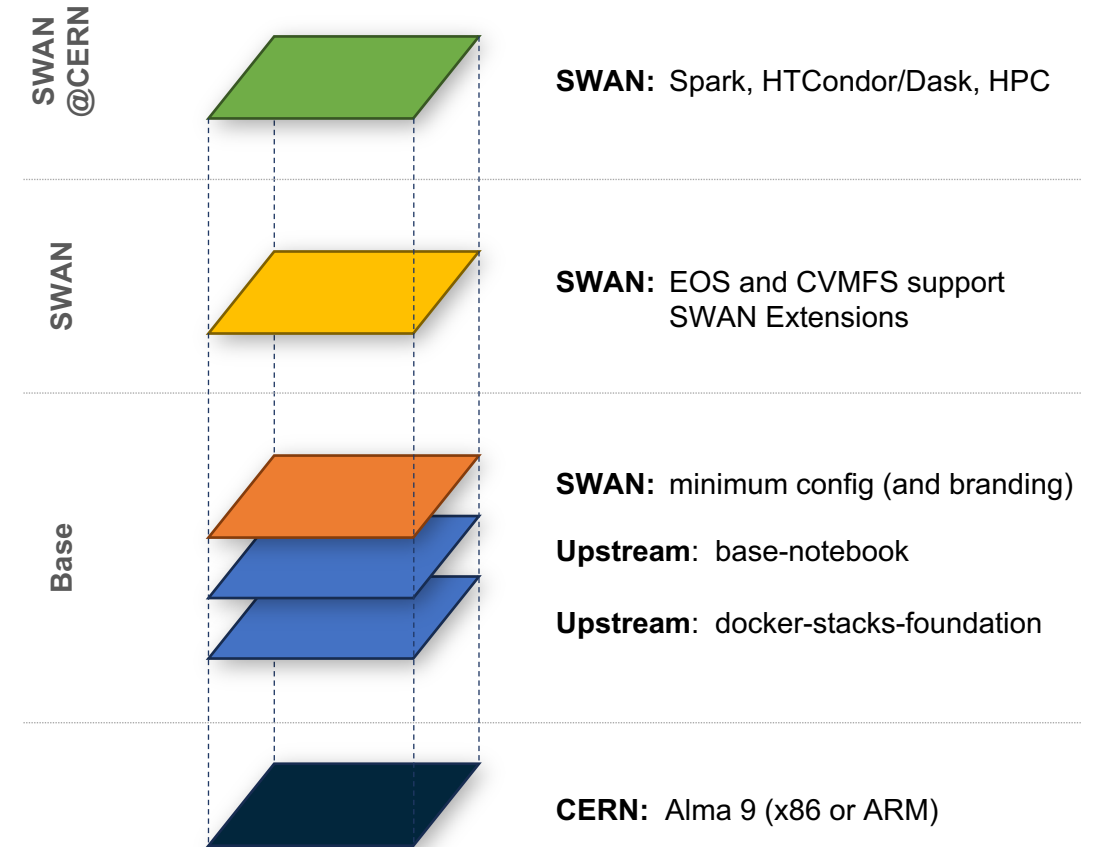  - Currently migrating to Lab v4 and making UI production ready

# Migration to Alma 9

> Key SWAN container images migrated from CERNCentOS7 to Alma9
- User session image (Jupyter server)
- JupyterHub image

> User images rewritten from scratch
- Like upstream images, but on top of Alma 9 instead of Ubuntu
- Same entry points and configuration options

> Modular components' configuration
- EOS, CVMFS, External resources, etc, are independently configured on separate scripts
- Easy to disable or add new components

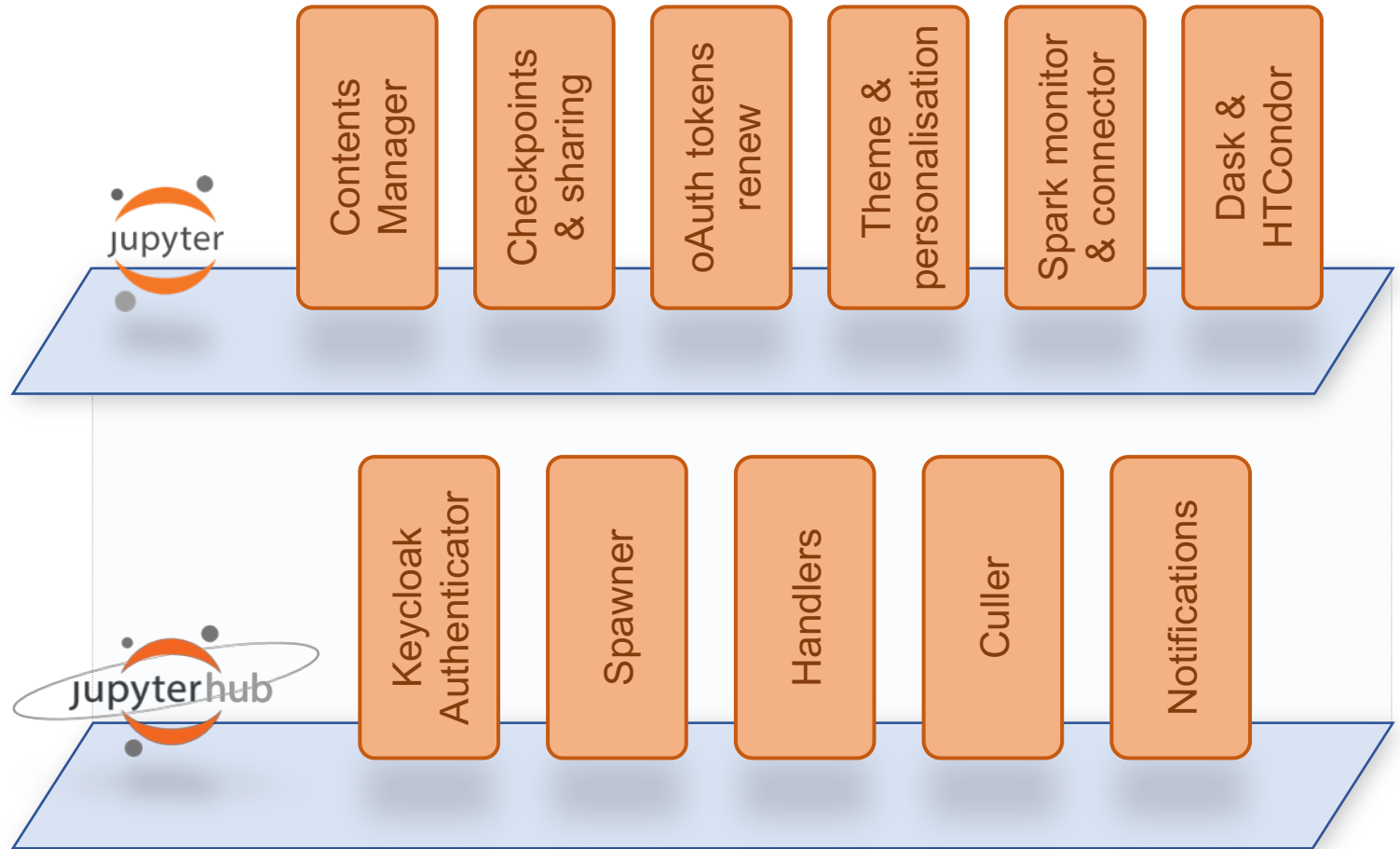SWAN @CERN

**SWAN:** Spark, HTCondor/Dask, HPC

SWAN

**SWAN:** EOS and CVMFS support
SWAN Extensions

Base

**SWAN:** minimum config (and branding)

**Upstream:** base-notebook

**Upstream:** docker-stacks-foundation

**CERN:** Alma 9 (x86 or ARM)

> More runtime freedom
  - They can be run within Jupyterhub (i.e. prod SWAN) but also independently (e.g. locally or headless in a CI)

> All dependencies updated to the latest versions
  - Some SWAN personalizations were replaced with upstream ones

jupyter
- Contents Manager
- Checkpoints & sharing
- oAuth tokens renew
- Theme & personalisation
- Spark monitor & connector
- Dask & HTCondor

jupyterhub
- Keycloak Authenticator
- Spawner
- Handlers
- Culler
- Notifications

> Spark infrastructure is being updated to Alma 9
  - Coordinating with Hadoop service the updates to Alma 9 for Spark on Hadoop
  - Investigating deploying Spark on YARN using container images to streamline the update

> All Spark Jupyter extensions have been updated to Lab 4

> CERN HPC integration is now in QA
  - Applications and use cases that do not fit the standard batch HTC model, typically parallel MPI applications
  - Uses CEPH as shared storage between submission and worker nodes
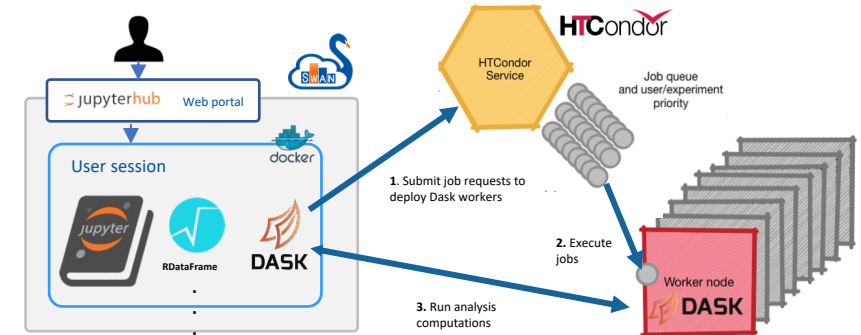  - CEPH FS integrated as PVCs in SWAN, mounted only for allowed users

# Analysis facility pilot

> Support interactive distributed analysis for High Energy Physics
  - Address the future analysis needs due to foreseen increase in data volumes.

> Dask as the connector to batch resources
  - The two main HEP analysis frameworks, ROOT and coffea, rely on Dask for running analysis distributedly

> For now, it uses overcommitted "static" slots on HTCondor
  - Optimizes usage of batch resources
  - A well-stacked batch farm with a good job mix can get to 80% CPU utilization
  - Known analysis jobs potential to stack nicely with other workloads to drive up utilization
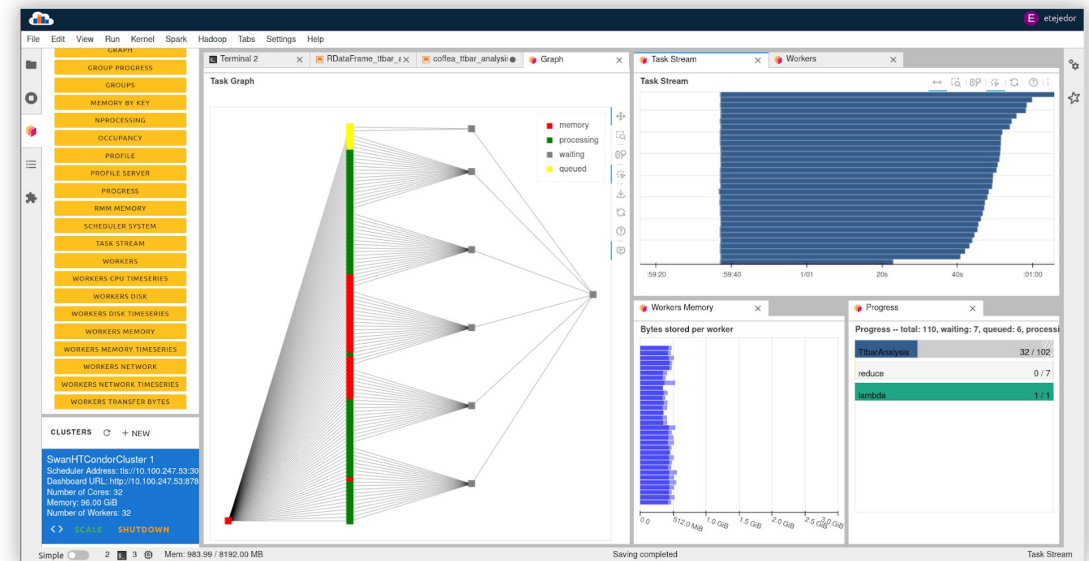
# Analysis facility pilot

> A Pilot has been approved
> - Validate demand
> - It will validate with real users its usefulness to CERN use cases and the necessity of improvements

> Future improvements
> - Use tokens throughout the workflow (currently a Kerberos auth is required)
> - Allow users to close the notebook UI and still be able to retrieve the jobs' status
> - Potentially different interactive jobs allocation model
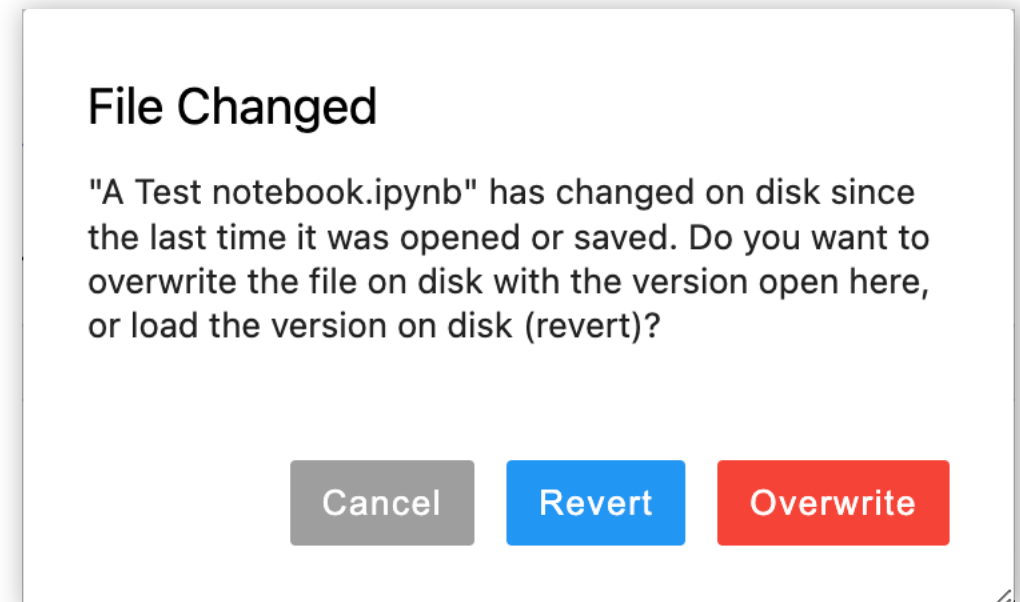> - Improve custom software environment integration

# Virtual environments

> Project kickstarted to allow the integration of LHC Control tools in SWAN

> Objective: Make it easier to manage project dependencies and ensure consistent execution across different environments
  - So that users can publish projects via Gitlab, and still recreate the environment consistently

> For performance reasons, for now, environments are created locally in the container storage
  - We are investigating the feasibility of the EOS Squash FS feature to persist across sessions

> We will try to integrate GUI tools to help add/remove packages
  - Similar to a PoC/GSoC project presented in previous years at CS3 Conference

> The integration of BinderHub is on hold for now

# A note on collaboration

# Current collaboration model for Jupyterlab

> 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 (?)

> Jupyterhub proposes "collaboration accounts" instead
  - "Real-time collaboration without impersonation"

**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)?
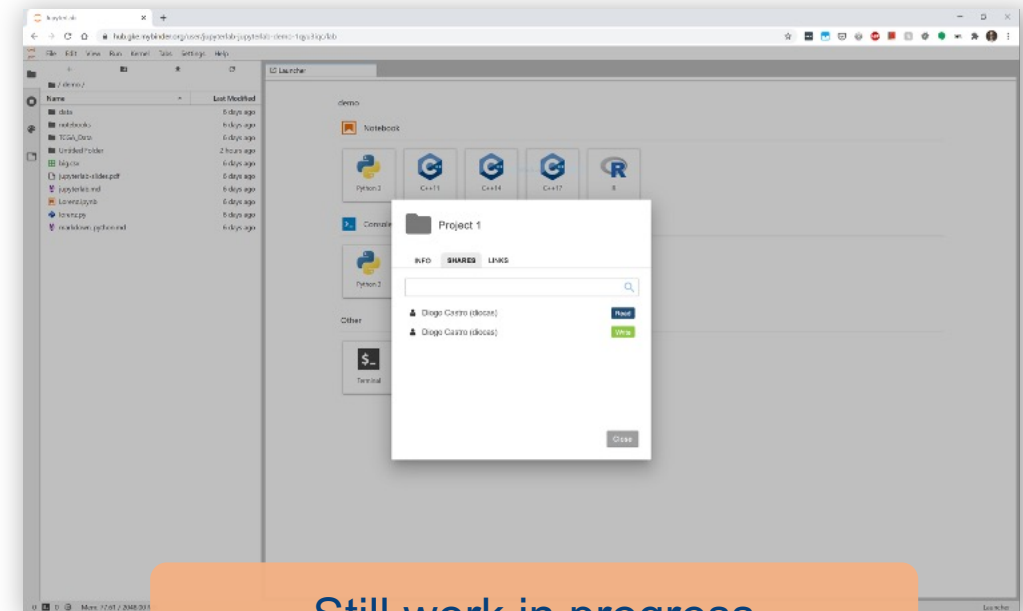
Cancel    Revert    Overwrite

> 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
  - This 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
  - Collaboration accounts help, but might be harder to coordinate or integrate with deployment

> We're not aware of use cases that would benefit from true concurrent editing

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

# Collaboration model of the CS3Mesh project

> ## 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



**Still work in progress**

https://github.com/sciencemesh/cs3api4lab

# Conclusion

# Takeaways

> With more configurable options in the upstream Jupyter project, SWAN is being simplified
  - It results in a project that is better to manage and operate at CERN
  - But also easier to deploy outside of CERN
  - The new docker images and full Kubernetes deployment are examples of that

> SWAN continues to work on its integration with heterogeneous and external resources
  - From GPUs to Htcondor via Dask or HPC
  - A Pilot Analysis Facility is ongoing to validate the demand and applicability to CERN use cases

> The collaboration model of Jupyter would benefit from our input
  - As deployers and developers of large sync&share services/products, we have relevant know-how
  - But we need to organise bahind a single voice

# Where to find us

> Contacts
> - swan-admins@cern.ch
> - http://cern.ch/swan
> - https://swan-community.web.cern.ch/

> Repository
> - https://github.com/swan-cern/

> Science Box
> - (deploys the SWAN Helm Chart)
> - https://cern.ch/sciencebox

# Evolving SWAN through simplification

Thank you

Diogo Castro
diogo.castro@cern.ch