# Evolving SWAN towards an Analysis Facility system

**Diogo Castro**

On behalf of the SWAN team

https://cern.ch/swan

# A reminder on SWAN

And its current status

# Integrating (CERN) services



UI/Core

Software

Analysis platforms

Storage

Compute

Infrastructure

- Service for Web-based Analysis
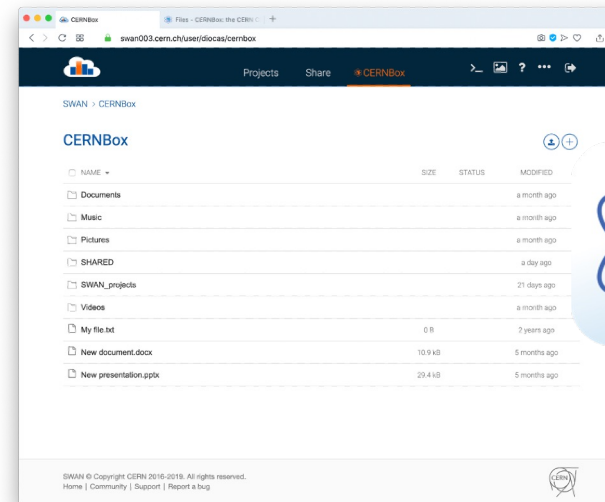  - Created in **2016**
  - Used by **200-250** people daily

# 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
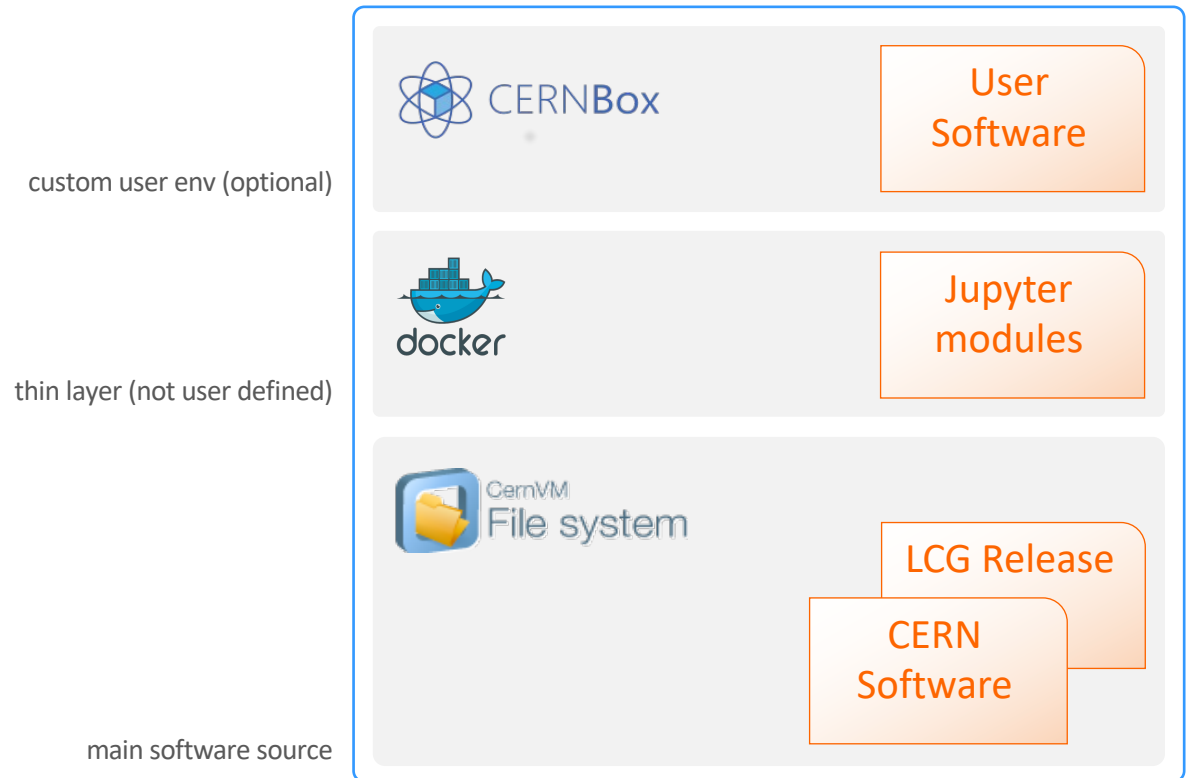- (Custom share API)

# Software

> 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

# Current project priorities

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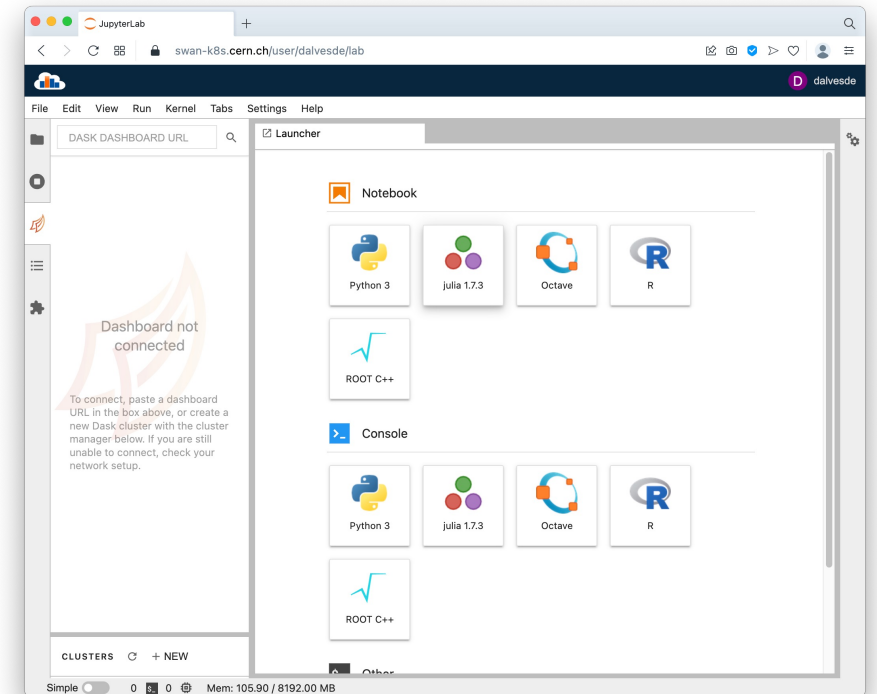
**5**
- New ways to manage software
- Binder

# Migration to Lab

> ## Final stages of migration

- Spark (monitor, connector) extensions migrated
- Other small extensions+branding, also done
- Looking into other integrations (i.e Dask)

> ## Deeper Sync&Share integration

- Ongoing integration with CERNBox using the CS3 APIs Jupyterlab extension [1]
- Full sharing and collaborative capabilities
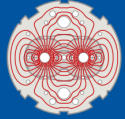- Easier integration of SWAN with other CS3 enabled EFSS



[1] See "Science Mesh demos" (Cs3api4lab), Tue @ 14:50
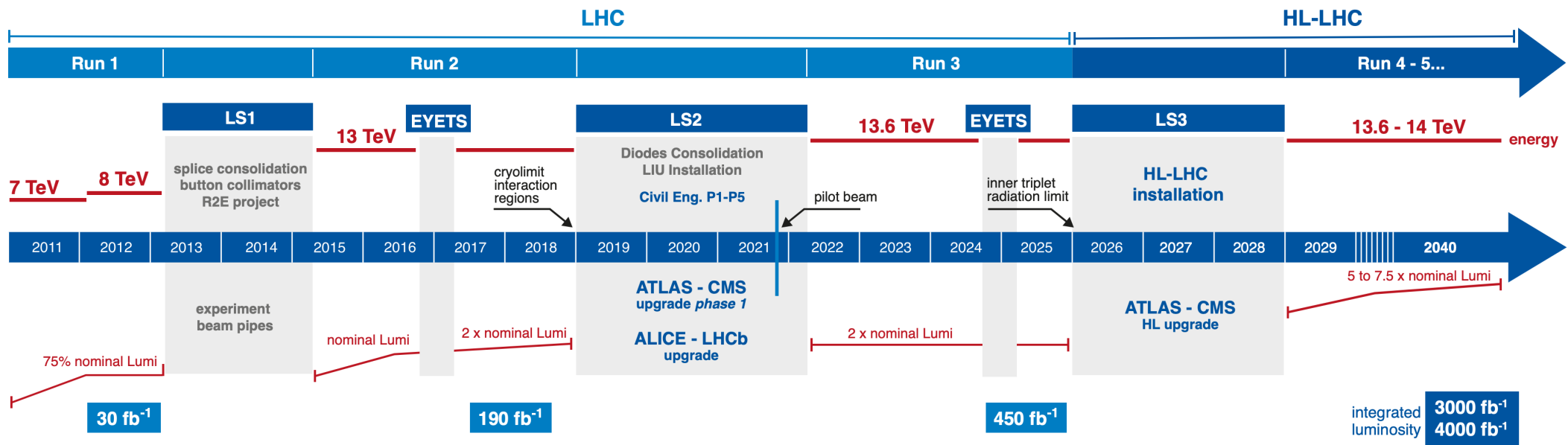
# Why an analysis facility?

And why SWAN needs to evolve
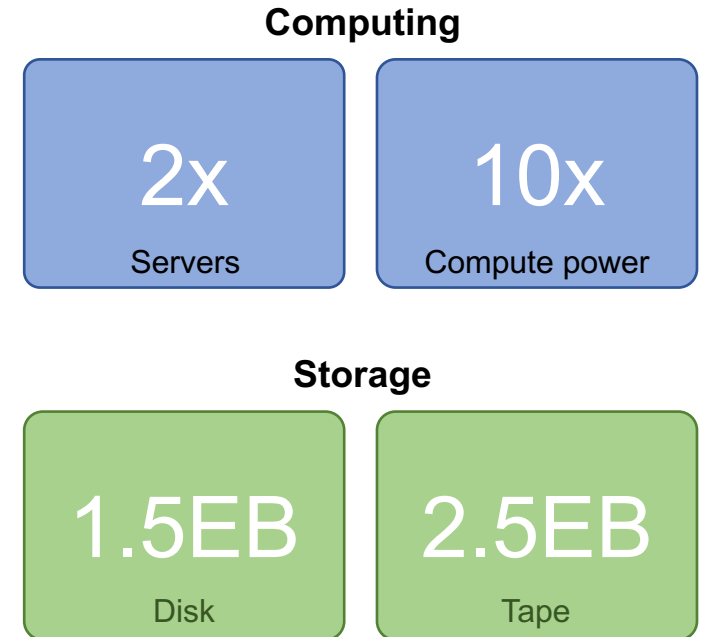
# Why an analysis facility?

> HL-LHC needs are pushing us to build modern Analysis Facilities
- Traditional batch processing
- Interactive computing on big datasets, with new interfaces (Jupyter)

> An AF should facilitate access to:
- Software
- Storage (+ sharing)
- Computing resources (elastic)

> Ongoing effort to provide an AF @ CERN
- Interdepartmental collaboration

**Computing**

| | |
|---|---|
| **2x** Servers | **10x** Compute power |

**Storage**

| | |
|---|---|
| **1.5EB** Disk | **2.5EB** Tape |

# Build on what already exists

> Evolution as opposed to a new and dedicated AF
>
> - Less costs and quicker to deploy
> - SWAN is a good candidate for the interface of an AF @ CERN!

> SWAN needs to be an entry point to external and heterogeneous resources
>
> - Multiple services already available at CERN
> - We don't want to run the full infrastructure
> - More freedom to users to chose what best fits their use cases

# Resources integration

Past, ongoing and future work

# Integration model

> ## A Platform for physics analysis
> - With support for both *single-node* and *distributed* analysis
> - Keep only lightweight resources "local"

> ## SWAN acts as a client to other resources
> - Allows connection to multiple types of resources, instead of creating a session on them directly
> - Doesn't sacrifice the usablity (users have access to the system without delay)
> - Test small and local (quick) and run big and distributed
> - Keeps the independence between services (i.e upgrade schedules, dependencies, etc)

# Spark

> SWAN is connected to the Spark clusters at CERN

    Physical: ~3800 cores, some dedicated
    Virtual: ~250 cores, on demand (kubernetes)

> Jupyter extensions available to:

    Connect to a certain cluster
    Monitor the execution

# GPUs

> SWAN allows to attach a GPU to a user session
  - Feature of the new SWAN k8s deployment
  - ~12 GPUS (Tesla T4)

> The GPUs are used interactively
  - When starting their session, the user selects a CUDA software stack and gets a GPU
  - GPU-enabled packages (e.g. tensorflow, PyTorch) can then be used in a notebook and offload to the GPU by default

> Currently looking into GPU concurrency

```
In [1]:   import tensorflow as tf

          tf.debugging.set_log_device_placement(True)

          # Create some tensors
          a = tf.constant([[1.0, 2.0, 3.0], [4.0, 5.0, 6.0]])
          b = tf.constant([[1.0, 2.0], [3.0, 4.0], [5.0, 6.0]])
          c = tf.matmul(a, b)

          Executing op MatMul in device /job:localhost/replica:0/task:0/device:GPU:0
```

# HTCondor

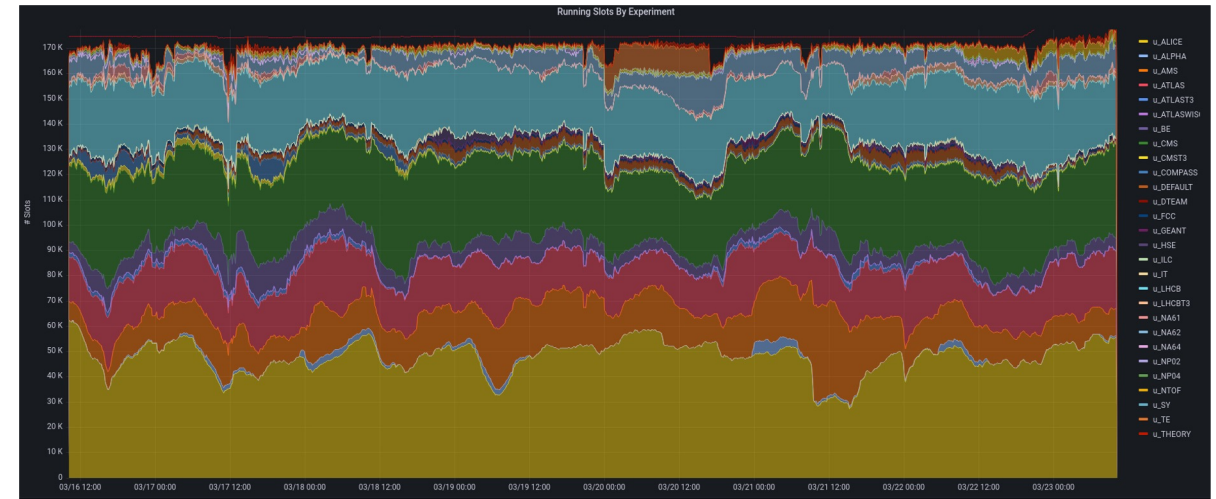> ## Goal: leverage HTCondor pools at CERN from SWAN

- Up to ~175k cores in shared pools at CERN – limited by the quotas assigned depending on experiment affiliation
- Already used for analysis

> ## Batch submission: already supported

- Condor packages available on CVMFS

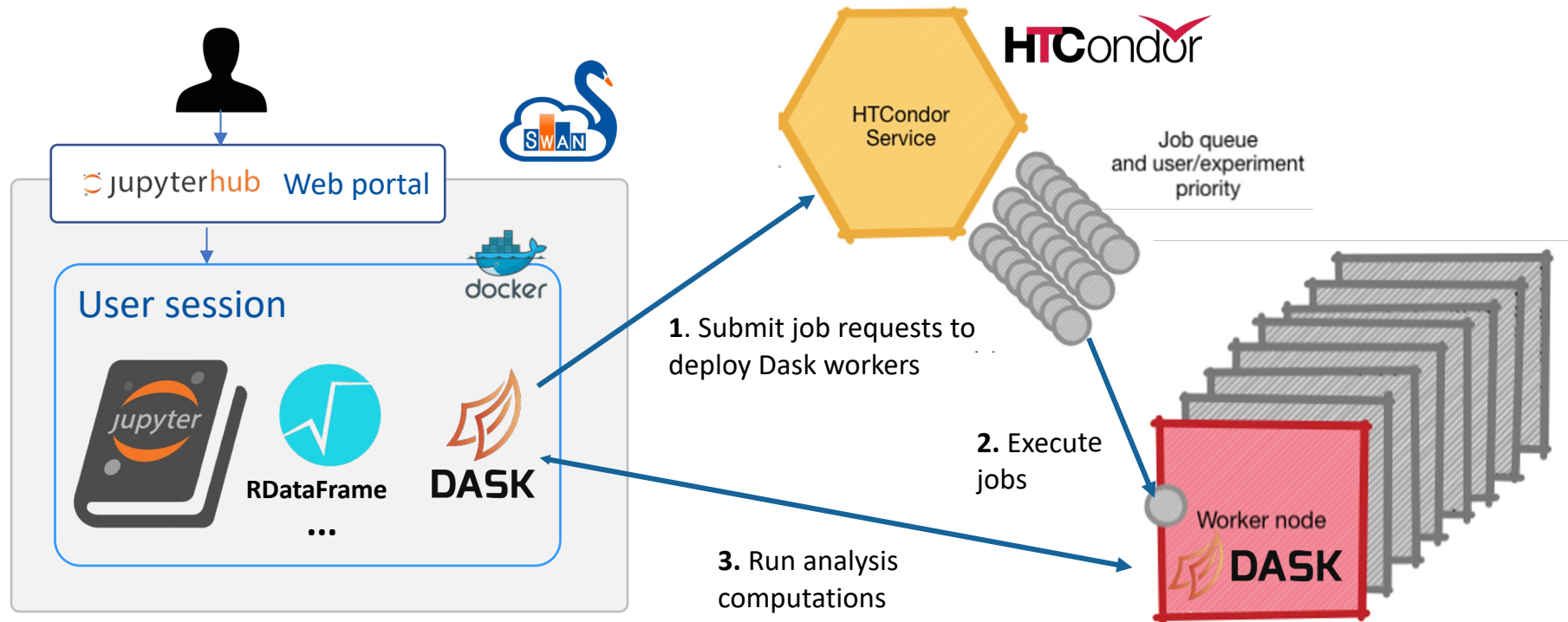> ## Interactive usage: in pilot phase

- Collaboration with Batch Service@CERN
- Dask packages available on CVMFS
- Will be exposed to users when migration to JupyterLab is finished (Q2 2023)

# HTCondor + Dask = interactivity

1. Submit job requests to deploy Dask workers

2. Execute jobs

3. Run analysis computations

# HTCondor + Dask = interactivity

> ## Upstream Dask extension

- A cluster shared across multiple "clients" (notebooks)
- Better resource utilization
- Spark connection is evolving into this model as well

# HPC

> ## HPC service at CERN
  - Applications and use cases that do not fit the standard batch HTC model, typically parallel MPI applications.
    - Ex: Computation Fluid Dynamics, Beam simulation, plasma simulation, …)
  - Uses Ceph FS between submission and work nodes

> ## SWAN integration ongoing
  - Testing authentication and software requirements
  - Missing storage integration: CERNBox/REVA/CEPH integration?
    - We don't want shared secrets from other services

> ## Reva Ceph FS "simplified" storage backend
  - PoC available, some features (sharing/ACLs, snapshotting) not available but not needed

# Other future possibilities

> Reana
- CERN's Reusable Analysis Platform
- Move from exploratory analysis into a reproducible format

> Kubeflow
- Machine learning and MLOps service at CERN
- Create and deploy pipelines directly from SWAN

# Conclusion

# Takeaways

> The HL-LHC upgrade will bring many challenges
  - In terms of storage and computing
  - But also on the complexity of the analysis

> The CERN infrastucture is evolving to cope with the load
  - And a new Analysis Facility workgroup was formed

> SWAN is a good candidate for the interface of an AF @ CERN
  - Gives access to the software, storage and UX expected by the users
  - More services need to be integrated

# 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 [1]
- (deploys the SWAN Helm Chart)
- https://cern.ch/sciencebox

[1] See "Driving the ScienceBox package into the future", Wed @ 11:45

# Evolving SWAN towards an Analysis Facility system

Thank you

Diogo Castro

diogo.castro@cern.ch