# Facilitating Collaborative Analysis in SWAN

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

https://swan.cern.ch

July 4th, 2018
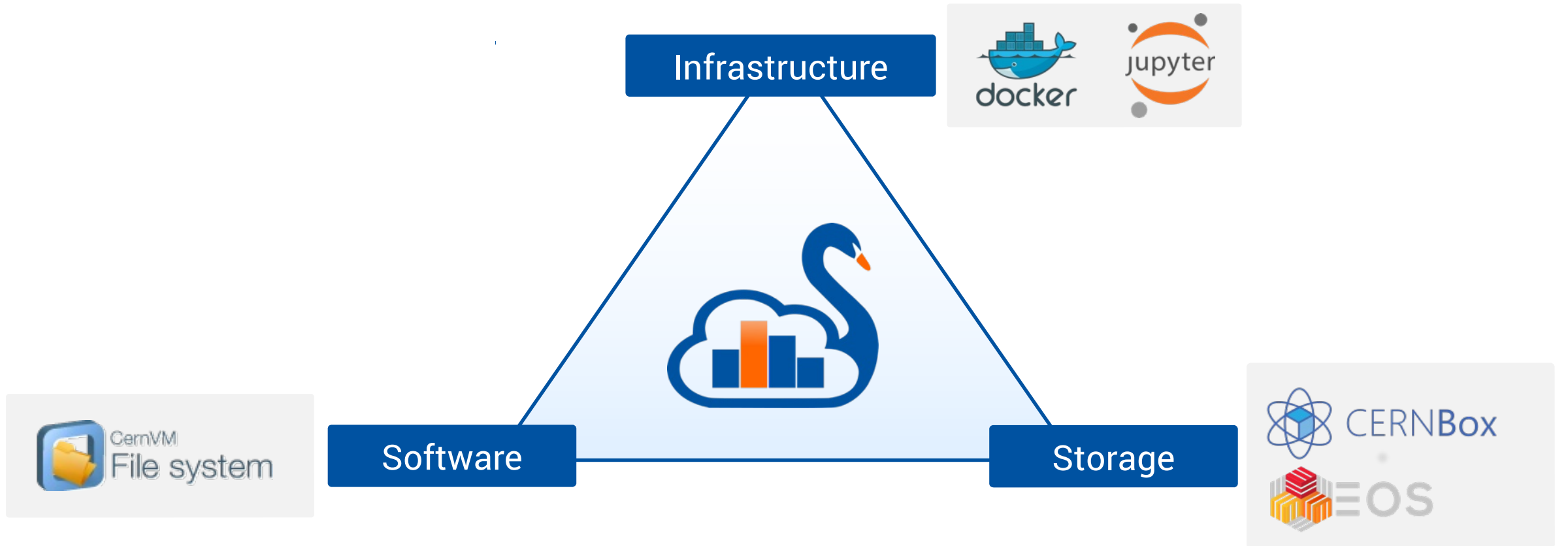**CHEP 2018, Sofia (Bulgaria)**

# Introduction

# SWAN in a Nutshell

> Analysis only with a web browser

  - No local installation needed
  - Based on Jupyter Notebooks – interactive computing
  - Calculations, input data and results "in the Cloud"

> Support for multiple analysis ecosystems

  - ROOT, Python, R, …

> Easy sharing of scientific results: plots, data, code

> Integration with CERN resources
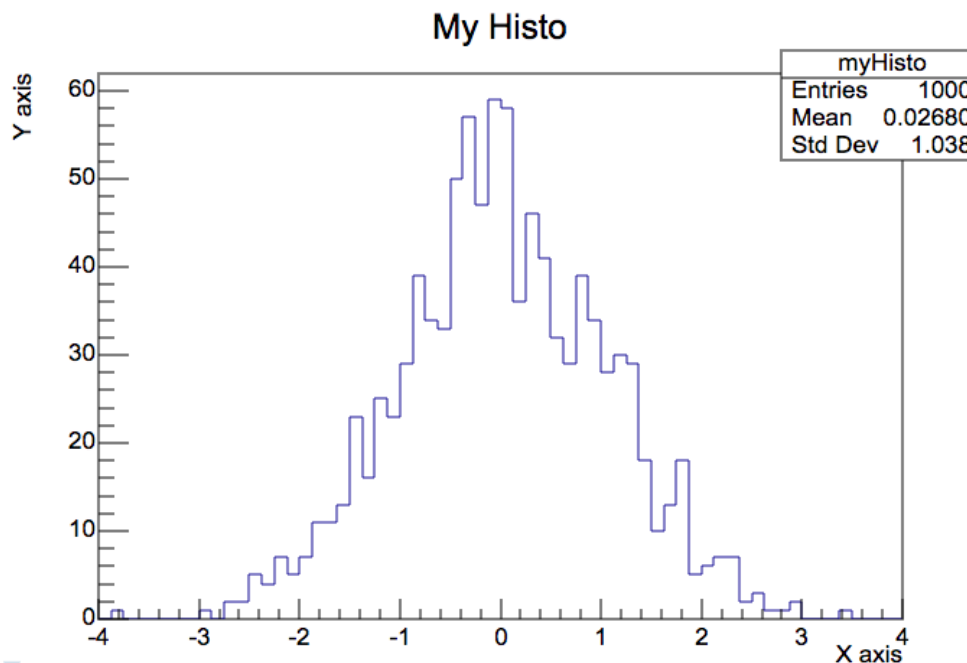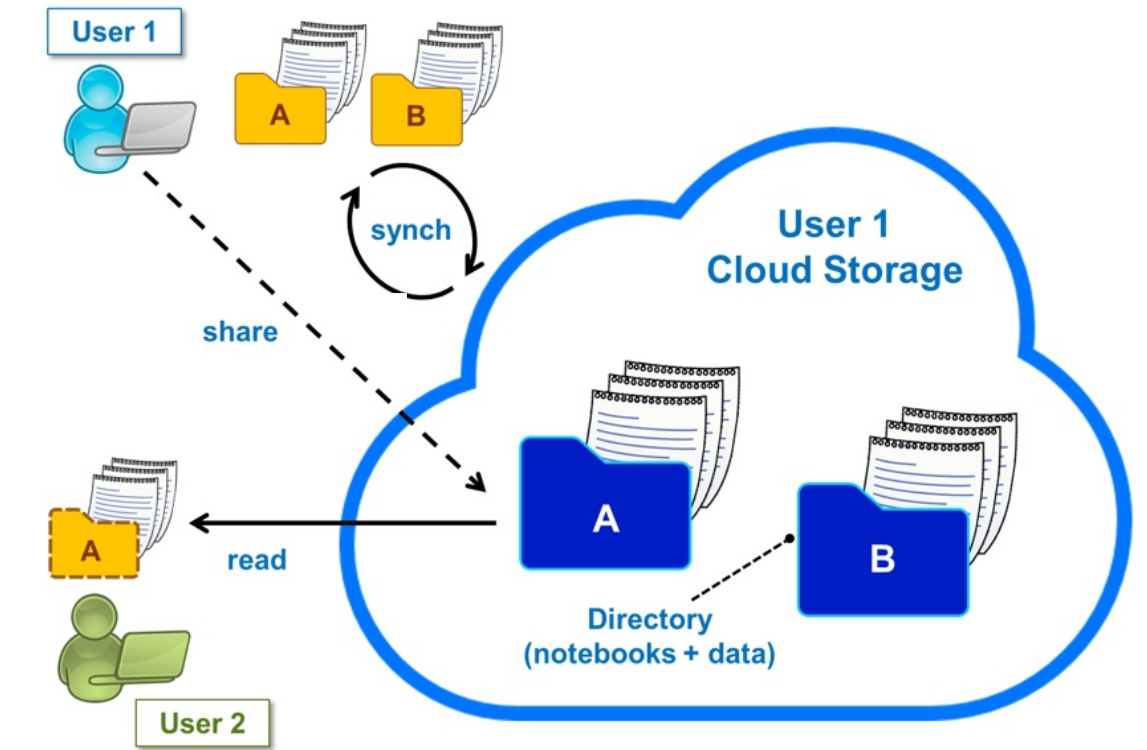
  - Sofware, storage, mass processing power

# Collaborative Analysis

# Cloud Storage as your Home

> ## CERNBox is SWAN's home directory
> - Storage for your notebooks and data

> ## Automatic synchronization
> - Files synced across devices and the Cloud

> ## Provides foundations for sharing
> - Collaborative analysis

# New User Interface

# Sharing Made Easy

> ## Sharing from within the SWAN interface
  - Integration with CERNBox

> ## Users can share "Projects"
  - Special kind of folder that contains notebooks and other files, like input data
  - Self contained

# The Share Tab

> Users can list which projects...
  - they have shared
  - others have shared with them

> Project information
  - Sharer
  - Size
  - Date

# Inspecting a Project

> By clicking on a shared project, a user can inspect its contents
  - Browsing of the project files
  - Static rendering of notebooks

> Useful to decide whether to accept or not the shared project

# Accepting a Shared Project

> When accepting a shared project, its contents are cloned to the receiver's CERNBox

- The receiver will work on their own copy

> Concurrent editing not supported by Jupyter

- Safer to clone

# Sharing Spark Projects

# SWAN for Education

> UP2University European Project
>   - Bridge the gap between secondary schools, higher education and the research domain

> SWAN used by students to learn physics and other sciences
>   - Let them use the very same tools & services used by scientists at CERN

> SWAN Boxed: distribution easily deployable on premises
>   - https://github.com/cernbox/uboxed

Poster: The EU Up to University Project, E.Bocchi, J. Moscicki

# Conclusion

# Summary

> The interface of SWAN has been completely redesigned to foster collaboration and sharing of results among scientists

> Sharing is now fully integrated in SWAN, where users can share their work in the form of Projects (notebooks + data)

> New functionality to share a project, list shared projects, inspect and clone shared projects into users CERNBox

> Concurrent edition of notebooks in Jupyter would give another dimension to sharing

# Facilitating Collaborative Analysis in SWAN

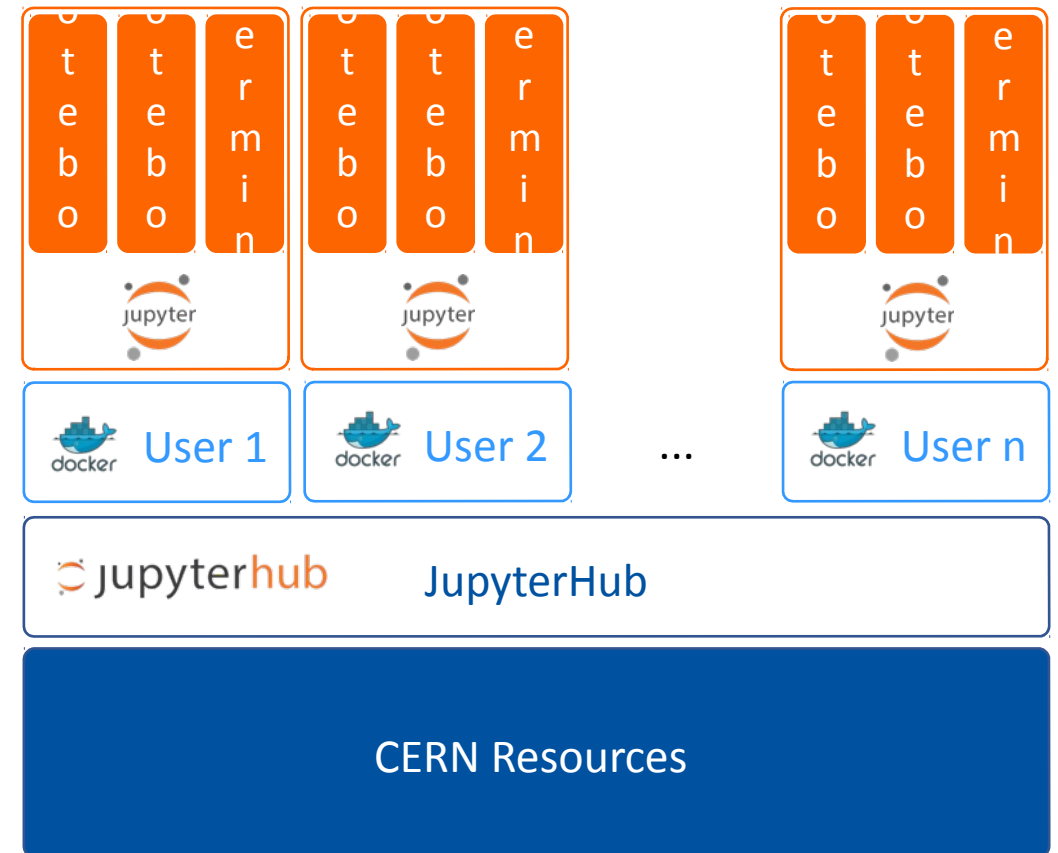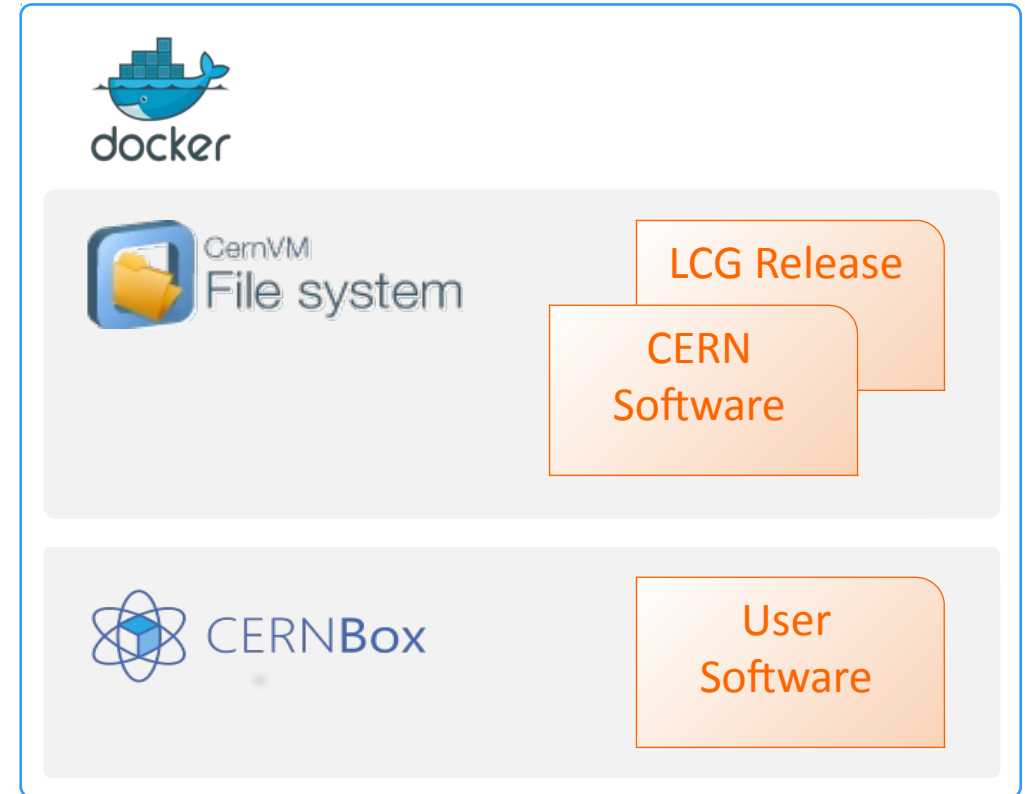Thank you

# Backup Slides

# Integrating Jupyter

> Configurable environments through user defined scripts

> Jupyterhub to allow multiple Jupyter instances
  - Single instance of Jupyter per user

> User sessions spawned as Docker containers
  - Enforces resource limits per user
  - To isolate users work
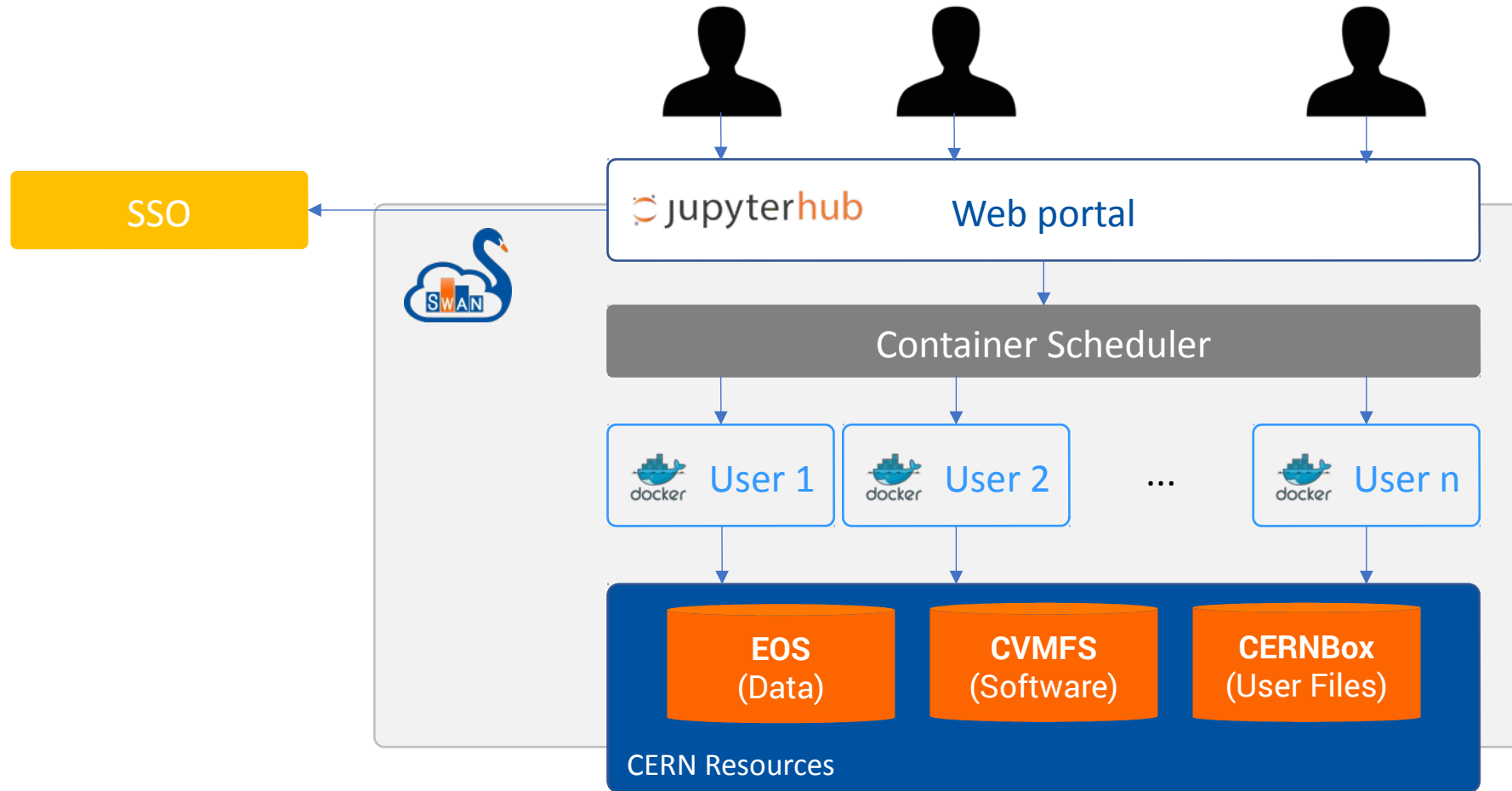
# Software

> ## Software distributed through CVMFS
> - LCG Releases: distribute a series of compatible packages
> - Software used by researchers is available

> ## Possibility to install other libraries in user storage (CERNBox)

# The Notebook as Interface

> A web-based interactive interface and platform that combines code, equations, text and visualisations
  - Ideal for sharing/collaboration
  - A "shell opened within the browser"

> Many supported languages (kernels)
  - In SWAN: Python, ROOT C++, R and Octave

> Interactive, usually lightweight computations

> Very useful for some use cases
  - Final steps of an analysis, exploration, teaching, documentation and reproducibility
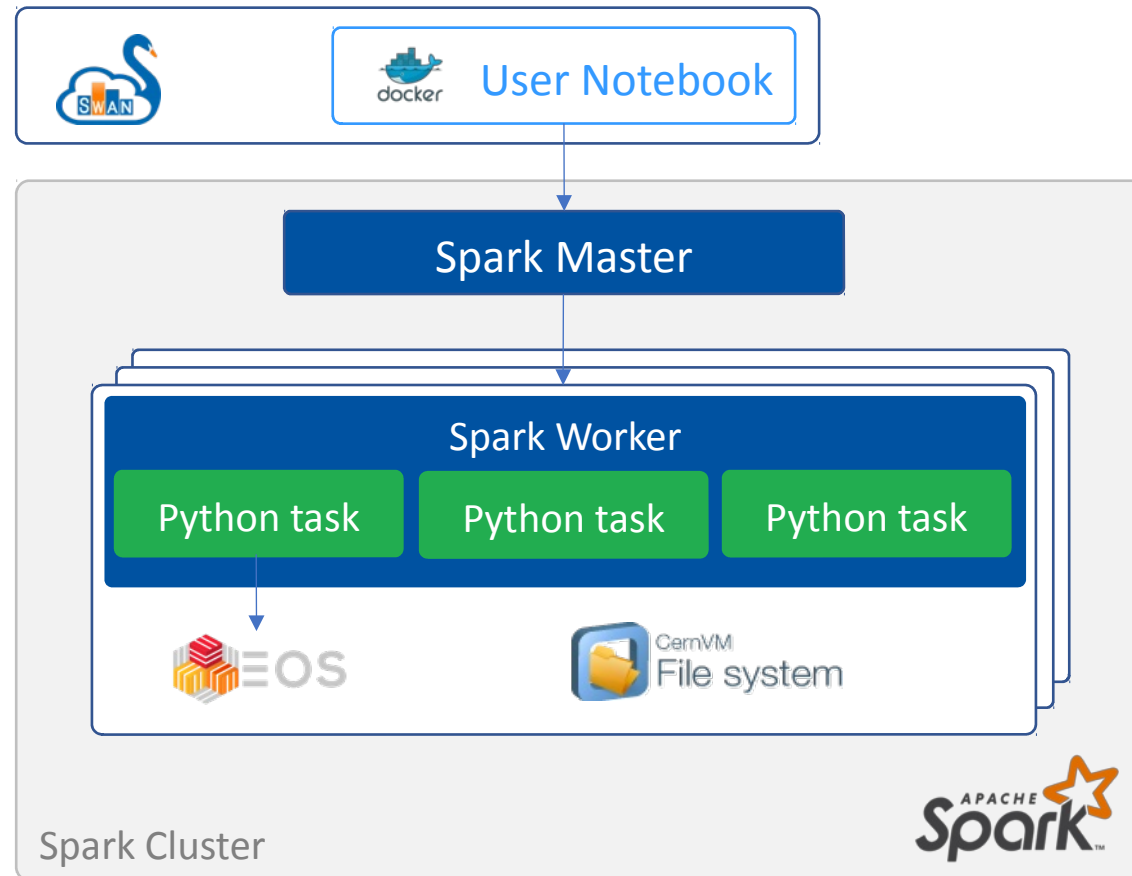
# New User Interface

# Boxed

> Containerized version of all the infrastructure

- Includes EOS, CERNBox, CVMFS and all Swan services (Jupyter Docker image, JupyterHub)
- Available in https://github.com/cernbox/uboxed

> Easily deployable on premises

- Installable in Linux systems
- Based on Docker Compose

# Integration with Spark

> Connection to CERN Spark Clusters

> User data accessed through EOS

> Graphical Jupyter extensions developed
  - Spark Connector
  - Spark Monitor
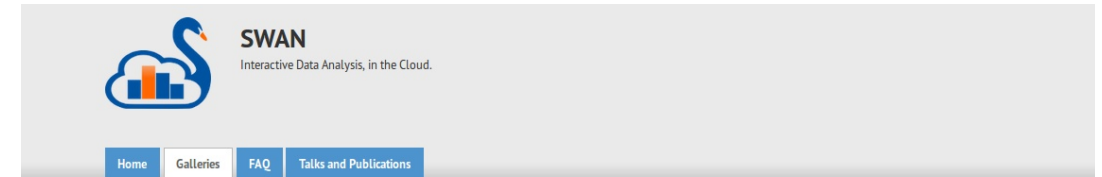
# HEP User Community

> SWAN development is guided by our user community

- New features (libs, kernels, ...) are requested by users from their real usage needs

> Gallery of examples

- Made in collaboration with our users
- Almost 50 notebooks in 7 categories



**Access with only a click**