# SWAN: service for web based analysis

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

https://swan.web.cern.ch

# Introduction

# Introduction

> There are analysis tools developed in CERN, but they require installation and configuration

> Some resources are only available from within CERN network
  - And remote connection might not be ideal

> External services, like IBM Bluemix or Mybinder, provide analysis services but lack some advanced features
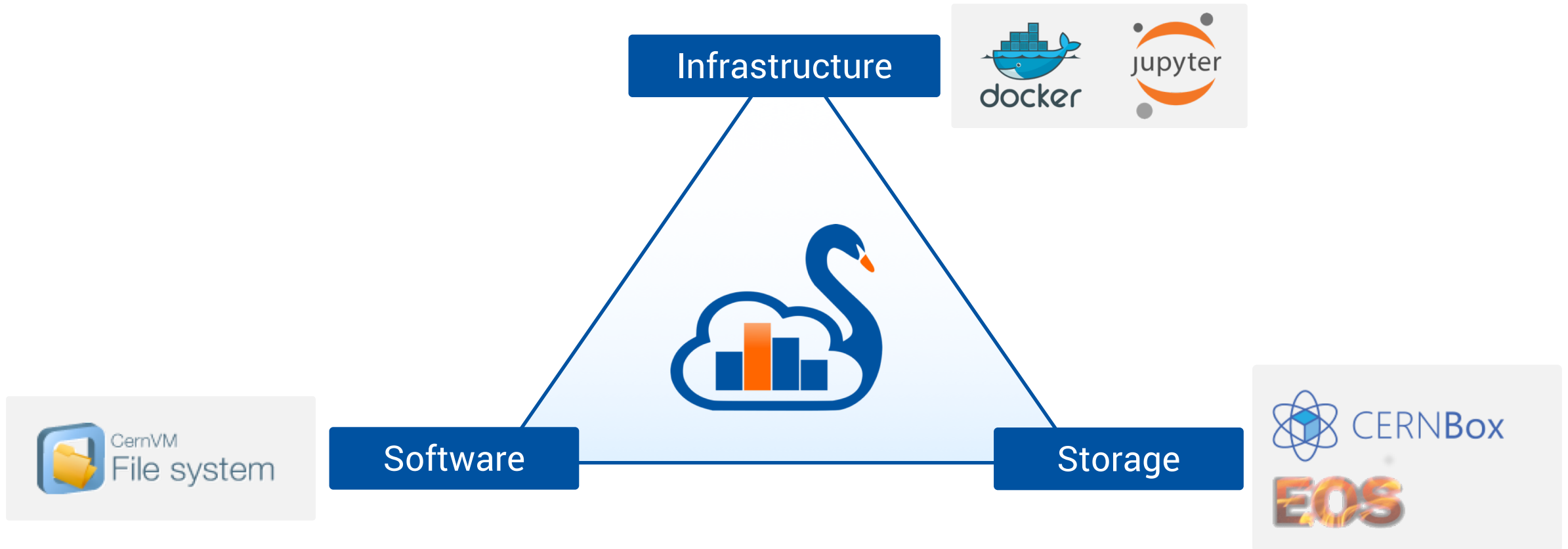  - Like software packages and integration with advanced and non volatile storage

# Motivation

> Analysis only with a web browser
  - Available everywhere and at anytime

> Easy to use (but powerful)
  - No local installation and configuration  needed

> Create easily sharable scientific results: plots, data, code
  - Storage is crucial: mass & synchronized

> Integration with CERN resources
  - Access software, user/experiments data, mass processing power

> Integration with other analysis ecosystems : ROOT, R, Python, ...

# SWAN

Infrastructure

Software

Storage

# Jupyter - The Notebook as Interface

> A web-based interactive interface and platform that combines code, equations, text and visualisations
  - Ideal for sharing/collaboration

> Many supported languages (kernels): Python, C++, Haskell, Julia, R …

> Very well received Project with major contributions and implementations from big names (IBM, Google,...)

> … In a nutshell: an "interactive shell opened within the browser"

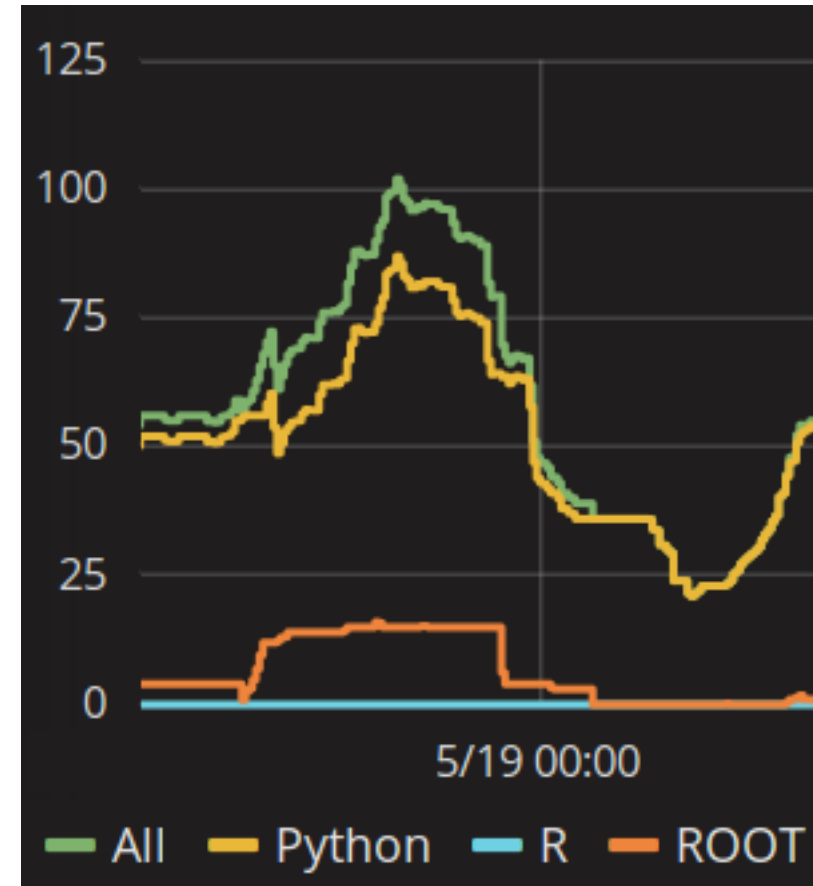# Jupyter - The Notebook as Interface

> Very useful for some use cases
  - Final steps of an analysis
  - Exploration
  - Teaching, documentation
  - Reproducibility

> Interactive, usually lightweight computations
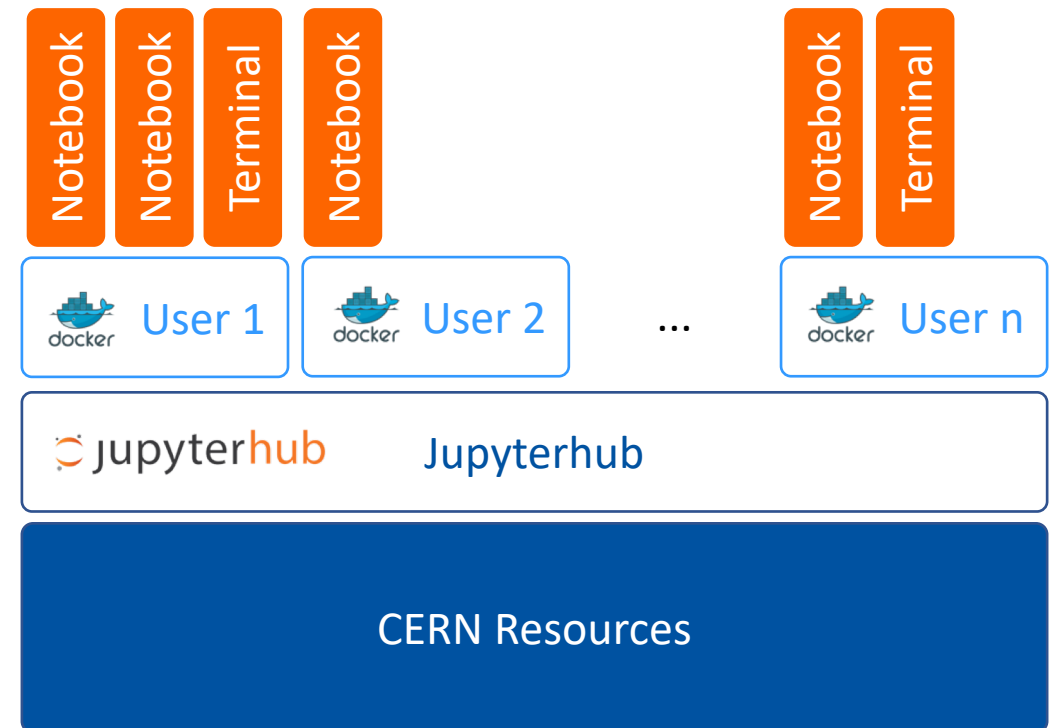
> Languages
  - Not restricted to any
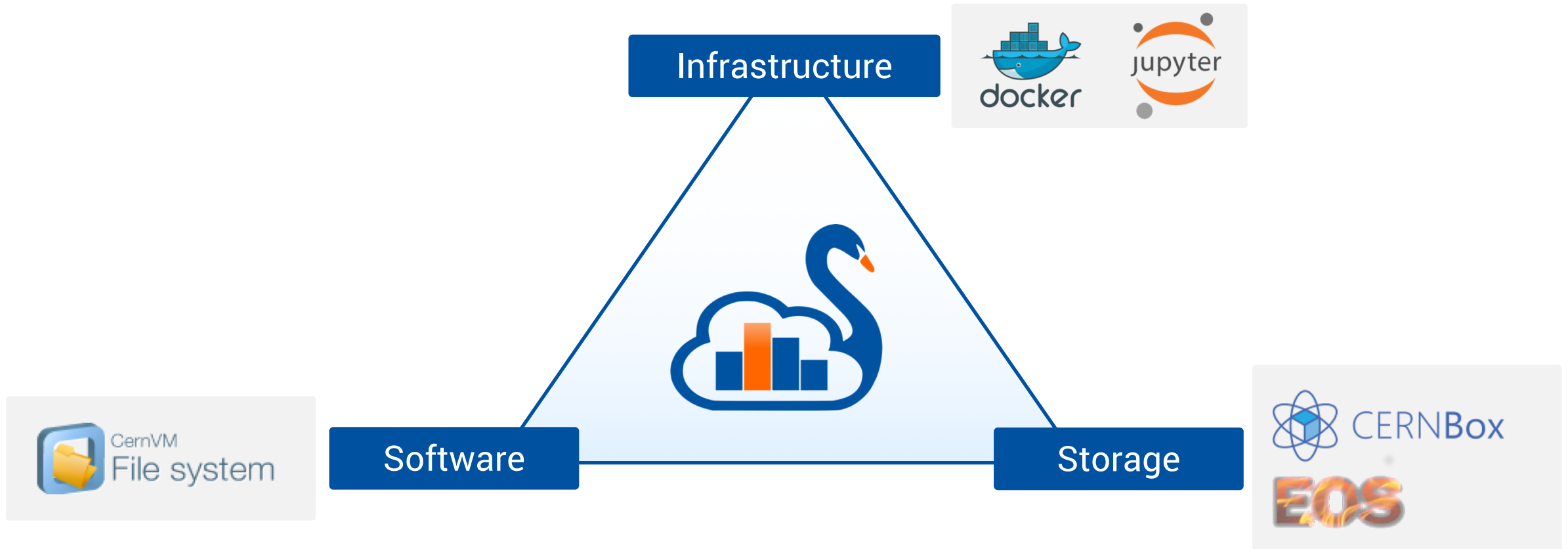  - Python (2&3) clearly dominating in SWAN

# Integrating Jupyter

> Jupyterhub to allow multiple Jupyter instances

> User sessions spawned as Docker containers
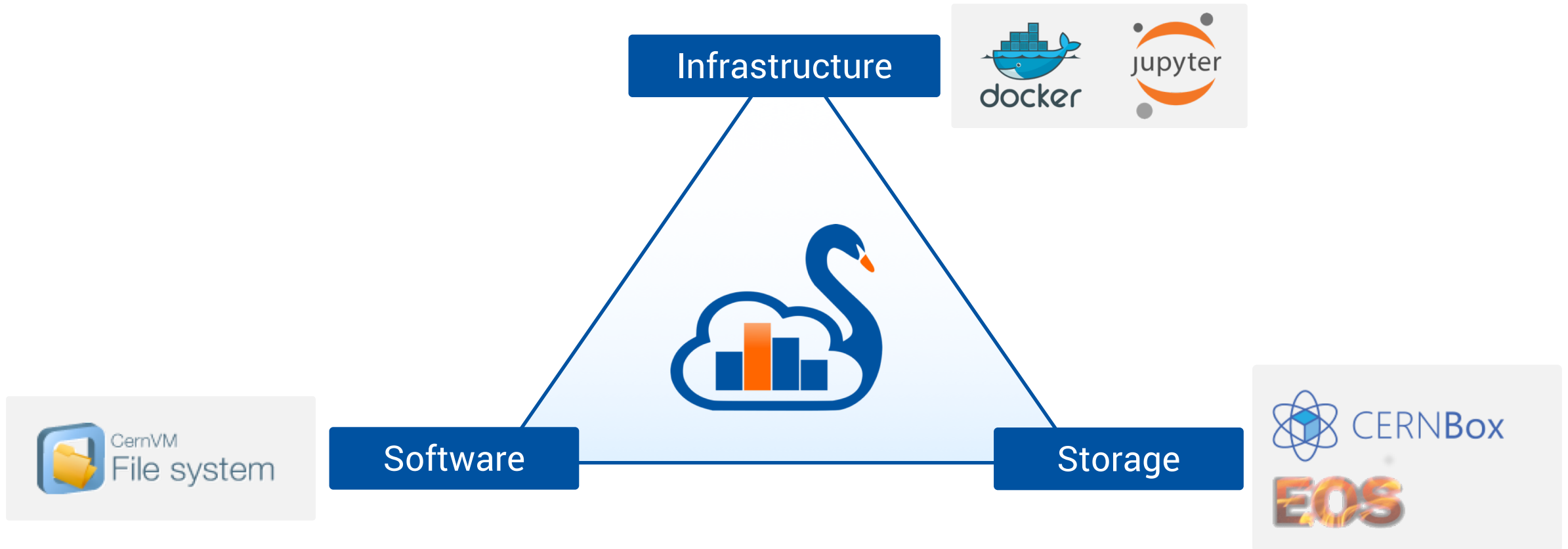  - To guarantee that resources allocated to users are honoured
  - To isolate users work

> Uses EOS mass storage system
  - All experiment data potentially available

> User personal space, synchronized through CERNBox
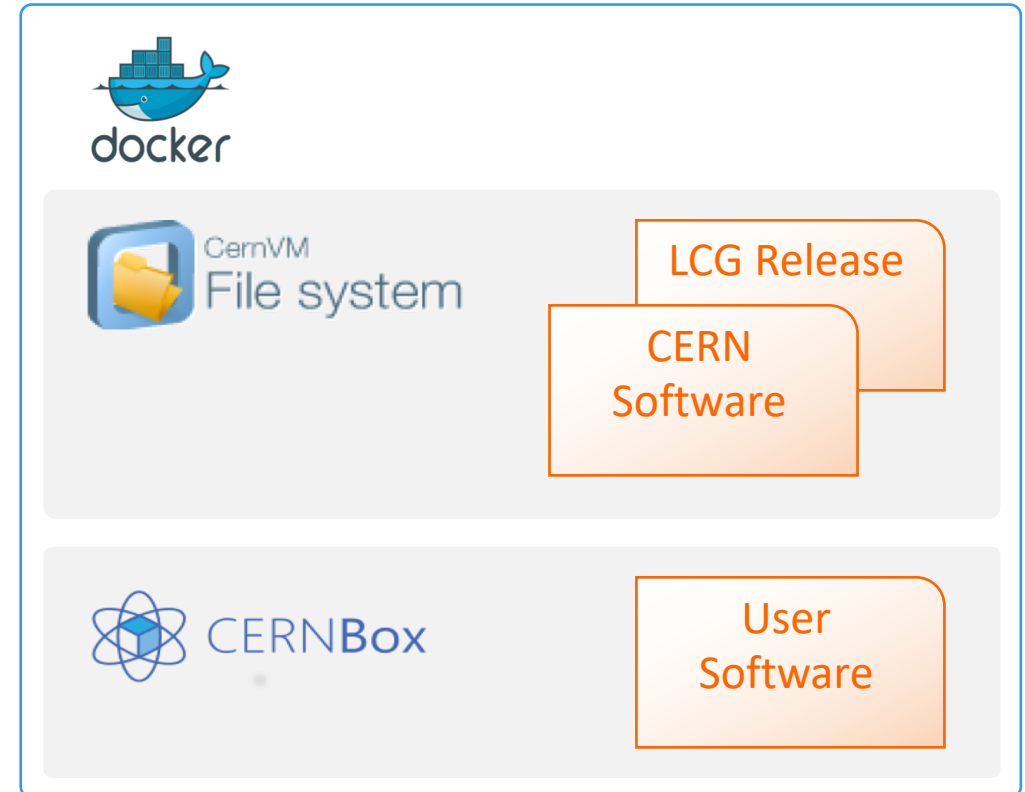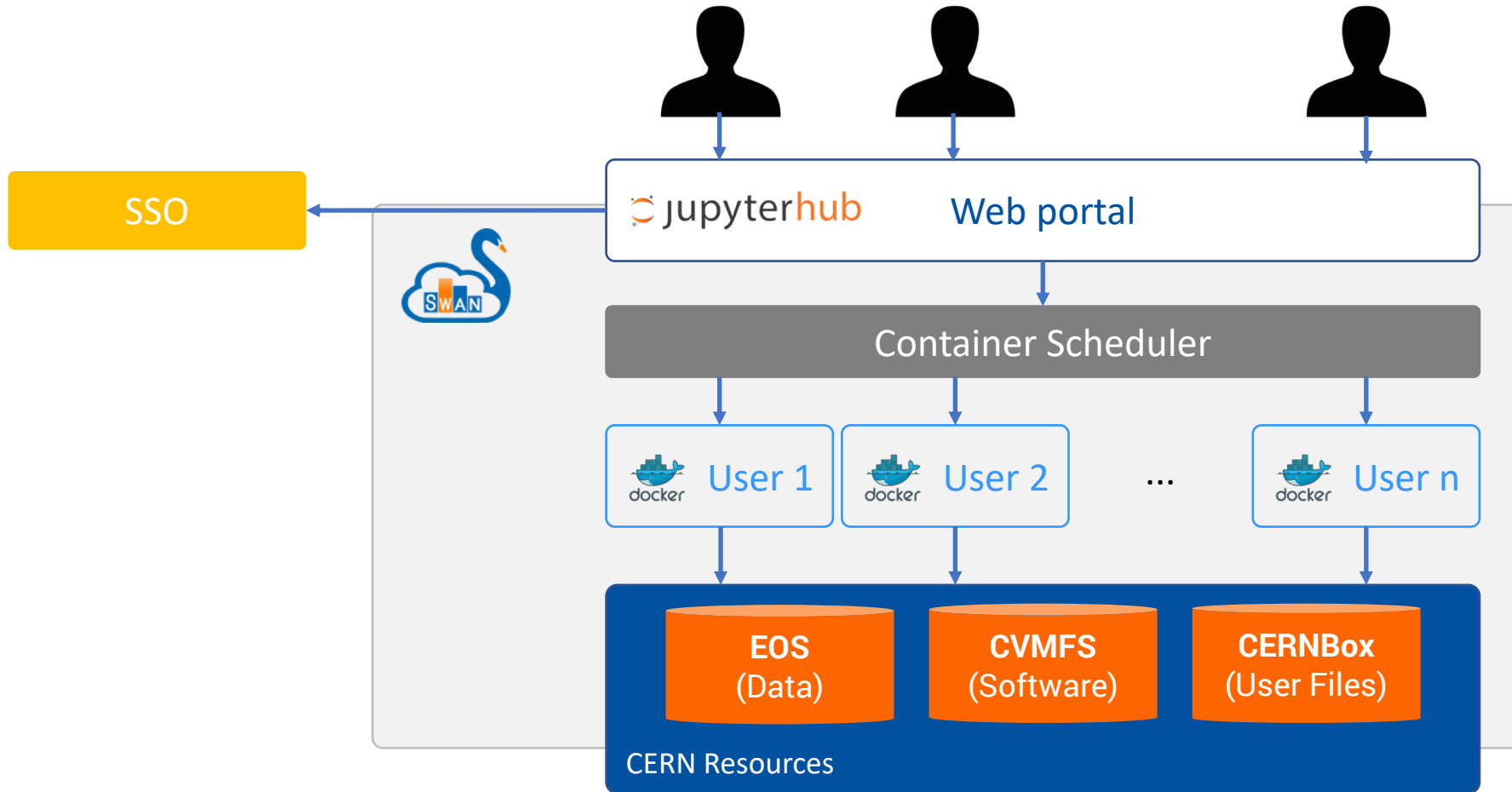  - All files synced across devices, the cloud and other users

# Software

> Software distributed through CVMFS as "LCG releases"
  - A release packs a series of compatible packages

> CVMFS also used by experiments to distribute software
  - Software used by researchers is available

> Multiple languages available
  - C++ (ROOT), Python, R

> Possibility to install other libraries in user local storage

# Architecture



SSO

jupyterhub    Web portal

Container Scheduler

docker User 1    docker User 2    ...    docker User n

**CERN Resources**

EOS (Data)    CVMFS (Software)    CERNBox (User Files)

# Why SWAN matters

# SWAN 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



**Access with only a click**

Almost 50 notebooks in 7 categories

# NXCALS & SWAN

- Spark Web Notebooks (like Jupyter):
  - **Web interface** with built-in Spark integration
  - Data visualisation (**tables, charts**, etc.)
  - Dynamic input forms and **data widgets**
  - Support **work in collaboration** and publishing results online
  - **Natural and very user-friendly for Data Scientists**

SWAN (Service for Web based ANalysis) is a platform to perform interactive data analysis in the cloud.

*Very productive collaboration. Big THANK YOU to our EP-SFT and IT-DB colleagues !*

*To avoid work duplication NXCALS bet on SWAN!*

9/25/2017                    CERN - BE/CO                    5

Some users started saying SWAN is fundamental for their work

20

# The numbers



Grand Total since May 2016 (beginning of monitoring):
- Number of sessions (containers): **~7k**
- Number of notebooks opened: **~14k**

# Other collaborations

> Building block in UP2University European Project

- Bridge the gap between secondary schools, higher education, and the research domain
- SWAN used by students to learn Physics and other Natural Sciences
- Let the kids use the very same tools & services used by real researchers doing Big Science at CERN

> Will be talked in CS3 Workshop in Krakow

- http://cs3.cyfronet.pl/

# Recent developments

# Integration with Spark

> One of the features requested by the community
- Team from the Beams department

> Allow users to connect to CERN Spark Clusters to submit jobs

> In collaboration with CERN Database and Storage groups

# New User Interface

# New User Interface

# New User Interface

```
In [5]:  sc.parallelize(range(0,10)).count()
         sc.parallelize(range(0,20)).count()
```

**Apache Spark:** 1 EXECUTORS  4 CORES  **Jobs:** 2 COMPLETED

| Job ID | Job Name | Status | Stages | Tasks | Submission Time | Duration |
|--------|----------|-----------|--------|-------|-----------------|----------|
| ▶ 3 | count | COMPLETED | 1/1 | 4 / 4 | a few seconds ago | 0s |
| ▶ 4 | count | COMPLETED | 1/1 | 4 / 4 | a few seconds ago | 0s |

```
Out[5]:  20
```

# Sharing made easy

> ## Sharing from inside SWAN interface

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

# Sharing made easy

> Users can clone a shared Project directly from the interface

> New Octave kernel

> Possible integration with HTCondor

# Conclusion

# Conclusion

> SWAN is a CERN service that provides Jupyter Notebooks on demand

> SWAN promotes a cloud based analysis model where users can do analysis only with their browser

> SWAN federates CERN services for software, storage and infrastructure so that users can find what they need in the service

> SWAN fosters collaboration and results sharing between scientists

> SWAN is an Interface for Mass Processing Resources (Spark)

# SWAN: service for web based analysis

Thank you