

SWAN: service for web-based analysis

D. Castro

On behalf of the SWAN team

<https://cern.ch/swan>



Apr 4th, 2019

Architects Forum Meeting



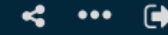
Introduction



SWAN in a Nutshell

- › Analysis only with a web browser
 - No local installation needed
 - Based on **Jupyter Notebooks**
 - Calculations, input data and results “in the Cloud”
- › Support for multiple analysis ecosystems
 - ROOT, Python, R, Octave
- › Easy sharing of scientific results: plots, data, code
- › Integration with CERN resources





Integration of SWAN with Spark clusters

This notebook demonstrates the functionality provided by a SWAN prototype machine that allows to offload computations to an external Spark cluster. The Spark version we are going to use is 2.1.0 and we are going to connect to the analytix cluster (as previously selected in the SWAN web form).

Step 1 - Acquire the necessary credentials to access the Spark cluster.

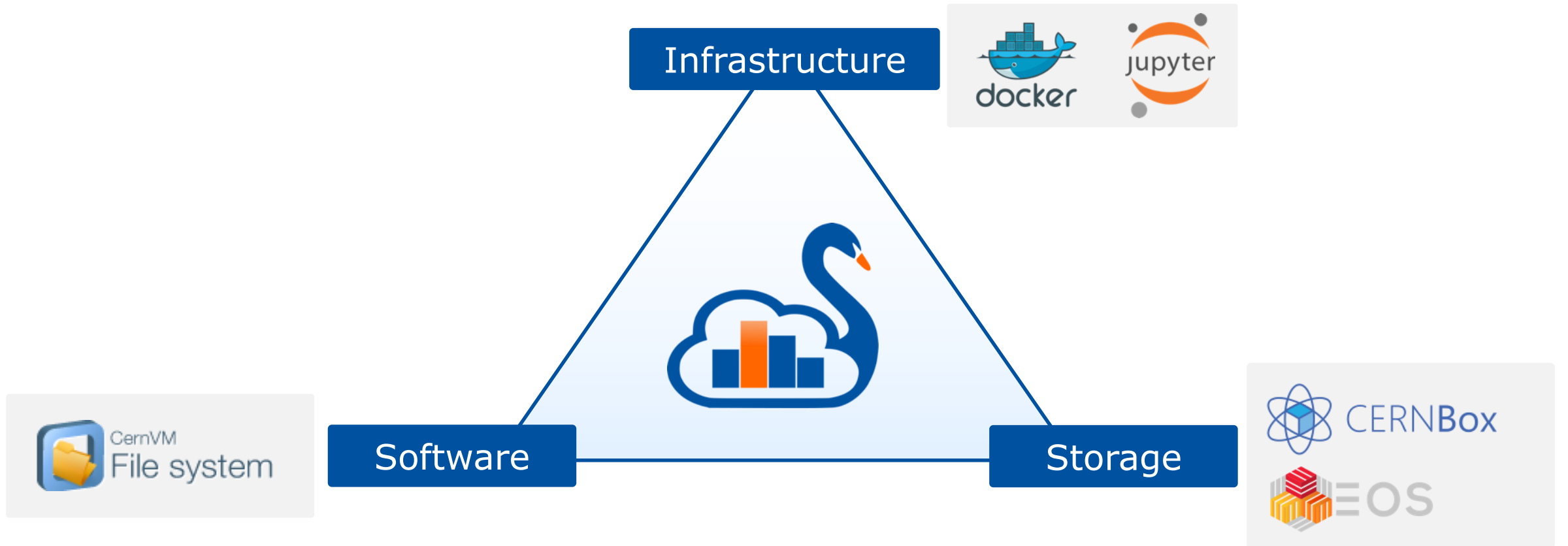
```
In [1]: import getpass
import os, sys, re

print("Please enter your password")
ret = os.system("echo \"%s\" | kinit" % re.escape(getpass.getpass()))

if ret == 0: print("Credentials created successfully")
else:      sys.stderr.write('Error creating credentials, return code: %s\n' % ret)
```



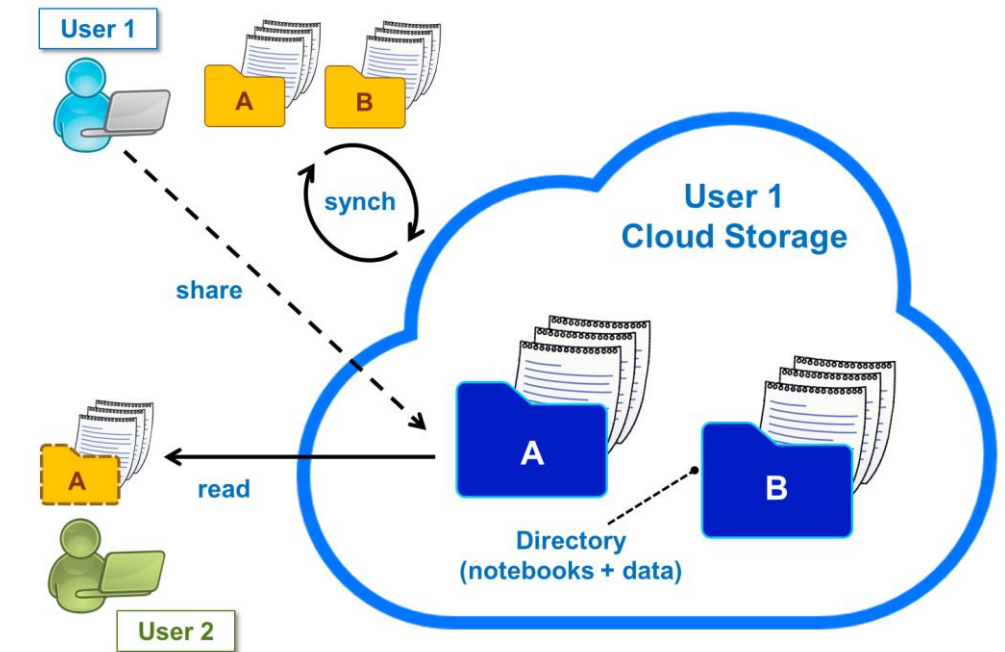
Integrating services





Storage

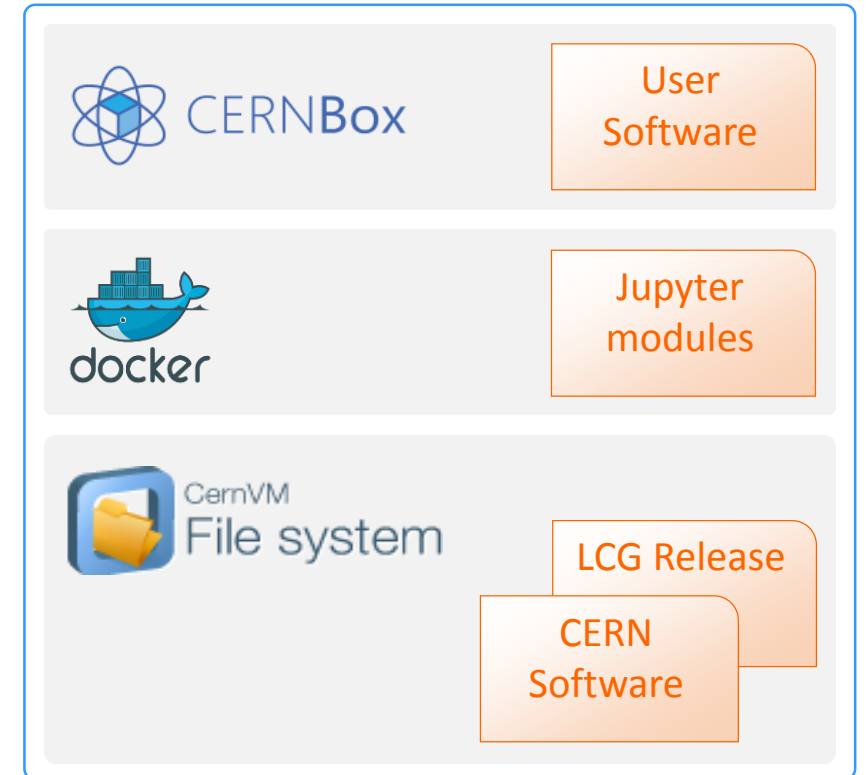
- > Uses EOS disk storage system
 - Same mountpoints as LXPlus
- > CERNBox is SWAN's home directory
 - Storage for your notebooks and data
- > Sync&Share
 - Files synced across devices and the Cloud
 - Collaborative analysis





Software

- > Software distributed through CVMFS
 - LCG Releases
 - Same mountpoints as LXPlus
 - Step towards reproducibility (across time and people)
- > Possibility to install libraries in user cloud storage
 - Good way to use custom/not mainstream packages
 - Configurable environment



SWAN features



Sharing made easy

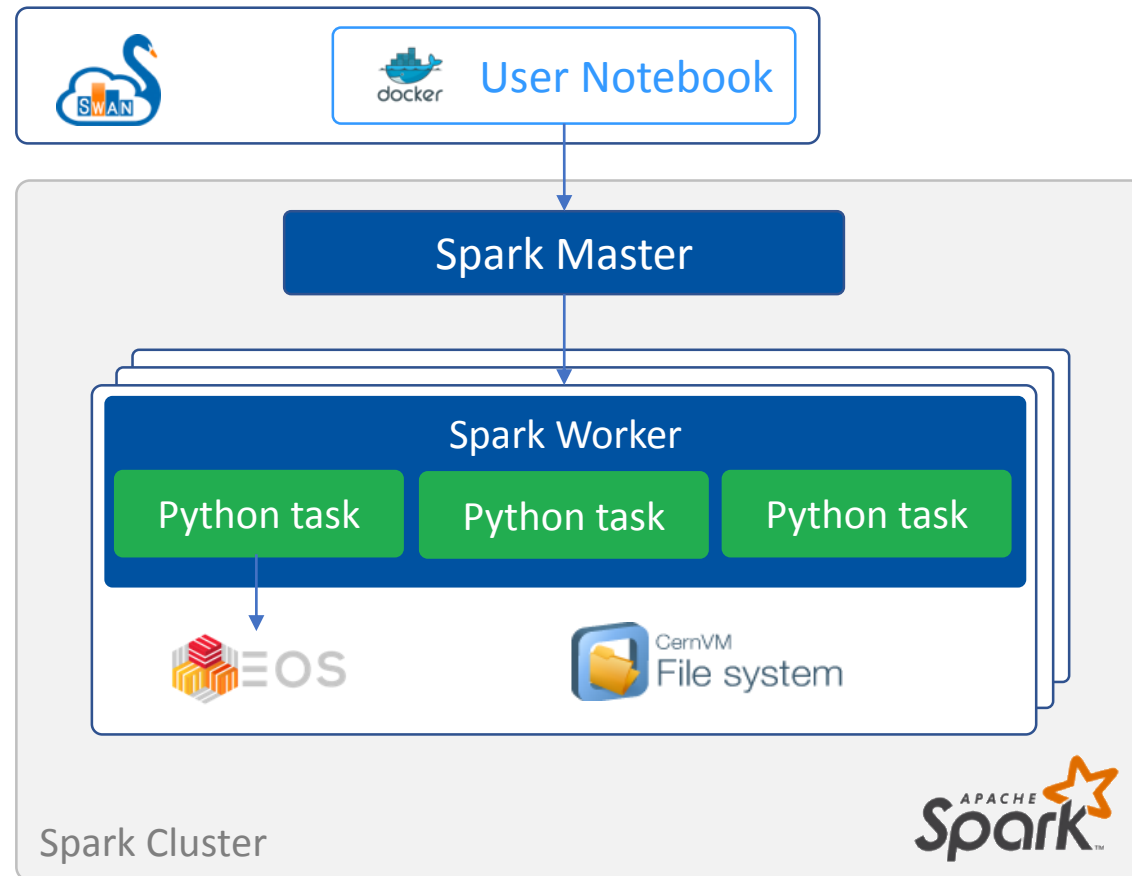
- > Sharing from inside SWAN interface
 - Integration with CERNBox
 - List shares from other users
- > Users can share “Projects”
 - Special kind of folder that contains notebooks and other files, like input data
 - Self contained
- > Concurrent editing not supported *yet* by Jupyter
 - Safer to clone
 - Will be available with Jupyterlab

The screenshot shows the SWAN interface with a 'Share Project' dialog box open. The dialog box has a title bar with a close button. Below the title bar, it says 'You are sharing: Super Real Analysis with TOTEM data'. There is a search bar with the placeholder text 'Start typing to add names...'. Below the search bar, it says 'Shared with' and lists two users: Danilo Piparo (danilo) and Enric Tejedor Saavedra (enric). At the bottom of the dialog, there are two buttons: 'Stop Sharing' (red) and 'Update' (blue). The background of the dialog is dimmed, showing the SWAN interface with the project name 'Super Real Analysis with TOTEM data' and a list of files: 'DistillDistribution.ipynb' and 'dataset.root'.



Integration with Spark

- > Connection to CERN Spark Clusters
- > Same environment across platforms
 - User data – EOS (xrootd)
 - Software – CVMFS
- > Graphical Jupyter extensions developed
 - Spark Connector
 - Spark Monitor





Spark Connector/Monitor



The screenshot shows a Jupyter notebook titled "Spark > Spark_Simple (autosaved)". The notebook content includes a section "Simple example with Spark" and a code cell with the following text:

```
In [ ]: from pyspark import SparkContext
```

Overlaid on the notebook is a "Spark clusters connection" dialog box. It shows the connection target is "hadalytic". Below this, there is a section for "Add a new option" with a text input field containing "Write the option name...". There is also a "Bundled configurations" section with a checkbox for "Include NXCALs options" which is currently unchecked. A "Selected configuration" section lists the following settings:

- spark.shuffle.service.enabled: false
- spark.driver.memory: 2g
- spark.executor.instances: 4

A green "Connect" button is located at the bottom of the dialog.

```
In [ ]: def f(x):  
        global a  
        a+=x  
        RDD9.foreach(f)  
        RDD9.foreach(f)  
        print(a.value)  
        #Display should appear automatically
```



Use cases



Use cases

> Final steps of analysis

- NanoAOD processing with RDataFrame
- Reconstruction of Higgs boson decaying to two Z bosons from events with four leptons

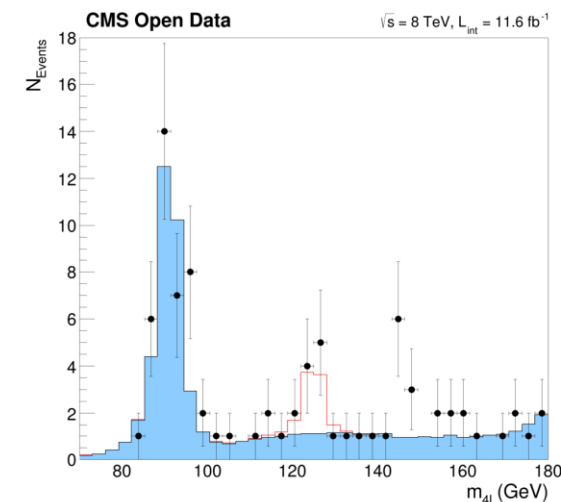
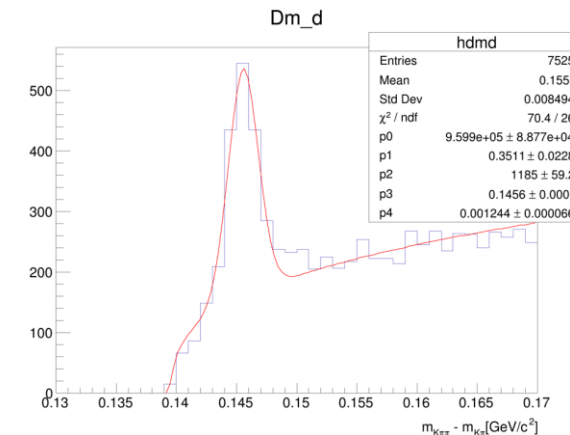
> Exploration

- H1 analysis with RDataFrame

> Teaching

- CERN Summer student courses: ~150 students, data analysis with ROOT
- CERN School of Computing exercises: ~70 students, parallelism
- CERN ATLAS PhD student courses: ~50 students, declarative data analysis
- Machine learning tutorials at CERN

https://root.cern/doc/master/group_tutorial_dataframe.html





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
 - More than 50 notebooks in 8 categories

Example notebooks at cern.ch/swan or inside the service

The screenshot shows the SWAN Gallery interface. At the top, there's a navigation bar with the SWAN logo and a menu icon. Below it, the breadcrumb path is "SWAN > Gallery > Basic Examples". The main heading is "Gallery" with a sub-heading "Basic Examples". A list of categories is shown on the left: Basic Examples, ROOT Primer, Accelerator Complex, FCC, Beam Dynamics, Machine Learning, Apache Spark, and Outreach. The main content area displays a grid of notebook thumbnails. Each thumbnail has a title and a small preview image. The thumbnails are: "Simple ROOTbook (Python)", "Simple ROOTbook (C++)", "Simple Fitting", "Simple I/O", "C++ from Python w/o bindings", "3D Visualisation", "CMS Opendata: di-muon analysis", "Fastjet (Interactive usage of 3rd)", and "Pandas". A dashed line points from the "CMS Opendata: di-muon analysis" thumbnail to an orange callout box.

Access with only a click



Science Box: SWAN on Premises

- › Packaged deployment of SWAN
 - Includes all SWAN components: CERNBox/EOS, CVMFS, JupyterHub
 - Deployable through Kubernetes or docker-compose
- › Some successful external/community installations
 - CERN UP2U Pilot
 - Used by students to learn physics and other sciences
 - PSNC (UP2U EU project)
 - SWAN and CERNBox for students and teachers
 - Open Telekom Cloud (Helix Nebula)
 - TOTEM analysis
 - AARNet
 - Australia's Academic and Research Network

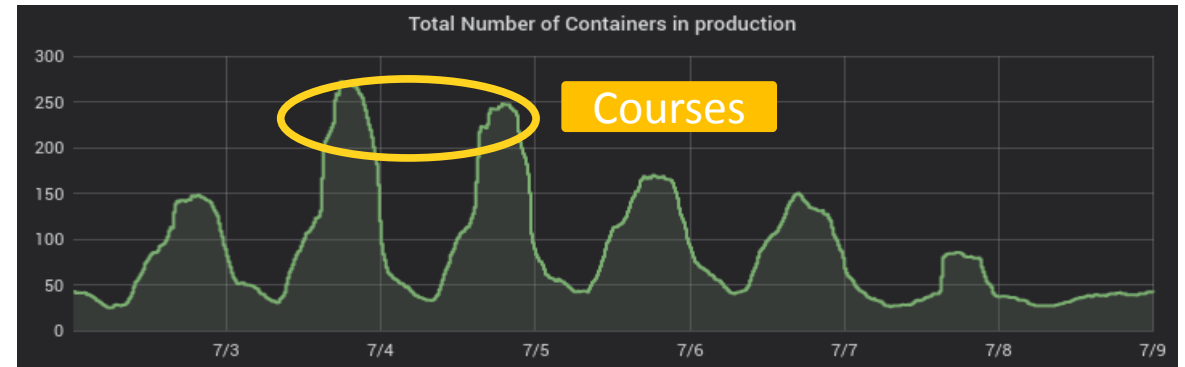


Usage numbers



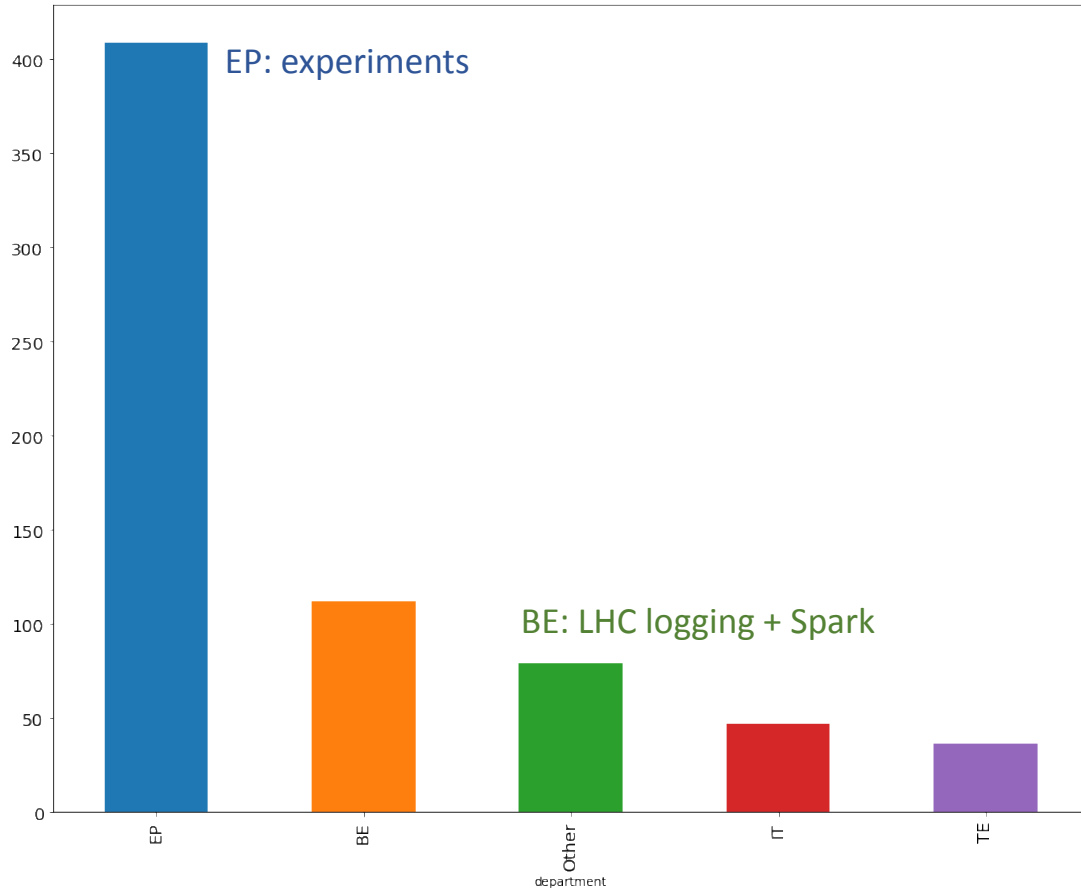
Usage numbers

- > **~200 user sessions** a day on average
 - Users doubled last year with new SWAN interface
- > **~1700 unique users** in 7 months
 - 728 in March
- > Spark cluster connection: 15 – 20 % of users
 - SWAN as entry point for accessing computational resources
 - Used for monitoring LHC accelerator hardware devices (NXCal)

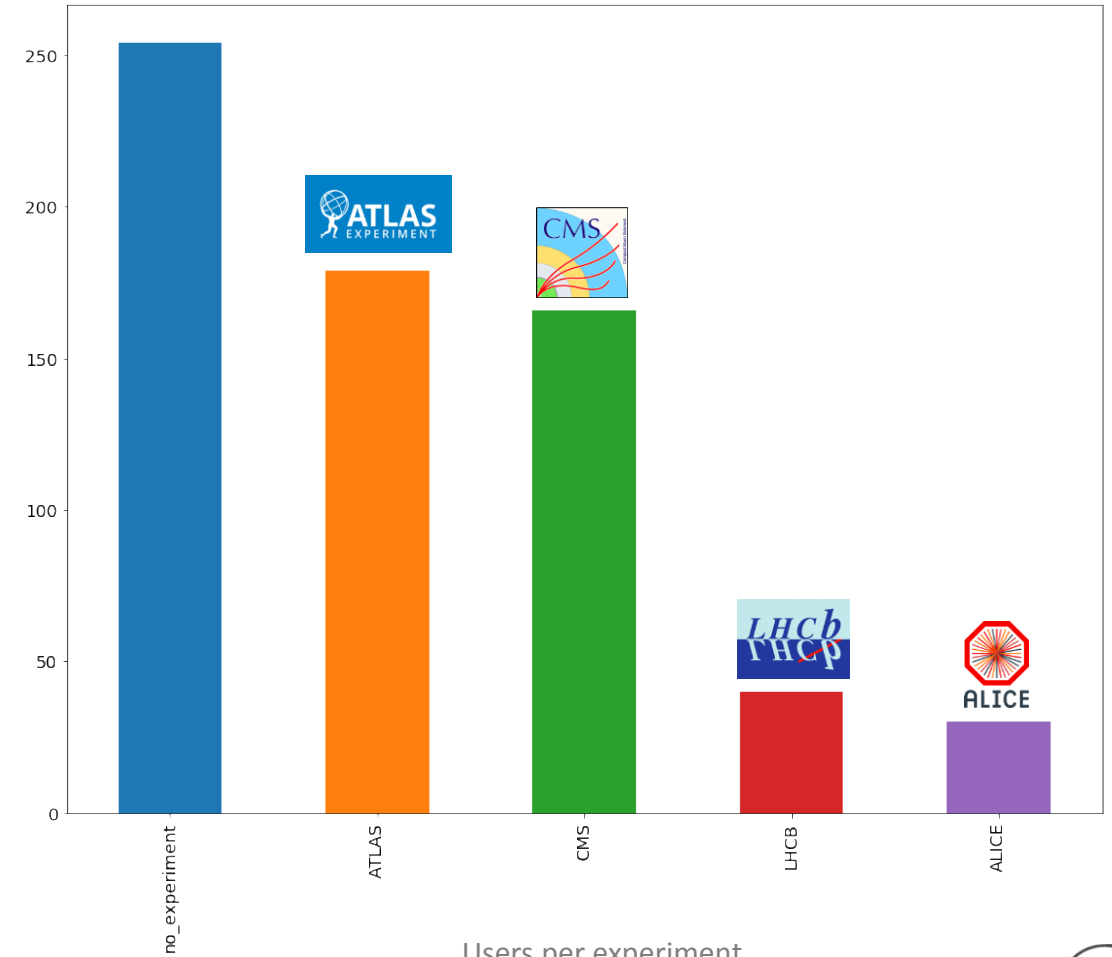




March (top 5)



Users per department

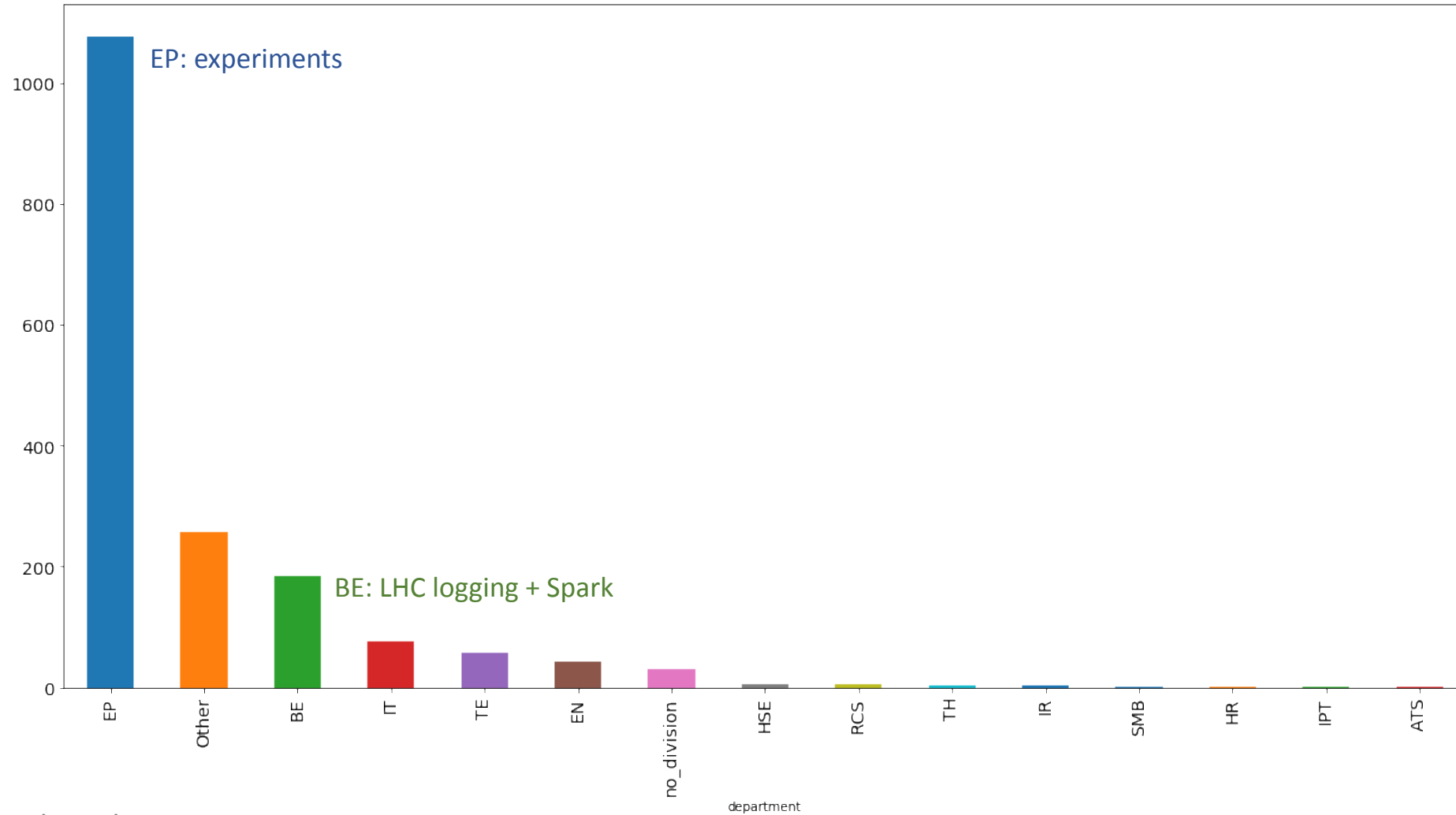


Users per experiment





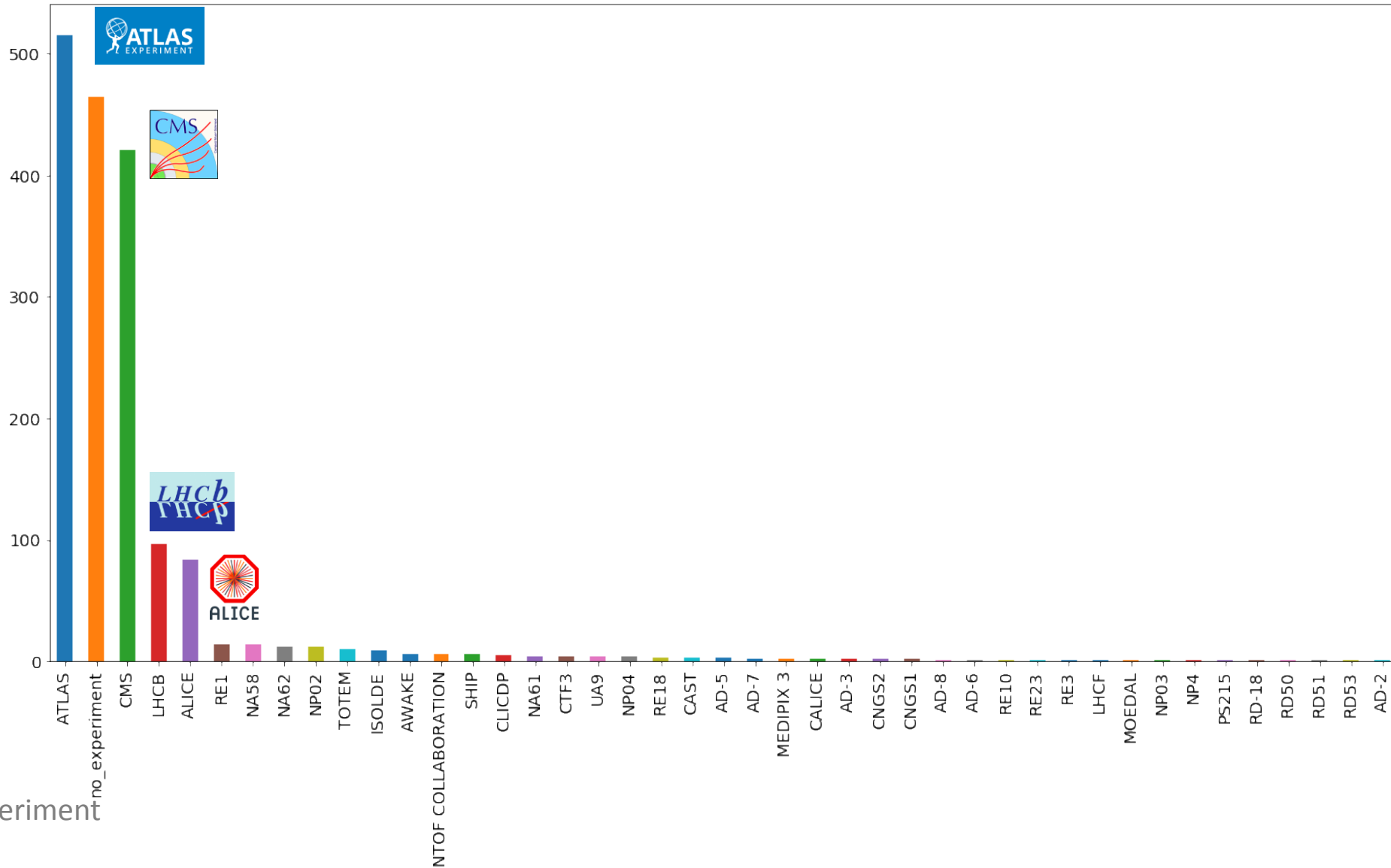
7 months period



Users per department



7 months period



Users per experiment



Looking ahead



Future work/challenges

- > Move to Jupyterlab
 - Porting the current extensions
 - Concurrent editing
- > New architecture
 - Based on Kubernetes
- > Exploitation of GPUs
 - Speed up computation of GPU-ready libraries (e.g. TensorFlow)
 - Interactive access to GPU resources (complementary to the batch access to GPU)
- > Ongoing effort:
 - Submit batch jobs from the notebook

Conclusion



Conclusions

- › SWAN is a CERN service that provides Jupyter Notebooks on demand
 - Promotes a cloud-based analysis model
 - Federates CERN services for software, storage and infrastructure
 - Deployable on premises
- › SWAN fosters collaboration and results sharing between scientists
- › The new Jupyterlab interface will bring new possibilities for collaborative analysis
 - With the introduction of concurrent editing of notebooks
- › SWAN became a fundamental Interface for Mass Processing Resources (Spark)
 - Not only for Physics analysis but also for monitoring the LHC hardware
- › Usage is growing
 - Missing dedicated manpower to ensure its sustainability

SWAN: service for web-based analysis

Thank you

Diogo Castro
diogo.castro@cern.ch