# SWAN: interactive data analysis on the web

**Prasanth KOTHURI, IT-DB**
On behalf of the SWAN team

https://cern.ch/swan
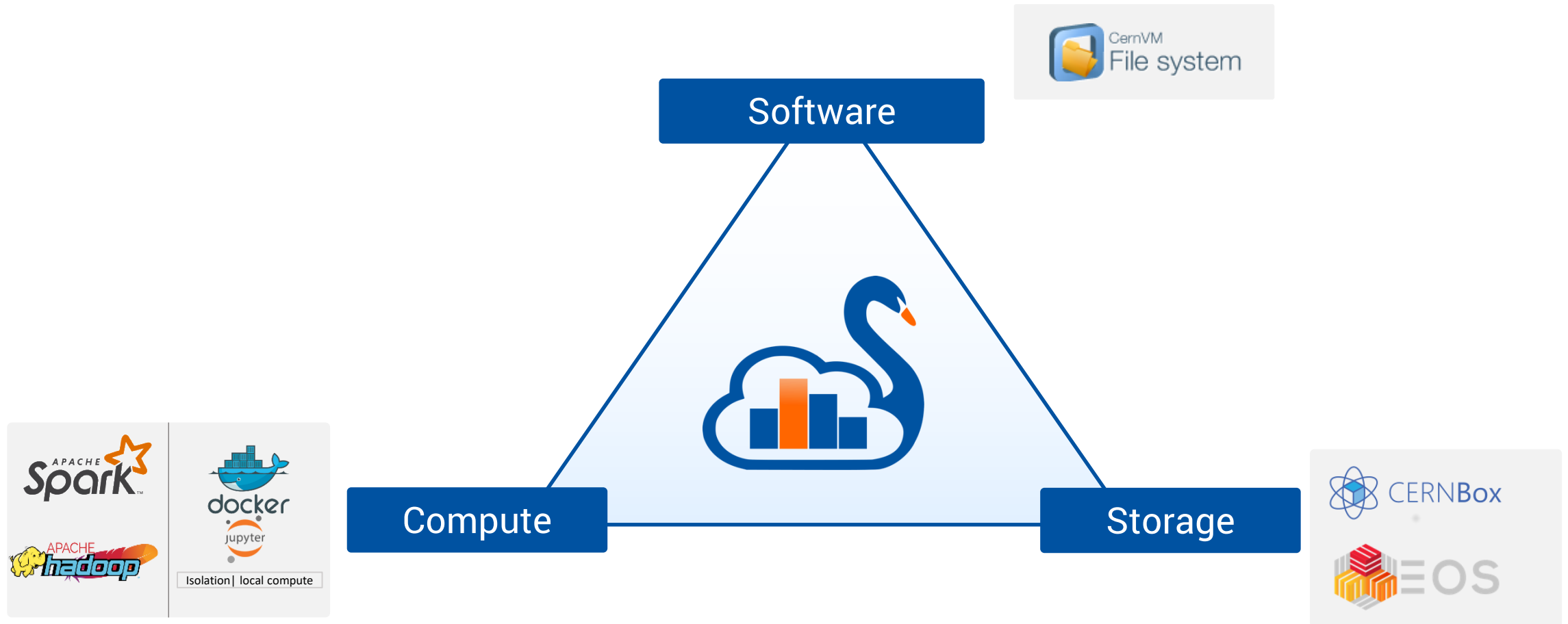
**Sep 18th, 2019**
PE Mini Lectures

# 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 and languages
>> - Python, ROOT C++, R and Octave

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

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

# Integrating services

# Jupyter - The Notebook as Interface

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

> Interactive, usually lightweight computations and distributed parallel processing capability with the integration of mass processing system (Apache Spark)

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

# User Interface

FILE   EDIT   VIEW   INSERT   CELL   KERNEL   NAVIGATE   WIDGETS   HELP

Not Trusted          Python 2

Code

## 2 Displaying graphics

We can now draw the histogram. We will at first create a canvas, the entity which in ROOT holds graphics primitives. Note that thanks to JSROOT, this is not a static plot but an interactive visualisation. Try to play with it and save it as image when you are satisfied!

```
In [5]: c = ROOT.TCanvas()
        h.Draw()
        c.Draw()
```

My Histo

| myHisto | |
| --- | --- |
| Entries | 1000 |
| Mean | 0.02680 |
| Std Dev | 1.038 |

Y axis

X axis

We'll try now to beautify the plot a bit, for example filling the histogram with a colour and setting a grid on the canvas.

```
In [6]: h.SetFillColor(ROOT.kBlue-10)
        c.SetGrid()
        h.Draw()
        c.Draw()
```

My Histo

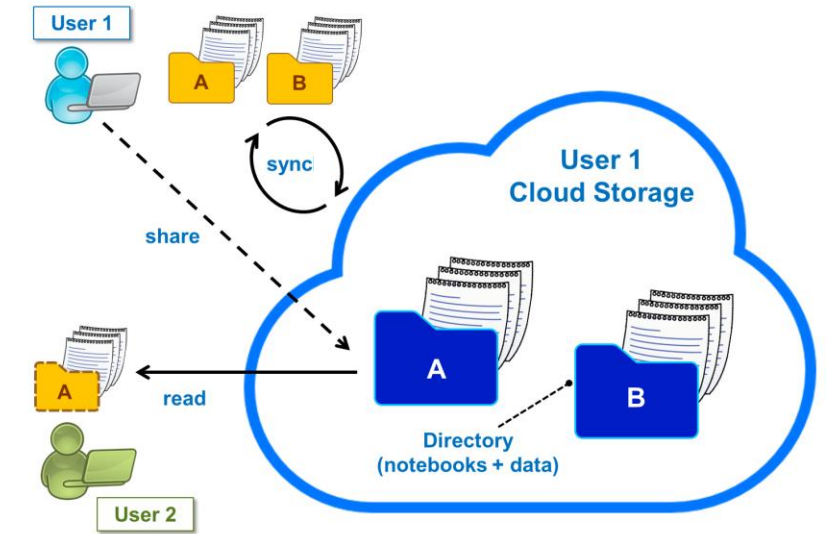| myHisto | |
| --- | --- |
| Entries | 1000 |
| Mean | 0.02680 |

7

# Cloud storage as your Home

> CERNBox is SWAN's home directory
- Storage for your notebooks and data
- Based on ownCloud
- 6PB of user data, 16k users

> Uses EOS disk storage system
- All experiment data potentially available
- 250PB of experimental data at CERN (LHC and others)

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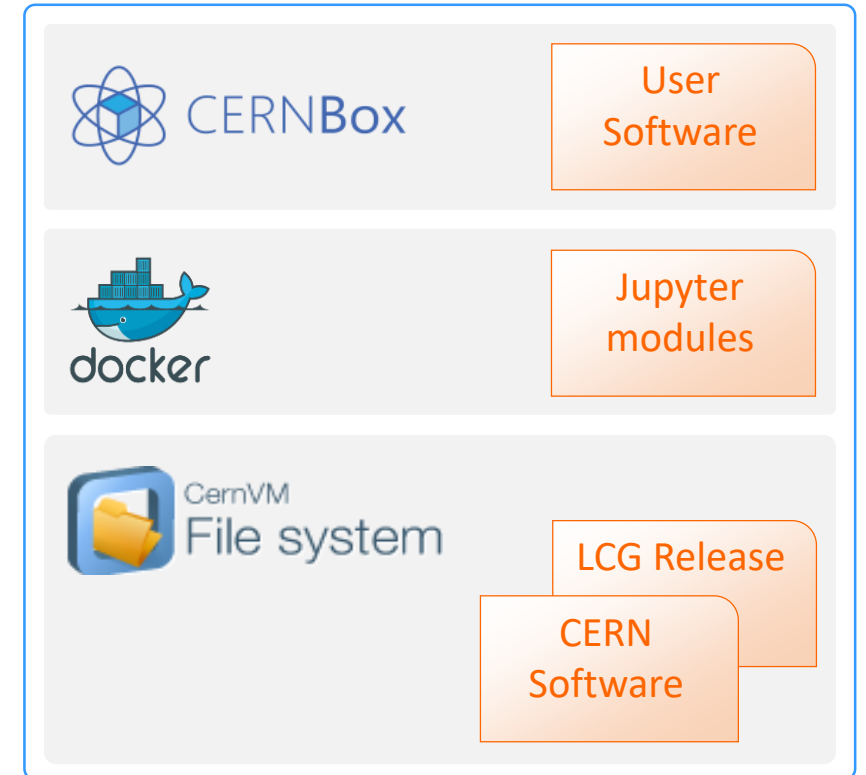- Collaborative analysis

# 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

# Software

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

CERNBox — User Software

docker — Jupyter modules

CernVM File system — LCG Release / CERN Software

# Architecture

# Access to Computing Resources

# Integration with Spark

> Connection to CERN Spark Clusters
>> ▪ Spark: general purpose distributed computing framework
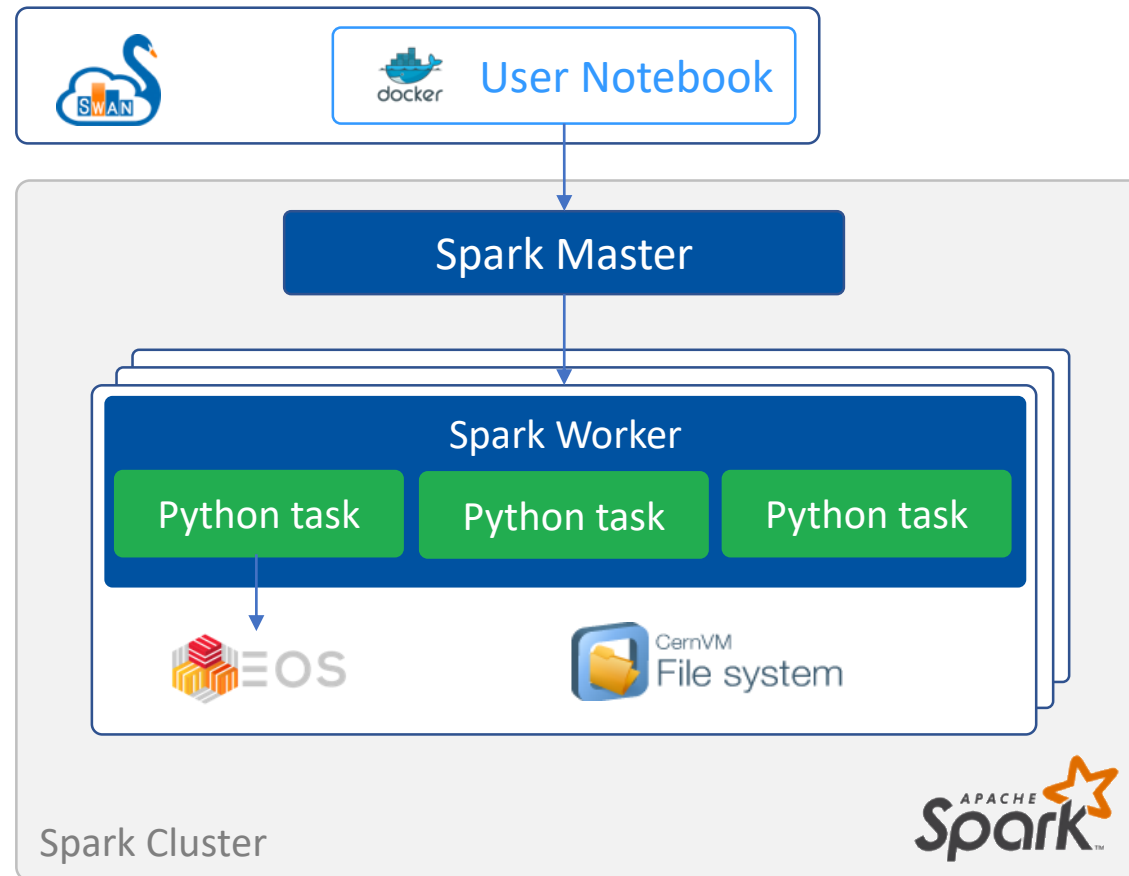
> Same environment across platforms (local/remote)
>> ▪ User data - EOS
>> ▪ Software - CVMFS

> Graphical Jupyter extensions developed
>> ▪ Spark Connector
>> ▪ Spark Monitor

> Not only used for Physics Analysis at CERN

> Spark Clusters
>> ▪ nxcals – Dedicated cluster for accelerator logging
>> ▪ Analytix – General purpose YARN cluster
>> ▪ Cloud Containers – General purpose Kubernetes cluster

# Spark Connector



- <u>Spark Connector</u> – handling the spark configuration complexity
  - User is presented with Spark Session (Spark) and Spark Context (sc)
  - Ability to bundle configurations specific to user communities
  - Ability to specify additional configuration

# A configurable system

# 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 community installations…
  - AARNet
  - PSNC
  - Open Telekom Cloud (Helix Nebula)

> …with different storage integrations
  - OwnCloud WebDav access (AARNet)

# Looking ahead

# Future work/challenges

> Move to Jupyterlab
  - Porting the current extensions
  - Concurrent editing

> New architecture
  - Based on Kubernetes

> Exploitation of GPUs
  - HEP is looking to ML
  - Speed up computation of GPU-ready libraries (e.g. TensorFlow)

# Demo

# Where to find us

# Where to find us

> Contacts
  - [swan-admins@cern.ch](mailto:swan-admins@cern.ch)
  - [http://cern.ch/swan](http://cern.ch/swan)

> Repository
  - [https://github.com/swan-cern/](https://github.com/swan-cern/)

> Science Box
  - [https://cern.ch/sciencebox](https://cern.ch/sciencebox)

# Conclusion

# Conclusion

> SWAN is a CERN service that provides Jupyter Notebooks on demand
- Promotes a cloud-based analysis model
- 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) at CERN
- Not only for Physics analysis but also for monitoring the LHC hardware

> Successfully deployed outside of CERN premises
- Personalized to fit the local infrastructure
- Ongoing effort to allow interoperability (CS3APIs, Science Mesh)

# SWAN and its analysis ecosystem

Thank you

Prasanth KOTHURI

prasanth.kothuri@cern.ch