Boosting productivity of researches integrating IT services into an interactive analysis platform



D. Piparo (CERN EP-SFT) for the SWAN team

https://cern.ch/swan

Aug 31st, 2018 GridKa School



The Work Behind These Slides

- I am here on behalf of the SWAN team, spanning two CERN Departments: Experimental Physics and Information Technology
- > The team: E. Tejedor, D. Castro, P. Mato, E. Bocchi, J. Moscicki, M. Lamanna, P. Kothuri, D. P.
- > I'd like to thank Diogo Castro and Enric Tejedor for contributing so much to this talk.

Introduction





Accelerating Science







https://kt.cern/success-stories/cern-virtual-machin e-file-system-euclid-science-data-centres





https://root.cern.ch/code-snippet/learning-root

Research















://root.cern.ch/code-snippet/learning-root







Distribution of All CERN Users by Nationality on 24 January 2018



7

~17.5k people in total ~12.8K International collaboration

m

The steps towards data analysis



Service for Web based Analysis



Not Only HEP Analysis: LHC Controls Data



CERM

Not Only HEP Analysis: LHC Controls Data



Service for Web based ANalysis: SWAN





- The physics Data Analysis Framework
 Created at CERN
- > Used for data analysis, I/O, ML, advanced graphics...
- > Integrated with Jupyter notebooks
 - First C++ kernel based on ROOT C++
- > The origin story of SWAN
- > <u>http://root.cern</u>





Service for Web based Analysis

> Analysis only with a web browser

- Blur boundaries between web and local applications
- Available everywhere and at anytime
- Providing access also to Jupyter notebooks
- > Easy to use (but powerful)
 - No local installation and configuration needed
- > Create easily shareable scientific results: plots, data, code
 - Storage is crucial: mass & synchronized
- > Integration with CERN resources
 - Access software, user/experiments data, mass processing power
- > Synergy among analysis ecosystems : ROOT, R, Python, ...



https://en.wikipedia.org/wiki/Argentina









Open Source Technologies









CERN

Does it Have a Terminal?





bash-4.1\$ root -b Welcome to ROOT 6.12/06 http://root.cern.ch Built for linuxx8664gcc From tag v6-12-06, 9 February 2018 Try '.help', '.demo', '.license', '.credits', '.quit'/'.q' root [0] .q bash-4.1\$ ls /cvmfs alice.cern.ch alice-ocdb.cern.ch atlas-condb.cern.ch fcc.cern.ch lhcb.cern.ch sft-nightlies.cern.ch alice.nightlies.cern.ch atlas.cern.ch cms.cern.ch geant4.cern.ch sft.cern.ch bash-4.1\$ ls /eos experiment fcc hepdata opendata project user bash-4.1\$





Uses EOS disk storage system >

All experiment data potentially available

CERNBox is SWAN's home directory >

- Storage for your notebooks and data
- Sync&Share >
 - Files synced across devices and the Cloud
 - Collaborative analysis









> Software distributed through CVMFS

- "LCG Releases" pack a series of compatible packages
- Single Docker image, many software stacks
- Reduced Docker Image size
- Lazy fetching of software
- > Possibility to install libraries in user cloud storage
 - Good way to use custom/not mainstream packages
 - Configurable environment
- Multiple "kernels" ("links" between Jupyter and a language) available
 - Python 2 or 3, ROOT C++, R, Octave







CERN

Leveraging the power of Jupyter

Service for Web based Analysis



Extensions, extensions everywhere





New User Interface

a	Configure Environment	a	۵
	Specify the parameters that will be used to contextualise the container which is created for you. See the online SWAN guide for more details.		
	91 \$		
	Platform more		
	x86_64-slc6-gcc62-opt \$		
	e.g. \$CERNBOX_HOME/MySWAN/myscript.sh		
SWAN	Number of cores more 2	SWAN	
	Memory more		
	8 GB		
	Hadalytic \$	Ctarting your appaign	
		Starting your session	
	Always start with this configuration		
	Start my Session	Waiting for swan-qa004,cern.ch	



CERN



a	Projects	Share	CERNBox		>_ ••• 🕩
SWAN > My Projects					
My Projects					(\pm)
NAME .				STATUS	MODIFIED
Proj1				4	5 days ago
Proj2					15 days ago
Project					21 days ago
Project 1					2 months ago
Project 2					4 months ago
ProjTest					15 days ago
Spark					7 days ago
SWAN-Spark_NXCALS_Example					20 days ago
🗑 teste					19 days ago

SWAN Copyright CERN 2017. All rights reserved. Home | Contacts | Support | Report a bug | Imprint





Collaborative Analysis

Service for Web based Analysis



Sharing made easy

- Sharing from inside
 SWAN interface
 - Integration with CERNBox
- > Users can share "Projects"
 - Special kind of folder that contains notebooks and other files, like input data
 - Self contained

&	Projects Share	Share Project	×
SWAN > My Projects > Super Real Analy	sis with TOTEM data	You are sharing: Super Real Analysis with TOTEM data	
Super Real Analysis wi	th TOTEM data 🗠	Search by name or username. Use "a:" for secondary accounts.	
□ NAME ▼		Start typing to add names Shared with	
DistillDistibution.ipynb		Danilo Piparo (danilo)	
SWAN © Copyright CERN 2016-2018. All rights re Home Contact Support Report a bug	iserved.	Stop Sharing Upd	ate

The Share tab

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

a	Projects	< Share	CERNBox		>_ ••• 🕞
SWAN > Share					
Projects shared with me	~				
NAME -			SIZE	SHARED BY	DATE
UP2University Pilot			Empty	jupytercon2	7 minutes ago
NAME * Higgs Boson discovery Super Real Analysis with TOTEM data	<u></u>			SHARED WITH 2 people/groups diogo	DATE 18 hours ago 19 hours ago
SWAN © Copyright CERN 2016-2018. All rights rese Home Contact Support Report a bug	rved.				CERN



Inspecting a Project

- > Users can inspect shared project contents
 - Browsing of the files
 - Static rendering of notebooks
- Useful to decide whether to accept or not the shared project

	Simple ROOTbook (C++)
	This simple ROOTbook shows how to create a histogram, fill it and draw it. The language chosen is C++.
	In order to activate the interactive visualsisation we can use the <u>JSROOT</u> magic:
In [1]:	%jsroot on
	Now we will create a histogram specifying its title and axes titles:
In [2]:	THIF h("myHisto", "My Histo; X axis; Y axis", 64, -4, 4)
	(THIF &) Name: myHisto Title: My Histo NbinsX: 64
	If you are wondering what this output represents, it is what we call a "printed value". The ROOT interpreter can indeed be instructed to "print" according to
	certain rules instances of a particular class. Time to create a random generator and fill our histogram:
In [3]:	<pre>Trandom3 rndmGenerator; for (auto i : ROOT::TSeqI(1000)){ auto rndm = rndmGenerator.Gaus(); h.Fill(rndm); }</pre>
In [3]:	<pre>We can now draw the histogram. We will at first create a <u>canvas</u>, the entity which in ROOT holds graphics primitives.</pre>
In [3]: In [4]:	<pre>injoid include instances of a particular class. Time to create a random generator and fill our histogram: TRandom3 rndmGenerator; for (auto i: ROOT:rSseqI(1000)){ auto rndm = rndmGenerator.Gaus(); h.Fill(rndm); } We can now draw the histogram. We will at first create a <u>canvas</u>, the entity which in ROOT holds graphics primitives. TCanvas c;</pre>
In [3]: In [4]: In [5]:	<pre>injoid include instances of a particular class. Time to create a random generator and fill our histogram: TRandom3 rndmGenerator; for (auto i: ROOT:TSEqI(1000)){ auto rndm = rndmGenerator.Gaus(); h.Fill(rndm); } We can now draw the histogram. We will at first create a <u>canvas</u>, the entity which in ROOT holds graphics primitives. TCanvas c; h.Draw(); c.Draw(); }</pre>

Accepting a Shared Project

- If accepted, project is cloned to the receiver's CERNBox
 - The receiver will work on his own copy
- > Concurrent editing not (yet?) supported by Jupyter
 - Safer to clone

	Projects	< Share	CERNBox		>_ ••• 🕞
SWAN > Share					
Projects shared with me	^				
NAME 👻			SIZE	SHARED BY	DATE
UP2University Pilot		Clone	Empty	jupytercon2	7 minutes ago
NAME - Higgs Boson discovery				SHARED WITH 2 people/groups	DATE 18 hours ago
NAME 👻				SHARED WITH	DATE
Super Real Analysis with TOTEM data				diogo	19 hours ago
SWAN © Copyright CERN 2016-2018. All rights resen Home Contact Support Report a bug ript:	ved.				CERN



Access to Computing Resources

Spark integration



> Beams department at CERN built and operates the LHC

- Data: monitoring LHC accelerator hardware devices
- Analysis of logs
- Different from physics data of the experiments

> With increase in data, they adopted Spark

- Spark: a modern cluster management system allowing distributed and interactive data analysis
- But it was missing a unified platform for analytics visualization



Integration with Spark

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









> Allows to monitor Spark jobs from within the notebook

 Interactive display appears automatically in the output of the cell

> Multiple views

- Progress bars for jobs and stages
- Task timeline
- Resource utilization



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



Access to Computing Resources

Batch jobs



Connecting More Resources

- Ongoing effort: submit
 batch jobs from the notebook
 - Monitoring display
 - Jobs tab



				Proje	cts Share	≡ Jobs	CERNBox		>_ ••• 🗭
S	Select jobs to p	perform actio	n on them.						Job ID Jobs per page -
	0.	Job ID	Job Name	Backend	Application	File Name	Status	Submission Time	Runtime
		3		Condor	Executable		SUBMITTING	Jul 24th, 3:26 pm	2
		2		Condor	Executable		NEW	121	2
		1		Localhost	Executable		NEW	(#.	Tre
		0		Localhost	Executable		COMPLETED	Jul 24th, 3:21 pm	00 seconds

In [6]:	۴%g j=g j.s	anga anga. ubmit	Jop()									
		S	Backend	LOCALHOST	Application:	EXECUTABLE	Splitter:	None	0 SUBJOBS			×
		Job	ID .	lob Name	State	JS		Subjo	bs	Submission Time	Runtime	
		0			COMPLI	ETED		No Sub	jobs	Jul 24th, 3:21 pm	00 seconds	



Education, Outreach



Notebooks for Teaching

- > SWAN and notebooks are very useful for teaching
 - Teachers prepare exercises in the form of notebooks
 - Students complete and run the exercises
- > Used frequently as support for:
 - 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





Science Box: SWAN on Premises

> UP2University European Project

- Bridge the gap between secondary schools, higher education and the research domain
- Partner universities (OU, UROMA, NTUA, ...), pilot schools
- <u>http://up2university.eu</u>
- > SWAN used by students to learn physics and other sciences
 - Let them use the very same tools & services used by scientists at CERN
 - Integrating Jupyter in teaching environment
 - Complete SWAN package, with CERNBox, EOS and CVMFS











> ~200 user sessions a day on average

- Users doubled this year with new SWAN interface
- > Spark cluster connection: 15 20 % of users
 - SWAN as entry point for accessing computational resources





- > 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

Example notebooks at cern.ch/swan



Basic Examples

This is a gallery of basic example notebooks: click on the images to inspect the underlying document, open in SWAN the single notebooks or the full git repository!

pen in 🔬 SWAN

Many of the notebooks are ROOTbooks, based on the ROOT framework. To know more about ROOT, visit root.cern.ch.



Education, Too



Outreach

Reach out with SWAN! This section collect a series of outreach efforts involving SWAN.

Particle open data teaching (Hiukkasfysiikan avoin data opetuksessa)



an introductory course about experimental HEP for future high school teachers. The result is great: check it out in SWAN!



Github repository: https://github.com/cmsopendata-finland/kurssimateriaali

Example notebooks at cern.ch/swan

 Finnish High School program



Integrate Software Packages



Example notebooks at root.cern

These examples show the functionalities of the RDataFrame class.

Files

file	df001_introduction.C View Notebook Open in Swam This tutorial illustrates the basic features of the RDataFrame class, a utility which allows to interact with data stored in TTrees following a functional-chain like approach.
file	df001_introduction.py View Notebook Open in SWAN This tutorial illustrates the basic features of the RDataFrame class, a utility which allows to interact with data stored in TTrees following a functional-chain like approach.
file	df002_dataModel.C View Notebook Open in SWAN This tutorial shows the possibility to use data models which are more complex than flat ntuples with RDataFrame
file	df002_dataModel.py View Notebook Open in SWAN This tutorial shows the possibility to use data models which are more complex than flat ntuples with RDataFrame
file	df003_profiles.C View Notebook Open in SWAN This tutorial illustrates how to use TProfiles in combination with the RDataFrame.
file	df003_profiles.py View Notebook Open in Swall This tutorial illustrates how to use TProfiles in combination with the RDataFrame.



Looking ahead



Future work/challenges

- > Move to Jupyterlab
 - Porting the current extensions
 - ROOT is already migrating its graphics to JavaScript only
 - Ultimate blur between local/cloud resources and apps
- > New architecture made to scale
 - Based on Kubernetes, leveraging CERN Container Service
- > Exploitation of GPUs
 - ML found its place in HEP
 - Speed up computation of GPU-ready libraries (e.g. TensorFlow)
- > Moving beyond LHC
 - New challenges from HL-LHC
 - Even small dataset will not fit inside users laptops: need R&D Analysis Facilities?

Where to find us





- > Contacts
 - swan-talk@cern.ch
 - <u>http://cern.ch/swan</u>
- > Repository
 - https://github.com/swan-cern/
- > Science Box: SWAN distribution **deployable on premises**
 - <u>https://cern.ch/sciencebox</u>
 - Includes all SWAN components: CERNBox/EOS, CVMFS, JupyterHub
 - Deployable through Kubernetes or docker-compose









> We successfully integrated CERN computing services with Jupyter

- Increase value of existing services' portfolio (storage, batch, auth)
- Successfully boosting productivity of scientists
- Jupyter became a new entry-point to our computing resources and data

> Successful model for teaching activities

- Tutorials for software, data analysis, ML, Spark...
- European project focused on higher education
- > Blending desktop and web application can help reach more users
 - Can Jupyterlab help?

