Data Mining as a Service DMaaS

P. Mato, <u>D. Piparo</u>, E. Tejedor – EP-SFT M. Lamanna, L. Mascetti, J. Moscicki – IT-ST

Cloud Services for Synchronisation and Services (CS3)

18/01/2016





- Describe a distributed service offering a web interface for data analysis based on Jupyter Notebooks
- Demonstrate how provision of CPU and storage resources as well as software are its building blocks
 - Focus on sync'd and mass storage
- Illustrate with a demo analysis how it can boost productivity and give access to innovative workflows
- Give <u>you</u> the possibility to try it out!

Prelude: The "Notebook"



A web-based interactive computing interface and platform that combines code, equations, text and visualisations.



http://www.jupyter.org

In a nutshell: an "interactive shell opened within the Also called:

"Jupyter Notebook" or "IPython Notebook"

Data Mining As a Service



Kernels are processes that run interactive code in a particular programming language and return output to the user. Kernels also respond to tab completion and introspection requests.

















… And capture their output

Shell Commands

```
In [1]: def thisFunction():
           return 42
In [2]: thisFunction()
Out[2]: 42
In [3]: %%bash
       curl rootaasdemo.web.cern.ch/rootaasdemo/SaasFee.jpg \
       > SF.jpg
                   % Received % Xferd Average Speed
                                                   Time
         % Total
         Time
                 Time Current
                                     Dload Upload
                                                   Total
         Spent Left Speed
       100 128k 100 128k 0
                                  0 2731k
                                               0 --:--:--
       --:-- 2787k
In [4]: from IPython.display import Image
       Image(filename="./SF.jpg",width=225)
```





A Distributed Service Building on top of CERN Services Portfolio



Data Mining as a Service

- Platform independent: only with a web browser
 - Analyse data via the Notebook web interface
- Calculations, input and results "in the cloud"
- Allow easy sharing of scientific results: plots, data, code
 Storage is crucial
- Simplify teaching of data processing and programming
 Not HEP specific, not only for cutting edge fundamental research
- C++, Python and other languages or analysis "ecosystems"
 - Also interface to widely adopted scientific libraries (e.g. ROOT*)

Novel Application, Existing Components

- The DMaaS project relies on technologies provided by CERN
- Scientific libraries Notebook integration (EP-SFT)
- Software distribution (EP-SFT, IT-ST): CVMFS
 - All software potentially available
- Virtualised CPU resources in OpenStack Cloud (IT-CM)
 - Interactive and batch usage
- Synergy with document sharing and publication (IT-CDA)
- Security, e.g. CERN credentials (IT-DI-CSO)
- Storage access (IT-ST): CERNBox, EOS
 - All data potentially available









EOS

Disk-based low latency storage infrastructure for physics users. Main target: physics data analysis. Storage backend for CERNBox.

Indico

Manage complex conferences, workshops and meetings.

CVMFS

HTTP based network FS, optimized to deliver experiment software Files aggressively cached and downloaded on demand.

ROOT

Software framework for data mining, visualisation and storage. Hundreds of PB of HEP data saved in ROOT format. Try it in your browser (notebooks!):

mybinder.org/repo/cernphsft/rootbinder











Other Technologies

Jupyterhub: Server application - manages login of users and redirection to notebook

- Existing solution
- Allows encapsulation: spawn Docker container at logon

Docker

- Isolation of users
- Boot faster than Virtual Machines
- Openstack support



Both have large user bases and an active community behind

The Big Picture



CÈRN



Potential Daily Usage

- Launch jobs on the batch farm
- Access notebook running on a container in the OpenStack instance
- Inspect produced data via CERNBox/EOS from the notebook
- Create plots and output data
- Share, access plots (and output data!) on the web with CERNBox web interface
- Security guaranteed by the usual CERN standards

Added value: remote users cannot open graphical connections to CERN (latency): Problem automatically solved in the above workflow

e.g.



A Taste of the User Experience

- Time for a demo:
- Download data
- Produce a plot after a simple analysis
- Share it via CERNBox

Test Node at CERN being used





Timeline

Intermediate steps accomplished:

- I) Single node, CERNBox, no CERN credentials
- 2) Single node on Openstack, CVMFS, CERNBox, CERN authentication (just demoed) See backu

See backup for more details about these setups

TODO:

 Distributed setup on Openstack, CVMFS, CERNBox, CERN authentication

DMaaS accessible to CERN users: 2nd quarter 2016

CERN Summer Student Program, ROOT lectures: Interactive notebooks offered

- 50 participants, perfect scaling, a success!
 - https://indico.cern.ch/event/407519

Data Science @ LHC Workshop, Multivariate analysis tutorial:

http://indico.cern.ch/event/395374/

- E-Planet exchange @ UERJ, Brazil
- 30 participants, every day for a week, 3h a day

https://indico.cern.ch/event/402660/

In addition, clear signs of appreciation of Notebook technology: see backup

Data Mining As a Service











Conclusions

- We will provide a service for data analysis in the cloud via a web interface
 - Platform independent: no need to install software
 - Rely on the robust services already provided by CERN
- Sync'd storage and access to the mass storage are crucial
 - Share data, code, documentation, results
 - CERNBox + EOS An optimal solution
- New ways of approaching data mining made accessible: boost productivity thanks to sync and file sharing services
- Give you the possibility to try it out!



Try it now!

• Access from the conference site until tomorrow

dmaasdemo.web.cern.ch

From the ETH "public" network only!

- Take a look to the provided notebooks, modify them, run them
 - Produce results!
 - Access them via CERNBox (https://cernbox.cern.ch)

Ask me for your user name and password!



Sign in		Files Running Clu	Isters
Username:		Select items to perform action	s on them.
rw15u098		☐	2
Password:		۵.	Select items to percentions on them.
Sign In	-	tutorials	
		B HowTo_ROOT-N	۵
		My First Notebo	□ □ hist
Т		myOutputFile.ro	□ □ hsimple
	CERNBOX	myPlot.pdf	🗆 🗅 io
			🗆 🗅 math
Hunning Clusters		Upload New - C	🗆 🗅 roofit
- 4			Co roostats

Backup Slides







Pilot Service Single Node





Interface to Job Submission Tools

Large volume of data – complex analysis: need to use many cores

- Single node: <u>TProcPool</u>, <u>IPython Parallel</u>, etherogeneous/ multithreaded code
- 2) Many nodes: Batch/Grid jobs

<u>CERN Batch Service</u> being considered in the full picture!

Several production grade, Python based job submission tools available:

- Ganga, GridControl, Panda, ...
- See A. Richards Presentation



Opportunity: Steer job submission to WLCG or local batch resources from the notebook.



CVMFS: Software Environment

Define a custom software environment via a web form

• Same mechanism for selecting hardware (e.g. GPU, N CPU cores, SSD disk)





And Notebooks





Integration of ROOT & Jupyter Notebooks delivered

- Python flavour
 - import ROOT: all goodies activated
 - %%cpp magic
- ROOT C++ flavour
 - Kernel distributed with ROOT itself
- Goodies
 - Tab completion
 - Display of graphics
 - Syntax highlighting
 - Asynchronous output capturing

ROOT comes with a C++ II/I4 compatible interpreter based on LLVM Technology



Try it Out: Local Server

Follow some simple instructions at:

https://root.cern.ch/how/how-create-rootbook

(basically build ROOT) and...

\$ root --notebook

This command:

- I. Starts a local notebook server
- 2. Connects to it via the browser

Provides a ROOT C++ kernel and the rest of ROOTbook goodies



Documentation and Links

ROOTbooks How-Tos

https://root.cern.ch/howtos#Jupyter%20Notebooks

ROOT bindings for Jupyter

https://github.com/root-mirror/root/tree/master/bindings/pyroot/ JupyROOT

ROOT C++ Kernel

https://github.com/ipython/ipython/wiki/IPython-kernels-for-other-

<u>languages</u>

Examples (15 already) from the *new* ROOT Tutorials can be found at:

https://root.cern.ch/code-examples#notebooks

both in Python and C++ (and mixed!)



ROOT: Binder Usage

Binder is a software package and a webservice (100% free and open source) to turn a GitHub repo into a collection of interactive notebooks powered by Jupyter and <u>Kubernetes</u>. <u>ROOT is on Binder</u>: you can try it at <u>ROOTBinder</u>.

On <u>ROOTBinder</u> you can find a collection of Notebooks aiming to illustrate the potential of the ROOT Framework.

> View, Create and Run ROOTbooks!



Anonymous access, no persistent storage

Binder Usage





hcontz.Draw("CONTZ"); c1.Draw();



Option CONTZ example

mybinder.org/repo/cernphsft/rootbinder



A C++ ROOTbook



Data Mining As a Service