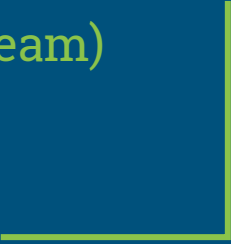# Containers@LHCb

ISGC Containers Workshop
Ben Couturier (for the LHCb Computing team)

# LHCb Collaboration

~1200 members,
69 Institutes,
16 Countries

# Containers in LHCb

**Appealing way to package/run experimental software**

- Need to build/run on SLC5/SLC6/Centos7
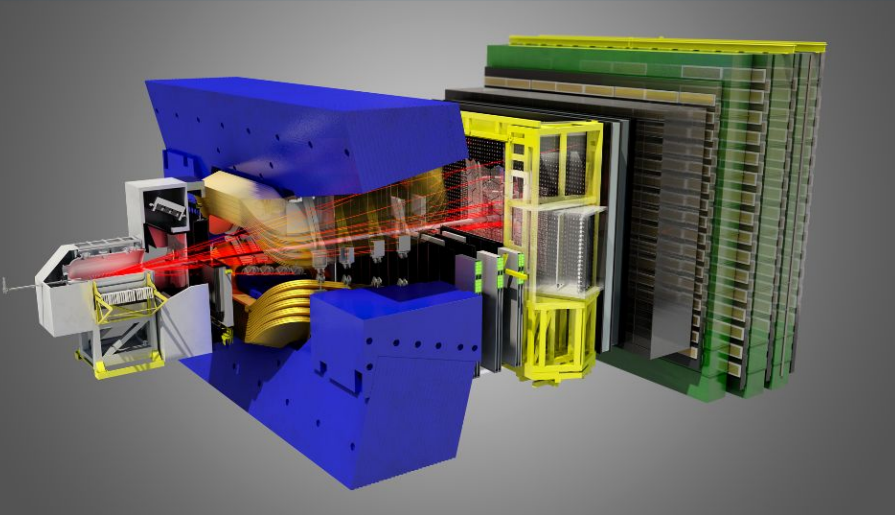- S.Binet prepared images in since 2014

*Boundary conditions:*

**LHCb application stack is rather large**

- 12GB for the application software
- slc6-build image: ~1.2GB

**Integration with CVMFS is therefore crucial**

- Either inside the docker image (not practical)
- Or on the system, shared between all images

# Multiple teams interested

**The computing team saw potential early on**

- **in line with microservices and builds**
- **Using production software**

**Analysts are also very keen, for different reasons**

- **Easy way to gather/configure all the tools needed for analysis**
- **Really helps reproducibility, in conjunction with gitlab/gitlabCI**

# Containers
for developers

# Containers for continuous build environment

**Container with correct OS started on demand by the Jenkins jobs**

- Great way to decouple build VMs from the builds
- Extremely useful for old stacks…

**But not without consequences**

- Restarting docker daemon can have disastrous effects
- Ditto for Jenkins agent

# Containers as development environment

**Practical way to allow development on developers' laptops**

- Used in the LHCb upgrade hackathons
- Provided CVMFS is already mounted
- Tricks needed for graphical applications…

# Containers as development environment



**Unfortunately containers cannot be used on shared clusters**

- Is Openshift a solution ?

**Not convincing on MacOS yet**

- I/O performance issues
- CernVM probably more convincing still

# Rerunning productions

**Using the LHCb Software Preservation DB**

- Can rerun any LHCb data processing
- Uses the Preservation DB to choose software version
- Runs on the appropriate image, loading applications from CVMFS

# Reproducible
## Analysis

# Reproducible Research@LHCb

**Different issue than production**

- Analysis software not necessarily on CVMFS
- Wider range of tools and methods
- Does not run on same data volumes

**Containers are a key enabler of "Continuous Analysis"**

- Used in coordination with version control systems (git/gitlab in our case) and continuous integration



Event 251784647
Run 125013
Thu, 09 Aug 2012 05:53:58

$\Lambda_b$
decay point

$\Lambda_b$

pp
collision point

μ
K
p

μ

# Huge interest in the research community

**Have tried to follow various initiatives...**

*Obviously NOT an exhaustive list!*

- Biomed researchers are very interested
- Way to run custom code on HPC clusters

# Data Analysis tools following suit

**Containers are the enabling technology for complex environments in pipelines:**

- e.g. www.pachyderm.io

**But also for interactive analysis:**

- https://github.com/everware
- https://cern.ch/swan
- [...]

*How do we run containerized analysis on our compute resources ?*

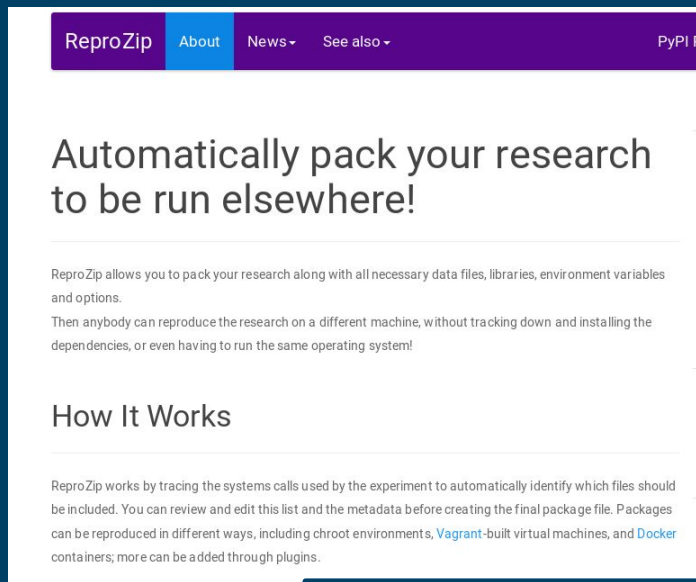# Containers
# v.s.
# VMs

**Maybe we don't have to choose**

- Some tools can use both e.g. ReproZip
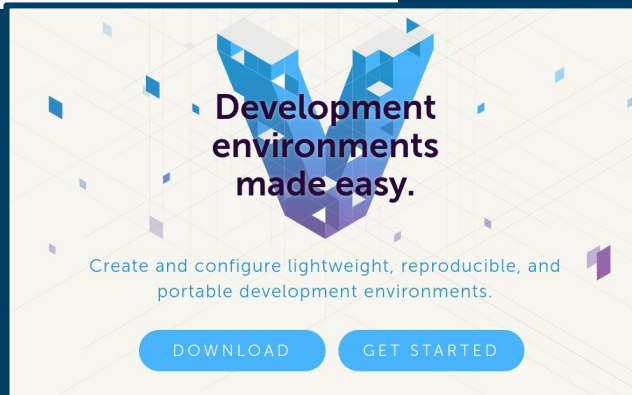
**They can be used in conjonction anyway**
- Openstack to start VMs running dockers images...

**Customizing VM images is not necessarily so complicated**

- Vagrant can be used with VMs

# Long term preservation of containers ?

**How long do we need to run them for ?**

- In the case of LHCb, we have to re-run old trigger versions for the duration of the experiment

**Are containers more/less preservable than VMs?**

- **Would standards help ?**
  https://www.opencontainers.org/



**OPEN CONTAINER INITIATIVE**

AN OPEN GOVERNANCE STRUCTURE FOR THE EXPRESS PURPOSE OF CREATING OPEN INDUSTRY STANDARDS AROUND CONTAINER FORMATS AND RUNTIME

# Exciting times!

- Very active field, loads of exciting developments, many questions open
- How to make sure our users can make full use of containers ?
- How to articulate VMs and containers ?
- How to organize all the images used ?
- What about security aspects ?

We are not alone, need for collaboration with other fields on the topic