



Long term data preservation demonstrator using CernVM

Jakob Blomer for the CernVM Team

ALICE Offline Week
March 2014



- ① Components and Functionality
- ② ALEPH software on CernVM / SL4
- ③ CMS Open Data Pilot on CernVM / SL5



Base Technology: Virtual Machines

Isolation; spawn (historic) software environment on any physical host.



Base Technology: Virtual Machines

Isolation; spawn (historic) software environment on any physical host.

Add-On 1: CernVM File System

CernVM-FS is a *versioning* and *snapshotting* file system used to make the virtual machine's content accountable.

It is used to distribute

- Operating System (Scientific Linux 4–6)
- Experiment software (/cvmfs/alice.cern.ch/...)
- Possibly: conditions data (/cvmfs/alice-ocdb.cern.ch/...)



Base Technology: Virtual Machines

Isolation; spawn (historic) software environment on any physical host.

Add-On 1: CernVM File System

CernVM-FS is a *versioning* and *snapshotting* file system used to make the virtual machine's content accountable.

It is used to distribute

- Operating System (Scientific Linux 4–6)
- Experiment software (`/cvmfs/alice.cern.ch/...`)
- Possibly: conditions data (`/cvmfs/alice-ocdb.cern.ch/...`)

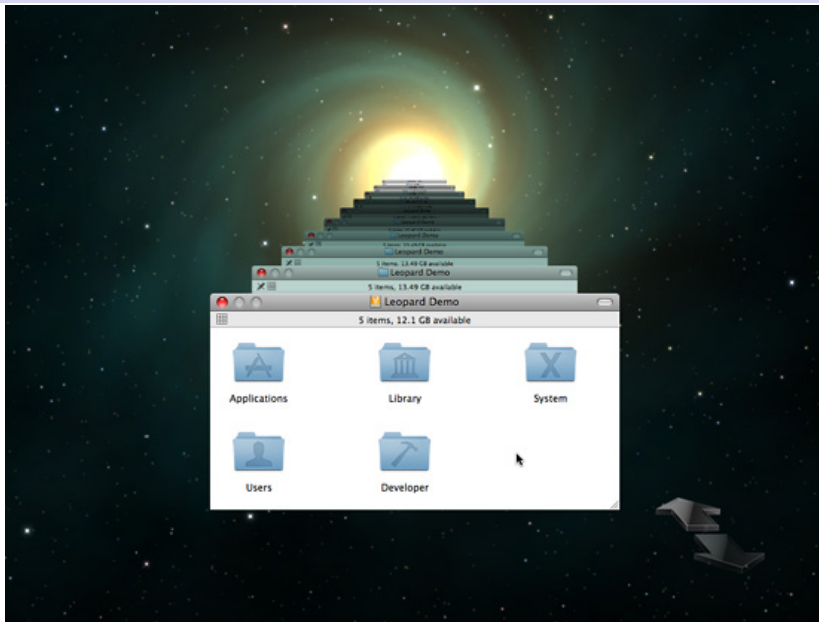
Add-On 2: CernVM Contextualization Agent

Interprets a textual specification for customizing CernVMs.

A generic operating system installation can convert into many roles.



A Time Machine for the Analysis Environment





A Time Machine for the Analysis Environment

Example

- For LHCb software stored in CernVM-FS, we can go back to essentially every day until October 2010
- This capability becomes more powerful since we can **associate meaningful tags with snapshots**

Tag list for the CernVM 3 operating system repository:

```
cernvm@cernvm002:~$ sudo cvmfs_server lstags cernvm-prod.cern.ch
```

NAME	HASH	SIZE	REVISION	TIMESTAMP	CHANNEL	DESCRIPTION
cernvm-system-3.1.0.0	fb17e39ca21729a9509fe836fc7f30d26cae1c82	14kB	11	28 Jan 2014 14:31:17	0	
cernvm-system-3.1.1.0	d855c3c05e4fcdb9d5c6f1d0b08c74094f4f5008	14kB	13	30 Jan 2014 00:11:10	0	
cernvm-system-3.1.1.1	3a06202aad3b3163b9c5bd36f48b25744f3f204	14kB	16	5 Feb 2014 21:03:00	0	
cernvm-system-3.1.1.2	fc2faf3bc87a2f74da7db22525189b5c582975de	14kB	18	16 Feb 2014 13:01:32	0	
cernvm-system-3.1.1.3	fc0d2515c9e79f9fd3cf8b01eac0a16746f4f6cb	14kB	20	4 Mar 2014 09:26:27	0	
cernvm-system-3.1.1.4	314d93015ce473d9a6c99a7365dd4ce38b4e7b13	14kB	22	17 Mar 2014 11:07:02	0	
HEAD	314d93015ce473d9a6c99a7365dd4ce38b4e7b13	14kB	22	17 Mar 2014 11:07:10	0	



Use Cases for a Preserved Software Environment

- ① Processing of legacy data
 - Software implicitly encodes knowledge about the correct interpretation of the data
 - **After substantial upgrades** and modifications of the detector, the new software might lose this legacy knowledge
 - **After experiment decommission**, porting and validation of software is likely to end
 - Porting and validation will at some point become prohibitively expensive or just impossible



Use Cases for a Preserved Software Environment

- ① Processing of legacy data
 - Software implicitly encodes knowledge about the correct interpretation of the data
 - **After substantial upgrades** and modifications of the detector, the new software might lose this legacy knowledge
 - **After experiment decommission**, porting and validation of software is likely to end
 - Porting and validation will at some point become prohibitively expensive or just impossible
- ② Validation of new software versions (see talk by S. Roiser)
 - Comparison with historic version provides input for validation

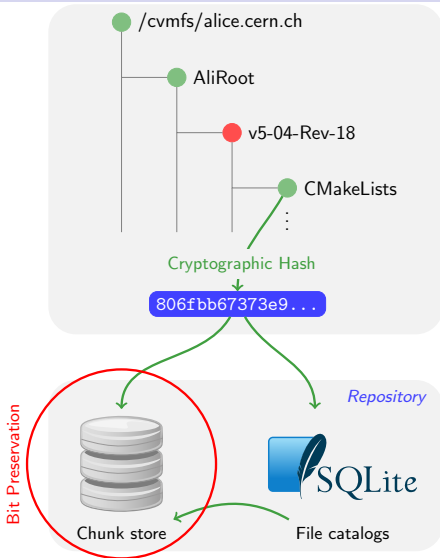


Use Cases for a Preserved Software Environment

- ① Processing of legacy data
 - Software implicitly encodes knowledge about the correct interpretation of the data
 - **After substantial upgrades** and modifications of the detector, the new software might lose this legacy knowledge
 - **After experiment decommission**, porting and validation of software is likely to end
 - Porting and validation will at some point become prohibitively expensive or just impossible
- ② Validation of new software versions (see talk by S. Roiser)
 - Comparison with historic version provides input for validation
- ③ Stable environment for education (cf. CMS Open Data Pilot)
 - Stable operating system and experiment software version accompanies “open data” set and well-defined analysis tasks
 - **Driver for data preservation:**
 - Opportunity to streamline data format and documentation
 - Disentangle from grid environment



Versioning and Snapshots in CernVM-FS



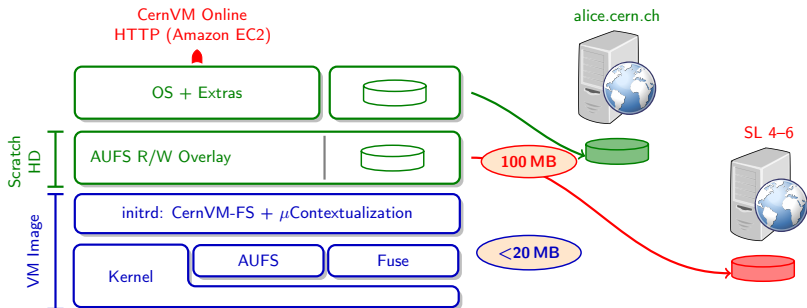
Data Store

- Eliminates duplicates
- Archiving

File Catalog

- Directory structure, symlinks
- Content hashes of regular files
- Digitally signed
- Plain files

Reduces preservation of software environment to bit preservation



Twofold system: μ CernVM boot loader + OS delivered by CernVM-FS

- μ CernVM has a modern Linux kernel, support for all relevant hypervisors and clouds
- The very same image can be *contextualized* to run Scientific Linux 4 32bit as well as the latest Scientific Linux 6 64bit
- ≈ 10 years with a single image



user-data.txt

```
[cernvm]
organisations=ALICE
repositories=alice,alice-ocdb,sft
shell=/bin/bash
config_url=http://cernvm.cern.ch/config
users=alice:alice:ion
edition=Desktop
keyboard=us
startXDM=on
auto_login=on

[ucernvm-begin]
cvmfs_tag=cernvm-system-3.1.1.4
[ucernvm-end]
```

Boot on CERN OpenStack

```
nova boot AliceVM -image "cvm3" -flavor m1.small \
  -key-name ssh-key -user-data user-data.txt
```



Exercise: resurrecting the ALEPH environment

- Can we use CernVM on current CERN OpenStack infrastructure to do ALEPH physics?
- Backport of CernVM-FS to Scientific Linux 4
- Template installation of Scientific Linux 4 for use with μ CernVM



Exercise: resurrecting the ALEPH environment

- Can we use CernVM on current CERN OpenStack infrastructure to do ALEPH physics?
- Backport of CernVM-FS to Scientific Linux 4
- Template installation of Scientific Linux 4 for use with μ CernVM

Instances

<input type="checkbox"/>	Instance Name	Image Name	IP Address	Size	Keypair	Status	Task	Power State	Uptime	Actions
<input type="checkbox"/>	cernvm-aleph01	ucernvm-slc4	188.184.134.26	m1.small 2GB RAM 1 VCPU 20.0GB Disk	-	Active	None	Running	3 months, 2 weeks	<input type="button" value="Create Snapshot"/> <input type="button" value="More"/>



```
pb-d-128-141-134-74:~ jakob$ ssh -X aleph@cernvm-aleph01
aleph@cernvm-aleph01's password:
[aleph@cernvm-aleph01 ~]$ source setaleph.sh
[aleph@cernvm-aleph01 ~]$ cd test/ALPHA/
[aleph@cernvm-aleph01 ALPHA]$ sh alpha.sh
*****
*****          ALPHA RUN          **** 11.6 ****
*****
*****

Wed Mar 19 16:10:27 CET 2014

*****
***   Compilation and creation of the makefile 6lep.mk
*****
gmake -f /home/aleph/test/ALPHA/6lep.mk
gmake: `6lep' is up to date.
```




Purpose: Provide an easy-to-use virtual machine with CMS computing environment for CMS Open Data

Data:

- Frozen data set
- Remote data access
Initially through XrootD, eventually DPHEP portal

Software:

- Frozen CMS software framework (CMSSW.4.2.8.patch7)
- *Complete* analysis environment required (compile + run)
- Requires Scientific Linux 5 compatible virtual machine

Virtual machine, user interface:

- Graphical environment
- Easy-to-install and easy-to-use



Deployment: as OVF/OVA bundle¹

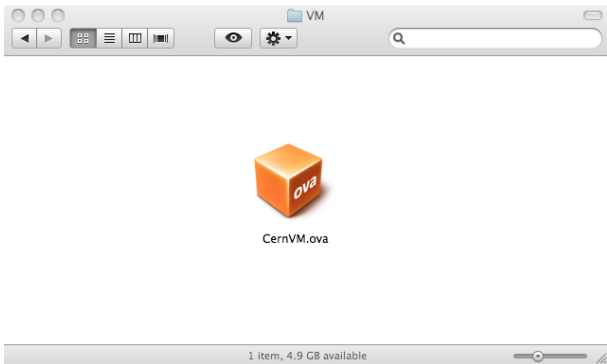
- Open specification for bundling VMs, stable since 2009
- OVA: tarball containing hard disk image and an XML specification

¹Open Virtualization Format / Open Virtual Appliance, <http://www.dmtf.org/standards/ovf>



Deployment: as OVF/OVA bundle¹

- Open specification for bundling VMs, stable since 2009
- OVA: tarball containing hard disk image and an XML specification

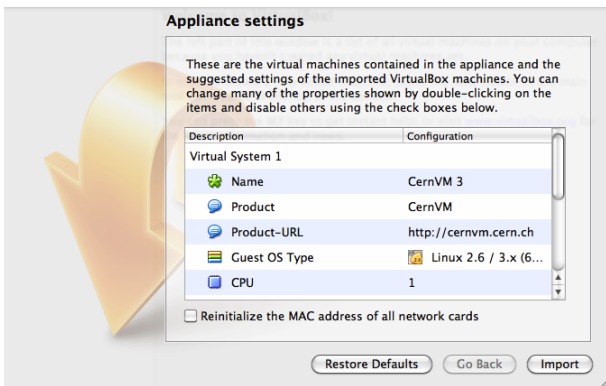


¹Open Virtualization Format / Open Virtual Appliance, <http://www.dmtf.org/standards/ovf>



Deployment: as OVF/OVA bundle¹

- Open specification for bundling VMs, stable since 2009
- OVA: tarball containing hard disk image and an XML specification

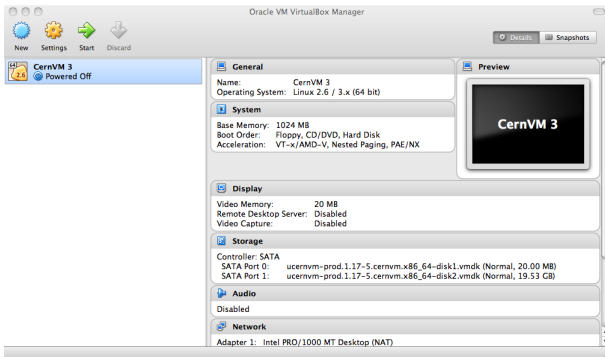


¹Open Virtualization Format / Open Virtual Appliance, <http://www.dmtf.org/standards/ovf>



Deployment: as OVF/OVA bundle¹

- Open specification for bundling VMs, stable since 2009
- OVA: tarball containing hard disk image and an XML specification

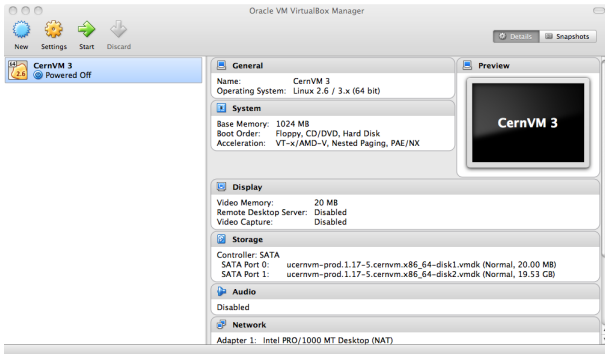


¹Open Virtualization Format / Open Virtual Appliance, <http://www.dmtf.org/standards/ovf>



Deployment: as OVF/OVA bundle¹

- Open specification for bundling VMs, stable since 2009
- OVA: tarball containing hard disk image and an XML specification



- OVA packaging fed back into CernVM baseline

¹Open Virtualization Format / Open Virtual Appliance, <http://www.dmtf.org/standards/ovf>



CMS Open Data Pilot on CernVM / SL5

Work in progress

The screenshot displays the CMS Open Data Pilot software interface. The main window title is "CMS Open Data Pilot [Running]" and the current file is "cmsShow: DoubleElectron.root [1/1], event [1/10]". The interface includes a menu bar (File, Edit, View, Window, Help), a toolbar with navigation buttons, and a status bar showing "Run 170348", "Lumi 57", "Event 30771350", and "Sun Jul 17 19:36:51 2011 CEST".

The main display area is divided into several panels:

- Summary View:** A sidebar on the left with a tree view of event objects including ECAL, HCAL, Jets, Tracks, Muons, Electrons, Vertices, BeamSpot, DT-segments, CSC-segments, Photons, and MET.
- Rho Phi:** A large central plot showing the distribution of particles in the Rho-Phi plane. It features a central cluster of particles with several tracks extending outwards, highlighted in red and blue.
- Rho Z:** A plot showing the distribution of particles in the Rho-Z plane, displaying a central cluster of particles.
- 3D Tower:** A 3D visualization of the particle distribution, showing a cylindrical structure with tracks extending from the center.

The interface also includes a "REWORKS" logo and a "CMS" logo in the top right corner. The bottom status bar shows the date "Wednesday 19 March 2014" and the system tray with various icons and the time "17:25".

Backup Slides

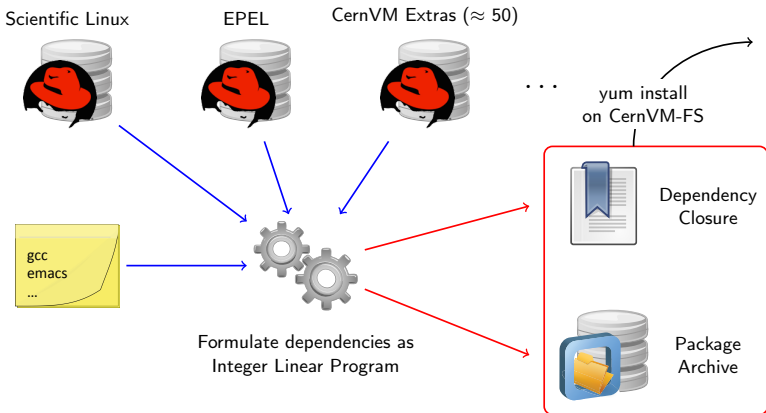


Build Process: Scientific Linux on CernVM-FS

Maintenance of the repository **should not** become a Linux distributor's job

But: should be reproducible and well-documented

Idea: automatically generate a **fully versioned, closed** package list from a "shopping list" of unversioned packages





Build Process: Package Dependency ILP

Normalized (Integer) Linear Program:

$$\text{Minimize } (c_1 \cdots c_n) \cdot \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \quad \text{subject to} \quad \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \cdots & a_{mn} \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \leq \begin{pmatrix} b_1 \\ \vdots \\ b_m \end{pmatrix}$$

Here: every available (package, version) is mapped to a $x_i \in \{0, 1\}$.

Cost vector: newer versions are cheaper than older versions.

(Obviously: less packages cheaper than more packages.)

Dependencies:

Package x_a requires x_b or x_c : $x_b + x_c - x_a \geq 0$.

Packages x_a and x_b conflict: $x_a + x_b \leq 1$.

(...)

Figures

$\approx 17\,000$ available packages ($n = 17000$), 500 packages on “shopping list”

$\approx 160\,000$ inequalities ($m = 160000$), solving time < 10 s (g1pk)

Meta RPM: $\approx 1\,000$ fully versioned packages, dependency closure

Idea: Mancinelli, Boender, di Cosmo, Vouillon, Durak (2006)



Hypervisor / Cloud Controller	Status
VirtualBox	✓
VMware	✓
KVM	✓
Xen	✓
Microsoft HyperV	✓
Parallels	⚡ ³
OpenStack	✓
OpenNebula	✓ ²
Amazon EC2	✓ ¹
Google Compute Engine	✓

¹ Only tested with ephemeral storage, not with EBS backed instances

² Only amiconfig contextualization

³ Unclear license of the guest additions