# CernVM – A Versatile Environment for High-Energy Physics Applications in the Cloud

Jakob Blomer

SLAC Computing Seminar
19th September 2014

# Virtualized and Cloud Resources

## IaaS Clouds

- Complement batch scheduler
- CERN OpenStack: 22 K cores

## Volunteer Cloud

- 75 % computing resources and > 1 trillion simulated events for CERN theory group
- Can be also seen as a "corridor grid"

## High-level Trigger Farms

- ATLAS: 28 K cores
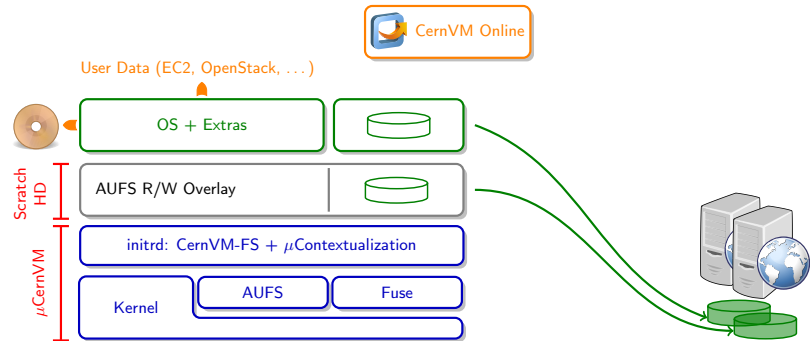- CMS: 12 K cores

## Long-term data preservation

- Preserve data analysis environment
- Outreach & education
- Becomes relevant for LHC experiments

Using these resources requires virtual machine image

## IaaS Clouds

- Complement batch scheduler
- CERN OpenStack: 22 K cores

## Volunteer Cloud

- 75 % computing resources and > 1 trillion simulated events for CERN theory group
- Can be also seen as a "corridor grid"

## High-level Trigger Farms

- ATLAS: 28 K cores
- CMS: 12 K cores

## Long-term data preservation

- Preserve data analysis environment
- Outreach & education
- Becomes relevant for LHC experiments
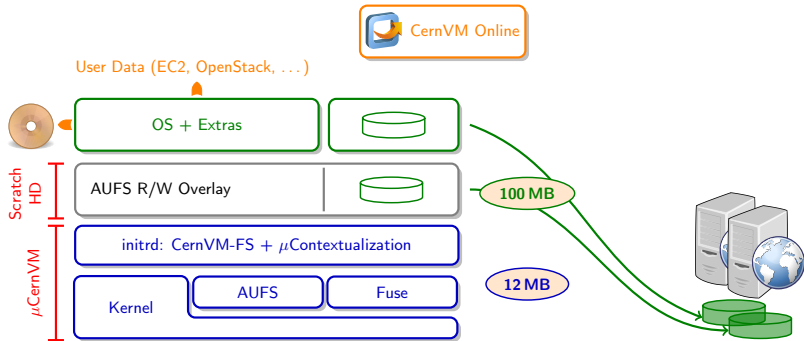
Using these resources requires virtual machine image

Twofold system: $\mu$CernVM boot loader + OS delivered by CernVM-FS

Twofold system: $\mu$CernVM boot loader + OS delivered by CernVM-FS

$\Rightarrow$ Drastic reduction in size

From "just enough operating system" to "operating system on demand"
400 MB image (compressed) $\mapsto$ 12 MB image + 100 MB cache

CernVM Kernel: 3.10 long-term kernel (2 year support)
Features: KSM, zRam, THP, cgroups, X32-ABI
Extra modules: AUFS, VMware drivers, VBox drivers, OpenAFS

**"Virtualization-friendly"**, minimal set of device drivers:
  100 modules / 8 MB as compared to 2 000 modules / 120 MB in SL6

---

1. Execute SYSLINUX boot loader
2. Decompress and load Linux kernel
3. Decompress init ramdisk, execute customized /init
   a) Start networking
   b) Contextualize (supports EC2, OpenStack, OpenNebula, vSphere)
   c) [Partition and] [format and] mount scratch space
   d) Mount CernVM-FS
   e) Mount AUFS root file system stack
   f) Change root file system and start operating system

| Hypervisor / Cloud Controller | Status |
|---|---|
| VirtualBox | ✓ |
| VMware | ✓ |
| KVM | ✓ |
| Xen | ✓ |
| Microsoft Hyper-V | ✓ |
| Parallels | ⚡[3] |
| Openstack | ✓ |
| OpenNebula | ✓ |
| Amazon EC2 | ✓[1] |
| Google Compute Engine | ✓[2] |
| Microsoft Azure | ? |
| Docker | ? |

[1] Only tested with ephemeral storage, not with EBS backed instances
[2] Only amiconfig contextualization
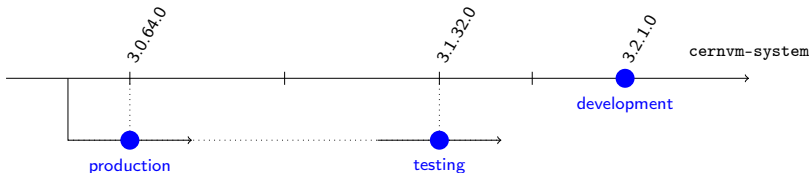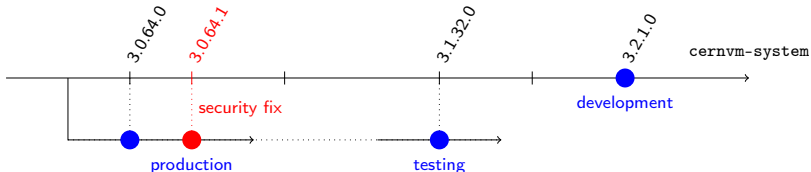[3] Unclear license of the guest additions

General idea: Install packages with yum into a CernVM-FS chroot jail
Problem: Typical package repositories are not designed
to *preserve* an environment

The CernVM 3 build process **ensures strong versioning** on three levels

1. `cernvm-system` meta RPM
   fully versioned dependency closure

2. Named branches in the CernVM-FS repository

3. Versioned snapshots provided by CernVM-FS
   allow the very same image to instantiate any `cernvm-system` version
   helpful for **long-term data preservation**

General idea: Install packages with yum into a CernVM-FS chroot jail
Problem: Typical package repositories are not designed
to *preserve* an environment

The CernVM 3 build process **ensures strong versioning** on three levels

1. `cernvm-system` meta RPM
   fully versioned dependency closure

2. Named branches in the CernVM-FS repository

3. Versioned snapshots provided by CernVM-FS
   allow the very same image to instantiate any `cernvm-system` version
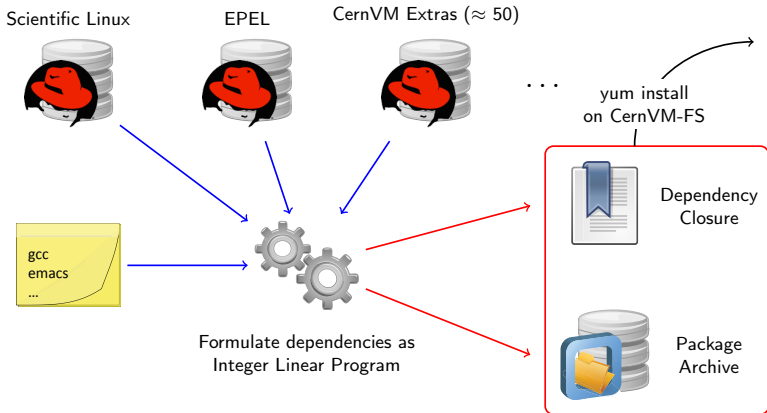   helpful for **long-term data preservation**

# 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

Scientific Linux          EPEL          CernVM Extras ($\approx 50$)

· · ·          yum install
on CernVM-FS

gcc
emacs
...

Dependency
Closure

Formulate dependencies as
Integer Linear Program

Package
Archive

Options for updating: **stay**, **diverge**, **rebase**

Rebase high-level perspective:

**1** On first boot, CernVM selects and pins newest available version

**2** Automatic update notifications

**3** Applying updates requires a reboot
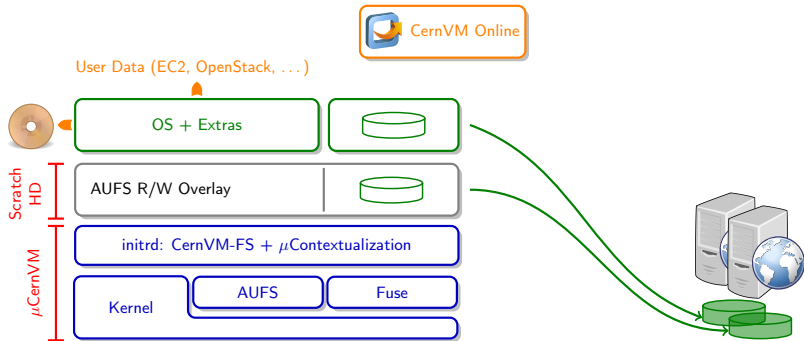Most security critical updates require a reboot anyway

## $\mu$CernVM Bootloader

- boot partition read-only, updates dropped on ephemeral storage

- 2 phase boot: start old kernel and ramdisk, **kexec** into updated version

## CernVM-FS OS Repository

- Mount updated CernVM-FS snapshot

- Conflict resolution wrt. local changes

  **1** keep local configuration
  **2** map user/group ids
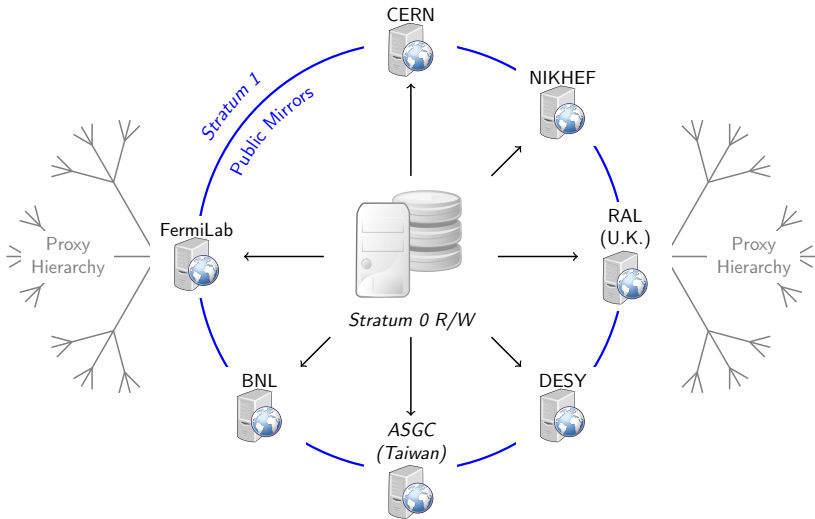  **3** merge rpm database

Options for updating: **stay**, **diverge**, **rebase**

Rebase high-level perspective:

1. On first boot, CernVM selects and pins newest available version

2. Automatic update notifications

3. Applying updates requires a reboot
   Most security critical updates require a reboot anyway

## $\mu$CernVM Bootloader

- boot partition read-only,
  updates dropped
  on ephemeral storage

- 2 phase boot:
  start old kernel and ramdisk,
  **kexec** into updated version

## CernVM-FS OS Repository

- Mount updated CernVM-FS
  snapshot

- Conflict resolution wrt.
  local changes

  1. keep local configuration
  2. map user/group ids
  3. merge rpm database

- Independent *repositories*, e. g. /cvmfs/atlas.cern.ch
- Single point of publishing
- HTTP Transport, access and caching on demand
- Has a life outside the CernVM virtual machine

**Server side: stateless services**

**Server side: stateless services**

**Server side: stateless services**

**Server side: stateless services**

**Server side: stateless services**

**Server side: stateless services**
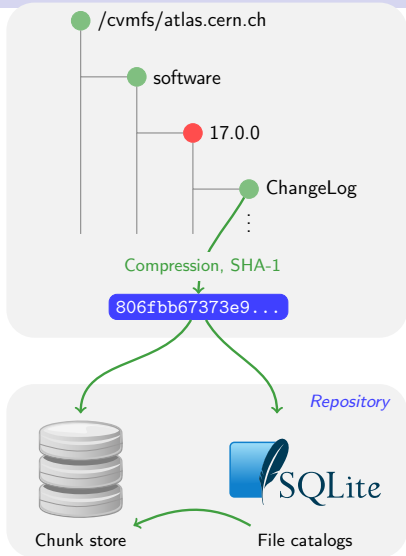
**Server side: stateless services**

- Kernel-level Union File System AUFS

- $< 5\%$ performance loss (untar)

---

- Fully POSIX-compliant read-write file system

- Encapsulated change set in scratch area

- In contrast to file-wise write: batch publishing of consistent snapshots

/cvmfs/atlas.cern.ch

software

17.0.0

ChangeLog

Compression, SHA-1

`806fbb67373e9...`

*Repository*

Chunk store          File catalogs

## Data Store

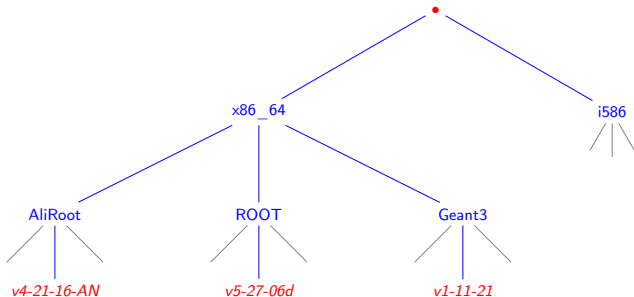- Compressed chunks (files)
- Eliminates duplicates

## File Catalog

- Directory structure, symlinks
- Content hashes of regular files
- Digitally signed
  ⇒integrity, authenticity
- Time to live
- Partitioned / Merkle hashes
  (user assisted)

⇒ Immutable files, trivial to check for corruption, consistency by snapshots
≈ 6× reduction in number of files / volume

## Automatic Approaches

File based (Tolia et al. 2004), directory based (Kutzner 2008),
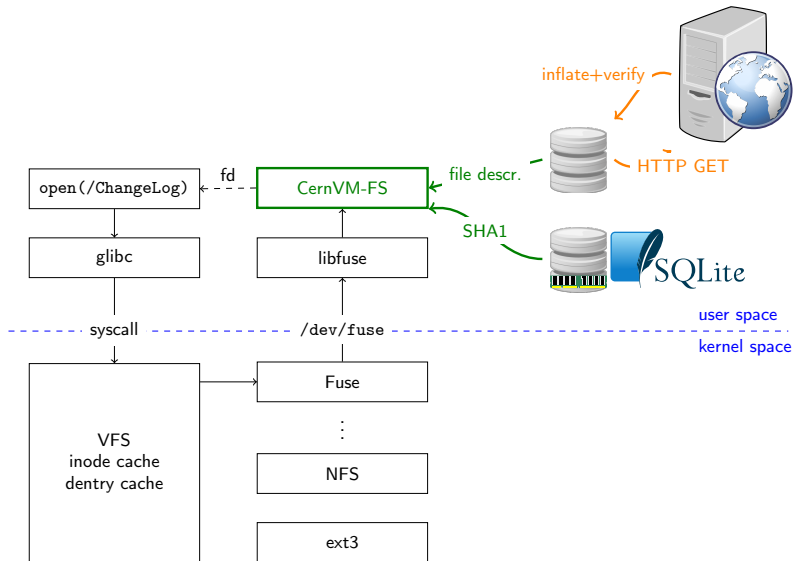global (Compostella et al. 2010)



## CernVM-FS: Semi-automatic, use of human knowledge about

- locality by software version
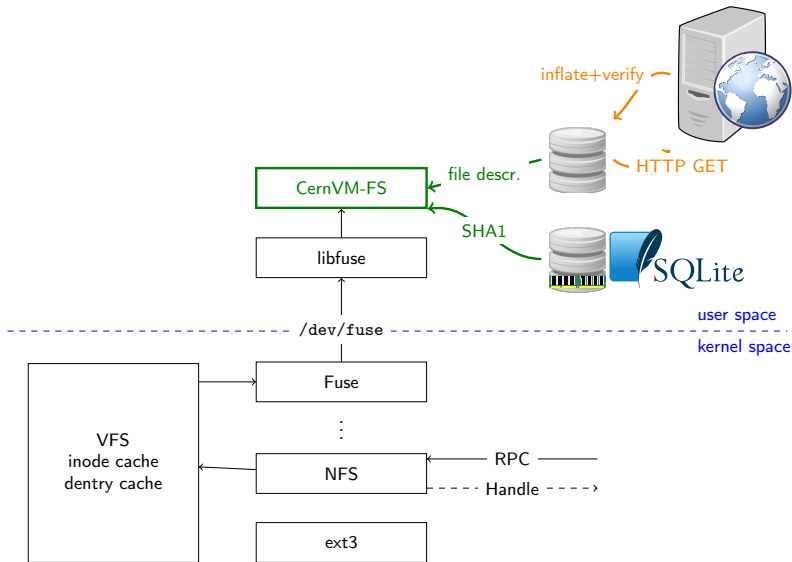- locality by frequency of changes

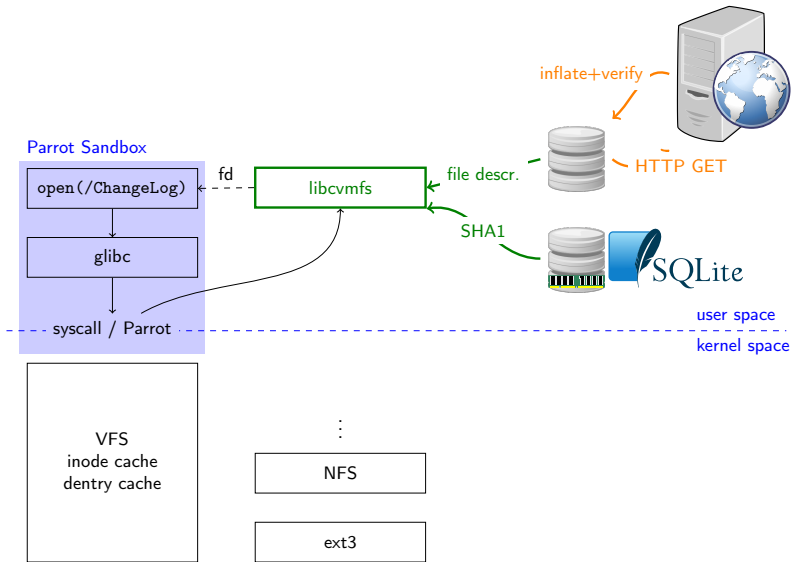http://ccl.cse.nd.edu/software/parrot

Goal Reduce cache duplication, local hard disk consumption

Environment Super computers, diskless clusters



External Connection

load

store

Memory Cache

**RAMCloud:** distributed key-value store in DRAM

Prototype based on RAMCloud

1 SQlite backend

2 Put/Get interface for CernVM-FS chunk storage

`https://ramcloud.stanford.edu`

1. As part of CernVM
   - Fuse module, auto-configured

2. As Fuse module on the worker node / workstation
   - Almost all of WLCG using it this way
   - Very little problems
   - "Private mounts" as normal user possible

3. NFS-exported Fuse module
   - Brings back infamous NFS bottleneck
   - Only solution for disk-less farms
   - DESY: 2 k nodes on CernVM-FS / NFS

4. As part of the Grid job
   - Using library interface + Parrot connector
   - Performance limitations, limited cache sharing
   - Runs as normal user

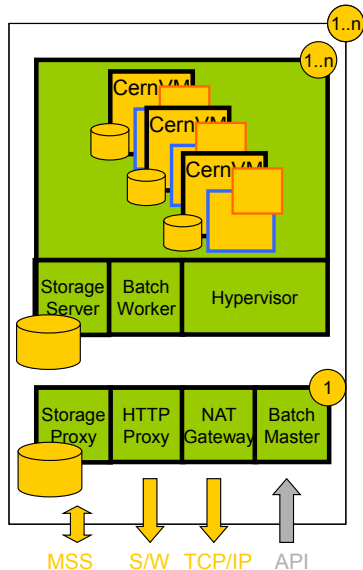Web proxies: Site-local, *Frontier* (CMS, ATLAS), Public Service

## Transfer and Conventions

- Payload in attached device (CD, floppy)

- Payload on web server (fixed IP, gateway)

- Script or "user data" string

- Interpreted by VM agent (e. g. cloud-init, amiconfig)

## Tasks

- Credentials (ssh, X.509)

- Condor head & batch services

- Squid server, XrootD proxy

- Network configuration & tuning



Source: Buncic

### 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
```

- Commonly known as
  . . . @Home projects

- Computing on *spare* resources
  of "interested citizens"

- Outreach program

- Opportunistic, volatile resources

- Big projects comparable to a
  TOP500 supercomputer

- LHC with 2 projects:
  LHC@Home & LHC@home 2.0

## SETI@Home

- Search for extraterrestrial life

- 130 000 active participants

- 630 Tera-FLOPS



## Einstein@Home

- Search for gravitational waves

- 34 000 active participants

- 470 Tera-FLOPS

Monte-Carlo simulations, parameter tuning
First BOINC project using virtual machines



## Numbers

- At any point in time
  600–700 VMs connected

- Overall: More than 1 trillion
  events created



Source: Harutyunyan

1 day geographic distribution, 2500 distinct IPs

## 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 | 3a06202aadc3b3163b9c5bd36f48b25744f3f204 | 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 |
```

**Exercise:** resurrecting the ALEPH environment

- 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 $\mu$CernVM

Instances   [Filter 🔍]  [Filter]  [➕ Launch Instance]  [Soft Reboot Instances]  [🗑 Terminate Instances]

| ☐ | Instance Name | Image Name | IP Address | Size | Keypair | Status | Task | Power State | Uptime | Actions |
|---|---|---|---|---|---|---|---|---|---|---|
| ☐ | cernvm-aleph01 | ucernvm-slc4 | 188.184.134.26 | m1.small \| 2GB RAM \| 1 VCPU \| 20.0GB Disk | - | Active | None | Running | 3 months, 2 weeks | Create Snapshot   More ▾ |

**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 bundle[1]

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

---

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

**Deployment**: as OVF/OVA bundle[1]

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



CernVM.ova

---

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

**Deployment**: as OVF/OVA bundle[1]

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



---

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

**Deployment**: as OVF/OVA bundle[1]

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



---

[1]**Open Virtualization Format / Open Virtual Appliance**, http://www.dmtf.org/standards/ovf

## CernVM

- $\mu$CernVM +
  OS template on CernVM-FS +
  Contextualization

- 20 MB image that adapts

- Image for IaaS clouds,
  volunteer clouds,
  long-term data preservation,
  development environment

- Next steps: SL7, Docker

## CernVM-FS

- Software Distribution

- HTTP transport, versioning,
  on-demand download & caching

- Limited applicability for data:
  common cluster working set
  required

- Next steps: better support for
  opportunistic resources,
  better support for small VOs

---

Source code: https://github.com/cernvm, https://github.com/cvmfs
Documentation and downloads: http://cernvm.cern.ch
Bug tracker: https://sft.its.cern.ch/jira/browse/CVM

Mailing lists: cernvm-talk@cern.ch, cvmfs-talk@cern.ch

**6** Backup Slides

CernVM-FS features targeted to loading the OS:

- Closing all read-write file descriptors in order to
  unravel file system stack on shutdown
- Redirect syslog messages
- In-flight change of DNS server
- GID / UID mappings
- $\Rightarrow$ This file system stack requires special support from a read-only Fuse
  branch since it is started before the operating system.

Normalized (Integer) Linear Program:

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

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:  $x_a + x_b \quad \leq 1$.
  $(\ldots)$

## Figures

$\approx 17\,000$ available packages ($n = 17000$),  500 packages on "shopping list"
$\approx 160\,000$ inequalities ($m = 160000$), solving time $< 10\,$s (`glpk`)
Meta RPM: $\approx 1\,000$ fully versioned packages, dependency closure

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

# 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

# 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

# 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

Based on ATLAS Figures 2012

| Software | Data |
|---|---|
| POSIX Interface | put, get, seek, streaming |
| File dependencies | Independent files |
| $10^7$ objects | $10^8$ objects |
| $10^{12}$ B volume | $10^{16}$ B volume |
| Whole files | File chunks |
| Low latency | High throughput |
| Absolute paths | Any mount point |
| Open source | Confidential |
| WORM ("write-once-read-many") Versioned | |

Based on ATLAS Figures 2012

CernVM-FS

Global
File System

Content-Addressable
Storage

Distributed
Caching

| Software | Data |
|----------|------|
| POSIX Interface | put, get, seek, streaming |
| File dependencies | Independent files |
| $10^7$ objects | $10^8$ objects |
| $10^{12}$ B volume | $10^{16}$ B volume |
| Whole files | File chunks |
| Variant symlinks | No Symlinks |
| Absolute paths | Any mountpoint |
| Open source | Confidential |
| WORM ("write-once-read-many") Versioned | |

# CernVM-FS Client in Heterogeneous Environments

In order to fully benefit from CernVM-FS, the file system has to be available on all relevant computing resources.

Range of Environments:

Scientific Linux, Fedora, Ubuntu, SuSE, OS X
1 core to 48+ cores
Tens of mounted repositories
Possibly no Fuse, disk-less server farms

**Portability:**

- Portable C++ / POSIX code

- Library interface, connector to *Parrot* (by Dan Bradley)
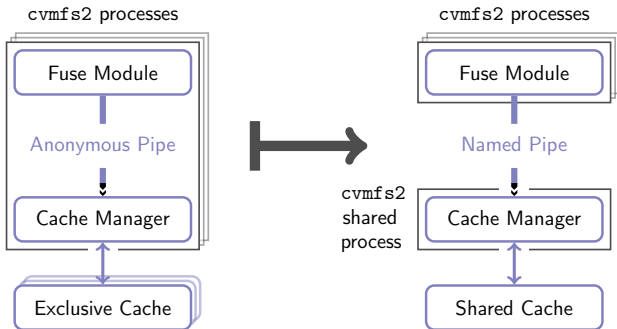
**Scalability:**

- Memory fragmentation
  open hash collision resolution $\mapsto$ linear probing
  path strings stored on the stack

- Concurrent file system access
  Fine-grained locking
  Asynchronous, parallel HTTP I/O

- Cache sharable among repositories

- Hotpatch functionality for Fuse client

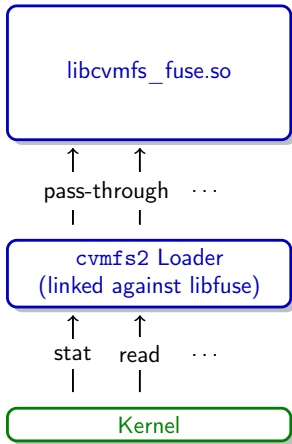Issue: Enforce shared *quota*, coordinated bookkeeping required
Idea: Turn the cache manager thread into a shared process



- No extra service: automatically spawned by first cvmfs mount point, automatically terminated by last unmount
- Named pipe can be turned into a network socket:
  Foundation for distributed shared memory cache

Addresses the issue of "worker node draining"
when there is a new version of CernVM-FS



libcvmfs_fuse.so

↑   ↑
pass-through   · · ·
|   |

cvmfs2 Loader
(linked against libfuse)

↑   ↑
stat   read   · · ·
|   |

Kernel

- A minimal loader implements the Fuse interface

- The logic is part of an (unloadable) shared library

- Very little state across file system calls: open files and open directories

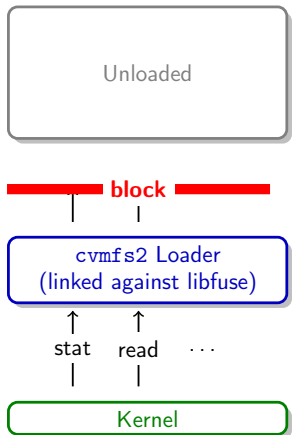- Can be also seen as a reload of parameters (like SIGHUP)

Addresses the issue of "worker node draining"
when there is a new version of CernVM-FS



- A minimal loader implements the Fuse interface

- The logic is part of an (unloadable) shared library

- Very little state across file system calls: open files and open directories

- Can be also seen as a reload of parameters (like SIGHUP)

Addresses the issue of "worker node draining"
when there is a new version of CernVM-FS



libcvmfs_fuse.so
(reloaded)

↑    ↑
pass-through  · · ·

cvmfs2 Loader
(linked against libfuse)

↑    ↑
stat  read  · · ·

Kernel

- A minimal loader implements the Fuse interface

- The logic is part of an (unloadable) shared library

- Very little state across file system calls: open files and open directories

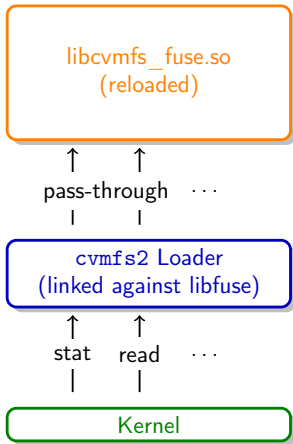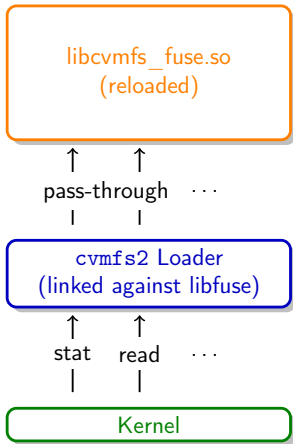- Can be also seen as a reload of parameters (like SIGHUP)

Addresses the issue of "worker node draining"
when there is a new version of CernVM-FS



libcvmfs_fuse.so
(reloaded)

↑    ↑
pass-through  · · ·
|    |

cvmfs2 Loader
(linked against libfuse)

↑    ↑
stat  read  · · ·
|    |

Kernel

- A minimal loader implements the Fuse interface

- The logic is part of an (unloadable) shared library

- Very little state across file system calls: open files and open directories

- Can be also seen as a reload of parameters (like SIGHUP)

**(Just released)**

## Working Set as seen with CernVM

- $\approx 10\%$ of all available files are requested at runtime
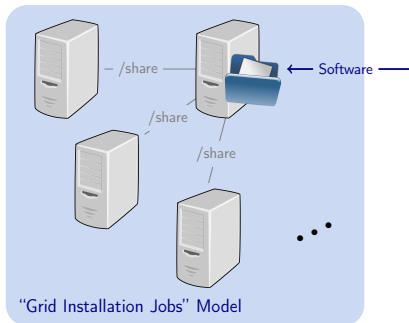- Median of file sizes: $< 4\,kB$

## Working Set as seen with CernVM

- $\approx 10\%$ of all available files are requested at runtime

- Median of file sizes: $< 4\,\mathrm{kB}$

## Flash Crowd Effect

- Up to $500\,\mathrm{kHz}$ meta data request rate

- Up to $1\,\mathrm{kHz}$ file open request rate



"Grid Installation Jobs" Model

## Working Set as seen with CernVM

- $\approx 10\%$ of all available files are requested at runtime

- Median of file sizes: $< 4\,\text{kB}$

## Flash Crowd Effect

- Up to 500 kHz meta data request rate
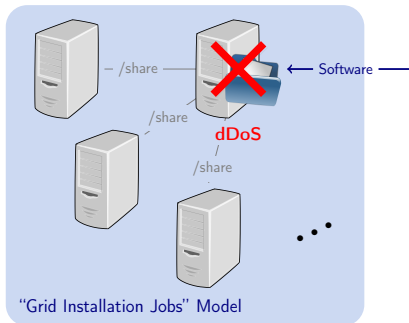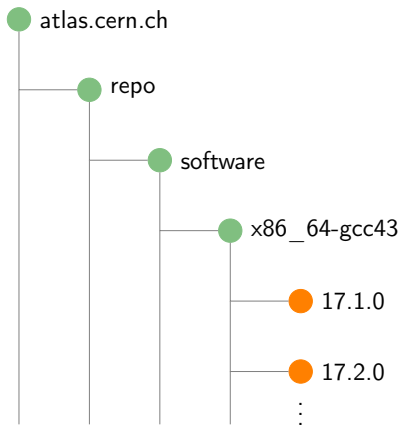
- Up to 1 kHz file open request rate



"Grid Installation Jobs" Model

Software Directory Tree

Software Directory Tree

Software Directory Tree

Between consecutive software versions: only $\approx 15\,\%$ new files

Fine-grained software structure (*Conway's law*)
Between consecutive software versions: only $\approx 15\%$ new files

Create tarball    Distribute tarball    Extract tarball

**Shared Software Area**

Create       Distribute     Extract      Create        Remount
tarball      tarball        tarball      file catalog  file system

**GROW-FS**

Release volume

**AFS** (using AFS replication in a hypothetical installation)

Create tarball               Seed tarball

**ALICE Installation Framework**

Transformation
in CAS              Replication     Wait for TTL expiry

**CernVM-FS**

New version tagged

Neue version available

Improvement: from days (Grid Installation Jobs) to hours

## Fuse Module

- Normal namespace:
  /cvmfs/<repository>
  e. g. /cvmfs/atlas.cern.ch
- Private mount as a user possible
- One process per fuse module +
  watchdog process
- Cache on local disk
- Cache LRU managed
- NFS Export Mode
- Hotpach functionality
  cvmfs_config reload

## Mount helpers

- Setup environment (number of file
  descriptors, access rights, . . . )
- Used by autofs on /cvmfs
- Used by /etc/fstab or mount as
  root
  mount -t cvmfs atlas.cern.ch
  /cvmfs/atlas.cern.ch

## Diagnostics

- Nagios check available
- cvmfs_config probe
- cvmfs_config chksetup
- cvmfs_fsck
- cvmfs_talk, connect to running
  instance

# Considerations about volunteer computing

- Resources are unmanaged, e. g. there are no grid services

- Resources are untrusted, e. g. no grid certificates can be stored
  Results are typically verified by

  - Sanity checker
  - Credit based trust relationship
  - Quorum of results

- Resources are volatile; small or resumable work units required

- Resources are heterogenous (poor and powerful, Linux, Windows, Mac)

- Resources are donated; credits and community support required

- Resources have poor I/O connection;
  works for CPU dominated workloads only

- You must never run out of work!

- Technology can be re-used for "corridor grid"