



CernVM-FS Turns 15

Jakob Blomer (CERN)

CernVM Workshop 2024

CERN, 16 September 2024

How CernVM-FS Came to Life

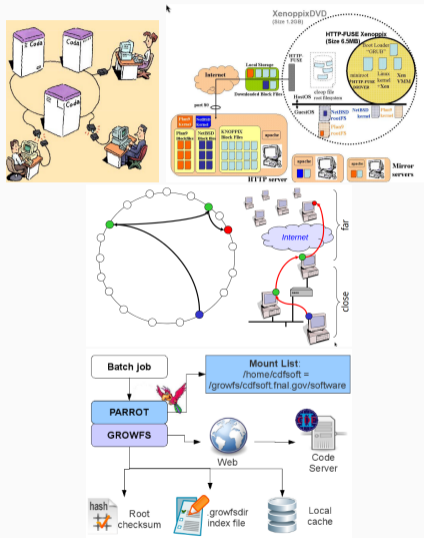
BUILDING N°: 001
Territorial Safety Officer
72504



How CernVM-FS Came to Life



- Part of the CernVM R&D project on virtualization
 - ▶ Predrag's 2018 Talk: CernVM 10 years after
- Decouple the experiment software from the virtual machine image using a **global network file system**
- Looked into several existing options
 - Coda: AFS with offline mode
 - HTTP-Fuse: on-demand bootable Linux image
 - Igor-FS: file system with P2P transport
 - GROWFS: CernVM-FS pre-cursor using the Parrot system call interception toolkit

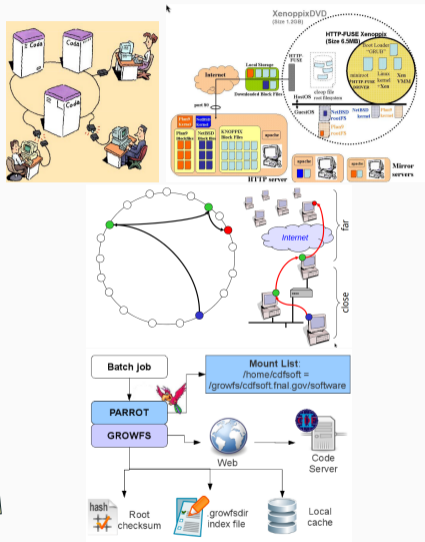


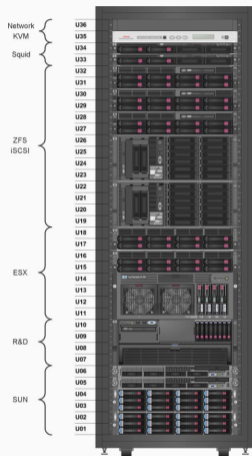
How CernVM-FS Came to Life



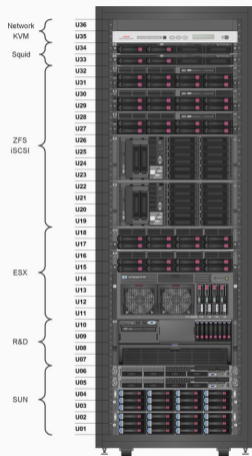
- Part of the CernVM R&D project on virtualization
 - ▶ Predrag's 2018 Talk: CernVM 10 years after
- Decouple the experiment software from the virtual machine image using a **global network file system**
- Looked into several existing options
 - Coda: AFS with offline mode
 - HTTP-Fuse: on-demand bootable Linux image
 - Igor-FS: file system with
 - GROWFS: CernVM-FS toolkit

CVMFS v1 was a GROWFS Fuse frontend written by Leandro Franco





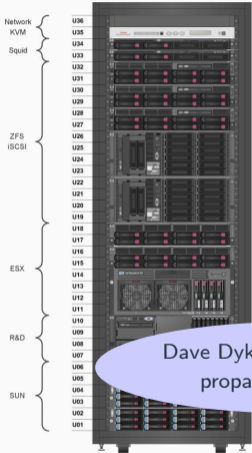
- Initial technology choices leading to CVMFS 2.0:
Fuse, C++, HTTP CDN, SQLite file catalogs, content-addressed storage
- CernVM infrastructure (including CVMFS storage, release managers etc.) operated from building 157
 - Fully virtualized with VMware ESX
 - Storage using Solaris/ZFS: our initial solution for snapshotting & replication
- Presented at CHEP 2010 in Taipei
(15 million files under management)
- Growing interest in using CernVM-FS on the Grid outside the VM
(virtualization came back later a few years later with OpenStack, Docker, k8s)
 - to address shortcomings of AFS, NFS, Grid installation jobs



- Initial technology choices (CernVM-FS, Fuse, C++, HTTP, Squid, ZFS, ISCSI, ESX, R&D, SUN) addressed storage
- CernVM infrastructure (VMware ESX, Solaris/ZFS, OpenStack, Docker, k8s, managers etc.) open source
- Fully virtualized with VMware ESX
- Storage using Solaris/ZFS: our initial solution for snapshotting & replication
- Presented at CHEP 2010 in Taipei (15 million files under management)
- Growing interest in using CernVM-FS on the Grid outside the VM (virtualization came back later a few years later with OpenStack, Docker, k8s)
 - to address shortcomings of AFS, NFS, Grid installation jobs

Meanwhile:

- Amazon EC2, S3, CloudFront
- iPhone, Android
- Higgs Discovery



- Initial technology choices (Fuse, C++, HTTP, storage)
- CernVM infrastructure (managers etc.) open source
 - Fully virtualized with VMware ESX
 - Storage using Solaris/ZFS: our initial solution for snapshotting & replication
- Presented at CHEP 2010 in Taipei (15 million files under management)
 - to address shortcomings of AFS, NFS, Grid installation jobs

Meanwhile:

- Amazon EC2, S3, CloudFront
- iPhone, Android
- Higgs Discovery

Dave Dykstra: "how will you propagate updates?"

Steve Traylen: "how about autofs for shared nodes?"



With the adoption of experiments and a global infrastructure, we faced several challenges:

1. Robust and scalable writing / publishing
2. We split development and operations (CERN IT, RAL, BNL, FNAL) and needed to find “standard” replacements for the ZFS storage → “Stratum 0” to “Stratum 1” replication
3. File system clients would suddenly run on Grid worker nodes instead VMs
 - Client had to become more robust and gentle on resources; internally, we needed to switch from libfuse’s path based interface to the low-level, inode based interface
 - No room for downtime for client updates → developed client hotpatch functionality
 - Eventually, we needed to split the Cern(VM) specific configuration from the core software, leading to the “cvmfs-config-xyz” packages and the config repository

Presented the roadmap to CVMFS 2.1 at CHEP 2012 in New York
(100 million files under management)



With the adoption of experiments and a global infrastructure, we faced several challenges:

1. Robust publishing
2. We split into sites (CERN IT, RAL, BNL, FNAL) and needed to find “standard” replacement → “Stratum 0” to “Stratum 1” replication
3. File system clients would suddenly run on Grid worker nodes instead VMs
 - Client had to become more robust and gentle on resources; internally, we needed to switch from libfuse’s path based interface to the low-level, inode based interface
 - No room for downtime for client updates → developed client hotpatch functionality
 - Eventually, we needed to split the Cern(VM) specific configuration from the core software, leading to the “cvmfs-config-xyz” packages and the config repository

Among the early, crucial supporters were Doug Benjamin and Ian Collier

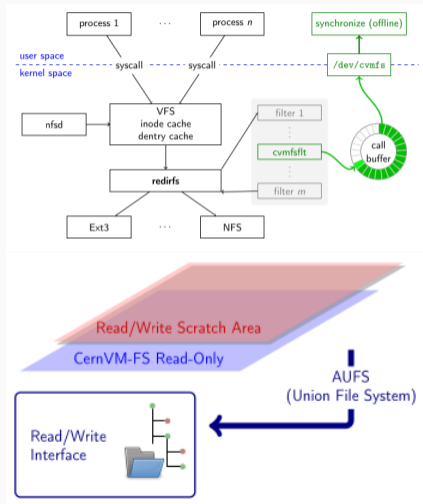
Presented the roadmap to CVMFS 2.1 at CHEP 2012 in New York
(100 million files under management)

How to solve writing



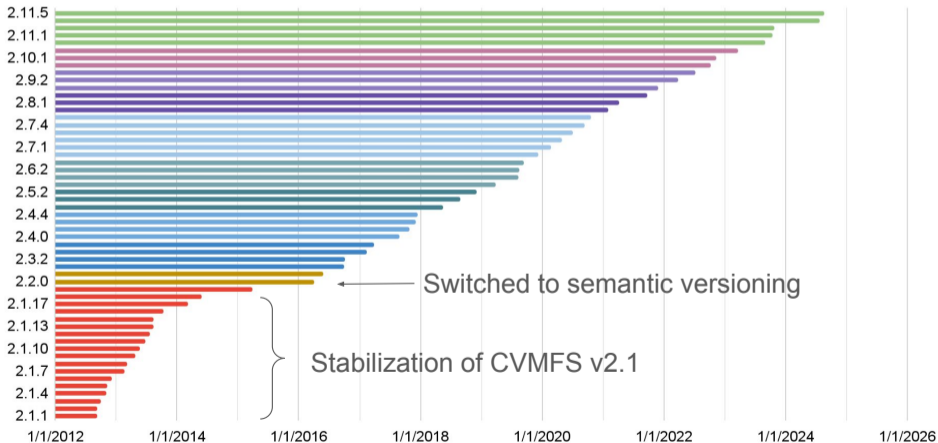
Approaches in chronological order

1. Traverse entire repository: does not scale
2. Fuse module for writing: slow and intricate
3. Change log: worked, brittle (until ~2011; propagation delay 0.5 days)
4. Union file system: AUFS for a while the only working kernel-level union file system, but not in the mainline kernel
 - Provided AUFS patched RHEL6 kernel
 - Moved to overlayfs in RHEL7 (initial contribution from Wellcome Sanger)
 - Propagation delay at <15 minutes
5. Special cases: grafting (since 2016), direct ingestion (since 2019)





CernVM-FS Release Dates



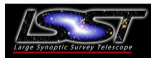


Off to new shores: 2015-2020



The production system stabilized, the environment changed

- Open source S3 compatible object stores became available (Ceph S3) fantastic alternative for CVMFS backend storage; took ~5 years from prototype to full production
- 1 TB per night of integration builds to be distributed to the grid: repository garbage collection
- Distributed publishing: long envisaged (since 2012) and deployed around 6 years later
- Plus: containers, HPC, data distribution
- New developments were summarized at CHEP 2018 in Sofia

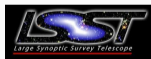




The production system stabilized, the environment changed

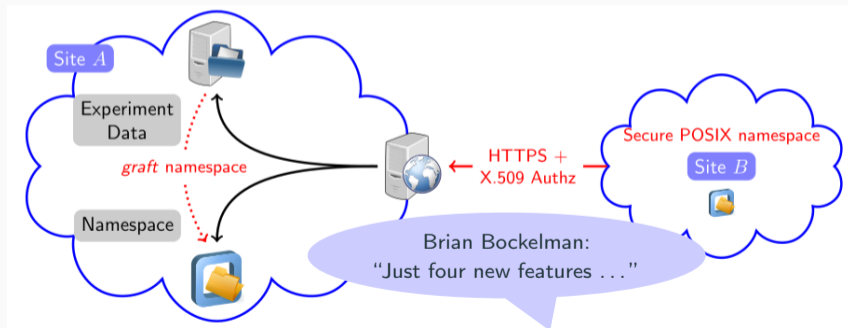
- Open source S3 compatible object stores became available (Ceph) a fantastic alternative for CVMFS backend storage; took ~5 years
- 1 TB per night of integration builds to be distributed to the grid
- Distributed publishing: long envisaged (since 2012) and deployed around 6 years later
- Plus: containers, HPC, data distribution
- New developments were summarized at CHEP 2018 in Sofia

Close collaboration with Ben Couturier and Enrico Bocchi





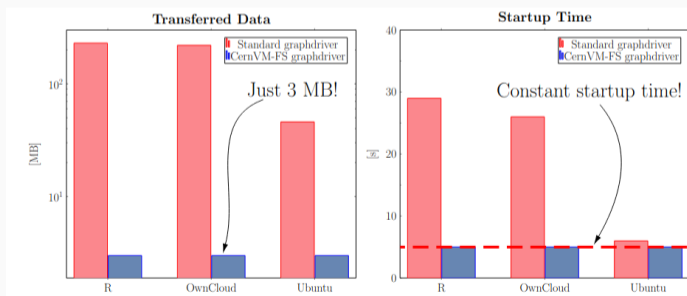
1. Namespace grafting
2. External data (split squid and XRootD HTTP traffic)
3. Uncompressed files
4. Authorization plugins





- With the advent of container virtualization (docker etc.), the platform isolation problem was finally solved without performance overhead.
- However, unlike full virtualization, we could (not yet) mount our file system in a container.
- The software distribution problem very much stayed the same; the docker “layers” offered modularization for image construction, but little advantage for distribution.

1. Tricked Docker into reading from CVMFS
2. Singularity/apptainer allowed for direct use of unpacked images
3. containerd snapshotter interface for file-based transfer





So close, yet so far:
no fuse, no internet connectivity, no local disk (cache), no site caches

Required development of a set of deployment options à la carte:

- preloader
- shrinkwrap
- parrot connector (now: cvmfsexec)
- cache plugins





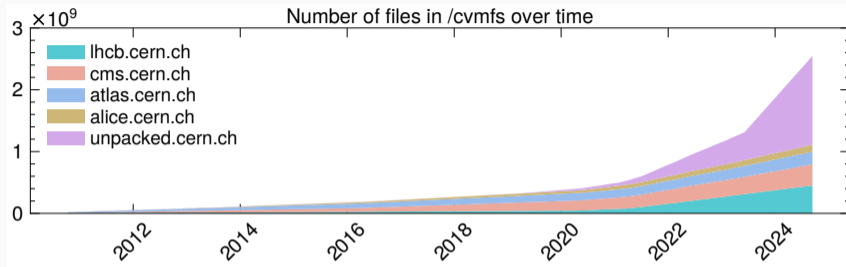
So close, yet so far:
no fuse, no internet connectivity, no local disk (cache), no site caches

Required development of a set of deployment options à la carte:

- preloader
- shrinkwrap
- parrot connector (now: cvmfs)
- cache plugins

Close collaboration
with Douglas Thain's
team at Notre Dame

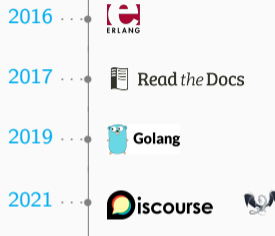




- Accelerating repository size: more platforms, more CI builds, containers, Python, ...
 - New, exciting use cases: $\mathcal{O}(100\text{ PB})$ data archive distributed at **jumptrading**
 - **EESSI**: a new approach to compute environment on European HPCs
 - Continuous improvement in Linux and fuse; most importantly: unprivileged fuse mounts (cvmfsexec)
- **Subject of this workshop**



Internal Affairs



- git replaced svn
- JIRA replaced Savannah
- Perl was phased out in 2018
- Jenkins replaced Electric Commander
- golang replaced erlang
- Github issues replaced JIRA in 2020
- clang-tidy replaced Coverity in 2021

Former Influential CernVM-FS Contributors



Predrag Buncic
CernVM Founder



Artem Harutyunyan
CernVM Original Team
→ Qualys, Mesosphere, Bardeen



Carlos Aguado Sanchez
CernVM Original Team
→ AWS



René Meusel
Union-FS based server
→ Rohde & Schwarz



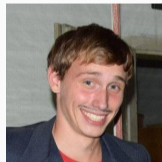
Nikola Hardi
Docker graph driver
→ Swisscom



Simone Mosciatti
unpacked.cern.ch
→ AWS



Radu Popescu
Gateway
→ Logitech



Jan Priessnitz
Parallel GC
→ MPI



Andrea Valenzuela
Containerized publisher
→ CMS



Thanks to Catalin Condurache



Thanks to Dennis van Dok

- Thank you all for the trust you put in the project and for patiently working with us
- Thank you, Predrag, for having run this wonderful incubator
- Roles changed – Valentin Völkl now at the helm of the CernVM project!

