# Software Distribution with CernVM-FS

Jakob Blomer    CERN, EP-SFT

✉ jblomer@cern.ch
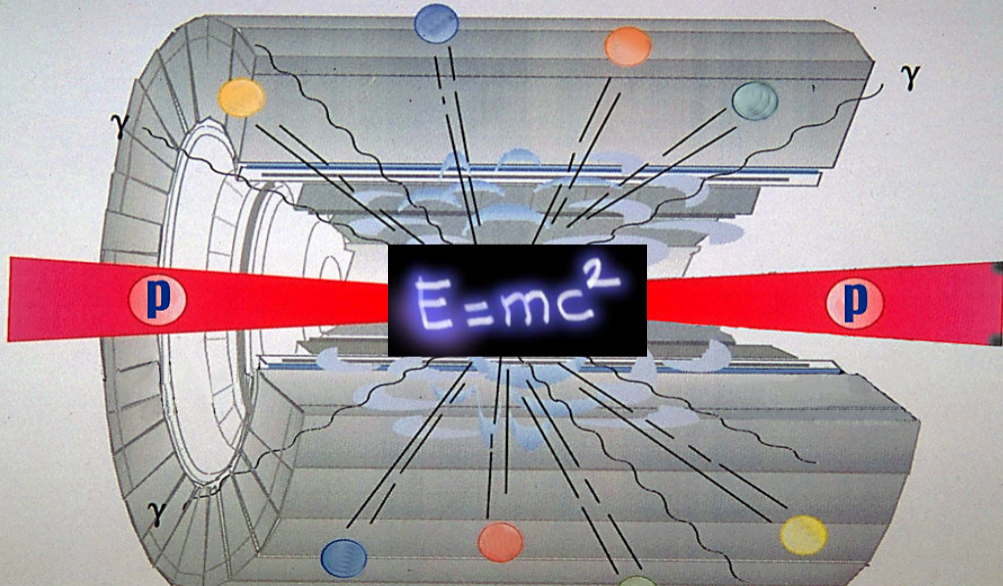
⌂ https://github.com/cvmfs/cvmfs
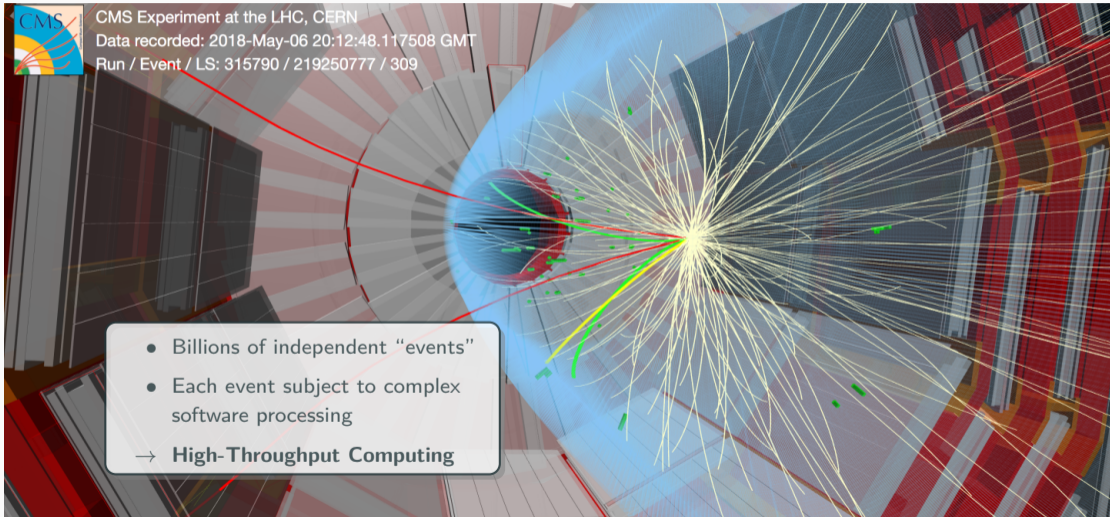
2019-05-17

- Billions of independent "events"
- Each event subject to complex software processing
- → **High-Throughput Computing**

**Distribution of All CERN Users by Location of Institute on 24 January 2018**



- **Physics and computing: international collaborations**
- "The Grid": ~160 data centers
- Approx.: global batch system

**MEMBER STATES**

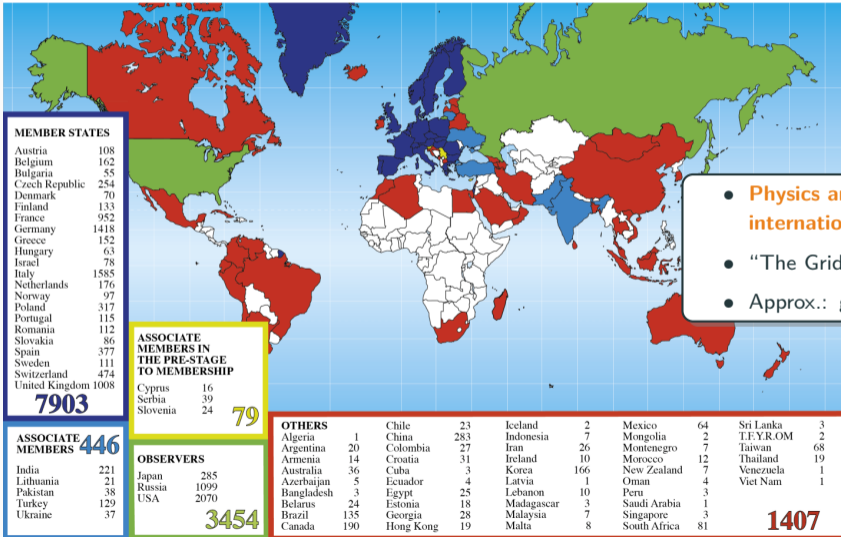| | |
|---|---|
| Austria | 108 |
| Belgium | 162 |
| Bulgaria | 55 |
| Czech Republic | 254 |
| Denmark | 70 |
| Finland | 133 |
| France | 952 |
| Germany | 1418 |
| Greece | 152 |
| Hungary | 63 |
| Israel | 78 |
| Italy | 1585 |
| Netherlands | 176 |
| Norway | 97 |
| Poland | 317 |
| Portugal | 115 |
| Romania | 112 |
| Slovakia | 86 |
| Spain | 377 |
| Sweden | 111 |
| Switzerland | 474 |
| United Kingdom | 1008 |

**7903**

**ASSOCIATE MEMBERS IN THE PRE-STAGE TO MEMBERSHIP**

| | |
|---|---|
| Cyprus | 16 |
| Serbia | 39 |
| Slovenia | 24 |

**79**

**ASSOCIATE MEMBERS 446**

| | |
|---|---|
| India | 221 |
| Lithuania | 21 |
| Pakistan | 38 |
| Turkey | 129 |
| Ukraine | 37 |

**OBSERVERS**

| | |
|---|---|
| Japan | 285 |
| Russia | 1099 |
| USA | 2070 |

**3454**

**OTHERS**

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Algeria | 1 | Chile | 23 | Iceland | 2 | Mexico | 64 |
| Argentina | 20 | China | 283 | Indonesia | 7 | Mongolia | 2 |
| Armenia | 14 | Colombia | 27 | Iran | 26 | Montenegro | 7 |
| Australia | 36 | Croatia | 31 | Ireland | 10 | Morocco | 12 |
| Azerbaijan | 5 | Cuba | 3 | Korea | 166 | New Zealand | 7 |
| Bangladesh | 3 | Ecuador | 4 | Latvia | 1 | Oman | 4 |
| Belarus | 24 | Egypt | 25 | Lebanon | 10 | Peru | 3 |
| Brazil | 135 | Estonia | 18 | Madagascar | 3 | Saudi Arabia | 1 |
| Canada | 190 | Georgia | 28 | Malaysia | 7 | Singapore | 3 |
| | | Hong Kong | 19 | Malta | 8 | South Africa | 81 |

| | |
|---|---|
| Sri Lanka | 3 |
| T.F.Y.R.OM | 2 |
| Taiwan | 68 |
| Thailand | 19 |
| Venezuela | 1 |
| Viet Nam | 1 |

**1407**

5

❶ Provide uniform, consistent, and versioned POSIX file system access to `/cvmfs`

```
$ ls /cvmfs/cms.cern.ch
slc7_amd64_gcc700   slc7_ppc64le_gcc530   slc7_aarch64_gcc700   slc6_mic_gcc481
...
```
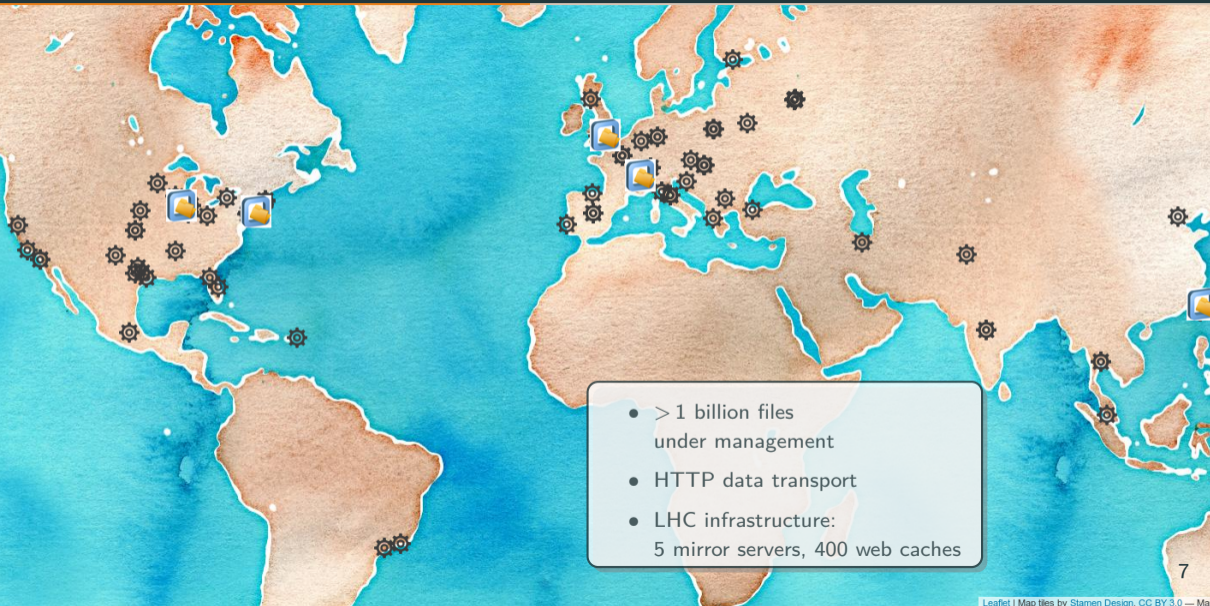
on **grids**, **clouds**, **supercomputers** and **end user laptops**
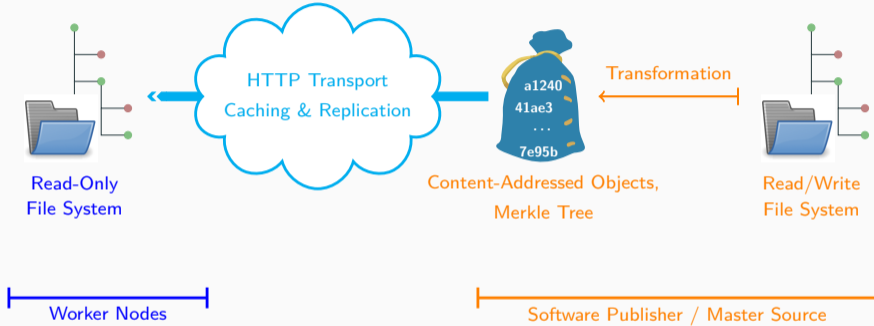
*read*

*publish*

❷ Populate and propagate new and updated content

- A few "software librarians" can publish into `/cvmfs`
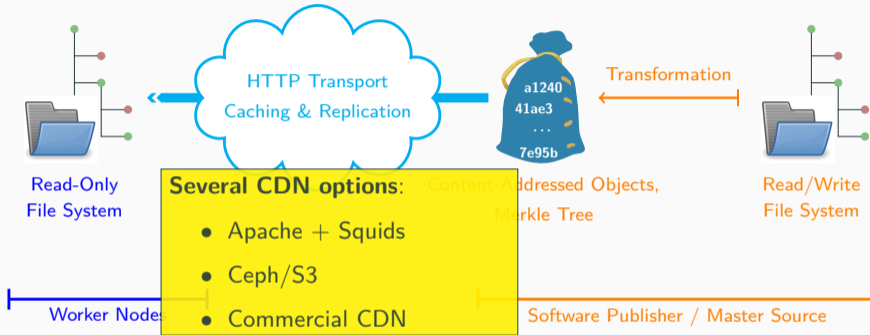- Transactional writes as in `git commit/push`

- > 1 billion files under management
- HTTP data transport
- LHC infrastructure: 5 mirror servers, 400 web caches
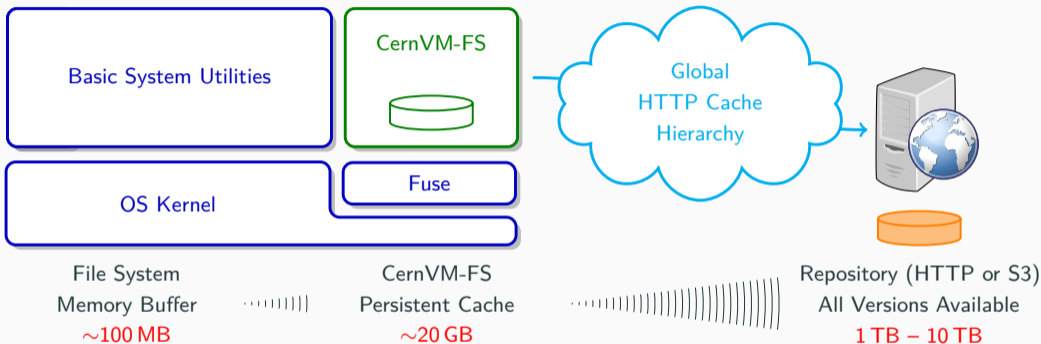
- Reading and writing treated asymmetrically
- Immutable objects, stateless services

- HTTP transport + caching
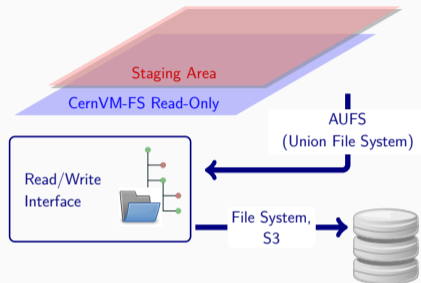- Consistency over availability

HTTP Transport
Caching & Replication

Transformation

a1240
41ae3
. . .
7e95b

Read-Only
File System

Read/Write
File System

Content-Addressed Objects,
Merkle Tree

**Several CDN options**:

- Apache + Squids
- Ceph/S3
- Commercial CDN

Worker Nodes

Software Publisher / Master Source

- Reading and writing treated asymmetrically
- Immutable objects, stateless services

- HTTP transport + caching
- Consistency over availability

8
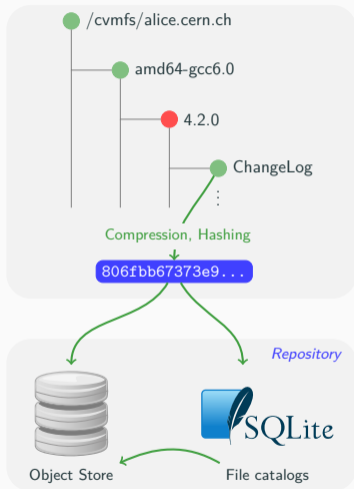
- Fuse based, independent mount points, e. g. /cvmfs/atlas.cern.ch
- High cache effiency because entire cluster likely to use same software

- Kernel-level union file system:
  AUFS, OverlayFS

**Publishing new content**

```
[ ~ ]# cvmfs_server transaction containers.cern.ch
[ ~ ]# cd /cvmfs/containers.cern.ch && tar xvf ubuntu1610.tar.gz
[ ~ ]# cvmfs_server publish containers.cern.ch
```

# Use of Content-Addressable Storage



**Object Store**
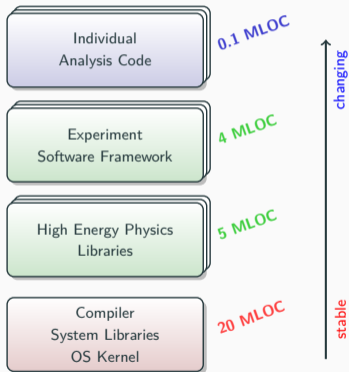- Compressed files and chunks
- De-duplicated

**File Catalog**
- Directory structure, symlinks
- Content hashes of regular files
- Digitally signed
  $\Rightarrow$ integrity, authenticity
- Time to live
- Partitioned / Merkle hashes
  (possibility of sub catalogs)

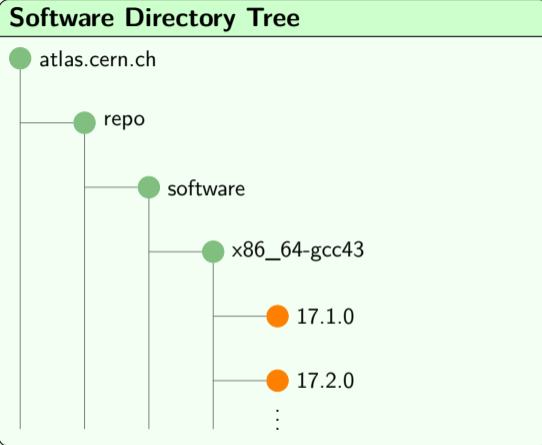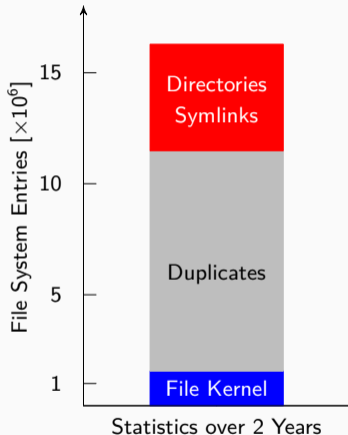$\Rightarrow$ Immutable files, trivial to check for corruption, versioning

Why a file system?

```
$ cmsRun DiPhoton_Analysis.py
```

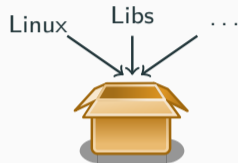| | |
|---|---|
| Individual Analysis Code | 0.1 MLOC |
| Experiment Software Framework | 4 MLOC |
| High Energy Physics Libraries | 5 MLOC |
| Compiler System Libraries OS Kernel | 20 MLOC |

changing

stable

**Key Figures**

- Hundreds of (novice) developers
- Hundred million binaries
- 1 TB / day of nightly builds
- ~100 000 machines world-wide
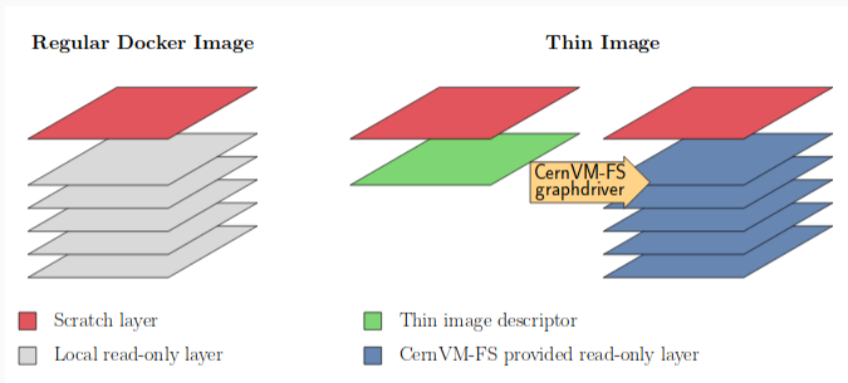- Daily production releases, remain available "eternally"

$\rightarrow$ **too much for a packaging approach**

13

**Software Directory Tree**

atlas.cern.ch

repo

software

x86_64-gcc43

17.1.0

17.2.0

File System Entries [×10⁶]

Directories Symlinks

Duplicates

File Kernel

Statistics over 2 Years

Between consecutive software versions: only ∼15 % new files
At runtime only tiny fraction of files actually accessed

14

Linux    Libs    . . .



- Containers are easier to create than to role-out at scale
- Ideally: containers for isolation and a software file system for distribution

**Regular Docker Image**

**Thin Image**

CernVM-FS
graphdriver

- 🟥 Scratch layer
- ⬜ Local read-only layer
- 🟩 Thin image descriptor
- 🟦 CernVM-FS provided read-only layer

Ideally: native support for (some) unpacked layers on a read-only file sytem