

CernVM 5

A versatile container-based platform
to run HEP applications



CernVM Appliance

- Started in 2008
- Virtual Appliance for CernVM-FS
→ A minimal yet complete run time environment for scientific software
- Version ≤ 4
 - Minimal Linux kernel + μ CernVM bootloader
 - CernVM-FS mounted during boot process; system loaded on-demand
 - Volume of ~20MB, but designed for full virtualization
- **Problem: Limited usability as a container**



Motivation

- **Operators are shifting infrastructure towards containerized environments**
 - Smaller images
 - Easier to distribute / retain
 - Less performance overhead
 - Easier and faster to deploy
 - Powerful orchestration tools (e. g. Kubernetes)



CernVM 5

- Still a complete platform to develop and run HEP applications
→ **But: Focus on container virtualization**
- Provide access to CernVM-FS in containerized environments
→ By mounting it in unprivileged containers running in various run times
- Equally practical as a container and full virtual machine
→ Running in versatile container run times and hypervisors

CernVM 5: Goals

- Minimal container image for a variety of applications
- Mounting CernVM-FS inside a container
- Stable usability in the presence of multiple versions of shared libraries
- Graphical user interface
- CernVM-FS cache management
 - Reducing start latency by building images with pre-loaded cache

CernVM 5: Goals

- ▶ • Minimal container image for a variety of applications
- Mounting CernVM-FS inside a container
- Stable usability in the presence of multiple versions of shared libraries
- Graphical user interface
- CernVM-FS cache management
 - Reducing start latency by building images with pre-loaded cache

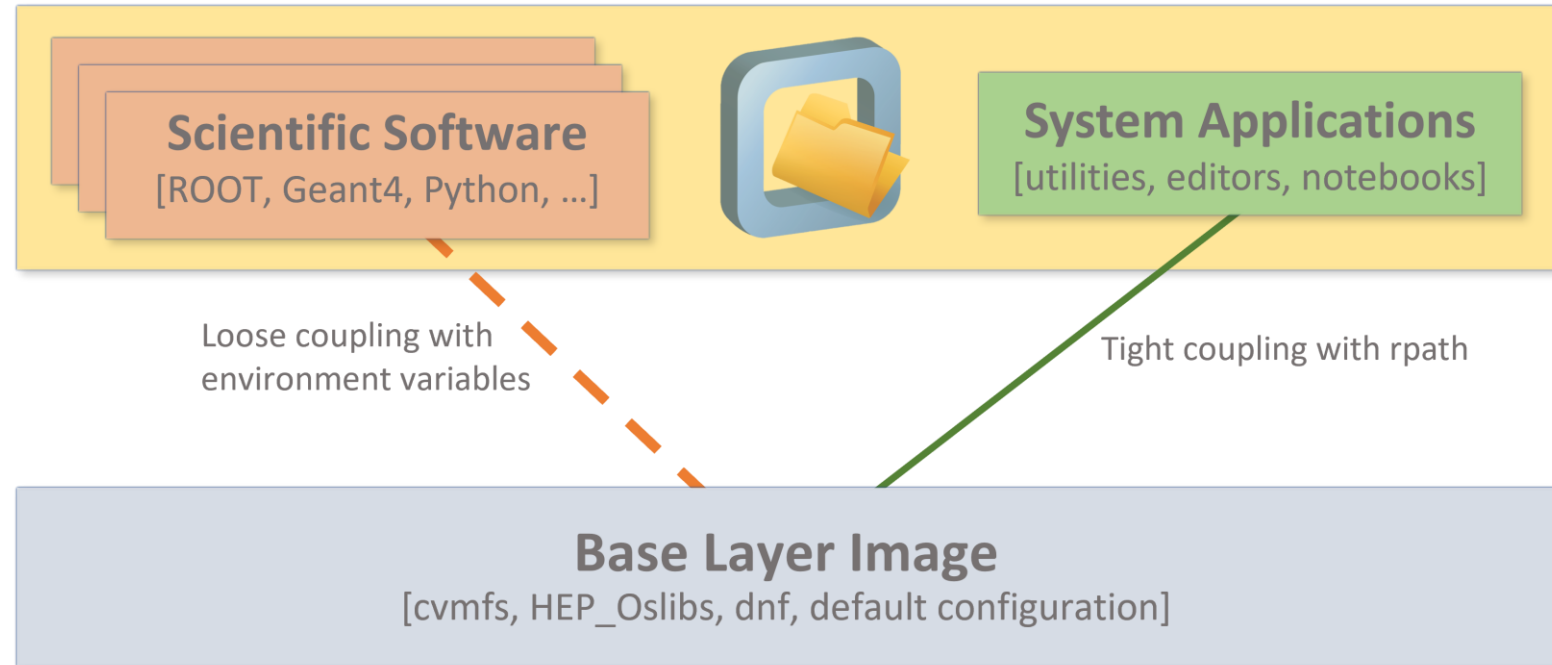
CernVM 5: Design

- Base Layer Image
 - CernVM-FS client
 - HEP_OSlibs
 - dnf package manager

→ The actual image, smallest unit
- System Applications
 - CernVM 5 system specific
 - Interactive applications

→ What is expected on-top
- Scientific Software Stacks
 - E. g. analysis frameworks
 - Set up using environment variables

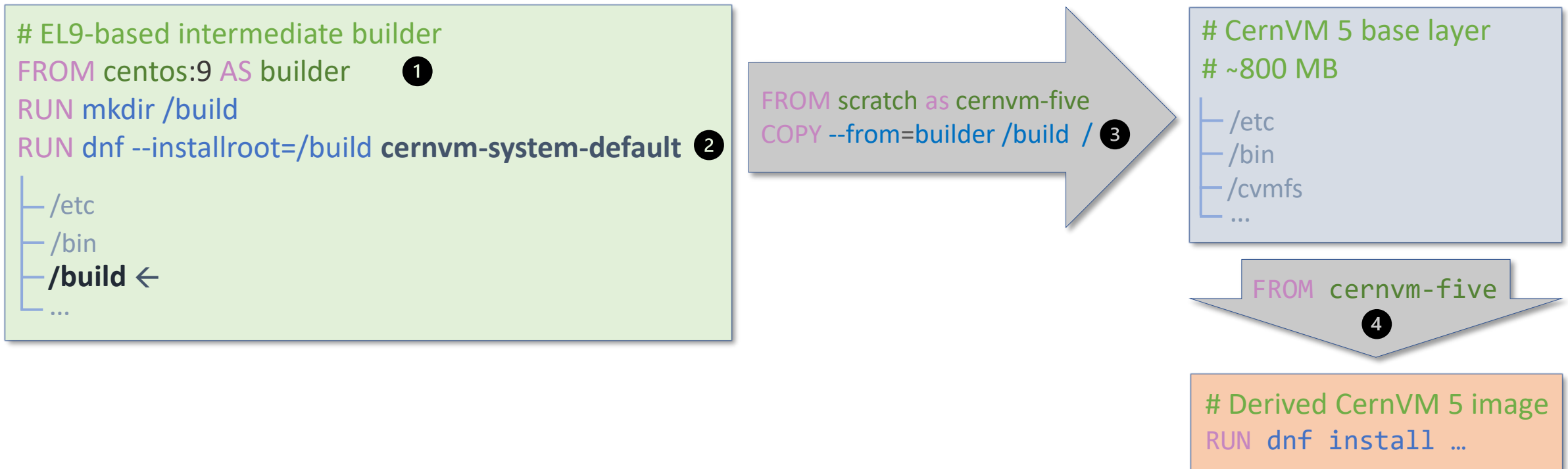
→ Externally managed



CernVM 5: Image Build Process



- Custom multi-stage build process
 - 1.) Set up intermediate build container
 - 2.) Construct CernVM 5 root file system in build directory
 - Defined in **cernvm-system-default** package (HEP_OSlibs, CernVM-FS, dnf, configuration files)
 - 3.) and 4.) Export build directory to standalone image



CernVM 5: System Applications on CernVM-FS



Outsourced system applications on CernVM-FS

- + Smaller image size
- + General advantages of CernVM-FS like e.g., lazy loading
- + Versioned by CernVM-FS

Subsystem-like installation

- + Easily extendable using standard package managers

System Applications Repository on CernVM-FS

[vim, nano, emacs, findutils, patchelf, ping, wget, strace, tree, diff, ...]

CernVM 5: Resulting Base Layer Image

	CernVM 5 Base Layer Image	Naively Derived Image	Naively Derived Image + System Applications
Volume (uncompressed)	805MB	1030MB	1450MB
Volume (compressed)	284MB	382MB	523MB
Installed Packages	457	502	612
Image Layers	single	multiple	multiple
Standard Derivability	yes	yes (but not a true base layer)	yes (but not a true base layer)

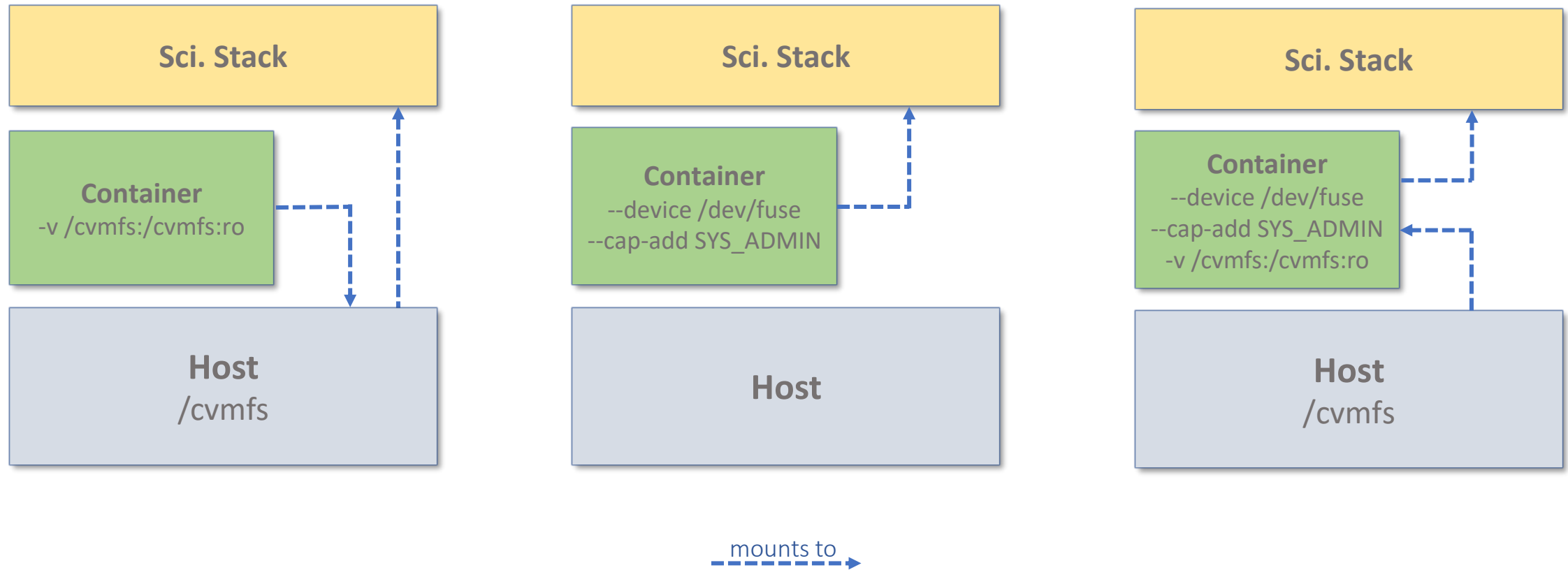
- + Better control over installed packages → smaller image, easier to distribute
- + Improved build time and caching
- + Single layered → less complex, faster to construct
- + Compared to ≤ 4 : Derivable image

CernVM 5: Goals

- Minimal container image for a variety of applications
- ▶ • Mounting CernVM-FS inside a container
- Stable usability in the presence of multiple versions of shared libraries
- Graphical user interface
- CernVM-FS cache management
 - Reducing start latency by building images with pre-loaded cache



CernVM 5: Mounting CernVM-FS



CernVM 5: Deployment



	Docker	Apptainer	Podman	Kubernetes	containerd
Bind mount CernVM-FS <code>-v /cvmfs:/cvmfs:ro</code>	✓	✓	✓	✓	✓
Mount CernVM-FS <code>--device /dev/fuse</code>	✓	✓	✓	✓	✓
Pre-mounted	✗	✓	✗	✗	✗
Host integration <code>/workspace</code> → as a shared folder <code>/data</code> → physics data	✓	✓	✓	✓	✓

CernVM 5: Distribution

Pushed to registry

- CERN Harbor Registry

CernVM-FS

→ Large scale distribution

As a VM

- CERN OpenStack
- Isolation use cases
- Cloud infrastructure
- Data acquisition systems

```
$ docker pull registry.cern.ch/cernvm/five/cernvm-five:latest
```

```
$ aptainer shell /cvmfs/unpacked.cern.ch/registry.cern.ch/cernvm/five/cernvm-five\:latest
```

```
$ openstack server create --image 'CernVM 5.1.2 – x86_64' (soon)
```

CernVM 5: Goals

- Minimal container image for a variety of applications
- Mounting CernVM-FS inside a container
- ▶ • Stable usability in the presence of multiple versions of shared libraries
- Graphical user interface
- CernVM-FS cache management
 - Reducing start latency by building images with pre-loaded cache

Problem: Presence of multiple versions of shared libraries

```
bash-4.4# which nano
/cvmfs/cernvm-five.cern.ch/x86_64/1.0a/usr/bin/nano
bash-4.4# source /cvmfs/sw.hsf.org/spackages4/key4hep-stack/release-2021-10-29-ip7764o/x86_64-centos8-gcc8.4.1-opt/setup.sh
bash-4.4# nano
Segmentation fault (core dumped)
```

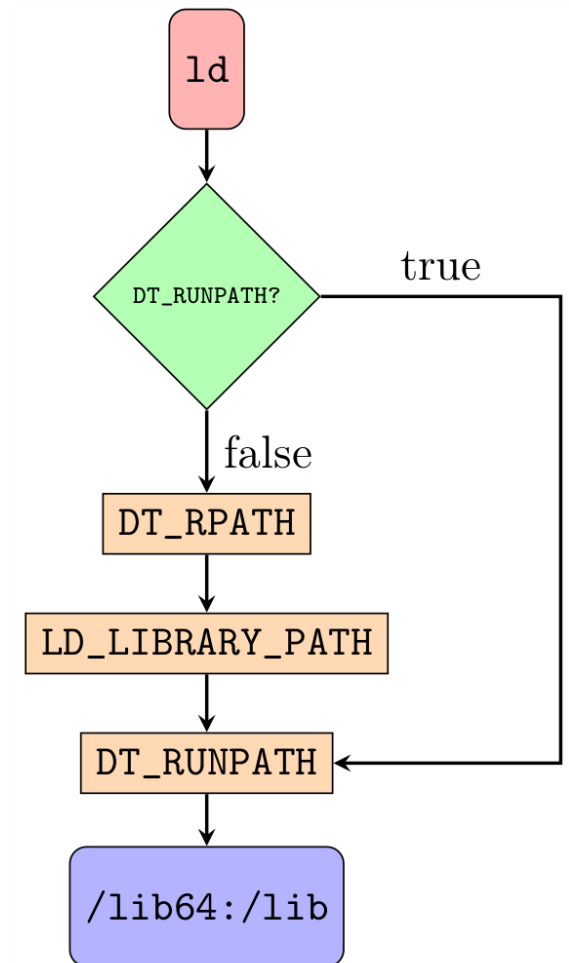
```
bash-5.1# ldd /cvmfs/cernvm-five.cern.ch/x86_64/1.0a/usr/bin/nano
linux-vdso.so.1 (0x00007f6d5ac3b000)
libmagic.so.1 => /cvmfs/sw.hsf.org/spackages4/file/5.40-qsol6zg/x86_64-centos8-gcc8.4.1-opt/lib/libmagic.so.1 (0x00007f6d5a7bf000)
libncursesw.so.6 => /cvmfs/sw.hsf.org/spackages4/ncurses/6.2-rxkbhc6/x86_64-centos8-gcc8.4.1-opt/lib/libncursesw.so.6 (0x00007f6d5a7bf000)
libtinfo.so.6 => /cvmfs/sw.hsf.org/spackages4/ncurses/6.2-rxkbhc6/x86_64-centos8-gcc8.4.1-opt/lib/libtinfo.so.6 (0x00007f6d5a34b000)
```

- Stack setup scripts override overwrite local default values
→ Locally and remotely installed executables break when linked with mismatching shared library

Dynamically Linked Executables

- Dynamically linked ELF executables
 - Use shared libraries installed on a system
 - Linked with dependencies during run time using a linker, e. g., ld
- Possible solutions
 - Static linking
 - Increases image size
 - rpath at compile-time
 - Requires custom packages / maintenance

→ Standard packages with post-build rpath processing



Standard Packages with Post-Build rpath



```
bash-5.1# cernvm-patch-rpath -h
Adds DL_RPATH to ELF executables using patchELF
cernvm-patch-rpath [-r <installation root (by default '/')> -l <location of shared libraries (by default '/lib64:/lib' resp. '/lib')>
```

- cernvm-patch-rpath:
 - 1.) Determine installed executables
 - 2.) Evaluate executable and original location of dependencies
 - E. g., /lib64 for a locally installed 64-bit, dynamically linked ELF executable
 - 3.) Set up DT_RPATH binary header accordingly

Pros:

- + Fast, repeatable
- + Isolated from LD_LIBRARY_PATH
- + Integrated in build chain
- + Can be used by users in their builds (**RUN** cernvm-patch-rpath)
- + Allows use of standard packages

Example: CernVM 5 executable with post-build applied rpath

```
bash-5.1# patchelf --print-rpath /cvmfs/cernvm-five.cern.ch/x86_64/1.2/usr/bin/nano
/cvmfs/cernvm-five.cern.ch/x86_64/1.2/lib64:/cvmfs/cernvm-five.cern.ch/x86_64/1.2/lib
```

```
bash-5.1# which nano
/cvmfs/cernvm-five.cern.ch/x86_64/1.2/usr/bin/nano
bash-5.1# source /cvmfs/sw.hsf.org/spackages4/key4hep-stack/release-2021-10-29-ip7764o/x86_64-centos8-gcc8.4.1-opt/setup.sh
bash-5.1# nano
bash-5.1# echo $?
0
```

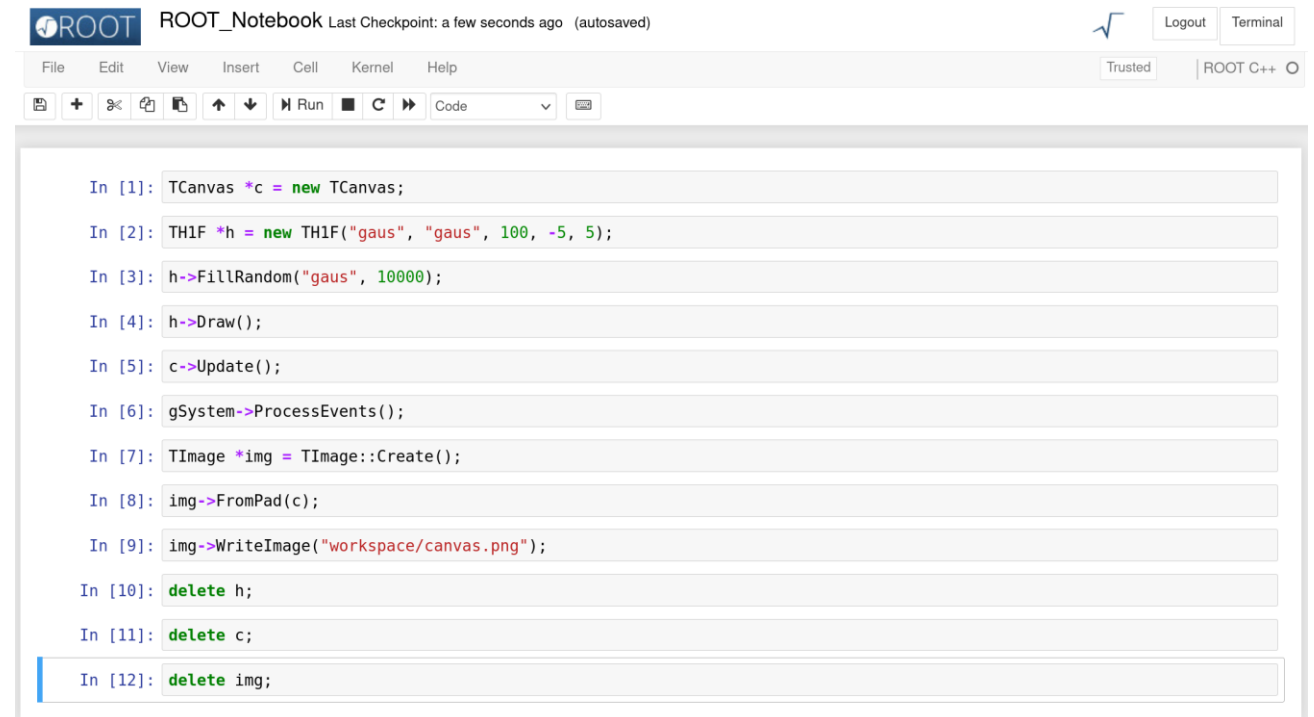
```
bash-5.1# ldd $(which nano)
linux-vdso.so.1 (0x00007ffd7abea000)
libmagic.so.1 => /cvmfs/cernvm-five.cern.ch/x86_64/1.2/lib64/libmagic.so.1 (0x00007fd7f1f1d000)
libncursesw.so.6 => /cvmfs/cernvm-five.cern.ch/x86_64/1.2/lib64/libncursesw.so.6 (0x00007fd7f1eda000)
libtinfo.so.6 => /cvmfs/cernvm-five.cern.ch/x86_64/1.2/lib64/libtinfo.so.6 (0x00007fd7f1ea9000)
```

CernVM 5: Goals

- Minimal container image for a variety of applications
- Mounting CernVM-FS inside a container
- Stable usability in the presence of multiple versions of shared libraries
- ▶ • Graphical user interface
- CernVM-FS cache management
 - Reducing start latency by building images with pre-loaded cache

CernVM 5: Graphical User Interface

- Web-based Jupyter Notebooks rather than a fully fledged GUI
 - Loaded from CernVM-FS
 - Hosted locally
(EXPOSE 8888)



The screenshot shows the ROOT Notebook interface. The title bar reads "ROOT Notebook Last Checkpoint: a few seconds ago (autosaved)". The menu bar includes "File", "Edit", "View", "Insert", "Cell", "Kernel", and "Help". The toolbar contains icons for file operations and execution. The code cell contains the following C++ code:

```
In [1]: TCanvas *c = new TCanvas;
In [2]: TH1F *h = new TH1F("gaus", "gaus", 100, -5, 5);
In [3]: h->FillRandom("gaus", 10000);
In [4]: h->Draw();
In [5]: c->Update();
In [6]: gSystem->ProcessEvents();
In [7]: TImage *img = TImage::Create();
In [8]: img->FromPad(c);
In [9]: img->WriteImage("workspace/canvas.png");
In [10]: delete h;
In [11]: delete c;
In [12]: delete img;
```

CernVM 5: Goals

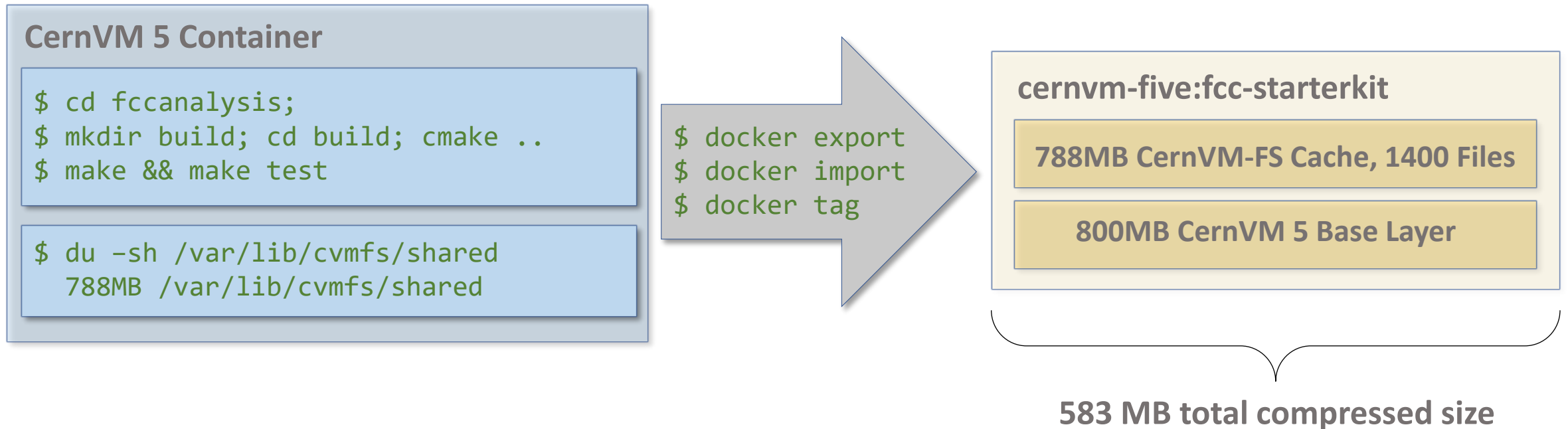
- Minimal container image for a variety of applications
- Mounting CernVM-FS inside a container
- Stable usability in the presence of multiple versions of shared libraries
- Graphical user interface
- ▶ • CernVM-FS cache management
 - Reducing start latency by building images with pre-loaded cache

CernVM-FS Cache Management

- Problem: High start latency of the CernVM-FS client when loading large packages for the first time
 - Because many small binaries need to be loaded
- Idea: Have container (and VM) images with pre-loaded cache for specific workloads
 - Tutorials
 - Computing campaigns
- But: Stacks change frequently
 - Automatic build chain required

Cache Pre-Loading – An Example

- Scenario / Workflow mapped in a script / testsuite
- Example FCC Starterkit tutorial



Speed Up (Inside CERN Network)

Command	CernVM 5 Base Image	CernVM 5 FCC Starterkit	Factor
KKMCee -h	1.609s	0.357s	4.5x
BHLUMI -h	1.569s	0.055s	28.0x
whizard Z_mumu.sin	50.233s	18.971s	2.6x
babayaga -h	1.635s	0.487s	3.4x
babayaga -f ¹	5.541s	5.289s	1.05x
Import ROOT	6.843s	1.544s	4.4x
fccanalysis run ²	11.230s	9.319s	1.2x

1 babayaga -f 15. -t 165. -e 91.2 -n 10000 -o bbyg_10000.LHE

2 fccanalysis run analysis_stage1.py --output p8_ee_ZH_ecm240.root --files-list ./p8_ee_ZH_ecm240_edm4hep.root

http://ecsft.cern.ch/dist/cernvm/five/vm/tutorials/cernvm-five-fcc-x86_64.qcow2.tar.gz

CernVM 5: VM Image Build Process

- Constructing the Root File System
 - Adding Kernel-enabled file system layer to CernVM base image

```
# Deriving from base layer (or CernVM 5 child images)  
FROM registry.cern.ch/cernvm/five/cernvm-five:latest
```

```
# Installing kernel and configurations, generating initrd  
RUN dnf install cernvm-kernel-default
```

```
|— /cvmfs  
|— /etc ←  
|— /boot ←  
|— ...
```

```
# Construct rootfs  
$ docker create vm  
  
# export rootfs  
$ docker export vm
```

RootFS (cernvm.tar)



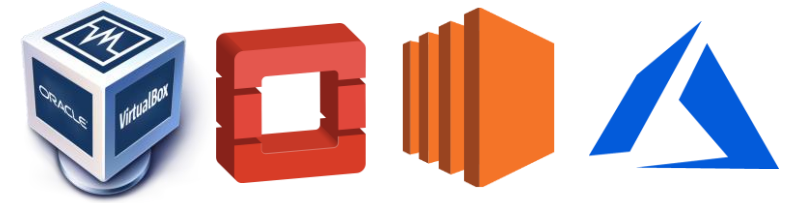
Kernel / initrd

Additional Packages

CernVM-FS Cache

CernVM 5 base layer

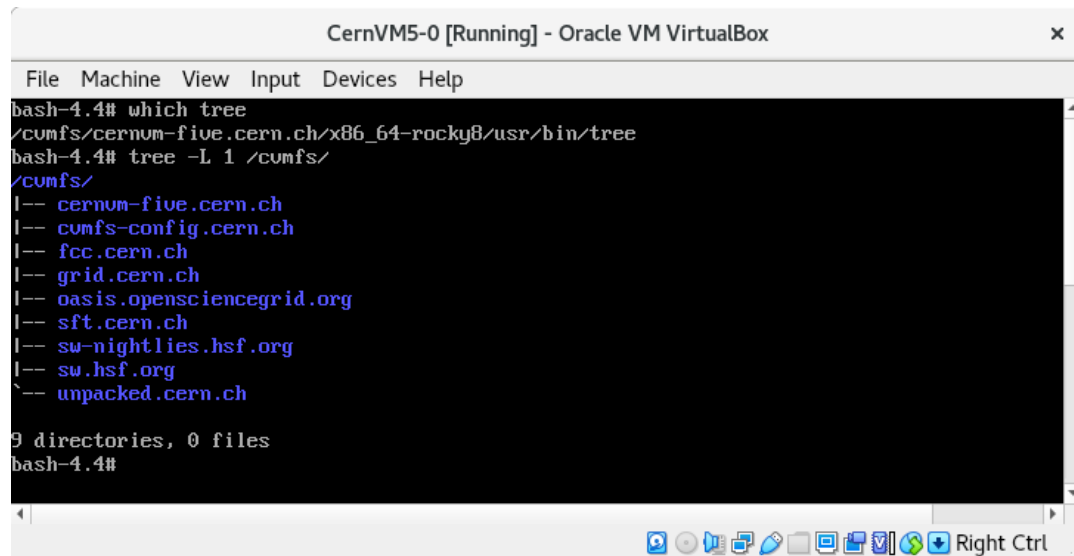
CernVM 5: Full Virtual Machine



- For cloud infrastructure
 - Such as KVM/CERN OpenStack, EC2, Azure
 - Contextualization (SSH keys, mountpoints, user data)
- For desktop hypervisor
 - Such as QEMU/KVM, VirtualBox
- Isolation-sensitive use cases, e. g. in data acquisition systems

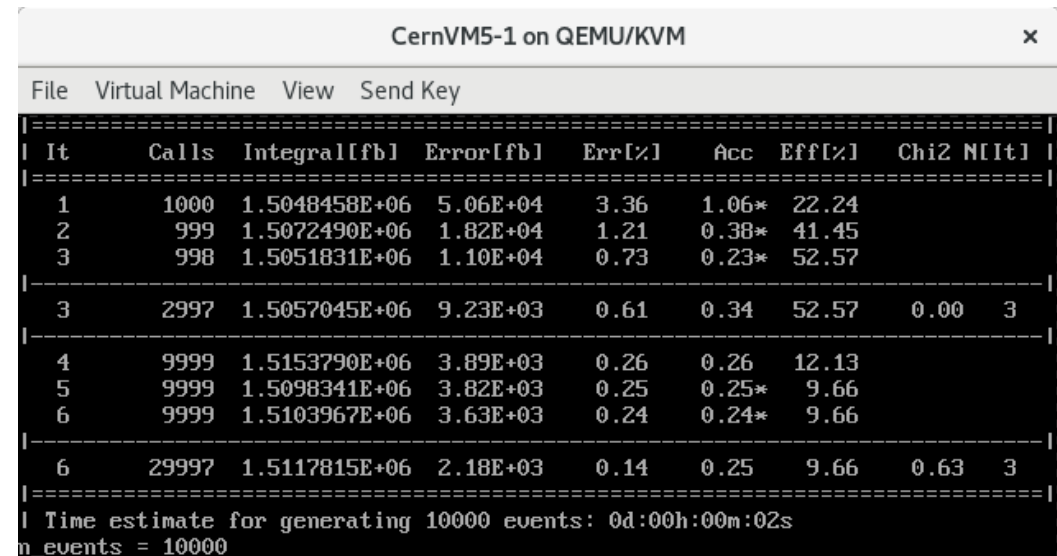
CernVM 5: Virtual Machine

- Available as raw and qcow2 and as cloud template image
- In line with the container build process
 - Every existing CernVM 5 container image and container extendable to full VM



```
CernVM5-0 [Running] - Oracle VM VirtualBox
File Machine View Input Devices Help
bash-4.4# which tree
/cumfs/cernvm-five.cern.ch/x86_64-rocky8/usr/bin/tree
bash-4.4# tree -L 1 /cumfs/
/cumfs/
|-- cernvm-five.cern.ch
|-- cumfs-config.cern.ch
|-- fcc.cern.ch
|-- grid.cern.ch
|-- oasis.opensciencegrid.org
|-- sft.cern.ch
|-- sw-nightlies.hsf.org
|-- sw.hsf.org
`-- unpacked.cern.ch

9 directories, 0 files
bash-4.4#
```



```
CernVM5-1 on QEMU/KVM
File Virtual Machine View Send Key
=====
| It      Calls  Integral[fb]  Error[fb]  Err[%]  Acc  Eff[%]  Chi2  N[It] |
=====
| 1       1000  1.5048458E+06  5.06E+04   3.36    1.06*  22.24  |
| 2        999  1.5072490E+06  1.82E+04   1.21    0.38*  41.45  |
| 3        998  1.5051831E+06  1.10E+04   0.73    0.23*  52.57  |
|-----|
| 3       2997  1.5057045E+06  9.23E+03   0.61    0.34   52.57  0.00  3  |
|-----|
| 4       9999  1.5153790E+06  3.89E+03   0.26    0.26   12.13  |
| 5       9999  1.5098341E+06  3.82E+03   0.25    0.25*   9.66  |
| 6       9999  1.5103967E+06  3.63E+03   0.24    0.24*   9.66  |
|-----|
| 6      29997  1.5117815E+06  2.18E+03   0.14    0.25   9.66  0.63  3  |
|-----|
| Time estimate for generating 10000 events: 0d:00h:00m:02s
n_events = 10000
```

CernVM 5: Conclusion

- Container image capable of serving arbitrary scientific stacks as a run time environment
- Built in a custom process using standard packages
- Processing of dynamically linked executables with DT_RPATH
- Natively derivable
- Mounting CernVM-FS inside various container runtimes via FUSE interface
- Extendable to a full virtual machine
- Jupyter Notebooks hosting
- Proof-of-concept for pre-loading CernVM-FS caches

Questions?