



Cloud Technologies in High-Energy Physics Data Processing

Jakob Blomer
CERN, EP-SFT

ICFDT 2016
1st April 2016

Rob Joyce, NSA, Head of Tailored Access Operations

*Cloud computing is a really fancy term for
someone else's computer.*

Rob Joyce, NSA, Head of Tailored Access Operations

*Cloud computing is a really fancy term for
someone else's computer.*

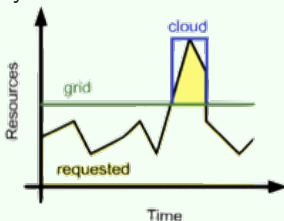
- ① High-energy physic's idea of cloud computing is mostly limited to "Infrastructure-as-a-Service" (IaaS)
- ② I will focus in this talk at cloud computing in the context of a means to improve our distributed systems ("the grid")

Drivers and Obstacles

- ① Cost ☺ (partially)
- ② Control & Trust ☹
- ③ Specialized applications ☹
- ④ New Technologies for Distributed Systems ☺
 - Virtualization
 - BLOB storage
 - NoSQL databases

Themes

- ① Hybrid academic-commercial clouds



Source: Mato

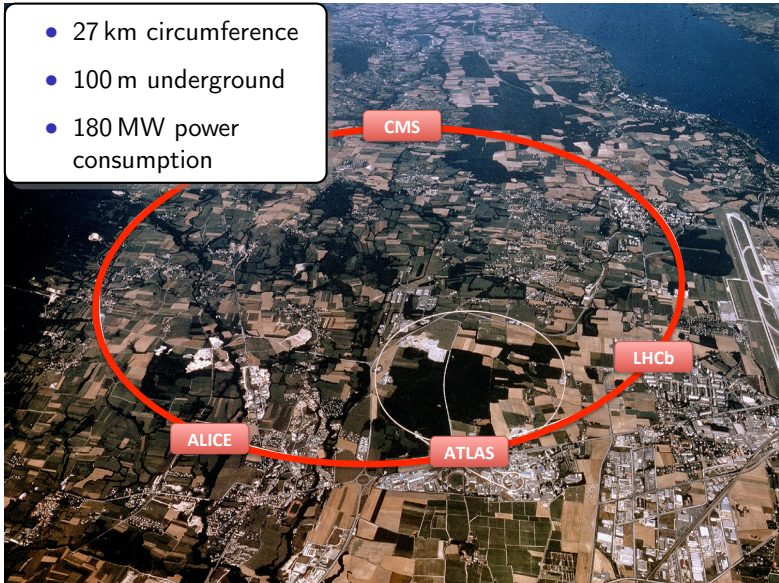
- ② Offload mainly simulation (up to 50%), i. e. no data lock-in
- ③ “Private” adoption of cloud technology
 - OpenStack for virtualization
 - Ceph/RADOS as a BLOB store
 - “Data Mining as a Service”

- ① HEP Computing Model
- ② Node Virtualization
- ③ Scientific Application Delivery
- ④ Volunteer Cloud
- ⑤ Cloud Storage
- ⑥ Summary & Outlook

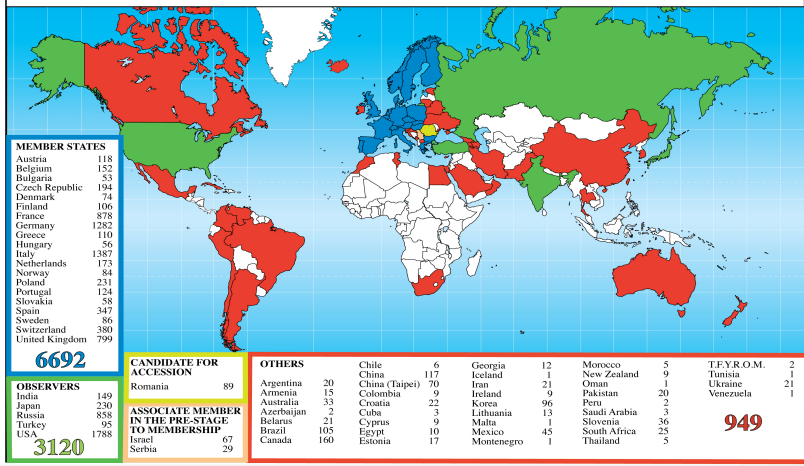
HEP Computing Model

Bird's Eye View of the Large Hadron Collider

- 27 km circumference
- 100 m underground
- 180 MW power consumption



Distribution of All CERN Users by Nation of Institute on 3 September 2012



Distribution of All CERN Users by Nation of Institute on 3 September 2012

- **Heterogeneous** computing landscape
- **Federation** of resources
Cloud: "Infinite Resources"
- **Fair-share** use model
Cloud: Pay-Per-Use

MEMBER STATES

Austria	118
Belgium	152
Bulgaria	53
Czech Republic	194
Denmark	74
Finland	106
France	878
Germany	1282
Greece	110
Hungary	56
Italy	1387
Netherlands	173
Norway	84
Poland	231
Portugal	124
Slovakia	58
Spain	347
Sweden	86
Switzerland	380
United Kingdom	799

6692

OBSERVERS

India	149
Japan	230
Russia	858
Turkey	95
USA	1788

3120

CANDIDATE FOR ACCESSION

Romania	89
---------	----

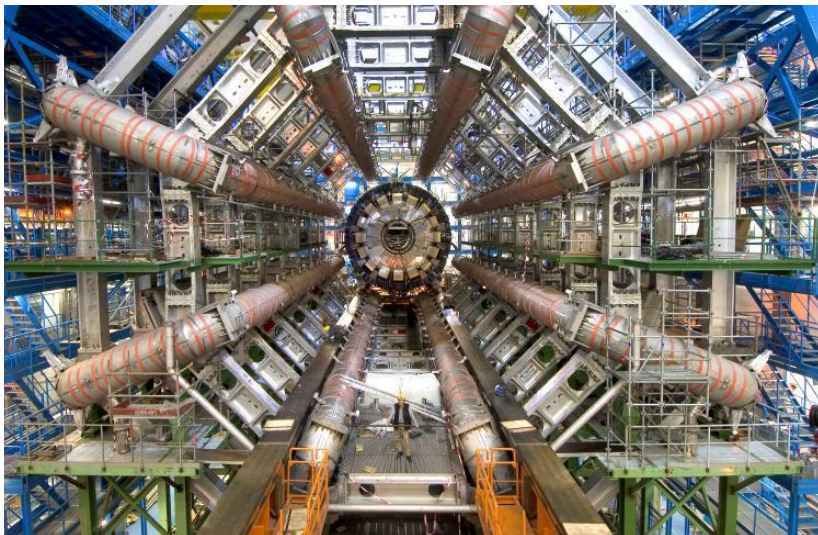
ASSOCIATE MEMBER IN THE PRE-STAGE TO MEMBERSHIP

Israel	67
Serbia	29

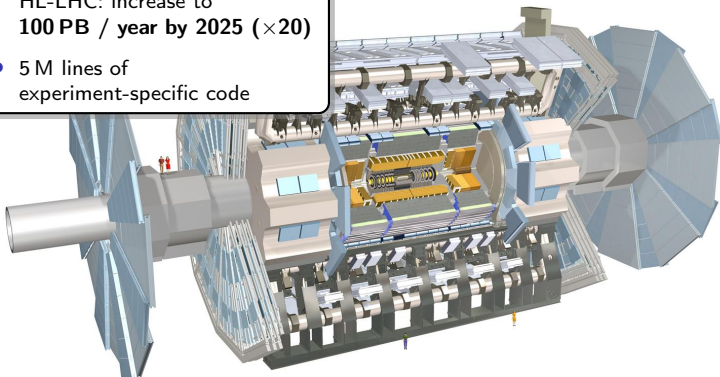
OTHERS

Argentina	20	Chile	6	Georgia	12	Morocco	5	T.F.Y.R.O.M.	2
Armenia	15	China	117	Iceland	1	New Zealand	9	Tunisia	1
Australia	33	China (Taipei)	70	Iran	21	Oman	1	Ukraine	21
Azerbaijan	2	Colombia	9	Ireland	9	Pakistan	20	Venezuela	1
Belarus	21	Croatia	22	Korea	96	Peru	2		
Brazil	105	Cuba	3	Lithuania	13	Saudi Arabia	3		
Canada	160	Cyprus	9	Malta	1	Slovenia	36		
		Egypt	10	Mexico	45	South Africa	25		
		Estonia	17	Montenegro	1	Thailand	5		

949



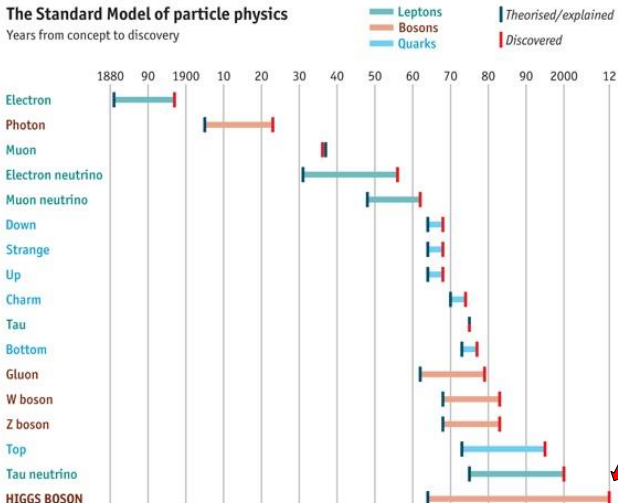
- 100 Million channels,
bunch crossing every 25 ns
- 1 PB/s internal data rate
- 5 PB data / year recorded
(without simulation)
HL-LHC: increase to
100 PB / year by 2025 ($\times 20$)
- 5 M lines of
experiment-specific code



Time Line of the Standard Model

The Standard Model of particle physics

Years from concept to discovery

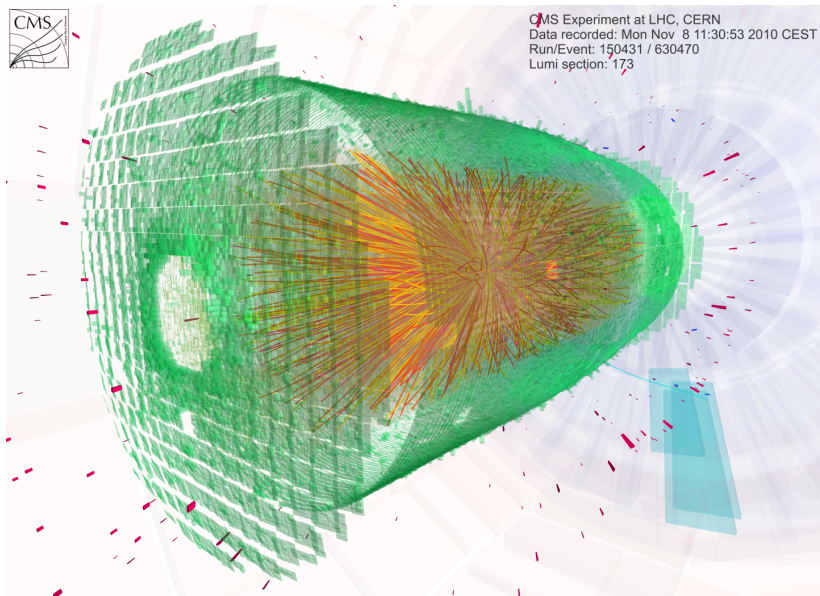


Source: *The Economist*

Is this the last particle to be discovered?



CMS Experiment at LHC, CERN
 Data recorded: Mon Nov 8 11:30:53 2010 CEST
 Run/Event: 150431 / 630470
 Lumi section: 173



- **Raw Event**

Megabytes or tens of megabytes per event

Sensor signals corresponding to a single crossing of beams

- **Event Summary Data (ESD)**

Hundreds of kilobytes per event

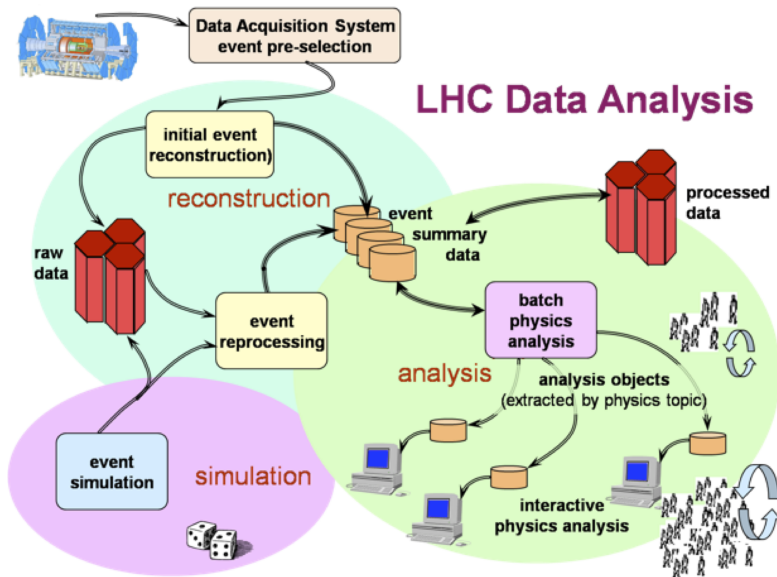
Reconstructed the physics information from digitized sensor signals.
For instance: transformation in global time and space coordinates,
reconstruction of trajectories, determination of particles' charge and
momentum

- **Analysis Object Data (AOD)**

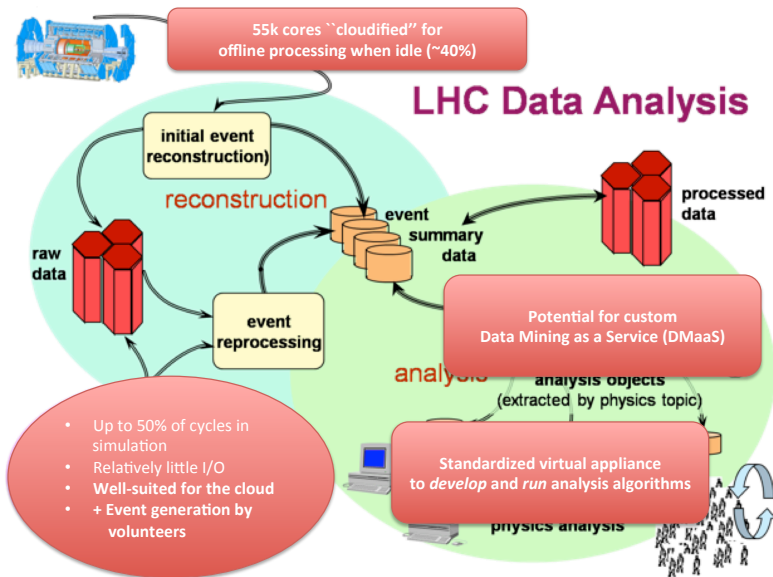
Tens of kilobytes per event

Stripped ESD tailored to specific analysis goals

Events are independent from each other,
embarrassingly parallel workload



Source: Harvey et al.



Distributed Processing of LHC Event Data



Source: WLCG

LHC Grid

- ~170 sites
- 10 – 10 000 servers / site
- 300 PB under management
- Commodity hardware

→ many small(ish) sites

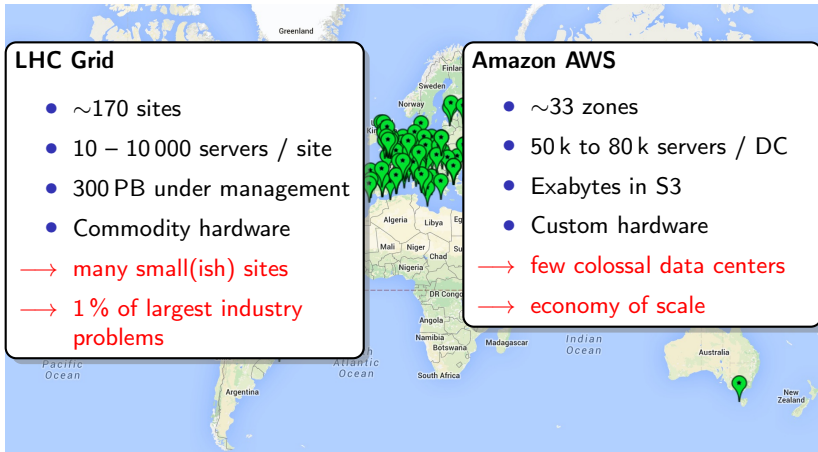
→ 1 % of largest industry problems

Amazon AWS

- ~33 zones
- 50 k to 80 k servers / DC
- Exabytes in S3
- Custom hardware

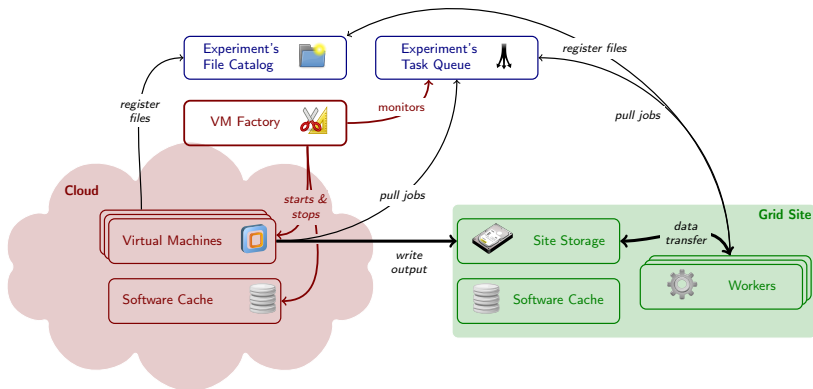
→ few colossal data centers

→ economy of scale



Node Virtualization

Cloud Resources in a Global Batch System



Different approaches to "VM Factory":

► CloudScheduler

► VAC

► WNoDeS

- Good VM performance requires careful configuration of hypervisor and guest
- High-performance I/O devices (e. g. 10 Gbit NIC, GPGPU, Xeon Phi, . . .) only usable with direct assignment to VMs
- Just beginning to explore non x86_64 architectures (ARM64, POWER8)

Performance engineering on CERN OpenStack:

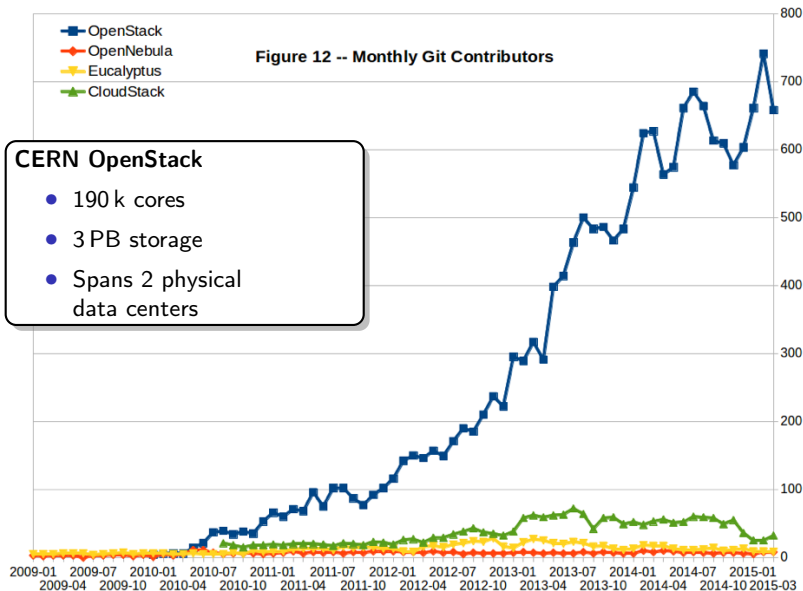
VM sizes (cores)	Before	After
4x 8	7.8%	3.3% (batch WN)
2x 16	16%	4.6% (batch WN)
1x 24	20%	5.0% (batch WN)
1x 32	20.4%	3-6% (bare SLC6 . . . batch WN)

► HEPiX'15

Source: Wiebalck et al.

Containers (e. g. Docker) provide isolated namespaces on top of the same kernel and have almost native performance.

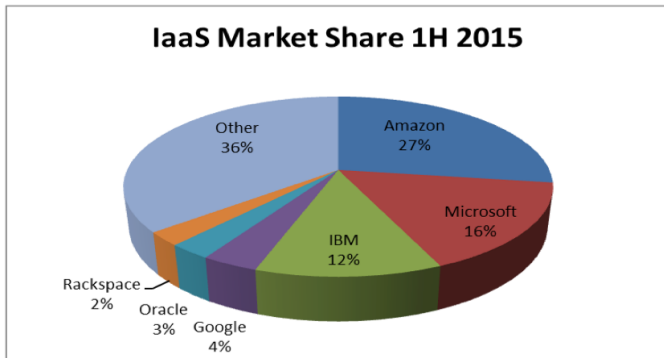
Figure 12 -- Monthly Git Contributors



CERN OpenStack

- 190 k cores
- 3 PB storage
- Spans 2 physical data centers

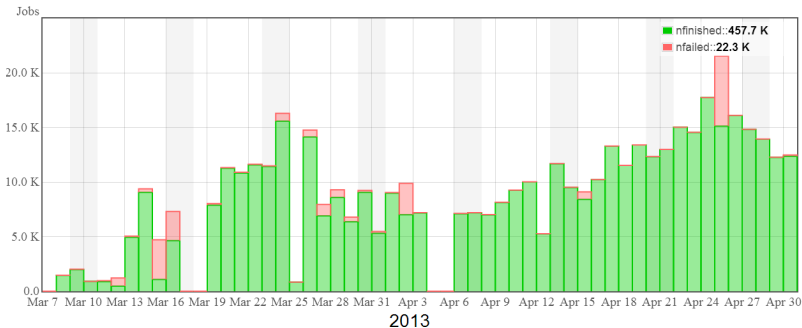
- Amazon (EC2), Microsoft (Azure), IBM (OpenStack,...), Google Cloud Platform, Rackspace (OpenStack/Azure/VMware), ...



- Ralph Finos, <http://wikibon.com/public-cloud-market-shares-2014-and-2015/>

Source: McNab

Failed and Finished Jobs



- ◆ Most of the job failures occurred during start up and scale up phase – as expected
- ◆ Reached throughput of 15k jobs per day

Source: Panitkin

► Google I/O

Brookhaven National Lab / Amazon EC2

► HEPiX'15

- 100 000 cores
- Amazon spot market biddings

CERN / Evaluation of Various Commercial Clouds

► HNSciCloud

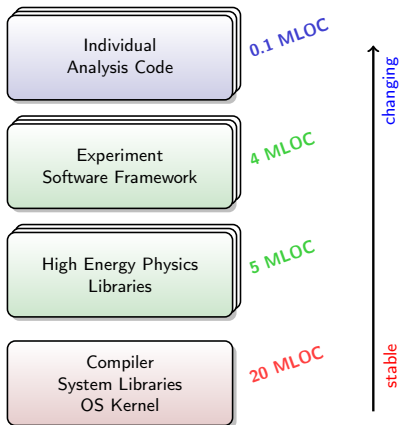
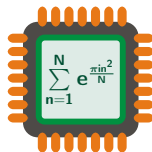
- HelixNebula Science Cloud:
EU funded project for federated clouds
- Current tests at a few thousand cores but ramping up
- ~5 k cores allocated on MS Azure
- A few thousand cores on Deutsche Börse Cloud Exchange

Lessons Learned

- Cloud resources require supporting services besides the VMs (caches, proxies)
- We need to figure out our cost of computing
- Cloud resources need to be continuously benchmarked
- More layers of complexity that can cause operational issues

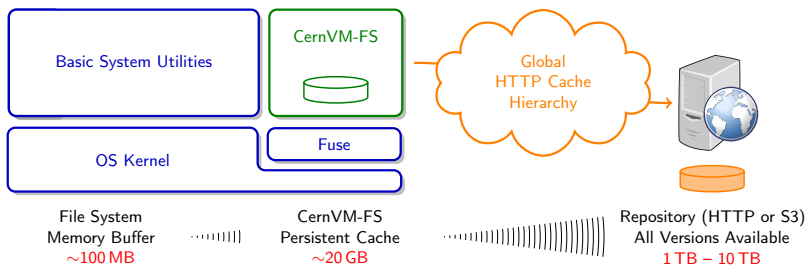
Application Delivery

```
> cmsRun DiPhoton_Analysis_cfg.py
```



Scientific Software is Special

- Frequent Updates
- Includes its dependencies
- Released versions need to stay available
- Quick turn-around desired from source code commit to deployment

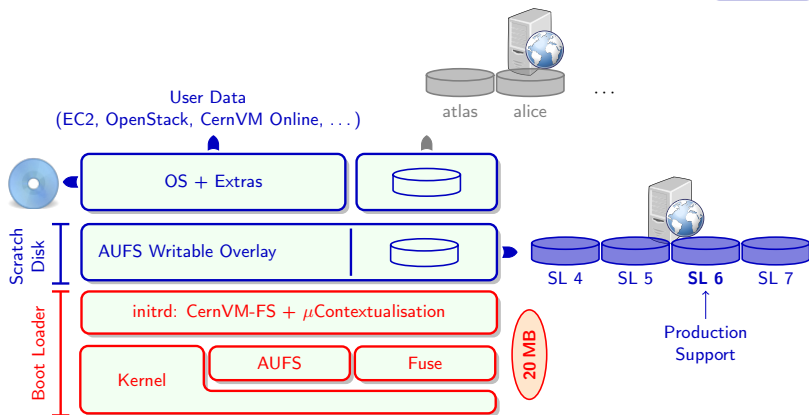


- Independent *repositories*, e. g. `/cvmfs/atlas.cern.ch`
- Single point of publishing
- HTTP Transport, access and caching on demand
- > 100 Million files, ~100,000 deployed clients (Grid, Cloud, HPC)

CernVM: A Tiny Yet Complete Appliance

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

► CHEP'13



~ 30 000 CernVM booted per day



CernVM Hypervisor Support Status

The success of CernVM is largely based on the fact that it runs in practically **all** cloud environments.

Hypervisor / Cloud Controller	Status
VirtualBox	✓
VMware	✓
KVM	✓
Xen	✓
Microsoft Hyper-V	✓
Vagrant	✓
OpenStack	✓
OpenNebula	✓
CloudStack	✓
Amazon EC2	✓
Google Compute Engine	✓
Microsoft Azure	✓
Docker	✓ ¹

¹ Release Candidate



**CernVM: complete and portable environment
for **developing** and **running** HEP data processing tasks**

**CernVM: complete and portable environment
for **developing** and **running** HEP data processing tasks**

Use Cases

- 1 **IaaS Clouds**
- 2 Development Environment
- 3 Volunteer Computing
- 4 Long-Term Analysis Preservation
- 5 Outreach & Education

Infrastructure-as-a-Service Cloud

Various clouds:

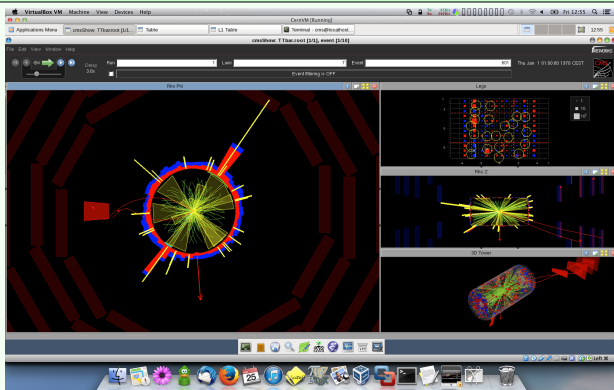
- ATLAS high-level trigger farm
- Cloud resources seamlessly integrated with experiment task queues (e. g. ATLAS CloudScheduler, LHCb VAC)
- Commercial providers
- ALICE software release testing on CERN OpenStack
- ...

**CernVM: complete and portable environment
for **developing** and **running** HEP data processing tasks**

Use Cases

- 1 IaaS Clouds
- 2 **Development Environment**
- 3 Volunteer Computing
- 4 Long-Term Analysis Preservation
- 5 Outreach & Education

Interactive Users: Linux Event Viewer on Mac



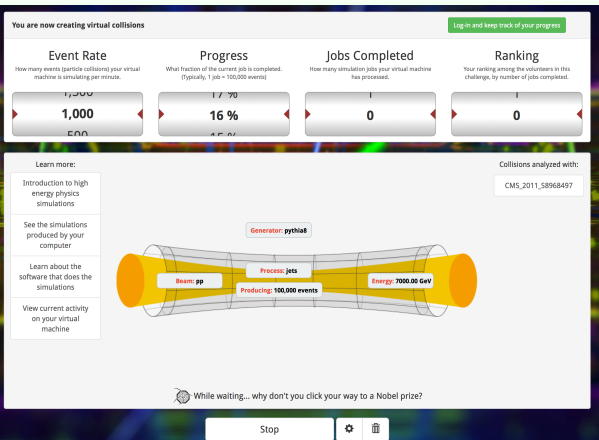
~300 users in 2015, out of which ~65 external users

**CernVM: complete and portable environment
for **developing** and **running** HEP data processing tasks**

Use Cases

- 1 IaaS Clouds
- 2 Development Environment
- 3 **Volunteer Computing**
- 4 Long-Term Analysis Preservation
- 5 Outreach & Education

Test4Theory and Computing Challenge



**CernVM: complete and portable environment
for **developing** and **running** HEP data processing tasks**

Use Cases

- 1 IaaS Clouds
- 2 Development Environment
- 3 Volunteer Computing
- 4 Long-Term Analysis Preservation
- 5 Outreach & Education

ALEPH (1989-2000) software in CernVM

```
jacob — aleph@cernvm-aleph01:~/test/ALPHA — ssh — 66x18
pb-d-128-141-134-74:~ jacob$ 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.
```

Demonstrates that VMs + Containers can bridge 10+ years

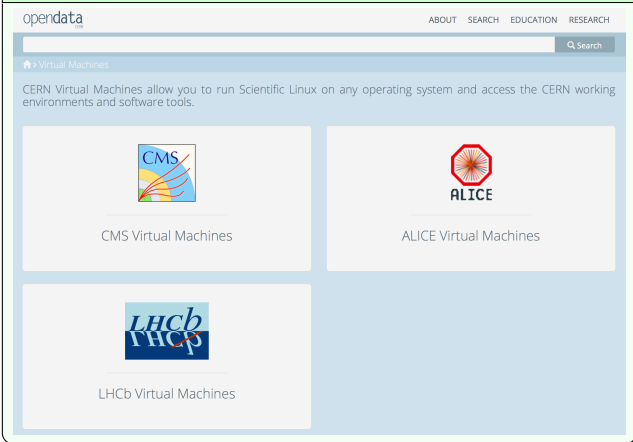
► DPHEP Status Report

**CernVM: complete and portable environment
for **developing** and **running** HEP data processing tasks**

Use Cases

- 1 IaaS Clouds
- 2 Development Environment
- 3 Volunteer Computing
- 4 Long-Term Analysis Preservation
- 5 Outreach & Education

CERN OpenData Portal, CERN@School



The screenshot shows the CERN OpenData Portal interface. At the top, there is a navigation bar with links for ABOUT, SEARCH, EDUCATION, and RESEARCH. Below this is a search bar with a magnifying glass icon and the word "Search". The main content area is titled "Virtual Machines" and includes a sub-header "CERN Virtual Machines allow you to run Scientific Linux on any operating system and access the CERN working environments and software tools." Below this, there are three large white boxes, each containing a logo and the name of a virtual machine environment: CMS Virtual Machines (with the CMS logo), ALICE Virtual Machines (with the ALICE logo), and LHCb Virtual Machines (with the LHCb logo).

Volunteer Computing

- Commonly known as ... @Home projects
- Computing on *spare* resources of “interested citizens”
- Opportunistic, volatile resources
- Big projects comparable to a **TOP500 supercomputer**

- Outreach program
- Requires **curation of volunteers**

SETI@Home

- Search for extraterrestrial life
- 130 000 active participants
- 630 Tera-FLOPS



Einstein@Home

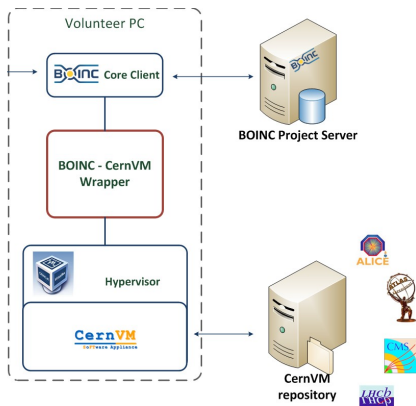
- Search for gravitational waves
- 34 000 active participants
- 470 Tera-FLOPS



LHC@Home 2.0

▶ ACAT'12

- Monte-Carlo simulations, parameter tuning
- **First BOINC project using virtual machines, avoids application porting**



Source: Harutyunyan

LHC@Home 2.0

- 600–700 VMs connected
- **> 2 trillion** events created

ATLAS@Home

- 10 k running jobs
- 2nd biggest simulation site

CMS@Home and Beauty@Home
under development

Cloud Storage



Thoughts behind Key-Value and BLOB Stores

~2000–2010

SQL DBs and file systems do not scale to the needs of large web services

- 1 Can we scale out storage (“horizontal scaling”) if we give up on the relational data model and the file system hierarchy?
 - Dictionary as a simple, easy to distribute yet useful data structure

Scaling



SQL DBs and file systems do not scale to the needs of large web services

- 1 Can we scale out storage (“horizontal scaling”) if we give up on the relational data model and the file system hierarchy?
 - Dictionary as a simple, easy to distribute yet useful data structure
- 2 In the presence of inevitable faults, can we still make progress (“availability”)?:
 - Boosted by Amazon Dynamo paper ▶ SOSP'07
 - Idea of “eventual consistency”:
heal data inconsistency once the system recovers from faults

Scaling

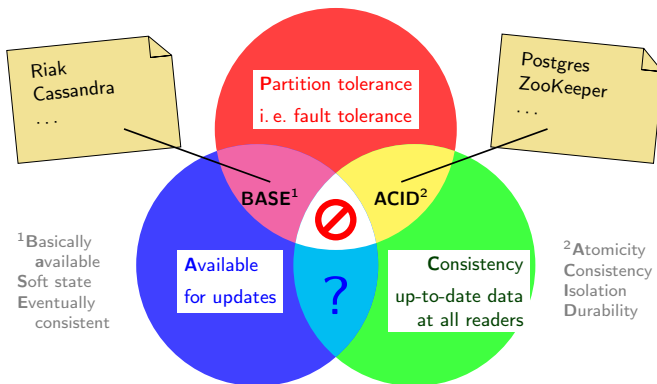
Fault-Tolerance

Even though the *interface* is simple, the *implementation* of a distributed key-value or BLOB store is highly non-trivial.

A distributed storage system can have at most two out of three desirable properties

► Brewer '97

► Brewer '12

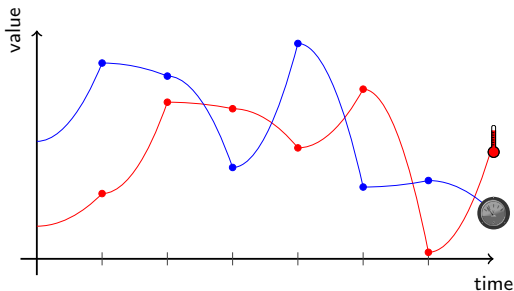


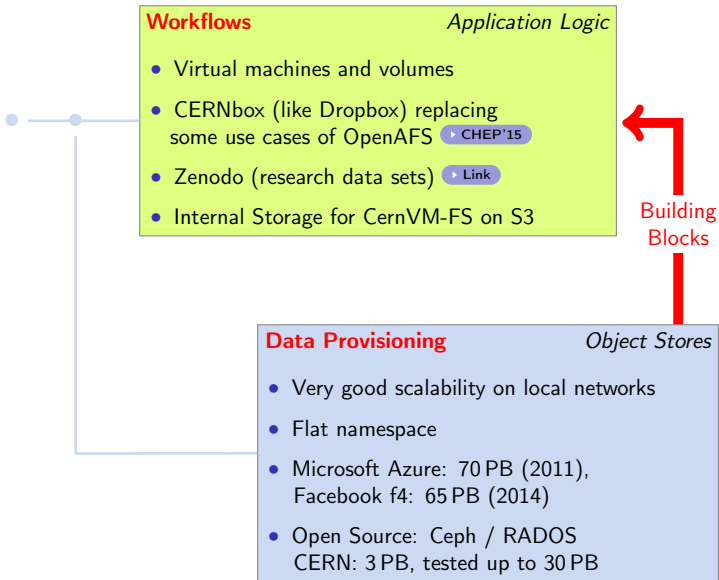
The tradeoffs between availability and consistency can be granular and subtle
For instance: a disconnected ATM might still allow small withdrawals

- CMS evaluated
Cassandra, MongoDB, Riak
▶ CHEP'15
- ALICE evaluated
RAMCloud and Riak
▶ Note

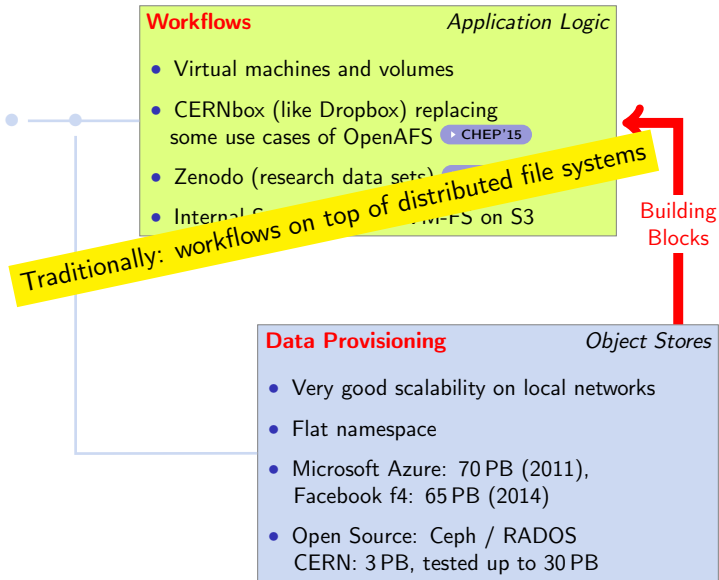
Conditions Data

- State of the detector at the time of data taking
- Small (terabytes / year) but critical to reconstruct data
- Traditionally kept in SQL databases

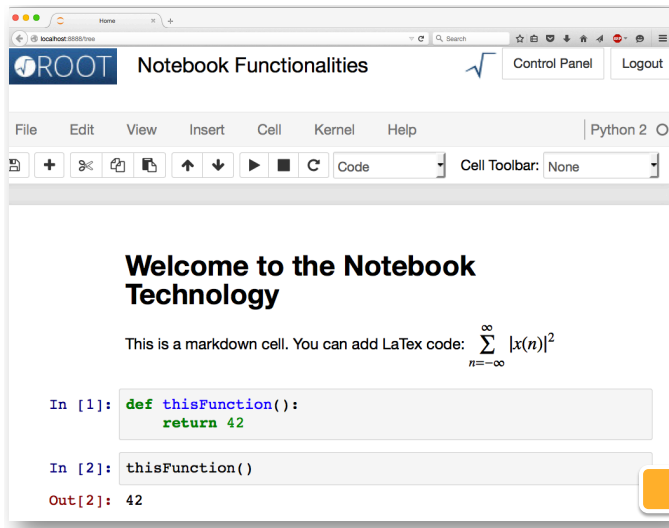




Workflows on Top of BLOB Storage



Summary & Outlook



ROOT Notebook Functionalities

Control Panel Logout

File Edit View Insert Cell Kernel Help Python 2

Cell Toolbar: None

Welcome to the Notebook Technology

This is a markdown cell. You can add LaTeX code: $\sum_{n=-\infty}^{\infty} |x(n)|^2$

```
In [1]: def thisFunction():
        return 42
```

```
In [2]: thisFunction()
```

Out[2]: 42

Code

Source: Piparo

The Future: Custom Cloud Apps?

In a browser

Text and
Formulas

Code

Shell Commands

Images

```
function():  
42
```

```
In [2]: thisFunction()
```

```
Out[2]: 42
```

```
> aaasdemo.web.cern.ch/rootaaasdemo/SaaSFee.jpg \  
> SF.jpg
```

% Total Time	% Received Time	% Xferd Current	Average Speed	Time
Spent	Left	Speed		
100	128k	100	128k	0
--:--:--	--:--:--	2787k		

```
In [4]: from IPython.display import Image  
Image(filename="./SF.jpg",width=225)
```

```
Out[4]:
```



The Future: Custom Cloud Apps?

```
In [1]: def thisFunction():
        return 42
```

```
In [2]: thisFunction()
```

```
Out[2]: 42
```

```
In [3]: %%bash
        curl rootaasdemo
        > SF.jpg
```

% Total	% R	Time	Time
Spent	Left		
100	128k	100	128k
--:--:-- --:--:-- 2787k			

```
In [4]: from IPython.display import Image
        Image(filename="./SF.jpg",width=225)
```

```
Out[4]:
```



- Notebook web interface
- Virtualized CPU resources
- Scientific software from CernVM-FS
- BLOB store backed data access

► CS3

Cloud Computing Resources in HEP

- ~5 % of resources provided by virtual machines, mostly for simulation
- Cost of commercial VMs almost on par with in-house cost
- Volunteer clouds: great opportunity for outreach

Cloud Storage Resources in HEP

- Storage: more expensive commercially because of data transfer costs
- Scalable key-value interface requires more application logic
- Very promising open source components: Ceph, Riak, RAMCloud, ...

Did utility computing arrive?

- ① On the infrastructure level:
x86_64 CPUs and *unstructured* storage are commodity
- ② HEP still has unique and specialised needs and applications:
 - Wide-area data management
 - Custom, optimized I/O format
 - Complex scientific software stacks
 - Legacy software
- ③ Our custom distributed systems must support our custom application workflows
- ④ We have today **much better building blocks** available than at the time when we designed the grid

