

Long term data preservation and virtualization

- Why preserving data? What to do with that old data?
- Possible usages of virtualization
- Some work done



Study Group for Data Preservation and
Long Term Analysis in High Energy Physics

Yves Kemp, DESY IT

2nd Workshop on adapting applications and
computing services to multi-core and virtualization
CERN, 22.6.2010

Inter-experiment Study Group on Data Preservation (DPHEP)

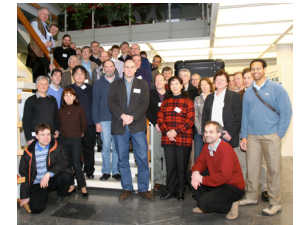
DPHEP-2009-001
July 31, 2009

www.dphep.org

Data Preservation in High-Energy Physics

Study Group for Data Preservation and Long-Term Analysis in High-Energy Physics

<http://dphep.org>



- **ICFA subgroup since 08/2009**

- Presented at HEPAP (DOE/NSF), FALC, ICFA
- Intermediate document released in November 2009

- **Next Steps:**

- Next workshop **KEK, July 8-10, 2010**
- New labs (**JLAB**), extension to other experiments (neutrino etc.)
- Investigate concrete models, technical proposal 2010
- Defined Resources, progress towards a common infrastructure

3 Workshops in 2009



Abstract
Data from high-energy physics (HEP) experiments are collected with significant financial and human effort and are mostly unique. At the same time, HEP has no coherent strategy for data preservation and re-use. An inter-experimental Study Group on HEP data preservation and long-term analysis was convened at the end of 2008 and held two workshops, at DESY (January 2009) and SLAC (May 2009). This document is an intermediate report to the International Committee for Future Accelerators (ICFA) of the reflections of this Study Group.

WIRED SUBSCRIBE SECTIONS BLOGS REVIEWS VIDEO HD Sign in / RSS Feeds

beyond the beyond

Dead Media Beat: High Energy Physics data

By Bruce Sterling | March 8, 2010 | 5:44 pm | Categories: Uncategorized

*When the data firehose stops gushing...

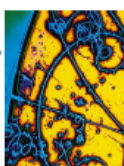
20 | WISSEN & BILDUNG

NACHRICHTEN Therapeutische Impfung gegen Leberkrebs denkbar

Thüringer Forscher entwickeln eine Impfung gegen Leberkrebs. Sie soll das Immunsystem auf die gefährlichen Zellen aufmerksam machen und so die Ausbreitung des Krebses verhindern. Erste Tests sind geplant.

Hieroglyphen im Teilchenlabor

Alle Forschungsergebnisse gehen verloren. Nichts davon ist im Teilchenlabor. Die Forscher suchen nach einer Möglichkeit, die Daten zu sichern. Erste Tests sind geplant.



Yves Kemp | Long Term Data Preservation

symmetry

dimensions of particle physics



Table of Contents volume 06 issue 06 december 09

Preserving the data harvest

Canning, pickling, drying, freezing—physicists wish there were an easy way to preserve their hard-won data so future generations of scientists, armed with more powerful tools, can take advantage of it. They've launched an international search for solutions.

By Nicholas Bock

The 2010 HEP Landscape (Colliders)

- > e^+e^- : LEP ended in 2000
 - No follow-up decided (ILC?) - after 2020
- > e^+p : HERA end of collisions at HERA in 2007
 - No follow-up decided (LHeC?) - after 2020
- > B-factories: BaBar ended in 2008, Belle → Belle II
 - Next generation in a few years (2013-2017)
- > pp: Tevatron ends soon (in 2011?)
 - The majority of the physics program will be taken over at the LHC
 - However: p-pbar is unique, no follow-up foreseen

*"LEP is scheduled to be dismantled soon so that its 27 km tunnel can become the home for the ambitious LHC proton collider, which is due to come into operation in 2005."
[CERN Courier, Dec. 1st, 2000]*

Data taking at HEP experiments takes 15-20 years, and some data are unique

- **What is the fate of the collected data?** (where "data" means the full experimental information..)



Data is needed:

- > My personal story: Around 2008, a retired professor asked the DESY data management group:
- > “Around 1975*, we had tapes from a bubble chamber experiment at the Computing Center. Are these still available, maybe copied to other media? I got a request from CERN concerning these tapes.”
- > (we did not have them...)
- > Honestly, I thought, no one would need these data anymore
 - New experiments, higher energy, better resolution, ...
 - No one able to read / understand the (scientific content of the) data
- > **First lesson learned here: Preserve data, and preserve the ability to perform some kind of meaningful operation on it.**



What is “meaningful operation”?

- Long-term completion and extension of scientific programs
 - Allow “late analyses” to be done: +5-10% more publications
 - “Late analyses” benefit from full statistics, best understanding of systematics

- Cross-collaboration analyses

- Often performed at the end of lifetime of collaborations
- Even among generations of experiments

- Re-use of the data

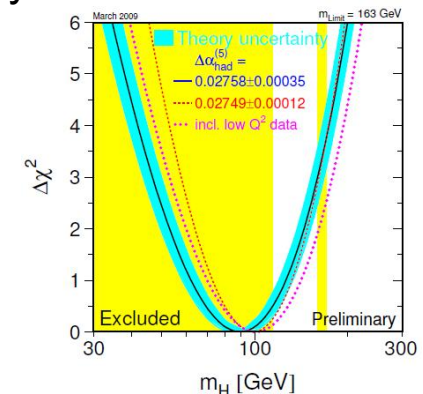
- Re-analyze the data with new theoretical models, new analysis techniques, ...

- Education, training and outreach

- E.g. analysis by students without restrictions (like collaboration membership...)

- Different goals – different solutions

- Both for the “data archival” and the “long term analysis” part



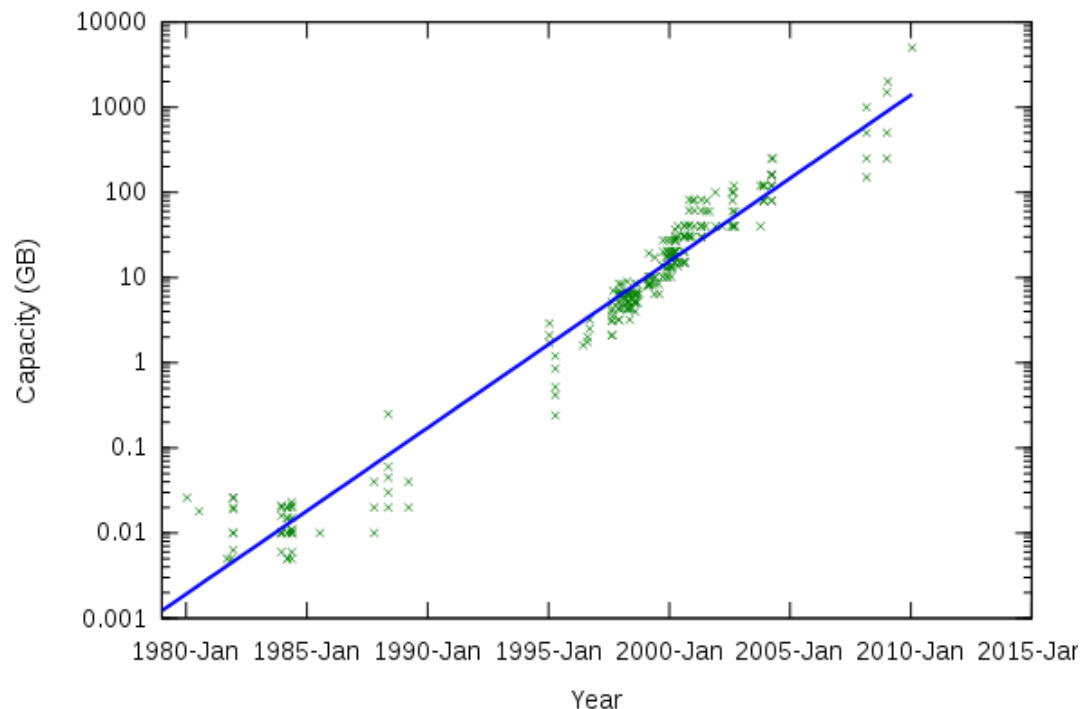
My first and very naïve ansatz

- OK, why don't we just put everything in a virtual machine?
 - Data archival is done elsewhere, just need "to plug that into the VM"
 - Your VM contains everything you need to develop and run code and analysis
- The problem would then be reduced to maintain virtual images, and maintain their ability to run. In the Cloud era, seems like a trivial task
- Problems: Everything in IT is a moving target:
 - Will your network always be the same?
 - Will your access protocol always be the same?
 - Are you sure you do not need new software (e.g. MC generators) that require a new OS?
 - Are you sure your i386/SL4 VM will produce the same results when emulated on a quantum computer in NN years?
 - What about services you need, like CondDB,...
- Naïve virtualization will not work... but still, virtualization can help



The Data Access Problem

- Data access for running or very recently finished experiments is huge
 - You know the LHC numbers. H1 e.g. (finished 2007) has ~1 PB of data
 - Need complicated systems to store and access data
- This problem will remain until you can “put all data in one machine”
 - For 1 PB of data, this could happen in 2025: All data fits in one hard drive (if you believe in the plot below)



Hard drive capacity over time
Wikimedia Commons



Some complex dependencies

Personal Analysis Code

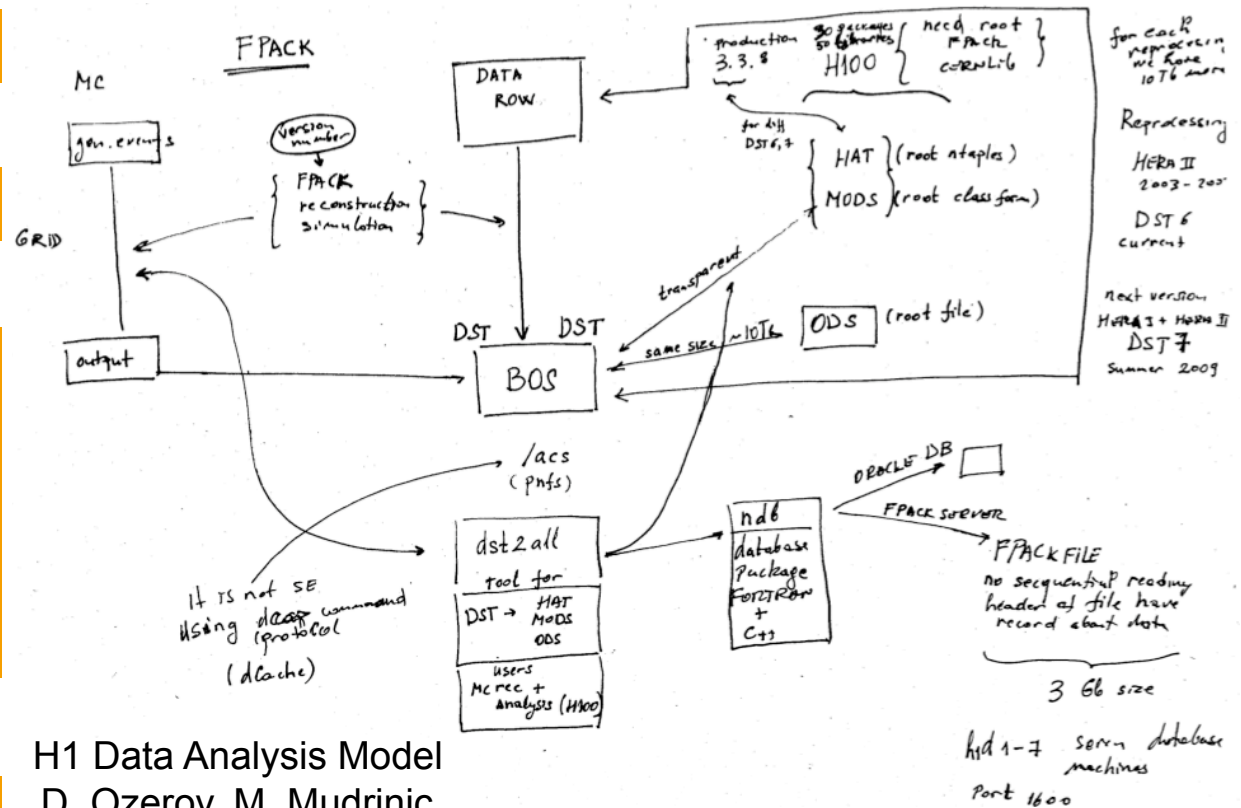
Experiment Framework

External Dependencies

- ROOT - Geant, Pythia, ...
- Compiler - DB (Oracle)
- Grid / Storage Middleware
- OS & Libs

Hardware Architecture

- x86 - x86-64
- Power
- Quantum Computers
- ...



- > Will we be able to reproduce the same results after many years of not looking at the data?
 - Two hypothetical antipodes involving Virtualization



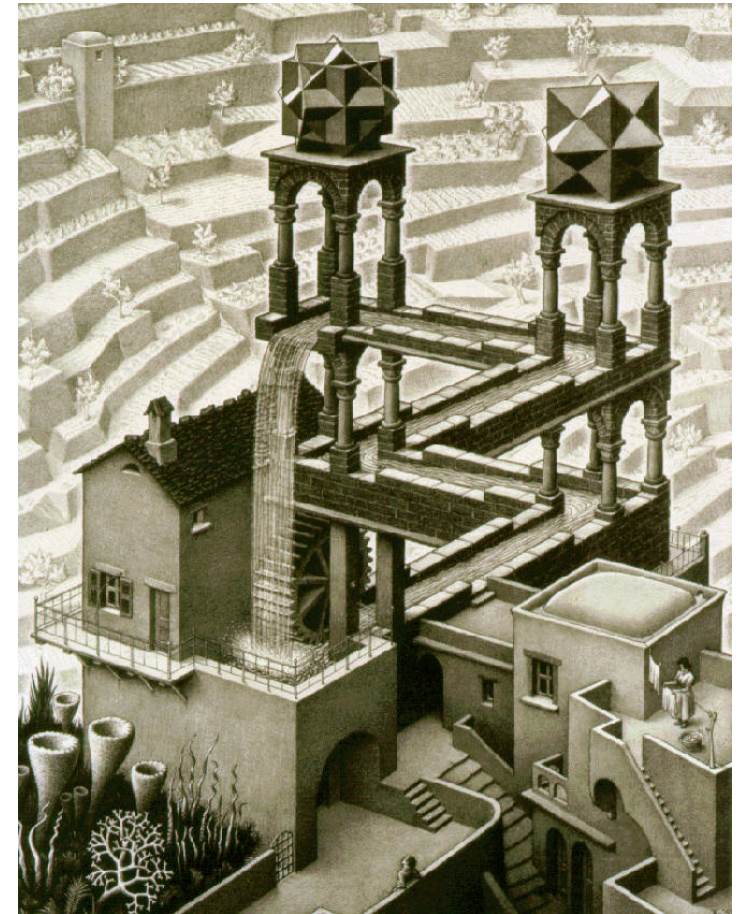
Scenario 1: "Freezing"

- At the end of the experiment:
 - Datasets closed, final reprocessing done
 - Software framework stable
- Virtual image of the OS with software is done
 - Important: Use a standardized format, like OVF
- Necessary services like Cond DB.:
 - Either integrated into images
 - Or also frozen into another image
- Data access:
 - Either maintain the old protocol/interface
 - Or use high-level protocols
- Running analysis in 20NN (with $NN \gg 10$):
 - Start the whole ensemble of VMs



Scenario 2: Continuous test-driven migration

- Start during running experiment
 - Or even before, when designing software framework
- Define tests
 - In the beginning on MC data, later real data
 - Certain code, running on certain data, yields certain result (e.g. $M_{\text{top}}=172.4 \text{ GeV}/c^2$)
- Have an automated machinery, which regularly compiles code for different OS / architectures, and runs the tests
- If test fails (e.g. compilation or execution fails, or result divergent)
 - Manual intervention: understand (and fix) problem
- ➔ Such automated tests are usually performed using virtualization techniques and workflows



M.C. Escher's "Waterfall"
(c) 2009 The M.C. Escher Company - theNetherlands.
All rights reserved. Used by permission. www.mcescher.com

Discussion “Freezing” / “test driven migration”

> Pro Freezing

- One-time effort, very small maintenance outside of analysis phase
- Also allows software w/o code (but might fail with DRM / licensing issues)

> Pro Test-driven migration

- Usability and correctness of code is guaranteed at every moment
- Data accessibility and integrity can be checked as well
- Fast reaction to standard/protocol changes
- General code quality can improve, as designed for portability and migration

> Cons Freezing

- Rely on certain standards and protocols that may evolve
- Potential performance problems

> Cons Test-driven migration

- Needs long-time intervention, more man-power and resources needed
- Some knowledge of the frameworks must be passed to maintainers



The BaBar Archival System

Lifetime:

It will take ~8 years for SuperB to be approved by the governments and funding agencies, constructed, commissioned and obtain a dataset more significant than BaBar's.

- 2012 to 2018 is the **minimum** required existence of the BaBar archival system

There may also be a need to validate initial results against those of BaBar and Belle₁₂

Status:

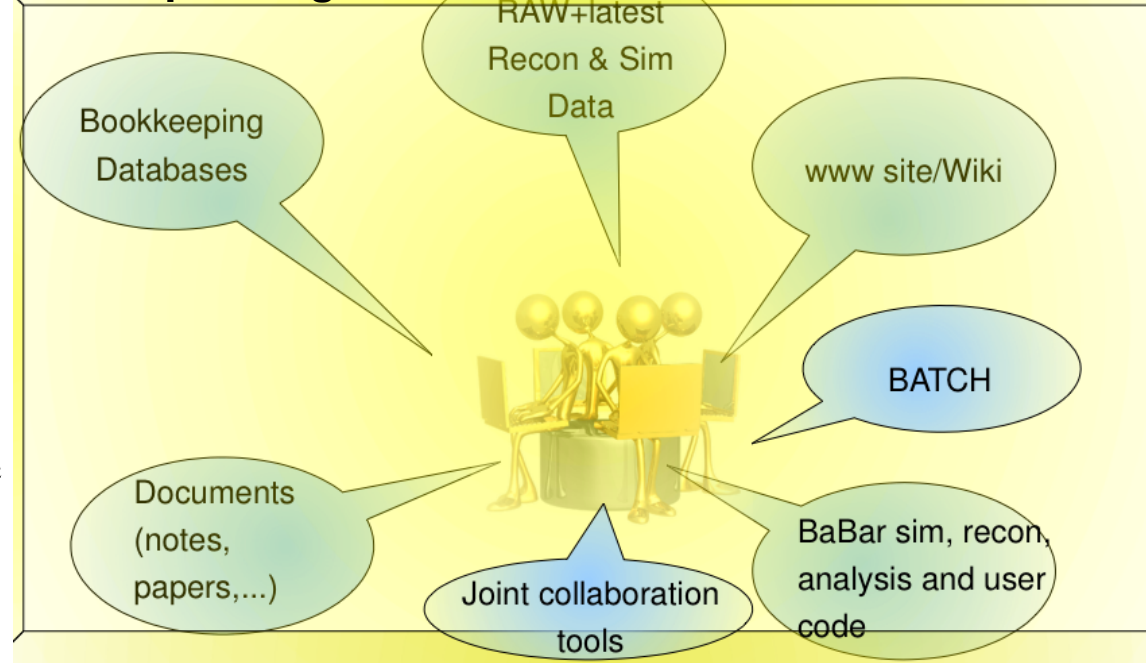
- Currently working on selecting prototype archival system hardware



32-cores, enough storage for micro data needed by most analyses

- Set it up and get testers to identify faults

Concept using Virtualization Containers:



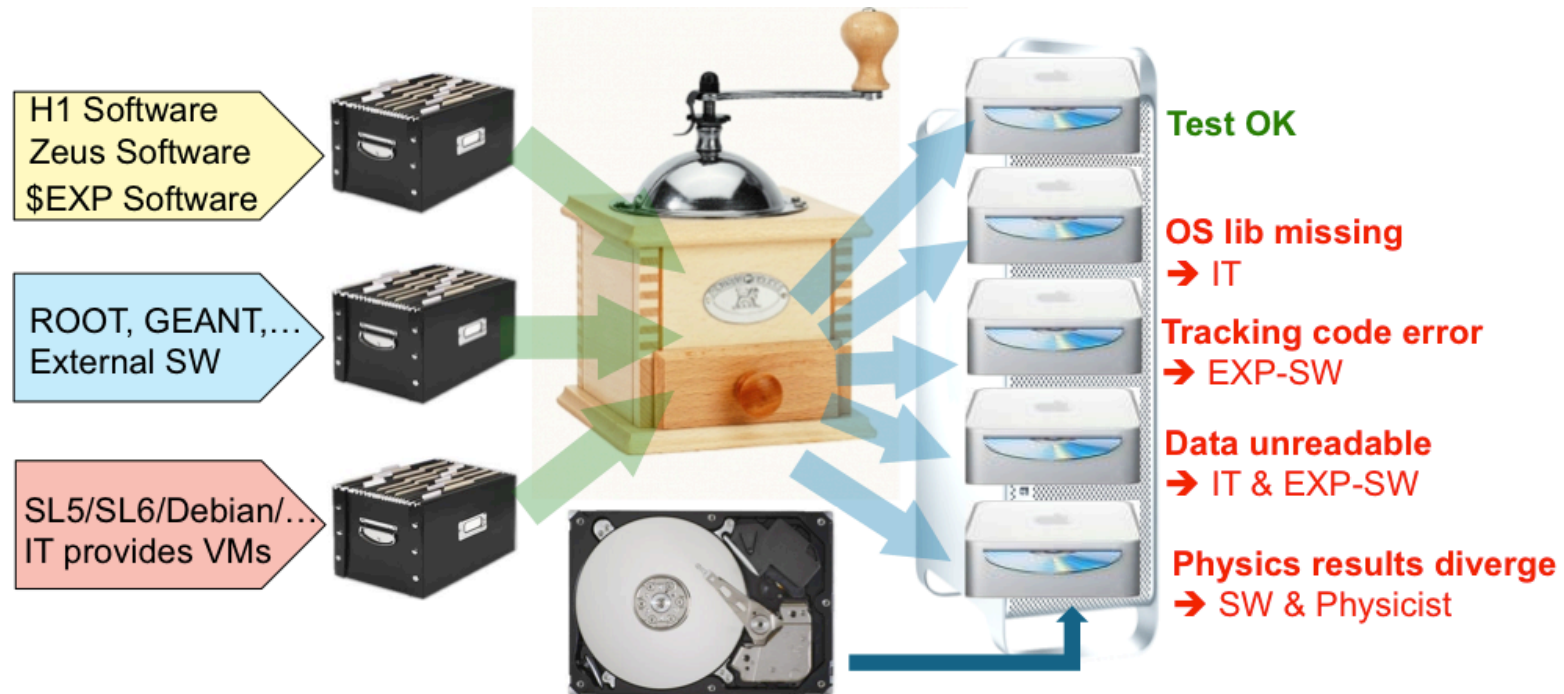
Slides from Homer Neal,
@CERN 8.12.2009

Remark: BaBar has migrated to SL5
Before “going virtual”



DESY: H1, Zeus & IT development

- > Automatically test SW & data
- > Migration helper (as automatic as possible)
- > Status: 5% mockup ready for evaluating man-power and needs
- > Generic effort: Open to other experiments



Clear separation between providers of input.

Automated VM image generator provided centrally.

Tests defined by \$EXP.
Test data store provided by IT.

Different VMs run SW and tests.
Depending on results, different action needed.



Combination of the two models

> Possible scenario:

- > Start with migrating software to most up-to-date OS. Best start this already during running experiment
 - ① Use an automated machinery for testing and migrating your software
 - ... as long as your data does not fit into one machine
 - ② If data fits in one machine: Freeze everything using the most up-to-date OS.
 - And preserve this VM
-
- > What about analysis?
 - **During phase 1:** Easy: You always have a living system, can add current SW, might need some manpower to do the large scale analysis, but success is guaranteed
 - **During phase 2:** New code: Might be difficult to incorporate in an (then old) VM
 - **During phase 2:** Reproduce results: Easy: Problem has become “small”, and is reduced to running an ensemble of VM

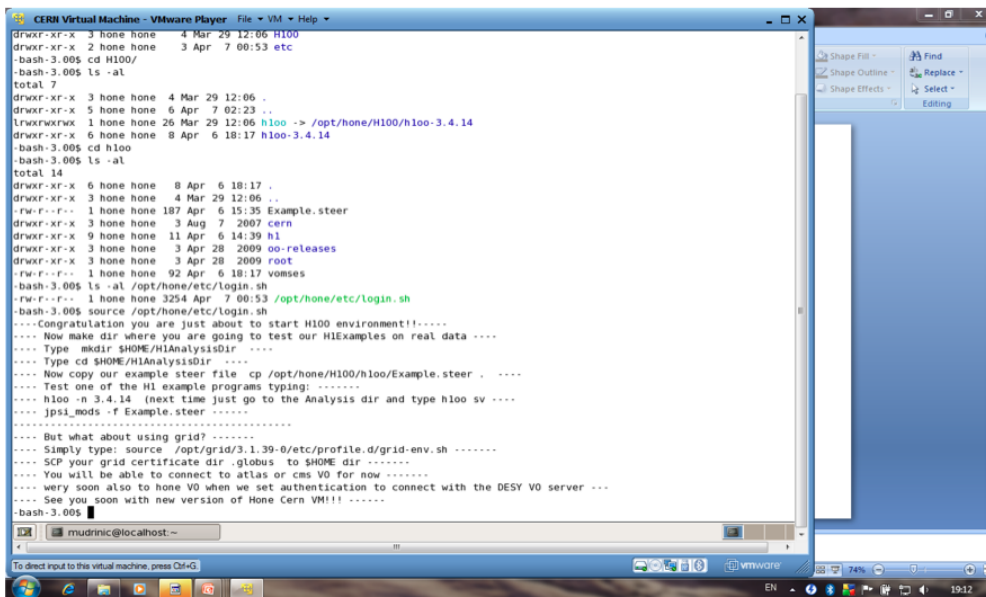


Acceptance problem of virtualization

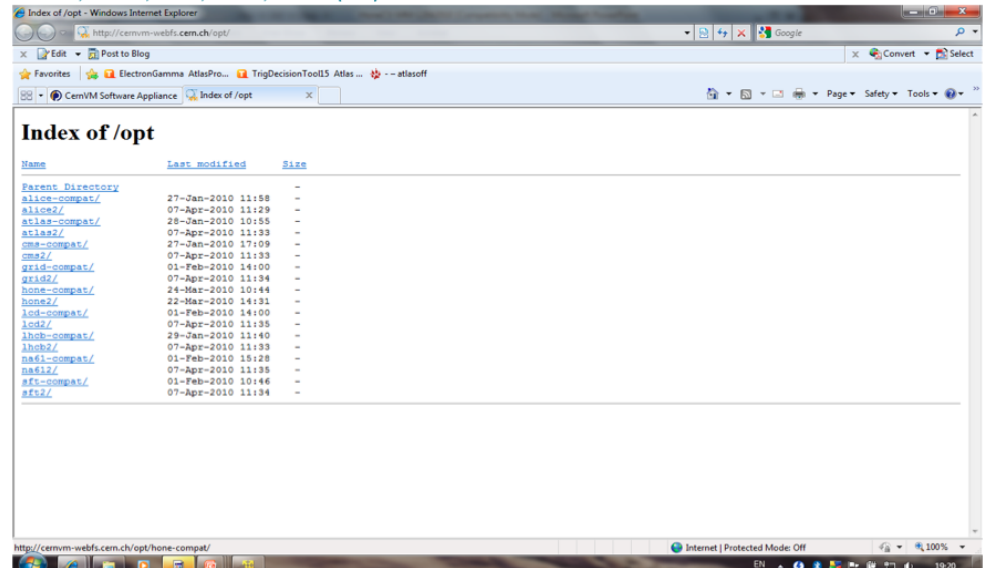
- Physicists are used to hardware. Virtualization is new, kind of a black box.
- Virtualization as a means of preserving ability to analyze precious data might have acceptance problems
- CernVM already now make VMs a reality for physicists
- At DESY, H1 (Mihajlo Mudrinic) has started some efforts together with Predrag:

At the beginning user only need to type:: `source /opt/hone/etc/login.sh`

Software Repository on the File Server:: <http://cernvm-webfs.cern.ch>
ATLAS,ALICE,CMS,LHCb,HONE(H1)



```
CERN Virtual Machine - VMware Player
drwxr-xr-x 3 hone hone 4 Mar 29 12:06 H100
drwxr-xr-x 2 hone hone 3 Apr 7 00:53 etc
-bash-3.005 cd H100/
total 7
drwxr-xr-x 3 hone hone 4 Mar 29 12:06 .
drwxr-xr-x 5 hone hone 6 Apr 7 02:23 ..
lrwxrwxrwx 1 hone hone 26 Mar 29 12:06 h100 -> /opt/hone/H100/h100-3.4.14
drwxr-xr-x 6 hone hone 8 Apr 6 18:17 h100-3.4.14
-bash-3.005 cd h100
-bash-3.005 ls -al
total 14
drwxr-xr-x 6 hone hone 8 Apr 6 18:17 .
drwxr-xr-x 3 hone hone 4 Mar 29 12:06 ..
-rw-r--r-- 1 hone hone 187 Apr 6 15:35 Example.steer
drwxr-xr-x 3 hone hone 3 Aug 7 2007 cern
drwxr-xr-x 9 hone hone 11 Apr 6 14:39 h1
drwxr-xr-x 3 hone hone 3 Apr 28 2009 oo-releases
drwxr-xr-x 3 hone hone 3 Apr 28 2009 root
-rw-r--r-- 1 hone hone 92 Apr 6 18:17 voemes
-bash-3.005 ls -al /opt/hone/etc/login.sh
-rw-r--r-- 1 hone hone 3254 Apr 7 00:53 /opt/hone/etc/login.sh
-bash-3.005 source /opt/hone/etc/login.sh
.... Congratulations you are just about to start our H1Examples on real data ....
.... Now make dir where you are going to test our H1Examples on real data ....
.... Type mkdir $HOME/H1AnalysisDir ....
.... Type cd $HOME/H1AnalysisDir ....
.... Now copy our example steer file cp /opt/hone/H100/h100/Example.steer . ....
.... Test one of the H1 example programs typing: ....
.... h100 -n 3.4.14 (next time just go to the Analysis dir and type h100 sv ....
.... jpsi_mods -f Example.steer ....
.....
.... But what about using grid? ....
.... Simply type: source /opt/grid/3.1.39-0/etc/profile.d/grid-env.sh ....
.... SCP your grid certificate dir .globus to $HOME dir ....
.... You will be able to connect to atlas or cms V0 for now ....
.... very soon also to hone V0 when we set authentication to connect with the DESY V0 server ...
.... See you soon with new version of Hone Cern VM!!! ....
-bash-3.005
```



- ZEUS are producing MC on VMs, have validated them against PhysM



Summary and outlook

- > Data preservation and long term analysis is an important project for HEP and an interesting field for computing science
 - HEP is not alone, but sheer amount of data makes it outstanding
 - This talk was far from complete, many more aspects to cover
- > Data preservation alone is worthless, if one does not preserve the ability to perform analysis
- > Different scenarios can be envisaged, two were presented that involve virtualization techniques
- > Two projects are in a prototype stage, will learn from them

- > ... more to come

- > Slides thanks to Dmitry Ozerov, David South, Mihajlo Mudrini, Cristinel Diacono and Homer Neal



Backup slides

Some backup slides with

- > H1 data model
- > More on H1 CernVM usage

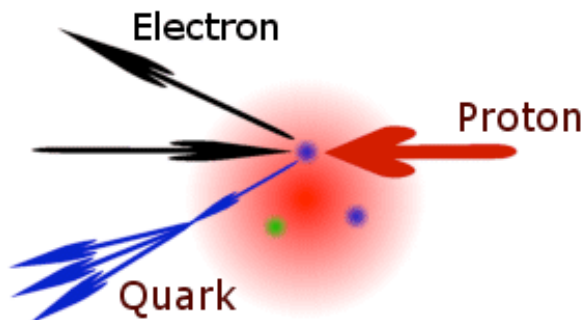


Hadron-Elektron-Ring-Anlage (HERA) at DESY

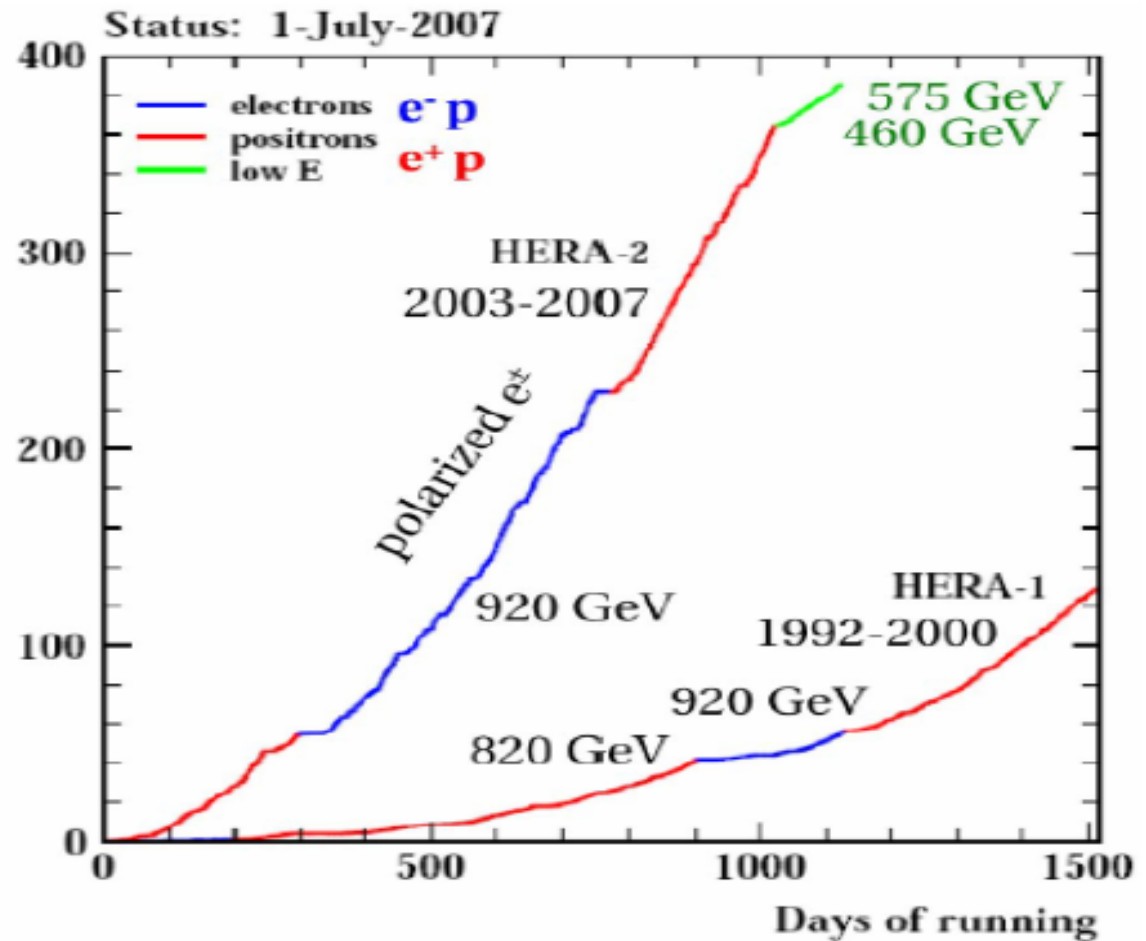
$e^\pm p$ collider, HERA at DESY, Hamburg, Germany.

Luminosity collected: 0.5fb^{-1} per experiment

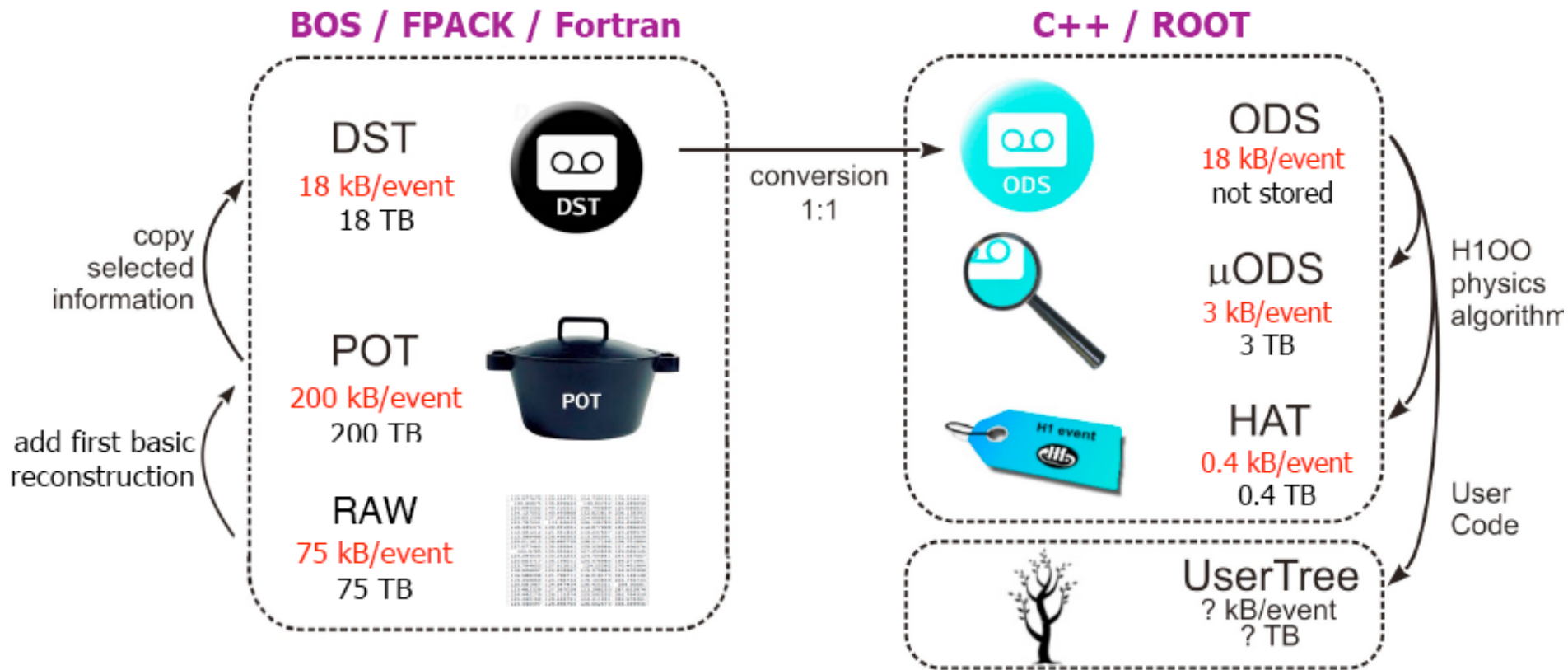
$E_e = 27.6\text{ GeV}$ $E_p = 920\text{ GeV}$ ($\sqrt{s} \approx 320\text{ GeV}$)



- Operated 1992-2007
- p: 460-920 GeV, 110 mA
- e: 27.6 GeV, 45 mA



H1 Data Event Model : Present and Past

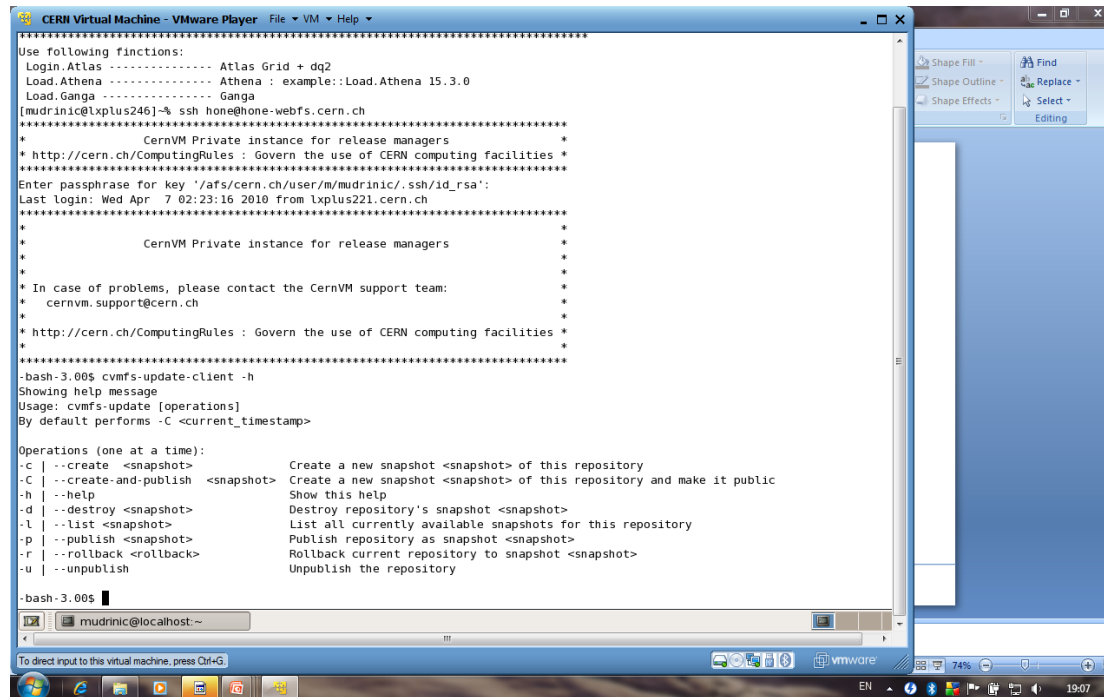


Example: HERA II data
~ 100Tb raw data => 13Tb DST => 1.3 Tb Mods/Hat (root ntuple)



H100 Setup and Maintenance

- > People from CernVM R&D project provided VM on CERN domain dedicated to H1 Experiment hone@hone-webfs.cern.ch
- > hone@hone-webfs.cern.ch is accessible through lxplus.cern.ch with ssh key, and is used for testing and publishing H100 software on CernVM File System.
- > After finishing with development, whole image of software can be published using command `cvmfs-update-client`
- > Many instances of image snapshot could be kept /deleted and published/unpublished.



```
CERN Virtual Machine - VMware Player File VM Help
Use following functions:
Login.Atlas ..... Atlas Grid + dq2
Load.Athena ..... Athena : example::Load.Athena 15.3.0
Load.Ganga ..... Ganga
[mudrnic@lxplus246]~% ssh hone@hone-webfs.cern.ch
*****
* CernVM Private instance for release managers *
* http://cern.ch/ComputingRules : Govern the use of CERN computing facilities *
*****
Enter passphrase for key '/afs/cern.ch/user/m/mudrnic/.ssh/id_rsa':
Last login: Wed Apr 7 02:23:16 2010 from lxplus221.cern.ch
*****
* CernVM Private instance for release managers *
*
* In case of problems, please contact the CernVM support team:
* cernvm.support@cern.ch
*
* http://cern.ch/ComputingRules : Govern the use of CERN computing facilities *
*****
-bash-3.005 cvmfs-update-client -h
Showing help message
Usage: cvmfs-update [operations]
By default performs -C <current_timestamp>

Operations (one at a time):
-C | --create <snapshot>          Create a new snapshot <snapshot> of this repository
-C | --create-and-publish <snapshot> Create a new snapshot <snapshot> of this repository and make it public
-h | --help                      Show this help
-d | --destroy <snapshot>       Destroy repository's snapshot <snapshot>
-l | --list <snapshot>          List all currently available snapshots for this repository
-p | --publish <snapshot>       Publish repository as snapshot <snapshot>
-r | --rollback <rollback>      Rollback current repository to snapshot <snapshot>
-u | --unpublish                 Unpublish the repository

-bash-3.005
```



Running H100 Analysis Framework on CernVM

The screenshot displays a CERN Virtual Machine environment. On the left, a terminal window shows the configuration and execution of the H1SteerTree framework. The terminal output includes:

```
Steering for H1SteerTree
-----
fModsFiles      = 1 entries
fHatFiles       = 1 entries
fTreeCacheSize = 1000000
-----
===== H1SteerManager: Done reading from file 'Exam
=====> Addfile for file /home/hluser/H1AnalysisDir/Data
=====> Addfile for file /home/hluser/H1AnalysisDir/Data
Consistency checks for list H1TreeEventList
=====
===== H1SteerManager: Using default values for class
H1SteerOdsEvent bank names:
HEAD CRME FRME TOFT TOFS CRPE DBPC DMIS DELE DBFC HRD
JDPX GHD GKI GTR GVX GEVC STR SVX SIPA TL23 NSF
BRUE DTAG DRP1 DRP2 DRP3 DL5W DTZS CSKH CSHY DFTS Y4T
DISM JMFI JMYX HD00 TEL1 HEAR FNCE FNCL FRSE FRXT DRM
TT10 TT20 TT30 TT1T TT2T TT2K TTEV BJKR
H1BankEvent uses all 3/4 letter branches in SetBranch
H1SteerOdsEvent: Call to cstrec disabled
===== H1SteerManager: End of defaults for class 'H1S
Creating transient ODS tree for chain number 1
Following subtrees exist:
Subtree nr 1 ODS (transient)
Subtree nr 2 MODS with 1 files and cachesize 1000000
Subtree nr 3 HAT with 1 files and cachesize 1000000
=====> Opening of H1Tree ok ===== Number of Events: 890
=====
=====> H1Tree::Open No HAT selection in H1SteerTree
=====> H1Tree::Open No Ascii List of Events given in H1
Info in <H1Tree::SelectHat>: Starting selection "NumJ
40064 events selected
event 5000 *****
[]
```

On the right, three histograms are displayed, each showing the distribution of $\mu\text{ODS: } m_{JPV}$ for different production and track types. The histograms are:

- elastic production, track-track ($\mu-\mu$)**: Shows a distribution with a peak around 3.0. Statistics: Entries: 63, Mean: 2.942, RMS: 0.4394.
- elastic production, track-cluster**: Shows a broader distribution with a peak around 3.0. Statistics: Entries: 405, Mean: 3.071, RMS: 0.5142.
- inelastic production, track-track ($\mu-\mu$)**: Shows a very narrow distribution with a peak around 3.161. Statistics: Entries: 4, Mean: 3.161, RMS: 0.6535.



