

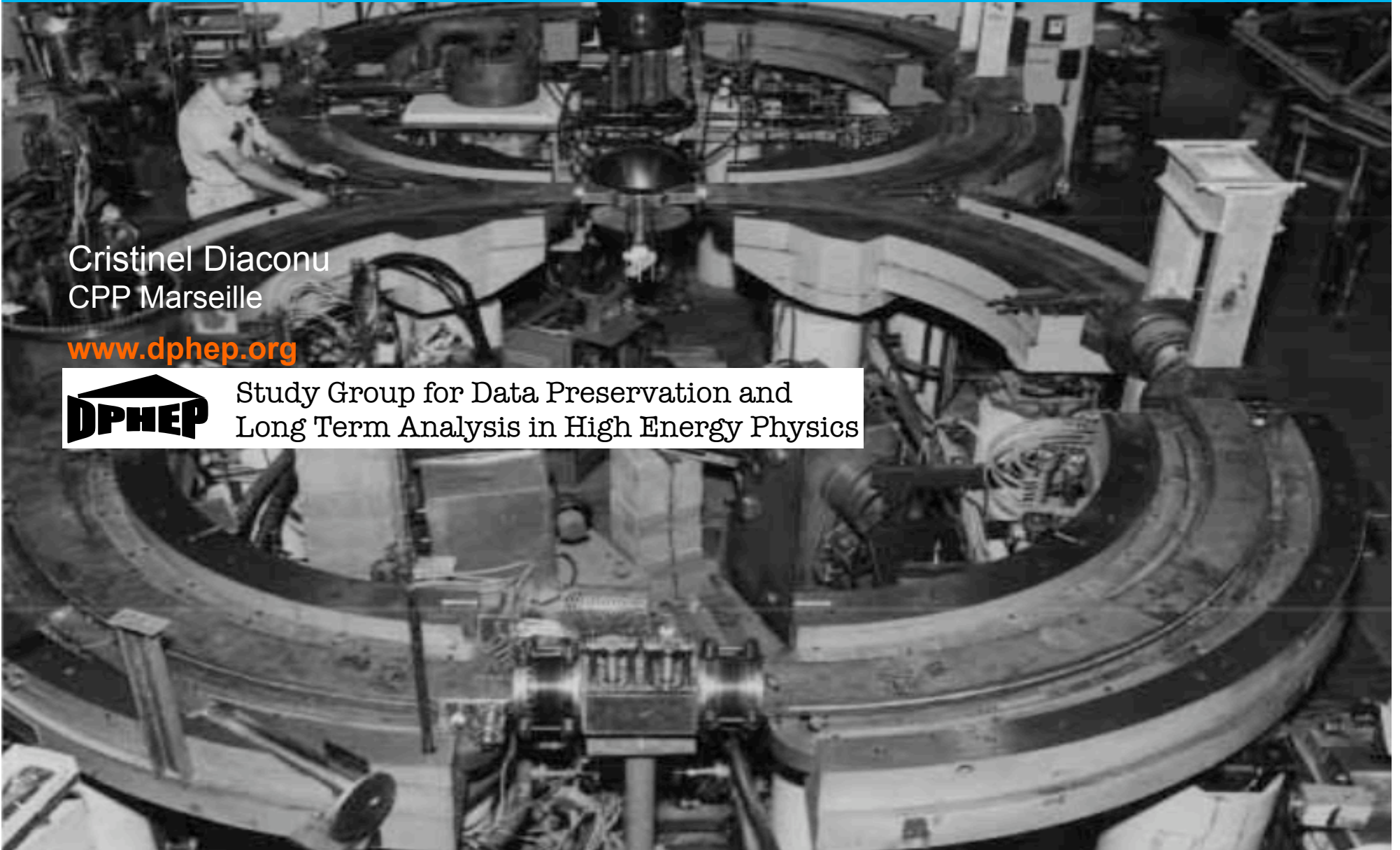
# Data Preservation in High Energy Physics

Cristinel Diaconu  
CPP Marseille

[www.dphep.org](http://www.dphep.org)



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



# Outline

## > Introduction

- Data preservation in HEP
- An international initiative: DPHEP
- The scientific potential of HEP data

## > DPHEP data preservation models

- Current strategies of the experiments
- Emerging projects in the DPHEP community

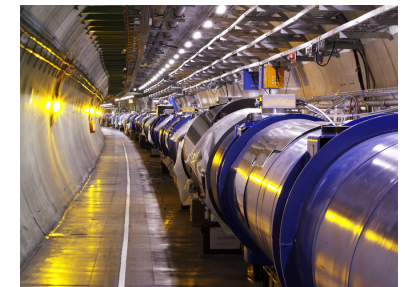
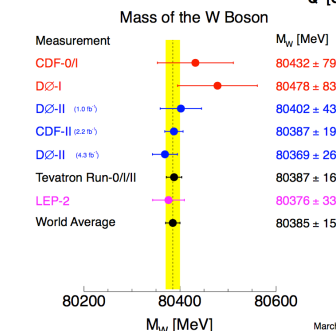
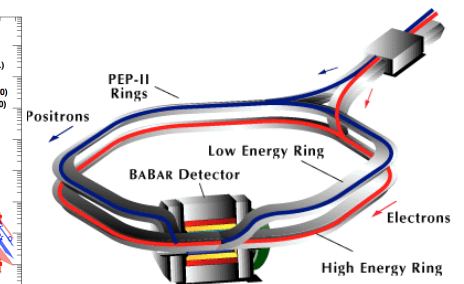
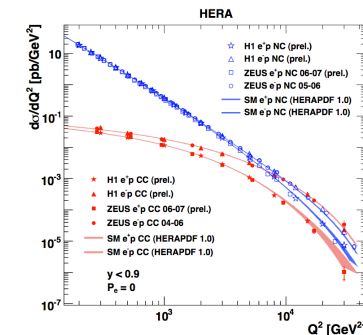
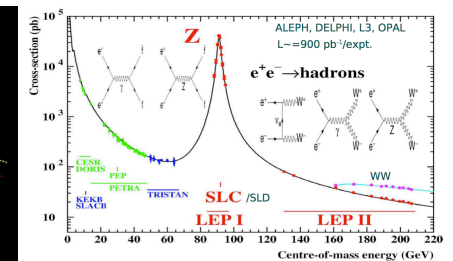
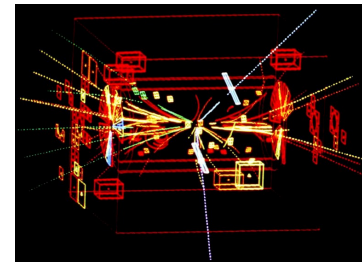
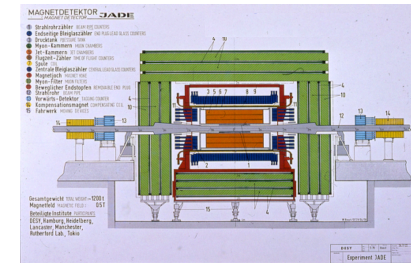
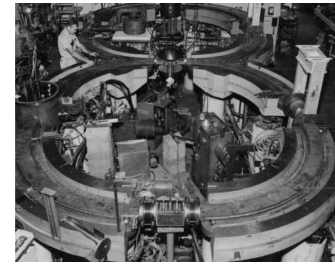
## > Future working directions

- Where we are now, where we need to go, and how that's going to happen

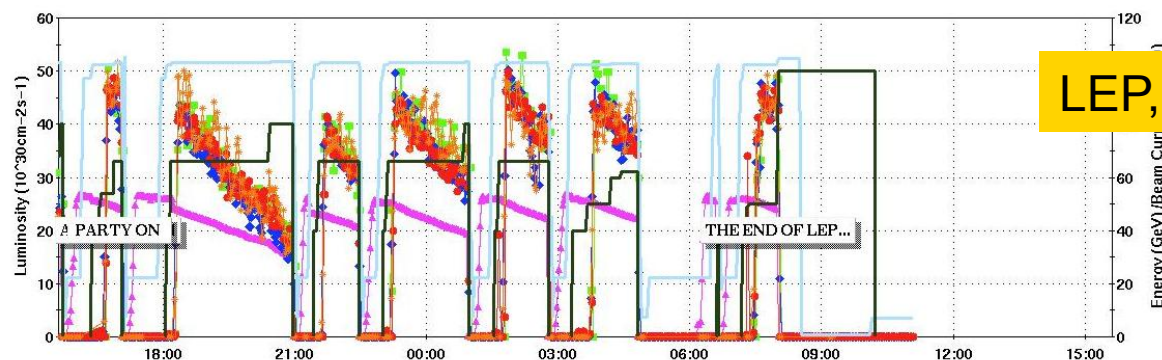


# Experimental particle physics in the collider era

- A wide variety of physics results from many, often very different experiments
- Energy frontier probed with increasingly complex accelerator installations
  - New experiments typically supersede previous, similar ones - but not always
- Growth in size of the necessary international collaborations, as well as the diversity of the data management
- The age of the LHC has truly arrived
  - The Super-B factories and other projects such as the ILC or next e-p(A) collider are to come



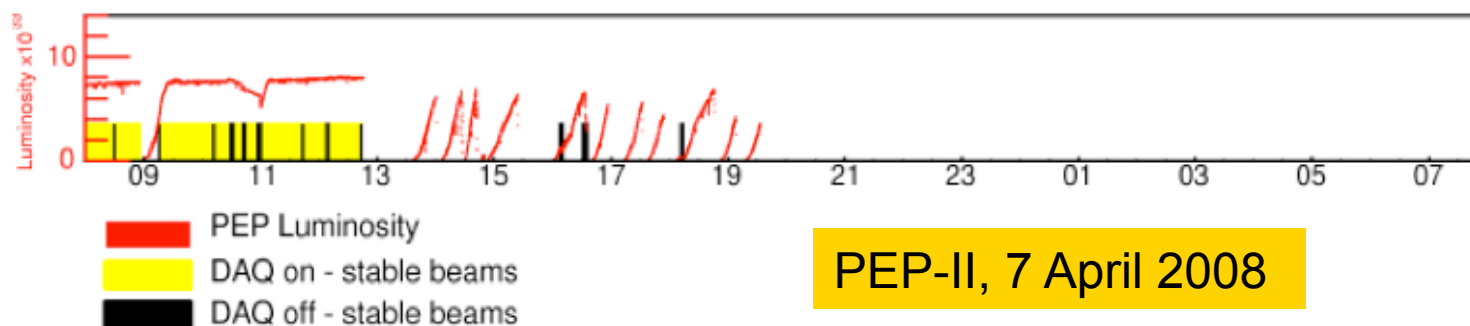
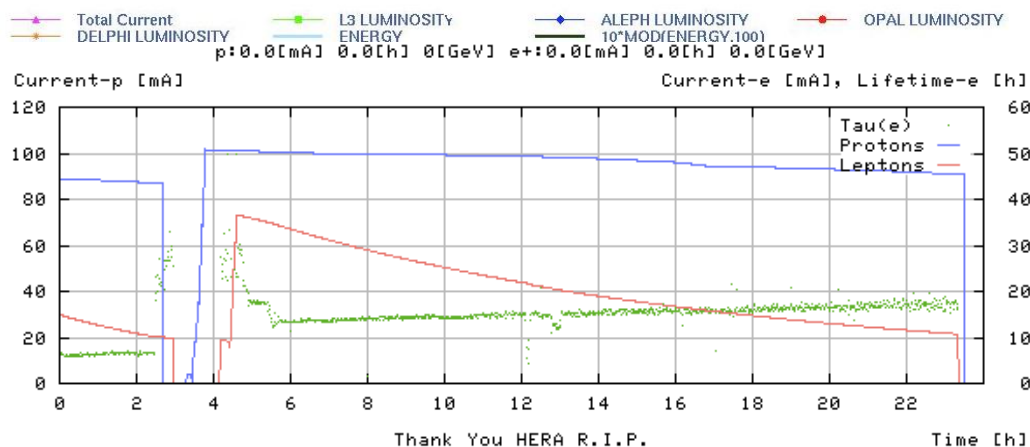
# The last years have seen the end of several experiments



LEP, 2 November 2000



HERA, 30 June 2007



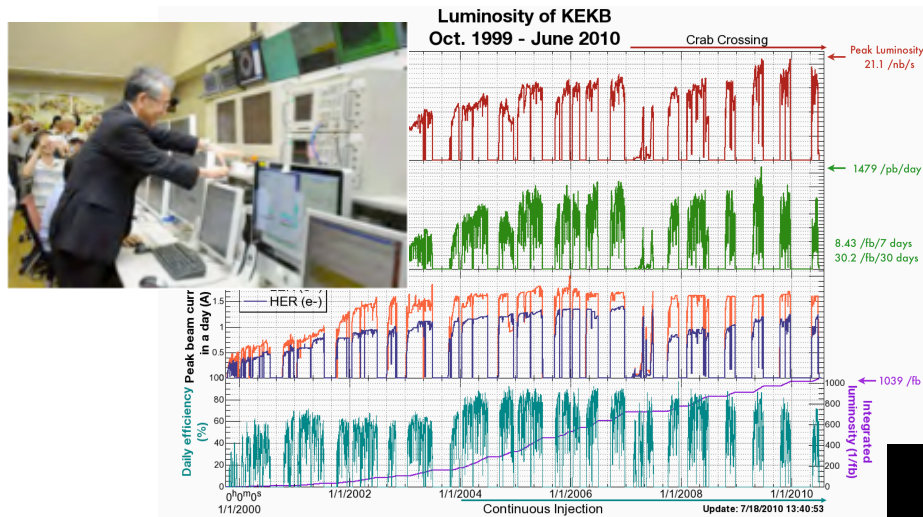
PEP-II, 7 April 2008

Mon Apr 7  
04:34h (19.0%) Stable Beams  
04:15h (93.1%) DAQ on  
84.6/pb Recorded Lumi  
1.3% DCH paused  
18s DCH paused (# 4)

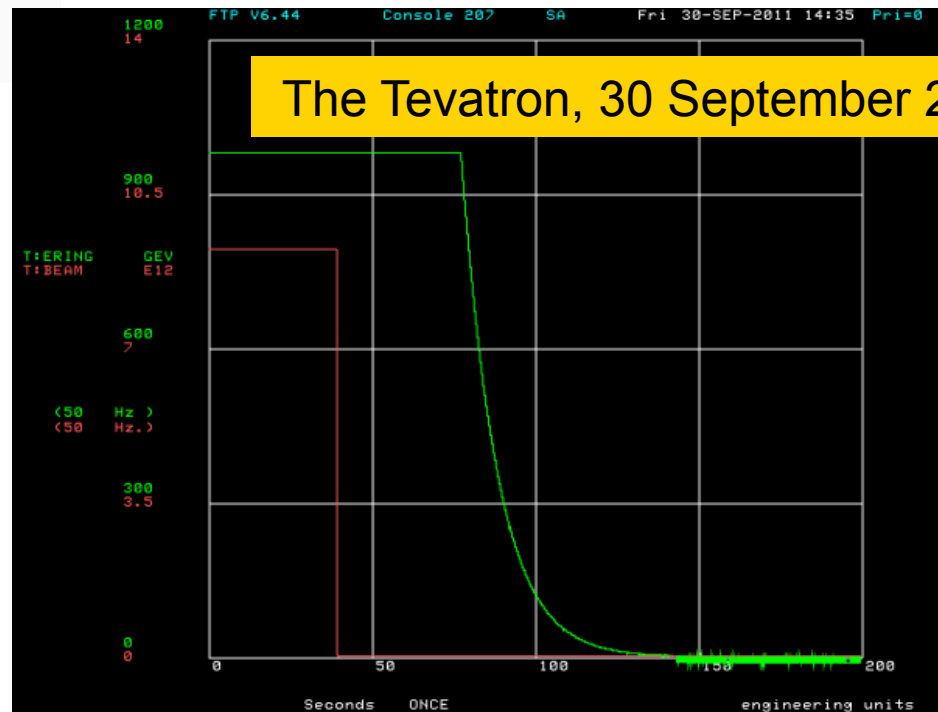


# The last years have seen the end of several experiments

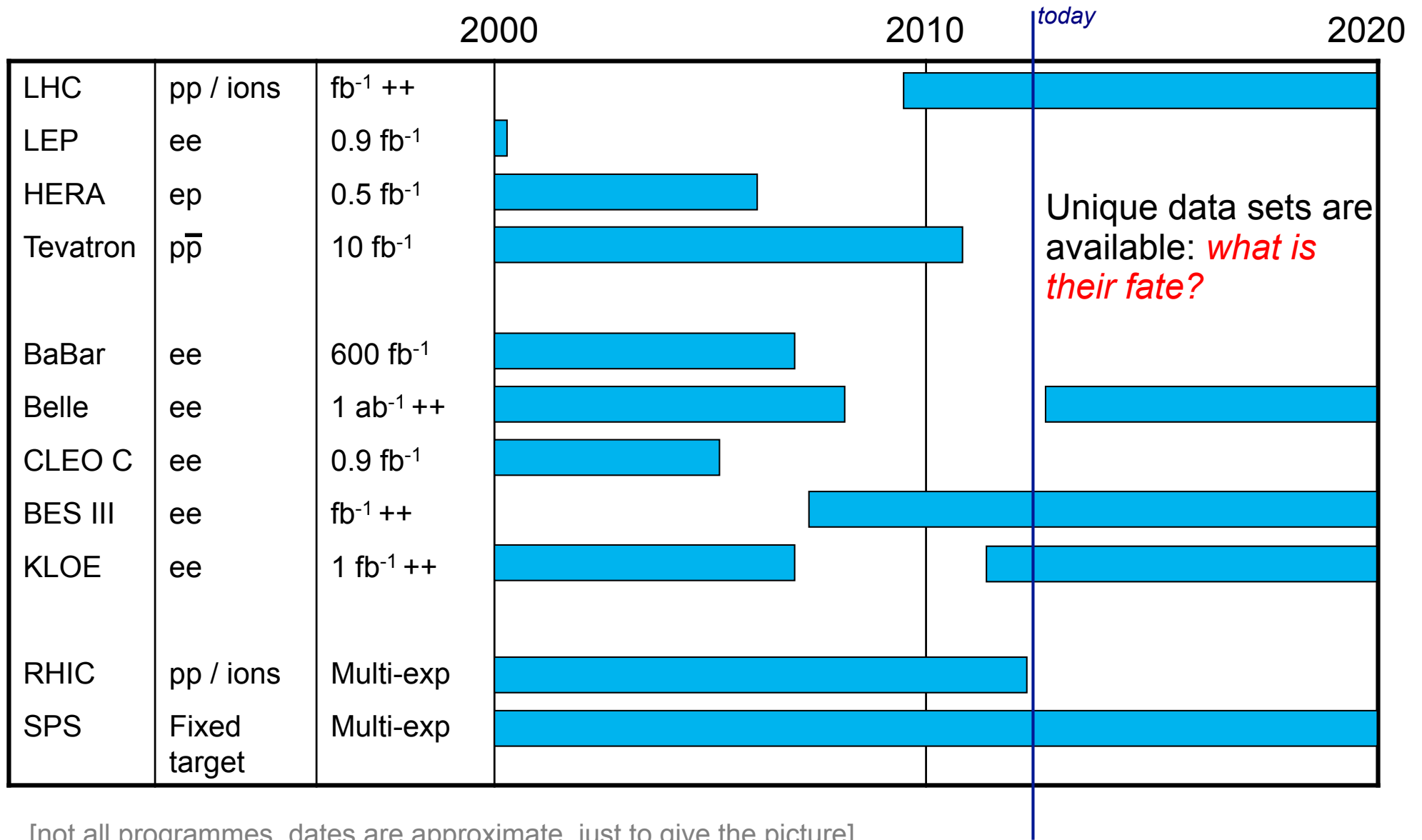
KEKB, 30 June 2010



The Tevatron, 30 September 2011



# HEP experimental programmes $\pm 10$ years



[not all programmes, dates are approximate, just to give the picture]

# The email you may receive one day (I did)

*To Whom it may concern,*

*In the tape storage area we still have 4132 tapes of type 3840 containing HERA data.*

*We do not have a functioning reading device anymore and the storage area was polluted recently, so it is likely that the tapes are damaged.*

*Would you like us to send you these tapes or should we **destroy them directly?***

*Yours Sincerely,*

*Tape admin. service [a large computing centre]*





# After the collisions have stopped

- Finish the analyses! But then what do you do with the data?
  - Until recently, there was no clear policy on this in the HEP community
  - It's possible that older HEP experiments have in fact simply lost the data
- Data preservation, including long term access, is generally not part of the planning, software design or budget of an experiment
  - So far, HEP data preservation initiatives have been in the main not planned by the original collaborations, but rather the effort a few knowledgeable people



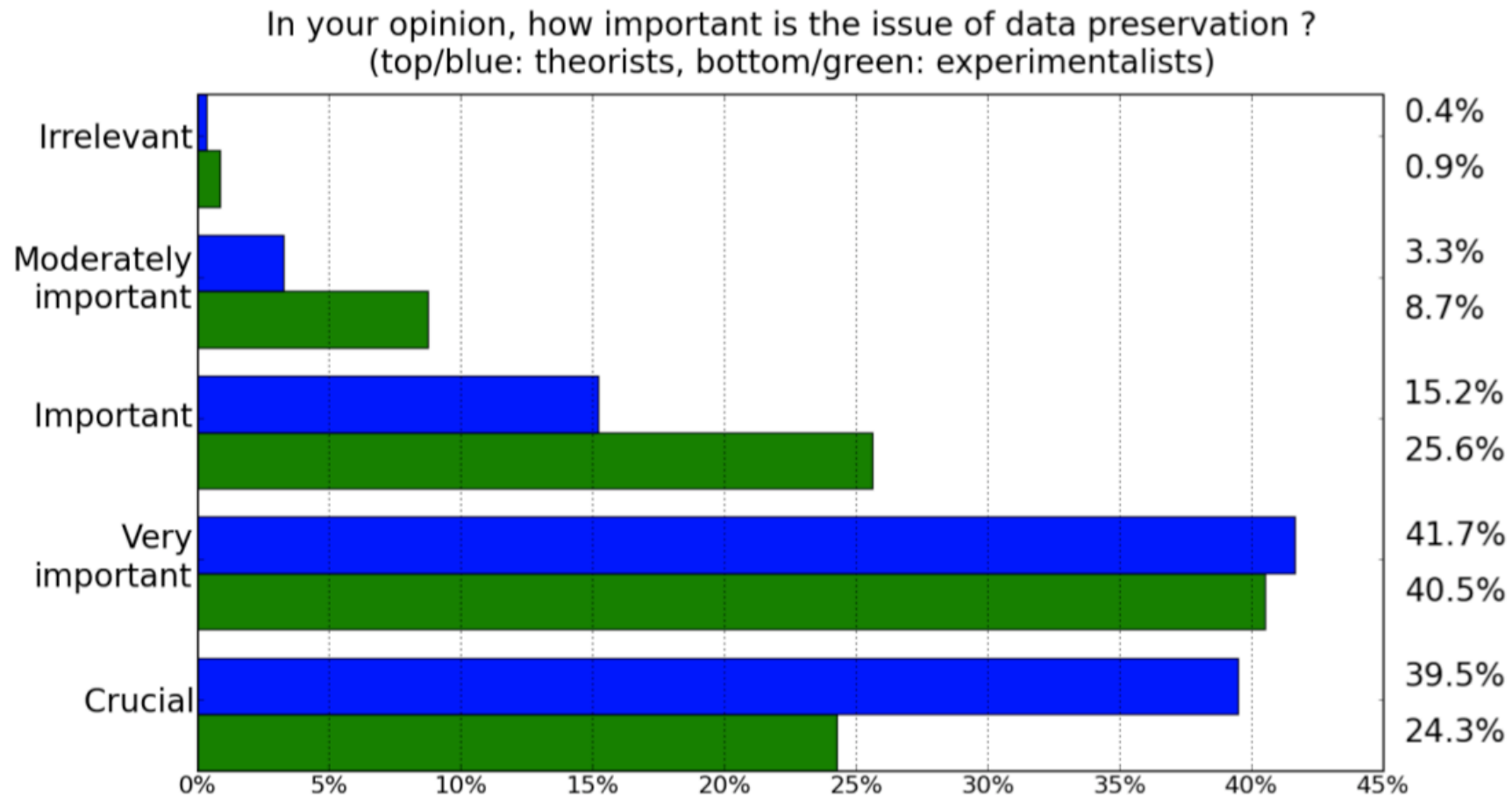
- The conservation of tapes is not equivalent to data preservation!

- *“We cannot ensure data is stored in file formats appropriate for long term preservation”*
- *“The software for exploiting the data is under the control of the experiments”*
- *“We are sure most of the data are not easily accessible!”*

# The difficulties of data preservation in HEP

- > Handling HEP data involves large scale traffic, storage and migration
  - The increasing scale of the distribution of HEP data can complicate the task
- > Who is responsible? The experiments? The computing centres?
  - Problem of older, unreliable hardware: unreadable tapes after 2-3 years
  - The software for accessing the data is usually under the control of the experiments
- > Key resources, both funding and person-power expertise, tend to decrease once the data taking stops
- > And a rather key ingredient to all this is: *why do it?*
  - Can the relevant physics cases be made?
  - Who says we want to do this anyway?
  - Is the benefit of all this really worth the cost and effort?

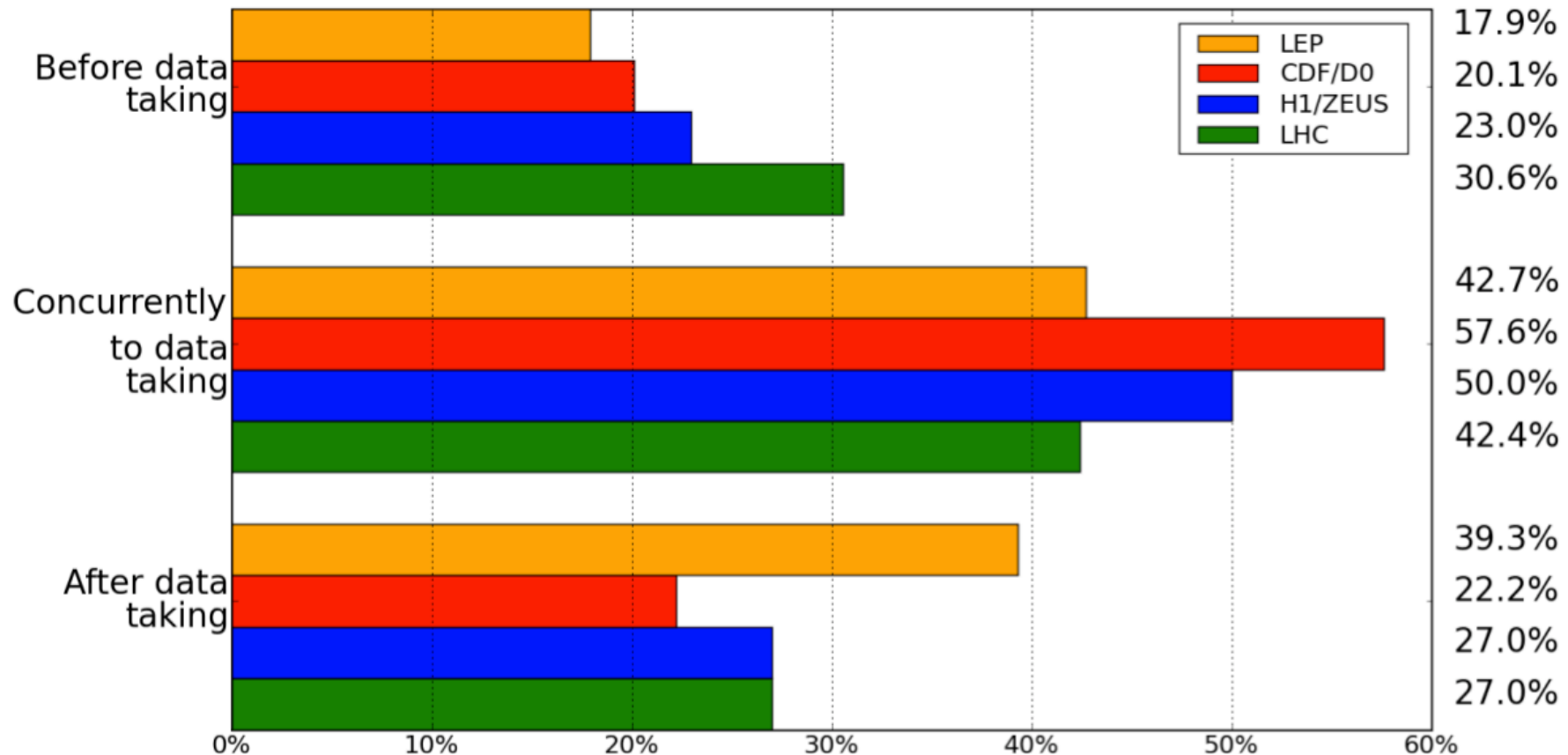
# Support for data preservation in the HEP community



arXiv:0906.0485

# Support for data preservation in the HEP community

In your opinion, when should this effort start in order to be the most effective ?

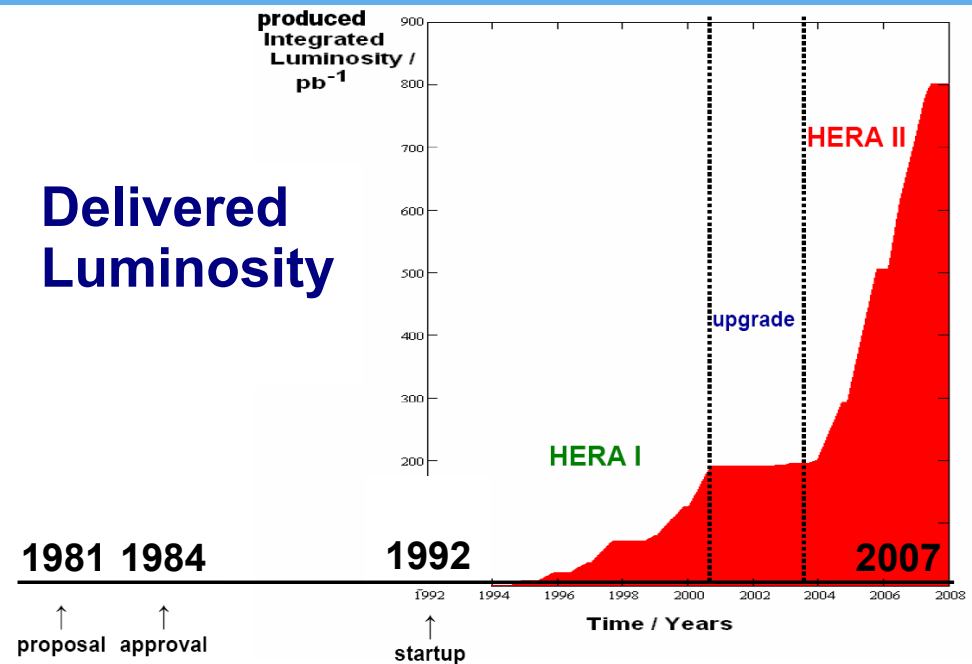


arXiv:0906.0485

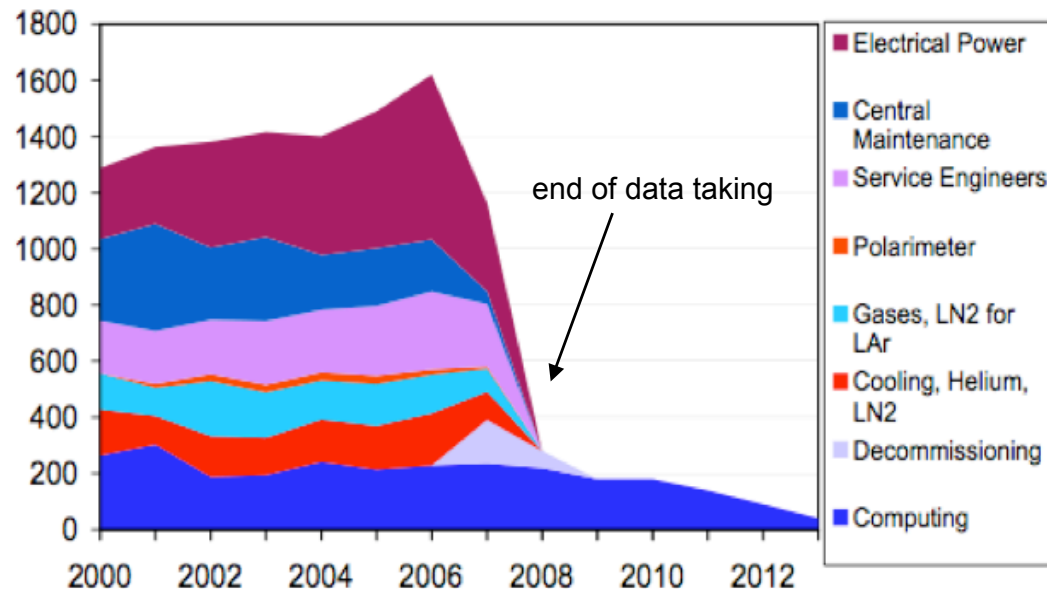
# Why is it Difficult to Preserve HEP Data?

- > Lots of data available to analyse at the end of collisions
- > The existing resources (funding and expertise) then decrease when the data taking stops

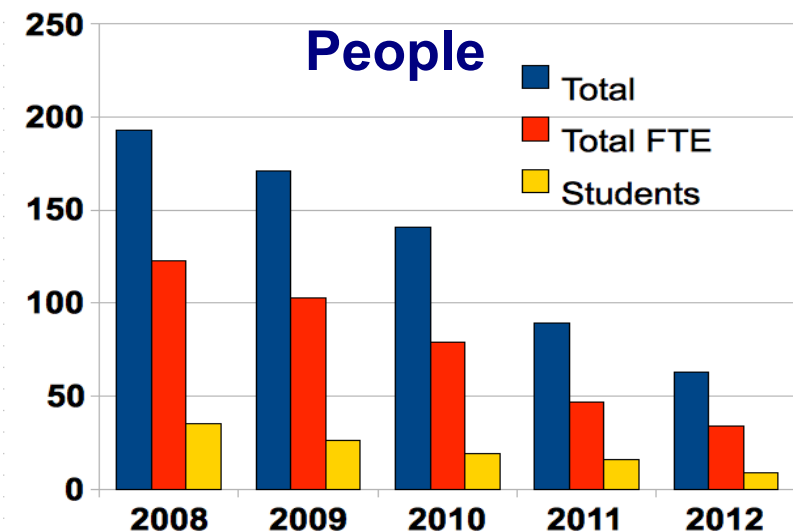
## Delivered Luminosity



## Funding

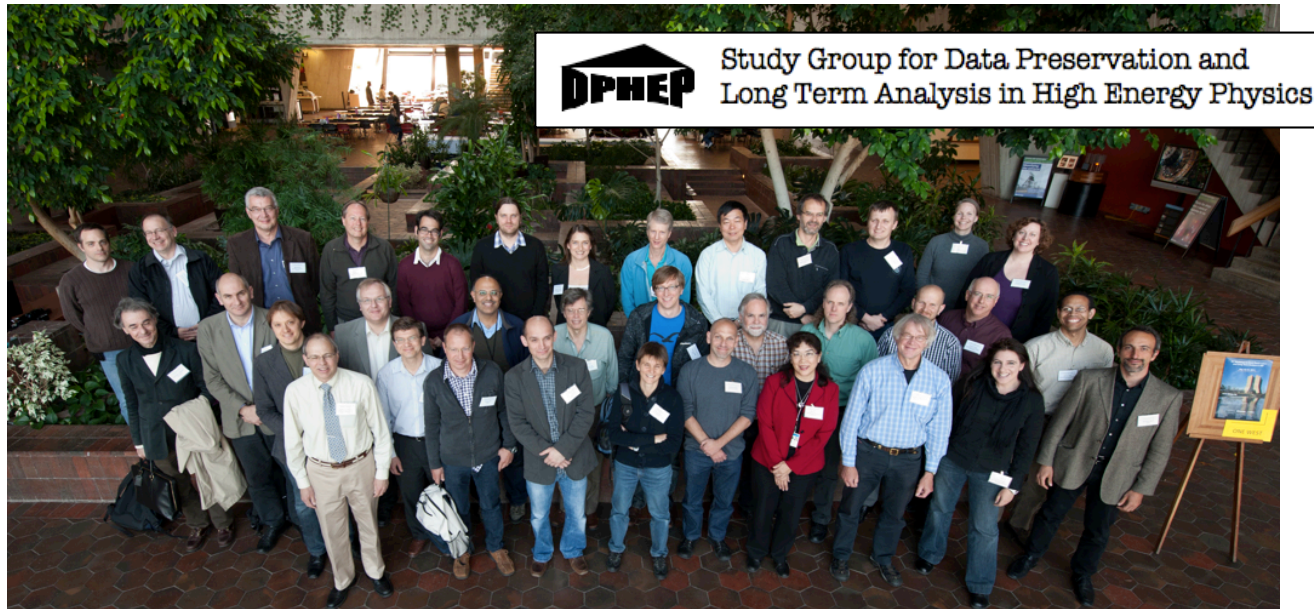


## People



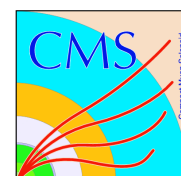
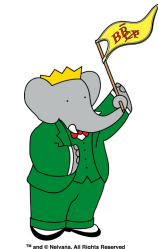
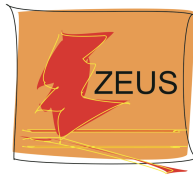


# DPHEP: An international study group on data preservation



- First contacts established in September 2008
  - Group since grown to over 100 contact persons (chair : CD)
  - Endorsed as an ICFA panel summer 2009
  - *All 4 LHC experiments joined in 2011*
- Steering Committee: representatives from all members
- International Advisory Committee:
  - Jonathan Dorfan (Chair, SLAC), Siegfried Bethke (Chair, MPIM), Gigi Rolandi (CERN), Michael Peskin (SLAC) Dominique Boutigny (IN2P3), Young-Kee Kim (FNAL), Hiroaki Aihara (IPMU/Tokyo), Alex Szalay (JHU)

# DPHEP: An international study group on data preservation



Institute of High Energy Physics  
Chinese Academy of Sciences



Jefferson Lab



Science & Technology  
Facilities Council

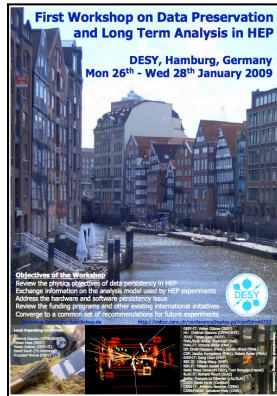


INSPIRE



# DPHEP: An international study group on data preservation

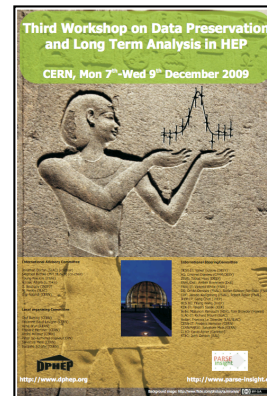
## > Series of DPHEP workshops held since 2009



Jan 2009: DESY



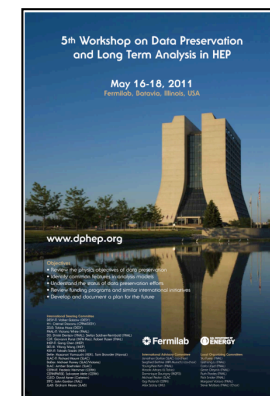
May 2009: SLAC



Dec 2009: CERN



Jul 2010: KEK



May 2011: Fermilab

## > The first task of the group was to establish the working directions

- “To confront data models, clarify the concepts, set a common language, investigate technical aspects, compare with other fields handling large data.”

## > Initial findings published in an interim report December 2009

- Focus on four key areas of the study group: **Physics Case for Data Preservation, Preservation Models, Technologies, Governance**

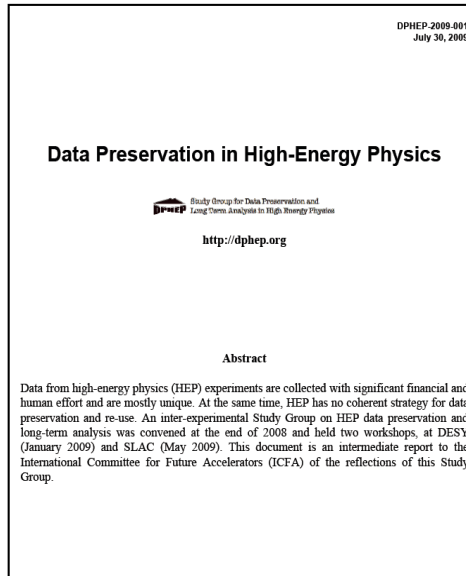
**arXiv:0912.0255**





# DPHEP Intermediate Recommendations (end 2009)

> **arXiv:0912.0255**



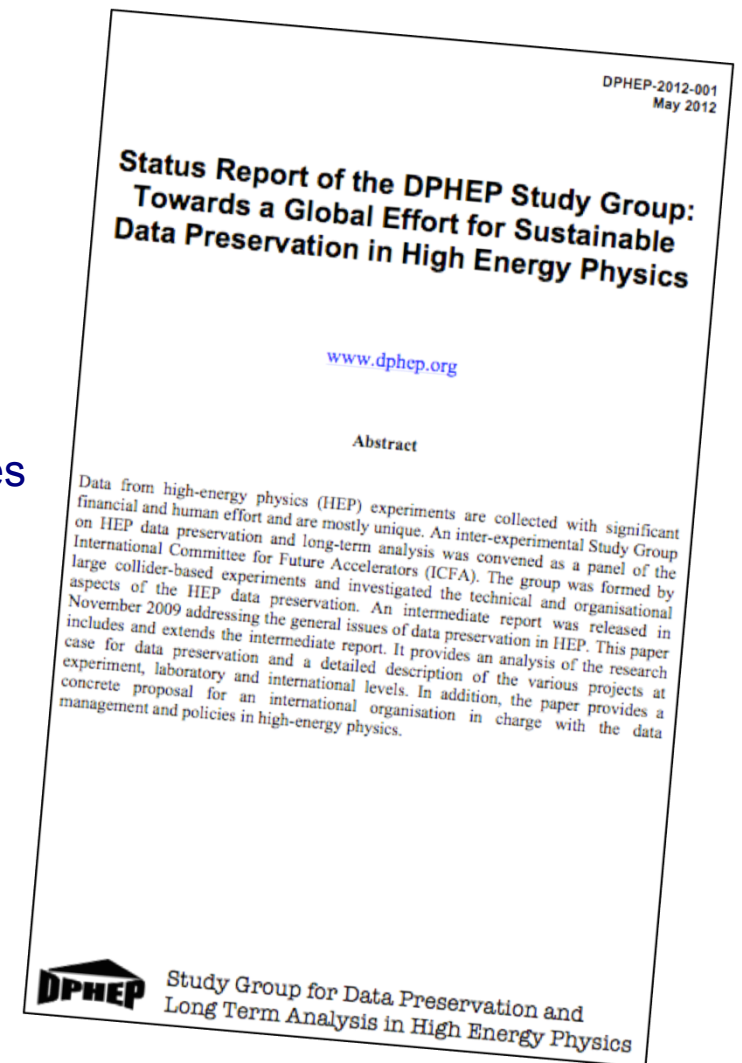
- > An urgent and vigorous action is needed to ensure data preservation in HEP
  - Many examples for the physics case explored
  - Data is rich and can be further exploited in most cases beyond the collaboration lifetime
- > The preservation of the full analysis capability of experiments is recommended, including the preservation of reconstruction and simulation software
- > An interface to the experiment know-how should be introduced: **data archivist** position in the computing centres
- > The preservation of HEP data requires a synergic action: collaborations, laboratories and funding agencies
- > An International Data Preservation Forum is proposed as a reference organisation. The Forum should represent experimental collaborations, laboratories and computing centres



# New DPHEP publication

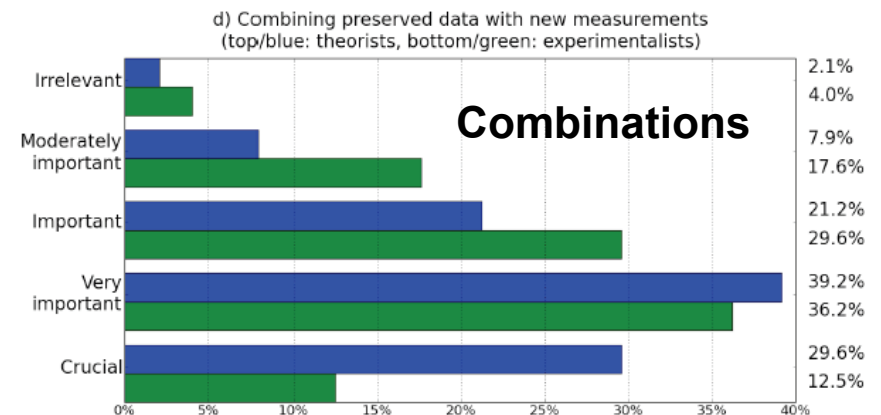
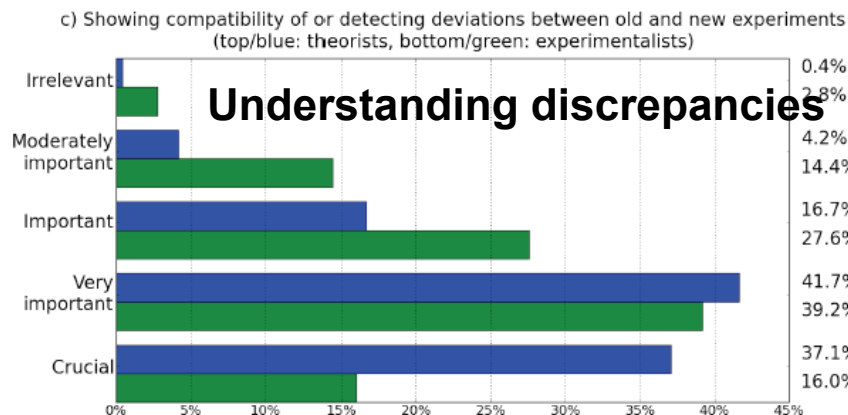
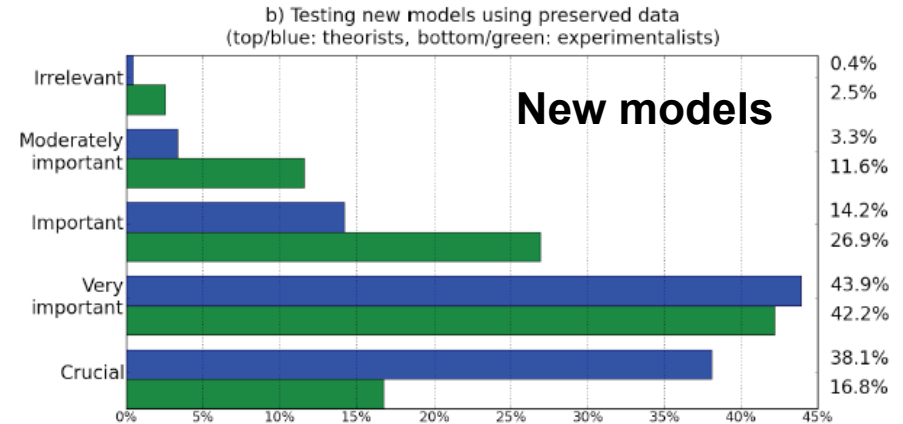
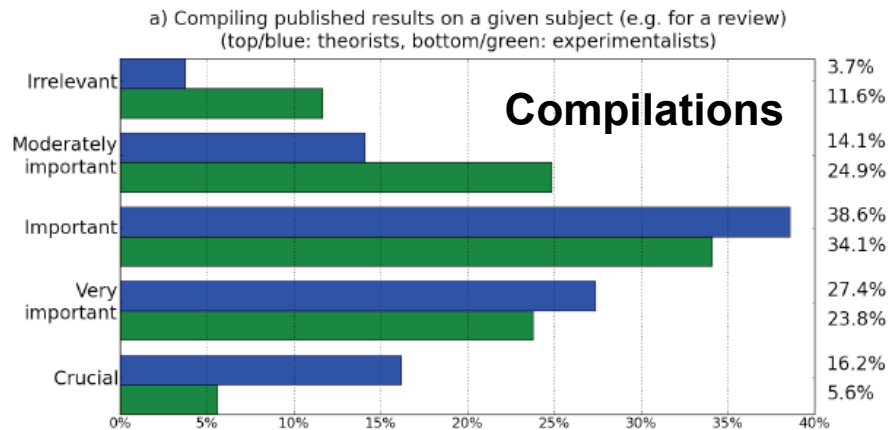
- Full status report of the activities of the DPHEP study group, including:
  - Tour of data preservation activities in other fields
  - An expanded description of the physics case
  - Defining and establishing data preservation principles
  - Updates from the experiments and joint projects
  - FTE estimates for these and future projects
  - Next steps to establish fully DPHEP in the field

**arXiv:1205.4667**



# Physics case: opinions in the HEP community

## Preserving HEP data is important for:



# Building the physics case: Reasons to preserve HEP data

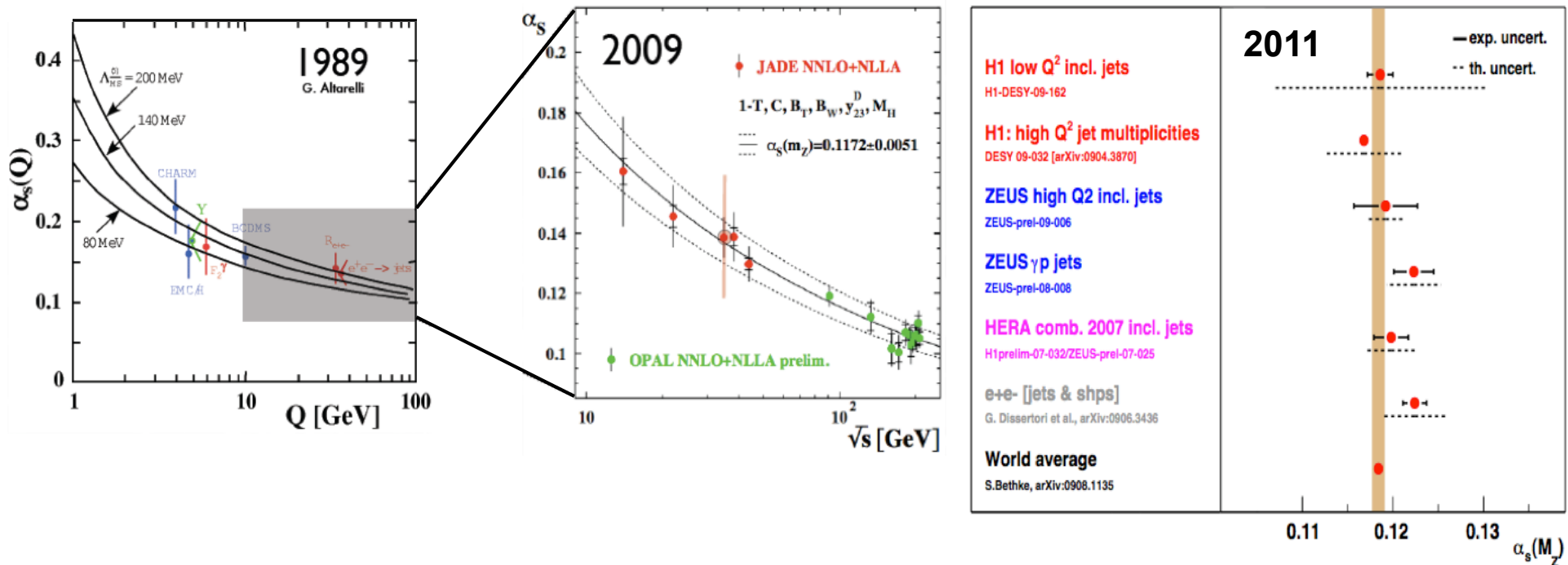
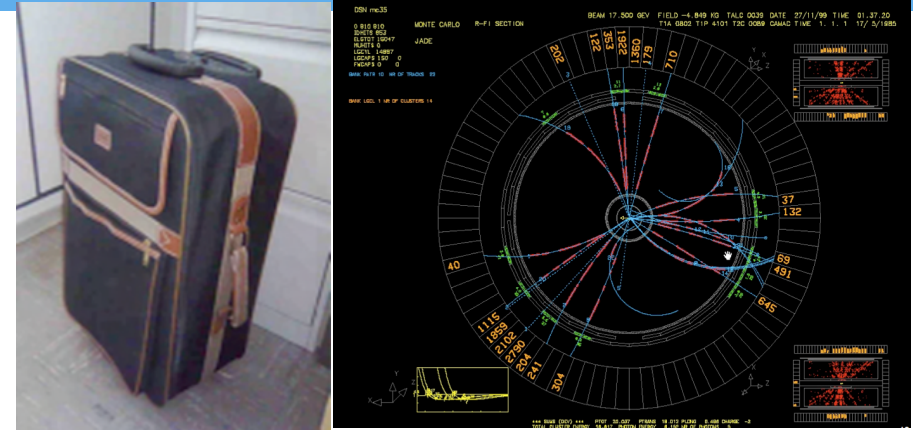
- > Long term completion and extension of an existing physics program
  - Up to 10% of papers are finalised in the “archival mode”
  - Gain in scientific output of the experiments
- > Cross-collaboration and combinations of physics results
  - During the active lifetime of similar experiments at one facility: LEP, HERA, TeVatron
  - And later across larger boundaries: Belle/BaBar, TeVatron/LHC
- > Revisit old measurements or perform new ones
  - Access to newly developed techniques, comparisons to new theoretical models
  - Unique data sets available in terms of energy, initial states
- > Use in scientific training, education, outreach



# Physics case: Improvement in theory and simulation

- JADE: Required full raw data preservation, software revitalisation, individual initiatives...

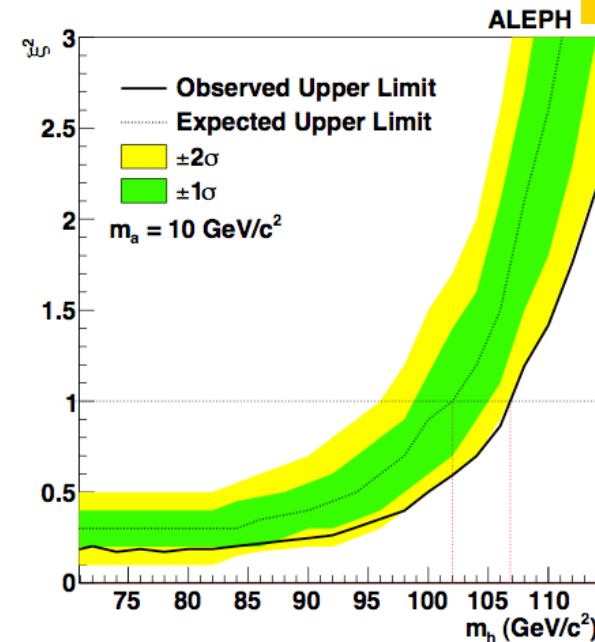
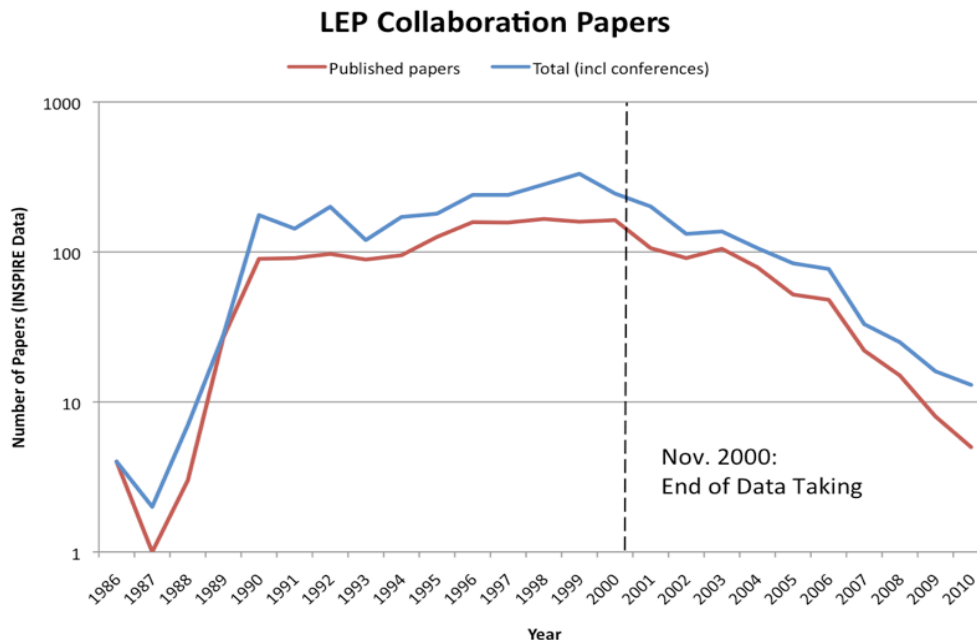
10 recent publications



- Around 10% of measurements are dominated by non-experimental errors: theory ( $N^{\text{th}}$ LO?) and simulation..

# Long term completion of the physics programme

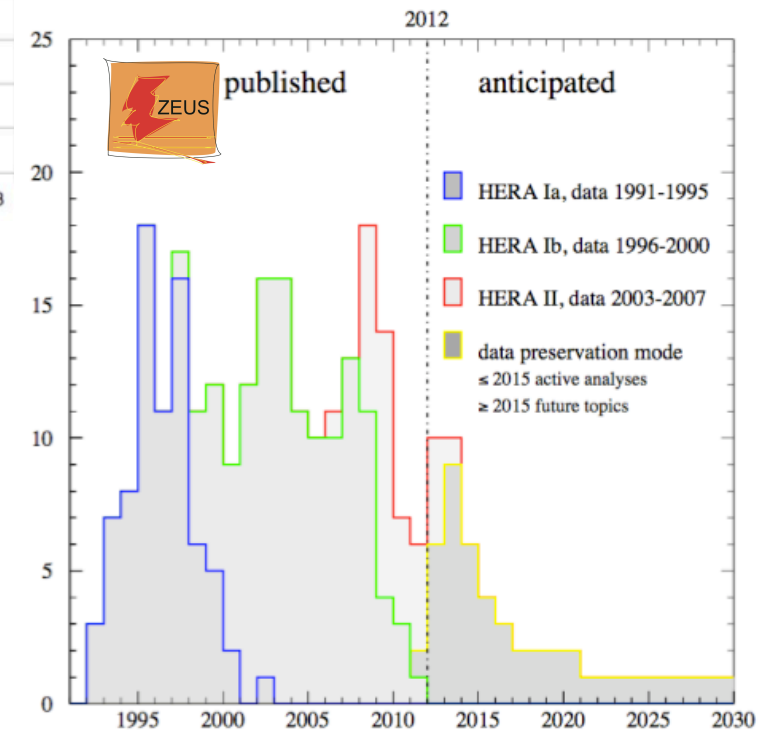
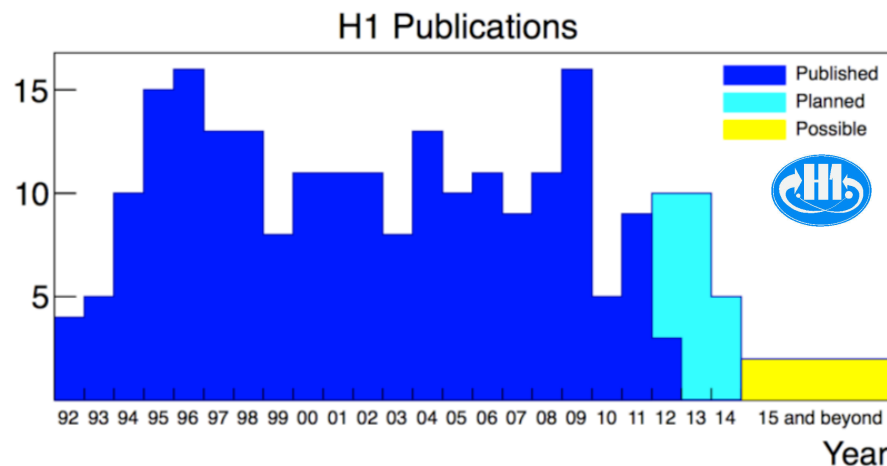
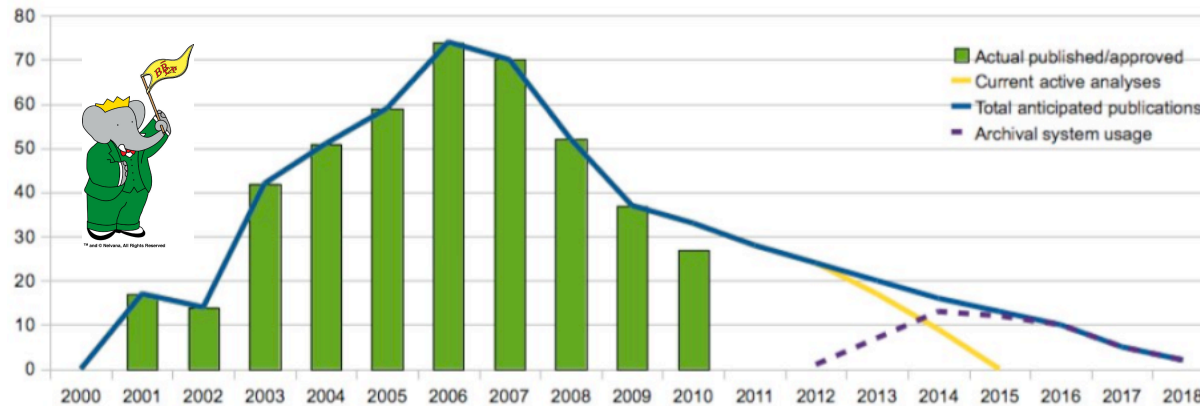
arXiv:1003.0705



- The publication tail of LEP is long, with new papers still appearing
- Well over 300 papers produced since the end of collisions in 2000
- Recent analysis of LEP data gave unique limits on a novel Higgs model
- Similar, if not longer publication tails predicted by the BaBar, H1 and ZEUS experiments, after taking into consideration the plans for data preservation



# Long term completion of the physics programme



- Similar publication tails predicted by the BaBar, H1 and ZEUS experiments, taking into consideration the plans for data preservation

# Cross-collaboration combinations of physics results

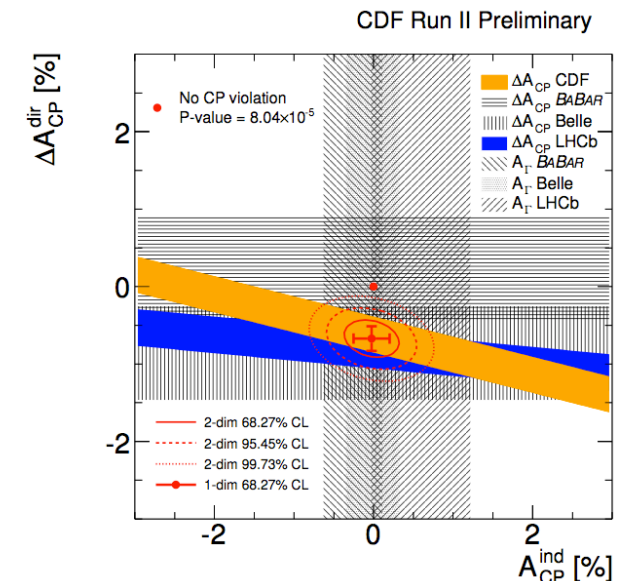
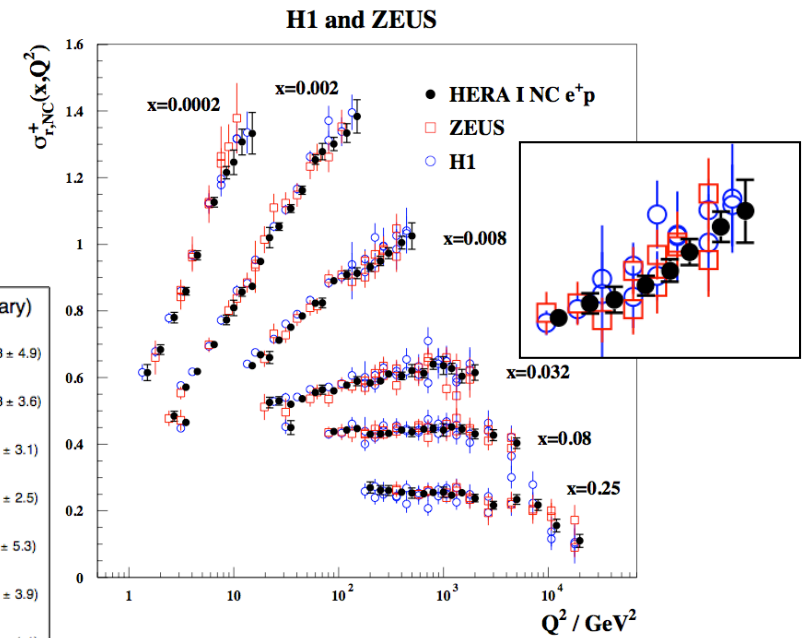
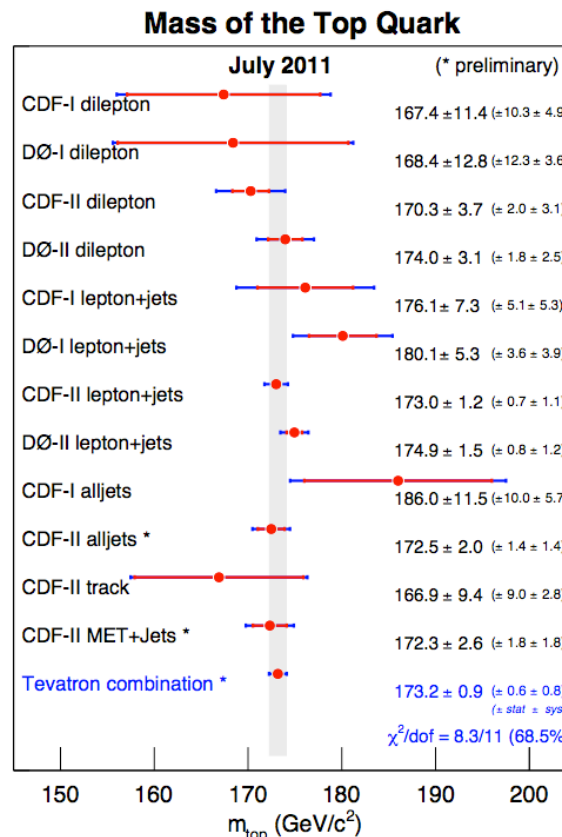
- > Combination of data from multiple experiments to produce new scientific results

- Improved precision and increased sensitivity

- > Comparison of experimental results

- Complimentary information from different physics
  - Verification of experimental observations

- > Both objectives facilitated by data preservation



# New theories, new interpretations

CERN-PH-EP-2011-080  
May 20, 2011

## Test of the $\tau$ -Model of Bose-Einstein Correlations and Reconstruction of the Source Function in Hadronic Z-boson Decay at LEP

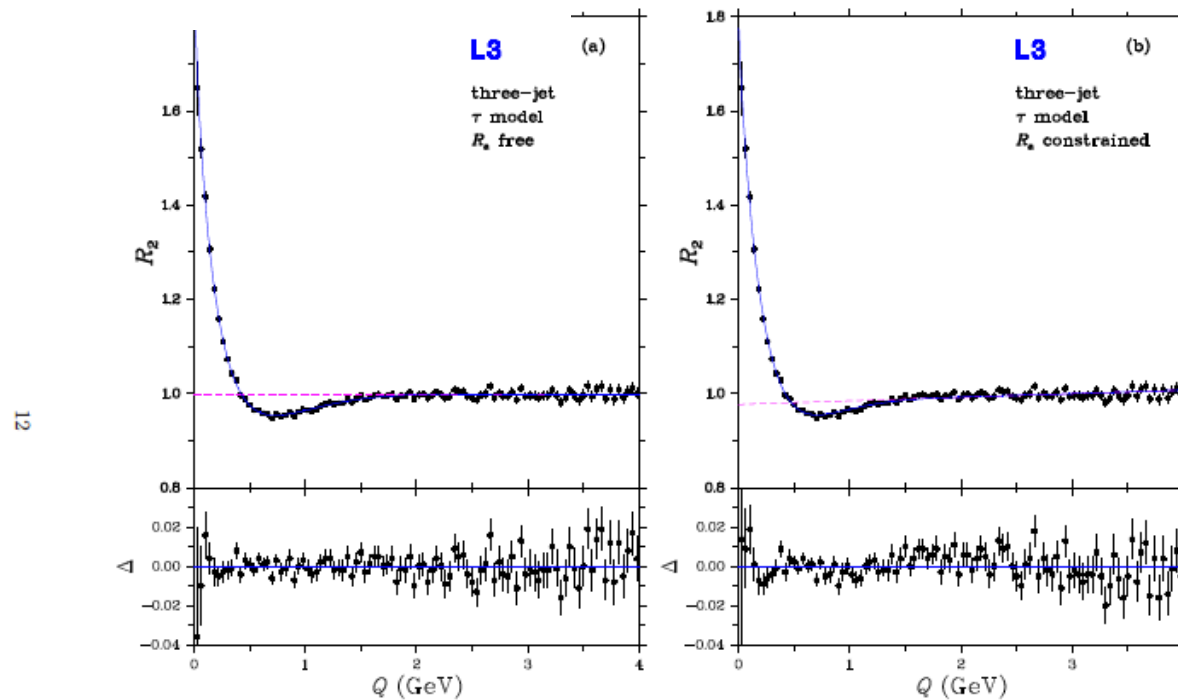


Figure 5: The Bose-Einstein correlation function  $R_2$  for three-jet events. The curve corresponds to the fit of the one-sided Lévy parametrization, Eq. (13), with the parameter  $R_a$  (a) free and (b) constrained by Eq. (14). The results of the fits are given in Tables 1 and 2, respectively. Also plotted is  $\Delta$ , the difference between the fit and the data. The dashed line represents the long-range part of the fit, *i.e.*,  $\gamma(1 + \epsilon Q)$ .

# Dark photons: subject is new, data is old

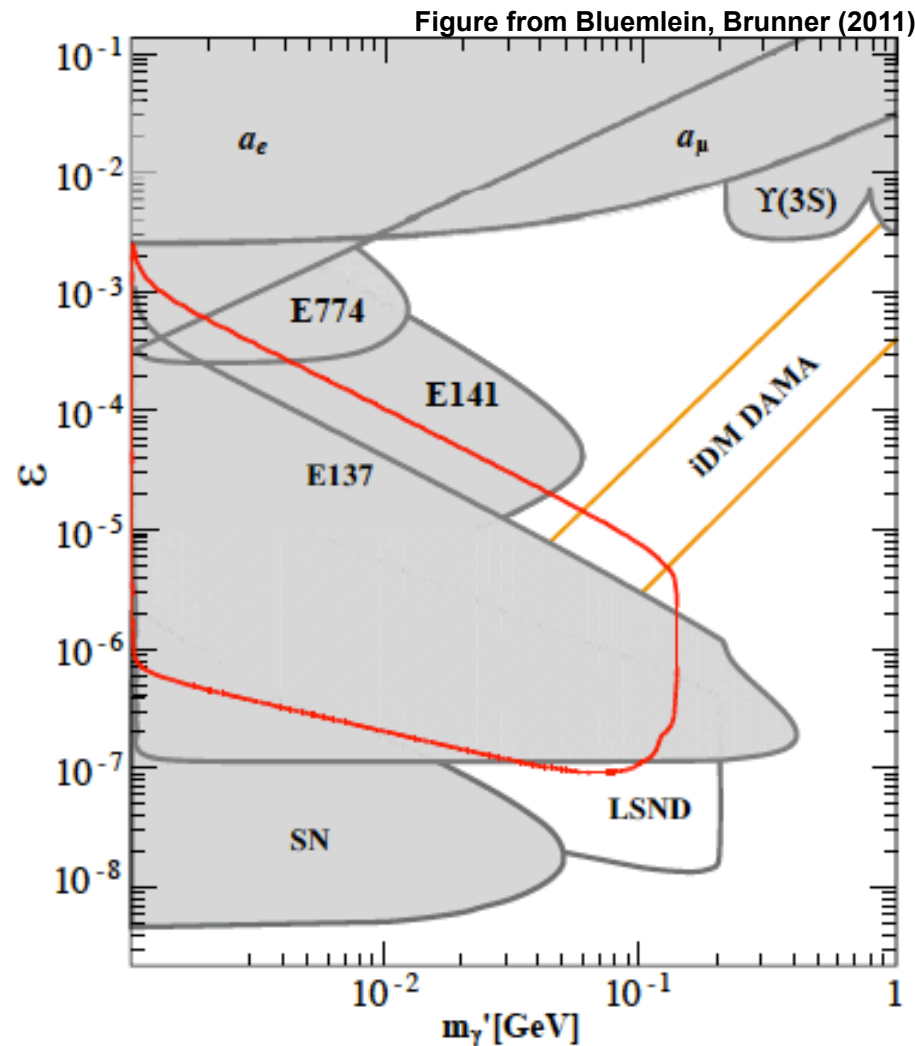
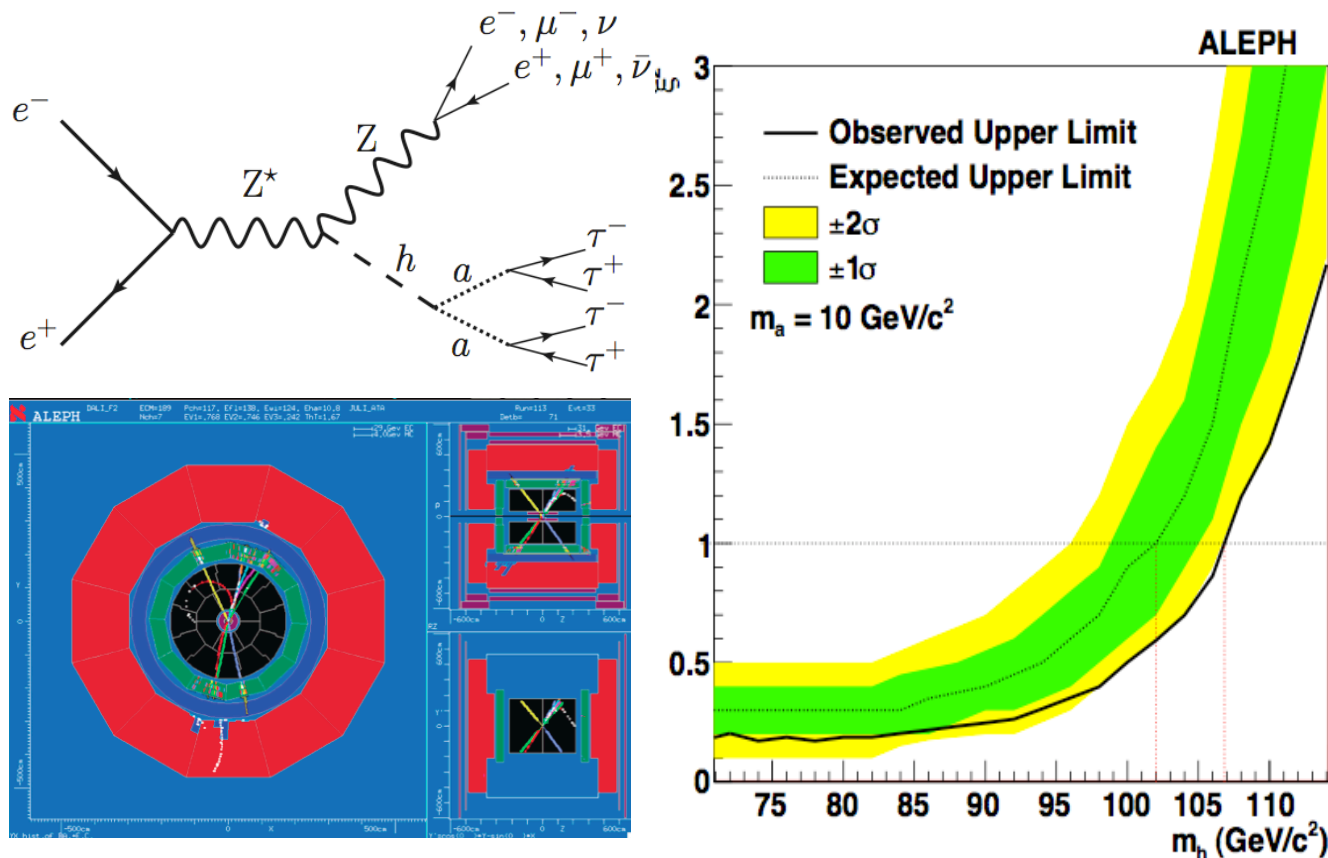


Figure 5: Comparison of the present exclusion bounds (red line) with other limits from the measurement of the anomalous magnetic moments  $a_e$  and  $a_\mu$  [19],  $\Upsilon(3S)$  decay [20], the beam dump experiments E137, E141, E774 [21–23], and supernovae cooling [4, 24]. We indicate the prospects for LSND [7, 25] (open grey-bounded area), and the DAMA/LIBRA region (open orange bounded area) [26]. The limits for  $\epsilon > 10^{-7}$  have been taken from Ref. [6].

# Physics case: searches in previous data sets

- Theory and “common sense” evolve
- ALEPH: Unique physics case analysed 10 years after the end of collisions
  - and 5 years after the official end of the collaboration



EUROPEAN ORGANISATION FOR NUCLEAR RESEARCH (CERN)

**Search for neutral Higgs bosons  
decaying into four taus at LEP2**

The ALEPH Collaboration\*)

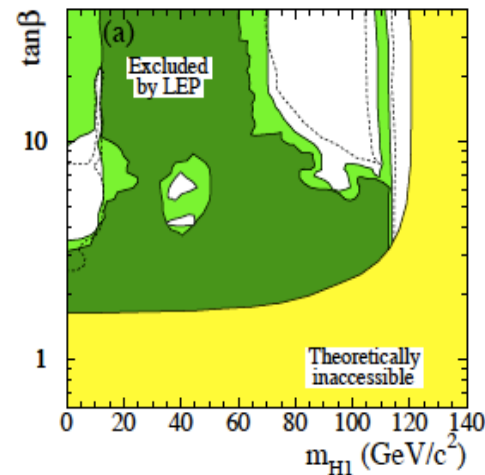
**Abstract**

A search for the production and non-standard decay of a Higgs boson,  $h$ , into four taus through intermediate pseudoscalars,  $a$ , is conducted on 683 pb<sup>-1</sup> of data collected by the ALEPH experiment at centre-of-mass energies from 183 to 209 GeV. No excess of events above background is observed, and exclusion limits are placed on the combined production cross section times branching ratio,  $\xi^2 = \frac{\sigma(e^+e^- \rightarrow Z h)}{\sigma_{\text{had}}(e^+e^- \rightarrow Z h)} \times B(h \rightarrow aa) \times B(a \rightarrow \tau^+ \tau^-)^2$ . For  $m_h < 107 \text{ GeV}/c^2$  and  $4 < m_a < 10 \text{ GeV}/c^2$ ,  $\xi^2 > 1$  is excluded at the 95% confidence level.

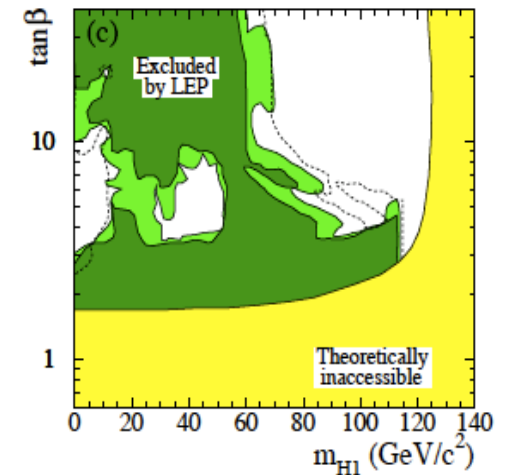
# Excluded?

- > Some external parameters may be not well known
- > Re-optimisation may be a case for re-analysis

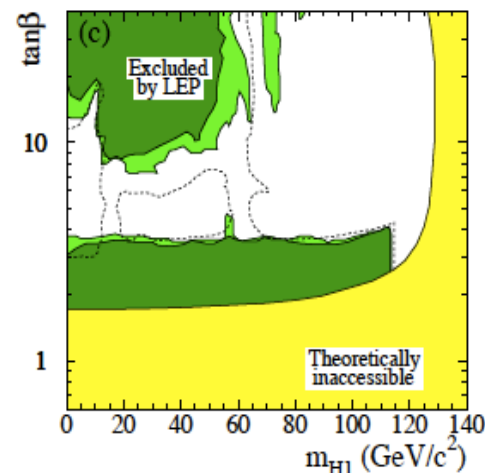
$$m_t = 169.3 \text{ GeV}/c^2$$



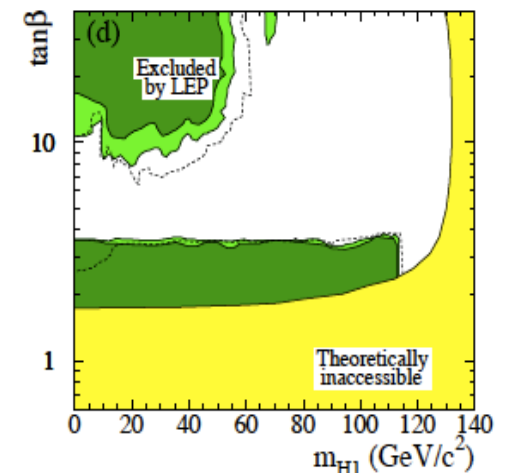
$$m_t = 174.3 \text{ GeV}/c^2$$



$$m_t = 179.3 \text{ GeV}/c^2$$



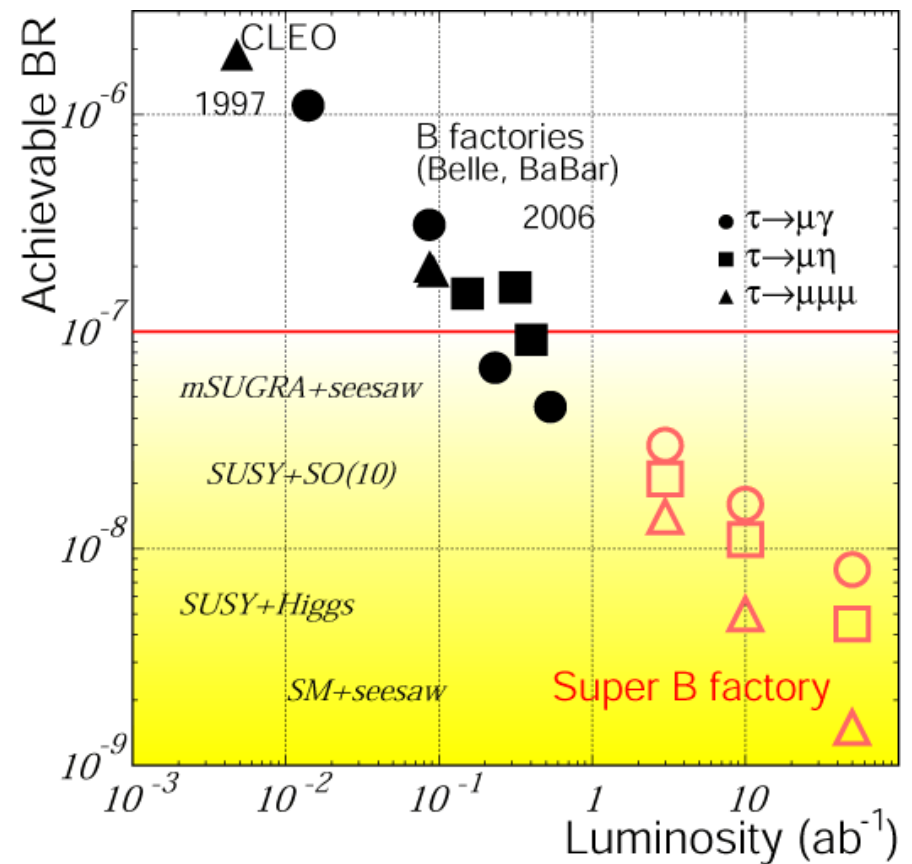
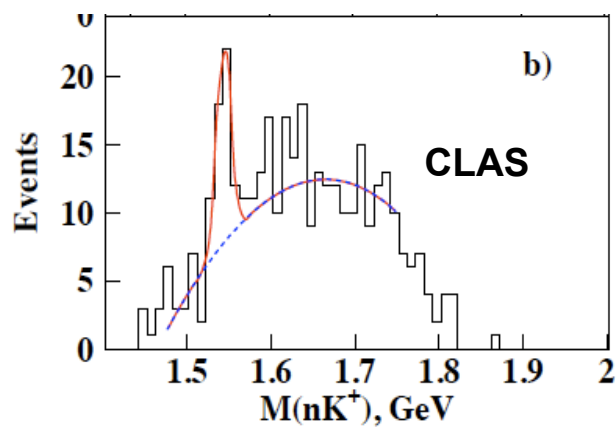
$$m_t = 183.0 \text{ GeV}/c^2$$



# More examples...

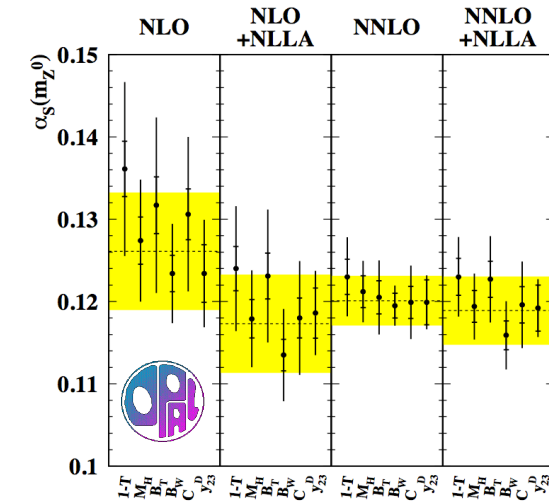
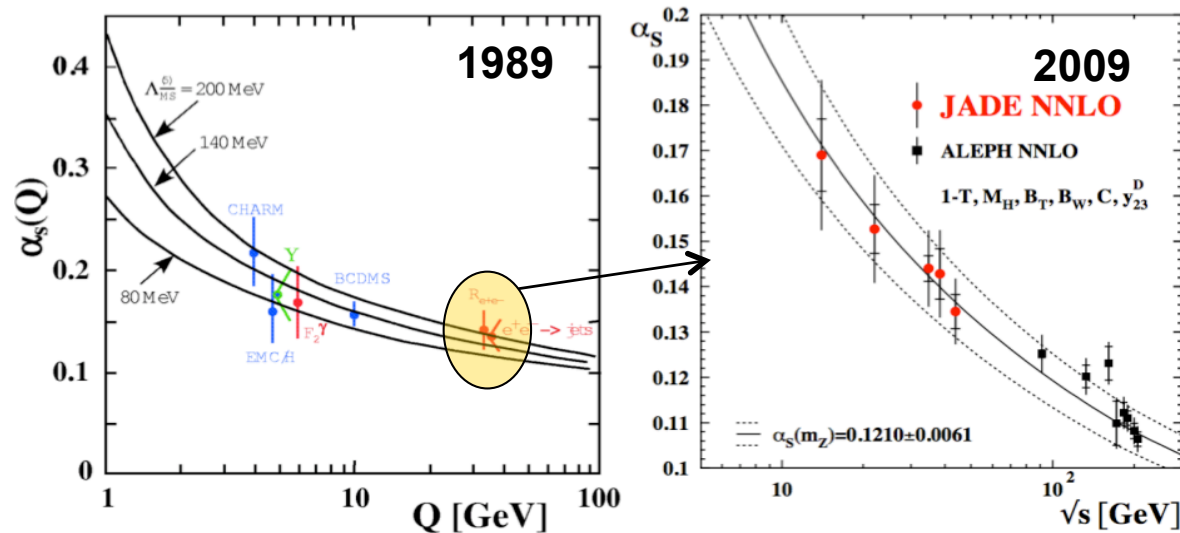
- B- and SuperB-factories
- Low energy
- ...and many others
  - your favourite?

...surprises can occur  
at lower energies too

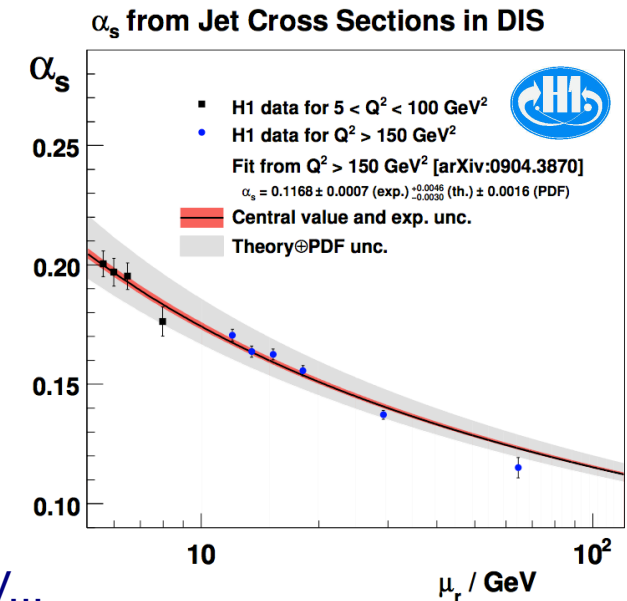




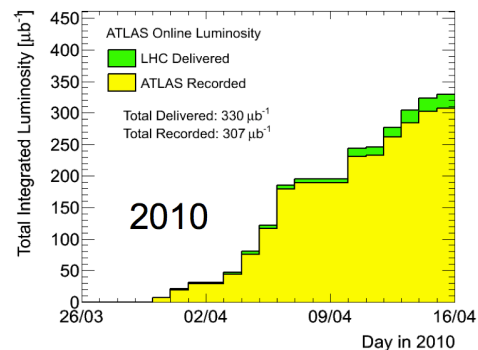
# Revisit old measurements or perform new ones



- > Access to newly developed techniques, comparisons to new theoretical models
  - History to be repeated with the HERA  $\alpha_s$  measurements
- > Unique data sets are available in terms of initial state particles and energy
  - HERA  $e^\pm p$ , Tevatron  $p\bar{p}$ , fixed target experiments..
  - Early LHC data: 900 and 2.36 TeV, 2010 low pile-up 7 TeV...

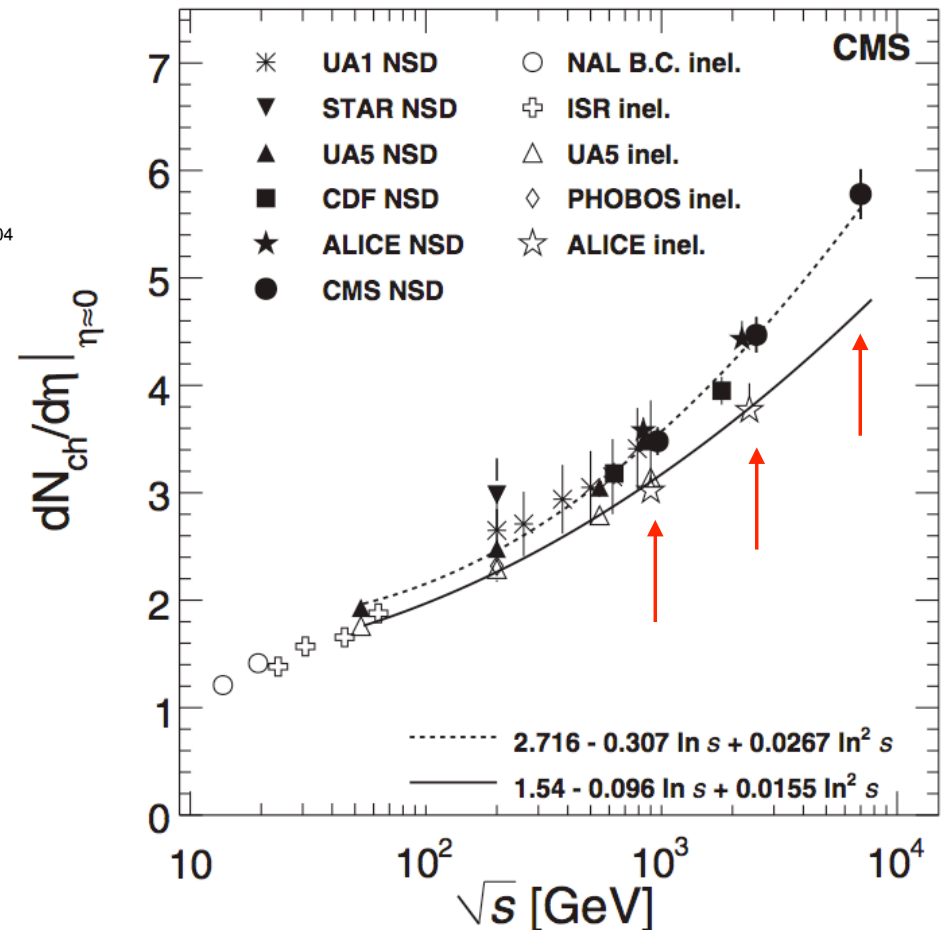


# What about LHC 900 GeV and 2.32 TeV data? And 7 TeV data?

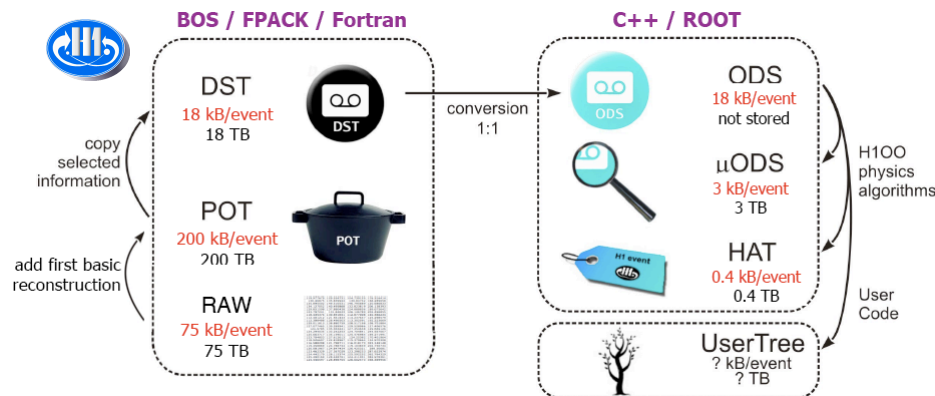
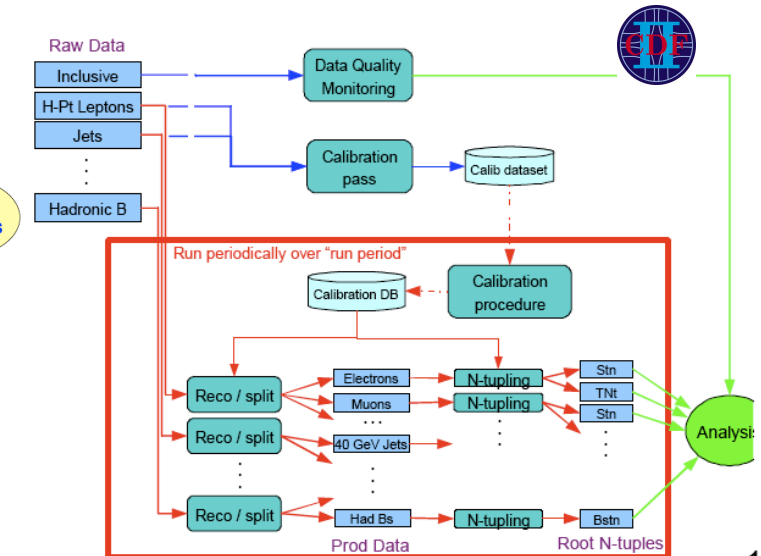
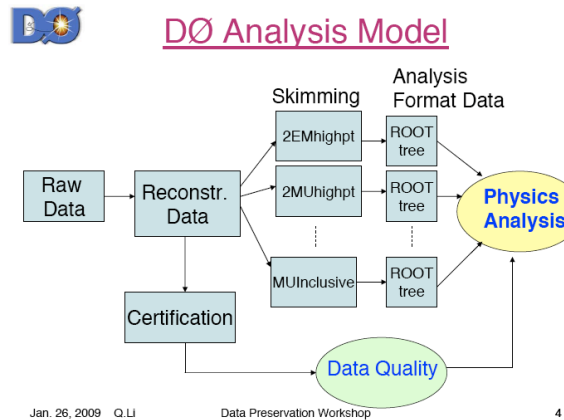
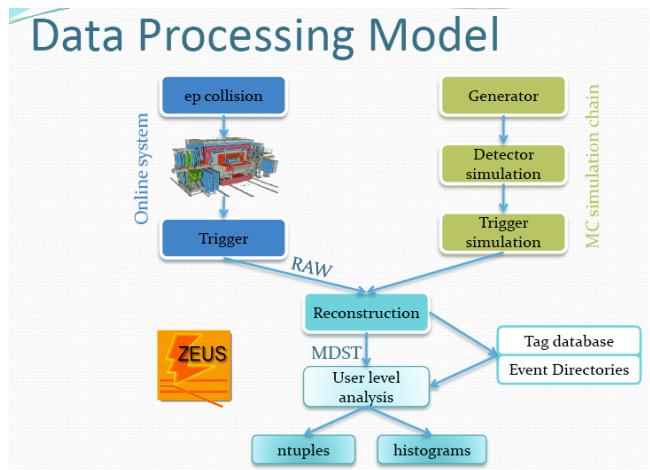


Centre-of-mass Energy	0.9 TeV	2.36 TeV
Selection	Number of Events	
BPTX Coincidence + one BSC Signal	72 637	18 074
One Pixel Track	51 308	13 029
HF Coincidence	40 781	10 948
Beam Halo Rejection	40 741	10 939
Beam Background Rejection	40 647	10 905
Valid Event Vertex	40 320	10 837

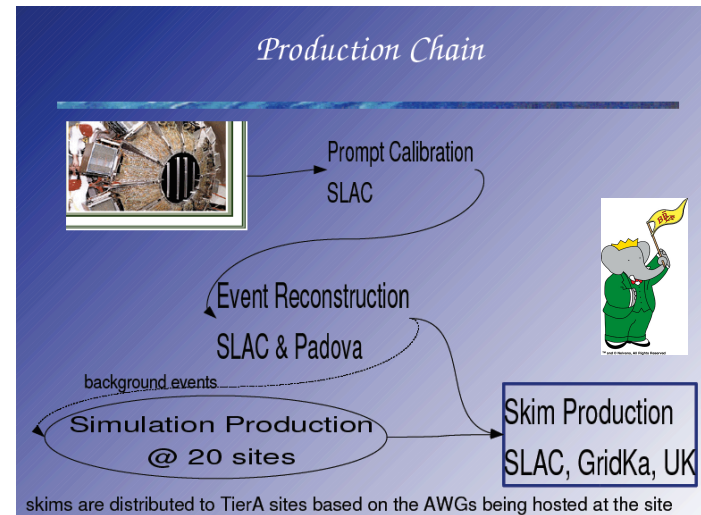
- Early LHC measurements made using data at a unique centre of masses
- 2010 low pile up 7 TeV data also at risk
- What happens when 14 TeV comes?



# Data analysis models in HEP



- Complicated, at first glance different
- Familiar descriptions of data analysis chain, from reconstruction to analysis level
  - RAW (→ POT) → DST → tuple → analysis

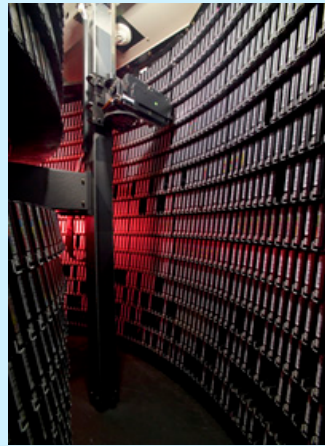


# Summary of information from the (pre-LHC) experiments

	BaBar	H1	ZEUS	HERMES	Belle	BESIII	CDF	DØ
<b>End of data taking</b>	07.04.08	30.06.07	30.06.07	30.06.07	30.06.10	2017	30.09.11	30.09.11
<b>Type of data to be preserved</b>	RAW data Sim/rec level Data skims in ROOT	RAW data Sim/rec level Analysis level ROOT data	Flat ROOT based ntuples	RAW data Sim/rec level Analysis level ROOT data	RAW data Sim/rec level	RAW data Sim/rec level ROOT data	RAW data Rec. level ROOT files (data+MC)	Raw data Rec. level ROOT files (data+MC)
<b>Data Volume</b>	2 PB	0.5 PB	0.2 PB	0.5 PB	4 PB	6 PB	9 PB	8.5 PB
<b>Desired longevity of long term analysis</b>	Unlimited	At least 10 years	At least 20 years	5-10 years	5 years	15 years	Unlimited	10 years
<b>Current operating system</b>	SL/RHEL3 SL/RHEL 5	SL5	SL5	SL3 SL5	SL5/RHEL5	SL5	SL5 SL6	SL5
<b>Languages</b>	C++ Java Python	C C++ Fortran Python	C++	C C++ Fortran Python	C C++ Fortran	C++	C C++ Python	C++
<b>Simulation</b>	GEANT 4	GEANT 3	GEANT 3	GEANT 3	GEANT 3	GEANT 4	GEANT 3	GEANT 3
<b>External dependencies</b>	ACE CERNLIB CLHEP CMLOG Flex GNU Bison MySQL Oracle ROOT TCL XRootD	CERNLIB FastJet NeuroBayes Oracle ROOT	ROOT	ADAMO CERNLIB ROOT	Boost CERNLIB NeuroBayes PostgresQL ROOT	CASTPR CERNLIB CLHEP HepMC ROOT	CERNLIB NeuroBayes Oracle ROOT	Oracle ROOT



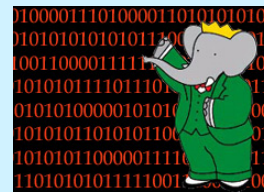
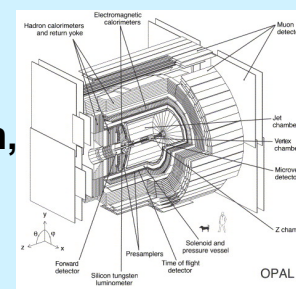
# What is HEP “data”?



**Digital information**  
The data themselves,  
volume estimates for  
preservation data of the  
order of **a few to 10 PB**

Other digital sources  
such as databases to  
also be considered

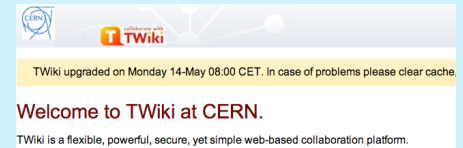
**Software**  
Simulation,  
reconstruction,  
analysis, user,  
in addition to  
any external  
dependencies



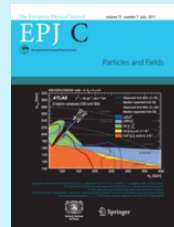
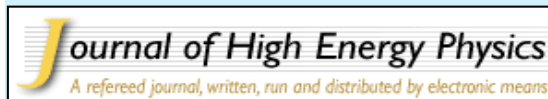
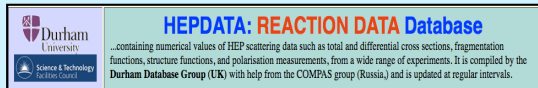
**CERNLIB Access**

- Access to the CERN Program Library is free of charge to all HEP users worldwide.
- Non-HEP academic and not-for-profit organizations: 1KSF/year

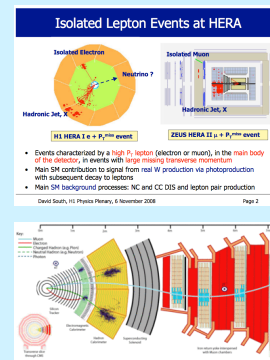
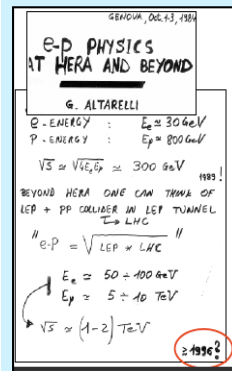
**Meta information**  
Hyper-news, messages,  
wikis, user forums..



**Publications** **arXiv.org**



**Documentation**  
Internal publications,  
notes, manuals, slides



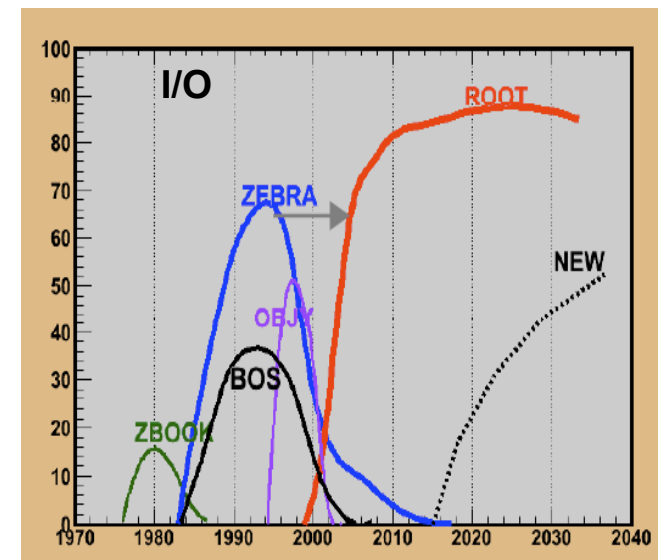
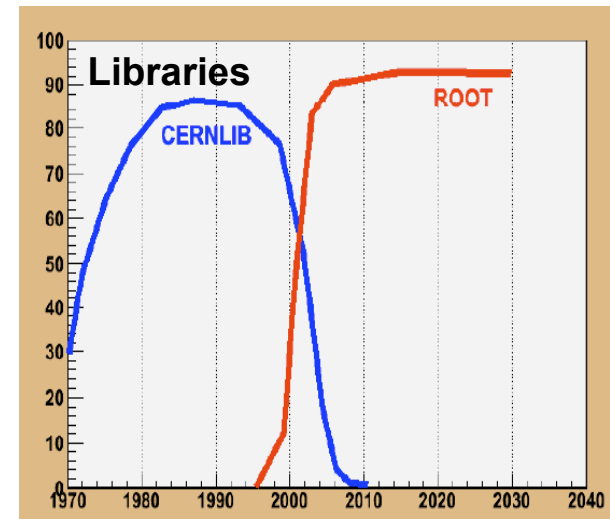
**Expertise and people**



# A serious issue: the software maintenance

R. Brun

- > Freezing: Technology preservation
  - Virtualisation techniques provide the software environment, freeze the hardware
  - Preparation step is not saved, lifetime limited as well
- > Better: Continuous migration
  - Follow technology changes, external software, new OS, redesign, recompile etc
  - Virtualisation can help here too
- > Preparation is not trivial
  - New operational model
  - Dependencies etc.
- > Supervision is needed for both data and software
  - Data archivist position



# DPHEP models of HEP data preservation

Preservation Model		Use Case
1	Provide additional documentation	Publication related info search
2	Preserve the data in a simplified format	Outreach, simple training analyses
3	Preserve the analysis level software and data format	Full scientific analysis, based on the existing reconstruction
4	Preserve the reconstruction and simulation software as well as the basic level data	Retain the full potential of the experimental data

Increasing cost,  
complexity and benefits

- > These are the original definitions of DPHEP preservation levels from the 2009 publication
  - Still valid now, although interaction between the levels now better understood



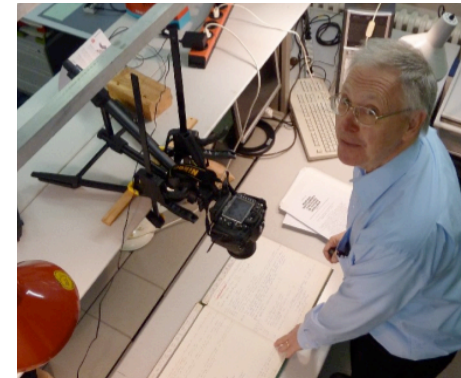
# DPHEP models of HEP data preservation

Preservation Model		Use Case	
1	Provide additional documentation	Publication related info search	Documentation
2	Preserve the data in a simplified format	Outreach, simple training analyses	Outreach
3	Preserve the analysis level software and data format	Full scientific analysis, based on the existing reconstruction	Technical Preservation Projects
4	Preserve the reconstruction and simulation software as well as the basic level data	Retain the full potential of the experimental data	

- > These are the original definitions of DPHEP preservation levels from the 2009 publication
  - Still valid now, although interaction between the levels now better understood
- > Originally idea was a progression, an inclusive level structure, but now seen as complementary initiatives
- > Three levels representing three areas:
  - Documentation, Outreach and Technical Preservation Projects

# Level 1: Documentation

- > The organisation of documentation turns out to be quite a task
  - Dedicated task forces set up by many of the experiments
  - Much material from pre-web days, or using all kinds of web applications
- > **Non-digital:** Cataloguing, organisation, scanning or photographing of appropriate of papers, notes, drawings, talks from pre-web days, detector schematics, blueprints, logbooks, ...
  - New *Virtual Archives* established by the experiments
- > **Digital:** Old online shift tools, detector configuration files, electronic logbooks, detailed run information, web content from out-dated servers with dead links, various wikis, meetings, talks, ...
  - Replacement of old web servers by VMs, hosted by the computer centres
  - Replacement of old pages to newer technologies such as wikis (use of (T)wikis much more prevalent in the LHC era)
  - Use of external services for hosting collaboration material



# Documentation projects with INSPIRE

- Internal notes from all HERA experiments now available on INSPIRE
  - Experiments no longer need to provide dedicated hardware for such things
  - Password protected now, simple to make publicly available in the future



# Documentation projects with INSPIRE

- Internal notes from all HERA experiments now available on INSPIRE
  - Experiments no longer need to provide dedicated hardware for such things
  - Password protected now, simple to make publicly available in the future



The screenshot shows the INSPIRE login interface. At the top, there is a blue header with the INSPIRE logo and a navigation bar with links: HEP, INST, HELP, SPIRES, and HEPNAMES. Below the header, a yellow banner reads: "Welcome to INSPIRE! INSPIRE is out of beta and ready to replace SPIRES. please email us at [feedback@inspirehep.net](mailto:feedback@inspirehep.net)". The main content area is titled "Login" and contains a message: "This collection is restricted. If you think you have right to access it, please authenticate yourself." Below this, there are input fields for "Username:" (containing "zeus") and "Password:". A checkbox labeled "Remember login on this computer." is present, followed by a "login" button and a link "(Lost your password?)". A "Note" states: "You can use your nickname or your email address to login." The footer includes links for "HEP :: Search :: Help", "Powered by Invenio v1.0.0-rc0+", and "Problems/Questions to [feedback@inspirehep.net](mailto:feedback@inspirehep.net)".



# Documentation projects with INSPIRE

- Internal notes from all HERA experiments now available on INSPIRE
  - Experiments no longer need to provide dedicated hardware for such things
  - Password protected now, simple to make publicly available in the future

The screenshot displays the INSPIRE website interface. At the top, there is a blue header with the INSPIRE logo and a welcome message: "Welcome to INSPIRE! INSPIRE is out of beta and ready to replace SPI". Below this, a search bar is visible. The main content area shows a search result for "ZEUS Internal Notes" with "10 records found". The results are listed as follows:

- 1. Inclusive-jet production in NC DIS with HERA II.**  
J. Terron C. Glasman. ZEUS-IN-09-004.  
[References](#) | [BibTeX](#) | [LaTeX\(US\)](#) | [LaTeX\(EU\)](#) | [Harvmac](#) | [EndNote](#)  
[Detailed record](#) - [Similar records](#)
- 2. Three-subjet distributions in neutral current deep inelastic scattering.**  
E. Ron C. Glasman, J. Terron. ZEUS-IN-09-003.  
[References](#) | [BibTeX](#) | [LaTeX\(US\)](#) | [LaTeX\(EU\)](#) | [Harvmac](#) | [EndNote](#)  
[Detailed record](#) - [Similar records](#)
- 3. 2009 Guide to Funnel: The ZEUS Monte Carlo Production Facility.**  
A. Parenti. ZEUS-IN-09-002.  
[References](#) | [BibTeX](#) | [LaTeX\(US\)](#) | [LaTeX\(EU\)](#) | [Harvmac](#) | [EndNote](#)  
[Detailed record](#) - [Similar records](#)
- 4. Automated calculation of radiative correction to electron-proton charged current DIS at HERA.**  
I. Marfin. ZEUS-IN-09-001.  
[References](#) | [BibTeX](#) | [LaTeX\(US\)](#) | [LaTeX\(EU\)](#) | [Harvmac](#) | [EndNote](#)  
[Detailed record](#) - [Similar records](#)

On the left side of the screenshot, there is a sidebar with a "Log" button and a "Note" section. At the bottom left, there is a small box with the text "HEP :: Search Powered by Problems/Queue Last update".

# Documentation projects with INSPIRE

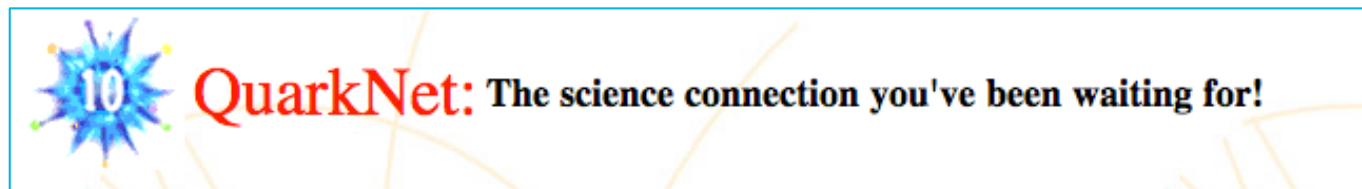
- Internal notes from all HERA experiments now available on INSPIRE
  - Experiments no longer need to provide dedicated hardware for such things
  - Password protected now, simple to make publicly available in the future



- The ingestion of other documents is under discussion, including theses, preliminary results, conference talks and proceedings, paper drafts, ...
  - More experiments working with INSPIRE, including CDF, D0 as well as BaBar

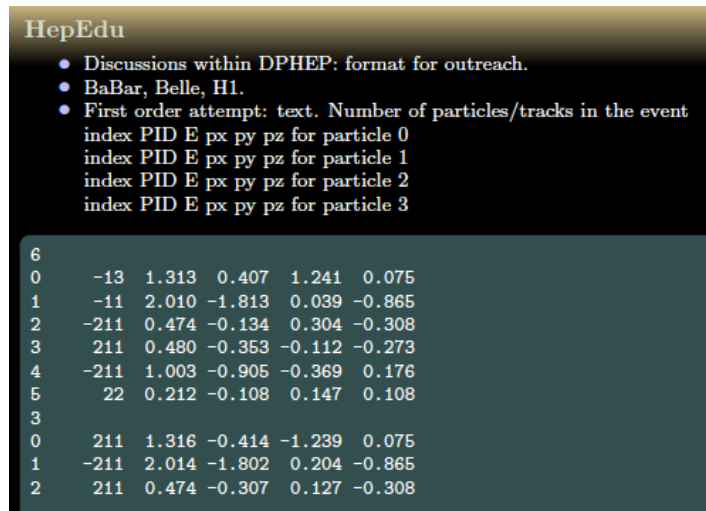
# HEP outreach initiatives

- Many initiatives promoting outreach efforts and to improve the public understanding of science in general



# Outreach

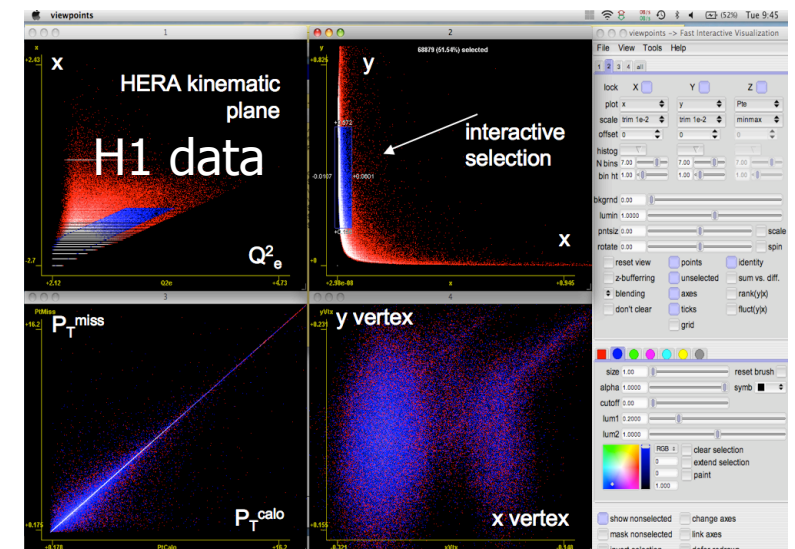
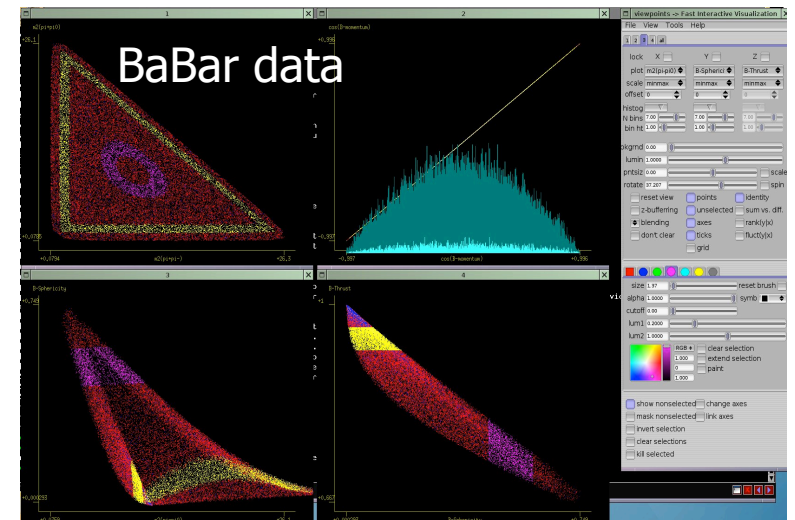
- > Use **real and preserved** data to enhance HEP education worldwide
- > Simple data format: input using text file of kinematics of HEP events



- > Discussions about common formats ongoing
  - B-lab (KEK) example considered
  - Experience at LHC
  - Connect to existing projects (master classes etc.)
- > A lot of potential, but little activity at multi-experiment level



## Viewpoints (NASA)





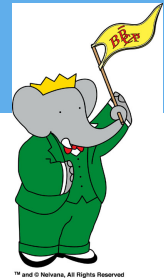
# Technical Projects: DPHEP preservation levels 3 and 4

- > This is really the main focus of the data preservation effort
  - Level 3: Access to analysis level data, MC and the analysis level software
  - Level 4: Access to reconstruction and simulation software, retain the full capability
- > Deciding on level 3 or 4 depends on the scope of your project
  - What do you want to be able to do in  $N$  years time?
  - Only level 4 gives full flexibility, but this also means not relying on frozen executables and binaries but rather retaining the ability to recompile: more work

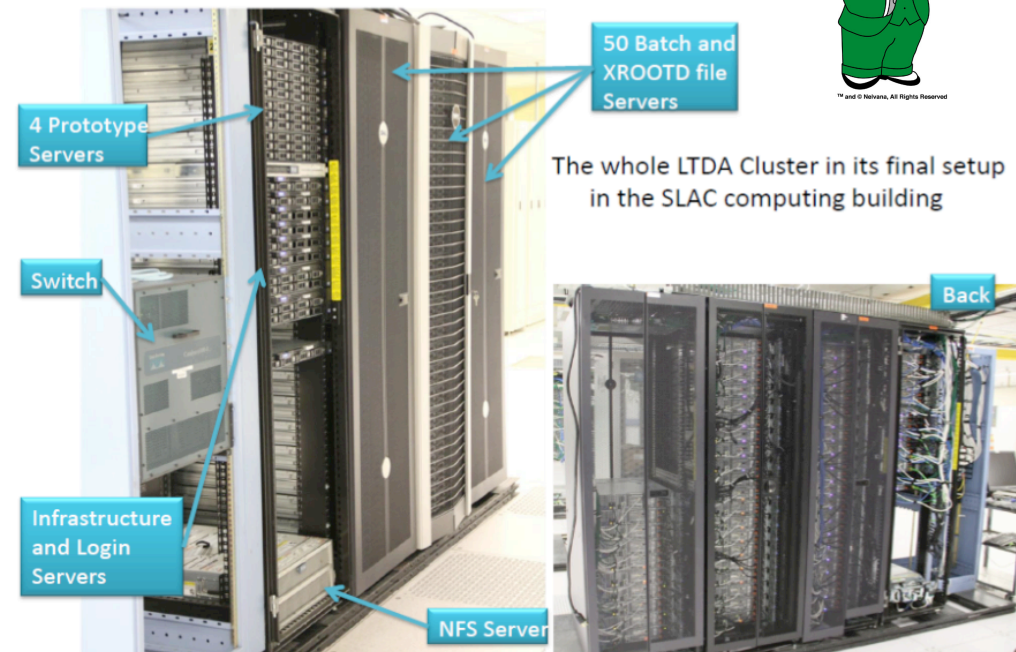
**The majority of DPHEP experiments aim for DPHEP level 4 preservation**

- > Remember: it's not about the data, but about still being able analyse it
  - Either keep your current environment alive as long as possible
  - Or adapt and validate your code to future changes as they happen
- > Two complimentary approaches taken at SLAC and DESY
  - Both employing virtualisation techniques, but in rather different ways

# The BaBar Long Term Data Access archival system

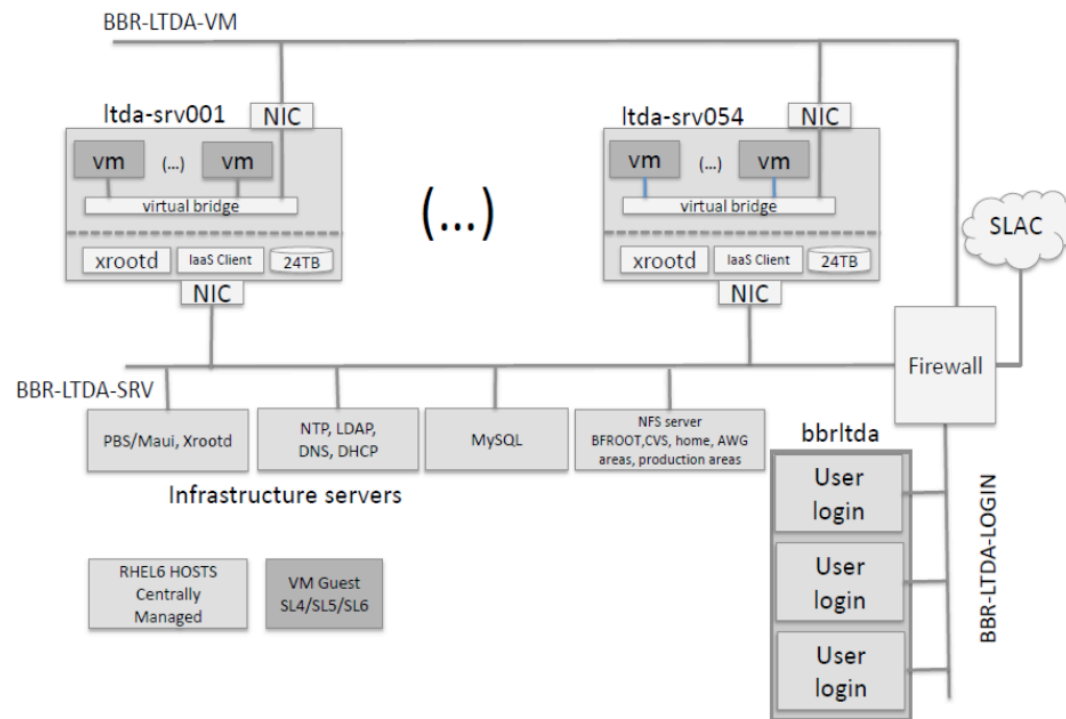


- > New BaBar system installed for analysis until at least 2018
- > Isolated from SLAC, and uses virtualisation techniques to preserve an existing, stable and validated platform
- > Complete data storage and user environment in one system



- > Required large scale investment: 54 R510 machines, primarily for data storage, as well as 18 other dedicated servers
  - Resources taken into account in experiment's funding model during analysis phase!
- > From the user's perspective, very similar to existing BaBar infrastructure

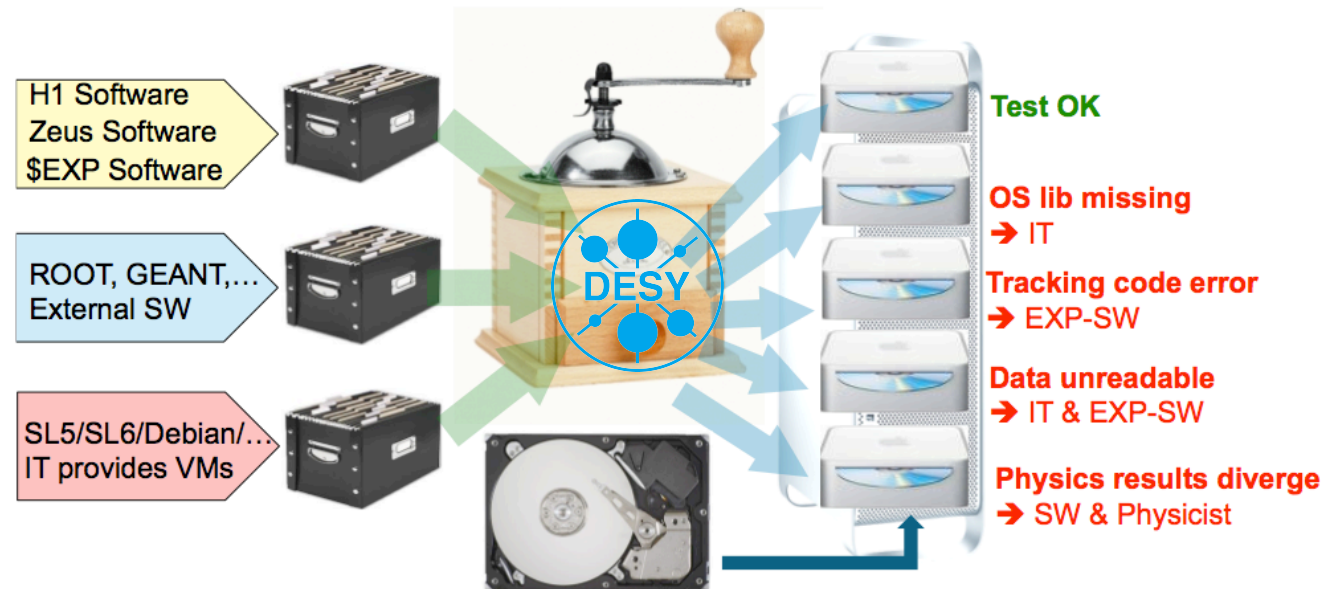
# The BaBar Long Term Data Access archival system



- > Crucial part of design is to allow frozen, older platforms to run in a secure computing environment
- > *Naïve* virtualisation strategy, not enough
  - Cannot support an OS *forever*
  - Security of system under threat using old versions

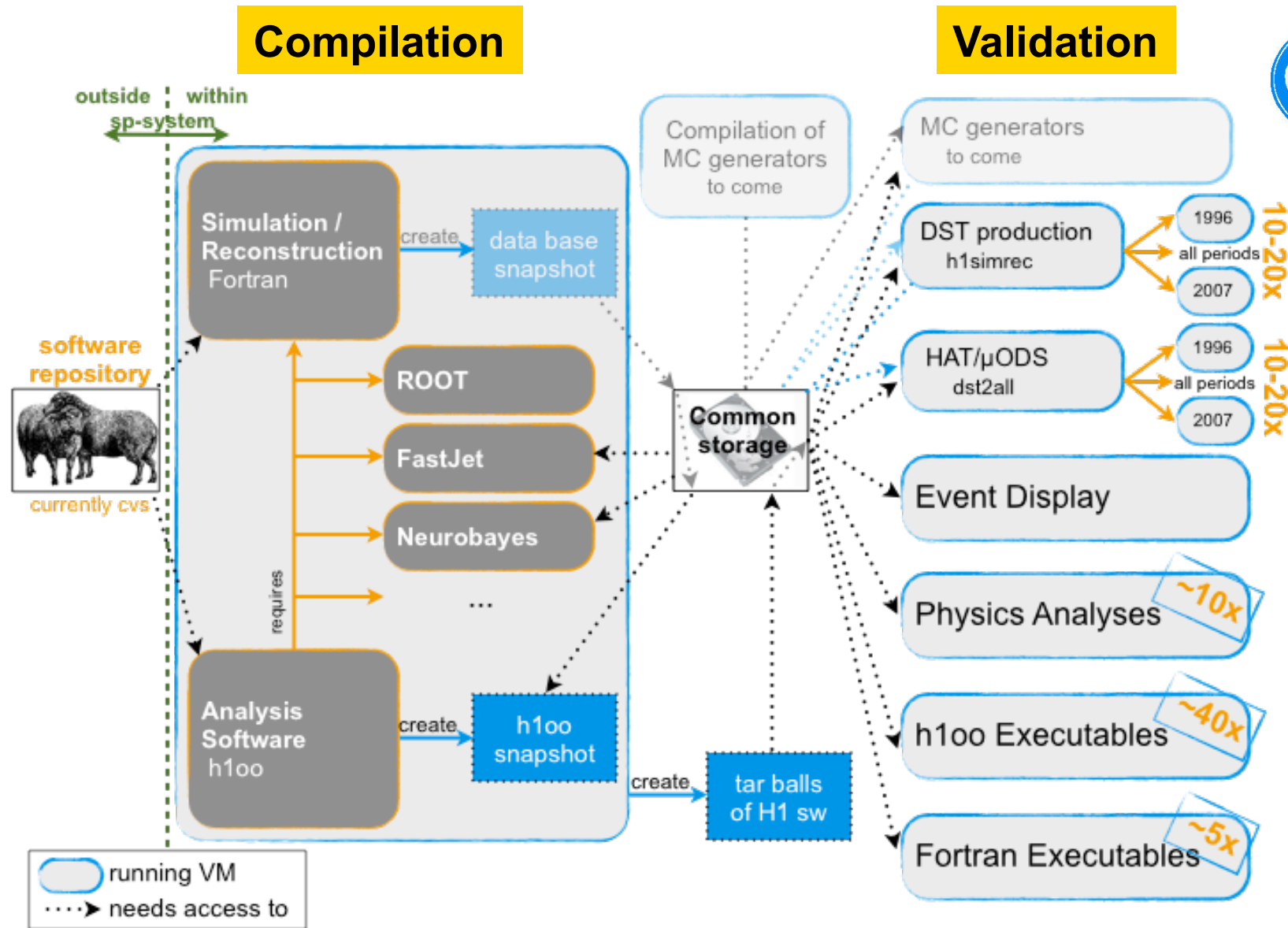
- > Achieved by clear network separation via firewalls of part storing the data (more modern OS) and part running analysis (the desired older OS)
- > Other BaBar infrastructure not included in VMs is taken from common NFS
- > More than 20 analyses now using the LTDA system as well as simulation

# The sp-system at DESY



- > Automated validation system to facilitate future software and OS transitions
  - Utilisation of virtual machines offers flexibility: OS and software configuration is chosen by experiment controlled parameter file
  - Successfully validated recipe to be deployed on future resource, e.g. Grid or IT cluster
  - Pilot project at CHEP 2010, full implementation now installed at DESY
- > Essential to have a robust definition of a complete set of experimental tests
  - Nature and number dependent on desired preservation level

# Example structure of the experimental tests: H1 (Level 4)

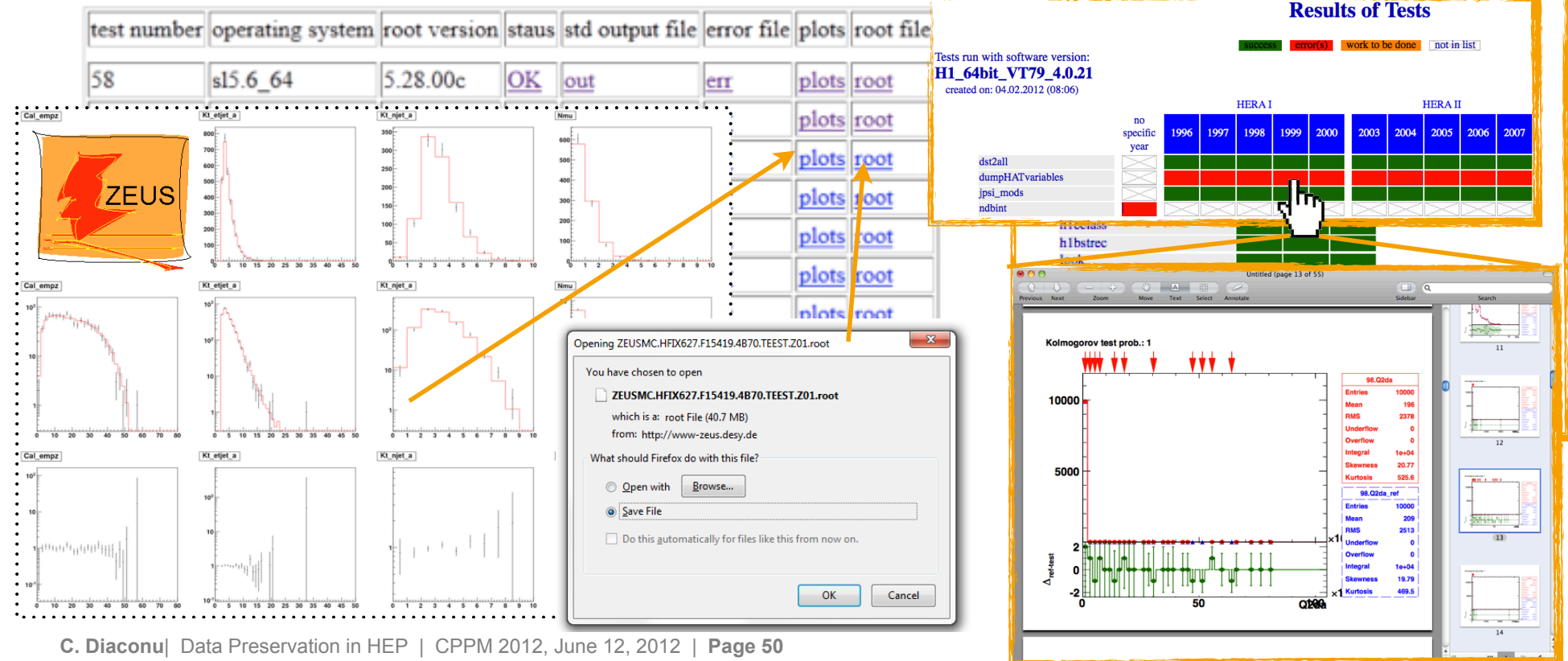


Including compilation of individual packages: about 250 tests planned by H1



# Digesting the validation results

- Display the results of the validation in a comprehensible way: web based interface
- The test determines the nature of the results
  - Could be simple yes/no, plots, ROOT files, text-files with keywords or length, ...






# Current status of the HERA experiments software

## > Common baseline of SLD5 / 32-bit achieved in 2011 by all experiments

- Validation of 64-bit systems is a major step towards migrations to future OS
- The system has already been useful in detecting problem visible only in newer software

## > Note that this system does not concern data integrity

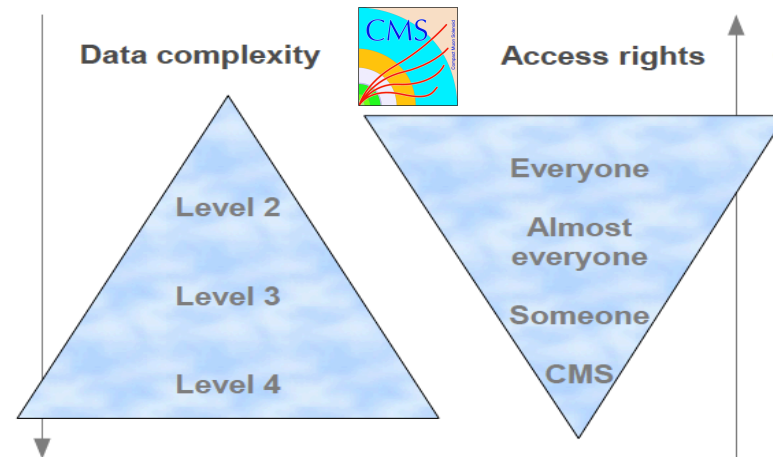
- The investigation into data archival options is underway

Process \ External Dependencies		SL5 32bit	ROOT				SL5 64bit						SL6 64bit		
		← Reference	5.26	5.28	5.30	5.32	Adamo	Cernlib		Fastjet	Neuro- 0312 bytes	Neuro- 3.3.0 bytes			
								2005	2006	2.3.3	2008				
	Accessing common ntuples						No dependence								
	ZMCSP (simulate/reconstruct MC)														
	Creating common ntuples														
	Validation														
	Compilation of s/w							No dependence							
	Generating MC files														
	Producing DST files														
	Producing h1oo files														
	Accessing h1oo files														
	Accessing ndb snapshot														
	Validation														
	Compilation of s/w														
	MC generation & digitisation														
	Reconstruction														
	Producing uDST														
	Analysing uDST (Fortran, HANNA++)														
	Validation														

ok
ongoing
not done
problem

# Data Preservation at the LHC

- > Reflection just started in ATLAS, ALICE, CMS, LHCb
  - Common understanding that starting earlier will consolidate the long term future
  - Strong wish to develop a common policy at CERN and within DPHEP
  - Specific cases already identified: Lower energy data, trigger configurations, pile up.
- > In terms of documentation, LHC experiments are in good shape
  - The electronic era: Twikis, accompanying notes, plans for extended use of INSPIRE
- > Outreach projects and open access explored
- > The distributed data model eases the worry of data loss
  - Although as previously stated: no successful preservation without associated long-term access
  - No concrete plans yet, but level 4 seen as the ultimate objective



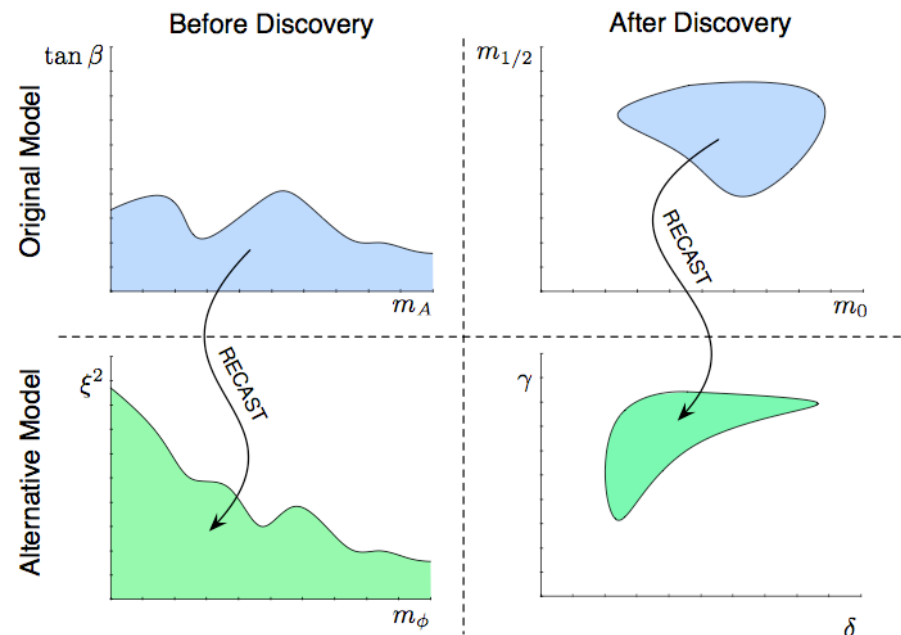
# A multi-preservation level tool: RECAST

arXiv:1010.2506

- > Framework developed to extend impact of existing analyses
- > Complementary approach of analysis archival, encapsulating the full event selection, data, backgrounds, systematics

- > Idea is to **recast** existing physics search results to constrain alternate model scenarios

- Complete information from original analysis contained in the data
- Already performed on ALEPH data, LHC experiments investigating



- > RECAST does not fit directly into the DPHEP preservation levels
  - Levels 3 and 4 are in the back-end, containing the complete archived analyses
  - However, only the selection in the publication is preserved, it could also be described as additional information, more like level 1



## ATLAS DMP Organization



- ◆ Data Preservation now included as part of the upgrade activity planning
  - ◆ May increase the funding options
  - ◆ Data Management Planning will be required by some funders for upgrade grants
  - ◆ Looking at the cost/benefit of various strategies
  - ◆ Resource tensioning with other upgrade activities





## ALICE DATA PRESERVATION POLICY



- ALICE recently initiated internal discussions on data preservation issues
  - Due to the operational requirement related to current data taking and processing, only a few people have been actively working on this issue.
  - The ALICE MB is being kept informed regarding this activity.
  - The current mandate is to define a policy that defines ALICE position with respect to the different preservation levels...
- ALICE participates now regularly to the internal LHC data preservation meetings
  - ALICE strongly supports the harmonization of the base principles and the development of common tools whenever possible...
  - These meeting are providing the starting points for our internal discussions
- We foresee the finalization of a draft for ALICE DP policy by the end of this year's running periods.

- ☐ Data preservation is upon us – and its common sense we initiate at least a requirements study – and preferably identify some effort to take it forward.
- ☐ LHCb is participating in the CERN+LHC wide discussions + DPHEP
- ☐ Open data access is out there – it may not be imminent in terms of requests – but it has the weight of funders and legislation
- ☐ It is far better to have a policy than not have one.
- ☐ We have presented a draft policy to focus the next step in iteration of the LHCb position.

## Collaboration with DPHEP, CERN and other LHC experiments

- DPHEP
  - Great benefit of having concepts (data levels etc.) already formulated.
  - Excellent forum to follow the experience from other experiments.
- CERN
  - CMS has relied on CERN expertise on data preservation and access in preparing this policy.
- Informal discussions with other LHC experiments at CERN:
  - Discussing and converging to a common view on this topic at CERN.
- CMS may have been the first LHC experiment to come up with a policy, but without the collaboration with others we would not yet be even there. **Thank you!**
- CMS hopes that this fruitful exchange of ideas and experience will continue and help us implementing this policy.

# Other experiments

- Tevatron experiments: transition to final analysis phase

The fate of the final data is being discussed (10 years retention is mentioned)  
Discussion on DP preservation intensified in experiments, formats, procedures etc.  
both CDF and D0 start task forces

- The issue is considered in Belle

For common work with Babar and transition of the analysis to Belle 2  
Goal: preserve Belle I data to 2017=> therefore a technological step needed!

- Other initiatives

CLEO archival initiative, BES3/IHEP as a part of the ongoing run  
JLAB/HallD: a physics case for combination explored  
BNL contact to experiments to be established

- LEP data:

Publications still produced in 2011, activity around data ongoing (Higgs results, RECAST)  
A recovery is still possible: but urgent action is needed

**> „If the same (loss of data as at PETRA) will happen with LEP data, I will sue the CERN DG“** (A well-known theorist after having seen reanalysed JADE results)

# Transition scenario and resources at the experimental level

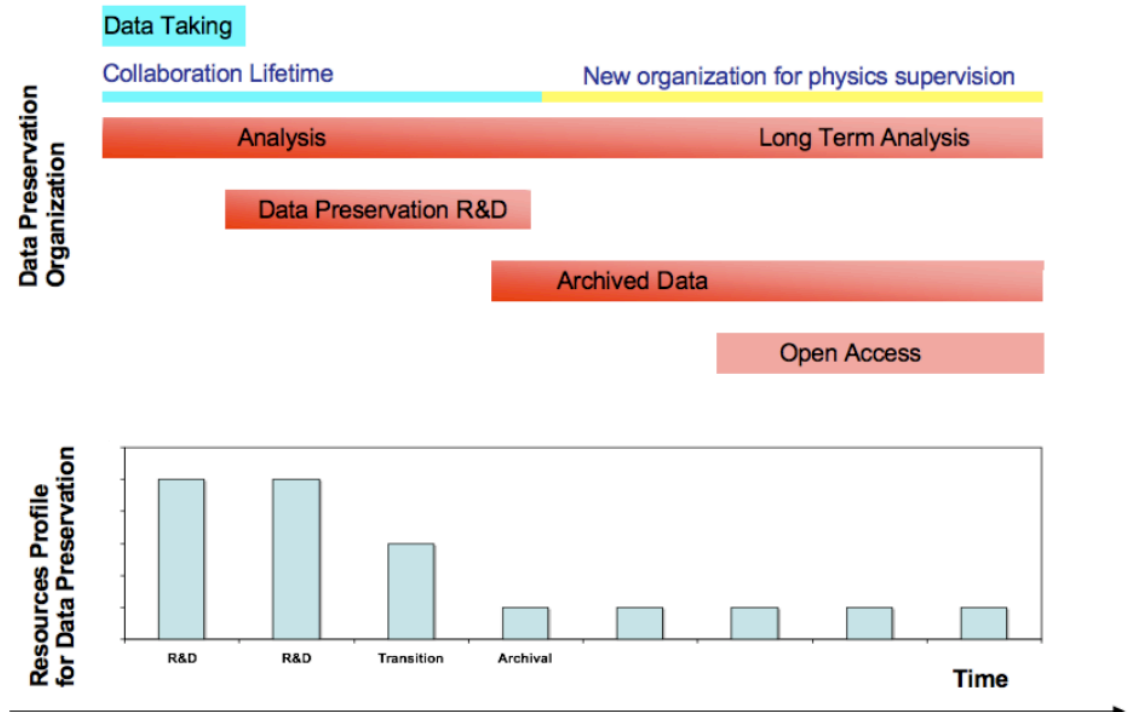
> Planning the transition to a long term analysis model

> R&D phase needed to develop the projects for the transition

> Long term custodianship of the physics data

> Resources / experiment

- Typically a surge of 2-3 FTEs for 2-3 years, followed by steady 0.5-1.0 FTE per experiment/lab
- This should be compared to 300-500 FTEs for many years per experiment!



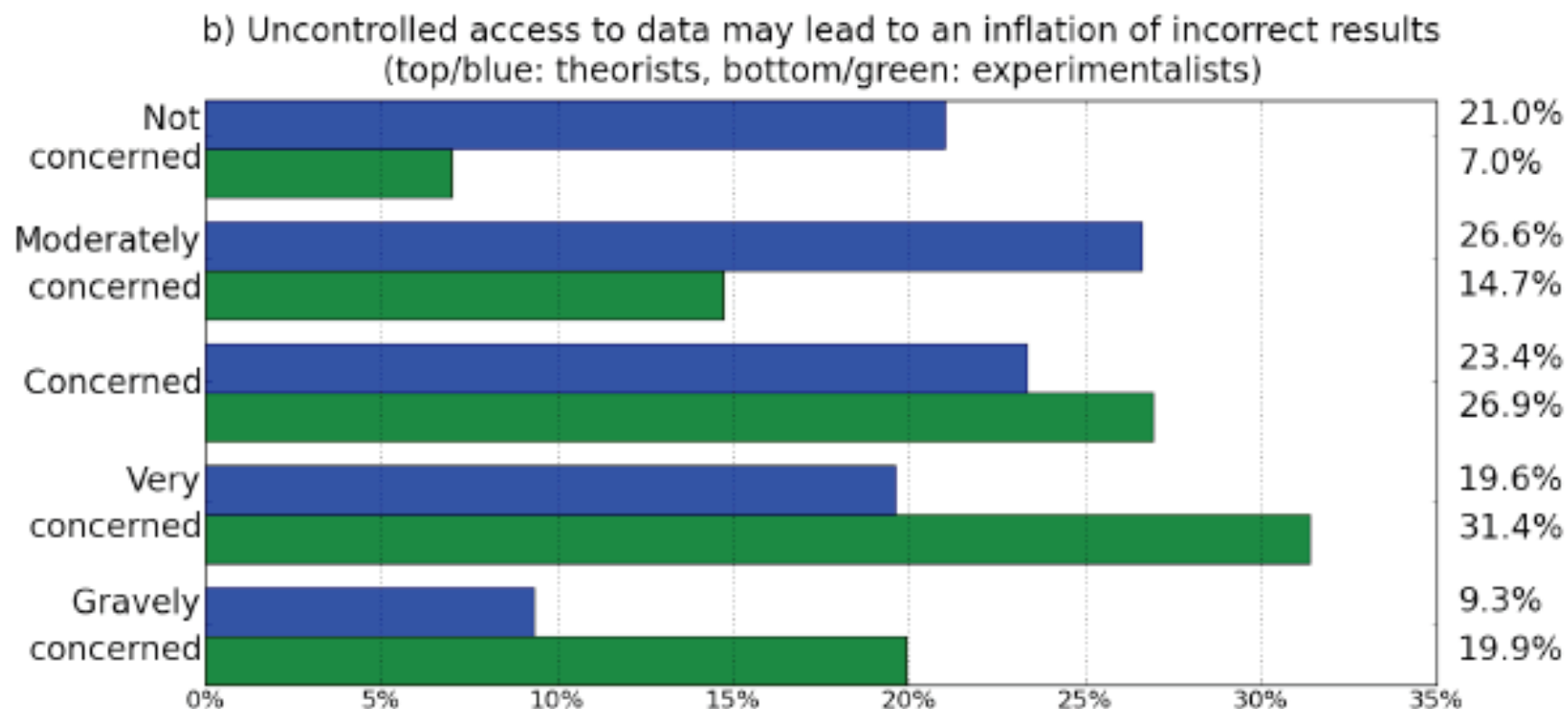
Cost estimates represent typically **much less than 1%** of the original investment

Scientific return: **O(10%)** in number of publications



# Risks of re-use?

Parse.insight



**Governance issues are very important**

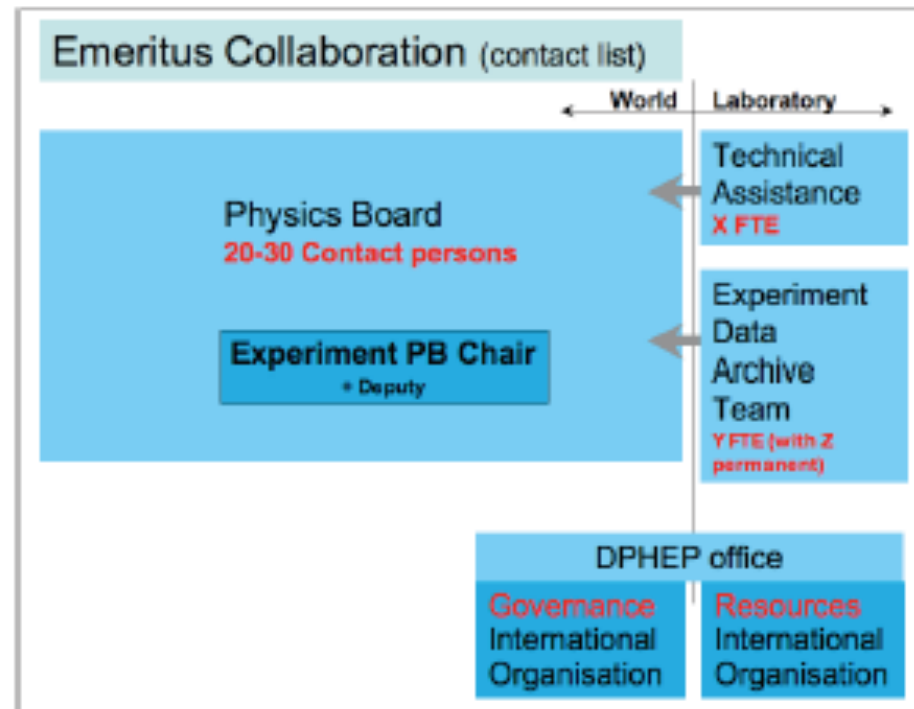
**"Errors using inadequate data are much less than those using no data at all."  
Charles Babbage**

# Long term organisation

Preservation project make sense if the scientific supervision is ensured

> Future structure of the collaboration should also be considered by HEP experiments

- Experimental organisation risks being left in an undefined state
- Transition should also be planned in advance of the projected end date

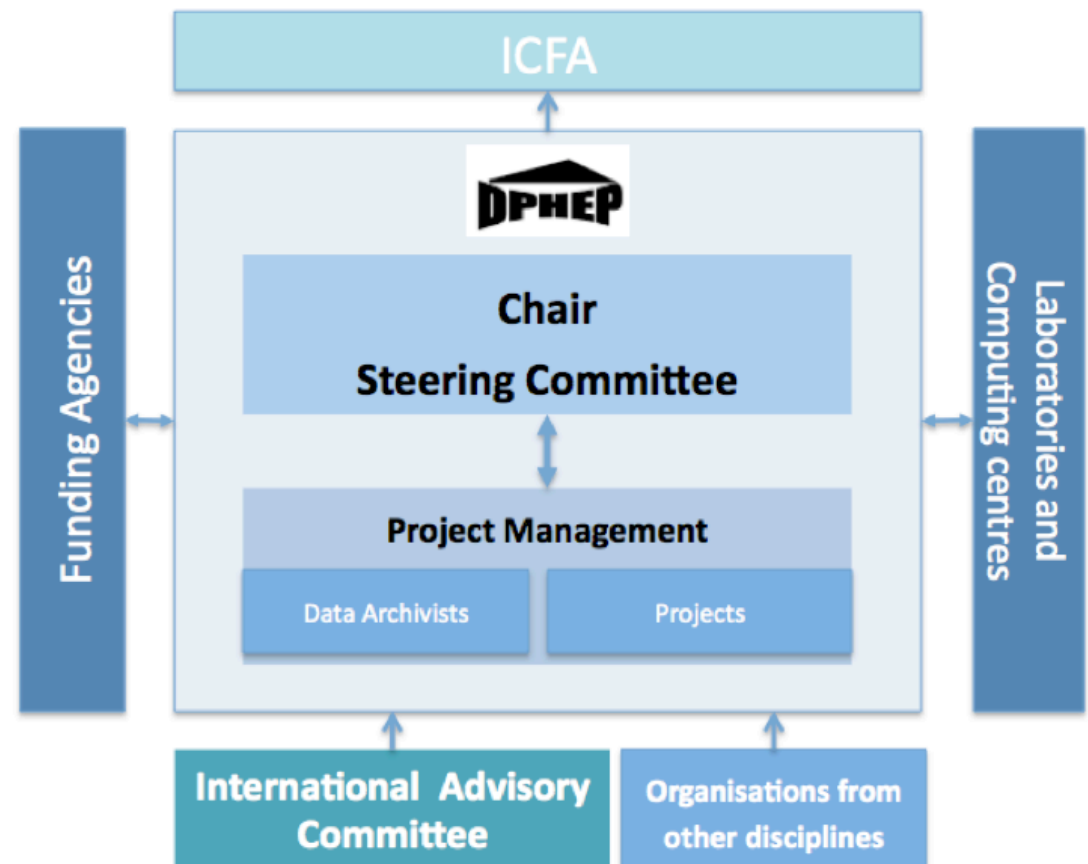


## Summing up: What has been achieved so far?

- > The DPHEP Study Group represents the first large scale effort to address data preservation in the field of high energy physics
- > The initial make up of the group was driven by the coincidence of the end of data taking at several large colliders, but had grown to include others including the LHC experiments
- > The activity of the group over the last three years has led to an increased understanding of the relevant issues, enabling problems to be addressed, recommendations to be formulated and multi-experiment projects to begin
- > To gain the most benefit from the work done so far, a transition from the current Study Group structure to a new, full time DPHEP Organisation

# The DPHEP Organisation

- Retain the basic structure of the Study Group, with links to the host experiments, labs, funding agencies, ICFA
- Installation of a full time DPHEP Project Manager, who acts as the main operational coordinator
- The DPHEP Chair (appointed by ICFA) coordinates the steering committee and represents DPHEP in relations with other bodies



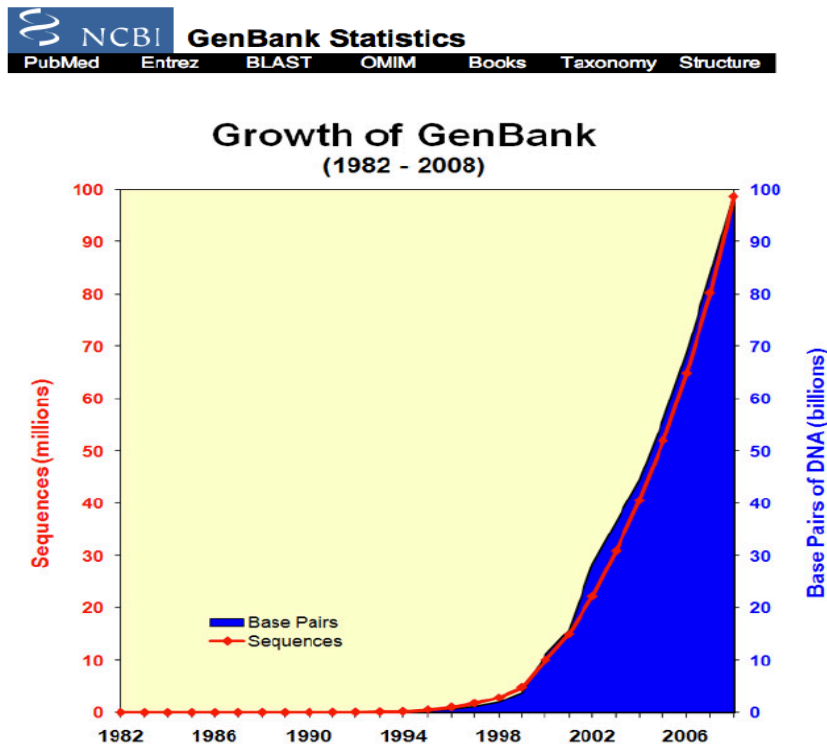
# DPHEP person power requirements

	Project	Goals and deliverables	Resources and timelines	Location, possible funding source, DPHEP allocation
Experiment and laboratory Priority: 1	Experimental Data Preservation Task Force	Install an experiment data preservation task force to define and implement data preservation goals.	1 FTE installed as soon as possible, and included in upgrade projects	Located within each computing team. Experiment funding agencies or host laboratories. DPHEP contact ensured, not necessarily as a displayed FTE.
	Facility or Laboratory Data Preservation Projects	Data archivist for facility, part of the R&D team or in charge with the running preservation system and designed as contact person for DPHEP.	1-2 FTE per laboratory, installed as a common resource.	Experiment common person power, support by the host labs or by the funding agencies as a part of the on-going experimental program. A fraction 0.2 FTE allocated to DPHEP for technical support and overall organisation.
Multi-experiment Priority: 3	General validation framework	Provide a common framework for HEP software validation, leading to a common repository for experiments software. Deployment on grid and contingency with LHC computing also part of the goals.	1 FTE	Installed in DESY, as present host of the corresponding initiative. Funding from common projects. Cooperation with upgrades at LHC can be envisaged. Part of DPHEP.
	Archival systems	Install secured data storage units able to maintain complex data in a functional form over long period of time without intensive usage.	0.5 FTE	Multi-lab project, cooperation with industry possible. Included in DPHEP person power.
	Virtual dedicated analysis farms	Provide a design for exporting regular analysis on farms to closed virtual farm able to ingest frozen analysis systems for a 5-10 years lifetime.	1 FTE	The host of this working group should be SLAC. Funding could come from central projects and can be considered as part of DPHEP.
	RECAST contact	Ensure contact with projects aiming at defining interfaces between high-level data and theory.	0.5 FTE	Installed with proximity to the LHC, the main consumer of this initiative, with strong connections to the data preservation initiatives that may adopt the paradigms.
	High level objects and INSPIRE	Extend INSPIRE service to documentation and high-level data object.	0.5-1.5 FTE	Installed at one of the INSPIRE partner laboratories.
	Outreach	Install a multi-experiment project on outreach using preserved data, define common formats for outreach and connect to the existing events.	1 FTE central + 0.2 FTE per experiment	A coordinating role can be played by DPHEP in connection with a large outreach project existing at CERN, DESY or FNAL. The outreach contributions from experiments and laboratories can be partially allocated to the common HEP data outreach project and steered by DPHEP.
Global Priority: 2	DPHEP Organisation	DPHEP Project Manager	1 FTE	A position jointly funded by a combination of laboratories and agencies.

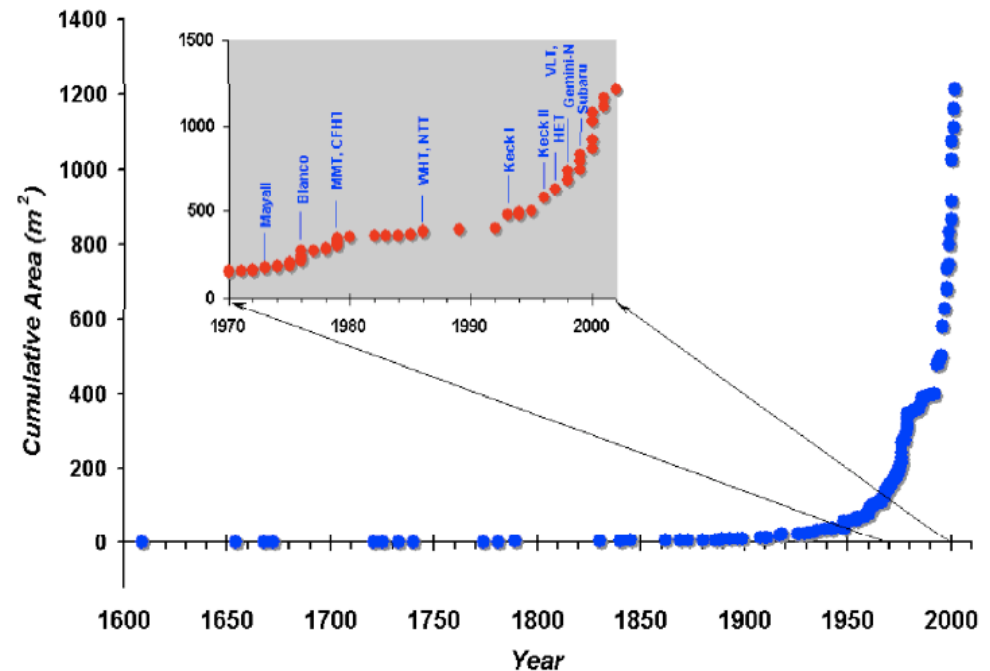


# We are not alone....

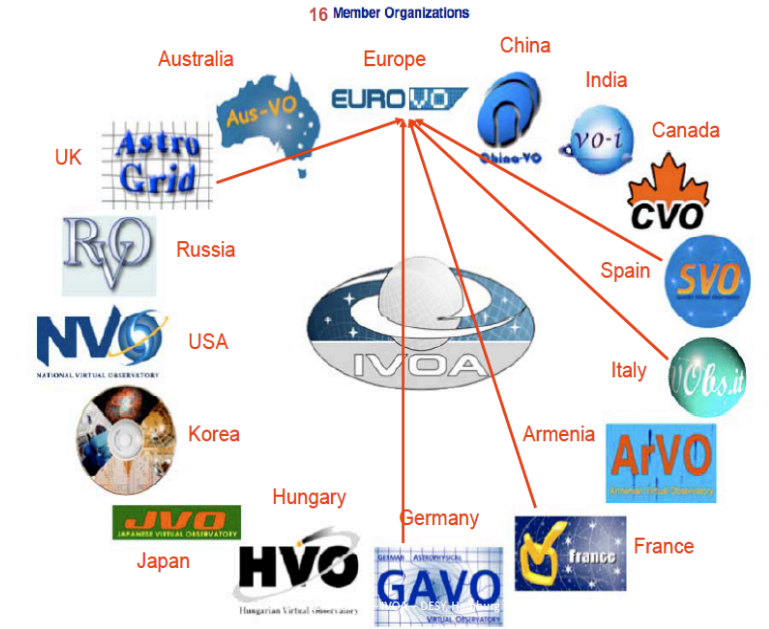
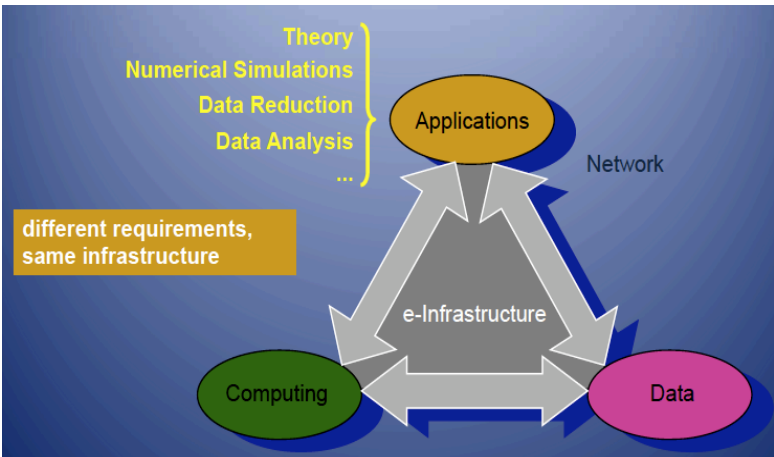
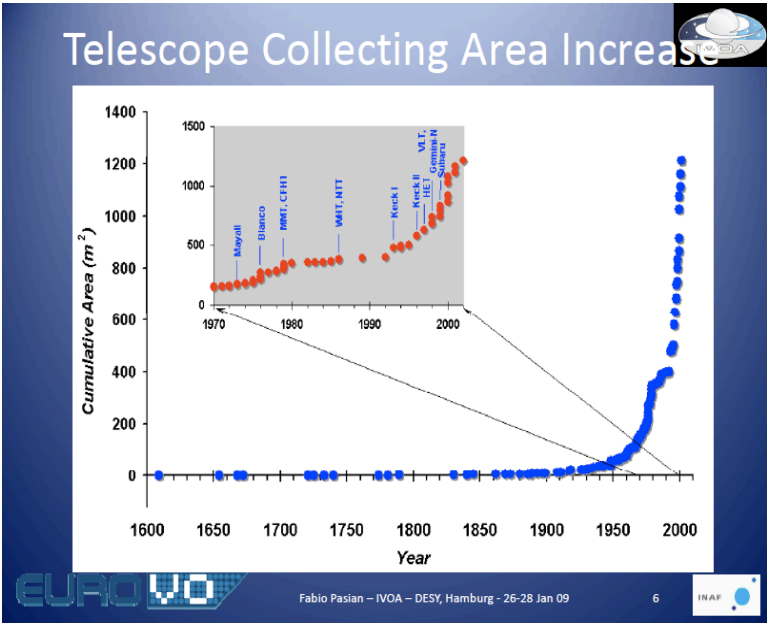
- Other fields observe a dramatic increase in data and are questioning the long term future of this data



## Telescope Collecting Area



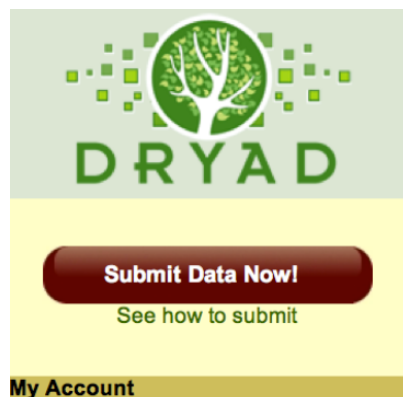
# Virtual Observatories in Astrophysics



- > Data Archives Inter-operable
- > Work on standards and access to
  - Data, simulation, mining techniques
- > International, multi-experiment
- > Agregated Person-power: about 100FTE

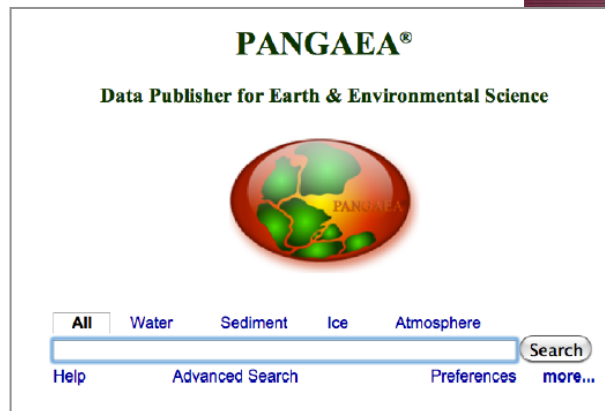
# Initiatives in other fields

- Data preservation and in particular open access and data sharing are present in other fields such as:
  - Astrophysics, molecular biology, earth sciences, humanities and social sciences



Blue Ribbon Task Force  
on Sustainable Digital Preservation and Access

[About Us](#) | [Members](#) | [Publications](#) | [Bibliography](#) | [News Center](#) | [Intranet](#)



[Home](#) | [News](#) | [Docs](#) | [WCS](#) | [Samples](#) | [Libraries](#) | [Viewers](#) | [Utilities](#) | [Keywords](#) | [Conventions](#) | [Resources](#)

The FITS Support Office

at NASA/GSFC



# Generic arguments

- Task forces already in place to address this issue in a generic way (standards)

- e.g. Blue Ribbon, APA, DPC, eSciDir, ...

<http://www.alliancepermanentaccess.eu>  
<http://brtf.sdsc.edu>  
(intermediate report and references)

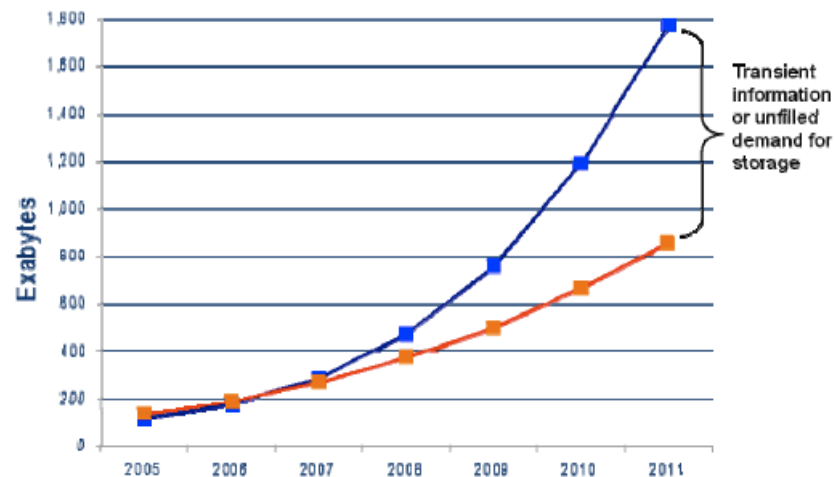


FIGURE 1.3: Information and Storage

Source: J. Gantz January 2008 (revised). Used with permission.

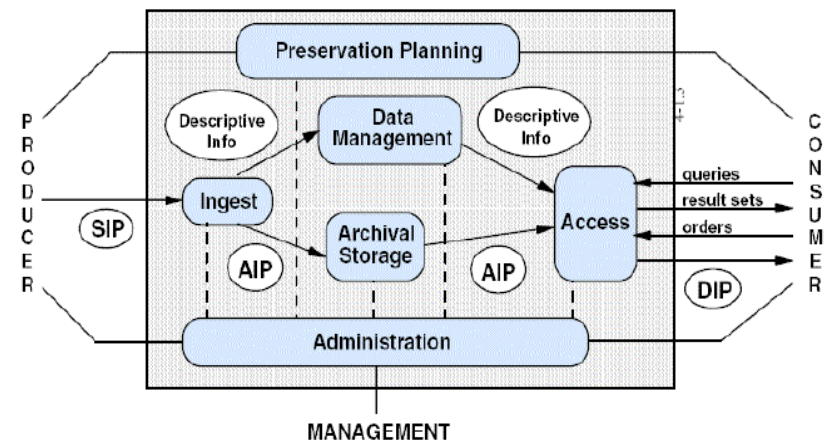


FIGURE 2.1: The OAIS Reference Model

<http://public.ccsds.org/publications/archive/650x0b1.pdf>, Page 4-1.

Source: Consultative Committee for Space Data Systems January 2002.

- Scientific Data is a major component of the ongoing efforts (complexity)
- Some scientific fields are well advanced : astrophysics

## ▶ Le 4<sup>ème</sup> Paradigme



Vision proposée par Jim Gray (Microsoft) en 2007

➔ Les données deviennent le vecteur principal de la Science

Astronomie

Génomique

Physique des particules

Sciences de la Terre

Réseaux de capteurs

Climat

Même les supercalculateurs deviennent des sources de données

Concept repris par Alex Szalay (Johns Hopkins University) qui propose la création d'un instrument nommé « Data-scope » pour scruter les données et en extraire la Science

Avec sa nouvelle salle et le projet CAPRI : « Cloud Académique Production Recherche Innovation », le CC-IN2P3 et ses 9 partenaires ont l'ambition de créer un tel Data-scope

27 septembre 2011

11

Presentation de D. Boutigny





- > **Report on current policies and practices of the High Energy Physics program for disseminating research results**
  - June 3, 2011
- > “To date no HEP experiment has provided large-scale open access to its raw form digital data, although limited access to processed data has sometimes been granted upon request. The size and complexity of these datasets present significant technological, governance, and support challenges. “
- > “DPHEP Study Group is an international effort working to develop solutions to these challenges and to provide common guidelines for use by future collaborations. “
- > “The preservation of HEP data and its dissemination requires organized action from the experimental collaborations, the participating laboratories, and the funding agencies.”
- > **BREAKING NEWS: NSF initiates a funding line for data preservation in HEP: proposal submitted**

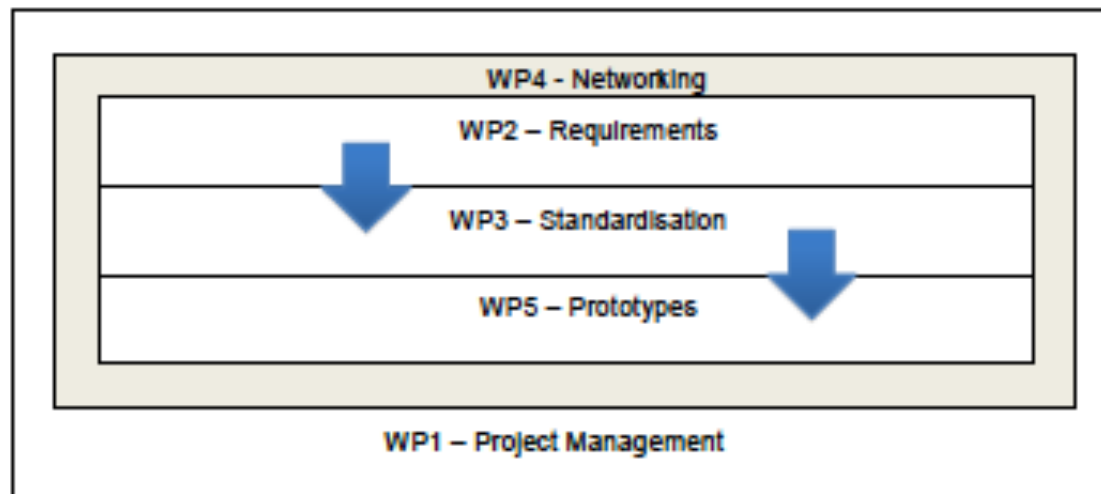
## > Multi-laboratory proposal within UE/FP7 INFRA-2012-3.2:

- The objective is to establish an EU/USA **coordination platform** aiming at full interoperability of scientific data infrastructures, and to demonstrate this coordination through several joint EU-USA prototypes that would ensure **persistent availability and effective sharing of data across scientific domains**, organisations and national boundaries. The platform should provide for: the collection of requirements and approaches for standardisation (development, promotion, adoption and maintenance); common ICT infrastructure approaches (technical, semantic, reference architecture, financing models, etc) in order to lower access barriers; harmonisation of intellectual property frameworks for scientific information; and mechanisms for international networking of experts and multidisciplinary communities. The joint prototypes should leverage and build upon **similar initiatives in Europe and USA**. The proposal should clearly describe synergies and collaboration with corresponding existing or potential NSF-funded initiatives.

## > CERN (coord.), DESY, CC-IN2P3 joint proposal

- with strong support from SLAC, FNAL, BNL (and letters of support from OSGrid, CMS)

## > A prototype for further developments, applications etc.



**More proposals  
in preparation:  
Get prepared for FP8**

# Quelque part en France...

## > Mastodons:

- La Mission Interdisciplinarité (MI) du CNRS lance **un défi sur la gestion, l'analyse et l'exploitation des très grandes masses de données scientifiques** (MASTODONS). Le but est d'identifier et de soutenir des actions de recherche dont les résultats ne pourraient être obtenus sans une fertilisation croisée des disciplines et sans une synergie effective entre chercheurs.

## > Projet: PREDON C. Diaconu (CPPM), G.Lammana (LAPP). S. Kraml (LPSC)

- le projet PREDON propose une approche nouvelle qui mélange les capacités scientifique, technique et organisationnelle des grandes collaborations en physique des particules et astrophysique pour définir et construire un system robuste de stockage et analyse des donnés à long terme.
- But pour 2012: montrer qu'il existe un intérêt a travers les disciplines et les instituts du CNRS
- Initiatives similaires MPI (Allemagne), INFN(Italie), STFC(UK)
- Workshops 2012:
  - Multi-disciplinary WS +DPHEP meeting(november 6-9 2012)
  - Finacing models for DP (Marseille?)

## Conclusion and outlook

**HEP data is potentially richer than the designed physics programs and the prolongation of its lifetime brings new and cost-effective science**

- > The DPHEP Study Group has established itself in the HEP community and has reached a milestone in the publication of the latest report, which contains a comprehensive appraisal of data preservation in HEP

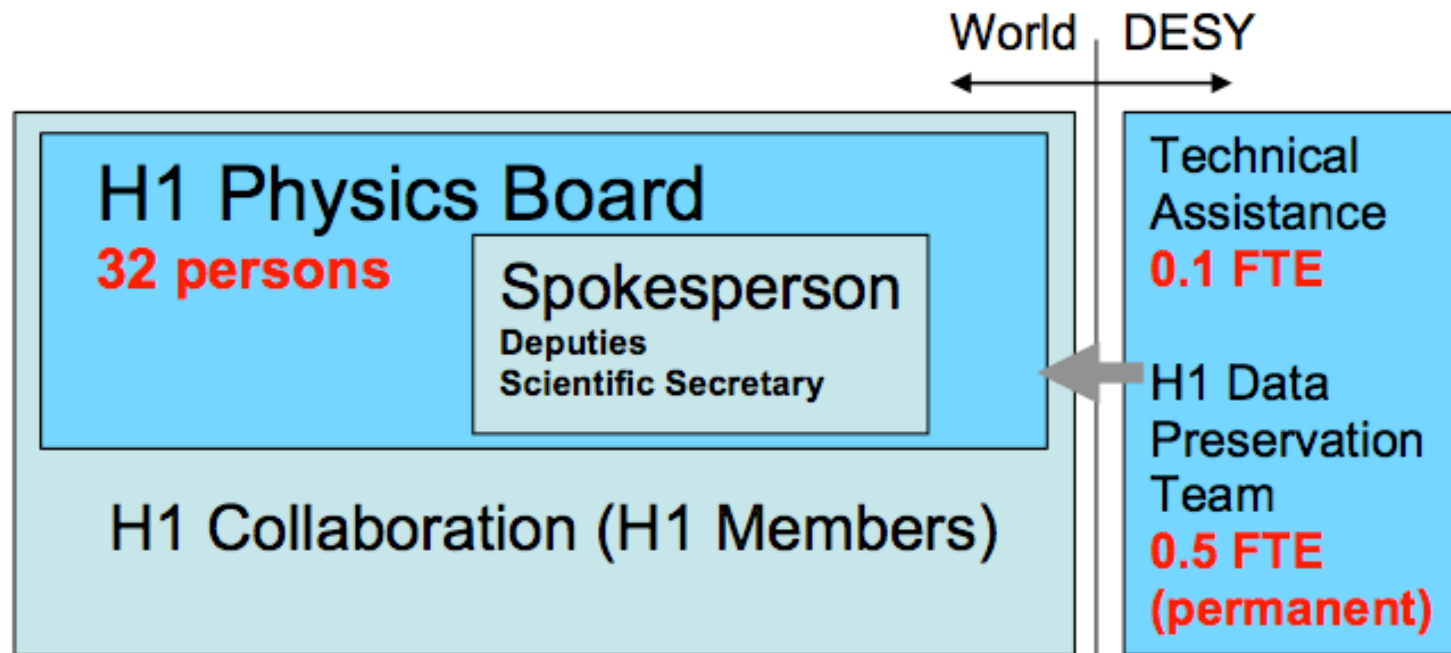
**arXiv:1205.4667**

- > The group will continue to investigate and take action in areas of coordination, preservation standards and technologies, as well as expanding the experimental reach and inter-disciplinary cooperation
- > In order to do this a transition of the Study Group to the more structured **DPHEP Organisation** should occur, with same or higher level of endorsement and a clear funding model

# EXTRAS

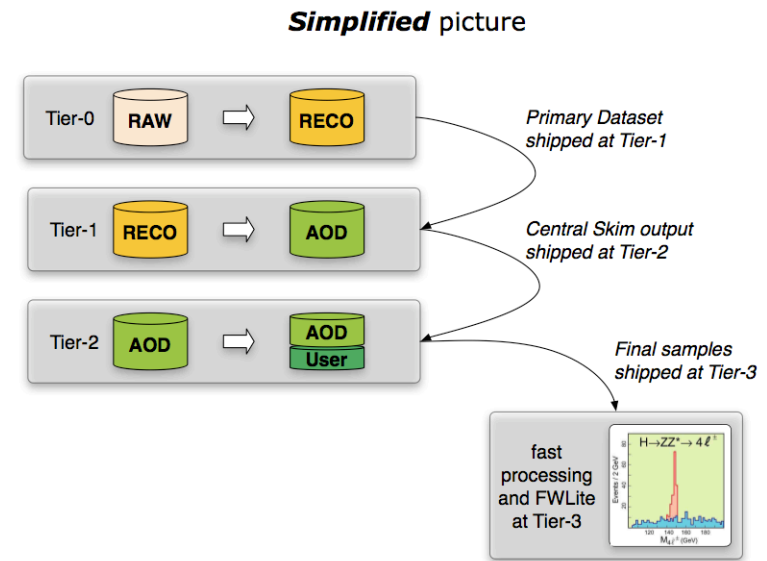
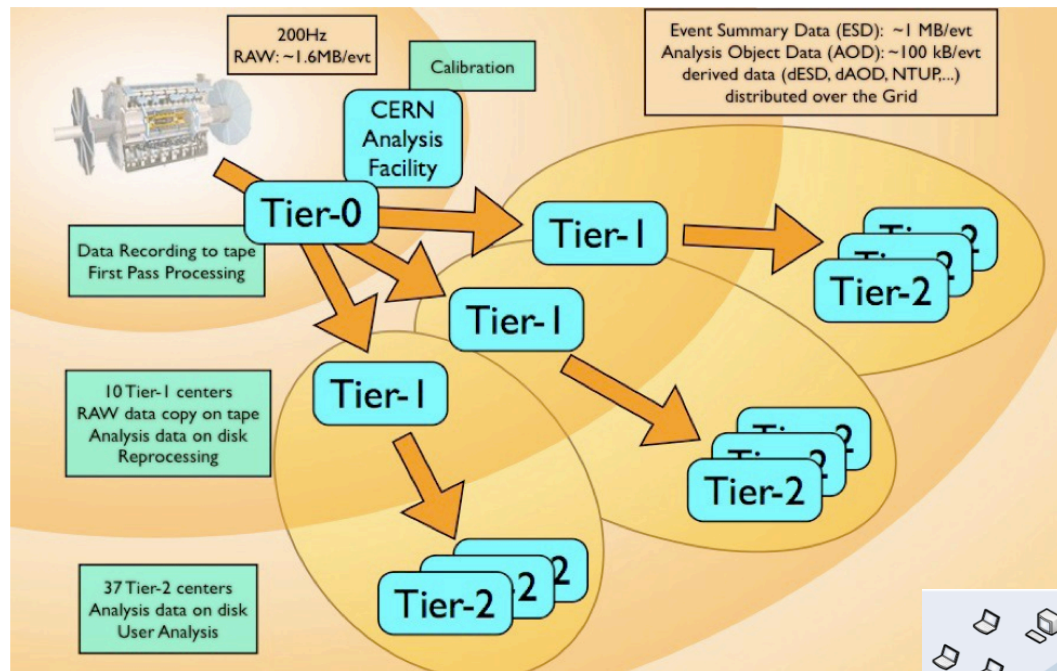


# Collaboration transitions

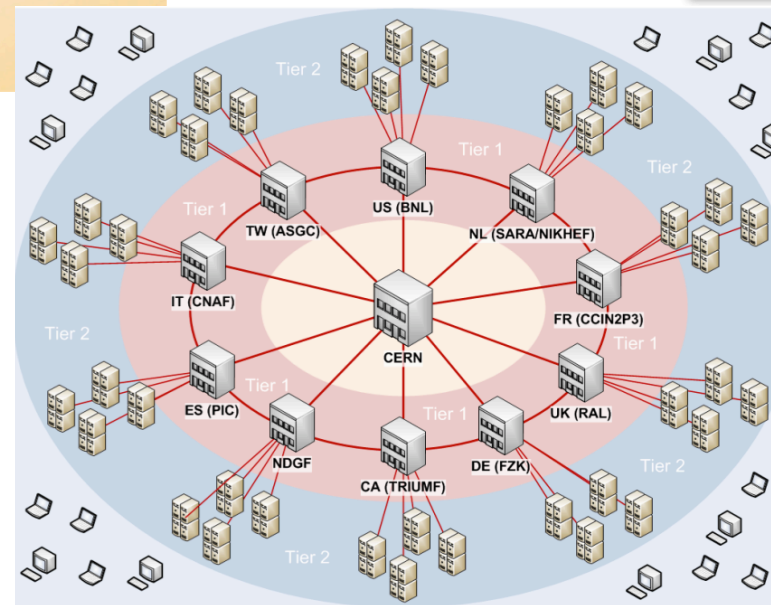


- Future structure collaborations should also be considered by experiments
  - Experimental organisation risks being left in an undefined state
  - Transition should also be planned in advance of the projected end date
  - Of particular note are authorship issues
  - Important when considering the future use of data and open access

# Data analysis models in HEP in the LHC era



- More skims - *yes*
- More distribution - *certainly*
- More complexity - *perhaps..*
- Data placement is key, but analysis-wise it's still very similar to what we had before



# INSPIRE: Paper histories



Welcome to INSPIRE  $\beta$ . Please go to SPIRES if you are here by mistake.  
Please send feedback on INSPIRE to [feedback@inspire-hep.net](mailto:feedback@inspire-hep.net)

HEP :: HELP :: ..... SPIRES HEPNAMES :: INST :: CONF :: EXP :: JOBS

[Home](#) > Events with Isolated Leptons and Missing Transverse Momentum and Measurement of W Production at HERA

[Information](#) [References \(52\)](#) [Citations \(8\)](#) **H1 internal**

## Events with Isolated Leptons and Missing Transverse Momentum and Measurement of W Production at HERA.

H1 Collaboration (F.D. Aaron (Bucharest, IFIN-HH & Bucharest U.) *et al.*) [Show all 256 authors](#).  
2009

**Eur.Phys.J. C64 (2009) 251-271**  
e-Print: [arXiv:0901.0488 \[hep-ex\]](#)

**Abstract:** Events with high energy isolated electrons, muons or tau leptons and missing transverse momentum are studied using the full  $e^+p$  data sample collected by the H1 experiment at HERA, corresponding to an integrated luminosity of  $474 \text{ pb}^{-1}$ . Within the Standard Model, events with isolated leptons and missing transverse momentum mainly originate from the production of single W bosons. The total single W boson production cross section is measured as  $1.14 \pm 0.25 \text{ (stat.)} \pm 0.14 \text{ (sys.) pb}$ , in agreement with the Standard Model expectation. The data are also used to establish limits on the  $WW\gamma$  gauge couplings and for a measurement of the W boson polarisation.

**Keyword(s):** INSPIRE: [W: production](#) | [transverse momentum: missing-energy](#) | [DESY HERA Stor](#) | [H1](#)



Record created 2009-01-05, last modified 2010-04-11 [Similar records](#)

[Abstract](#) and [Postscript](#) and [PDF](#) from arXiv.org  
[Journal Server](#)  
[Reaction Data \(Durham\)](#)

[Export](#)  
[BibTeX](#), [EndNote](#), [LaTeX\(US\)](#), [LaTeX\(EU\)](#), [NLM](#), [DC](#)

- > Envisage an additional link for H1 members only
- > Provides additional information such as preliminary results, earlier draft versions and documentation from the publication procedure

# INSPIRE: Paper histories



Welcome to [INSPIRE](#) ?. Please go to [SPIRES](#) if you are here by mistake.  
Please send feedback on INSPIRE to [feedback@inspire-hep.net](mailto:feedback@inspire-hep.net)

HEP :: [HELP](#) :: ..... [SPIRES](#) [HEPNAMES](#) :: [INST](#) :: [CONF](#) :: [EXP](#) :: [JOBS](#)

[Home](#) > [Events with Isolated Leptons and Missing Transverse Mo](#)

[Home](#) > > [Search Results](#)

Information | [References \(52\)](#) | [Citation](#)

**Events with Isolate**

[Abstract and Postscript](#)  
[Journal S](#)  
[Reaction Data](#)

Abs  
data  
with  
prod  
also

Key

Record created 2009-01-05, last mod

## Events with Isolated Leptons and Missing Transverse Momentum and Measurement of W Production at HERA

### PUBLICATION HISTORY

**Preliminary Results**  
[HEP-EPS 2007 conference paper](#) | July 2007  
[Prepared for Deep Inelastic Scattering 2007](#) | April 2007  
[Prepared for 42nd Rencontres de Moriond \(Electroweak\)](#) | January 2007  
[Prepared for the 62nd DESY PRC](#) | October 2006  
[ICHEP 2006 conference paper](#) | July 2006  
[Prepared for the 60th DESY PRC](#) | November 2005  
[HEP-EPS 2005 conference paper](#) | July 2005  
[Lepton Photon 2005 conference paper](#) | June 2005  
[Prepared for Deep Inelastic Scattering 2005](#) | April 2005  
[Prepared for the 58th DESY PRC](#) | October 2004  
[Analysis of High Pt HERA II Data](#) | [ICHEP 2004 conference paper](#) | August 2004  
[High Pt Analysis of the HERA II Data](#) | [Prepared for Deep Inelastic Scattering 2004](#) | April 2004

**T0 talks**  
[Pre-T0 Talk](#) | 08.02.2008  
[T0 Talk](#) | 24.07.2008  
[T0 Addendum](#) | 14.08.2008

**Paper Drafts**  
[First Draft](#) | [Answers to Draft](#) | 15.08.2008  
[Second Draft](#) | [Answers to Draft](#) | 19.11.2008  
[Referee Report](#) | 20.11.2008  
[Final Version](#) | 06.01.2009

C. Diaconu | Data Preservation

# For completeness, the HERA data summary



- > Final ZEUS data reprocessing to mDST completed in 2009
  - Basic preserved data format: ROOT based “Common Ntuples” (CN)
  - Ultimately RAW, MDST data and MC removed from robots, keep only CN
  - Reduces total amount to be preserved for ZEUS from the current 1 PB to ~ **200 TB**



- > Final H1 reprocessing of HERA II data 2009, HERA I repro almost there
  - Common analysis software H1OO started in 2000, uses ROOT based data format, used by all H1
  - In addition, a monthly MC production of up to 1/4 billion events
  - H1 to preserve RAW data, as well as one DST version and one analysis level version
  - Estimate total amount to be preserved for H1 to be ~ **200-500 TB**



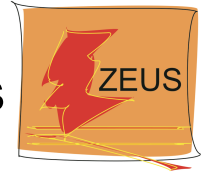
- > Main format for HERMES analyses is the mDST
  - New production planned before final freeze
  - Last years of data taking with recoil detector, still need improved calibrations
  - MC productions on Grid for on-going analyses
  - Total amount to preserve on tapes ~ **20-500 TB**



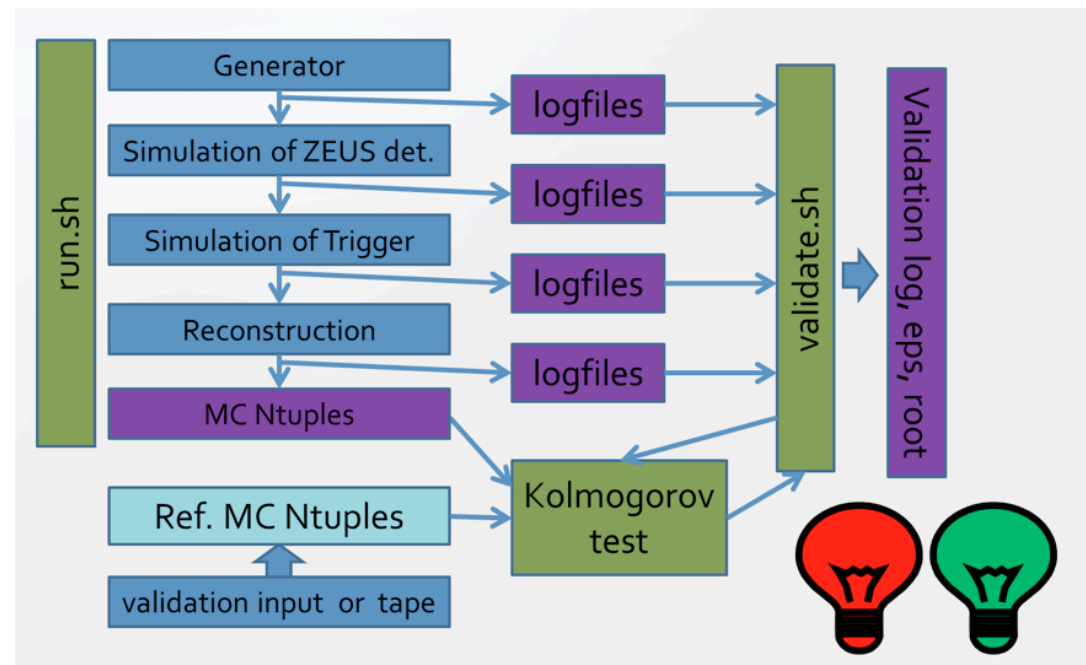
- > Preservation of HERA-B data under investigation within DESY-IT
  - Total amount of data currently ~ **250 TB**, decreases once preservation model established



## Example structure of experiment tests: ZEUS (Level 3 + MC chain)



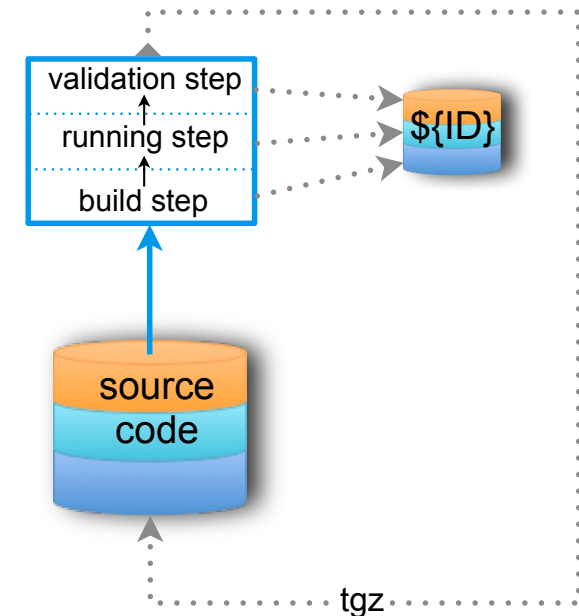
- > ZEUS strategy: use ROOT based analysis level Common Ntuples as data format for preservation – DPHEP level 3
- > Only external dependence is ROOT
  - Validation of new ROOT versions included as analysis level tests in the **sp-system**
- > However, the MC production chain executables will also be preserved as a standalone package
- > In addition, an interface for new generators is developed, which is also included in the validation system



# Running jobs in the sp-system

## > Initial step

- Compilation of analysis (level 3) and sim/rec (level 4) software
- **Or:** use tar-balls with pre-compiled software
- Provide access to software
  - Copy tar-balls to persistent storage
- All output kept in directory with unique name



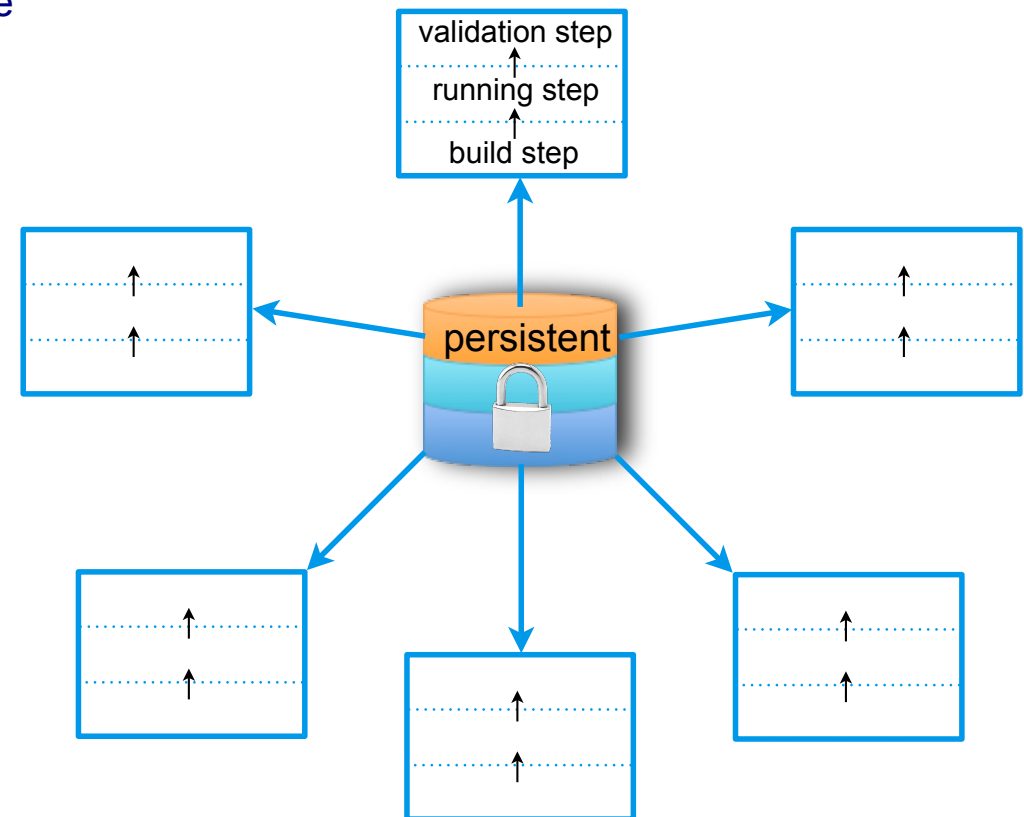
# Running jobs in the sp-system

## > Initial step

- Compilation of analysis (level 3) and sim/rec (level 4) software
- **Or:** use tar-balls with pre-compiled software
- Provide access to software
  - Copy tar-balls to persistent storage
- All output kept in directory with unique name

## > Run parallel tests

- Set up software environment
- Validate binaries with persistent input
  - e.g. event display, database access, ...



# Running jobs in the sp-system

## > Initial step

- Compilation of analysis (level 3) and sim/rec (level 4) software
- **Or:** use tar-balls with pre-compiled software
- Provide access to software
  - Copy tar-balls to persistent storage
- All output kept in directory with unique name

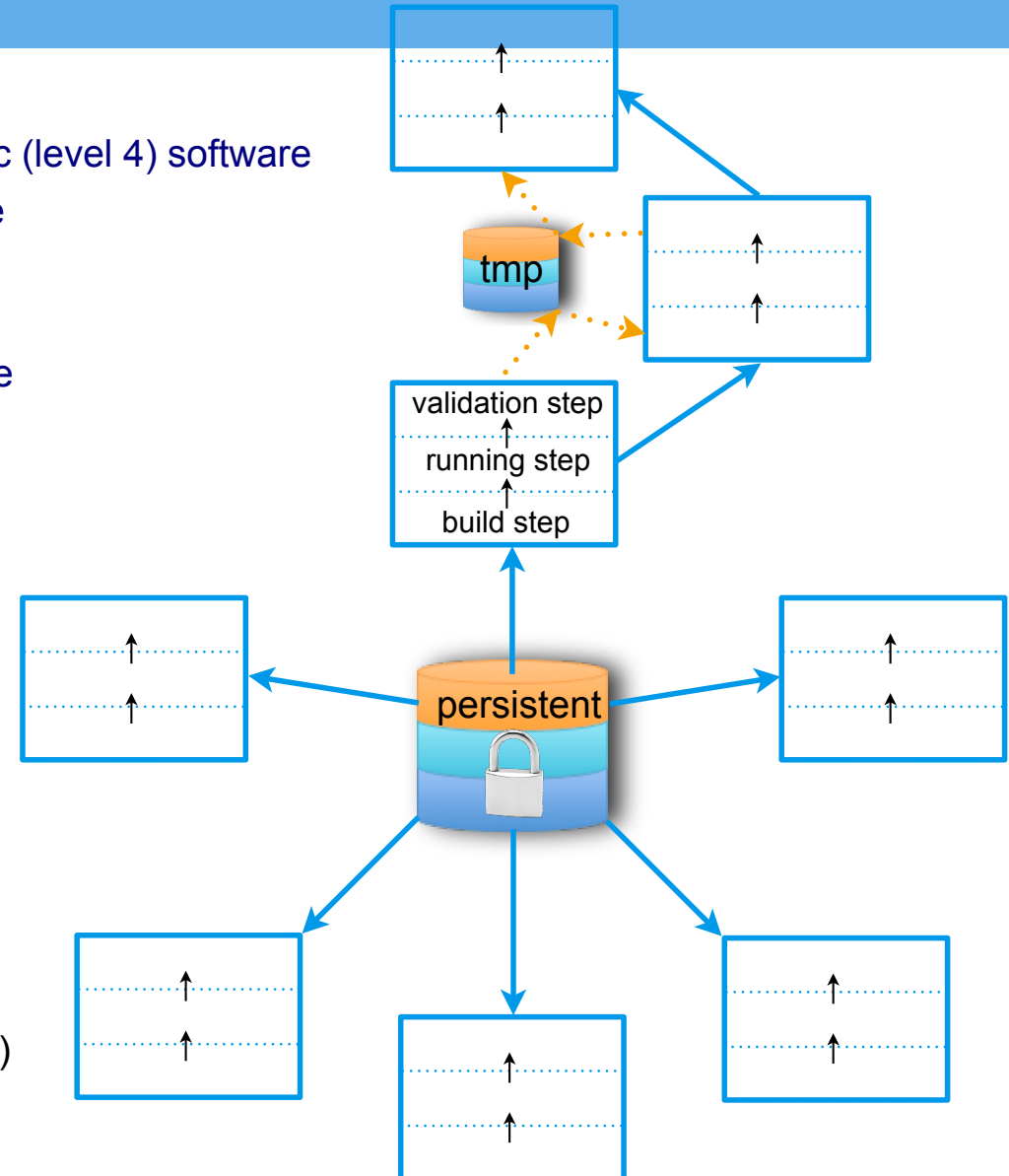
## > Run parallel tests

- Set up software environment
- Validate binaries with persistent input
  - e.g. event display, database access, ...

## > Run sequential tests

- Set up software environment
- Validate file production
  1. MC generation (produce gen files)
  2. Reconstruction (gen. files → DSTs)
  3. Analysis level (DSTs → ROOT files)
- Tests use output of previous test as input

## > Results remain accessible or can be reproduced with identical results



# Securing the resources

- The new DPHEP organisation will develop at least three levels:
  - Experiment / collaboration level projects
  - Multi-experiment level initiatives
  - Global DPHEP level projects or positions
- It is foreseen that funding must come from different sources, in particular for common DPHEP enterprises or positions
- The experiment and laboratory level projects are highest priority (1-2 FTE per site), followed by the appointment of the DPHEP Project Manager, which is a full time position
- Many potential multi-experiment projects also exist, including those shown today, which depend on additional funding, typically 0.5-1 FTE



# LEP Paper Tables

	2001	2002	2003	2004	2005	2006	2007	2008	2009	Total	2004-2009
<b>ALEPH</b>	46	42	24	34	12	9	4	4	2	<b>607</b>	65
<b>DELPHI</b>	64	30	31	58	21	19	7	7	2	<b>678</b>	114
<b>L3</b>	51	40	23	52	16	11	5	2	0	<b>578</b>	86
<b>OPAL</b>	61	38	32	55	9	11	4	3	2	<b>675</b>	84
<b>All</b>	<b>222</b>	<b>150</b>	<b>110</b>	<b>199</b>	<b>58</b>	<b>50</b>	<b>20</b>	<b>16</b>	<b>6</b>	<b>2538</b>	<b>349</b>

*Table 1: Statistics of peer-reviewed publications of the LEP collaborations.*

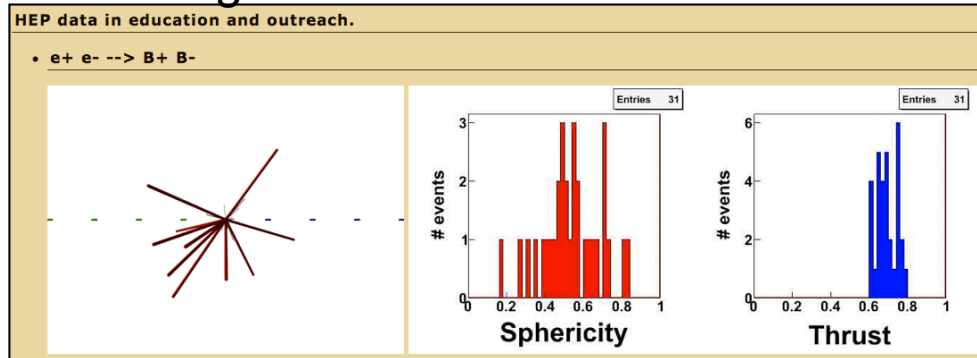
Papers 2004-2009	ALEPH	DELPHI	L3	OPAL	All
<b>Electroweak</b>	17	26	22	24	89
<b>QCD</b>	19	25	19	22	85
<b>Higgs Searches</b>	6	14	8	9	37
<b>SUSY Searches</b>	4	7	5	9	25
<b>Exotica Searches</b>	5	12	10	7	34
<b>Flavour Physics</b>	6	15	4	5	30
<b>Exclusive Channels</b>	3	8	8	2	21
<b>Cosmo-LEP</b>	3	3	6	0	12
<b>Other</b>	2	4	4	6	16
<b>Total</b>	<b>65</b>	<b>114</b>	<b>86</b>	<b>84</b>	<b>349</b>

*Table 2: Distribution of physics topics in LEP publications in the years 2004-2009.*

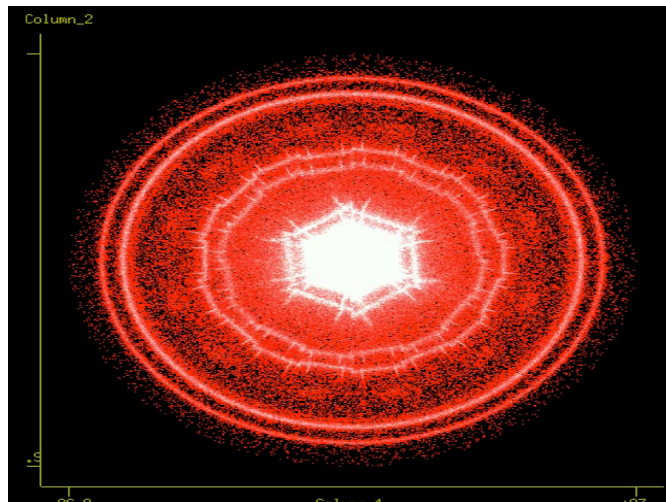
# Outreach Data and Tools

[http://www.slac.stanford.edu/~bellis/HEP\\_data.html](http://www.slac.stanford.edu/~bellis/HEP_data.html)

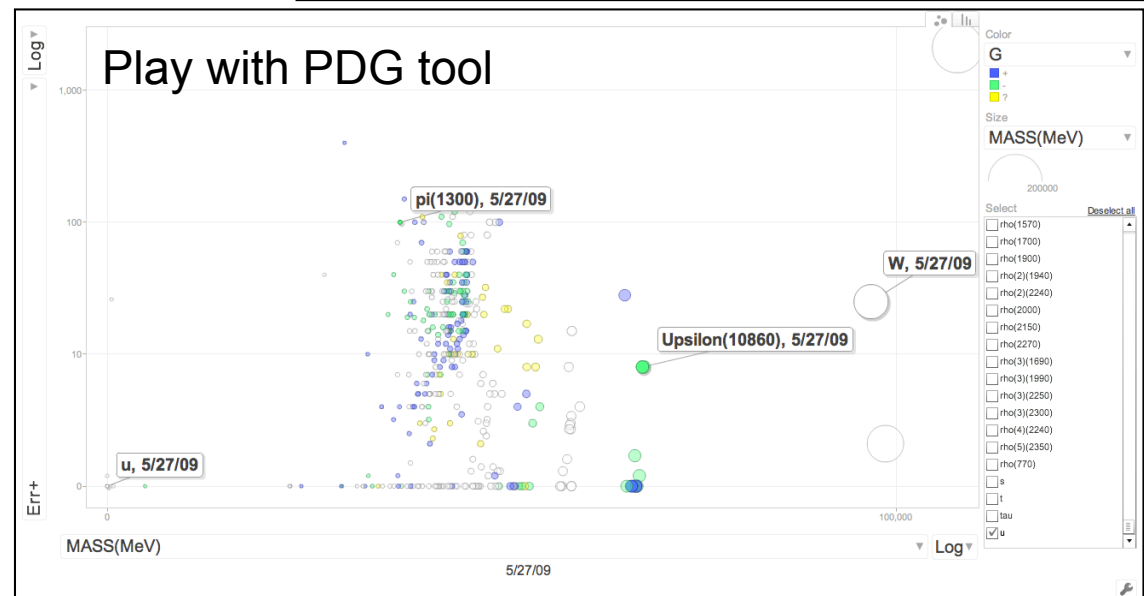
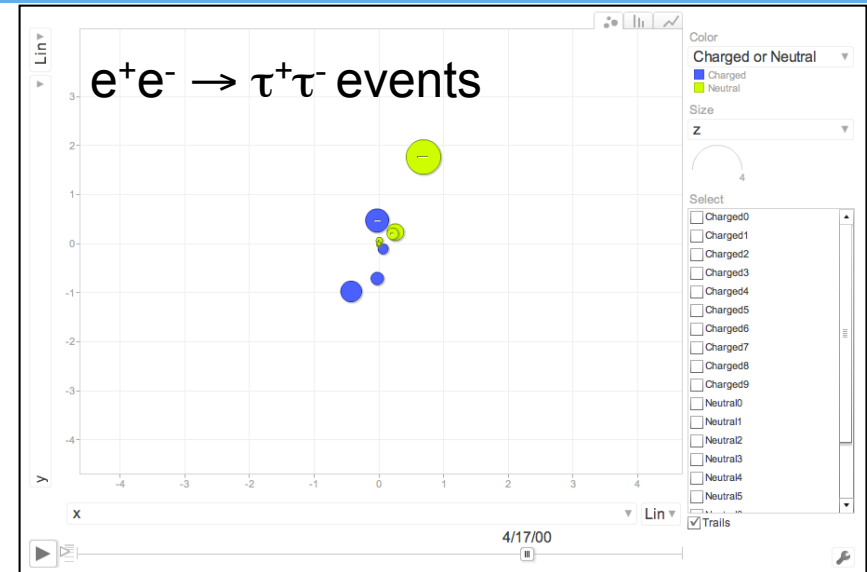
## Movie of generic $e^+e^- \rightarrow B^+B^-$ events



Several outreach tools already  
being used in classrooms



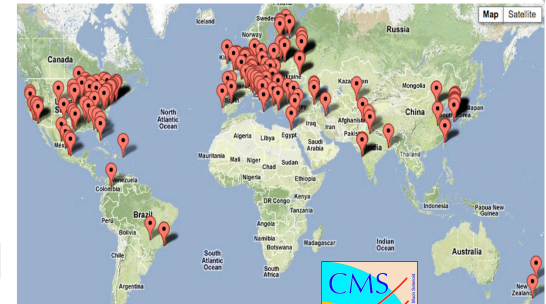
Tomography studies using converted  
photons in BaBar Silicon Vertex-Tracker



# Science Hack Day: Increasing the access to LHC data

<http://cms.web.cern.ch/news/cms-public-data-activity-scoops-prize-nairobi>

## CMS public data activity scoops prize in Nairobi



An application using real event data from CMS has won "Best Science" prize in a public "Science Hack Day" held in Nairobi between 13<sup>th</sup> and 15<sup>th</sup> April 2012. Science Hack days bring together a wide range of enthusiastic members of the public to create something completely new using existing scientific systems or data.

The winning application visualized real CMS di-muon events from the 2011 LHC run, which are made public for use in various educational programmes, such as the [IPPOG Masterclasses](#), [Quarknet](#) and [I2U2](#). The application showed an animation of muons produced in CMS superimposed on a map of the world, showing where they would go if they were to continue without stopping (which they don't in reality).

Other prizes were awarded to Leah Atieno, a 15-year-old high-school student, for a voice-controlled walking robot and Denis Munene for a crowd-mapping platform to help promote the fight against malaria.

The Nairobi event, involving 240 developers, is part of broader series of Science Hack Day events. CMS data previously featured in another very successful event in [San Francisco](#).

[News article by Gythan Munga, HumanIpo](#)

[See photos of the event](#)

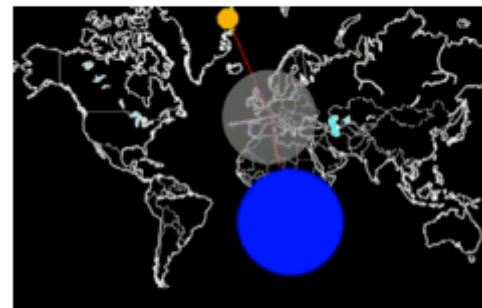
[Youtube film](#)

[Link to more Science hack events](#)

2012-04-20, by [Lucas Taylor](#)



CMS use of public data in a "Science Hack" event in Nairobi. Photo credit: [Matt Biddulph, via Flickr](#)



Application developed to visualise where muons from CMS would go if they continued forever

## Level 2: Simplified formats for outreach

- Within DPHEP and the member collaborations there are generic ideas, such as common formats and user interfaces
  - In terms formats, much can be learned from other fields such as astrophysics or life sciences
- Such outreach formats in HEP are typically based on ROOT, containing particle 4-vectors and simple event information
  - Composite-particle reconstruction, finding signals
  - Initiatives in place at BaBar, Belle and LHC experiments
- A multi-experimental project is desirable, coordinated via DPHEP, and based in several locations (CERN, FNAL, DESY..)
  - To include associated tutorials linked to preserved HEP data from several sources

