

La préservation des données scientifiques dans la physique des hautes énergies

C. Diaconu
CPPM



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

PREDON

A project for scientific data preservation in France

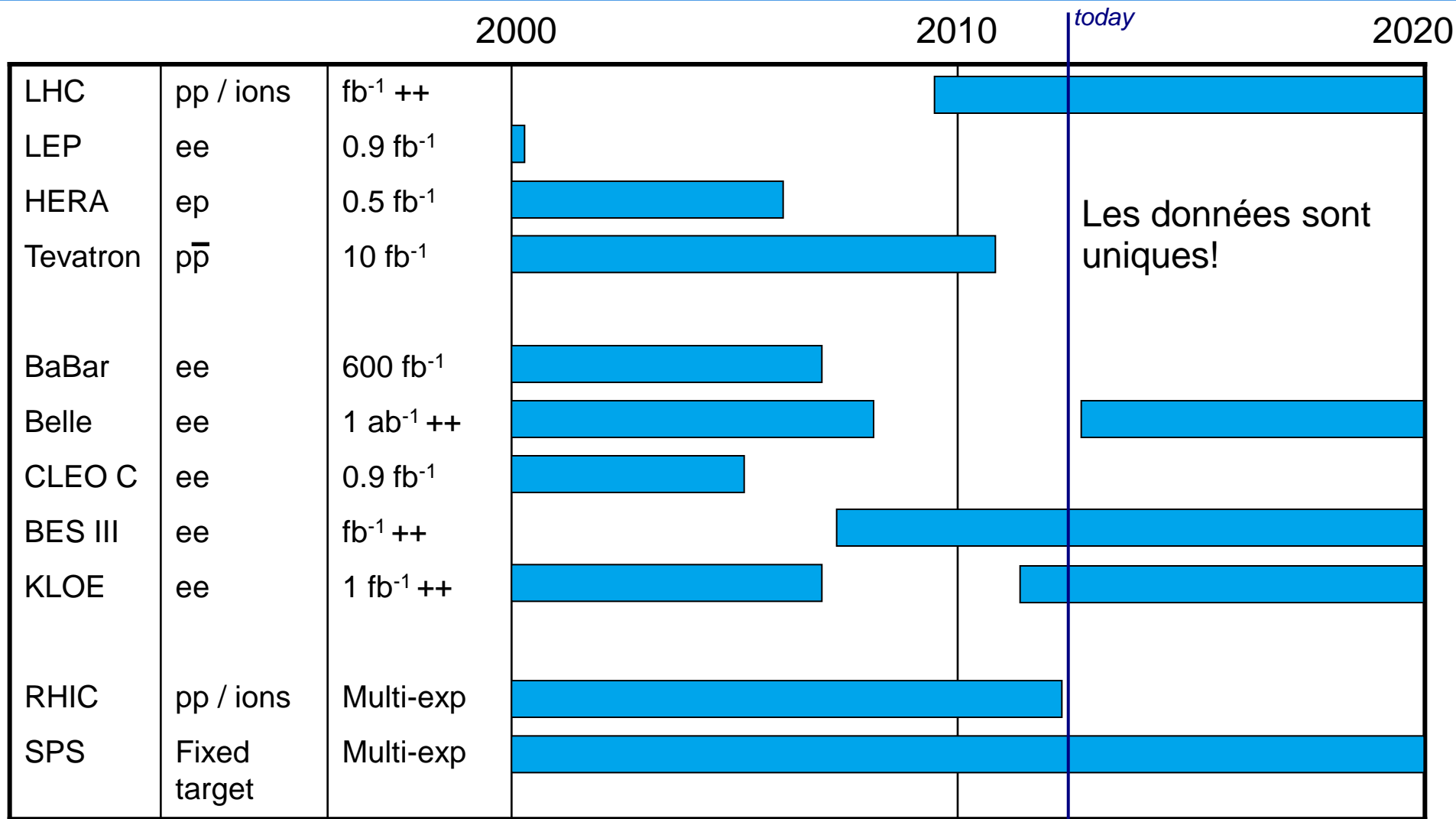


High Energy Physics

- Exemple: Physique des Hautes énergies
- Large Hadron Collider (27 Km, 13 TeV, 40MHz)



Exemple: Programme experimental de la physique des hautes énergies HEP \pm 10 ans



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

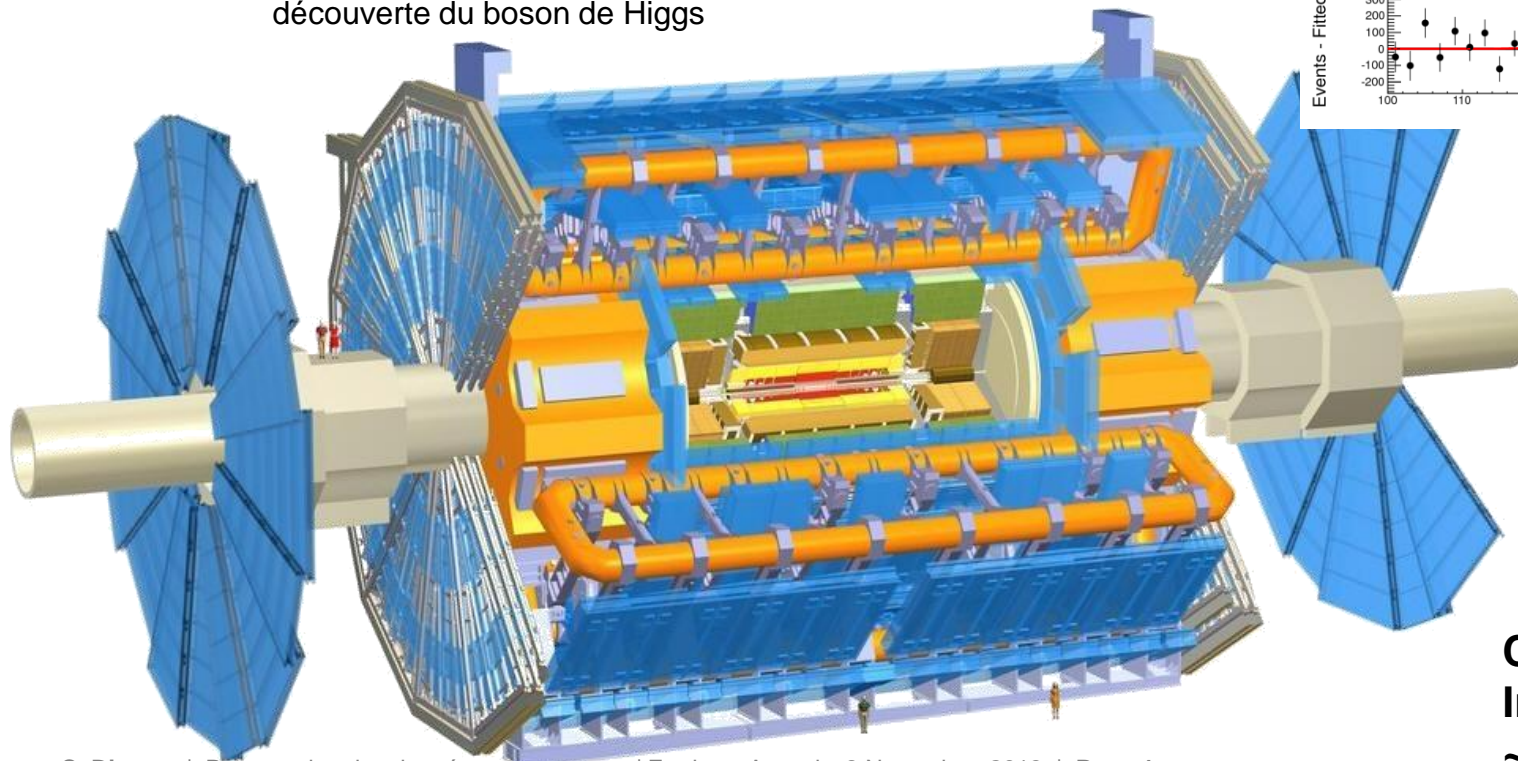
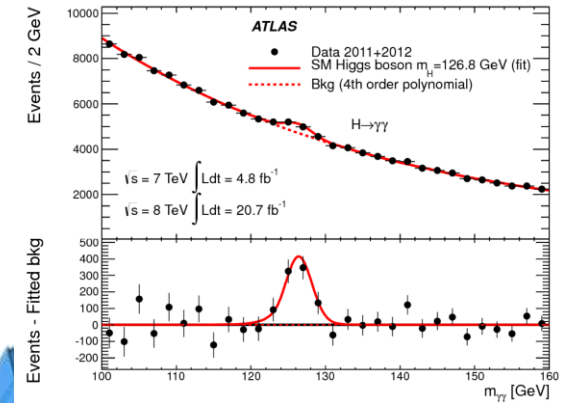
Des instruments gigantesques

➤ ATLAS: L'équivalent d'une camera avec 25Gpixels (avec une cinquantaine de technologies différentes) et 40 000 000 000 « photos » par seconde (100Pb)

➤ manips LHC:

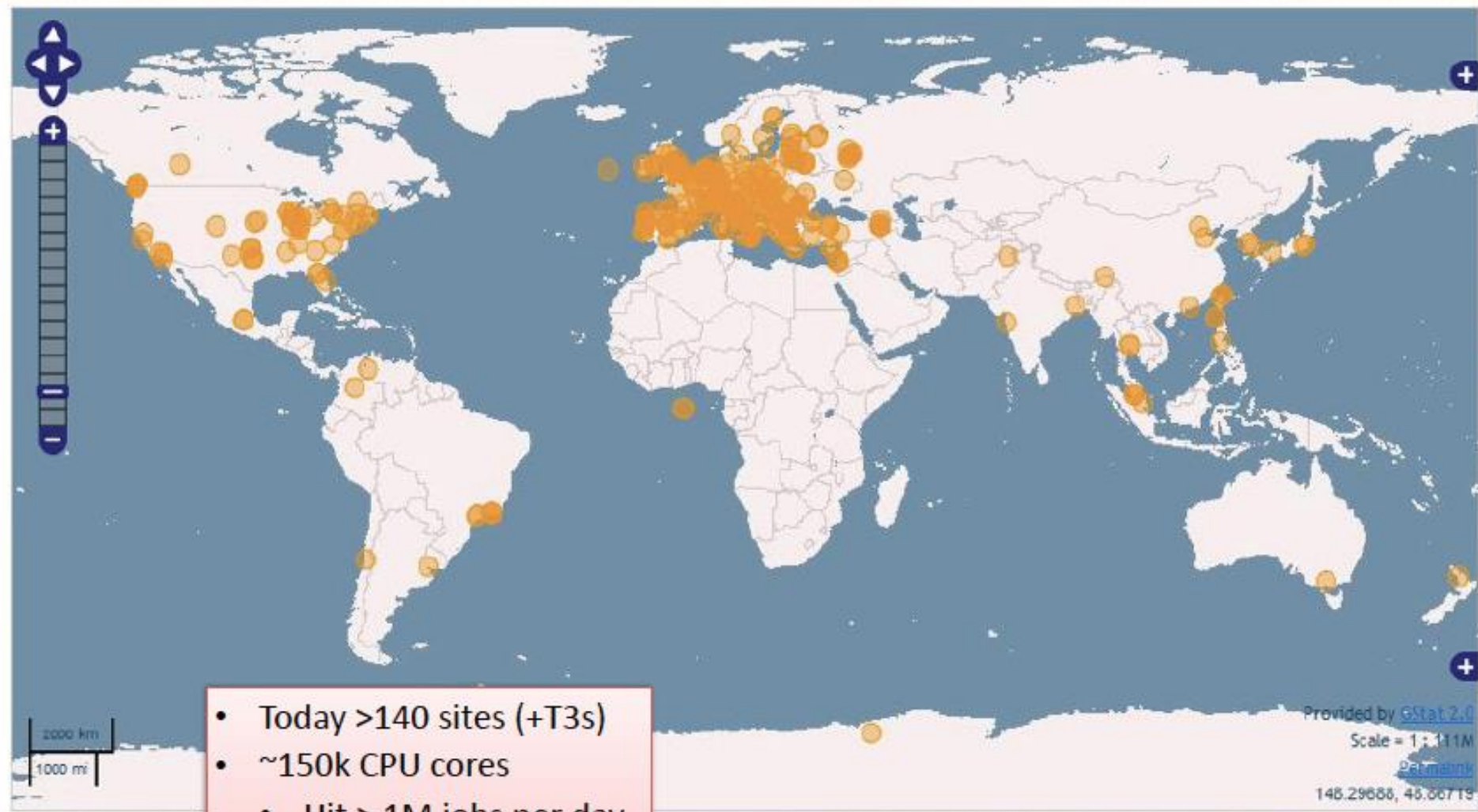
- 1000 articles scientifiques en 2 ans

découverte du boson de Higgs

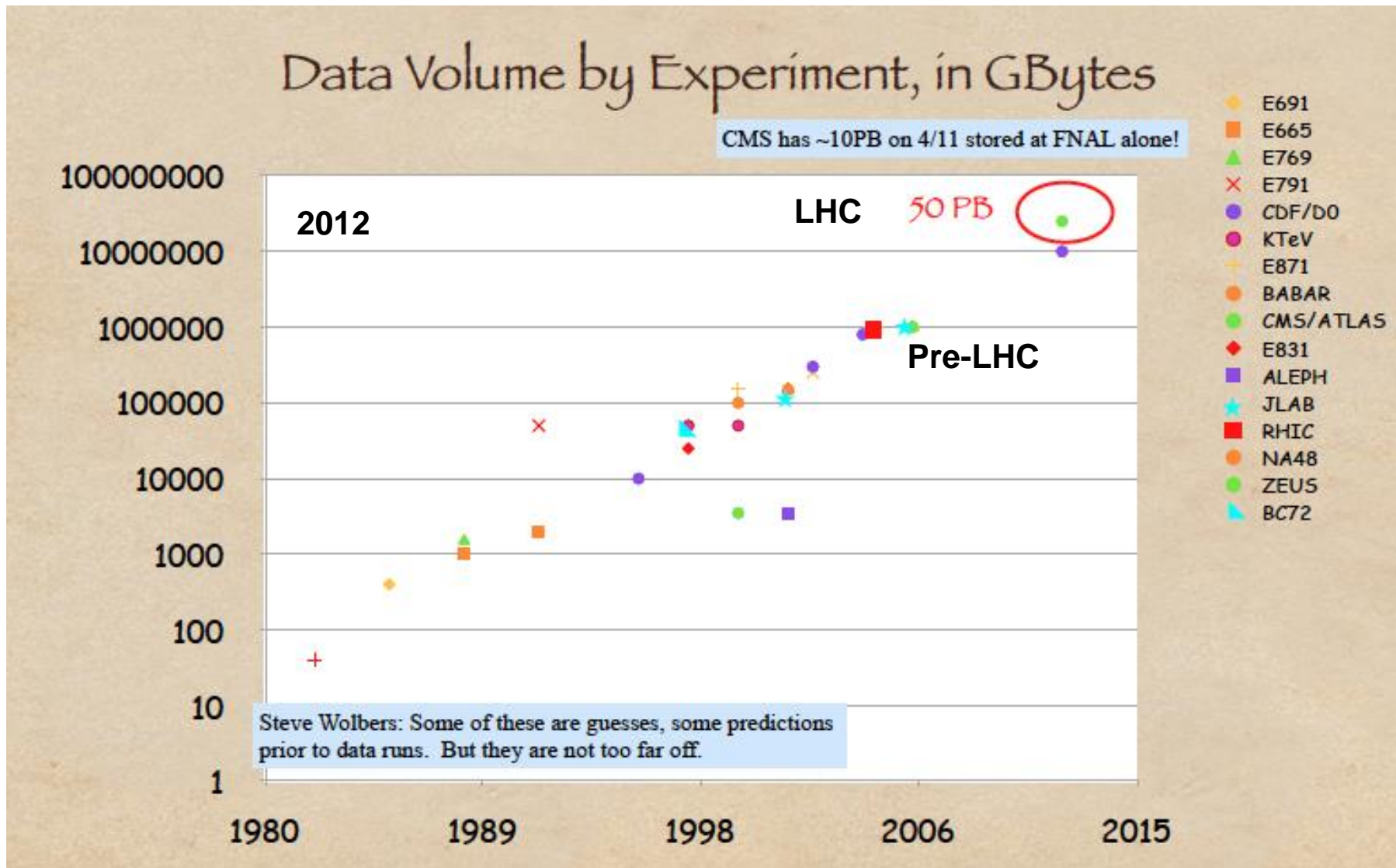


**Collaboration
Internationale
~3000 chercheurs**

LHC computing

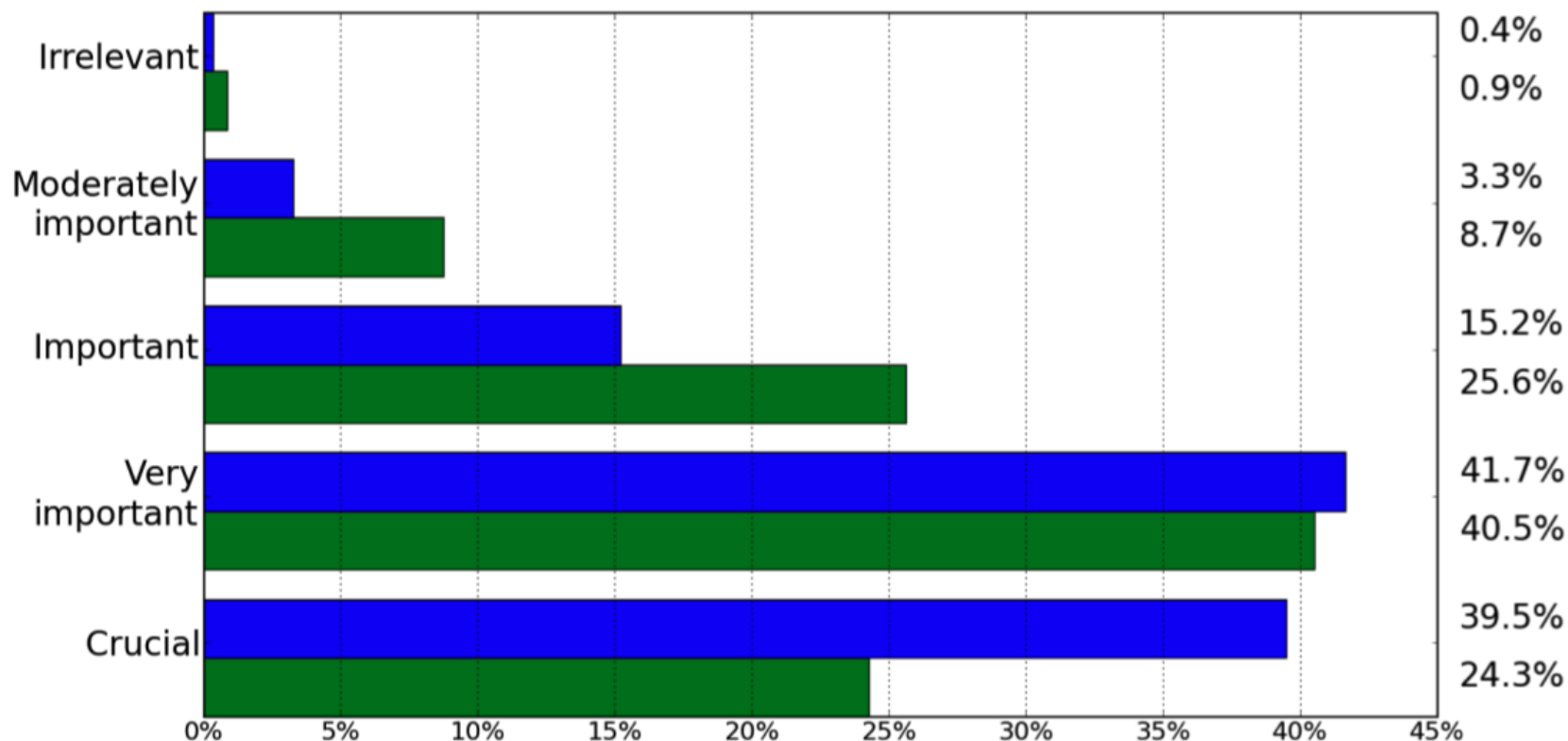


Quantité de données: 1PB -> 100PB->1EB



L'opinion de la communauté scientifique

In your opinion, how important is the issue of data preservation ?
(top/blue: theorists, bottom/green: experimentalists)

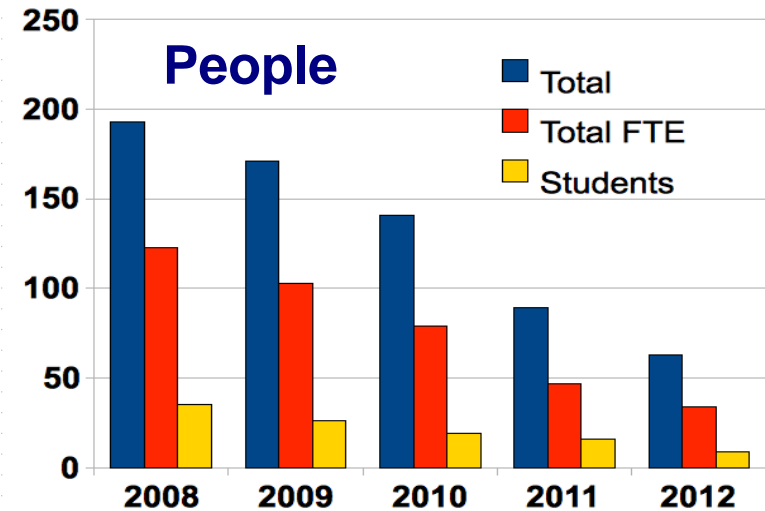
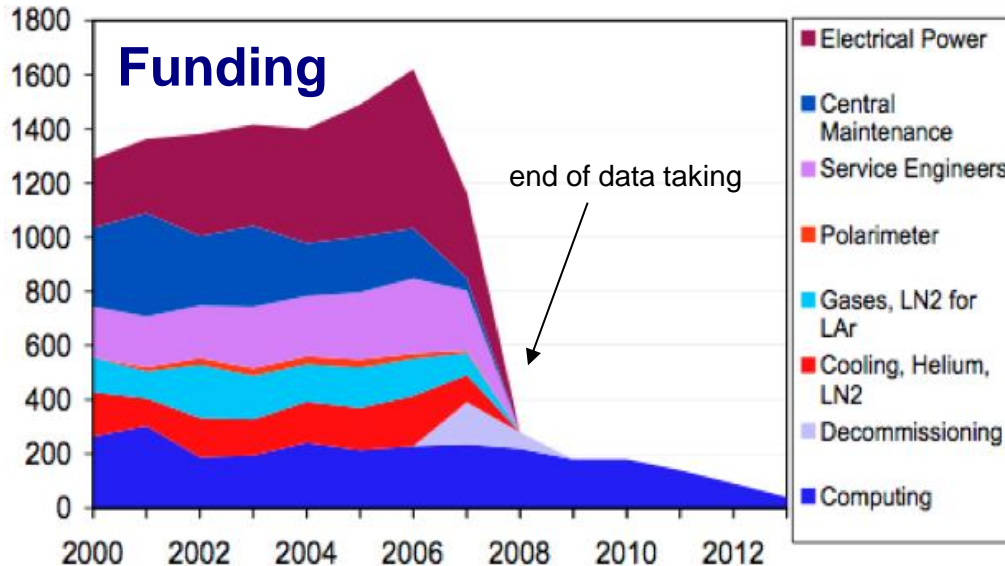
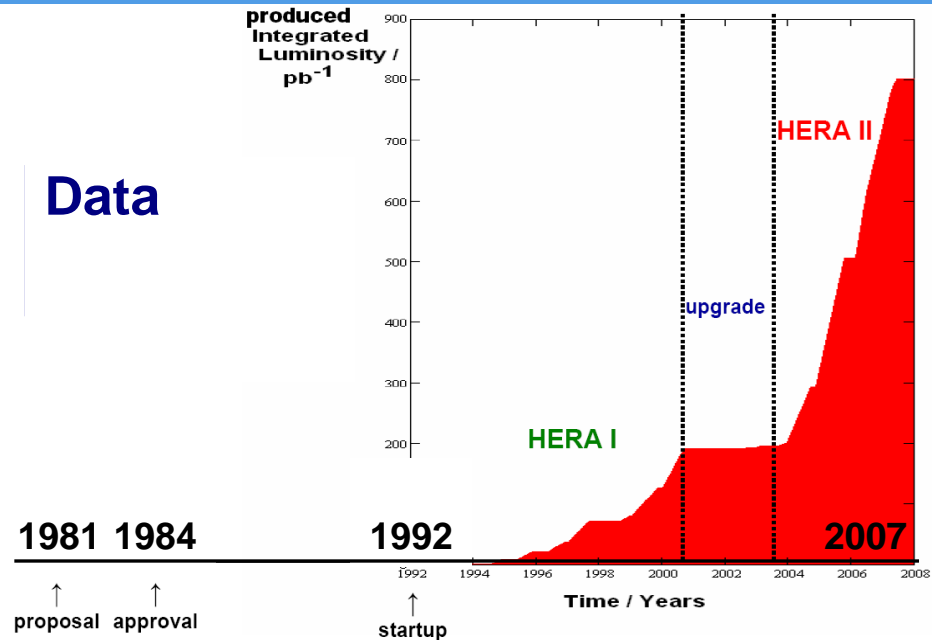


arXiv:0906.0485

Est-il difficile de préserver les données?

- > Les programmes accumulent la plupart des données vers la fin du programme
- > Les ressources (financières et humaines) décroissent rapidement après la fin des manip
- > En absence d'une planification rigoureuse les conditions pour des nouveaux projets ne sont pas idéales

Data



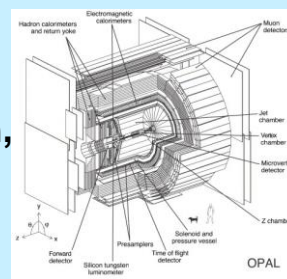
What is HEP "data"?



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

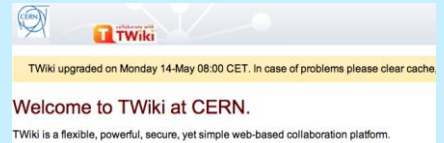
Software Simulation, reconstruction, analysis, user, in addition to any external dependencies



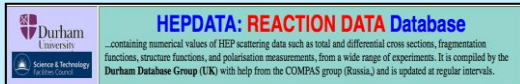
CERNLIB Access

- 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

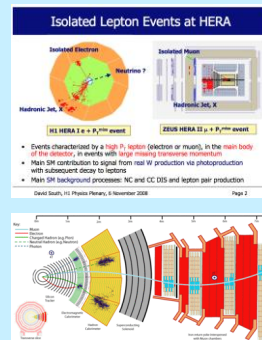
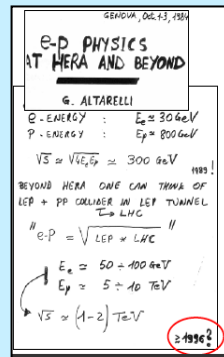
Meta information
Hyper-news, messages, wikis, user forums..



Publications arXiv.org



Documentation
Internal publications, notes, manuals, slides



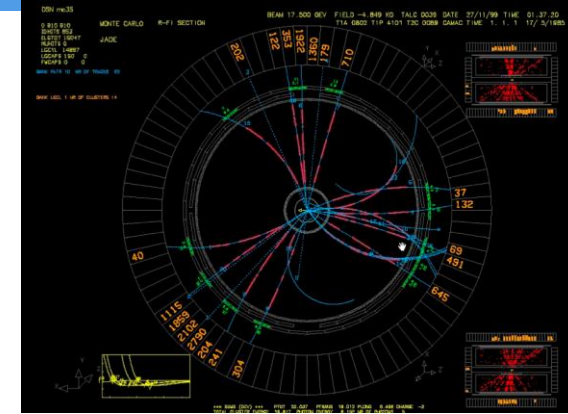
Expertise and people



Un exemple typique

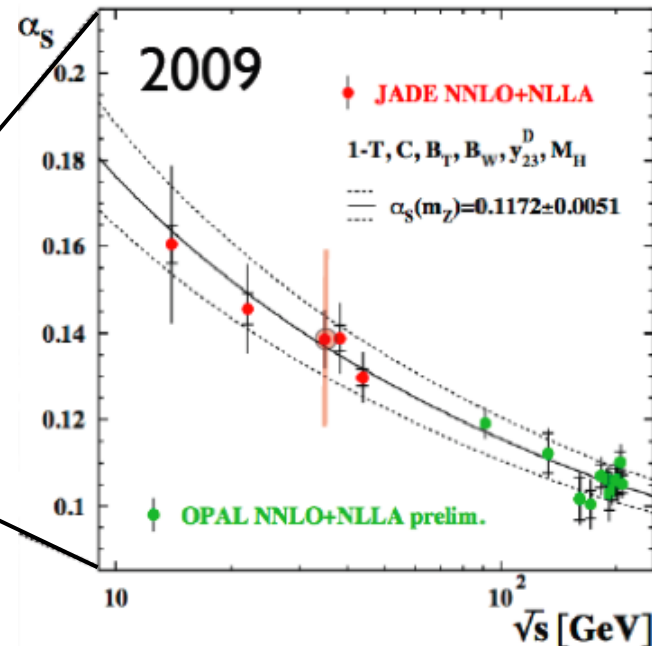
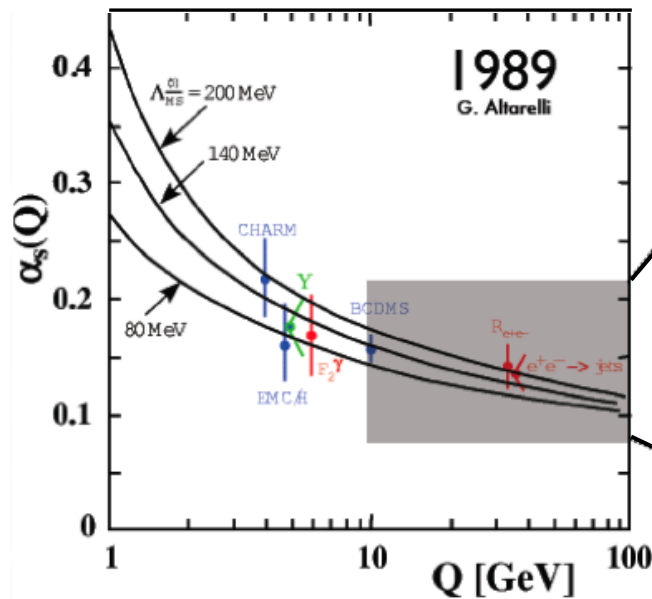
➤ Experience JADE

- Données sauvées par hasard
Nom de code: "la valise"
- Ré-analyse après 20 ans

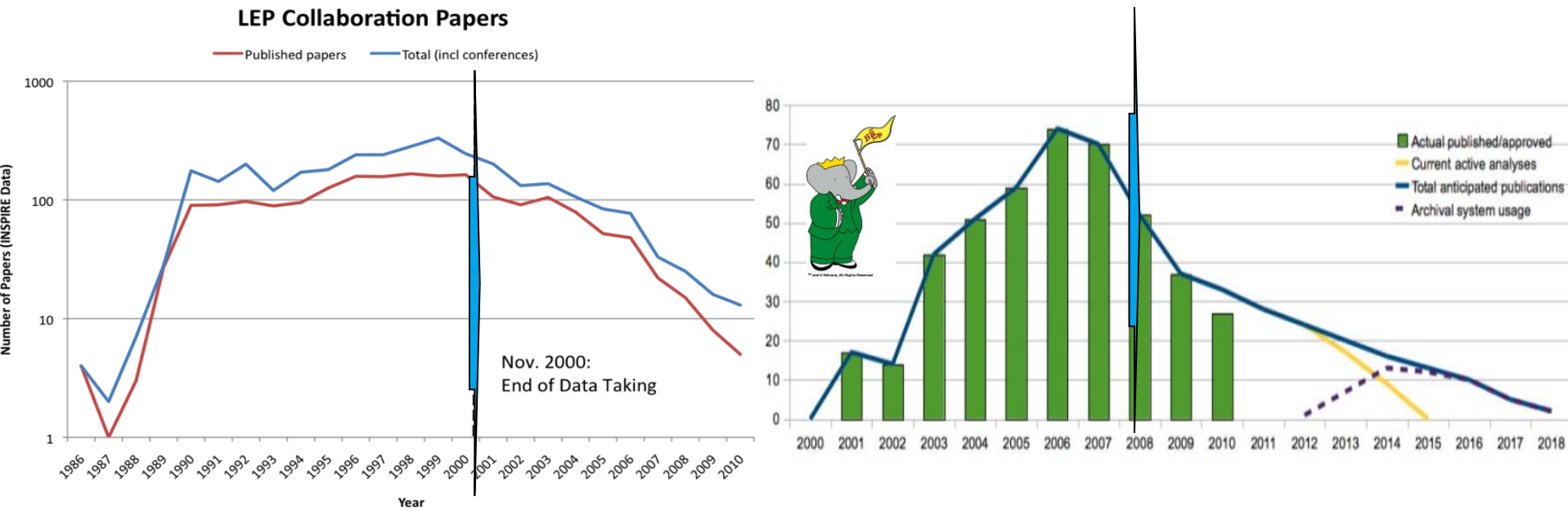


10 publications

2011



Publications à long terme



> LEP: 1989-2000

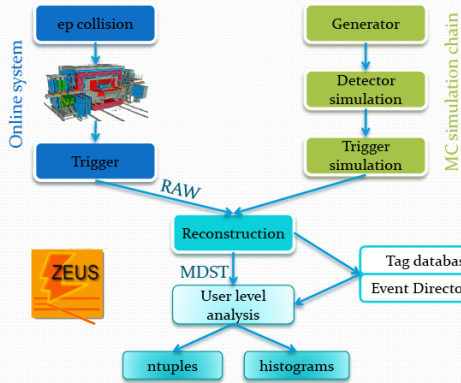
- mais 300 publications produites après 2000 et une centaine après 2005

> C'est systématique: les publications continuent long temps après la fin de l'expérience

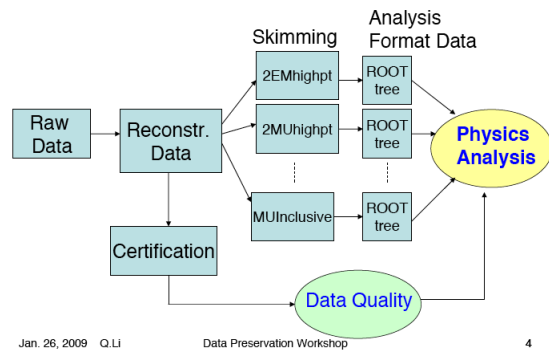
- Nouvelles idées, théories etc.

Modèles de données

Data Processing Model



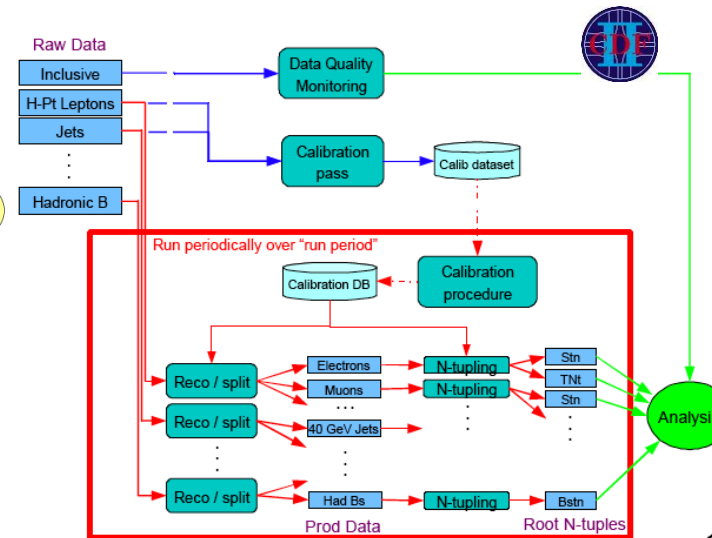
DØ Analysis Model



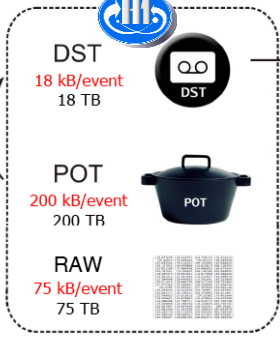
Jan. 26, 2009 Q.LI

Data Preservation Workshop

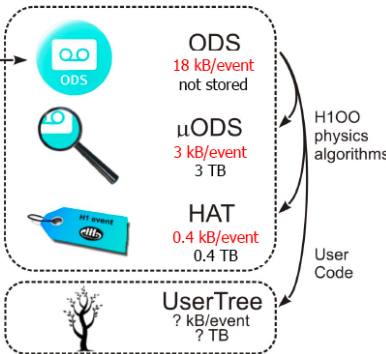
4



BOS / FPACK / Fortran

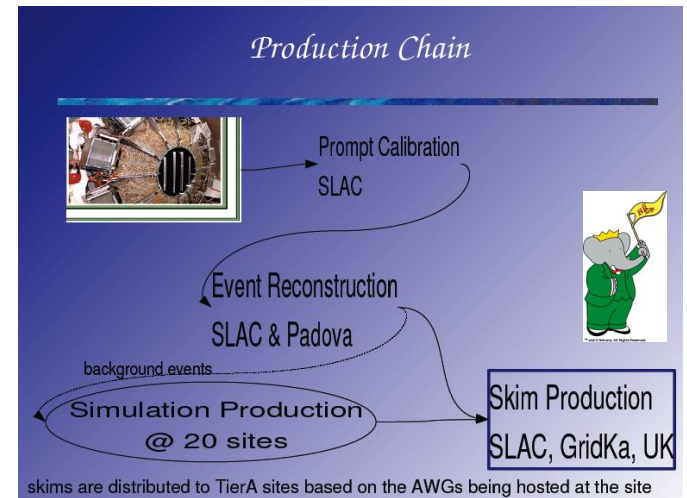


C++ / ROOT



Réduction et abstraction

- RAW (→ POT) → DST → ntuple → article
- 1.5Mb → 15Mb → 30 Kb (par "événement")
- Pas de format standardisé



DPHEP : définition des niveaux de préservation

> En progression de la complexité et les couts

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	

Level 1: Documentation

> Une tache considérable: des groupes de travail dédiées

> **Non-digital:** Cataloguing, organisation, scanning or photographing of appropriate of papers, notes, drawings, talks from pre-web days, detector schematics, blueprints, logbooks, ...

- *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 now available on INSPIRE

- Password protected now, simple to make publicly available in the future

The image shows a screenshot of the INSPIRE website. On the left, there is a login form for the 'ZEUS Intern' collection. The form includes fields for 'Username' (with 'zeus' entered) and 'Password'. Below the password field is a checkbox for 'Remember login on this computer.' and a 'login' button. A note below the form states: 'Note: You can use your nickname or your email address to login.' The top of the page features the INSPIRE logo and a navigation bar with links for 'HEP', 'INST', 'HELP', 'SPIRES', and 'HEPNAMES'. A welcome message reads: 'Welcome to INSPIRE! INSPIRE is out of beta and ready to replace SP. please email us at feedback@inspirehep.net'.

On the right, a yellow header indicates 'ZEUS Internal Notes' with '10 records found'. Below this, four records are listed:

- 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](#)

- > The ingestion of other documents is under discussion, including theses, preliminary results, conference talks and proceedings, paper drafts, ...
- > More on InSpire: reduced data?

HEP outreach initiatives

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



NETZWERK
TEILCHENWELT QUARKS, ELEKTRONEN & CO.



B - L a b

ビー・ラボ: 新しい素粒子発見のための公開データ解析プログラム
Open data analysis program to search for new particles

since 2004 copyright @ Belle collaboration



QuarkNet: The science connection you've been waiting for!

THE PARTICLE ADVENTURE
THE FUNDAMENTALS OF MATTER AND FORCE

LANGUAGES MIRROR SITES

Supported by the DOE and NSF

An award-winning interactive tour of quarks, neutrinos, antimatter, extra dimensions, dark matter, accelerators and particle detectors from the Particle Data Group of Lawrence Berkeley National Laboratory.

THE STANDARD MODEL
The theory of fundamental particles and forces
GO!

THE LARGE HADRON COLLIDER
EXPLORES UNSOLVED MYSTERIES
GO!

ACCELERATORS AND PARTICLE DETECTORS
GO!

International Particle Physics Outreach Group

INTERNATIONAL MASTERCLASSES
hands on particle physics

CPEP
Contemporary Physics Education Project

Home Fundamental Particles Plasma Physics and Fusion History and Fate of the Universe Nuclear Science | FUNDING CREDITS

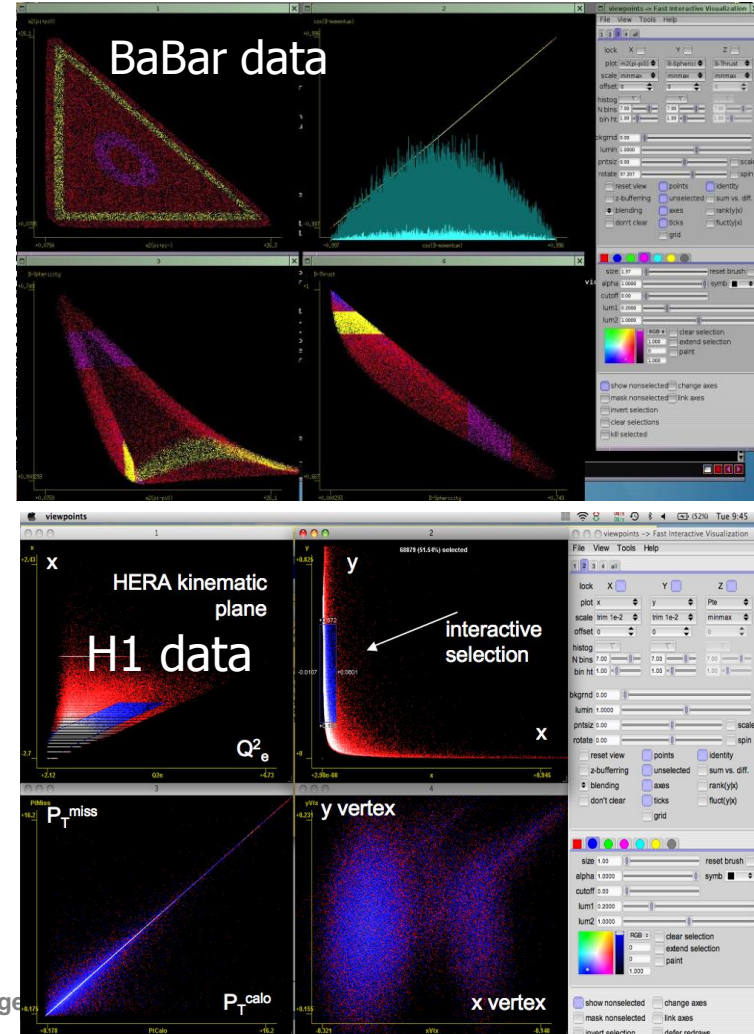
CPEP member Rush Holt elected to U.S. Congress

Outreach



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

Viewpoints (NASA)



HepEdu

- Discussions within DPHEP: format for outreach.
 - BaBar, Belle, H1.
 - First order attempt: text. Number of particles/tracks in the event
- ```
index PID E px py pz for particle 0
index PID E px py pz for particle 1
index PID E px py pz for particle 2
index PID E px py pz for particle 3
```

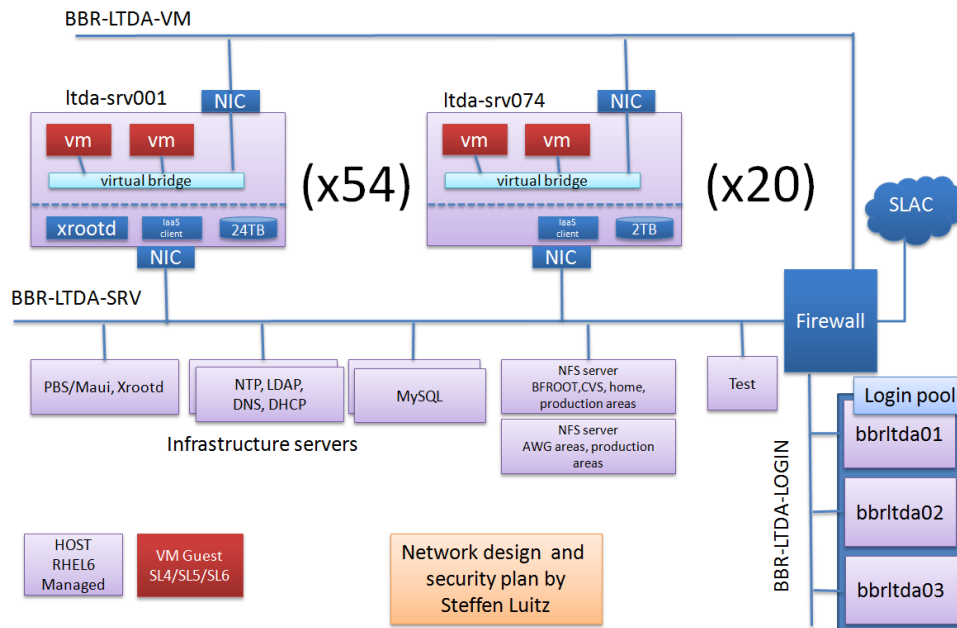
```
6
0 -13 1.313 0.407 1.241 0.075
1 -11 2.010 -1.813 0.039 -0.865
2 -211 0.474 -0.134 0.304 -0.308
3 211 0.480 -0.353 -0.112 -0.273
4 -211 1.003 -0.905 -0.369 0.176
5 22 0.212 -0.108 0.147 0.108
3
0 211 1.316 -0.414 -1.239 0.075
1 -211 2.014 -1.802 0.204 -0.865
2 211 0.474 -0.307 0.127 -0.308
```

- Discussions about common formats ongoing
  - B-lab (KEK) example considered
  - Experience at LHC
  - Connect to existing projects (master classes etc.)

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

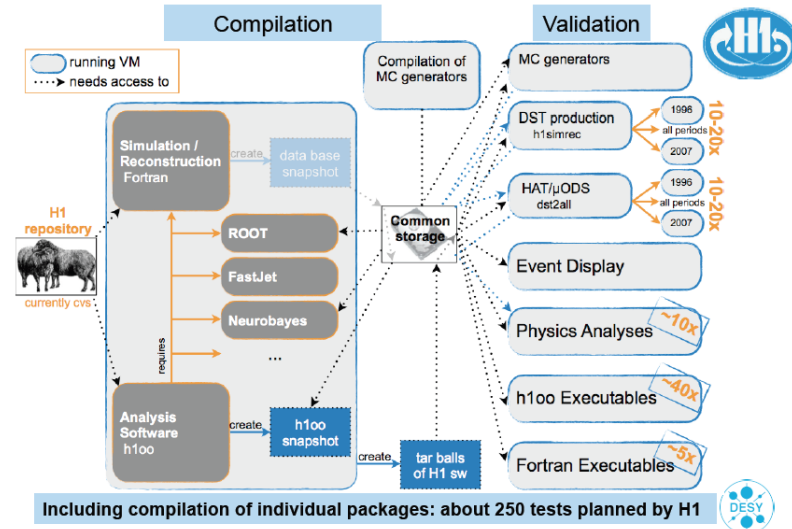
|                                                | BaBar                                                                                             | H1                                                    | ZEUS                    | HERMES                                                | Belle                                                | 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<br>Sim/rec level<br>Data skims in ROOT                                                   | RAW data<br>Sim/rec level<br>Analysis level ROOT data | Flat ROOT based ntuples | RAW data<br>Sim/rec level<br>Analysis level ROOT data | RAW data<br>Sim/rec level                            | RAW data<br>Sim/rec level<br>ROOT data      | RAW data<br>Rec. level<br>ROOT files (data+MC) | Raw data<br>Rec. level<br>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>Longévité recherchée: &gt; 10 ans</b>       |                                                                                                   |                                                       |                         |                                                       |                                                      |                                             |                                                |                                                |
| <b>Current operating system</b>                | SL/RHEL3<br>SL/RHEL 5                                                                             | SL5                                                   | SL5                     | SL3<br>SL5                                            | SL5/RHEL5                                            | SL5                                         | SL5<br>SL6                                     | SL5                                            |
| <b>Languages</b>                               | C++<br>Java<br>Python                                                                             | C<br>C++<br>Fortran<br>Python                         | C++                     | C<br>C++<br>Fortran<br>Python                         | C<br>C++<br>Fortran                                  | C++                                         | C<br>C++<br>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<br>CERNLIB<br>CLHEP<br>CMLOG<br>Flex<br>GNU Bison<br>MySQL<br>Oracle<br>ROOT<br>TCL<br>XRootD | CERNLIB<br>FastJet<br>NeuroBayes<br>Oracle<br>ROOT    | ROOT                    | ADAMO<br>CERNLIB<br>ROOT                              | Boost<br>CERNLIB<br>NeuroBayes<br>PostgreSQL<br>ROOT | CASTPR<br>CERNLIB<br>CLHEP<br>HepMC<br>ROOT | CERNLIB<br>NeuroBayes<br>Oracle<br>ROOT        | Oracle<br>ROOT                                 |

## Préservation d'un système d'accès et calcul à des données complexes (SLAC/Stanford USA)

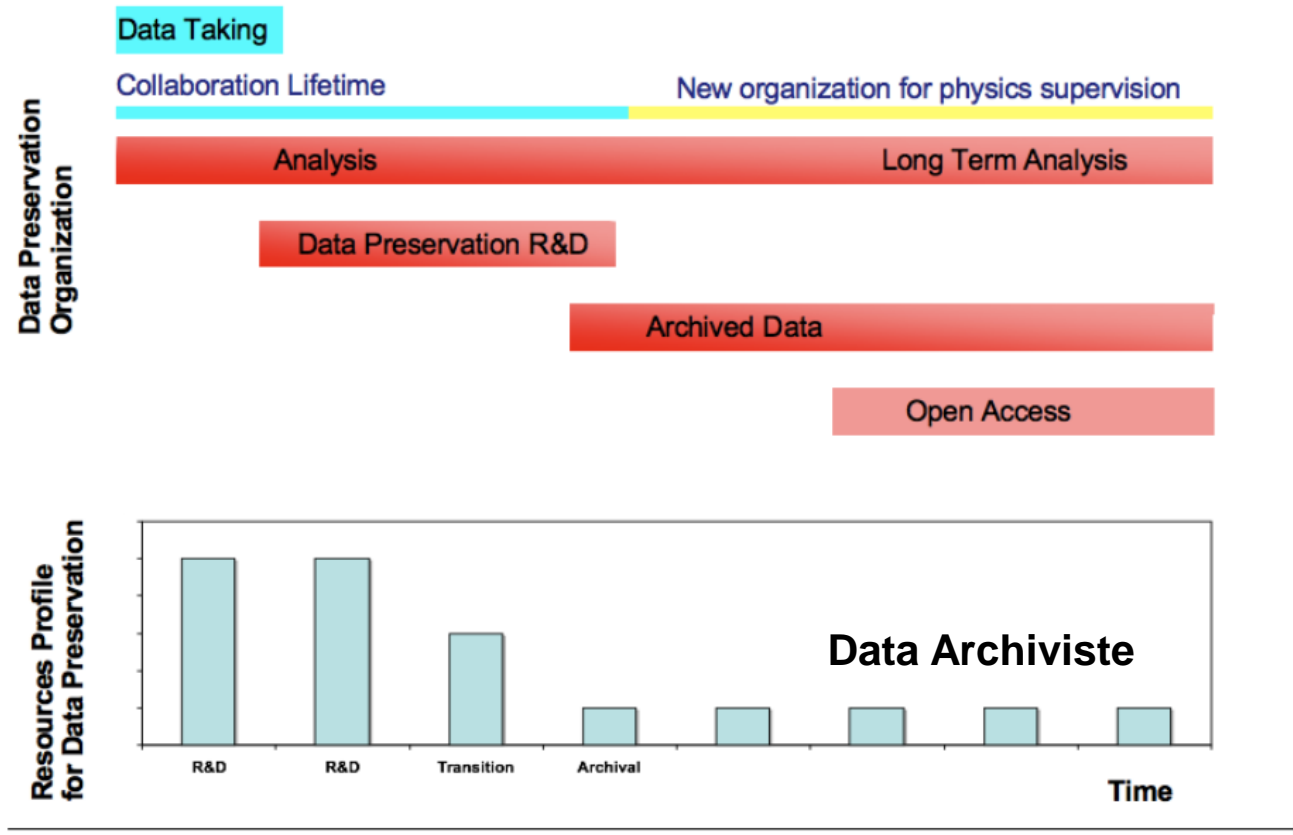


## Système de préservation et migration Virtualisation, validation intensive (DESY, Hambourg, Allemagne)

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



# Plan de sauvetage, vers un modèle économique



## Estimation du cout du projet:

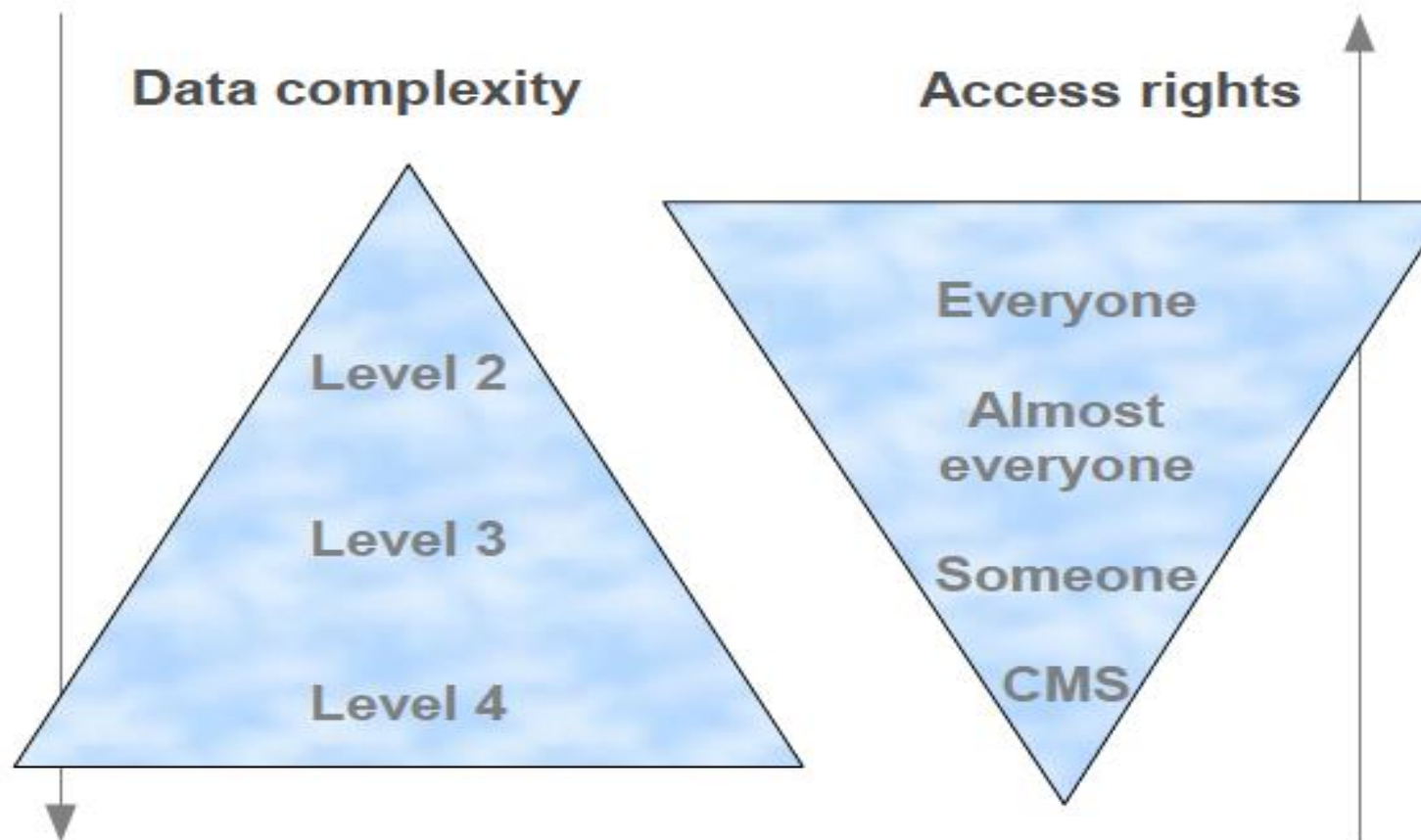
- >1% du cout total pour une production scientifique augmentée de 10%

## Des réflexions en cours sur:

- **Systèmes de stockage robustes, basse consommation etc.**
- **Couts de la préservation de données à long terme**

# Data Preservation and Open Access at the LHC

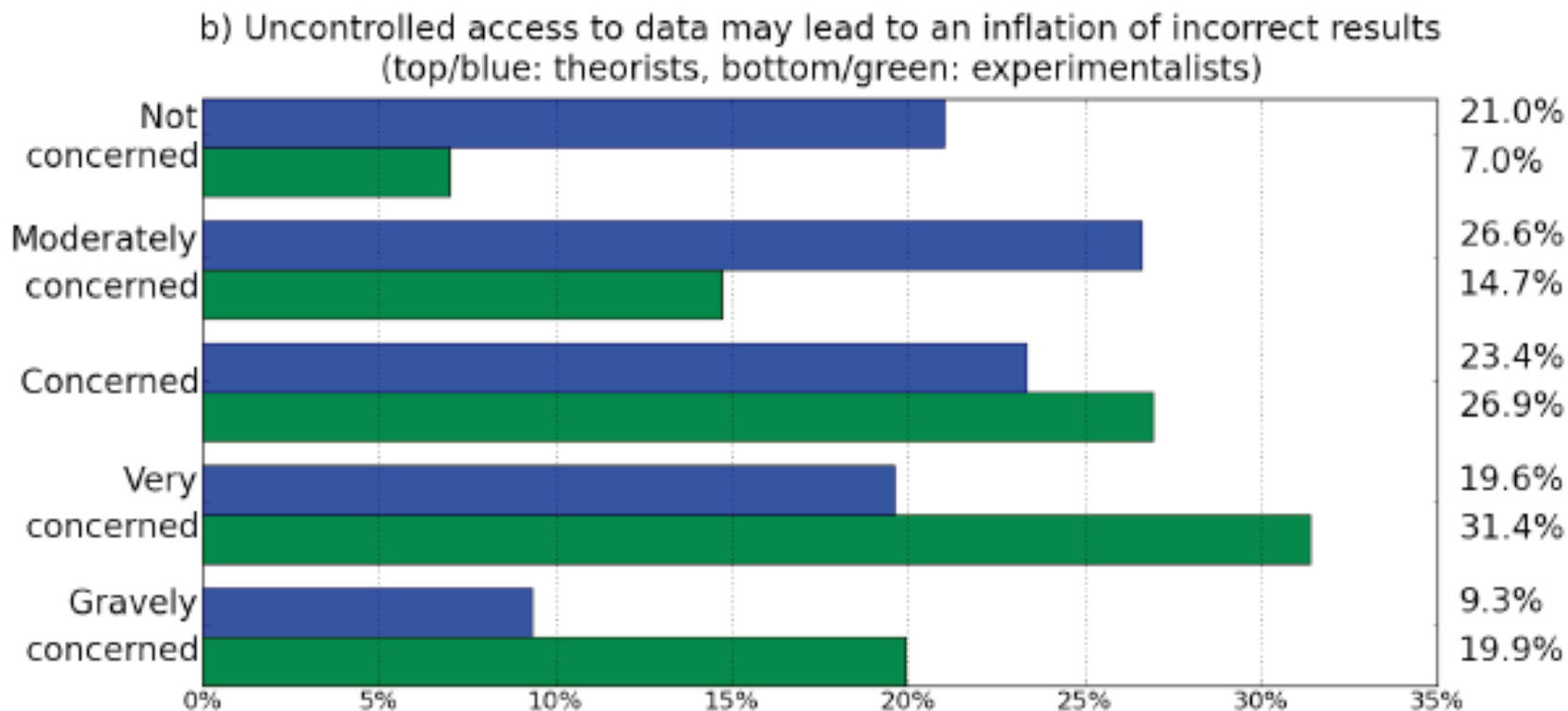
- Reflection just started in ATLAS, ALICE, CMS, LHCb
  - Common understanding that starting earlier will consolidate the long term future
  - Strong wish to develop a common policy at CERN and within DPHEP



# Est-ce que la réutilisation des données est risquée?

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

Parse.insight



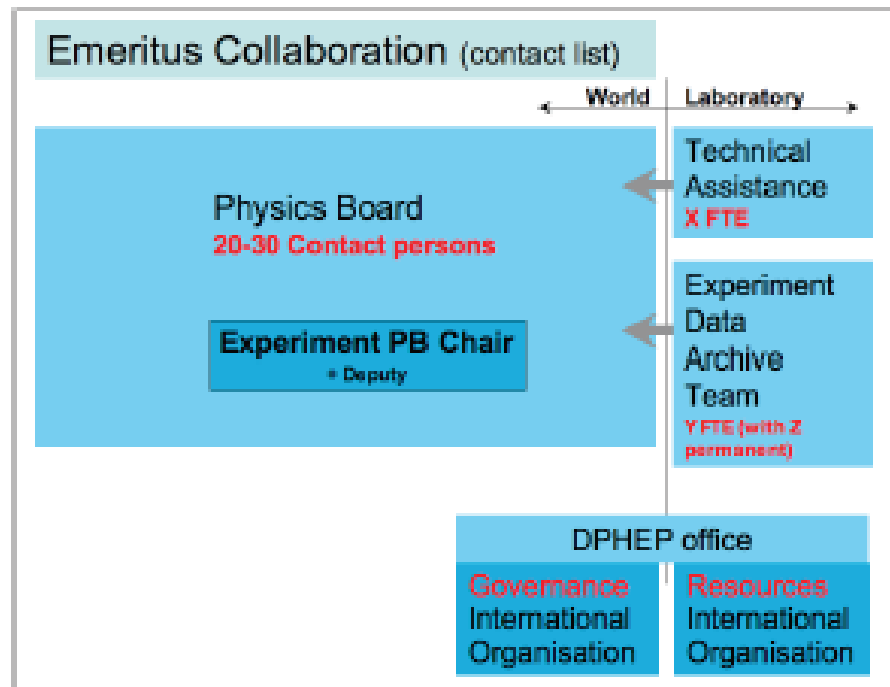
**Governance issues are very important to support data usage**

# 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

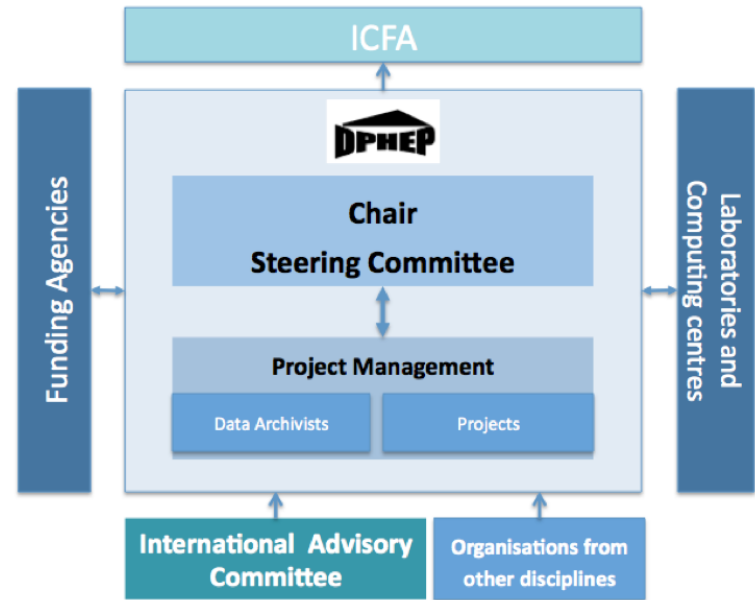
- Experimental organisation risks being left in an undefined state
- Transition should also be planned in advance of the projected end date



# DPHEP: Groupe d'études et organisation internationale



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



- > Study Group DPHEP:
  - > Participation des grands laboratoires (CERN, DESY, FERMILAB, SLAC, KEK, IHEP et experiences)
- > Organisation internationale en cours de mise en place
  - > 100 contact personnes de contact
  - > Chair: D. Diaconu Project Manager: Jamie Shiers (CERN)



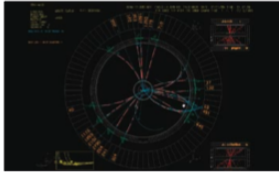
# DPHEP Visibility

CERN Courier, May 2009

DATA PRESERVATION

## Study group considers how to preserve data

For experimentalists in high-energy physics, the data are like treasure, but how can they be saved for the future? A study group is investigating data-preservation options.



A simulated event in the JADE detector, generated using a refined Monte Carlo program and reconstructed using revitalized software more than 10 years after the end of the experiment. (Courtesy Sigi Bethke.)

High-energy-physics experiments collect data over long time periods, while the associated collaborations of experimentalists exploit these data to produce their scientific publications. The scientific potential of an experiment is in principle defined and exhausted within the lifetime of such collaborations. However, the continuous improvement in areas of theory, experiment and simulation – as well as the need to re-analyse old data. Examples of such analyses already exist and they are likely to become more frequent in the future. As experimental complexity and the associated costs continue to increase, many present-day experiments, especially those based at colliders, will provide unique data sets that are unlikely to be improved upon in the short term. The close of the current decade will see the end of data-taking at several large experiments and scientists are now confronted with the question of how to preserve

the complexity of the hardware and a more dynamic part closer to the ROOT analysis environment and is in most cases done in C++ using the ROOT analysis environment and is mainly performed on local computing farms. Monte Carlo simulation also uses a farm-based approach but it is striking to see how popular the Grid is for the mass-produced events. The amount of data that should be saved varies between 0.5 PB and 10 PB for each not huge by today's standards but nonetheless large degree of preparation for long term data varies 1 but it is obvious that no preparation was fore- of the programs; any conservation initiatives ill with the end of the data analysis.



February 2011

## Rescue of Old Data Offers Lesson for Particle Physicists

Old data tends to get forgotten as physicists move on to new and better machines.



May 2011

Data Preservation

- ICFA Study Group on Data Preservation and Long Term Analysis in High Energy Physics. High Energy Physics experiments initiated with this Study Group a common reflection on data persistence and long term analysis in order to get a common vision on these issues and create a multi-experiment dynamics for further reference:

<https://www.dphep.org/>



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

symmetry dimensions of particle physics

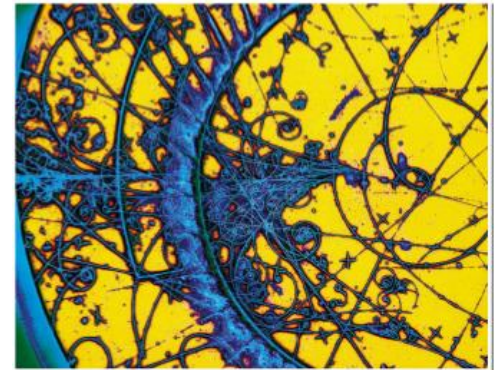
A joint Fermilab/SLAC publication

VOLUME 06 ISSUE 06 DECEMBER 09

Symmetry, December 2009

Berliner Zeitung, Nummer 50, Dienstag, 10. Februar 2010

## Wissenschaft



Beschleuniger gehören ab dem Herbst zu den wichtigsten Instrumenten der Teilchenphysik. Sie sind mit Millionen Wasserstoff- und Kohlenstoff-Teilchen, die durch die Beschleuniger in die Teilchenbeschleuniger beschleunigt werden, um die Teilchenphysik zu untersuchen.

## Die Hieroglyphen von morgen

An Beschleunigern sind immense Datenmengen entstanden – die Archivierung beginnt erst jetzt

von Thomas Kluge

Was der neue Teilchenbeschleuniger LHC die Elementarteilchenphysik in den nächsten Jahren bringen wird, ist noch unklar. Die Experimente werden aber in den nächsten Jahren die Teilchenphysik revolutionieren. Die Datenmengen werden dabei immens zunehmen. Die Datenmengen werden dabei immens zunehmen. Die Datenmengen werden dabei immens zunehmen.

### Der Teilchenzoo

Die Teilchenphysik ist ein Bereich der Physik, der sich mit den kleinsten Bausteinen der Materie beschäftigt. Die Teilchenphysik ist ein Bereich der Physik, der sich mit den kleinsten Bausteinen der Materie beschäftigt.

### Kosten der Beschleuniger werden steigen

Die Kosten für die Beschleuniger werden in den nächsten Jahren stark ansteigen. Die Kosten für die Beschleuniger werden in den nächsten Jahren stark ansteigen.

### Meilenstein Weg

Die Teilchenphysik hat in den letzten Jahren viele Meilensteine erreicht. Die Teilchenphysik hat in den letzten Jahren viele Meilensteine erreicht.

### Mit steigendem Nachfrager werden die Teilchenbeschleuniger LHC gebaut

Die Teilchenbeschleuniger LHC werden in den nächsten Jahren gebaut. Die Teilchenbeschleuniger LHC werden in den nächsten Jahren gebaut.

Die Teilchenphysik ist ein Bereich der Physik, der sich mit den kleinsten Bausteinen der Materie beschäftigt. Die Teilchenphysik ist ein Bereich der Physik, der sich mit den kleinsten Bausteinen der Materie beschäftigt.

### Die Datenmengen werden steigen

Die Datenmengen werden in den nächsten Jahren stark ansteigen. Die Datenmengen werden in den nächsten Jahren stark ansteigen.

### Meilenstein Weg

Die Teilchenphysik hat in den letzten Jahren viele Meilensteine erreicht. Die Teilchenphysik hat in den letzten Jahren viele Meilensteine erreicht.

### Mit steigendem Nachfrager werden die Teilchenbeschleuniger LHC gebaut

Die Teilchenbeschleuniger LHC werden in den nächsten Jahren gebaut. Die Teilchenbeschleuniger LHC werden in den nächsten Jahren gebaut.

Berliner Zeitung und Frankfurter Rundschau, February 2010

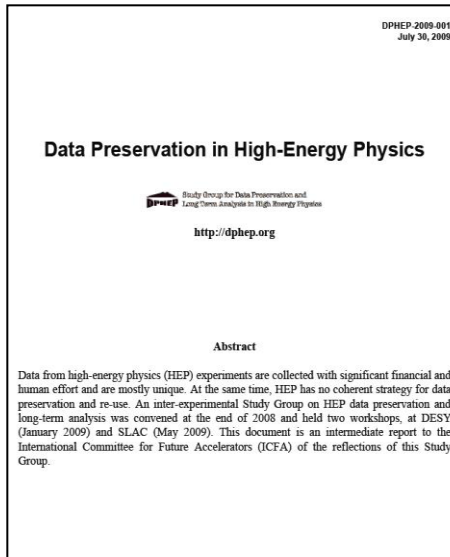


C. Diaconu, Preservation des dpm

age 25

# DPHEP Intermediate Recommendations (end 2009)

> [arXiv:0912.0255](https://arxiv.org/abs/0912.0255)

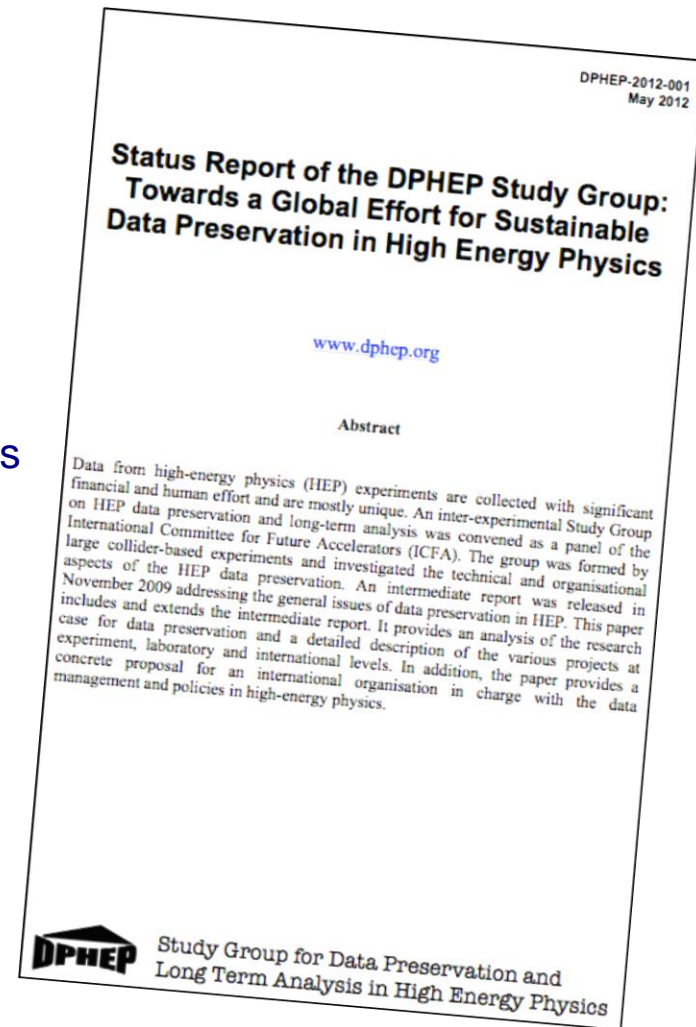


- > An urgent and vigorous action is needed to ensure data preservation in HEP
  - Many examples for the physics case explored
  - Data is rich and can be further exploited in most cases beyond the collaboration lifetime
- > The preservation of the full analysis capability of experiments is recommended, 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

- 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:1205.466

7



# ICFA Statement on LTDP

## ICFA: International Committee for Future Accelerators

- > *The International Committee for Future Accelerators (ICFA) supports the efforts of the Data Preservation in High Energy Physics (DPHEP) study group on long-term data preservation and welcomes its transition to an active international collaboration with a full-time project manager. **It encourages laboratories, institutes and experiments to review the draft DPHEP Collaboration Agreement with a view to joining by mid- to late-2013.***
- > *ICFA notes the lack of effort available to pursue these activities in the short-term and the possible consequences on data preservation in the medium to long-term. **We further note the opportunities in this area for international collaboration with other disciplines and encourage the DPHEP Collaboration to vigorously pursue its activities.** In particular, the effort required to prepare project proposals must be prioritized, in addition to supporting on-going data preservation activities.*
- > ***ICFA notes the important benefits of long-term data preservation to exploit the full scientific potential of the, often unique, datasets.** This potential includes not only future scientific publications but also educational outreach purposes, and the Open Access policies emerging from the funding agencies.*
- > 15 March 2013

# Data Preservation in a multidisciplinary context

- > **More Coordination:** The organisation should be brought to a long-term perspective by solid, commensurate and courageous decisions of the funding and coordination bodies responsible for the wealth of HEP experimental data produced so far.
- > **More Standards** An increased standardisation will increase the overall efficiency of HEP computing systems and it will also be beneficial in securing long-term data preservation.
- > **More Technology:** These new techniques (virtualisation etc.) seem to fit well within the context of large scale and long-term data preservation and access.
- > **More Experiments:** The expansion of the DPHEP organisation to include more experiments is one of the goals of the next period.
- > **More Cooperation: Cooperation with other fields in data management: access, mining, analysis and preservation; appears to be unavoidable and will also dramatically change the management of HEP data in the future.**

# Conclusions HEP

- > Les données scientifiques ont un potentiel qui dépasse le cadre de recherche initial et qui doit être exploité à long terme
- > La préservation de données scientifique est économiquement avantageuse:
  - Recherche à bas cout
- > Une technologies de frontière est nécessaire
  - Virtualisation, cloud computing, workflows....
  - Expertise IST essentielle
- > La collaboration internationale est essentielle
- > La collaboration multi-disciplinaire est nécessaire
  - Méthodes, technologies, approches (PREDON)