

Scientific Data Preservation

C. Diaconu

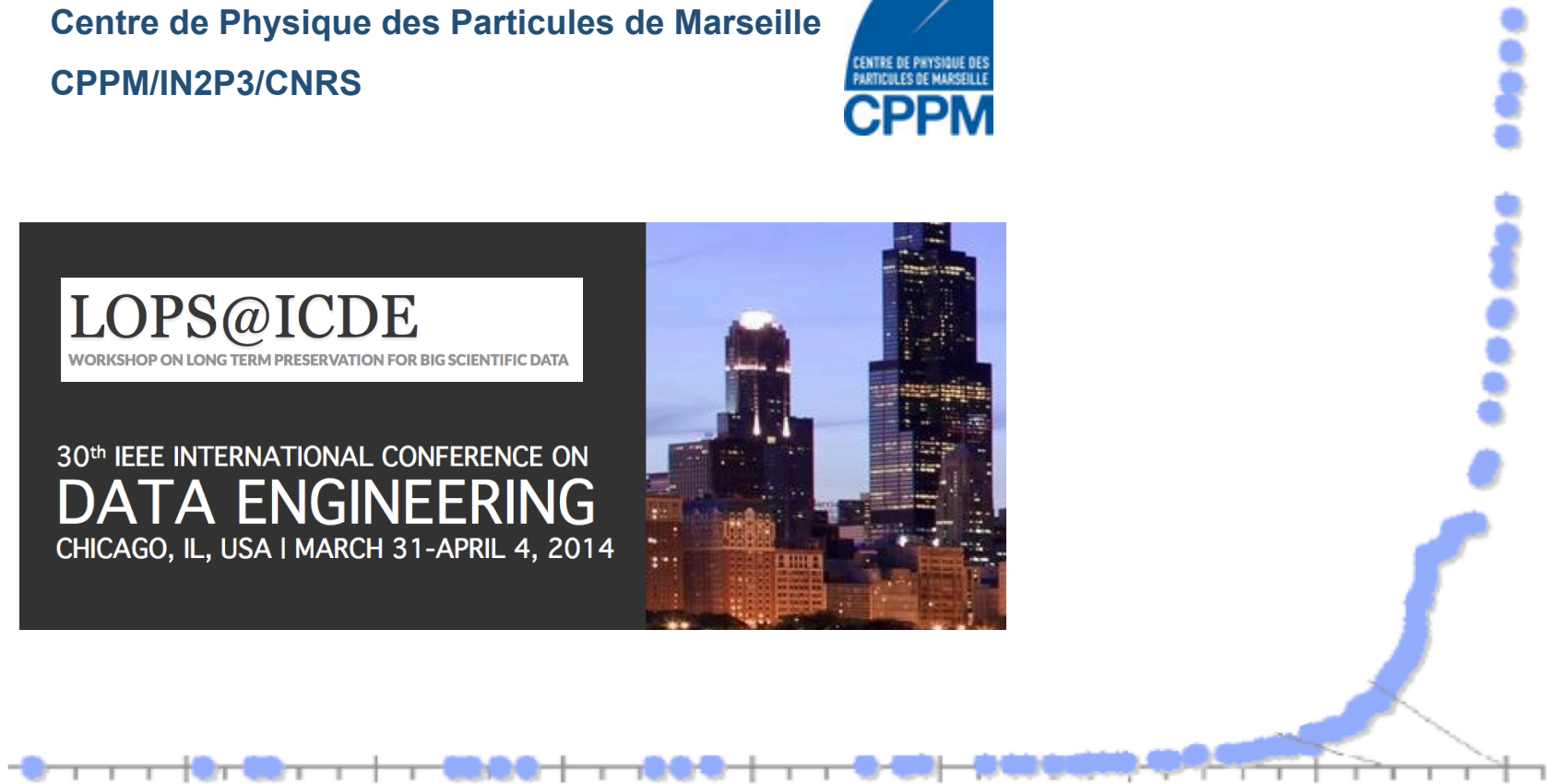
Centre de Physique des Particules de Marseille
CPPM/IN2P3/CNRS



LOPS@ICDE

WORKSHOP ON LONG TERM PRESERVATION FOR BIG SCIENTIFIC DATA

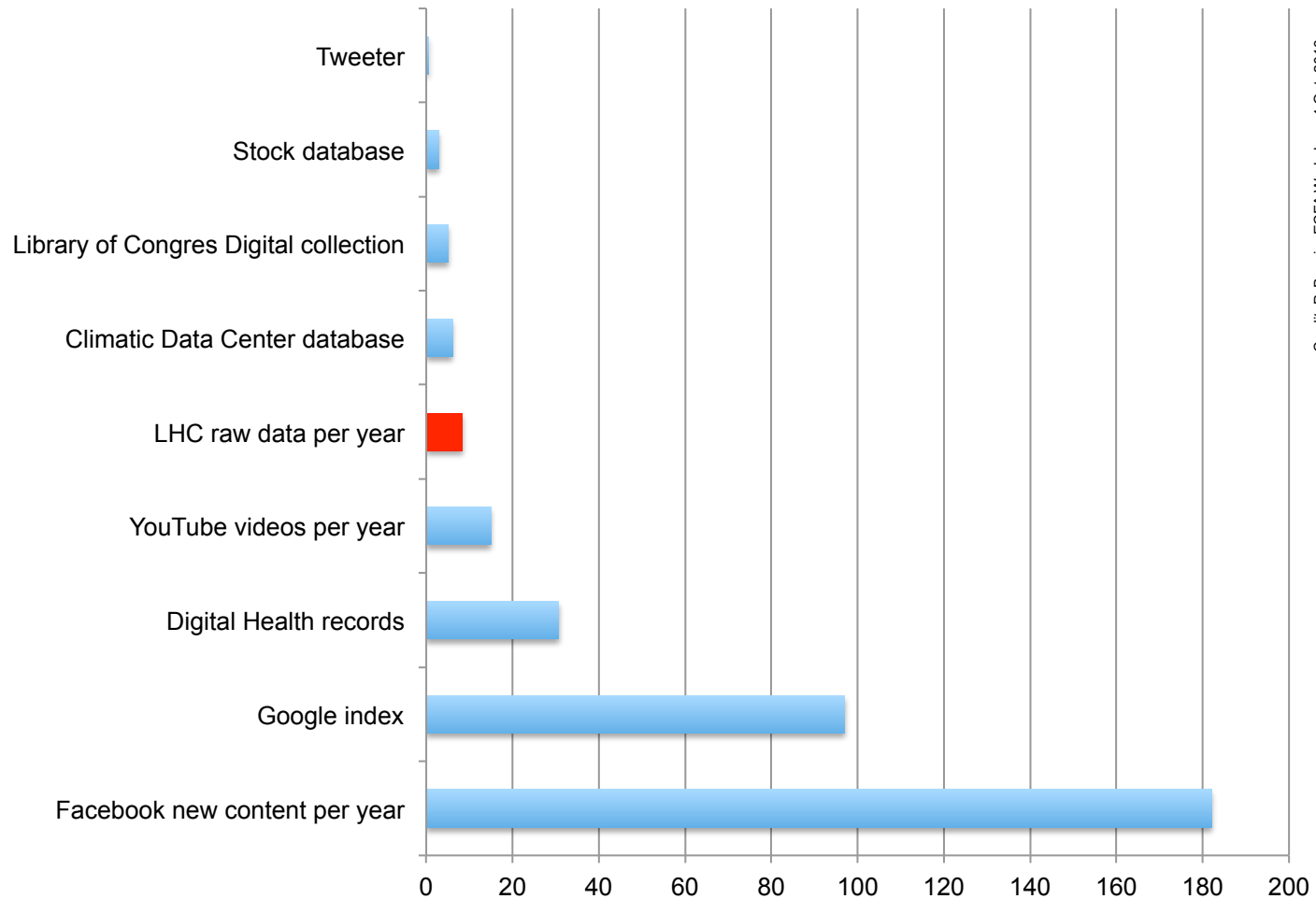
30th IEEE INTERNATIONAL CONFERENCE ON
DATA ENGINEERING
CHICAGO, IL, USA | MARCH 31-APRIL 4, 2014



Data Big Bang



Data collection accelerates



Credit: P. Buncic, ECFA Workshop, 4 Oct. 2013

PB

Digital data are fragile

- Storage capacity is physically exceeded
- Unattended/orphaned data vanishes quickly

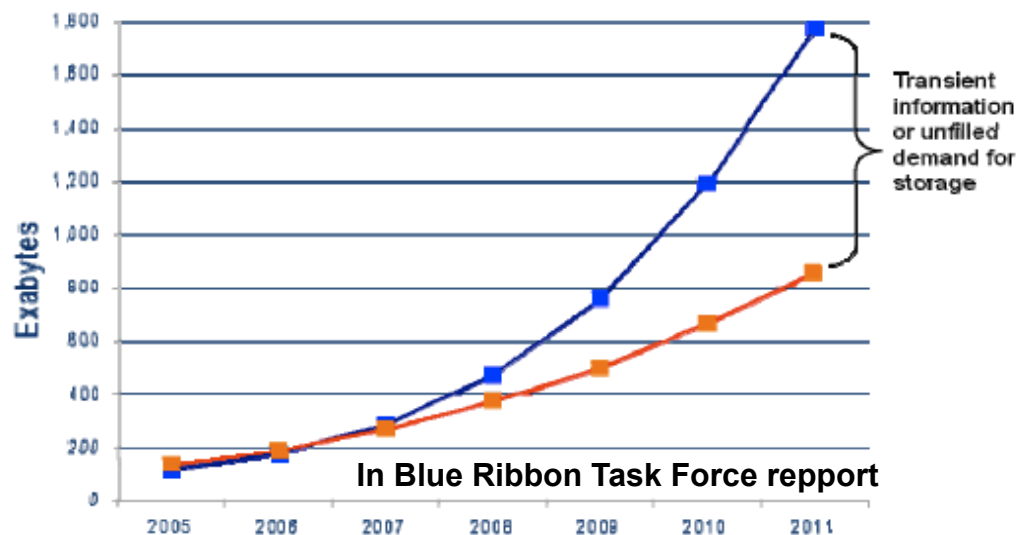


FIGURE 1.3: **Information and Storage**

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

Models of data preservation and access

> Collaborations address this issue in a generic way

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

<http://www.alliancepermanentaccess.eu>
<http://brtf.sdsc.edu>

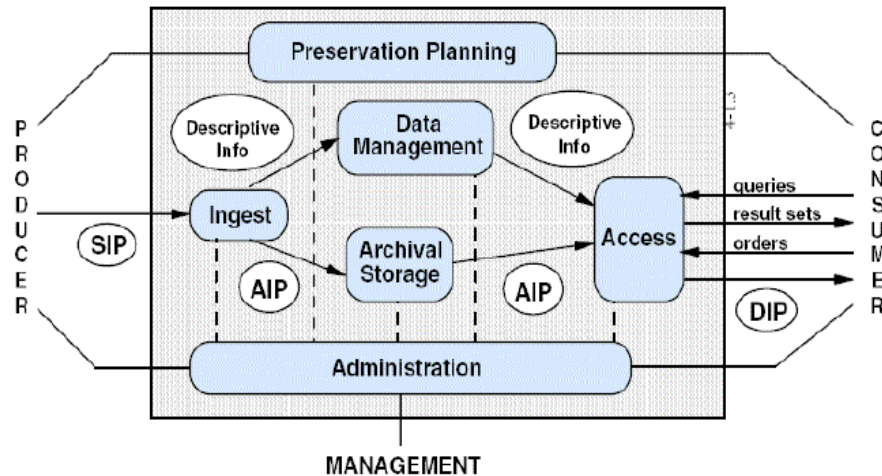


FIGURE 2.1: **The OAIS Reference Model**

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

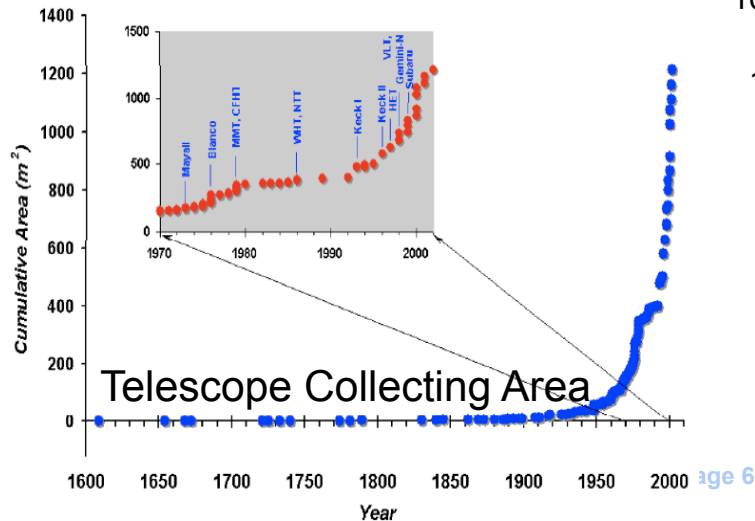
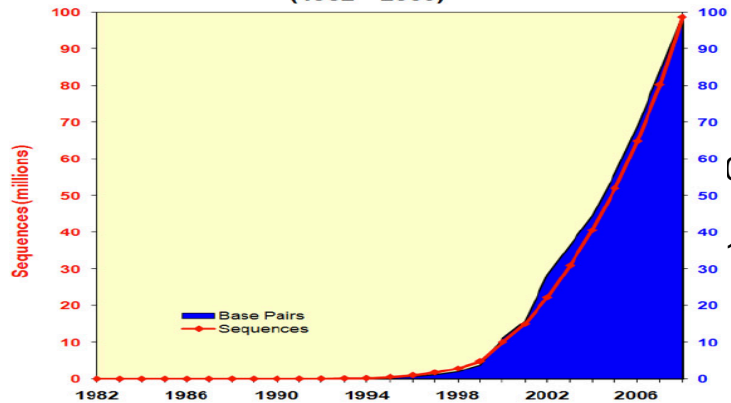
Source: Consultative Committee for Space Data Systems January 2002.

> Scientific Data is a major component of the ongoing efforts (complexity)

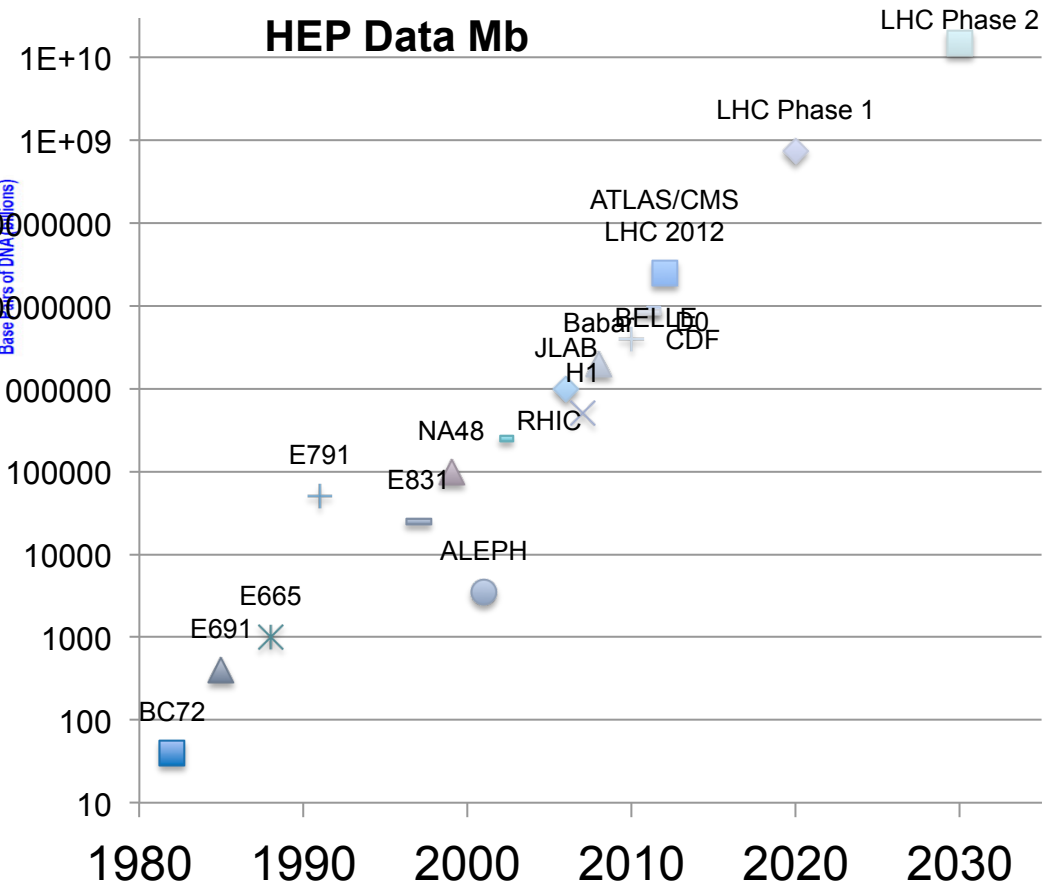
Big Scientific Data

- Scientific research observes a dramatic increase in data and are questioning the long term future of this data

Growth of GenBank
(1982 - 2008)



HEP Data Mb



Scientific Data

- > Structured following a scientific plan
- > Diverse sources
- > Large and expensive projects
 - Not easy to repeat
- > Contain unique knowledge
 - « Time stamped »
- > Data Observatories
 - Contain more information than initially needed

Scientific endeavours: big and complex

- > High energy physics projects
- > Large Hadron Collider (27 Km, 13 TeV, 40MHz)
- > 100Pb/year

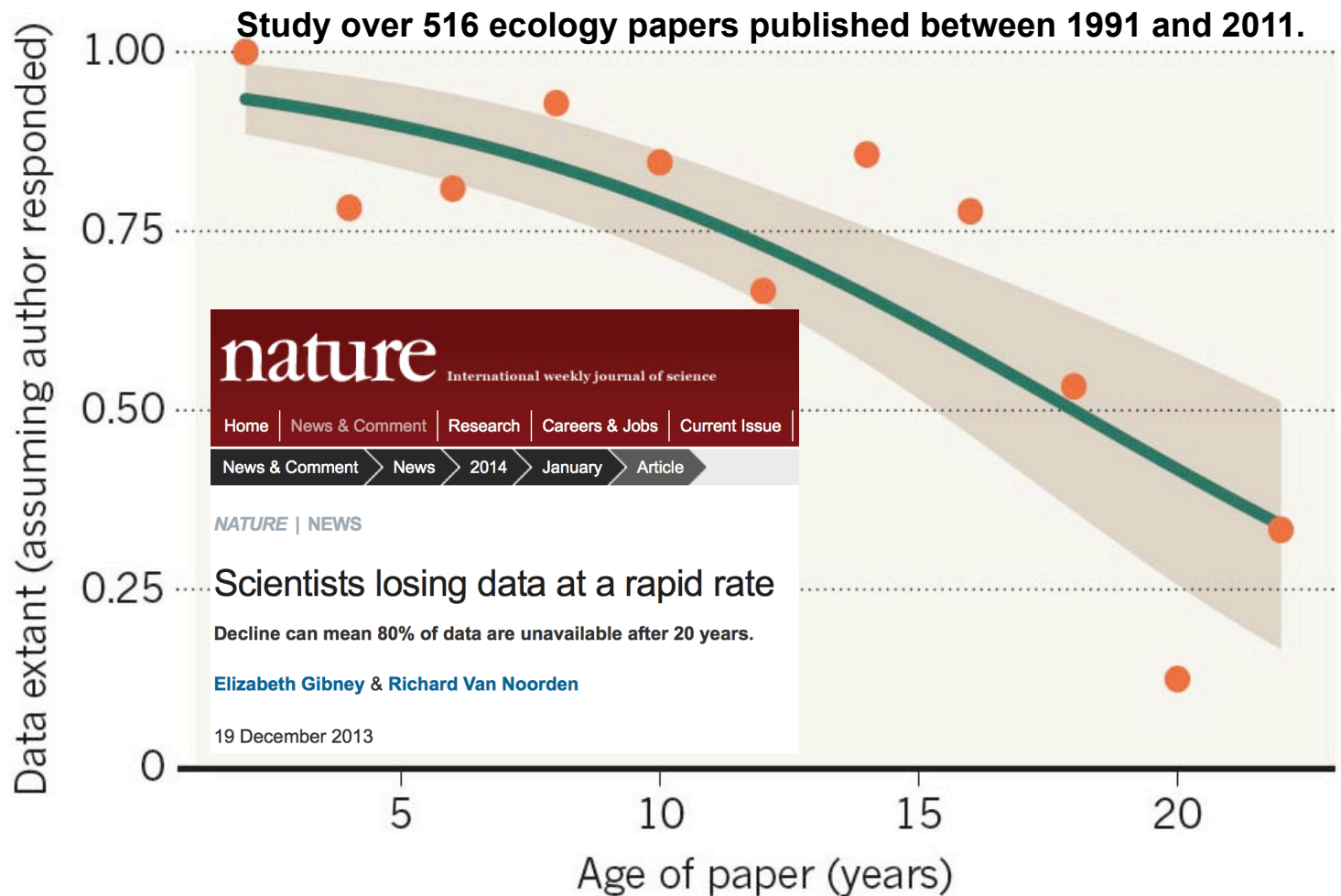


Another project

Where is your data?

MISSING DATA

As research articles age, the odds of their raw data being extant drop dramatically.



Preservation: where is the problem?



NATURE | NEWS

عربي

LHC plans for open data future

Researchers share results to keep them accessible.

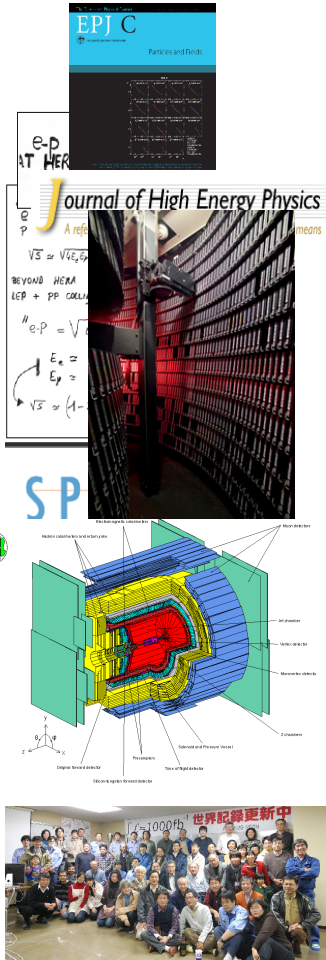
Elizabeth Gibney

26 November 2013

“When the LHC programme comes to an end, it will probably be the last data at this frontier for many years. We can’t afford to lose it.”

Storing the data is not a problem: hard drives are cheap and getting cheaper. The challenge is preserving knowledge that is less commonly stored — the software, algorithms and reference plots specific to each experiment. These often degrade or disappear with time, says Cristinel Diaconu of the Marseilles Centre for Particle Physics in France, who is chair of the international Data Preservation in Long Term Analysis in High Energy Physics (DPHEP) study group. He worries that if the data continue to be stored in their current state, physicists trying to decipher them in 10 years’ time will be unable to reconstruct the discovery of the Higgs boson. “When the LHC programme comes to an end, it will probably be the last data at this frontier for many years,” he says. “We can’t afford to lose it.”

Scientific Data: what is it?

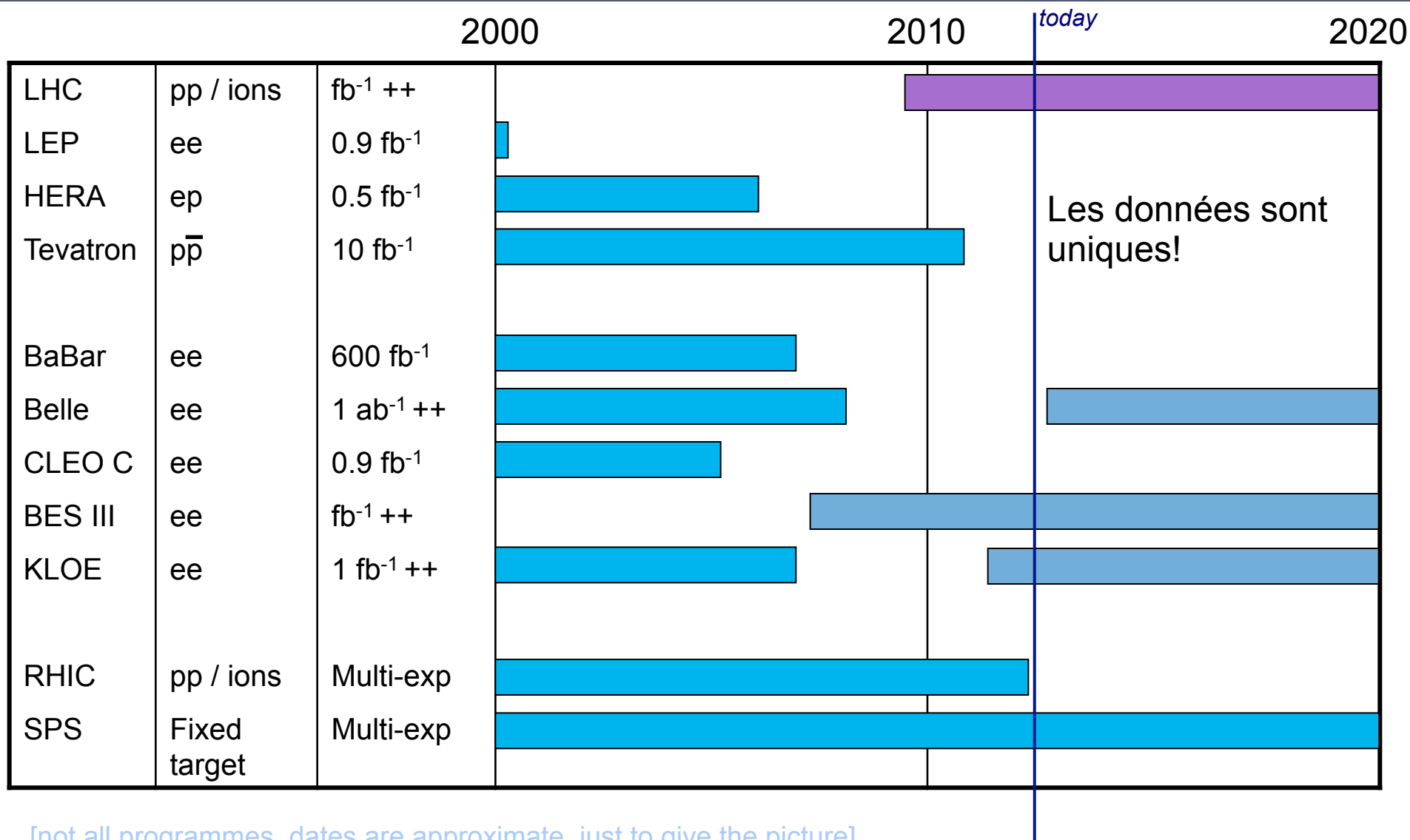


- Publications
- Documentation
- Raw
- Données Processées
- Meta-données
- Workflows
- Software
- Diffuse knowledge
-more...

Complexité, couts
Utilité

Technologie,
méthodologie
Organisation

Exemple: HEP experiments in ± 10 ans



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

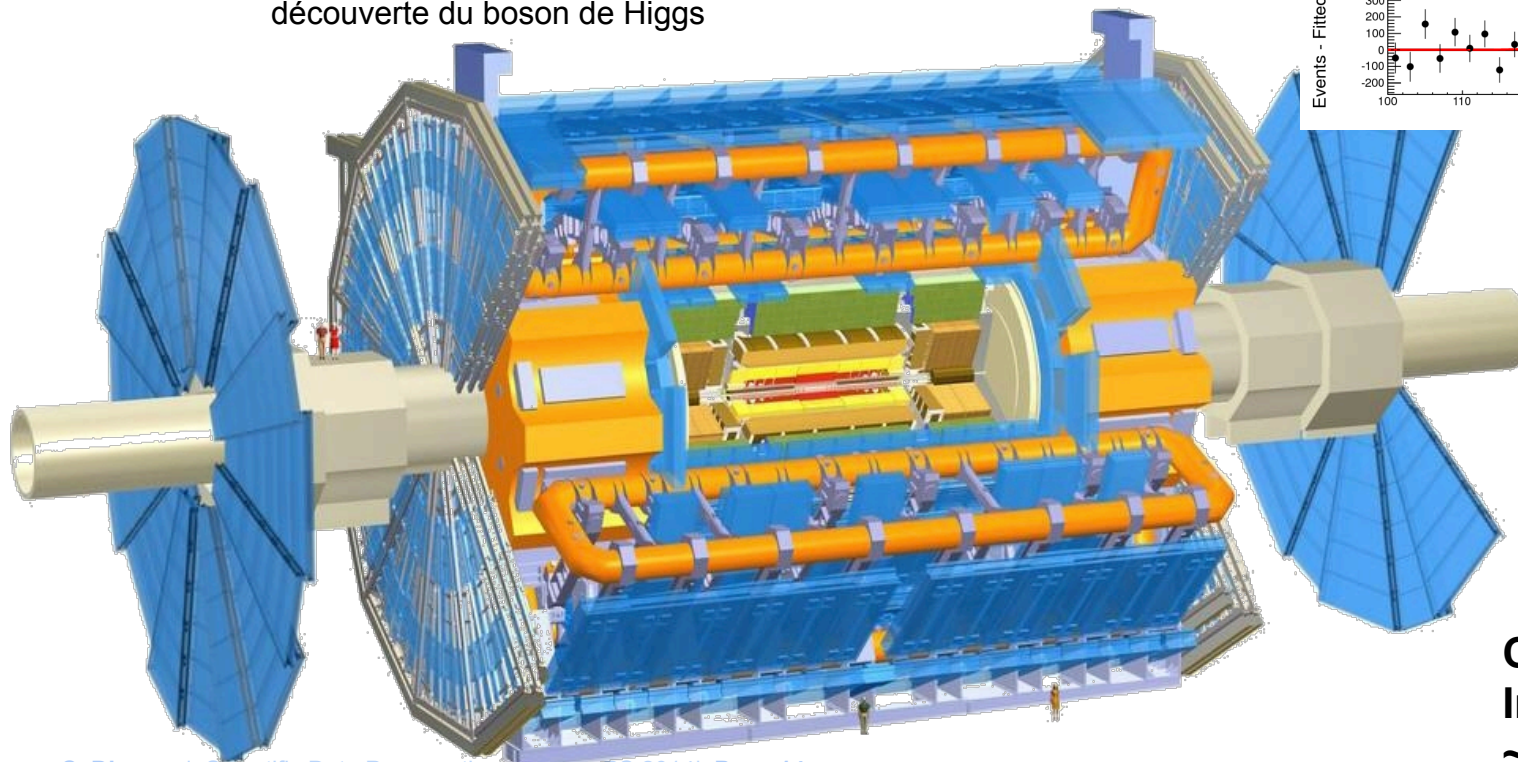
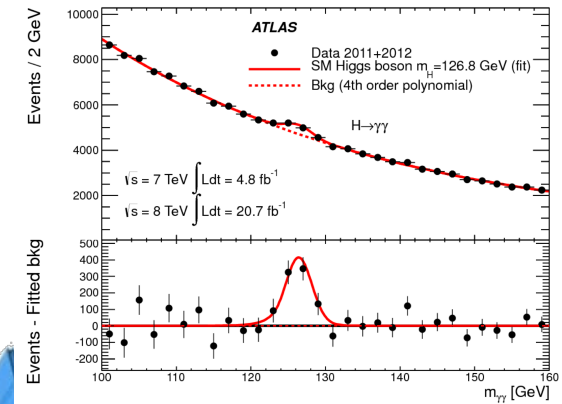
Complex experiments

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

> manip LHC:

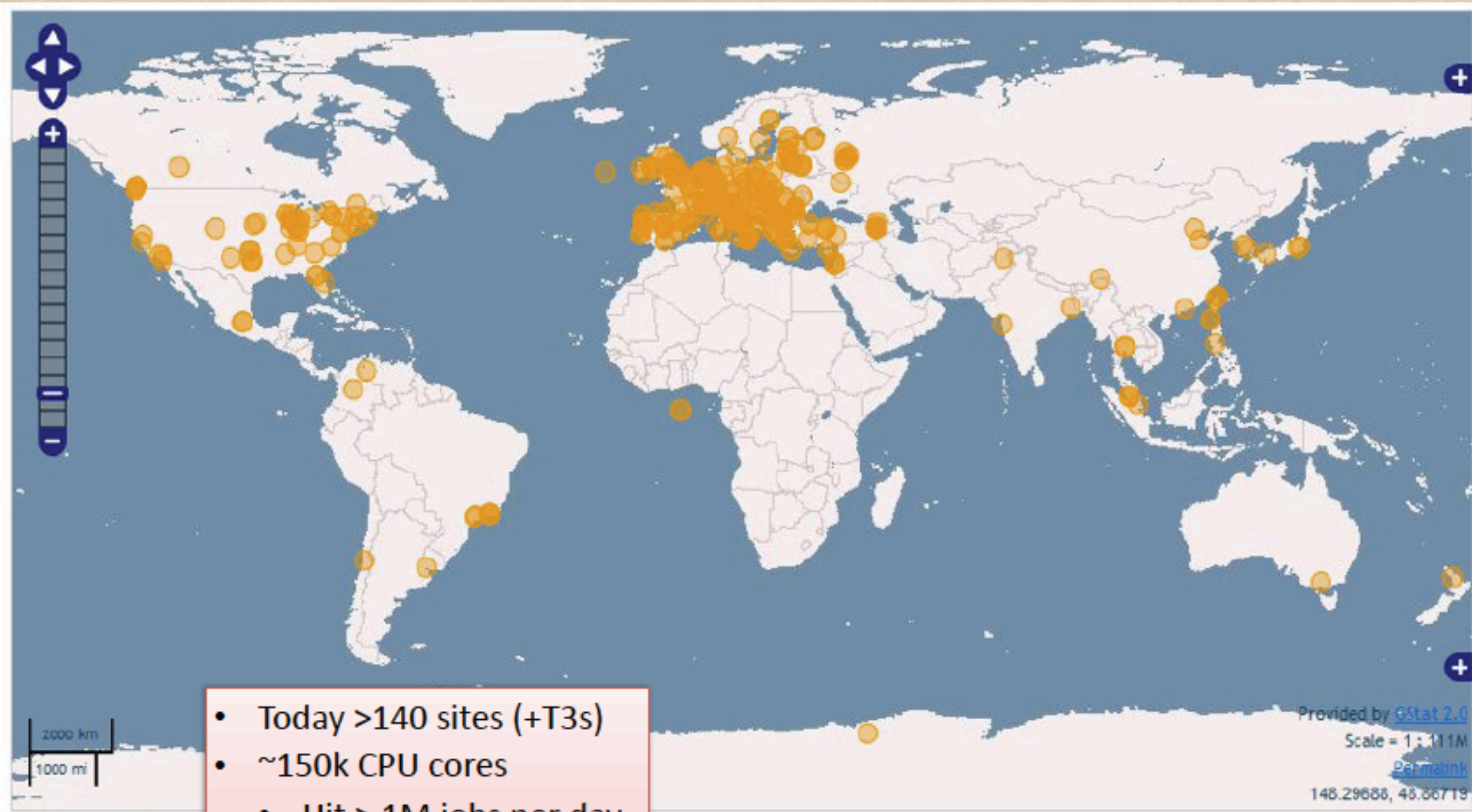
- 1000 articles scientifiques en 2 ans

découverte du boson de Higgs

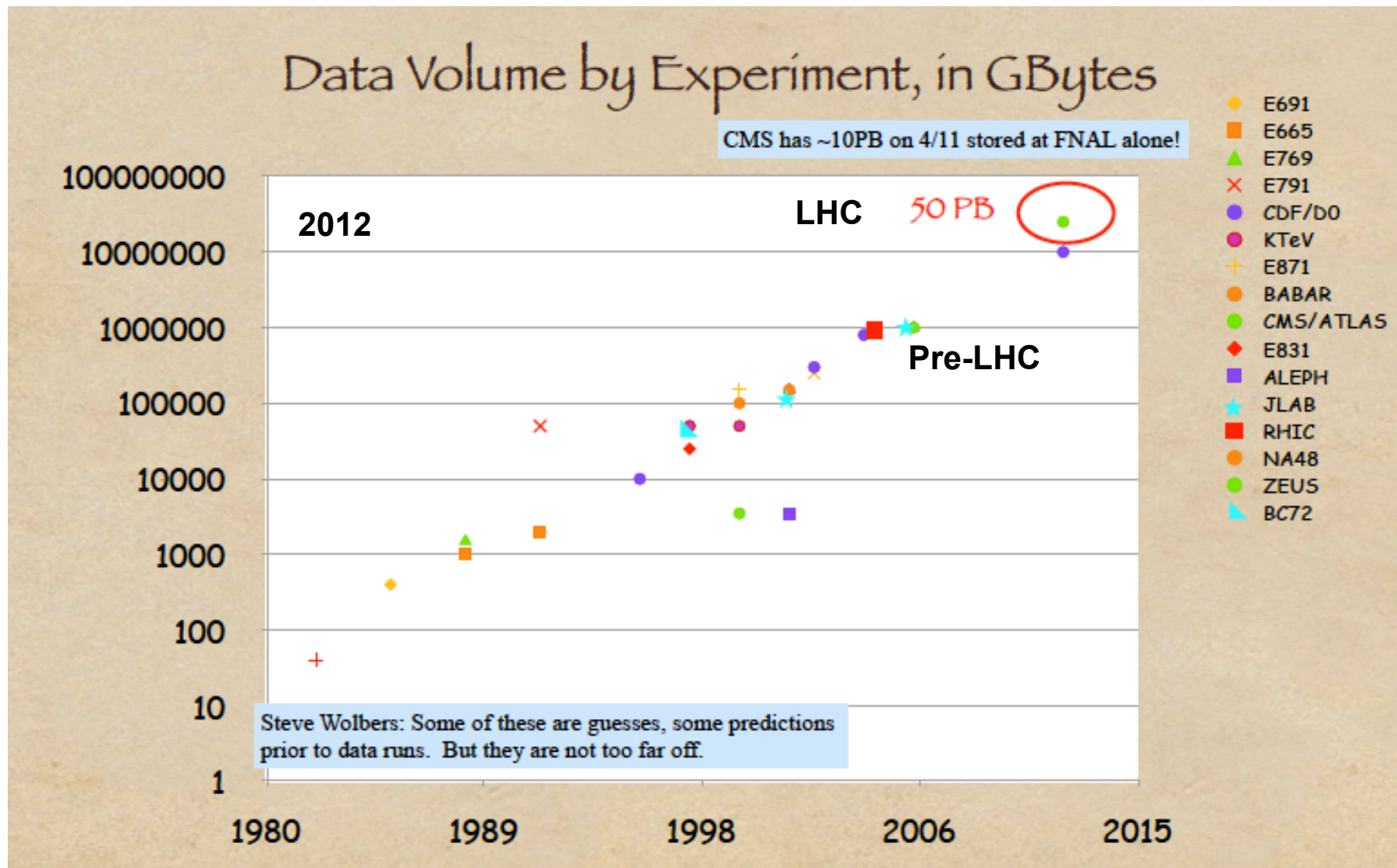


**Collaboration
Internationale
~3000 chercheurs**

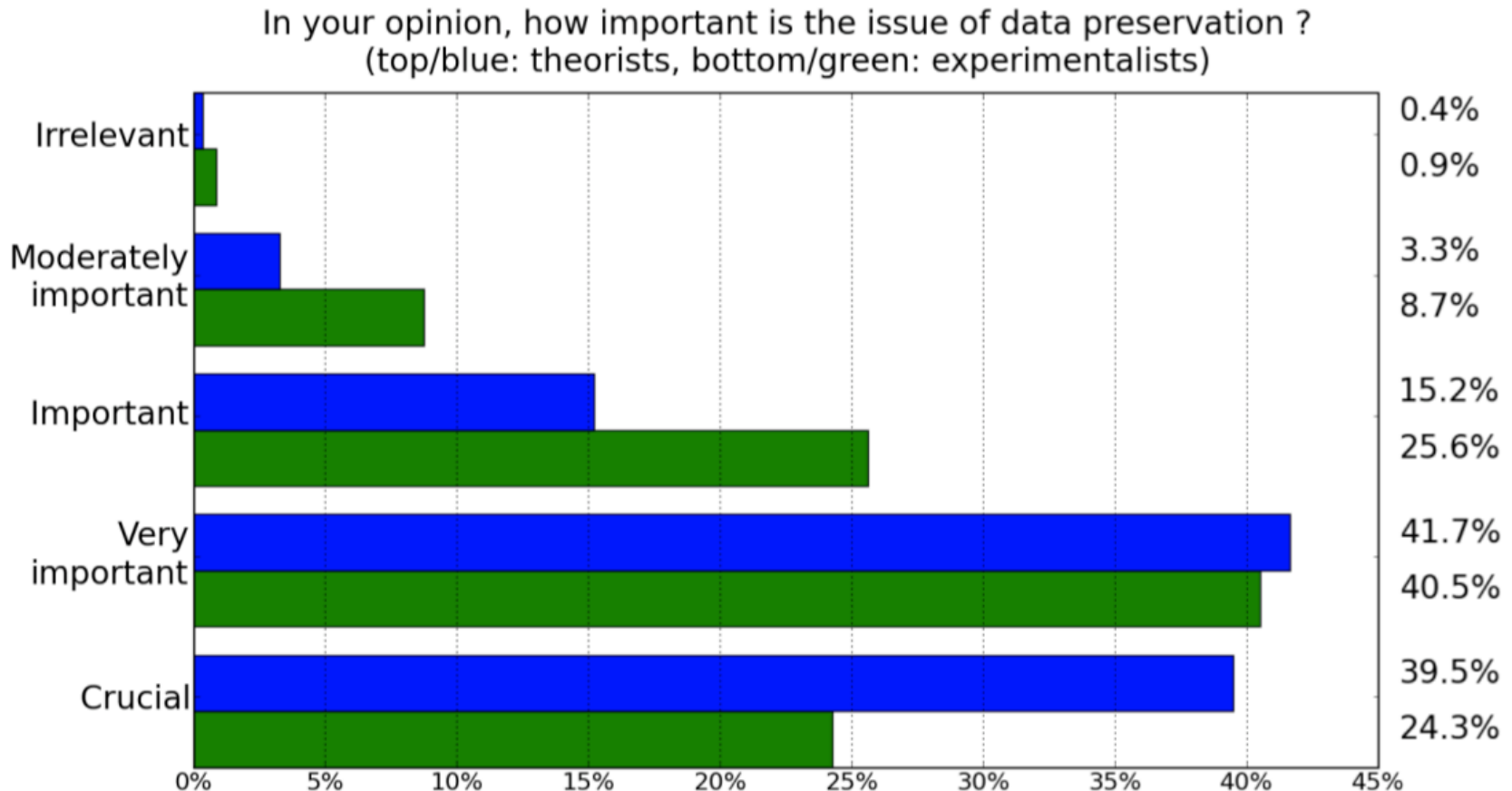
LHC computing



Data sets in time: 1PB -> 100PB->1EB



Should these data be preserved?

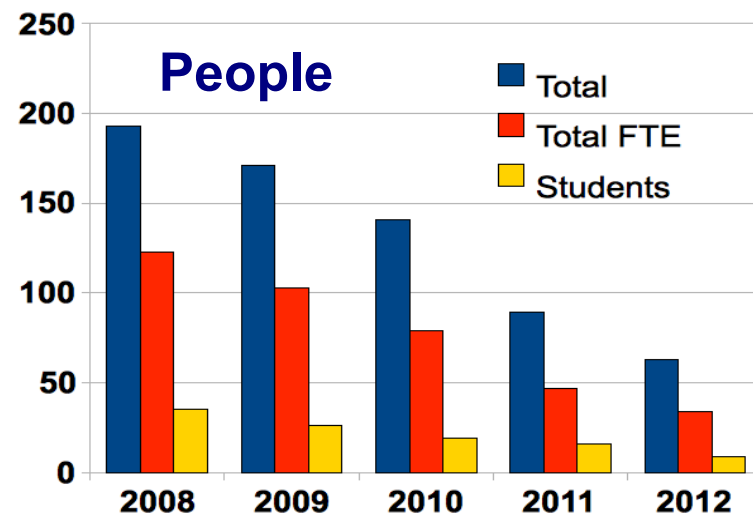
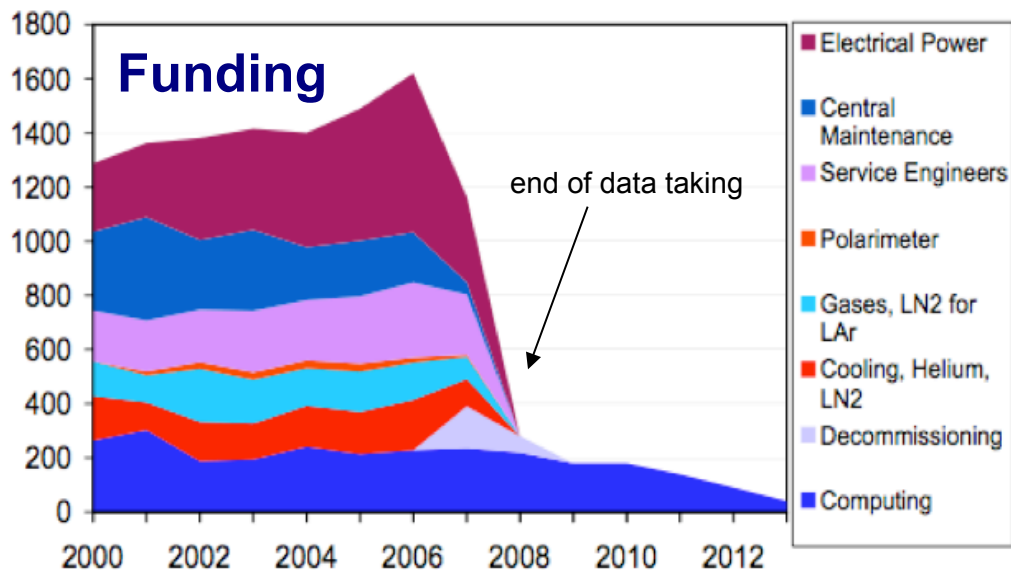
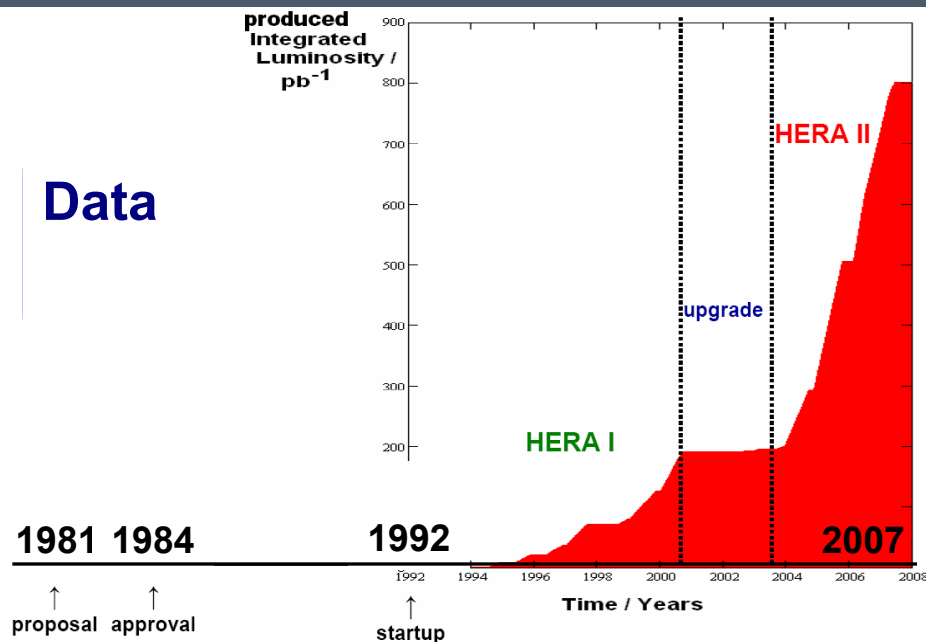


arXiv:0906.0485

Why there may be a problem?

- > End of the programs is difficult
 - > Resources
 - > Experts
- > Upfront planning is essential

Data



The email you may receive one day (I did)

Dear Dr. Diaconu,

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.

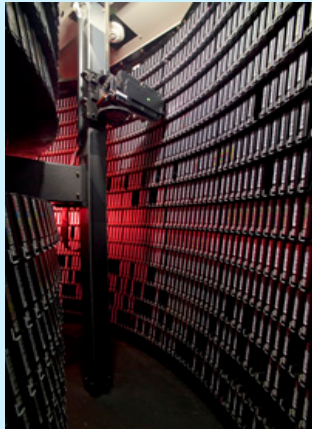
Would you like us to send you these tapes or should we **destroy them directly?**

Yours Sincerely,

Tape admin. service [a large computing centre]



What is “data”?



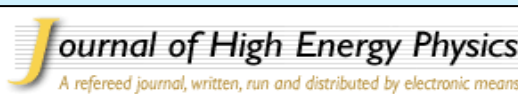
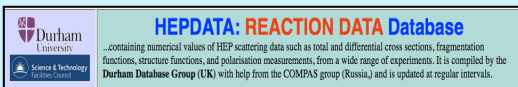
Digital information

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

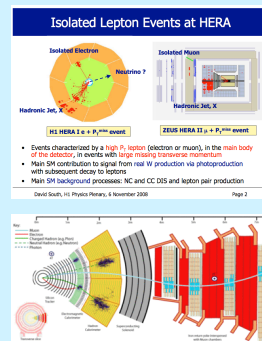
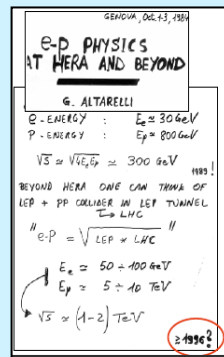
Publications

arXiv.org

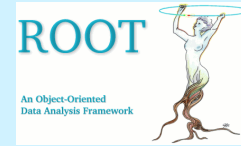
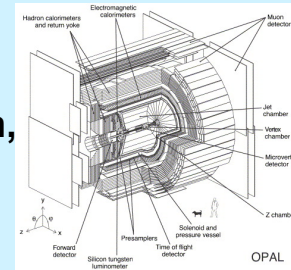


Documentation

Internal publications, notes, manuals, slides



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



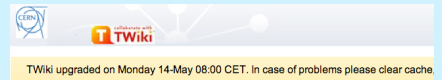
Software License for LHC++

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..



Welcome to TWiki at CERN.

TWiki is a flexible, powerful, secure, yet simple web-based collaboration platform.



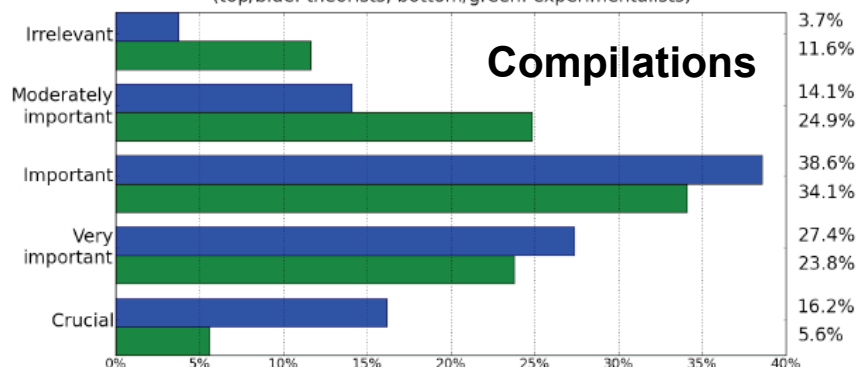
Expertise and people



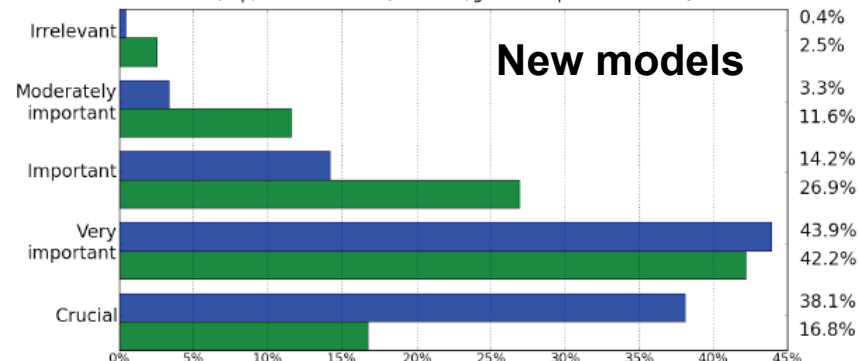
Data usage

Preserving HEP data is important for:

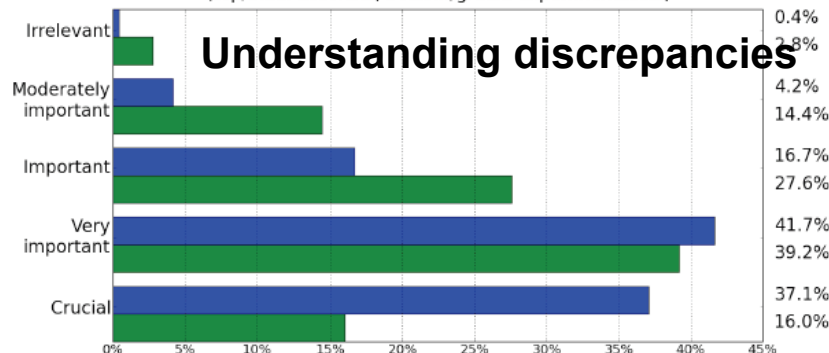
a) Compiling published results on a given subject (e.g. for a review)
(top/blue: theorists, bottom/green: experimentalists)



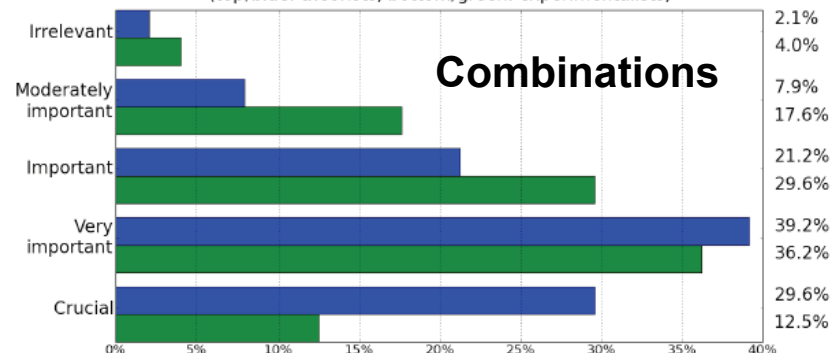
b) Testing new models using preserved data
(top/blue: theorists, bottom/green: experimentalists)



c) Showing compatibility of or detecting deviations between old and new experiments
(top/blue: theorists, bottom/green: experimentalists)



d) Combining preserved data with new measurements
(top/blue: theorists, bottom/green: experimentalists)



Rescued data used for fundamental results

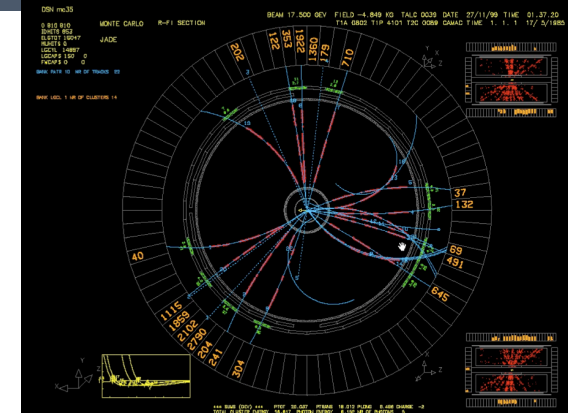
➤ Experience JADE

- Données sauvées par hasard

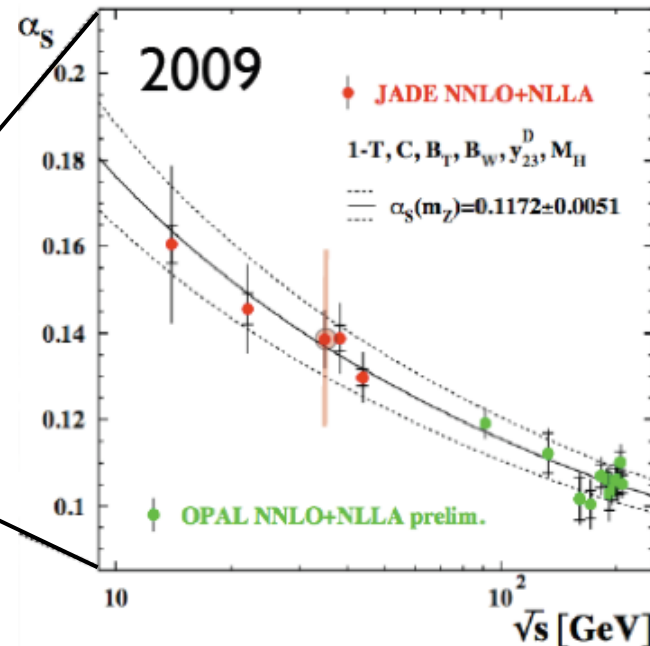
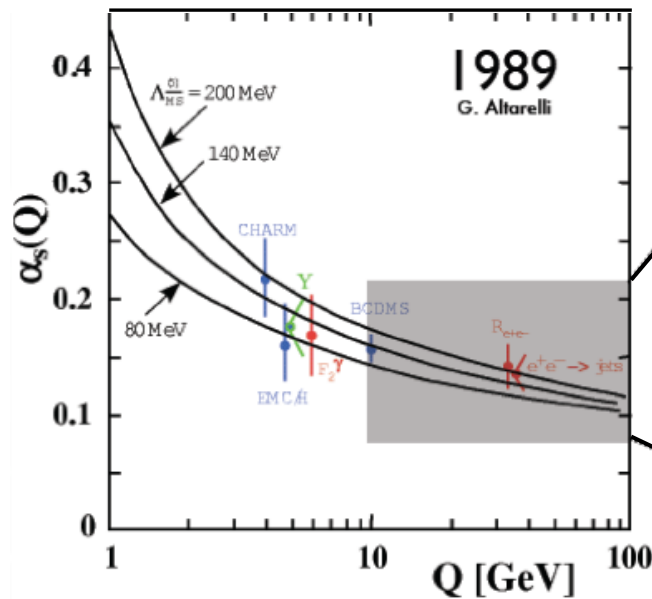
Nom de code: "la valise"

- Ré-analyse après 20 ans

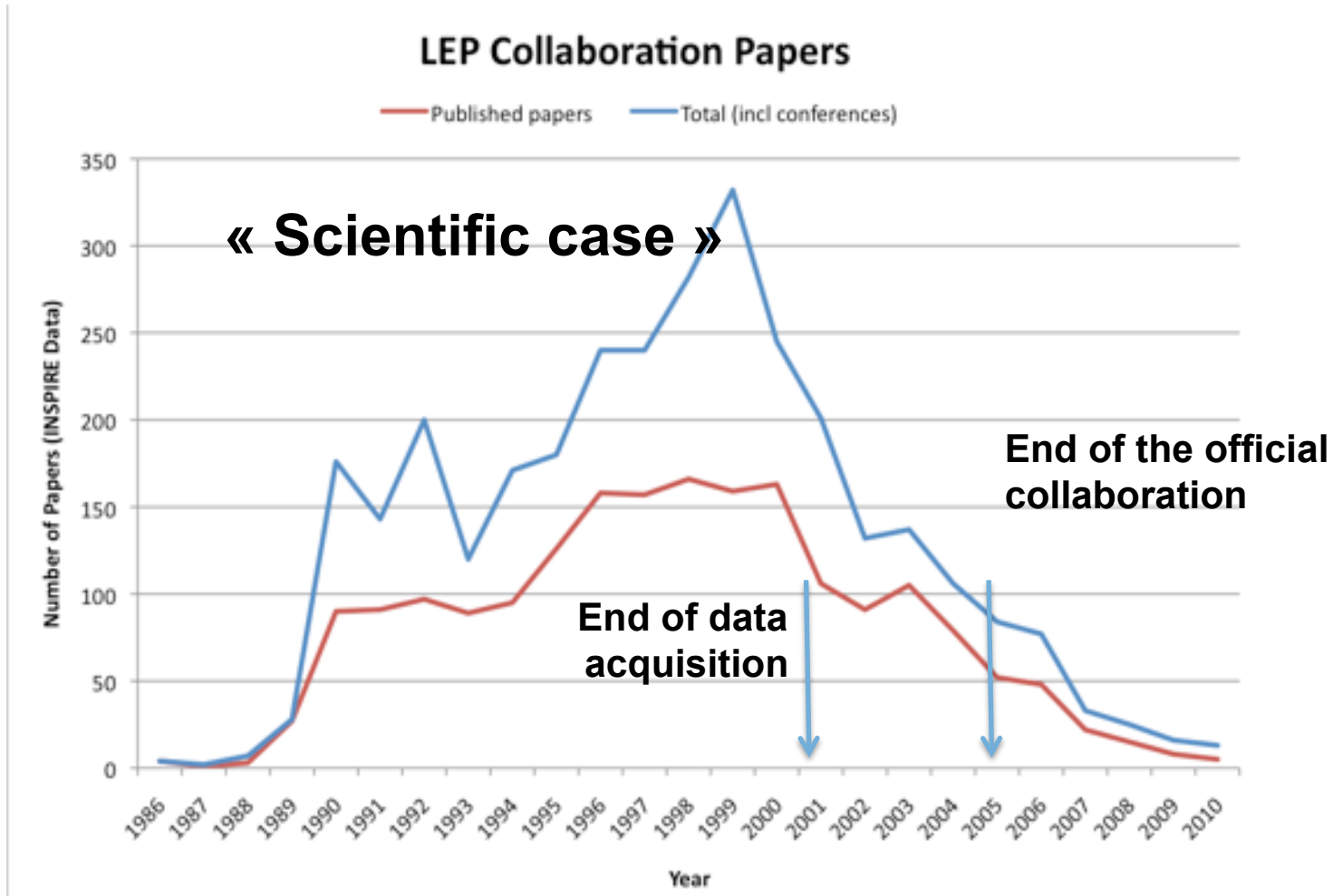
10 publications



2011

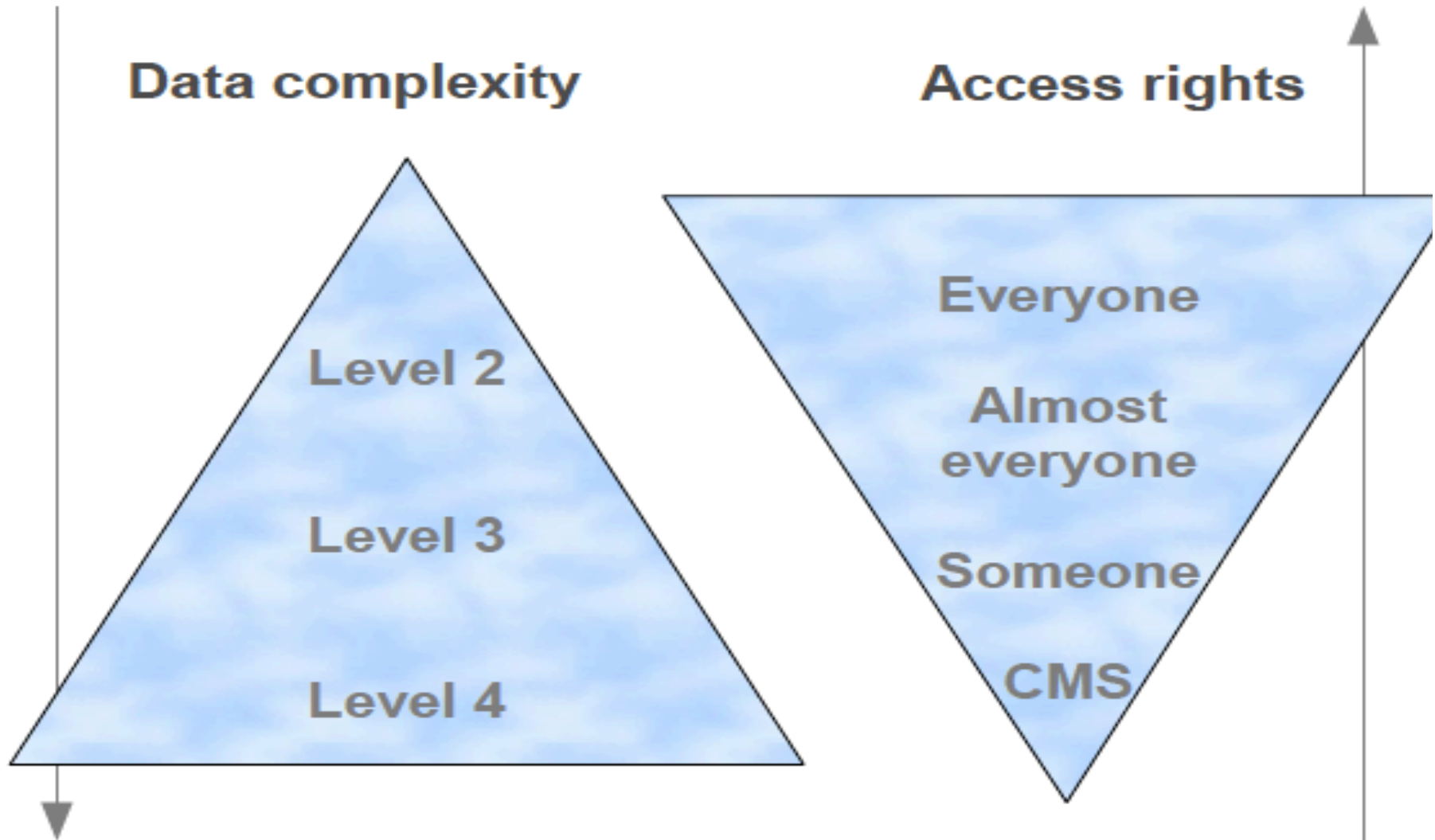


Long term publications yield



➤ Preserved data enables low cost science

Preservation complexity levels and access rights



Level 1: Documentation

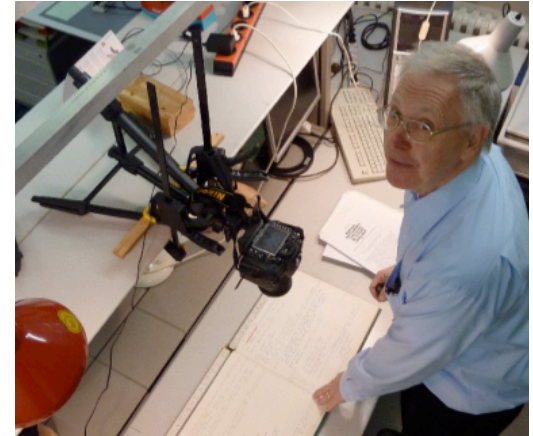
> Une tâche considérable: des groupes de travail dédiés

> **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

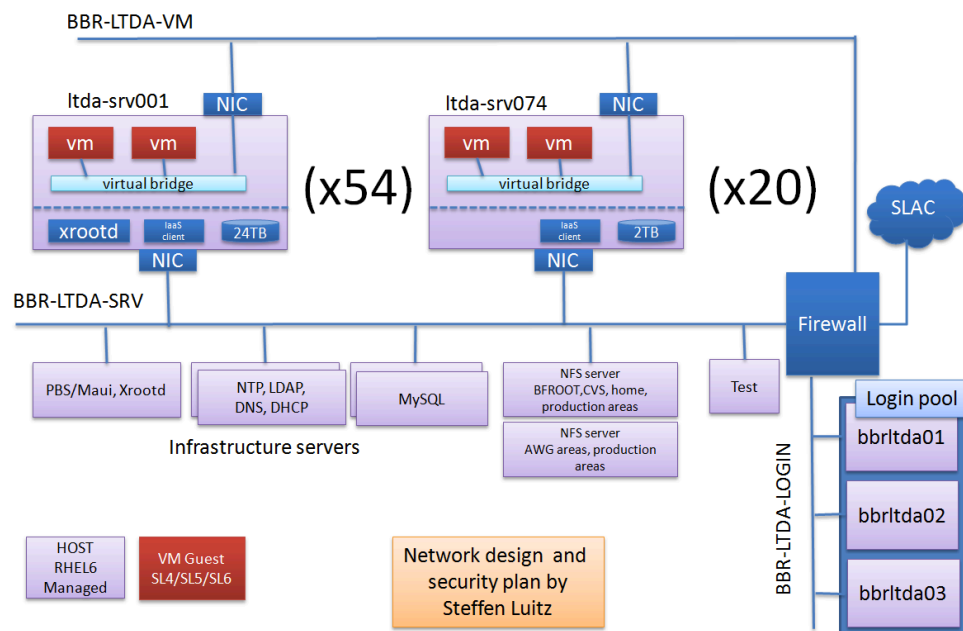
The screenshot shows the INSPIRE website interface. At the top, there's a navigation bar with links: HEP, INST, HELP, SPIRES, HEPNAMES. Below this, a banner for 'ZEUS Internal Notes' is visible. The main content area features a login form with fields for 'Username' (containing 'zeus') and 'Password'. There's 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 footer contains links for 'HEP', 'Search', 'Help', and 'Powered by Invenio v1.0.0-rc0+'. A message at the top right of the page says: 'Welcome to INSPIRE! INSPIRE is out of beta and ready to replace SP. please email us at feedback@inspirehep.net.'

ZEUS Internal Notes 10 records found

- 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](#)

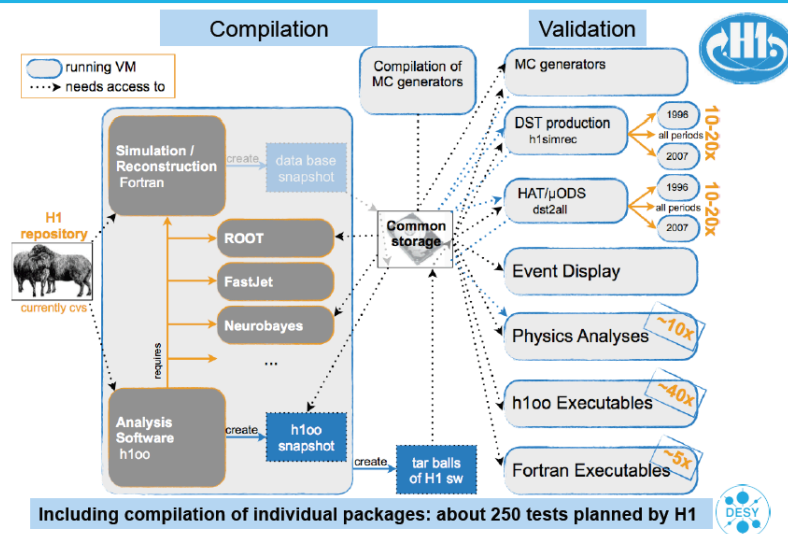
- The ingestion of other documents is under discussion, including theses, preliminary results, conference talks and proceedings, paper drafts, ...
- More on InSpire: reduced data?

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)

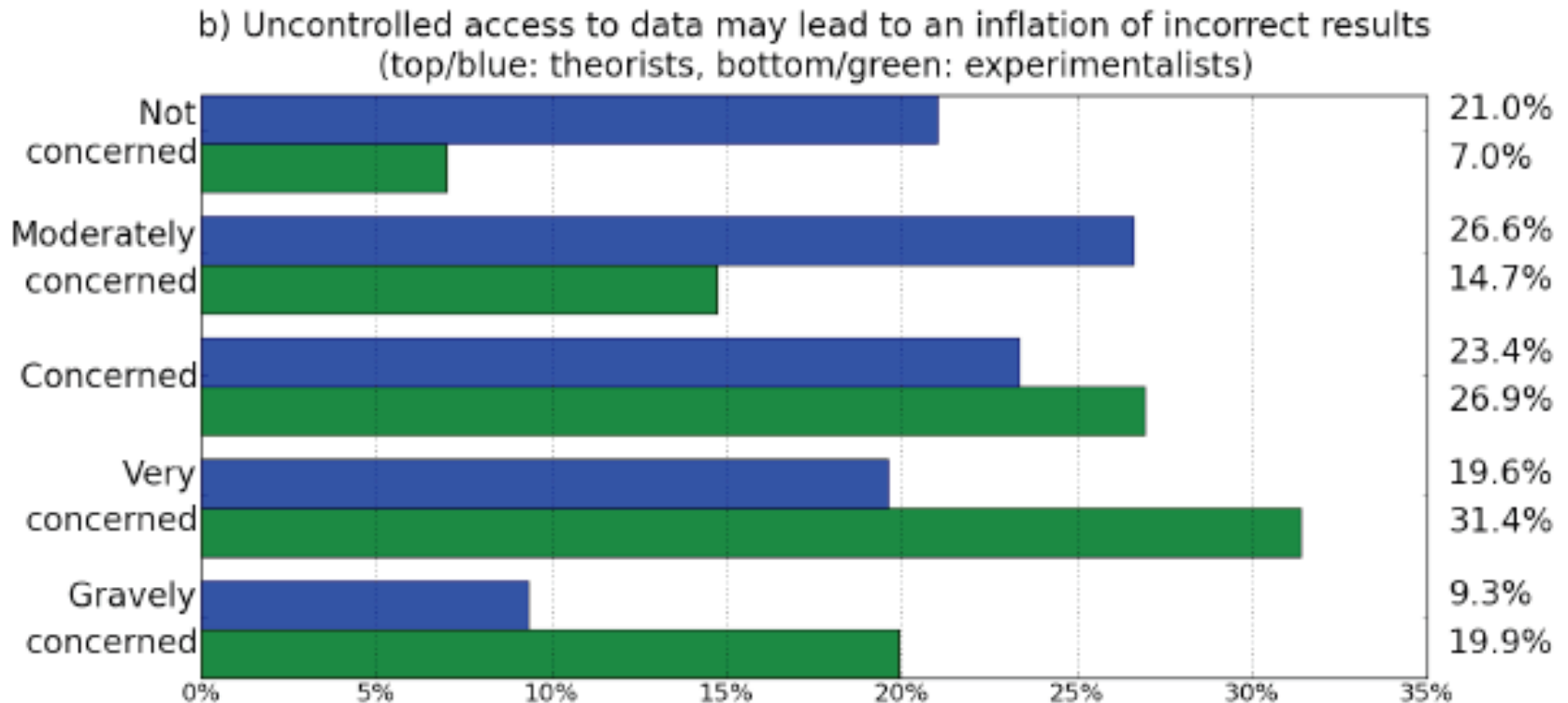
Example structure of experimental tests: H1 (Level 4)



Are complex preserved data « dangerous »?

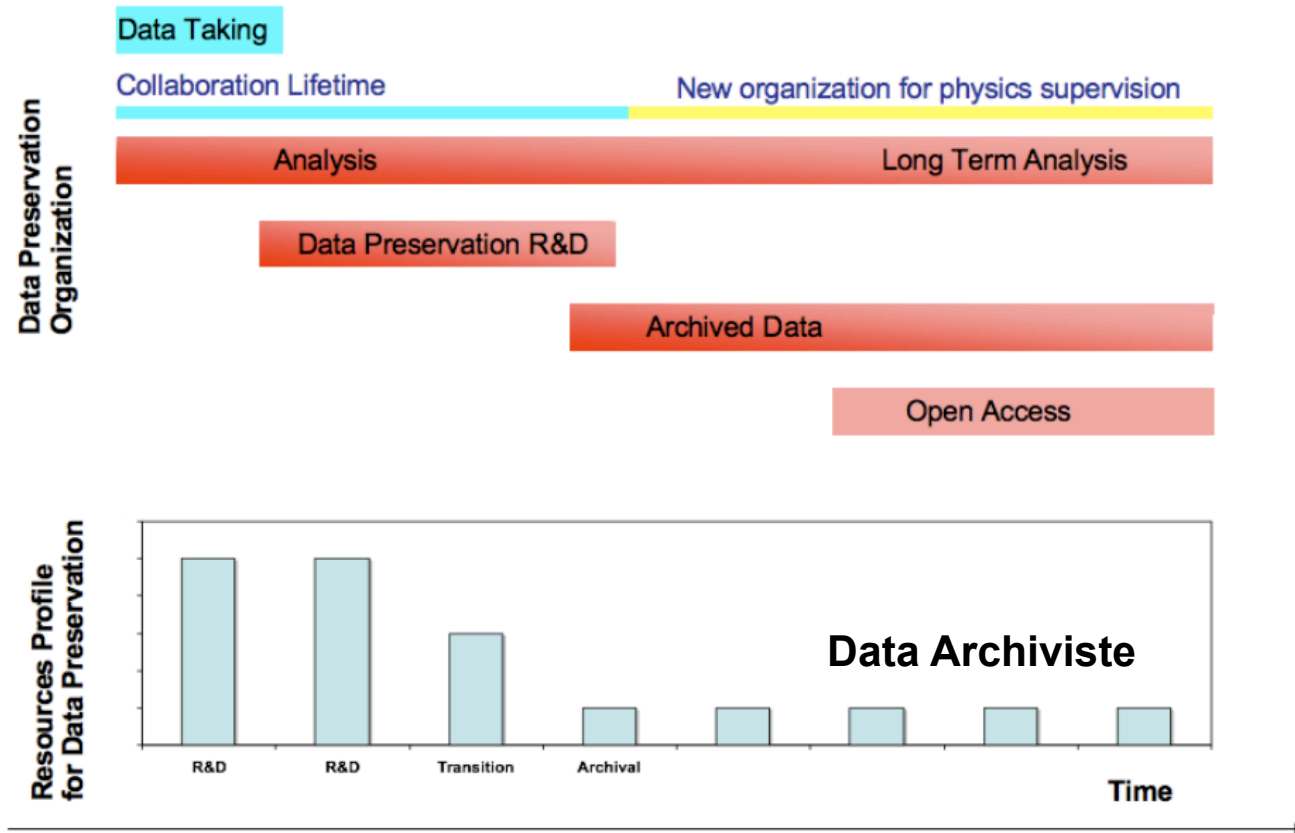
**"Errors using inadequate data are much less than those using no data at all."
Charles Babbage**

Parse.insight



Governance issues are very important to support data usage

Long term organisation and economical models

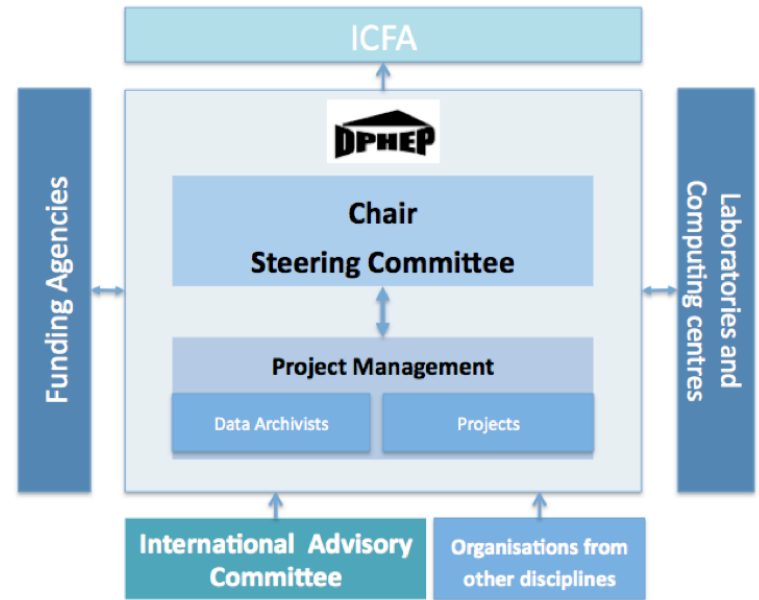


**The specific costs around 1% of the project
Scientific outcome around 10% more papers**

DPHEP: International organisation

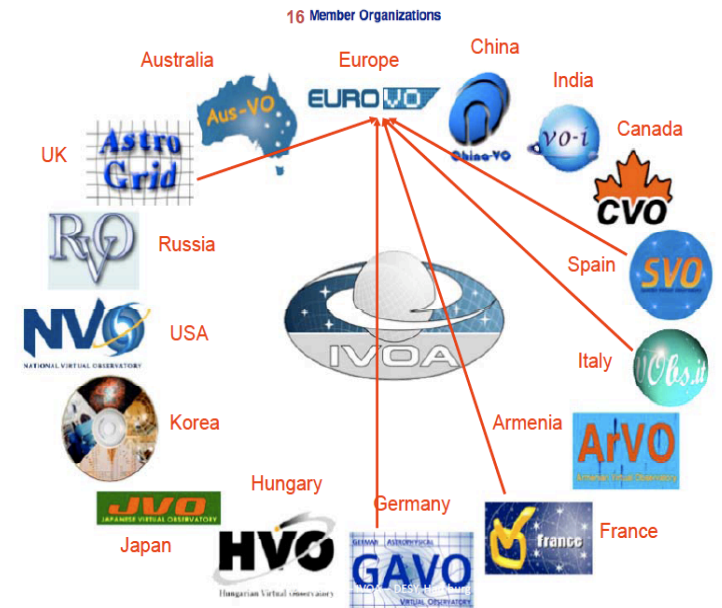
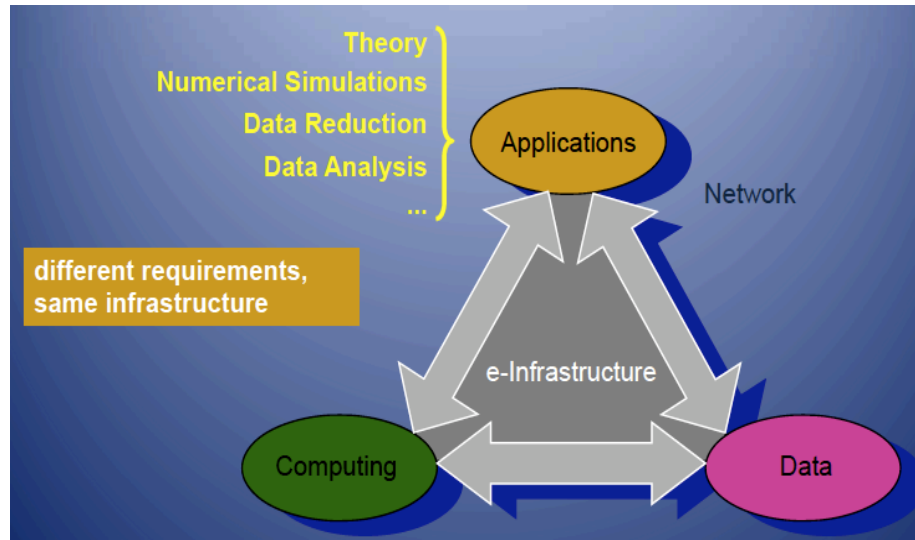


Study Group for Data Preservation and
Long Term Analysis in High Energy Physics



- Study Group DPHEP:
 - Large laboratories CERN, DESY, FERMILAB, SLAC, KEK, IHEP and experiments
- Organisation internationale en cours de mise en place
 - 100 contact personnes de contact
 - Chair: D. Diaconu Project Manager: Jamie Shiers (CERN)

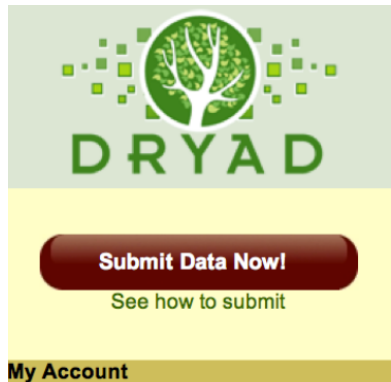
Virtual Observatories in Astrophysics



- > 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 science



Blue Ribbon Task Force
on Sustainable Digital Preservation and Access

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



PANGAEA®

Data Publisher for Earth & Environmental Science



[All](#) [Water](#) [Sediment](#) [Ice](#) [Atmosphere](#) [Search](#)
[Help](#) [Advanced Search](#) [Preferences](#) [more...](#)



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

The FITS Support Office

at NASA/GSFC



Scientific Data preservation in a multidisciplinary approach

> Challenges:

- **Scientific Potential:** these data sets contain unexploited information, which may give rise to highly useful for joint, multi-disciplinary project.
- **Complexity:** the data collected by experimental devices considered in the project are unique and encodes a large typology, well beyond the regular, well-structured data produced in large quantities in the industrial world.
- **Technological et methodological:** the installation of procedures, workflows, algorithms for long term data preservation, as well as the definition of suitable technological frameworks constitute novel investigation domains.

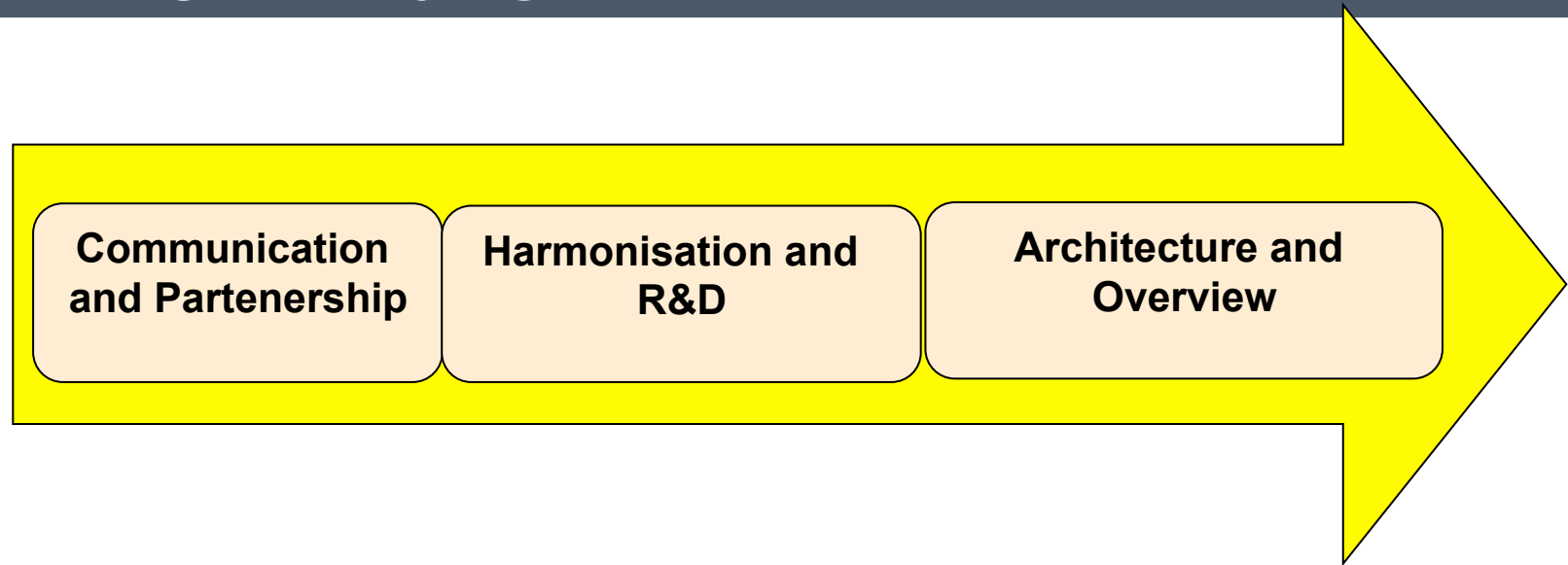
PREDON

A project for scientific data preservation in France



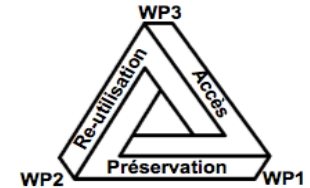
	Volume données	Complexité	Diversification des sources	Structuration au niveau international	Algorithmes et methodologies pour la preservation
IN2P3 HEP	+++	+++	+	++	+
INSU, IRD Astrophysics Earth Sciences	++	++	++	+++	++
CINES INS2I IT, Algorithms, workflows	+	++	+++	+	+++

PREDON: Plans



- > Short term (2013/2014): **Communication and partnerships**
 - Enlarge the community
- > Medium term (2014/2015) : **Harmonisation and R&D**
 - Communication: exchanges and workshops
 - Demonstrator access and preservation
- > Long term (2016) **Architecture and overview**
 - “Observatoire National des Données Scientifiques”

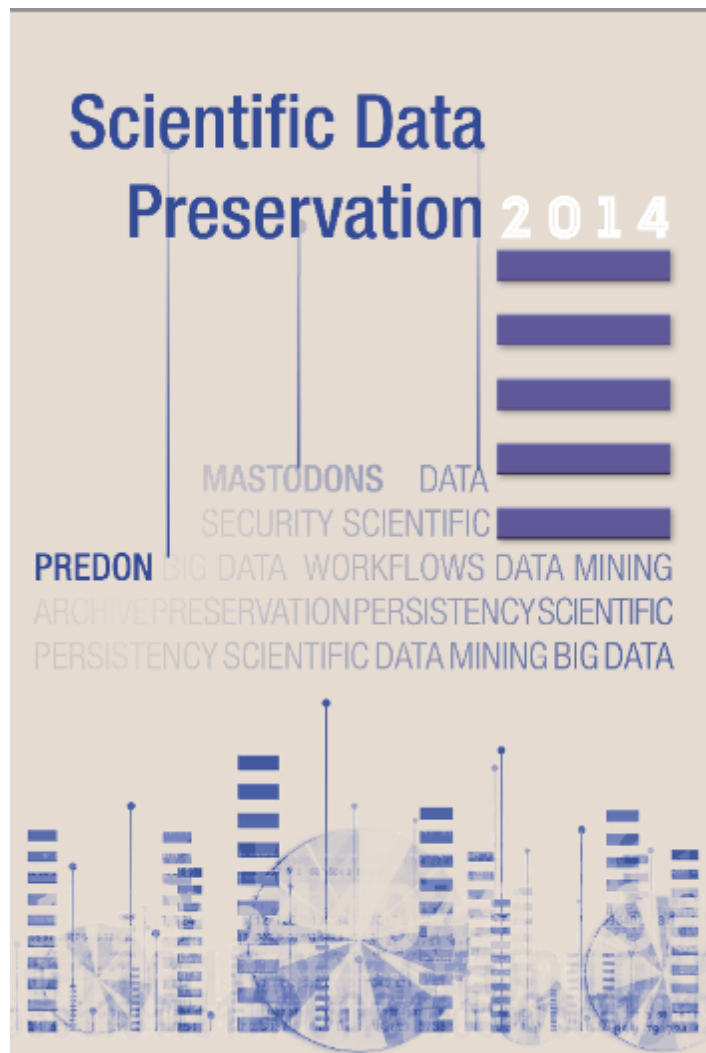
PREDON as a project



Working Package	Objectives	Participants (*coordinator)
WP1 Technologies and Methodologies	Explore methodologies and technologies suitable for a coherent and robust scientific data preservation in a multi-disciplinary context and on a multi-platform computing centre	CINES* APC
WP2 Algorithms and Workflows	Investigate generic and mathematically robust workflows and algorithms for data mining suited for data and workflow preservation; data- and process-based workflows and mining techniques to be used in a multi-disciplinary environment towards long term data preservation	LAM LIRMM LIPADE* LIPN
WP3 Data formats and interfaces	A parallel approach for data collection, storage, processing, analysis and preservation with the aim to achieve common standards for scientific data treatment	APC CPPM LAM* LPSC
WP4 General coordination	Program coordination, dissemination and international cooperation	CPPM*

Scientific data preservation: white paper

Available at <http://predon.org>



CHAPTER 1: SCIENTIFIC CASE

DATA PRESERVATION IN HIGH ENERGY PHYSICS	7
VIRTUAL OBSERVATORY IN ASTROPHYSICS	8
CRYSTALLOGRAPHY OPEN DATABASES AND PRESERVATION: A WORLD-WIDE INITIATIVE	15
SATELLITE DATA MANAGEMENT AND PRESERVATION	20
SEISMIC DATA PRESERVATION	26

CHAPTER 2: METHODOLOGIES

WORKFLOWS AND SCIENTIFIC BIG DATA PRESERVATION	37
LONG TERM ARCHIVING AND CCSDS STANDARDS	38
CLOUD AND GRID METHODOLOGIES FOR DATA MANAGEMENT AND PRESERVATION	42
SCIENTIFIC DATA PRESERVATION, COPYRIGHT AND <i>OPEN SCIENCE</i>	49

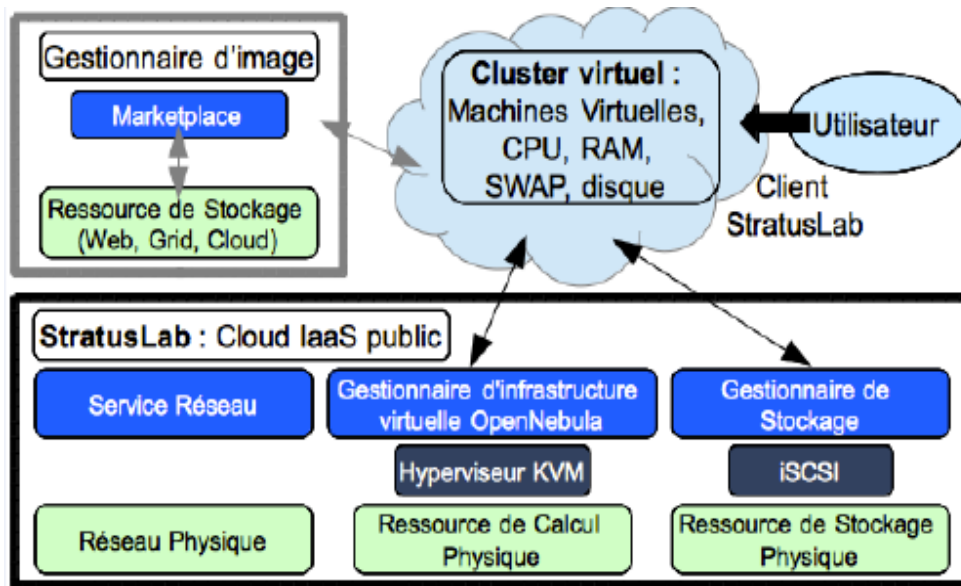
CHAPTER 3: TECHNOLOGIES

STORAGE TECHNOLOGY FOR DATA PRESERVATION	61
REQUIREMENTS AND SOLUTIONS FOR ARCHIVING SCIENTIFIC DATA AT CINES	62
VIRTUAL ENVIRONMENTS FOR DATA PRESERVATION	65

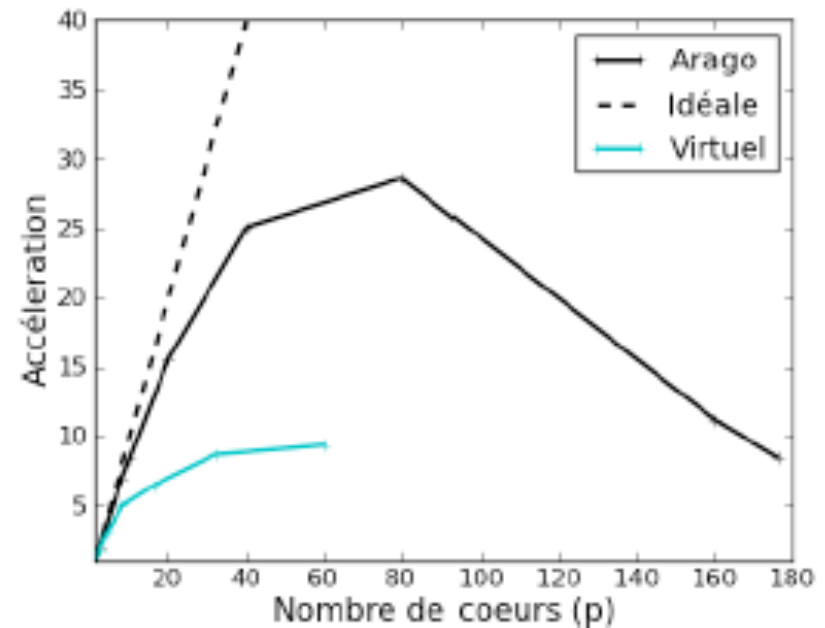
Exemple projet: Data processing & storage in the cloud

LabEx UnivEarths project at APC / François Arago Centre:

- potential of the cloud versus classical data processing and storage opportunities
- test processing on Francois Arago Centre cluster, compared with Cloud StratusLab
- questions: accessibility, data security, short-term and long-term cost



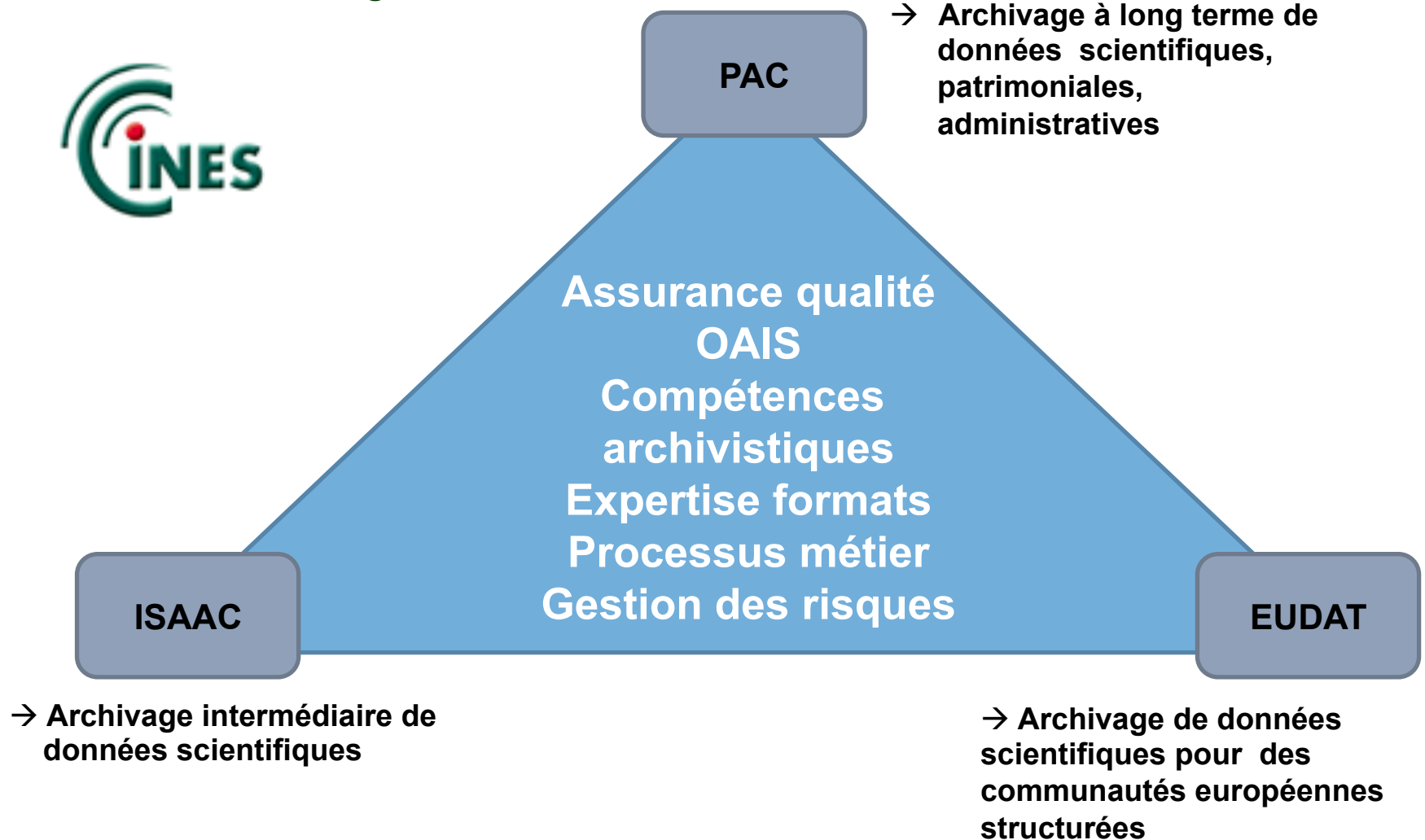
Schematic description of the cloud StratusLab, which is a European public cloud project IaaS which started in 2010.



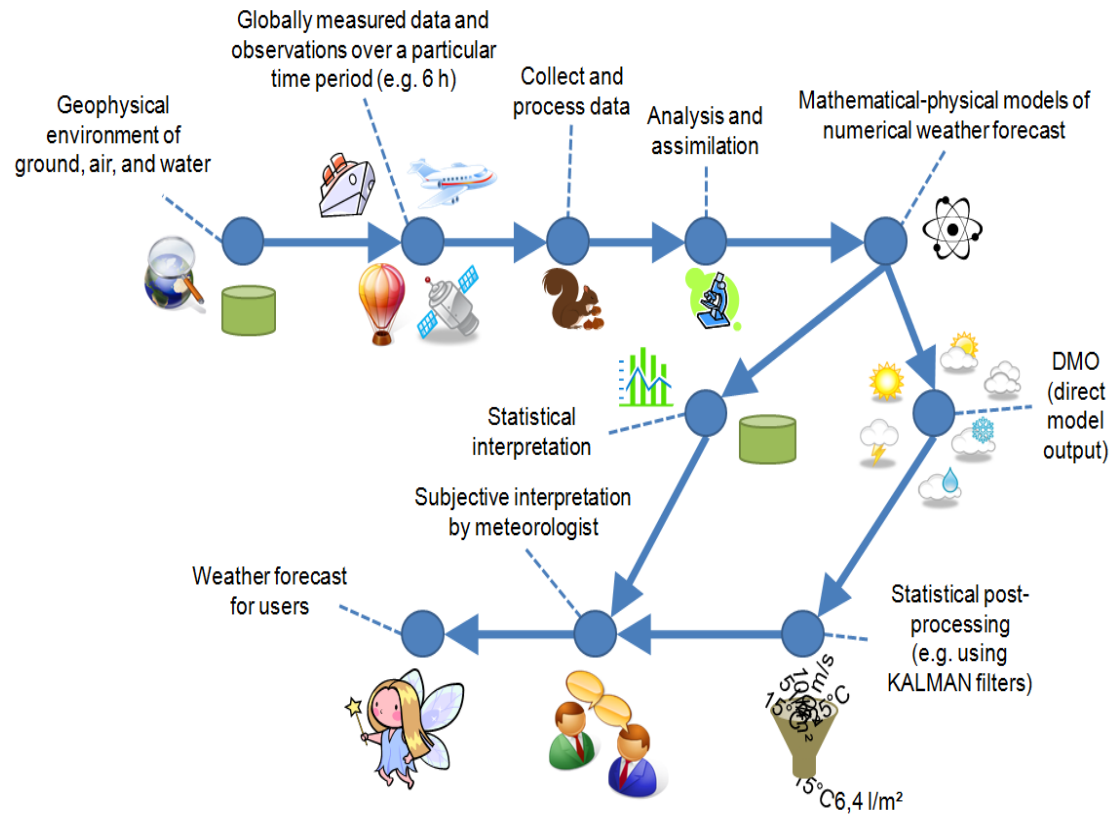
Processing speed does accelerate much faster on a classical computing cluster compared to cloud computing (Cavet et al. 2012)

Archival expertise CINES

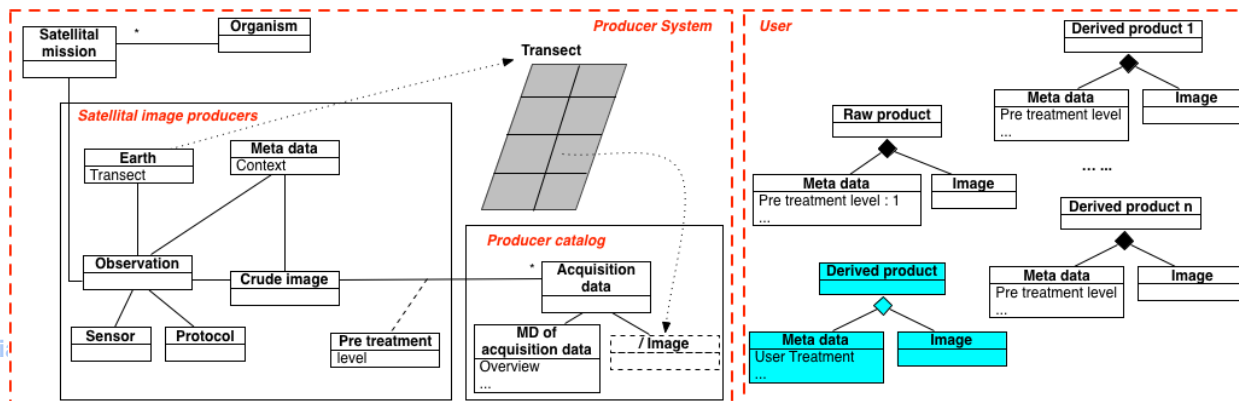
Les services d'archivage au CINES



Workflows for data preservation



Rigorous approach emerges



Long Term Archiving and CCSDS standards

Danièle Boucon, CNES

The primary objective of the Producer-Archive Interface Specification (PAIS) standard is to provide concrete XML files supporting the description and the control of transfers from a Producer to an Archive.

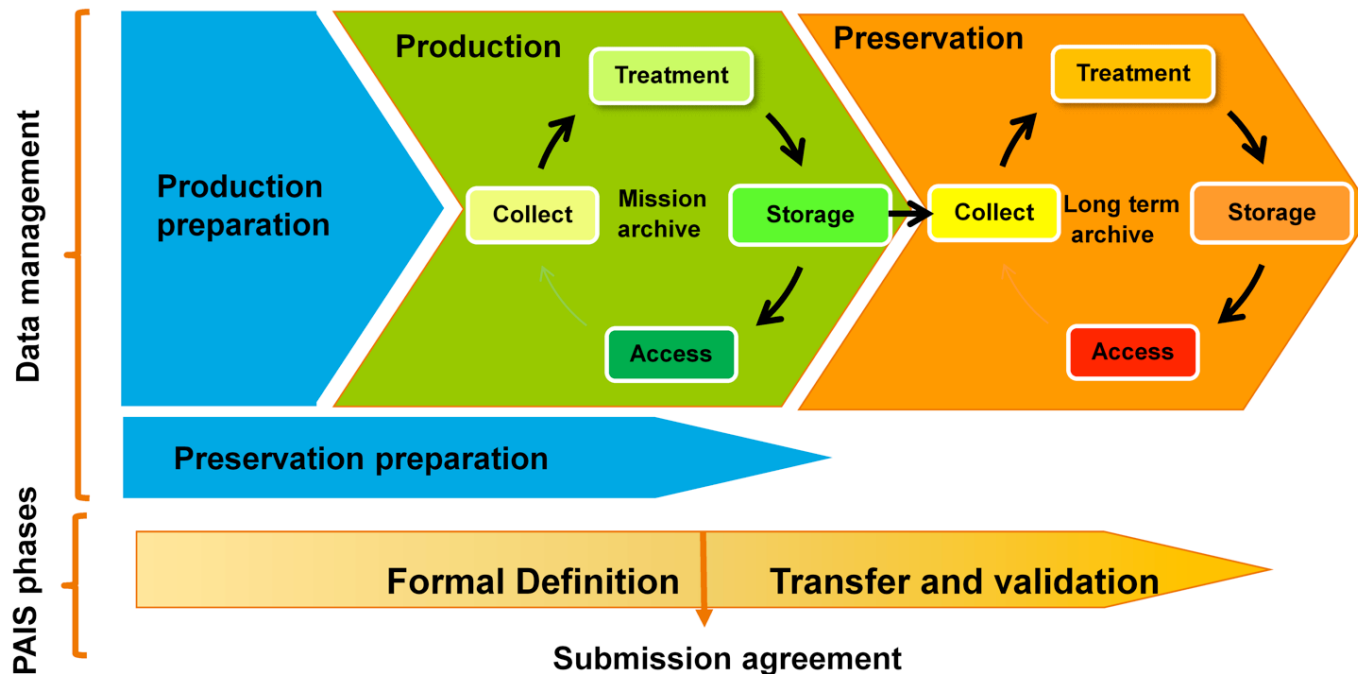
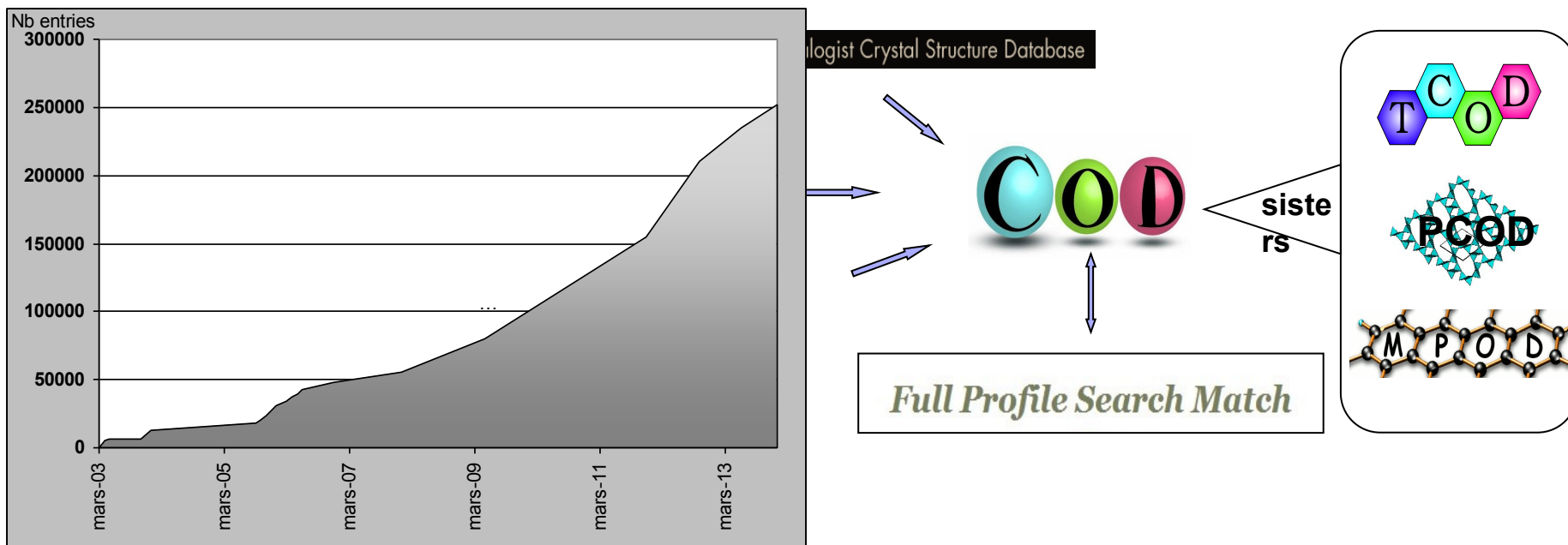


Figure 3: PAIS, preservation process and data lifecycle

Crystallography Open Databases and Preservation: a World-Wide Initiative

Daniel Chateigner (for the COD Advisory Board)



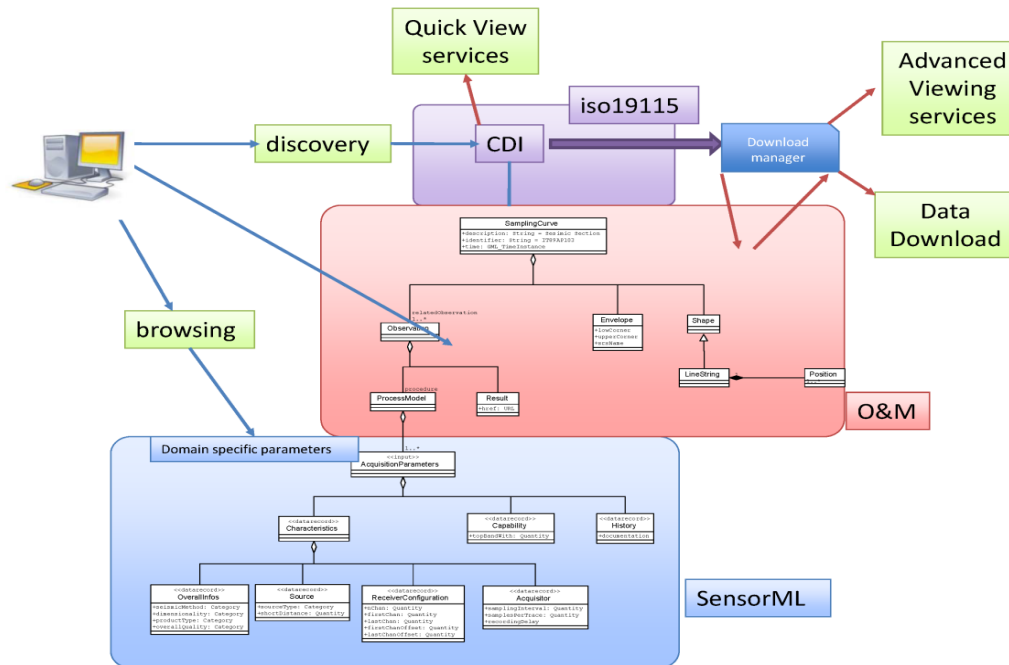
“...there is not yet sufficient coherence of experimental metadata standards or national policy to rely on instrumental facilities to act as permanent archives;
-there is not sufficient funding for existing crystallographic database organisations (which maintain curated archives of processed experimental data and derived structural data sets) to act as centralised stores of raw data, although they could effectively act as centralised metadata catalogues;
-few institutional data repositories yet have the expertise or resources to store the large quantities of data involved with the appropriate level of discoverability and linking to derived publications.”

Seismic Data Preservation

**Marc SCHAMING, Institut de Physique du Globe (CNRS/UNISTRA),
Strasbourg**

Conclusion

Preservation of seismic data is essential, but usually not considered by scientists, because it takes resources to document metadata, to read and copy tapes, to convert formats, etc. These tasks should be addressed at national and/or European level. Some European projects (Seiscan/Seiscanex, Geo-Seas) demonstrated that it is possible and useful. Repositories at national level should pursue this task with geophysical skills.



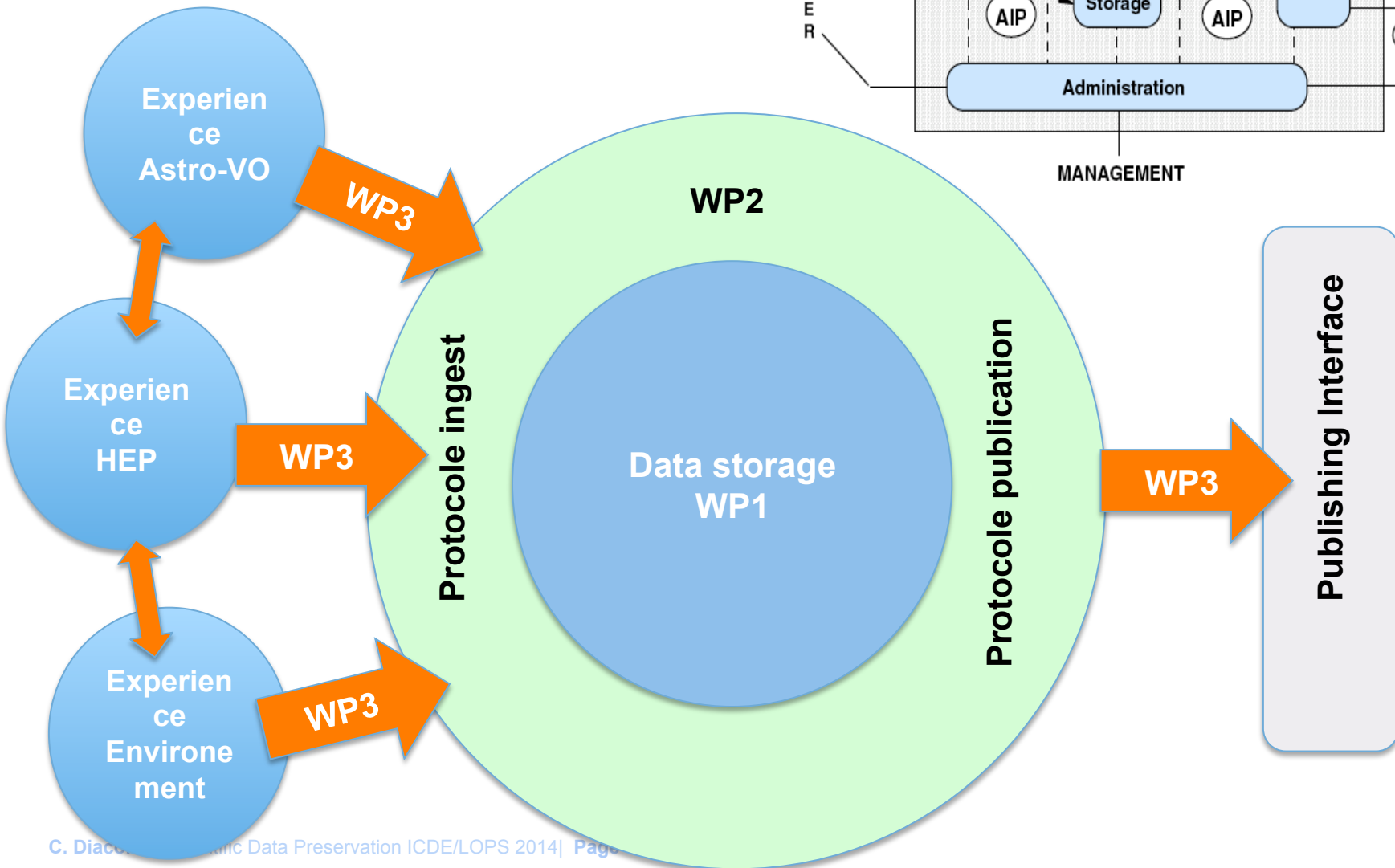
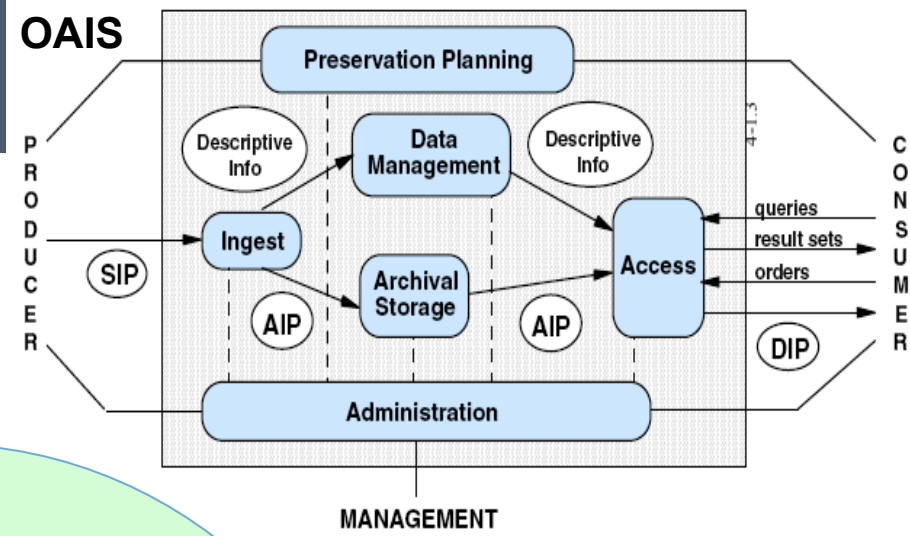
Scientific Data Preservation, Copyright and *Open Science*

Philippe Mouron, Aix-Marseille University, Faculté de droit et de science politique

- > **The best guarantee for ensuring the integrity of a resource is based on property.**
- > However, isn't there a public ownership of scientific research?
 - In truth, even if the public authorities may fundamentally participate in the scientific research, this does not mean, *ipso facto*, that they own its results.
- > ...any paper, article, report, record, thesis, book, graphic, map,... conducting personal choices of a researcher, or expressing his own personality, will be considered as a work of mind [...] are copyrightable
- > **The goal of digital preservation of scientific data must therefore be reconciled with intellectual property rights.**
- > Open model of management of intellectual property rights.
 - Tools: open access licenses (e.g. Creative Commons)

Towards a demonstrator

OAIS

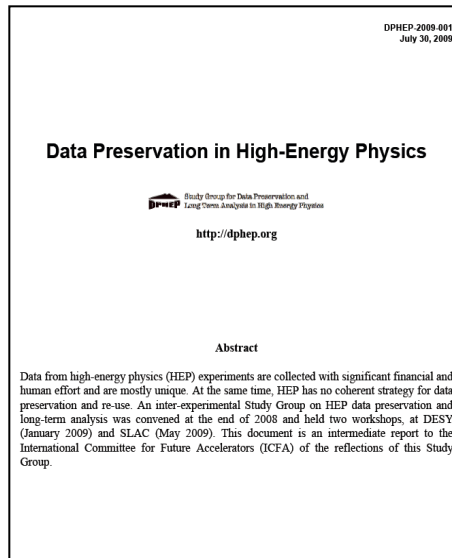


Conclusions

- > Research data in the recent period is a « gold mine » with unexploited research potential and therefore should be preserved long term
- > The preservation is worth the effort and enables low cost science
- > International cooperation is essential
- > Frontier technologies and research are involved
- > Innovative solutions can emerge from multidisciplinary initiatives

DPHEP Intermediate Recommendations (end 2009)

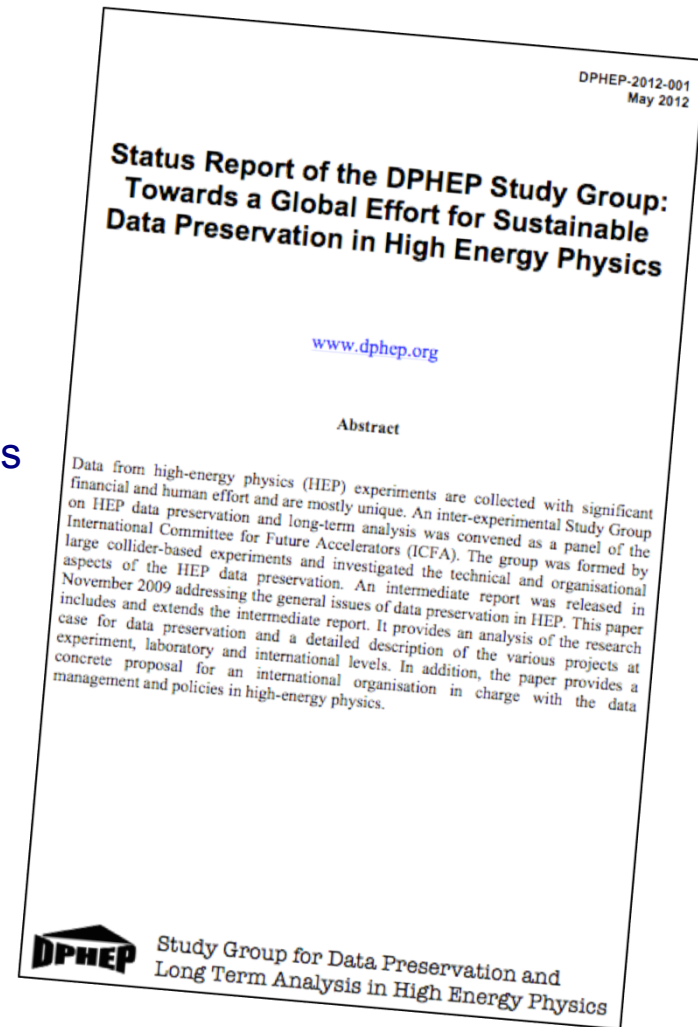
> [arXiv:0912.0255](https://arxiv.org/abs/0912.0255)



- > 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:1205.4667



A word on access and data preservation

Example: NSF Policy

Investigators are expected to share with other researchers, at no more than incremental cost and within a reasonable time, the primary data, samples, physical collections and other supporting materials created or gathered in the course of work under NSF grants. Grantees are expected to encourage and facilitate such sharing.

Proposals [...] must include a supplementary [...] "Data Management Plan" (DMP) [...] describ[ing] how the proposal will conform to NSF policy on the dissemination and sharing of research results.

<http://www.nsf.gov/bfa/dias/policy/dmp.jsp>

Very similar policies in other funding agencies (and growing interest for these aspects in the context of “big data” strategies)



Riding the wave

How Europe can gain from the rising tide of scientific data

Final report of the High Level Expert Group on Scientific Data
A submission to the European Commission
October 2010

**A myriad of projects/coalitions on
“data infrastructures”
either funded or
in preparation for FP8**

-APA, EUDAT, DPM, RDA...

Scientific e-infrastructure – a wish list

The ideal data infrastructure for science will have a long list of technical characteristics. Here are some suggestions.

- Open deposit, allowing user-community centres to store data easily
- Bit-stream preservation, ensuring that data authenticity will be guaranteed for a specified number of years
- Format and content migration, executing CPU-intensive transformations on large data sets at the command of the communities
- Persistent identification, allowing data centres to register a huge amount of markers to track the origins and characteristics of the information
- Metadata support to allow effective management, use and understanding
- Maintaining proper access rights as the basis of all trust
- A variety of access and curation services that will vary between scientific disciplines and over time
- Execution services that allow a large group of researchers to operate on the stored data
- High reliability, so researchers can count on its availability
- Regular quality assessment to ensure adherence to all agreements
- Distributed and collaborative authentication, authorisation and accounting
- A high degree of interoperability at format and semantic level

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.**