

La préservation des données scientifiques: une mine d'or pour la science de demain

Workshop PREDON
APC 5-6 Novembre 2014



C. Diaconu

Centre de Physique des Particules de Marseille
CNRS et Aix Marseille Université (AMU)

Les données digitales sont fragiles

- La capacité de stockage est physiquement dépassée depuis longtemps

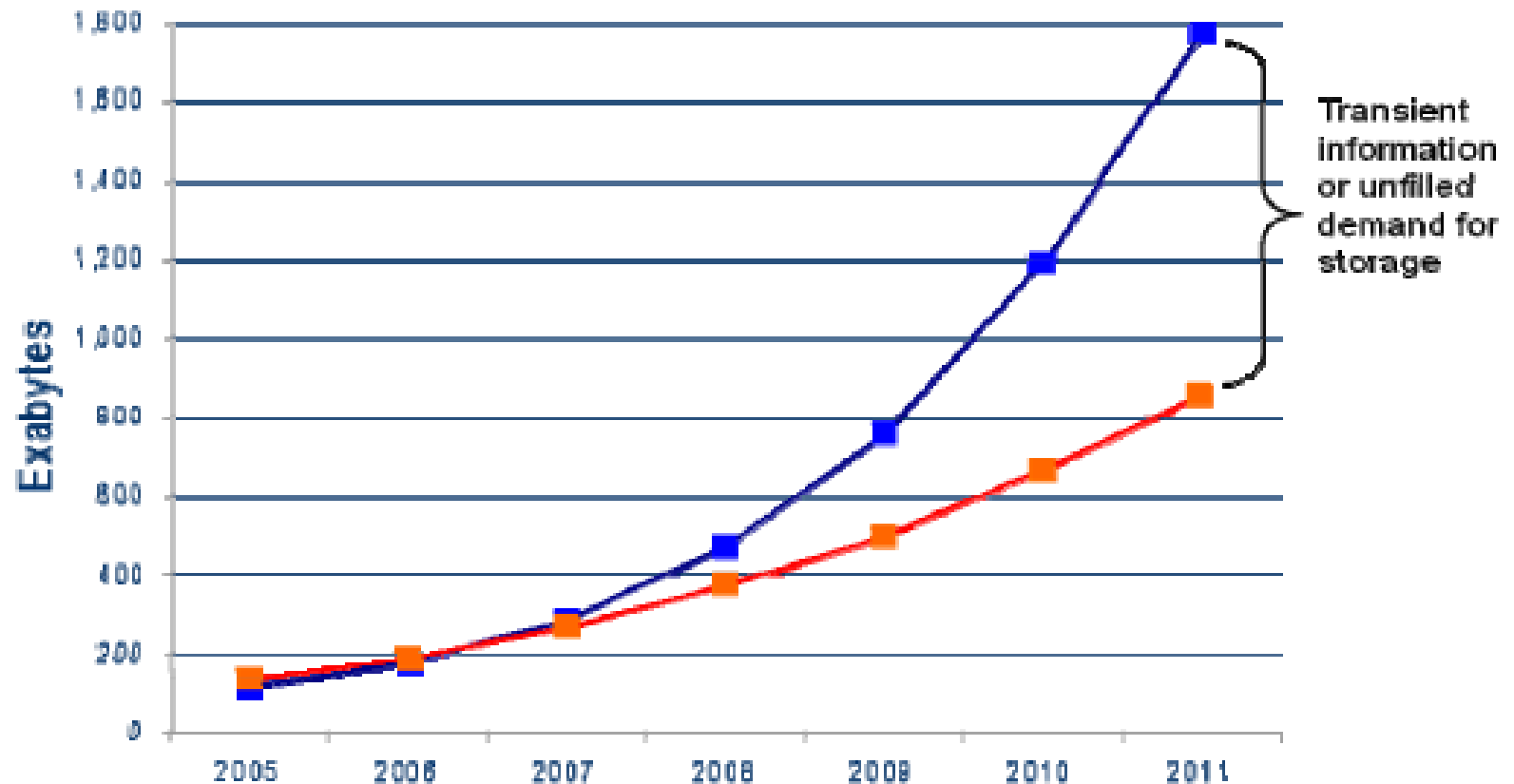
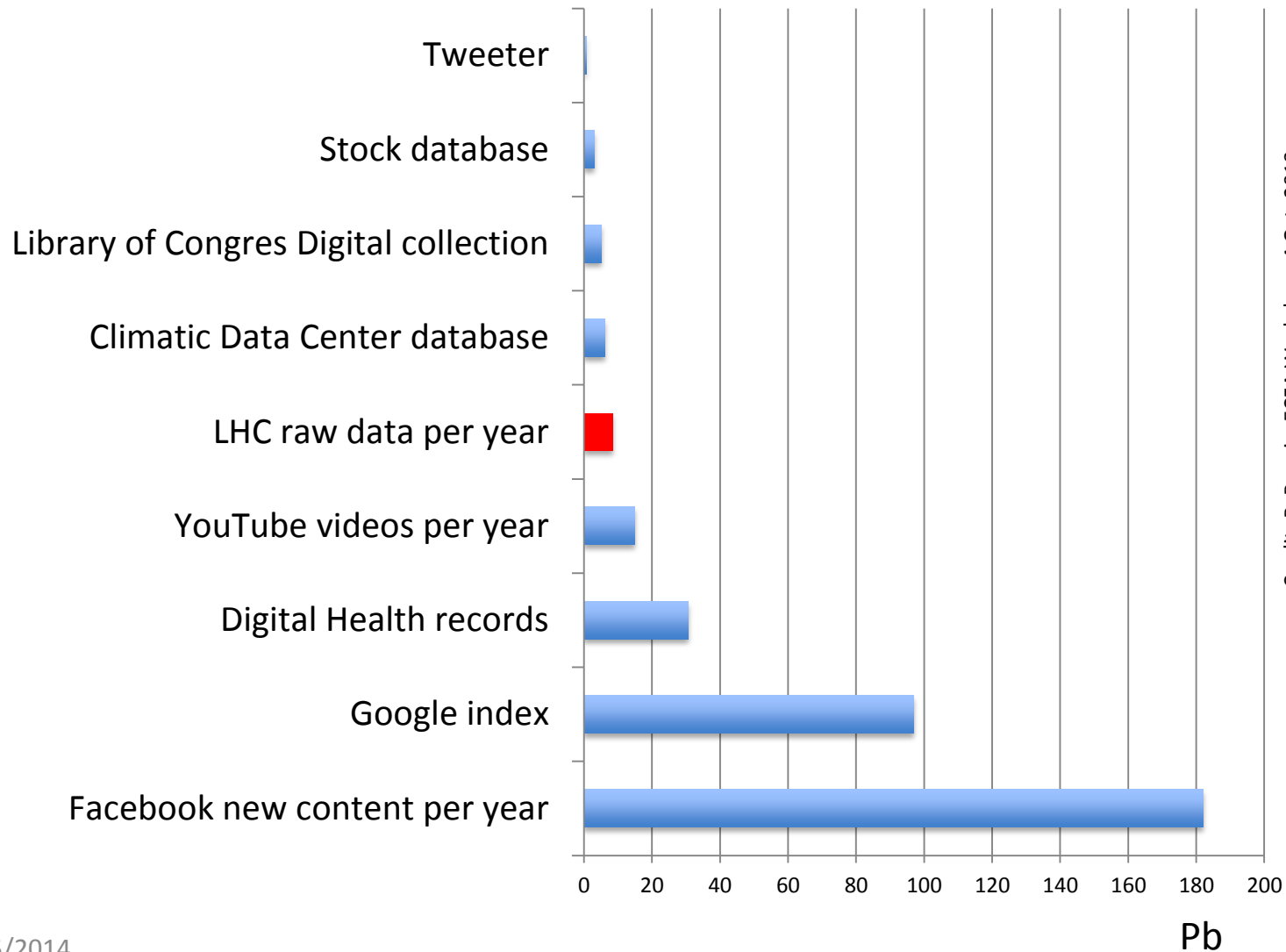


FIGURE 1.3: Information and Storage

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

Données digitales explosent (les données scientifiques aussi)

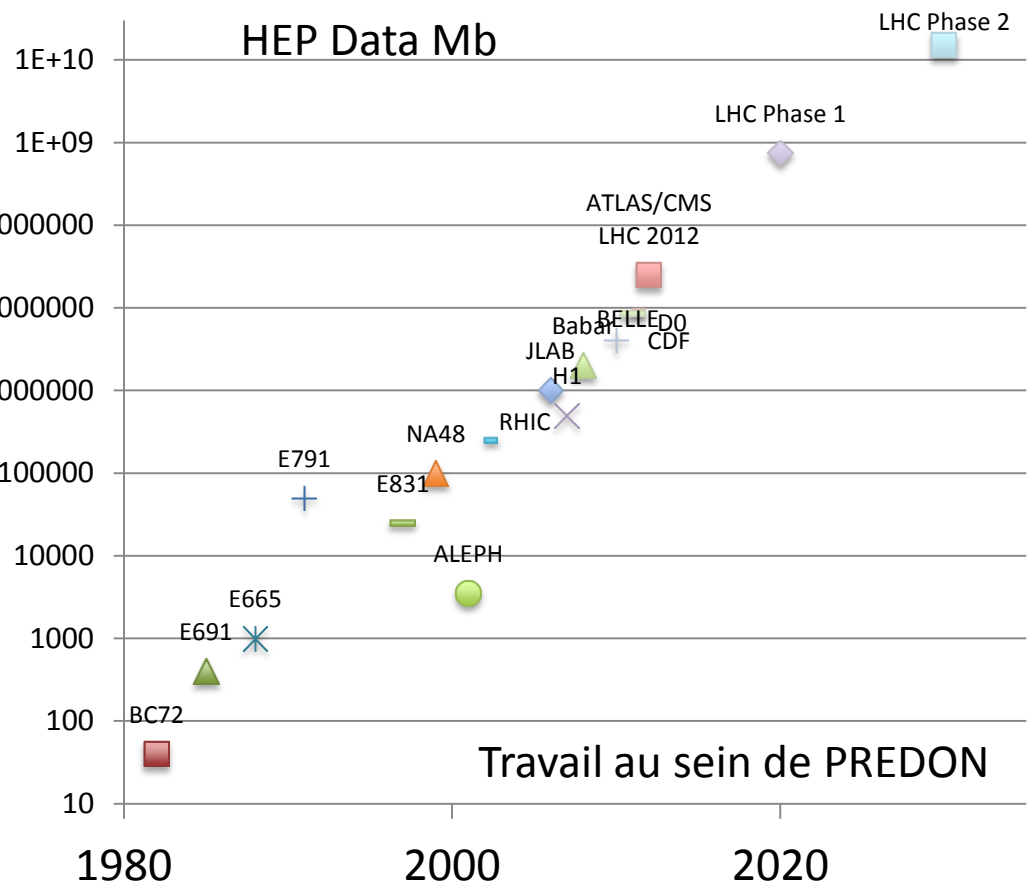
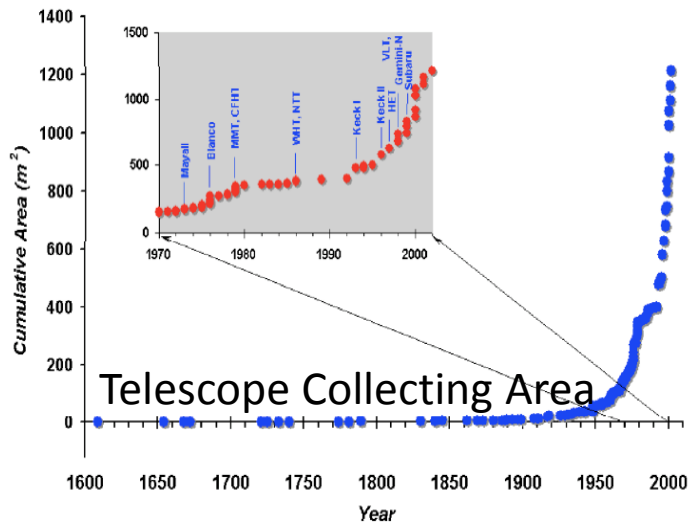
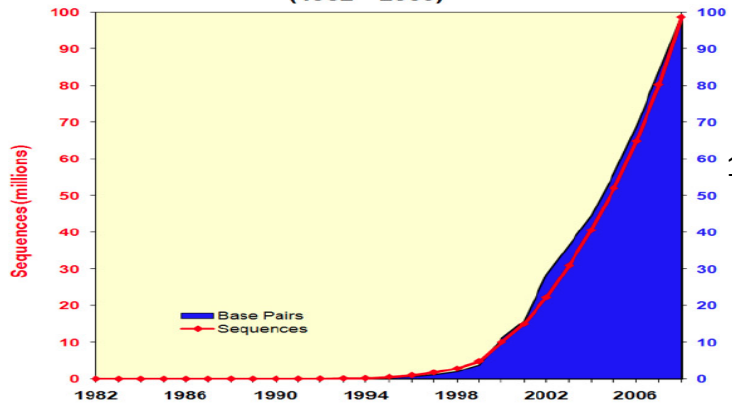


« Big Scientific Data »

- La recherche est « digitale »
 - Augmentation dramatique de la quantité/complexité des données



Growth of GenBank (1982 - 2008)

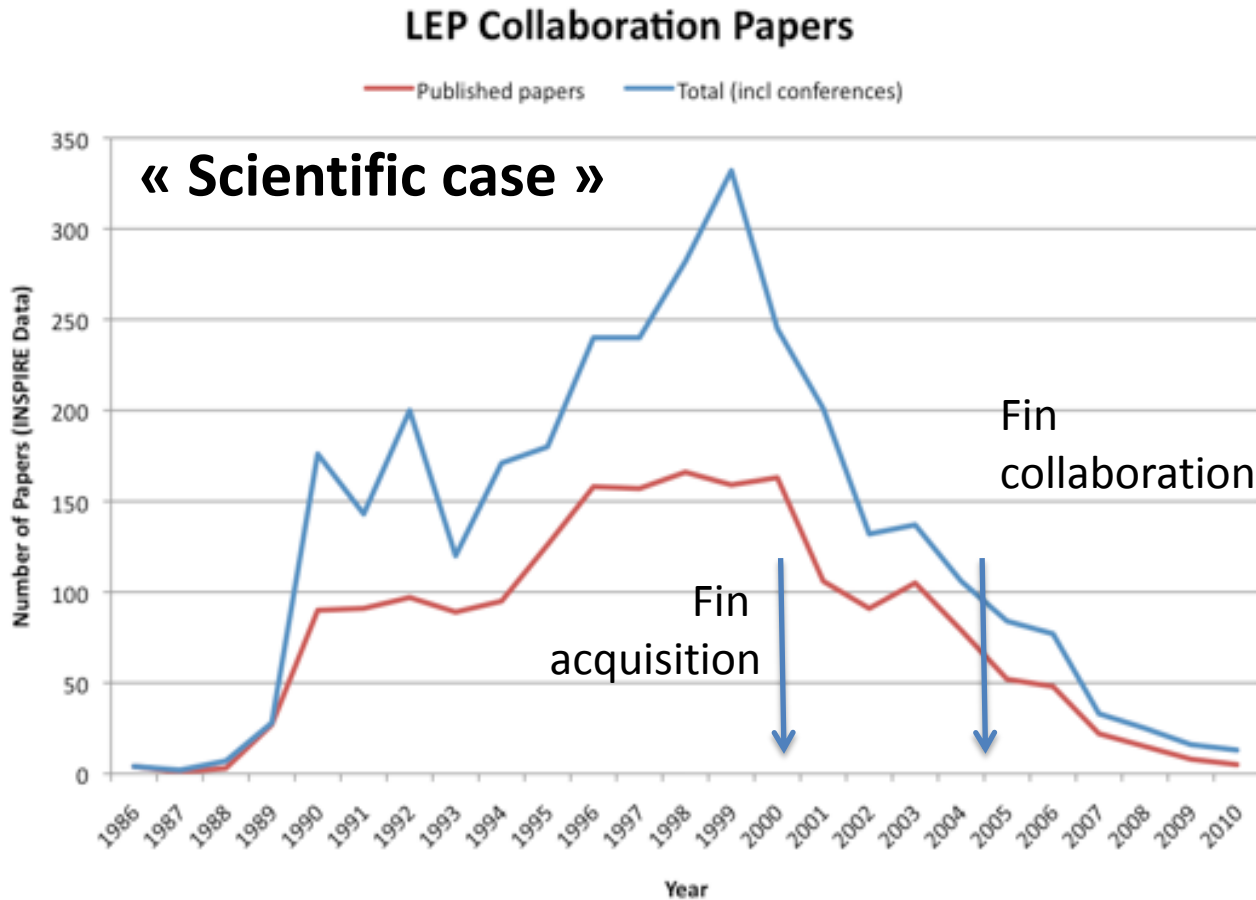


Travail au sein de PREDON

Est-ce que les données scientifiques sont spéciales (« big » à part)?

- Riches en information
 - structurées suivant un plan de recherche et une démarche scientifique
- De plus en plus diverses
 - la plupart des disciplines produisent massivement des données
- Souvent produites avec des efforts financiers et humains significatifs (voir gigantesques)
 - Plus ça coute cher, moins c'est reproductible
- Englobent des connaissances uniques
 - « Time stamped »
- De plus en plus dans une logique « observatoire »:
 - Les données contiennent plus que ce qu'on voulait au départ
 - Seulement l'information décantée est publiée de suite (1/10)
- **PRESERVATION!**

Est-ce que ça vaut le coup de garder des données « anciennes »?



nature International weekly journal of science

Home | News & Comment | Research | Careers & Jobs | Current Issues

Archive > Volume 503 > Issue 7477 > News > Article

NATURE | NEWS

عربي

LHC plans for open data future

Researchers share results to keep them accessible.

Elizabeth Gibney

26 November 2013

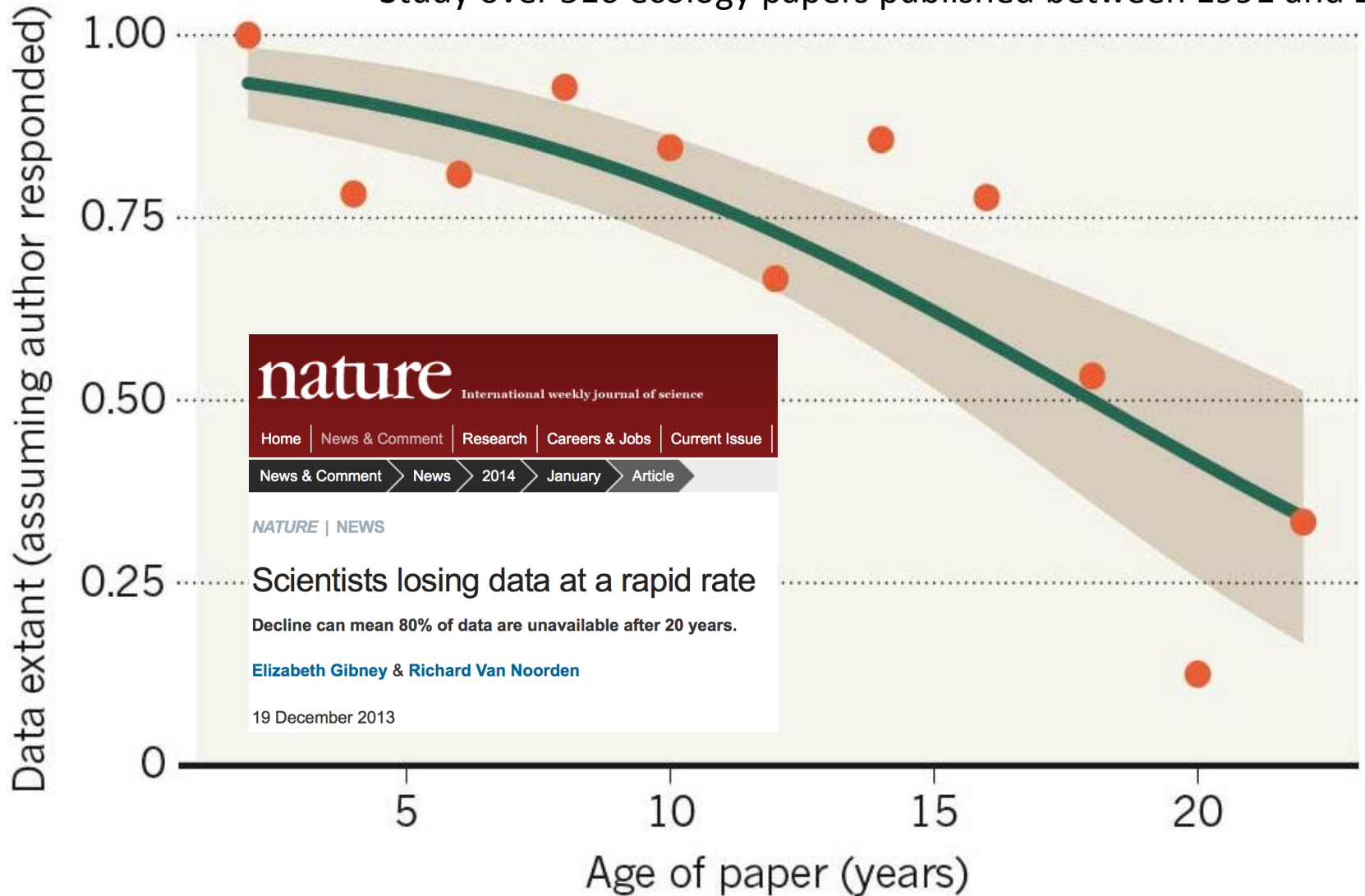
“When the LHC programme comes to an end, it will probably be the last data at this frontier for many years. We can’t afford to lose it.”

Estimation: gain scientifique de 10% pour un cout bien inférieur à 1%

MISSING DATA

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

Study over 516 ecology papers published between 1991 and 2011.



nature International weekly journal of science

Home | News & Comment | Research | Careers & Jobs | Current Issue

News & Comment > News > 2014 > January > Article

NATURE | NEWS

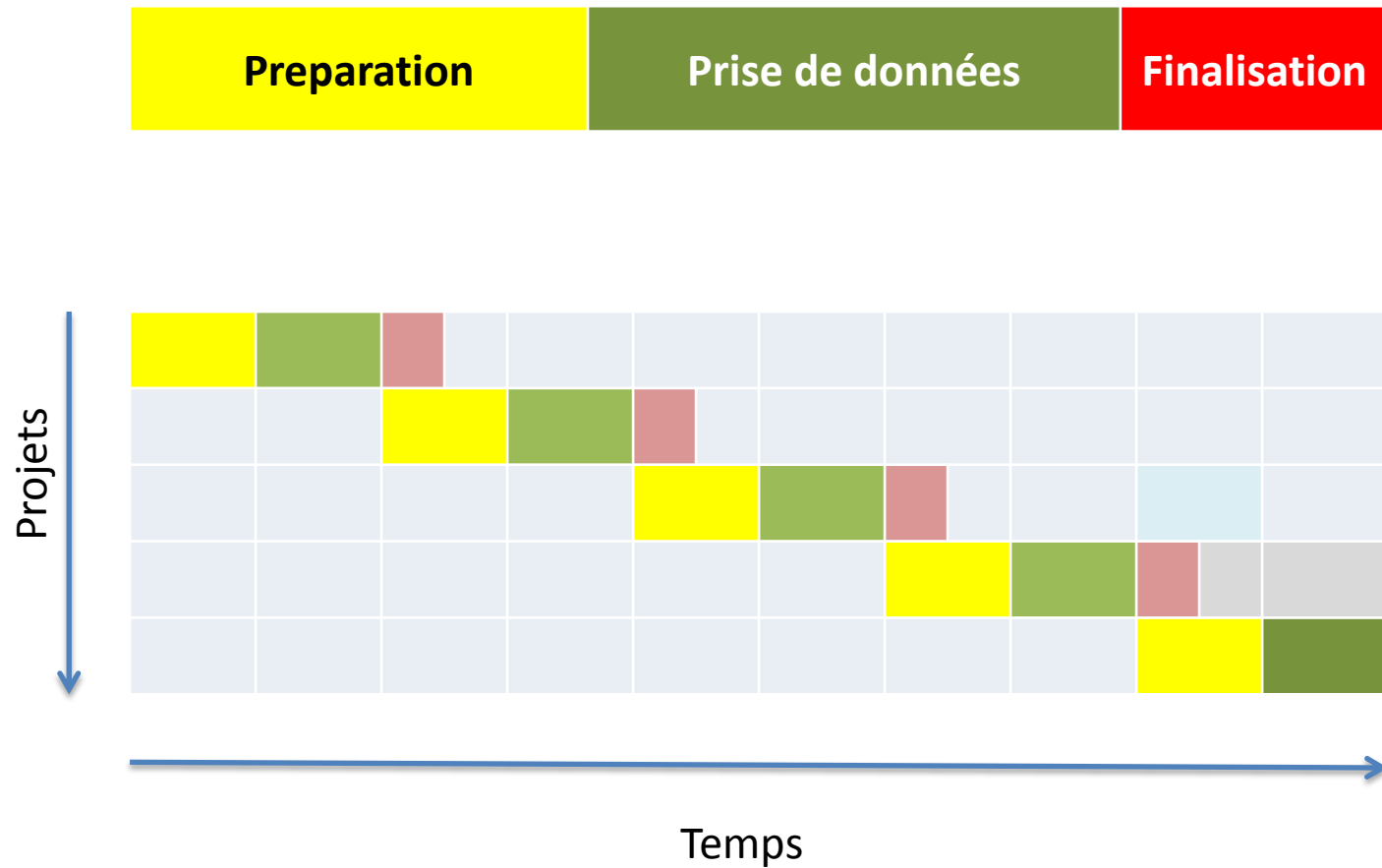
Scientists losing data at a rapid rate

Decline can mean 80% of data are unavailable after 20 years.

Elizabeth Gibney & Richard Van Noorden

19 December 2013

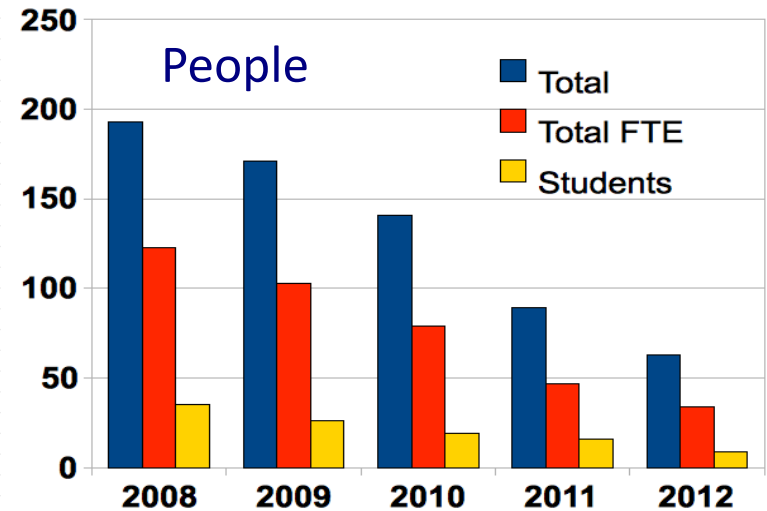
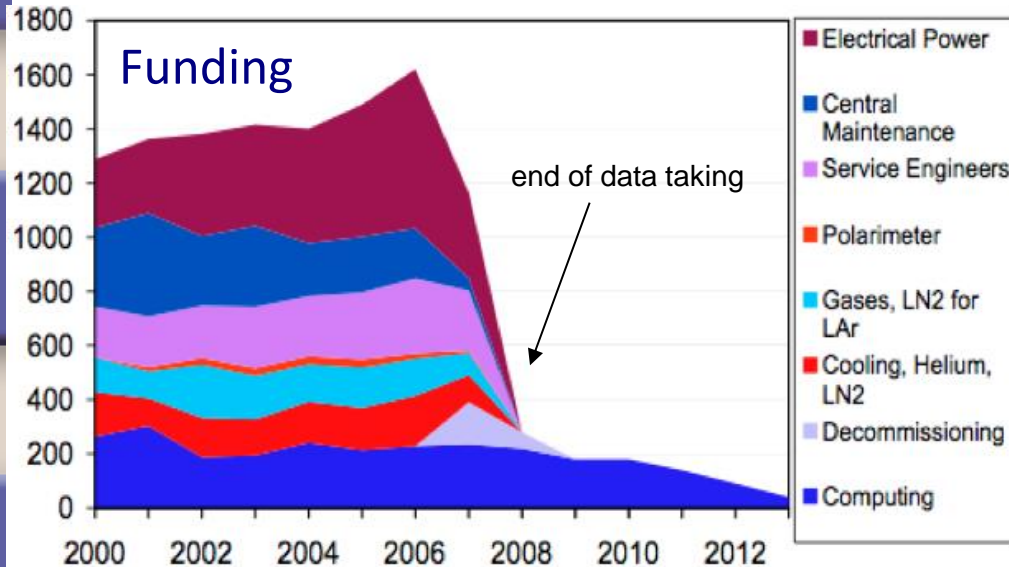
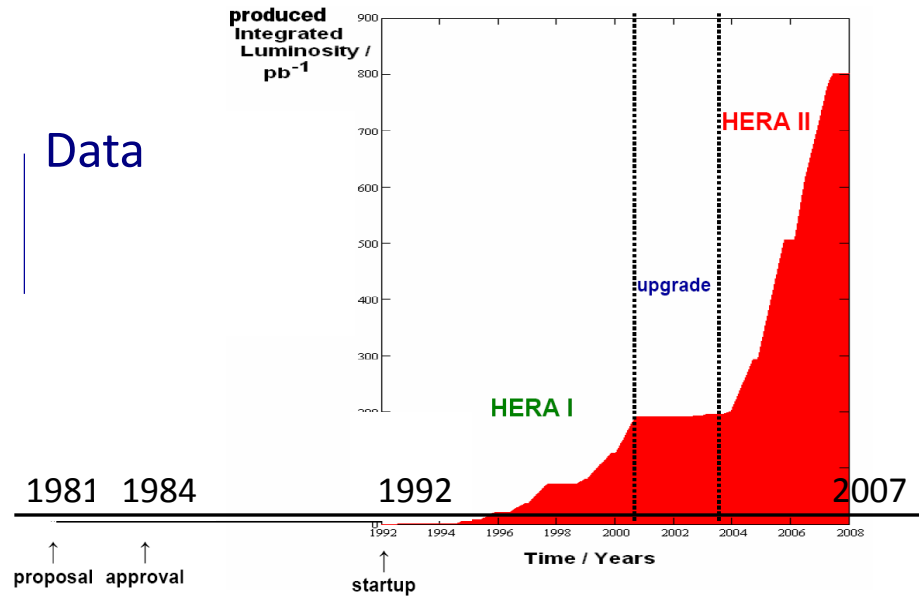
Pourquoi ca pose un problème?



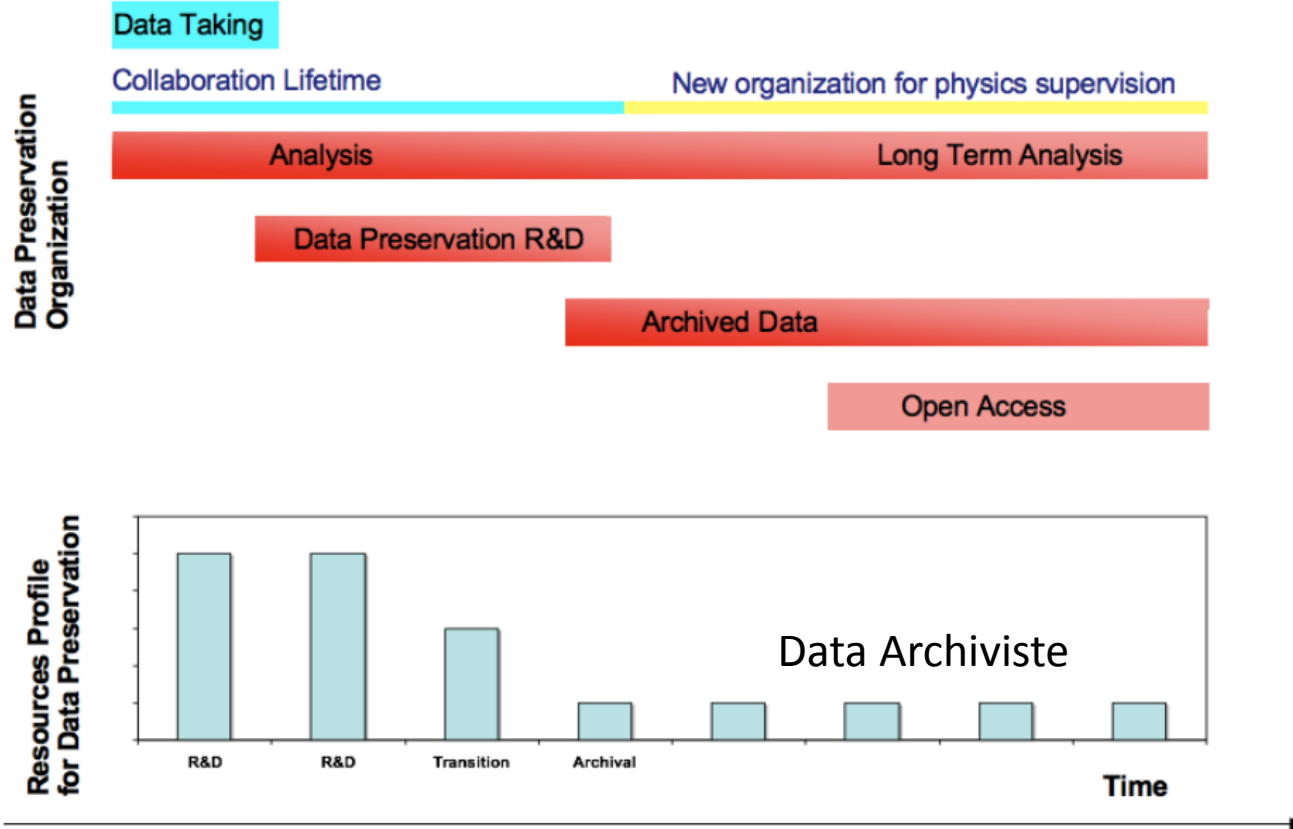
Le contexte de la préservation de données

- > La fin des programmes scientifiques est le plus gros danger pour les données
 - > Décrue des ressources et du personnel
 - > Organisation a long terme est cruciale

Data



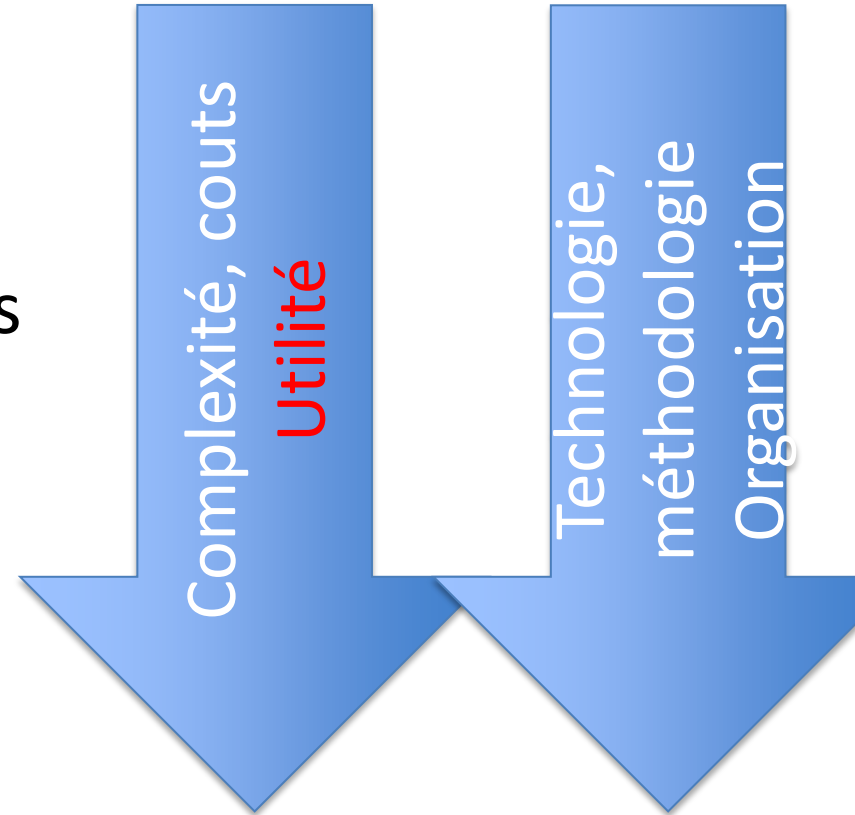
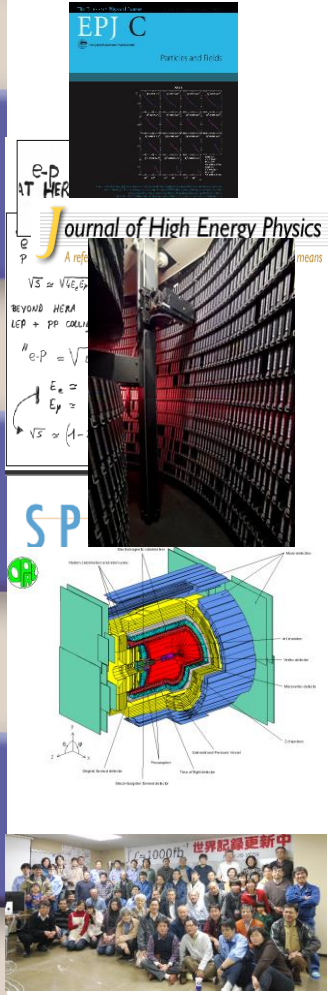
Préservation: modèle économique et organisationnel



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

Données Scientifiques

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

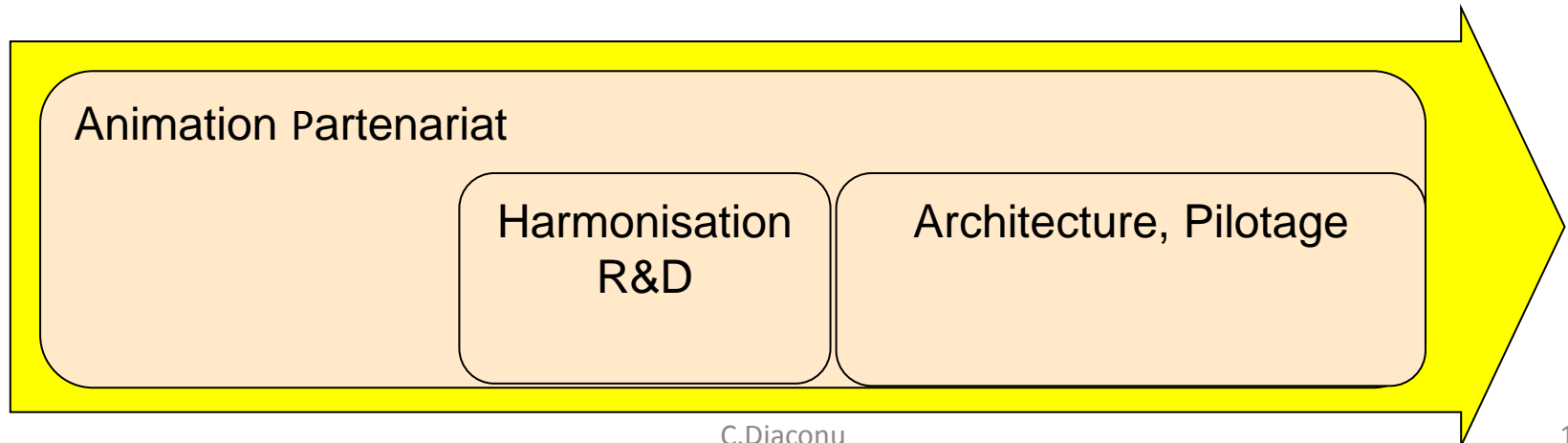


Quel modèle de préservation pour les données scientifiques?

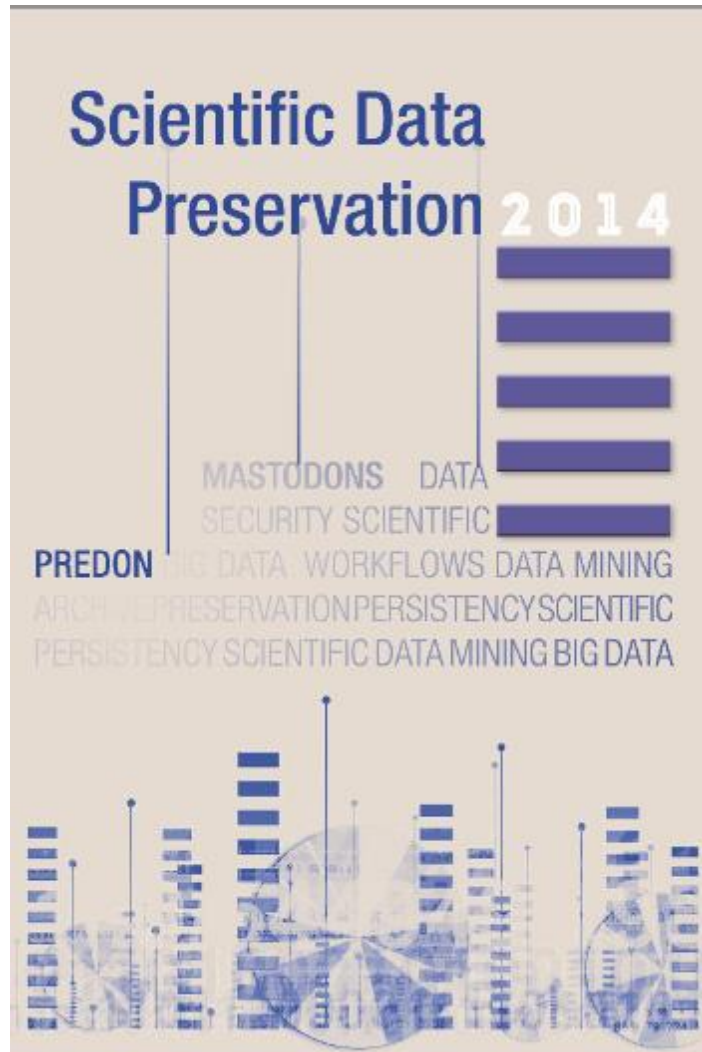
PREDON <http://predon.org>

- Projet dans le cadre « Mastodons/Big Data » de la MI/CNRS

	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	+	++	+++	+	+++



Livre blanc sur la préservation de données (« facts finding »)



CHAPTER 1: SCIENTIFIC CASE

DATA PRESERVATION IN HIGH ENERGY PHYSICS	7
VIRTUAL OBSERVATORY IN ASTROPHYSICS	15
CRYSTALLOGRAPHY OPEN DATABASES AND PRESERVATION: A WORLD-WIDE INITIATIVE	20
SATELLITE DATA MANAGEMENT AND PRESERVATION	26
SEISMIC DATA PRESERVATION	31

CHAPTER 2: METHODOLOGIES

WORKFLOWS AND SCIENTIFIC BIG DATA PRESERVATION	38
LONG TERM ARCHIVING AND CCSDS STANDARDS	42
CLOUD AND GRID METHODOLOGIES FOR DATA MANAGEMENT AND PRESERVATION	49
SCIENTIFIC DATA PRESERVATION, COPYRIGHT AND <i>OPEN SCIENCE</i>	55

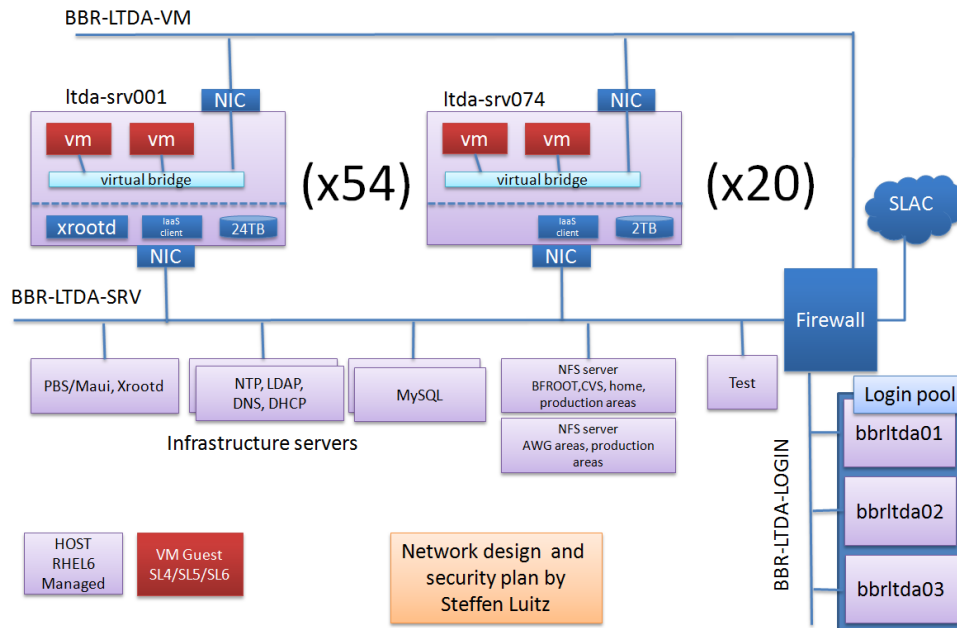
CHAPTER 3: TECHNOLOGIES

STORAGE TECHNOLOGY FOR DATA PRESERVATION	62
REQUIREMENTS AND SOLUTIONS FOR ARCHIVING SCIENTIFIC DATA AT CINES	65
VIRTUAL ENVIRONMENTS FOR DATA PRESERVATION	73

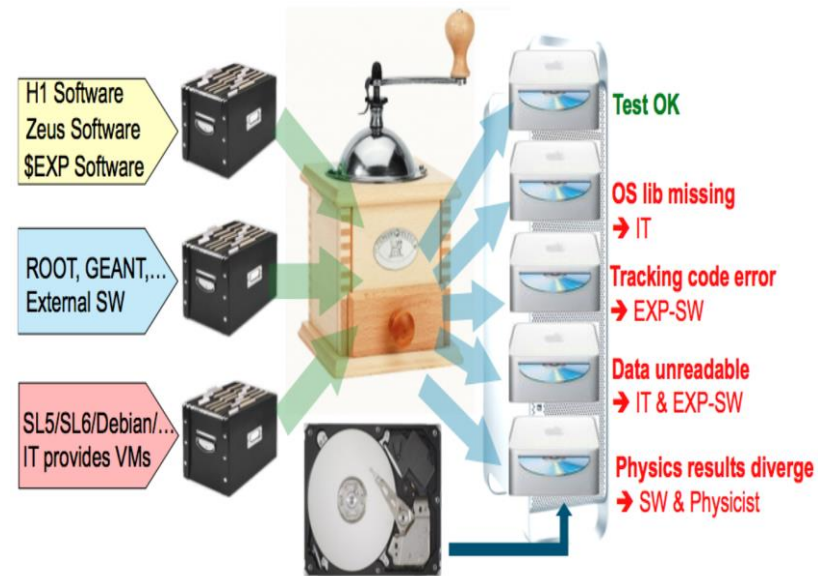
Physique des Particules

dphep.org

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)

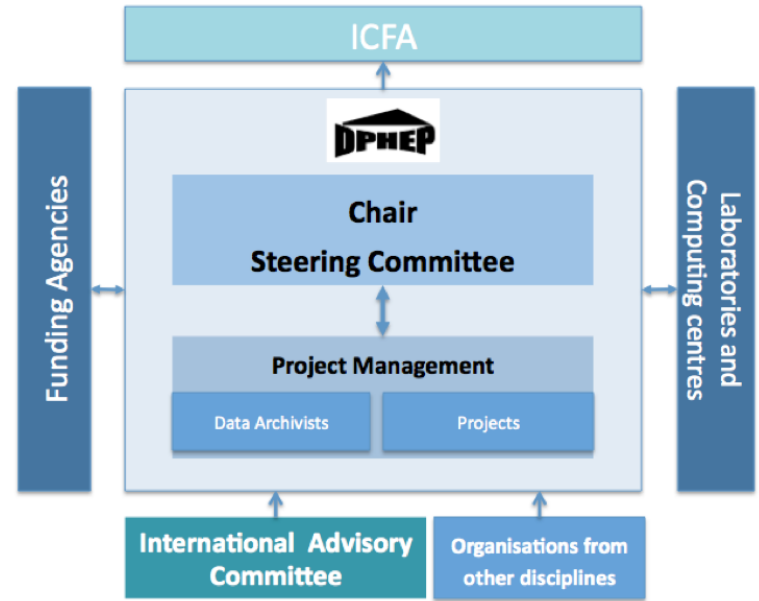


DPHEP: « Project Manager » nommé au CERN en Octobre 2012 (Scientific chair: CD)

DPHEP: International organisation

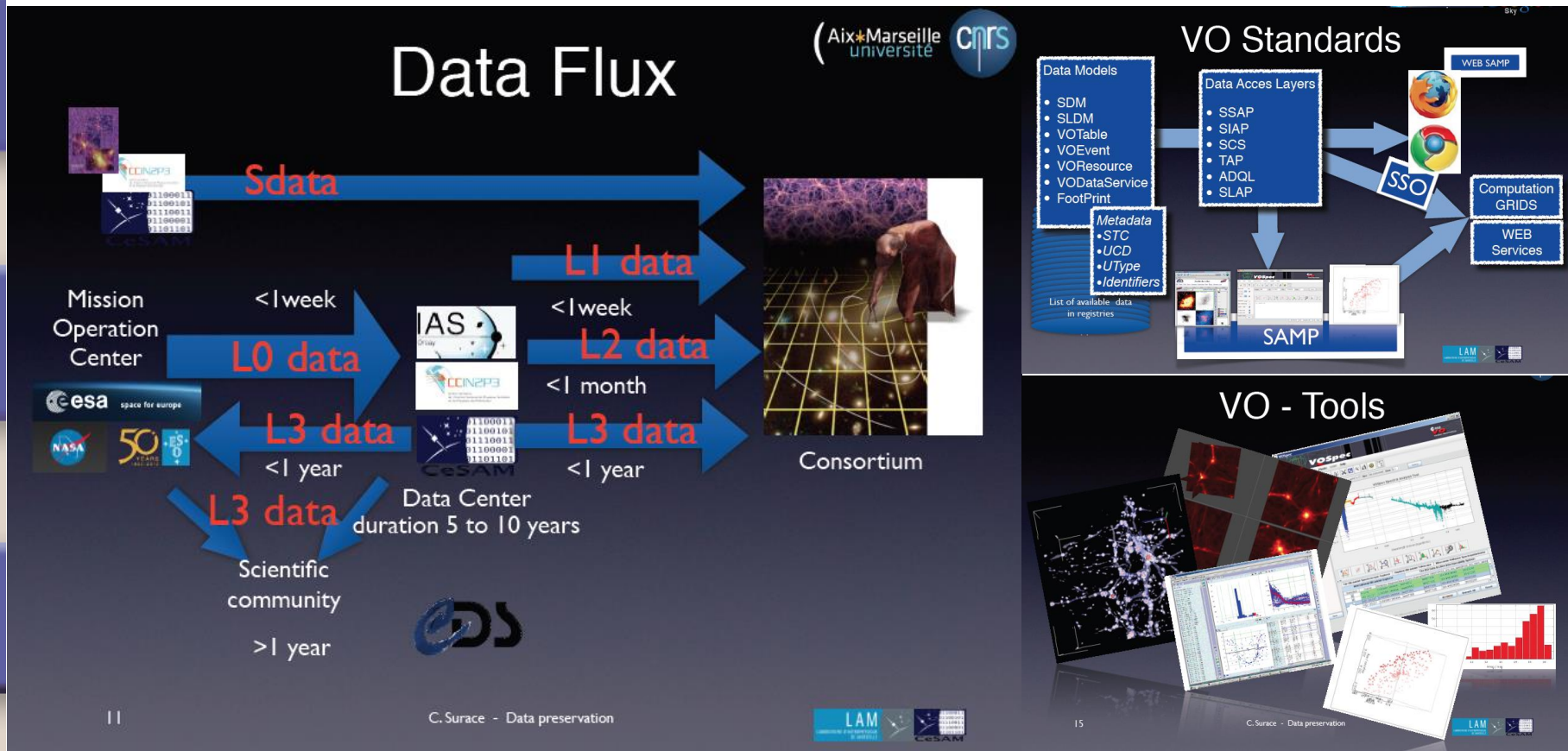


Study Group for Data Preservation and
Long Term Analysis in High Energy Physics



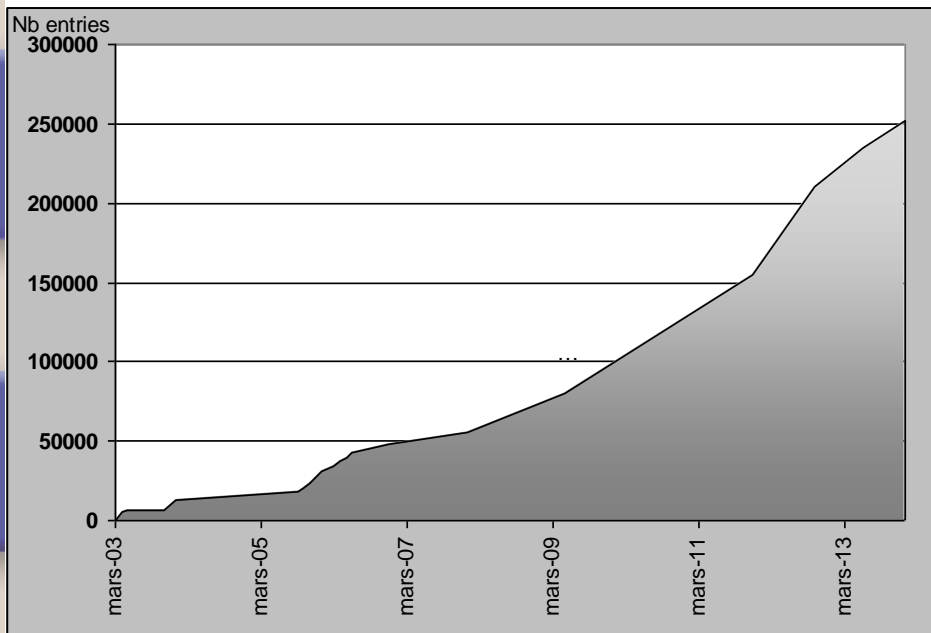
- > Study Group DPHEP:
 - > Large laboratories CERN, DESY, FERMILAB, SLAC, KEK, IHEP and experiments
- > **Breaking news: Organisation internationale mise en place**
 - > MoU signé en Juillet 2014
 - > 100 contact personnes de contact
 - > Chair: D. Diaconu Project Manager: Jamie Shiers (CERN)

Exemple projet astrophysique: Virtual Observatories

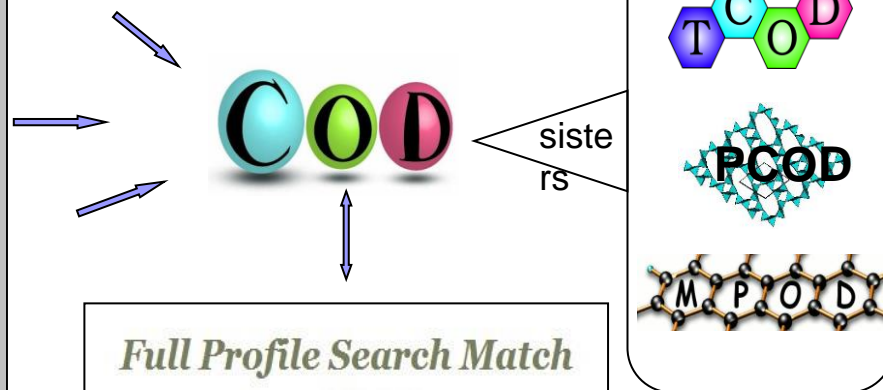


Crystallography Open Databases and Preservation: a World-Wide Initiative

Daniel Chateigner (for the COD Advisory Board)



Crystallography Open Database



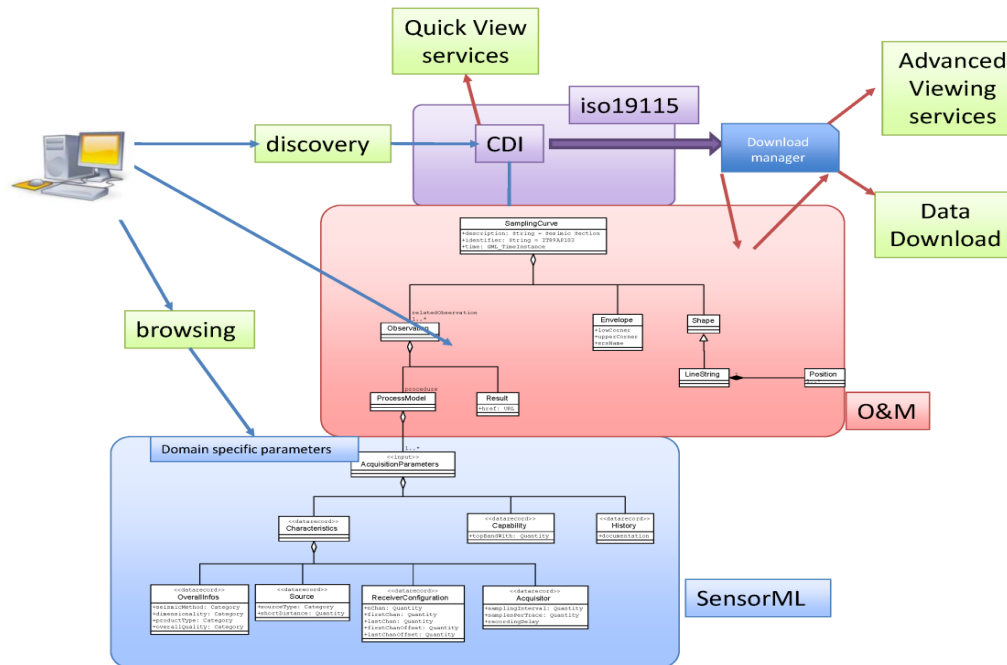
“...there **is not yet sufficient coherence** of experimental metadata standards or national policy to rely on instrumental facilities to act as permanent archives;
-there **is not sufficient funding** for existing crystallographic database organisations (which maintain curated archives of processed 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

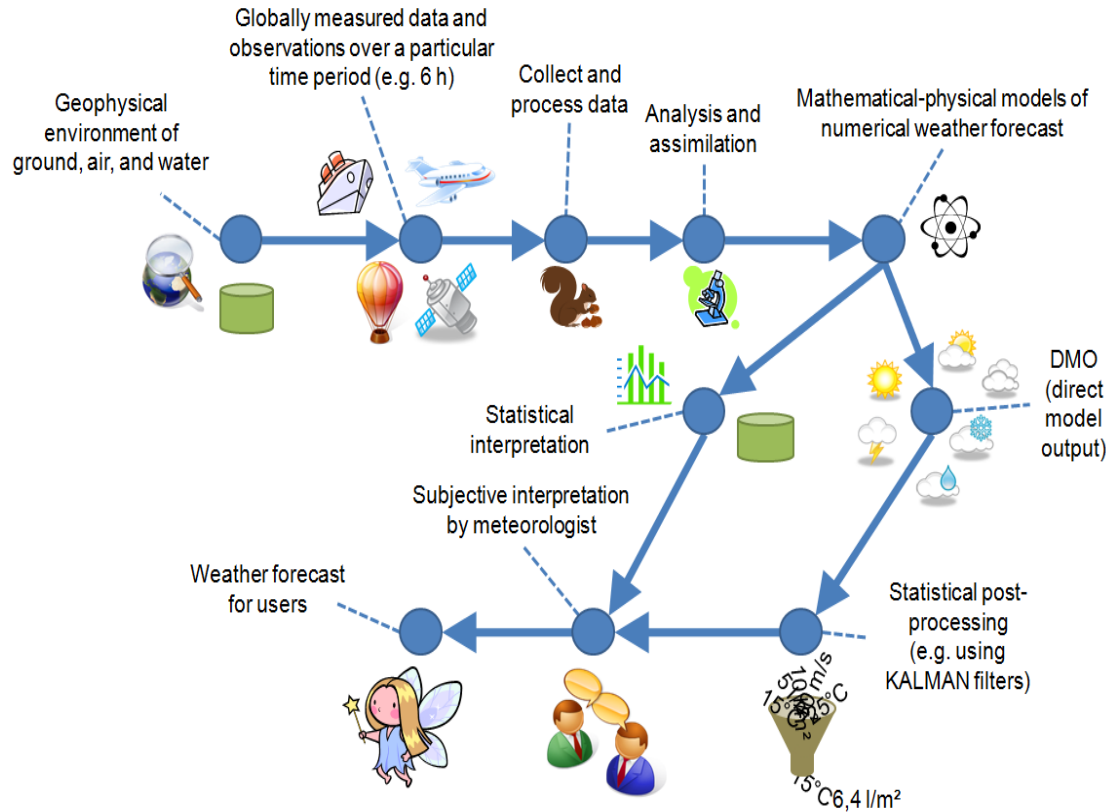
Marc SCHAMING, Institut de Physique du Globe (CNRS/UNISTRA), Strasbourg

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.

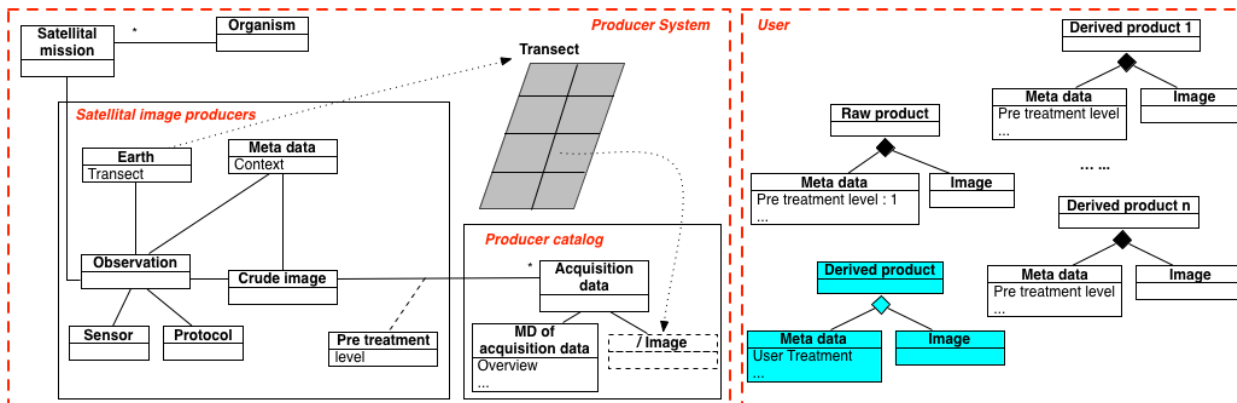


Workflows et préservation



Similarité entre les disciplines

Besoin d'une approche théorique rigoureuse



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.

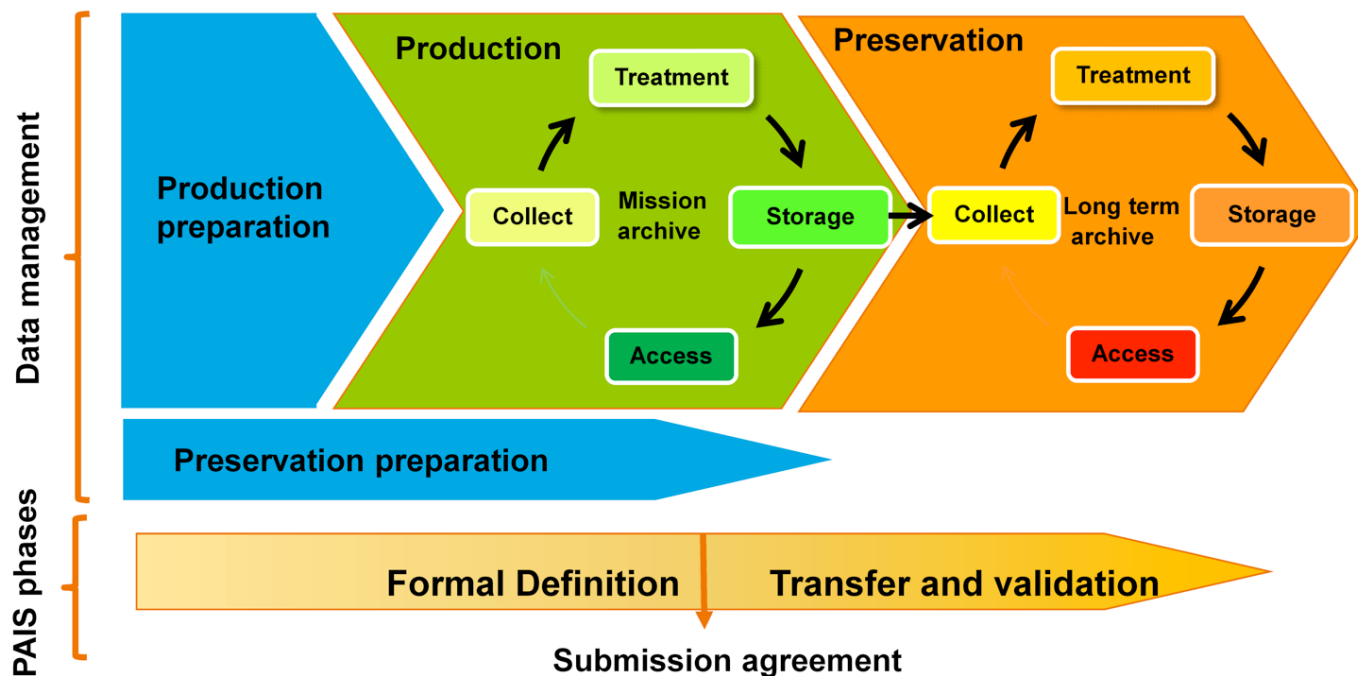


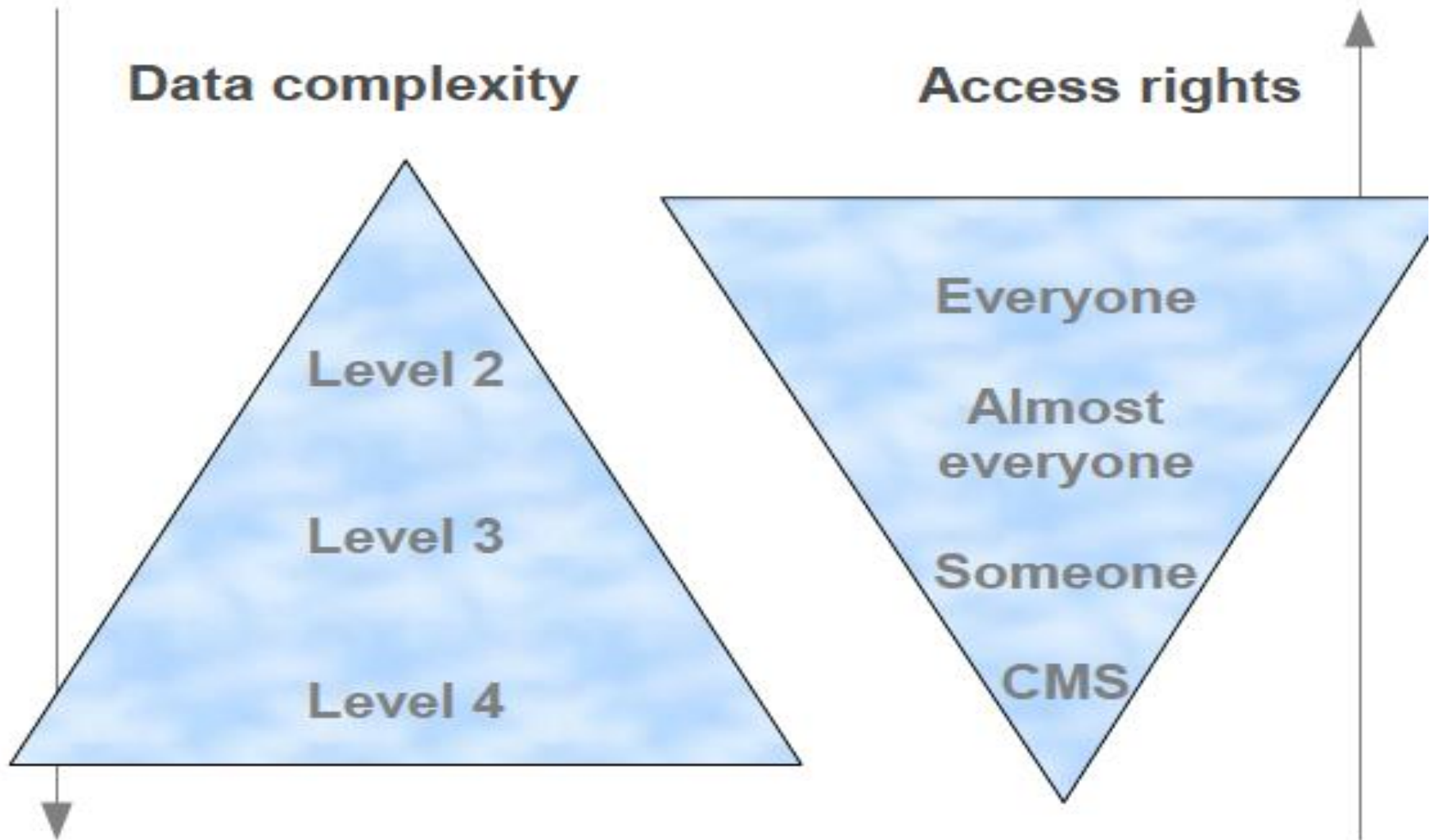
Figure 3: PAIS, preservation process and data lifecycle

Scientific Data Preservation, Copyright and *Open Science*

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

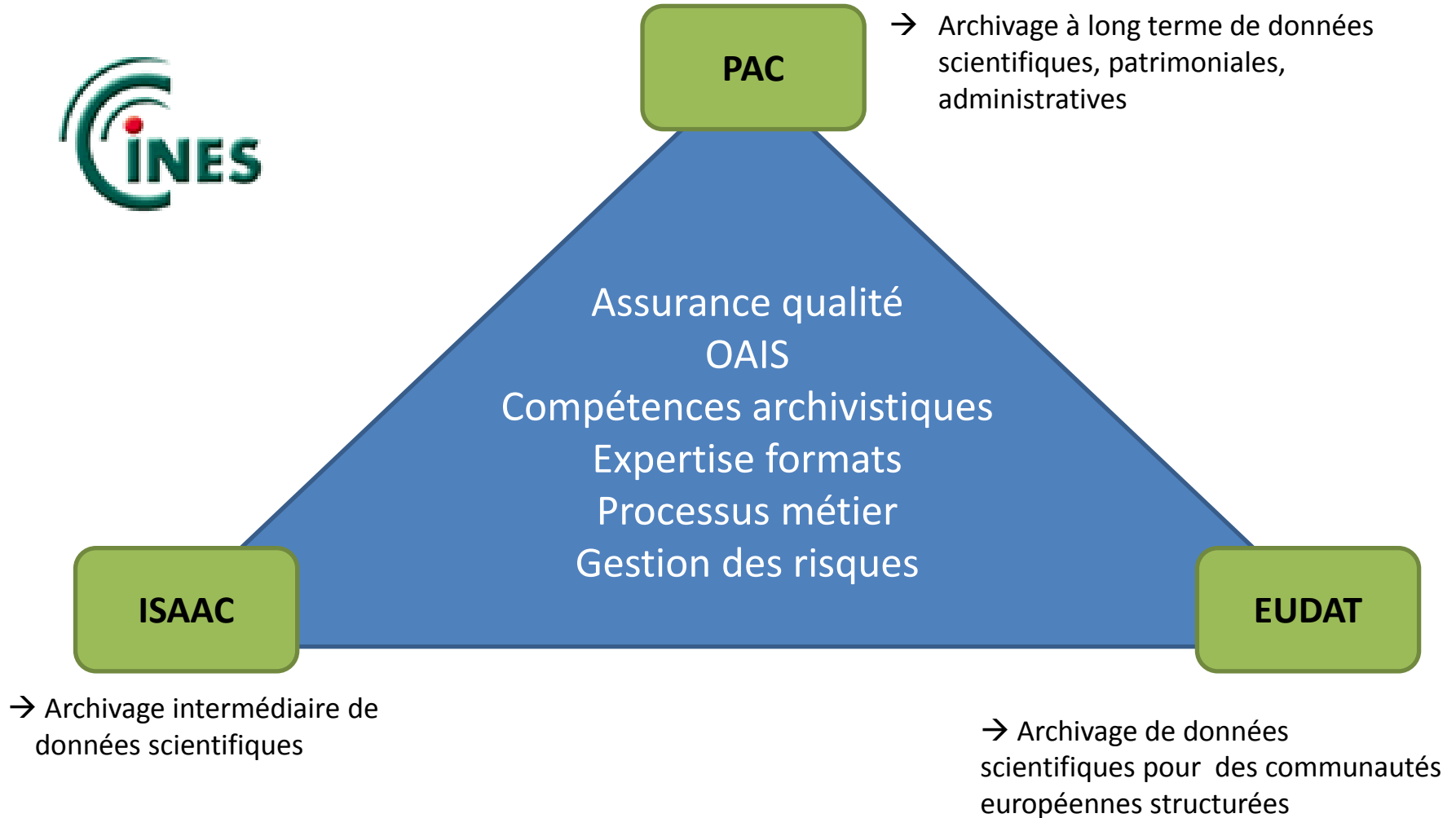
- **The best guarantee for ensuring the integrity of a resource is based on property.**
- However, isn't there a public ownership of scientific research?
 - In truth, even if the public authorities may fundamentally participate in the scientific research, this does not mean, *ipso facto*, that they own 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)

Preservation complexity levels and access



Archival expertise CINES

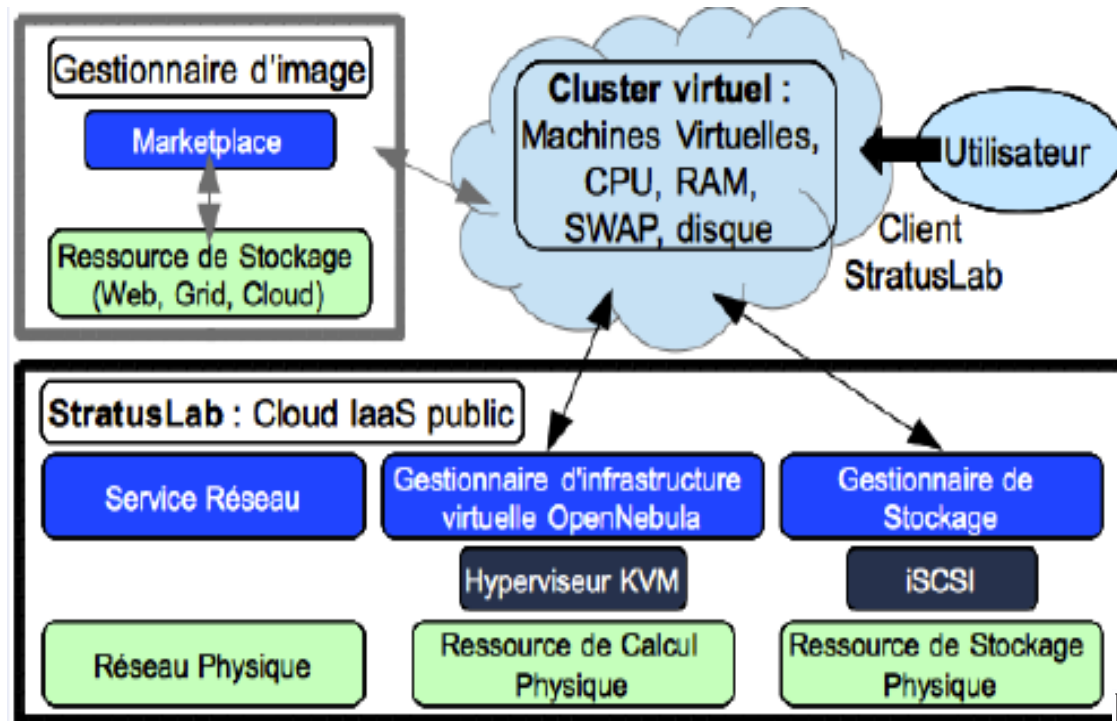
Les services d'archivage au CINES



Exemple projet: Data processing & storage in the cloud

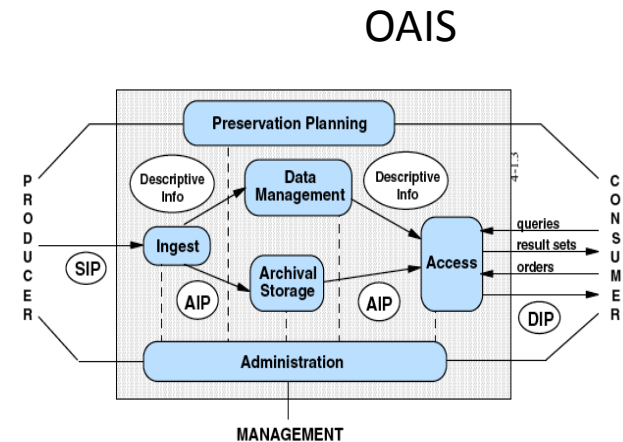
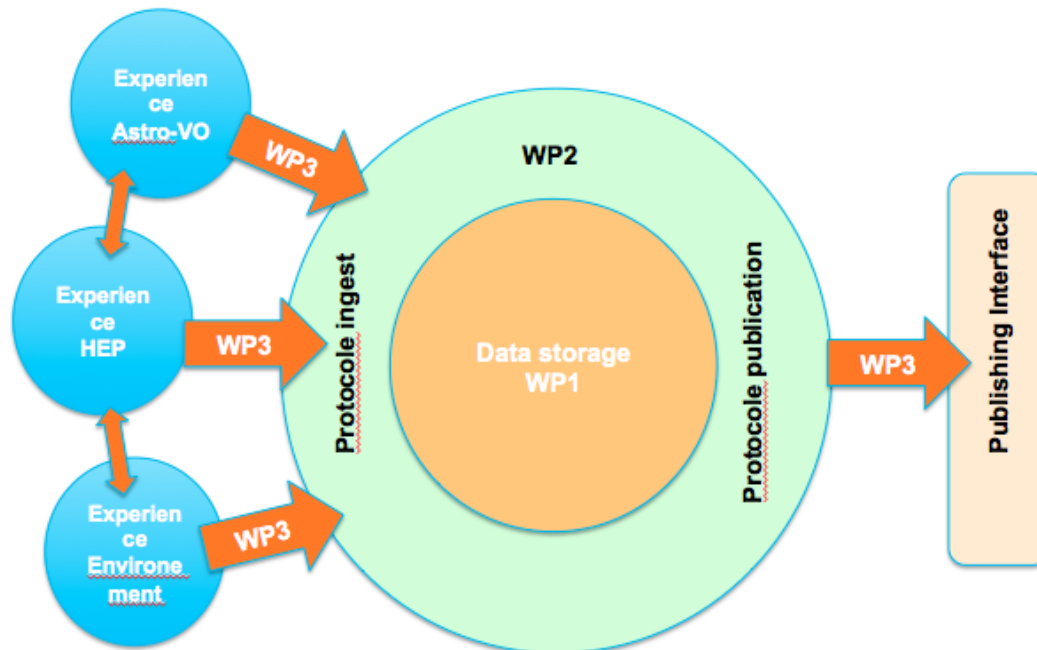
LabEx UnivEarths project at APC / François Arago Centre:

- potential of the cloud versus classical data processing and storage opportunities
- test processing on Francois Arago Centre cluster, compared with Cloud StratusLab



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

PREDON: Concept demonstrator



- But : « forcer » les frontières entre les disciplines, par exemple:
 - essayer des formats astrophysique (VOT) et des outils de visualisation (Tulip) sur des données HEP
 - Stocker des données complexes et très « custom » dans un projet de sauvegarde de données généraliste (ISAAC)

Interface données HEP – ISAAC (CINES)

PREDON 01/2013

Production Monte Carlo (données simulées) de la Collaboration H1/DESY : Interface vers le projet ISAAC-CINES

La production Monte Carlo s'appuie sur des logiciels nommés "générateurs", l'output de ces logiciels est encore **processé** en « simulation » et ensuite « reconstruction ».

Pour stocker (et ensuite être capable de trouver) ces données l'interface est la suivante :

Find Monte-Carlo Generator's File

Generator name:

Filename: (can be part of filename)

Physics Working Group: ID:

Lepton type: Radiative MC: NC/CC: Q2 min:

Analysis purpose: other:

Si je choisis par exemple le **générateur** « django.14 » je trouve une liste de « productions » qui correspondent à cette **requête**. Les productions sont l'unité **identifiées** par un « ID » de **générateur** unique

6 records found
Rows printed: 1 - 6

ID	Generator	File name	Lumi	Events	Q2 min	Date
1891	djangoh14	acs/mc/djangoh14/DE14.NCHERAZ.POSI.NGPOL.MRSH.Q21000.A00.A01	366.24	200000		100008-JAN-07
1892	djangoh14	acs/mc/djangoh14/DE14.NCHERAZ.POSI.NGPOL.MRSH.Q21000.A00.A01	3508.1	200000		100008-JAN-07
1913	djangoh14	acs/mc/djangoh14/DE14.CC.POSI.MRSH.Q21000.CDMtuned.A00.A03	15832.34	200000		100009-FEB-07
1914	djangoh14	acs/mc/djangoh14/DE14.CC.ELEC.MRSH.Q21000.CDMtuned.A00.A03	8579.25	200000		100009-FEB-07
1915	djangoh14	acs/mc/djangoh14/DE14.CC.POSI.MRSH.Q21000.B.CDMtuned.A00.A04	721.531	95	200000	100002-22-FEB-07
1916	djangoh14	acs/mc/djangoh14/DE14.CC.ELEC.MRSH.Q21000.B.CDMtuned.A00.A04	88596.87	200000		100002-22-FEB-07

[Find another generator's files](#)

Si je clique sur la première production, je retrouve des **meta-données**, en particulier, le nom de l'**ensemble** des fichiers, ainsi que un pointeur sur la « log file » qui contient

Niveau information package

• Proposition : l'unité archivée est la production

- Les métadonnées reprennent les informations de ce niveau

Generator's File

Id: 3956		Generator: django14		Date: 22-FEB-07		Working group: ReX	
File name: /acs/mc/djangoh14/DJANGO14.CC.ELEC.MRSH.Q210000.B.CDMtuned.A00-A04							
Analysis purpose: hadronic final state - jets							
Main cuts and comments: CC DIS, DJANGO + CDM QED radiative effects on, Q2 > 10000 GeV ² , no Weighting, PDF set 3036 = MRSH, Special CDM steering from HAQ di-jet CC analysis (F.Keil), tuned for High Q2.							
Lumi: 88596.87/pb		Events: 500003		E Energy: 27.6		P Energy: 920	
Lepton type: e-		Radiative MC: Y		NC/CC: CC		Q2 Min: >=10000	
Log file				simulated/reconstructed files			

[Add new comment!](#) [Find another generator's files](#)

- Et pour chaque fichier la composant...

2 records found

Simulated and reconstructed files from generator's file No. 3956

Records listed: 1 - 2

ID	p_id	Input file name	Events	Run period	Request Date	Status
4504	3956	acs/mc/djangoh14/DJANGO14.CC.ELEC.MRSH.Q210000.B.CDMtuned.A00-A04	500000	04.05.e-	22-FEB-07	done
4966	3956	acs/mc/djangoh14/DJANGO14.CC.ELEC.MRSH.Q210000.B.CDMtuned.A00-A04	500000	04.05.e-	30-OCT-07	done

[\(ctrl -> modify that job\)](#)

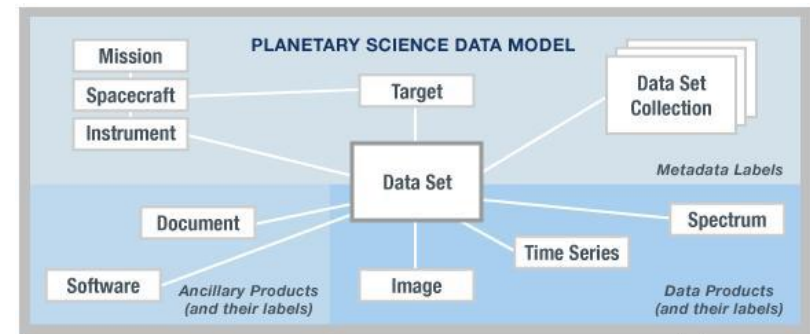
ICDE conference Chicago March 2014

http://lipade.math-info.univ-paris5.fr/lops/?page_id=96

- Workshop LoPS
 - organisé par PREDON
 - Journée dédiée à la préservation des données scientifiques
 - Agenda:

PDS- A Model-Driven Planetary Science Data Architecture for Long-Term Preservation

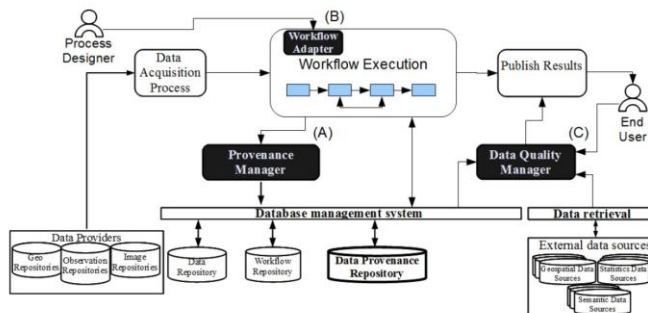
John Hughes (Jet Propulsion Laboratory, NASA)



A provenance-based approach to manage long term preservation of scientific data.

Claudia Medeiros (University of Campinas)

Quality aware workflows



Invest in capturing and maintaining data in **well-annotated, accessible, structured data repositories**

Computer Scientists, Statisticians/Data Scientists, Domain Experts (Scientists) must **systematize the analysis of massive data**

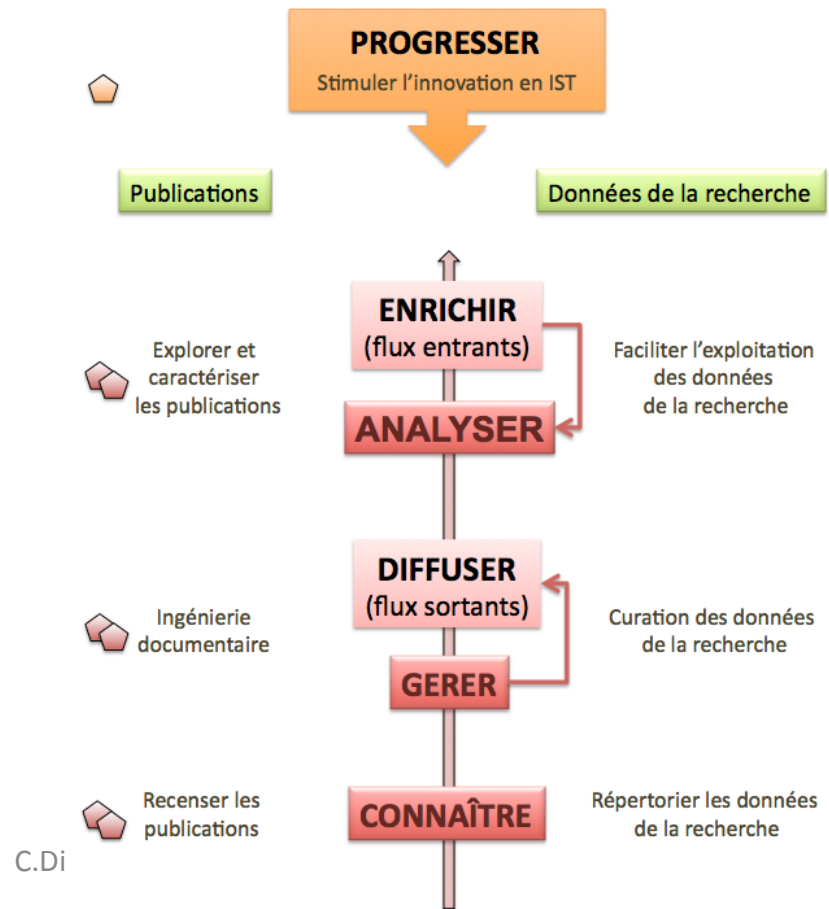
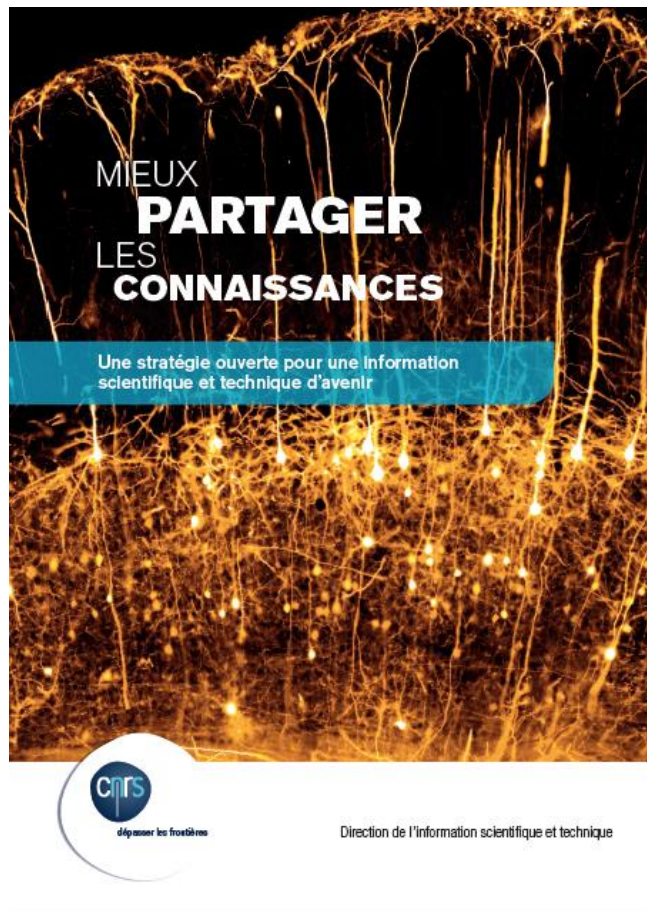
Develop **computing infrastructures** for sharing and analyzing highly distributed, heterogeneous data

Sustainability in both the data and the software infrastructures are critical

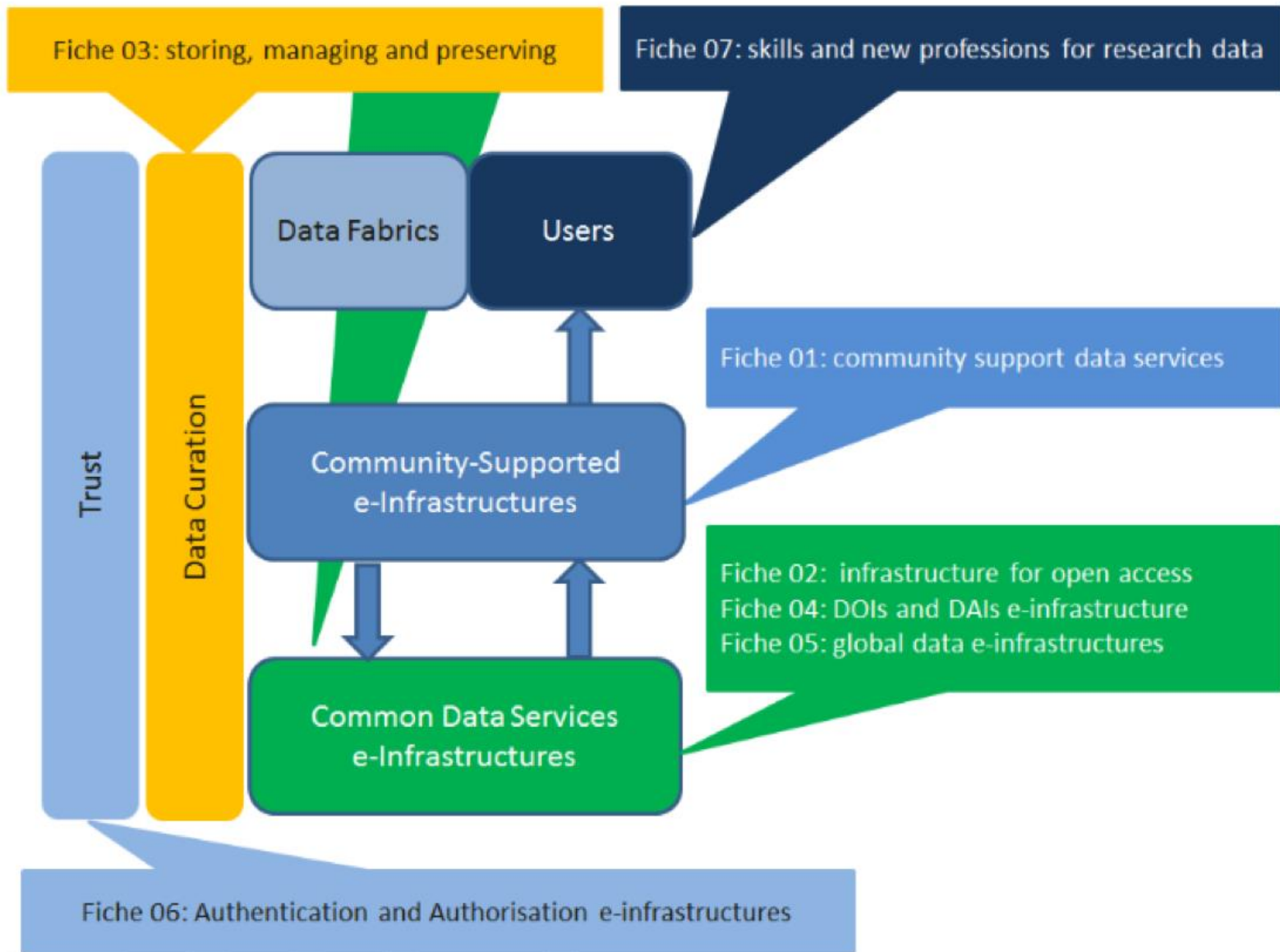
C.Diaconu

Strategie nationale IST


















- PAP3: Analyser et valoriser l'information



Opportunités H2020



Agenda

Wednesday, 5 November 2014	
09:00	Case for Preservation (until 10:30) 
09:10	The golden mine of the future: scientific data preservation - Cristinel Diaconu (CPPM, Aix-Marseille Université, CNRS/IN2P3 (FR)) 
09:40	Towards a public analysis database for LHC new physics searches - Sabine Kraml (Centre National de la Recherche Scientifique (FR)) 
11:00	Methodology for Data Preservation (until 13:05) 
11:00	Le pôle de recherche sur la conservation des données numériques (GIS-SPADON) - Pollack Jean-Dominique (Institut Jean le Rond d'Alembert, Université Pierre et Marie Curie/CNRS) 
11:35	Le projet pluri-disciplinaire IDV (Imagerie du vivant) de l'Université Sorbonne Paris Cité et quelques réflexions / méthodologies liées aux données. - Christophe Cérin (urn:Google) 
12:10	Data preservation methodology at CNES/ CCSDS - Danielle Boucon (CNES) 
14:00	Methodology for Data Preservation (until 16:00) 
14:00	Préservation de données dans le contexte IndexMed - Romain David (IMBE) 
14:45	Perenisation des donnees au CDS (Centre de données astronomiques de Strasbourg) - Gilles Landais (Centre de Données Astronomiques de Strasbourg) 
15:30	Préservation de données en imagerie médicale - Pierre Bourdoncle 
16:30	Technologies for Data Preservation (until 17:45) 
16:30	Cloud technology for algorithms preservation - Cécile Cavet (APC/ Univ. Paris 7)  <div style="float: right; margin-right: 10px;"> summary  </div>
17:00	EUDAT - STEPHANE COUTIN (C) 
17:45	Discussion (until 18:30) 
20:00	Social Dinner (until 23:00) () 

Objectifs du workshop

- Tour des projets au sein de PREDON
- Récevoir des points de vue d'autres disciplines
 - Nouvelles approches: documentation, juridique, économique
- Connexion aux projets similaires en France
- Document: « **Scientific Data Preservation 2015** »
 - Contributions suite aux workshops PREDON
 - Contributions invitées: propositions?

Conclusions

- Les données scientifiques ont un potentiel qui dépasse le cadre de recherche initial et qui doit être exploité à long terme
 - Preservation \Leftrightarrow Accès ouvert
- La préservation de données scientifique est économiquement avantageuse:
 - Recherche à bas cout
- Une technologies de frontière est nécessaire
 - Préservation de toute la chaine « grise »
 - Virtualisation, cloud computing, workflows....
- La collaboration multi-disciplinaire est essentielle
 - au niveaux national et international
- Projet PREDON: animation, R&D, architecture
 - Un GRAND merci aux organisateurs!

BACKUP

Groupe d'études PREDON

2013 2012 Prop.

> **IN2P3**

- Cristinel Diaconu, Dirk Hofmann, Angélique Pèpe, Magali Damoiseaux, D. Christofol (CPPM, Marseille)
- Sabine Kraml (LPSC, Grenoble)
- Giovanni Lamanna (LAPP, Annecy)
- Volker Beckmann (APC, Centre Francois Arago, Paris 7)

> **CCIN2P3**

- Ghita Rahal, Jean-Yves Nief (CC-IN2P3)

> **INSU**

- Christian Surace (LAM/OAMP Cesam, Marseille)

> **INS2I**

- Mustapha Lebbah (LIPN, Paris 13)
- Salima Benbernou (LIPADE, Paris 5)
- Anne Laurent, Sophie Nicoud (LIRMM, Montpellier)

> **CINES**

- Stéphane Coutin, Marion Massol (CINES, Montpellier)

> **IRD**

- Thérèse Libourel, Yuan Lin (Espace DEV)

Nouveau contacts en 2013/2014 suite aux workshops:

Daniel Chateigner, CRISMAT/ENSICAEN, données **cristallographie**

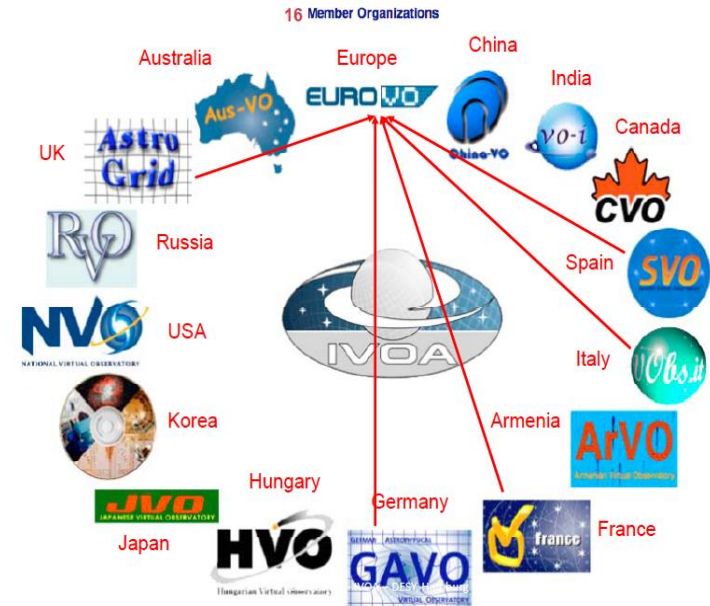
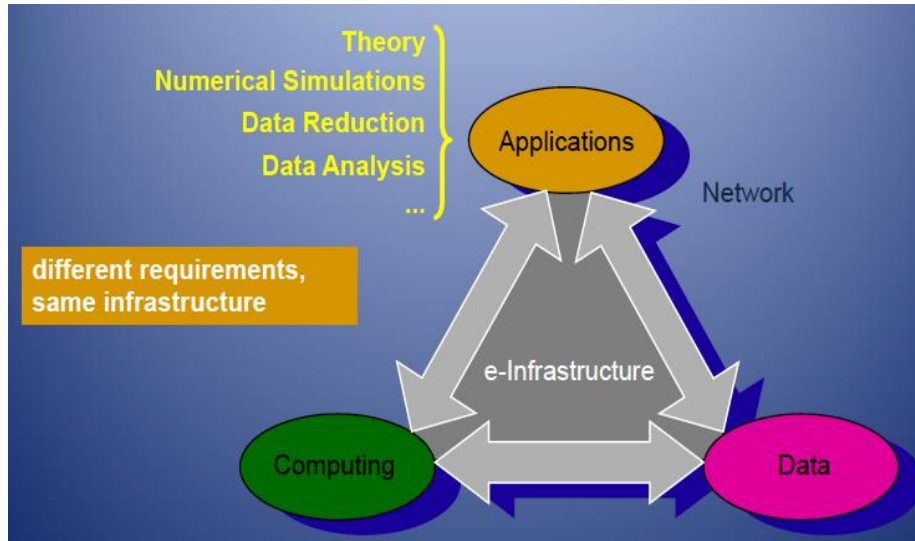
Marc Schaming, Institut de **Physique du Globe** (CNRS/UNISTRA), IPG Strasbourg

Catherine Boisson de l'Observatoire de Meudon / LUTH/**INSU** CTA

Danièle Boucon, expert en préservation de données **CNES**

Jean-Dominique Pollack, LAM, UPMC, GIS **SPADON**

Virtual Observatories in Astrophysics



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

Préservation des connaissances

- Le stockage des données à long terme demande une organisation rigoureuse
- Le vrai challenge technique est la préservation des connaissances « meta-digitales »

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

The logo for the journal 'nature', featuring the word 'nature' in a white serif font on a dark red rectangular background.

DPHEP : définition des niveaux de préservation

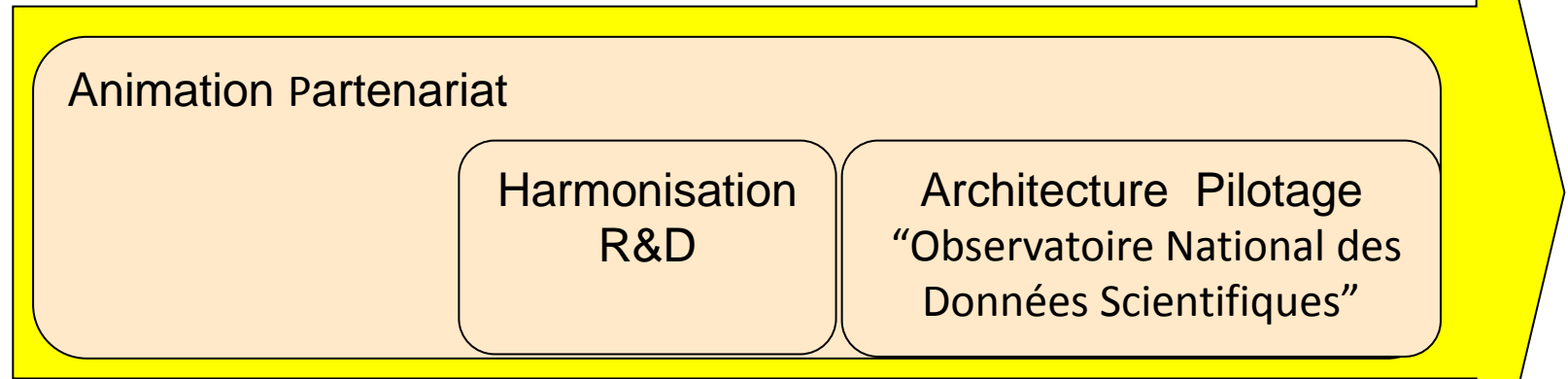
> En progression de la complexité et les couts

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	

MASTODONS

- Stockage et gestion de données (par exemple, dans le Cloud), sécurité, confidentialité.
- Calcul intensif sur des grands volumes de données, parallélisme dirigé par les données.
- Visualisation de grandes masses de données.
- Extraction de connaissances, datamining et apprentissage.
- Qualité des données, confidentialité et sécurité des données.
- Problèmes de propriété, de droit d'usage, droit à l'oubli.
- **Préservation/archivage des données pour les générations futures.**
 - **PREDON (PREservation des DONnees)**

PREDON



- Court terme (2012/2013 et après): **Animation et partenariat**
 - Elargir le champ de réflexion, constituer un consortium multi-disciplinaire
- Medium terme (2013/2014) : **Harmonisation et projets R&D**
 - Communication: exchanges and workshops
 - Livre blanc sur la préservation et la mise à disposition des données scientifiques dans un contexte multi-disciplinaire
 - Démonstrateur accès et préservation de données scientifiques complexes
- Long term (2015/2016) **Architecture et pilotage**
 - “Observatoire National des Données Scientifiques”
 - Coalition de grands centres de données et projets multi-disciplinaires
 - Support et suivi des lots de données scientifiques : accès et préservation

Workshop on Data Preservation at ICDE 2014

LOPS@ICDE

WORKSHOP ON LONG TERM PRESERVATION FOR BIG SCIENTIFIC DATA



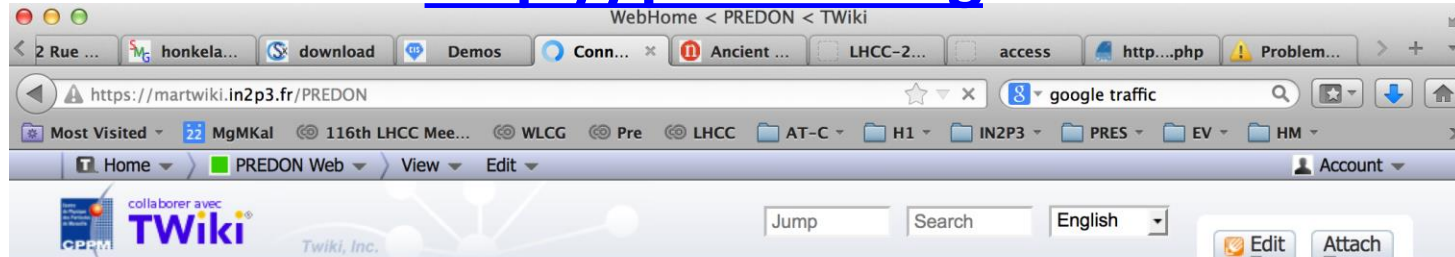
- Coordonnateurs workshop: S.Benbernou, C. Diaconu
- <http://lipade.math-info.univ-paris5.fr/lops/>
- LOPS will be held in conjunction with the 30th IEEE International Conference on Data Engineering. Chicago, IL, USA. March 31-April 4, 2014.

PREDON 2014

- Organisation Workshop LOPS@ ICDE2014
- Continuation et initiation de nouveaux mini-projets et démonstrateurs pour des cas spécifiques de préservation de données (stages)
 - HEP-Data @ ISAAC
 - Formats de données transdisciplinaires
- Réunions du groupe de travail : nouveau contacts, séminaires
 - Extensions possibles à d'autres domaines (bio, IST, économie)
 - Aborder des questions communes (cout, persistance, open access, éducation, outreach etc.)
- Organisation d'un Atelier sur la préservation des données scientifiques et en relation avec la thématique « Big Data »
 - Publication PREDON: 2015
- Participation aux groupes de travail au niveau international et aux projets et consortia en cours de constitution pour des programmes de financement H2020.

Site web PREDON

<http://predon.org>



Tags: [create new tag](#), [view all tags](#)



A project for scientific data preservation in France



Breaking News!



January 21st, 2014: PREDON document "Scientific Data Preservation", a facts finding white paper produced following 2012/2013 workshops is available.

"Data observatories, based on open access policies and coupled with multi-disciplinary techniques for indexing and mining may lead to truly new paradigms in science. It is therefore of outmost importance to pursue a coherent and vigorous approach to preserve the scientific data at long term. The preservation remains nevertheless a challenge due to the complexity of the data structure, the fragility of the custom-made software environments as well as the lack of rigorous approaches in workflows and algorithms. [...]"

"The present document includes contributions from the participants to the PREDON Study Group, as well as invited papers, related to the scientific case, methodology and technology. This document should be read as a "facts finding" resource pointing to a concrete and significant scientific interest for long term research data preservation, as well as to cutting edge methods and technologies to achieve this goal. A sustained and coherent and long term action in the area of scientific data preservation would be highly beneficial."

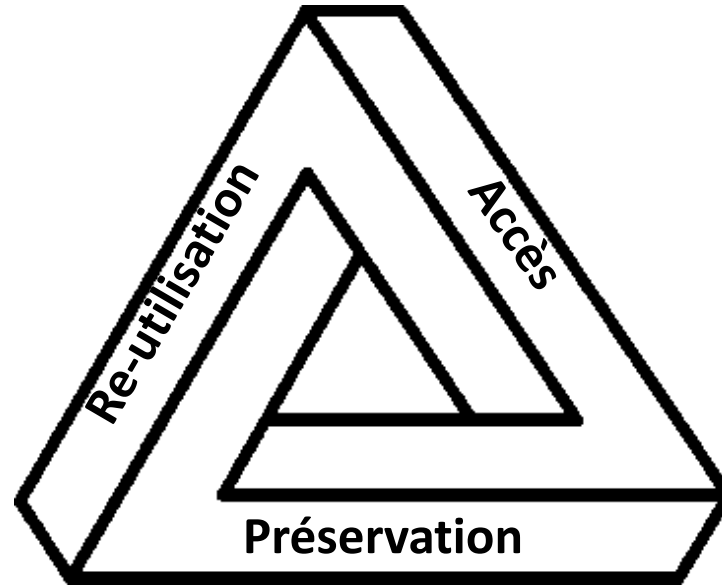
Challenges

Scientific data collected with modern sensors or dedicated detectors exceed very often the perimeter of the initial scientific design. These data are obtained more and more frequently with large material and human efforts. A large class of scientific experiments are in fact unique because of their large scale, with very small chances to be repeated or superseded by new experiments in the same domain: for instance high energy physics and astrophysics experiments involve multi-annual and even yearly repeatable. Other scientific experiments are in fact unique by nature: earth science, medical sciences etc. since the collected

Summary of information from the (pre-LHC) experiments

	BaBar	H1	ZEUS	HERMES	Belle	BESIII	CDF	DØ
End of data taking	07.04.08	30.06.07	30.06.07	30.06.07	30.06.10	2017	30.09.11	30.09.11
Type of data to be preserved	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)
Data Volume	2 PB	0.5 PB	0.2 PB	0.5 PB	4 PB	6 PB	9 PB	8.5 PB
Desired longevity of long term analysis	Unlimited	At least 10 years	At least 20 years	5-10 years	5 years	15 years	Unlimited	10 years
Longévité recherchée: > 10 ans								
Current operating system	SL/RHEL3 SL/RHEL 5	SL5	SL5	SL3 SL5	SL5/RHEL5	SL5	SL5 SL6	SL5
Languages	C++ Java Python	C C++ Fortran Python	C++	C C++ Fortran Python	C C++ Fortran	C++	C C++ Python	C++
Simulation	GEANT 4	GEANT 3	GEANT 3	GEANT 3	GEANT 3	GEANT 4	GEANT 3	GEANT 3
External dependencies	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

Préservation, réutilisation, libre-accès



- La préservation suppose la mise à disposition en accès libre
 - Maximiser le bénéfice
- ← Le libre-accès facilite la préservation à long terme
 - ← Elargir la communauté, multiplier les connaissances

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>

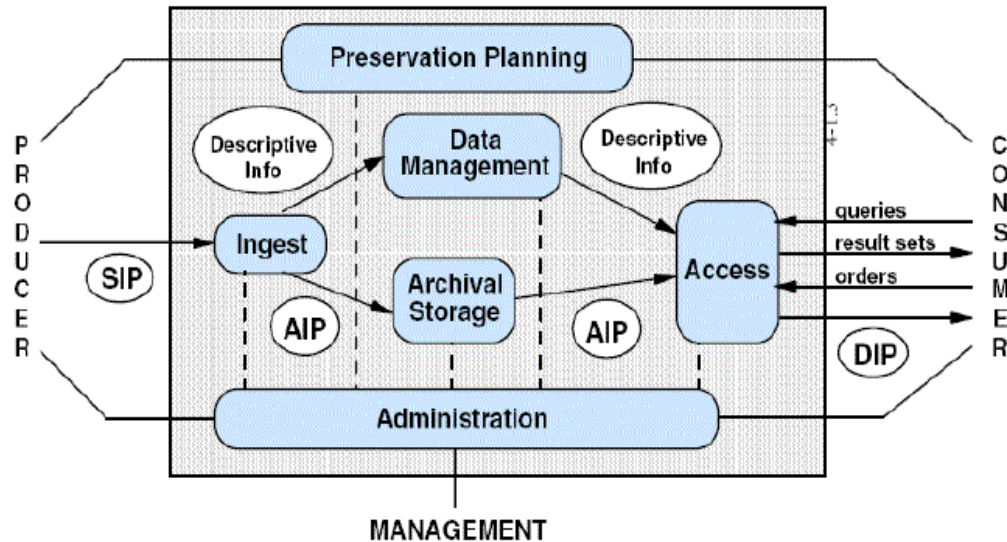


FIGURE 2.1: **The OAIS Reference Model**

<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)

Documentation

- Une tache considérable

- > **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



PREDON: Challenges

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