



# Data preservation (?) at CCIN2P3

*First workshop on data preservation and long term analysis in HEP*

**Fabio Hernandez**  
fabio@in2p3.fr

Hamburg, January 26th 2009



l r f u  
c e a  
s a c l a y

# ▶ Contents



- Brief introduction to the site
- Data preservation
  - What we do
  - What we don't/can't do
- Final remarks
- Questions & comments

# ▶ IN2P3 Computing Centre

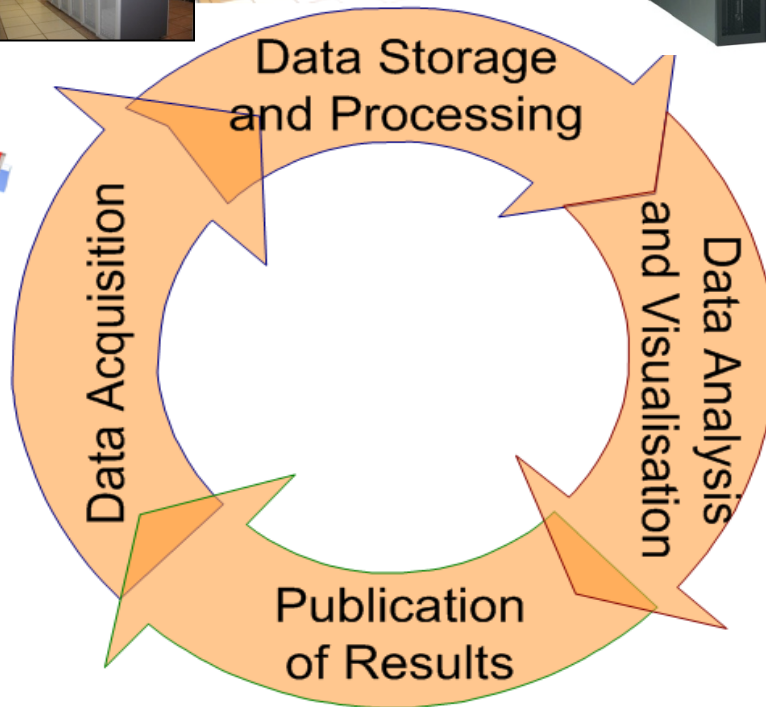
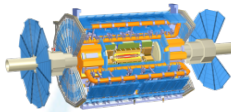
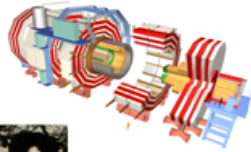
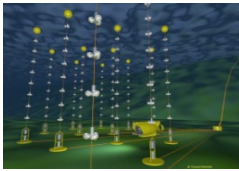
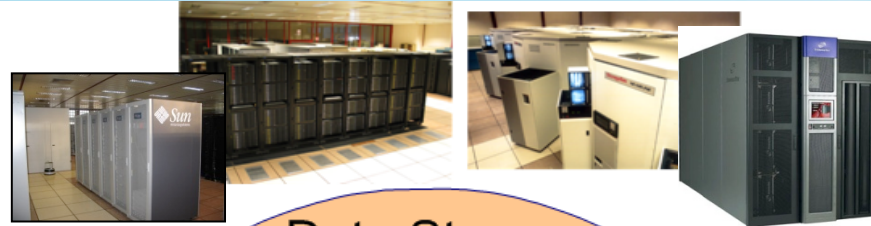


- National high-throughput data processing **shared** facility
  - not co-located with an experimental site
- Mission
  - mass storage repository
  - high-throughput computing facilities
  - technical consulting and training services for IN2P3 laboratories
  - web hosting, video-conferencing, e-mail infrastructure, news, wiki, webcast, ...
  - 24x7 service
- Users
  - ~70 research groups, mainly international collaborations in nuclear physics, particle physics and astro-particle physics
  - since 2002 serving also bio-medical applications and doing some technology transfer with industry



- Co-funded by CNRS and CEA/DSM
- ~70 FTEs

# IN2P3 Computing Centre (cont.)



## ▶ IN2P3 Computing Centre (cont.)



- Facility shared by several experiments
  - Operating a tier-1 for the 4 LHC experiments
- Connected to several grid infrastructures, serving several virtual organizations
- Users are not on site
- Vast majority of people involved in the centre's operation are not directly involved in (i.e. formally not members of) any experiment

# ▶ Who's storing data at CCIN2P3?



- On line/nearline data
  - **Disk: 49 groups**
    - *Group = experiment and/or laboratory*
  - **Mass storage system-managed data: 61 groups**
    - *These data are mostly on tape*
- Off line data
  - **Tape: 30 groups**
  - **More on this later**

# ▶ Who's storing data at CCIN2P3? (cont.)



## Top 10 users of data storage services

*(ranked by amount of stored data as of 31/12/2008)*

Rank	Tape (HPSS)	Disk
1	D0	Atlas
2	CMS	CMS
3	Babar	LHCb
4	Atlas	SuperNovae
5	QCD	Planck
6	Phenix	Alice
7	Virgo	Babar
8	Alice	EROS
9	Antares	D0
10	Pierre Auger	Phenix

# ▶ Data preservation



- Roles and responsibilities
  - We consider data preservation as a shared responsibility between our users (i.e. the experiments) and us
  - Those roles and responsibilities are informally defined and not documented
    - *For instance, when supporting a new experiment, we don't formally establish « who does what » on this subject*
- No explicit funding for this activity



# ▶ Data preservation: what we do



- Low-level services
  - Migration of data from one generation of storage media to the next
    - Applies both to data on disk and on tape
    - Migration process (both technically and organizationally) is agreed with the experiments
    - We trigger the process, regularly but not systematically
  - Migration of data as a result of changes in the data format
    - For instance, during the migration from mainframe-based computing to the UNIX world
  - Conservation of tape cartridges (in the vault)
    - In several cases, they contain raw data for some experiments (in particular those which the experimental site is not well connected to the network)
    - Catalogue of information on the contents of those cartridges is (expected to be) maintained by the experiment themselves
    - We store the cartridges in appropriate environmental conditions (temperature, humidity, reasonably low levels of dust, etc.)

# ▶ What we do: the vault (cont.)



© CC-IN2P3 / CNRS / f. de Stefanis

# ▶ What we do: the vault (cont.)



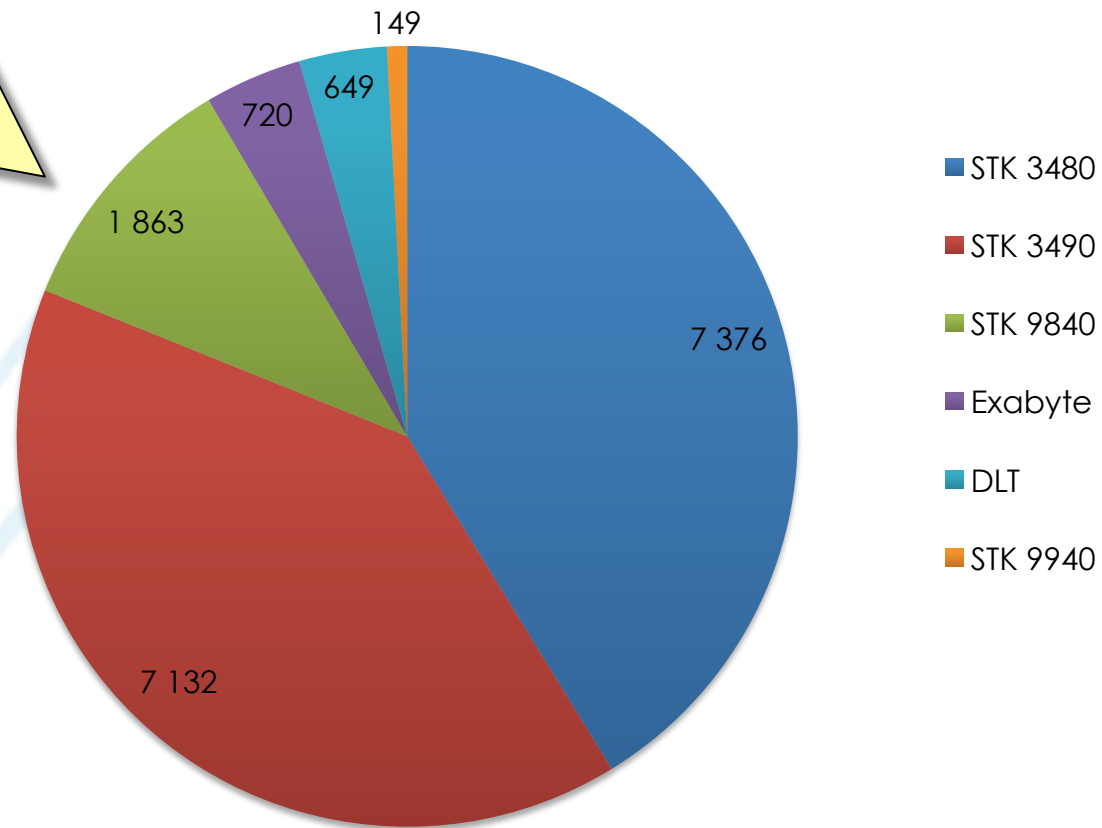
# ▶ What we do: the vault (cont.)



Number of tape cartridges in the vault

Around 18.000 cartridges currently in the vault

Was 120.000 5 years ago

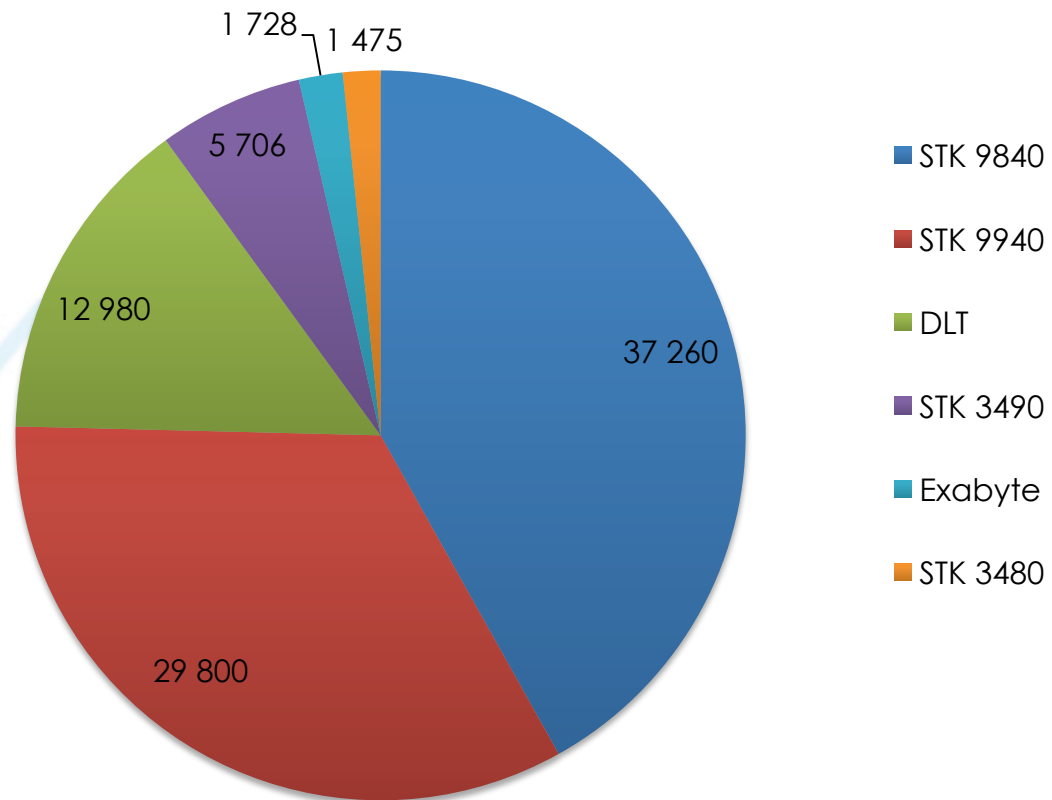


# ▶ What we do: the vault (cont.)



**Amount of data in the vault**  
(in GB)

Around 90 TB of data

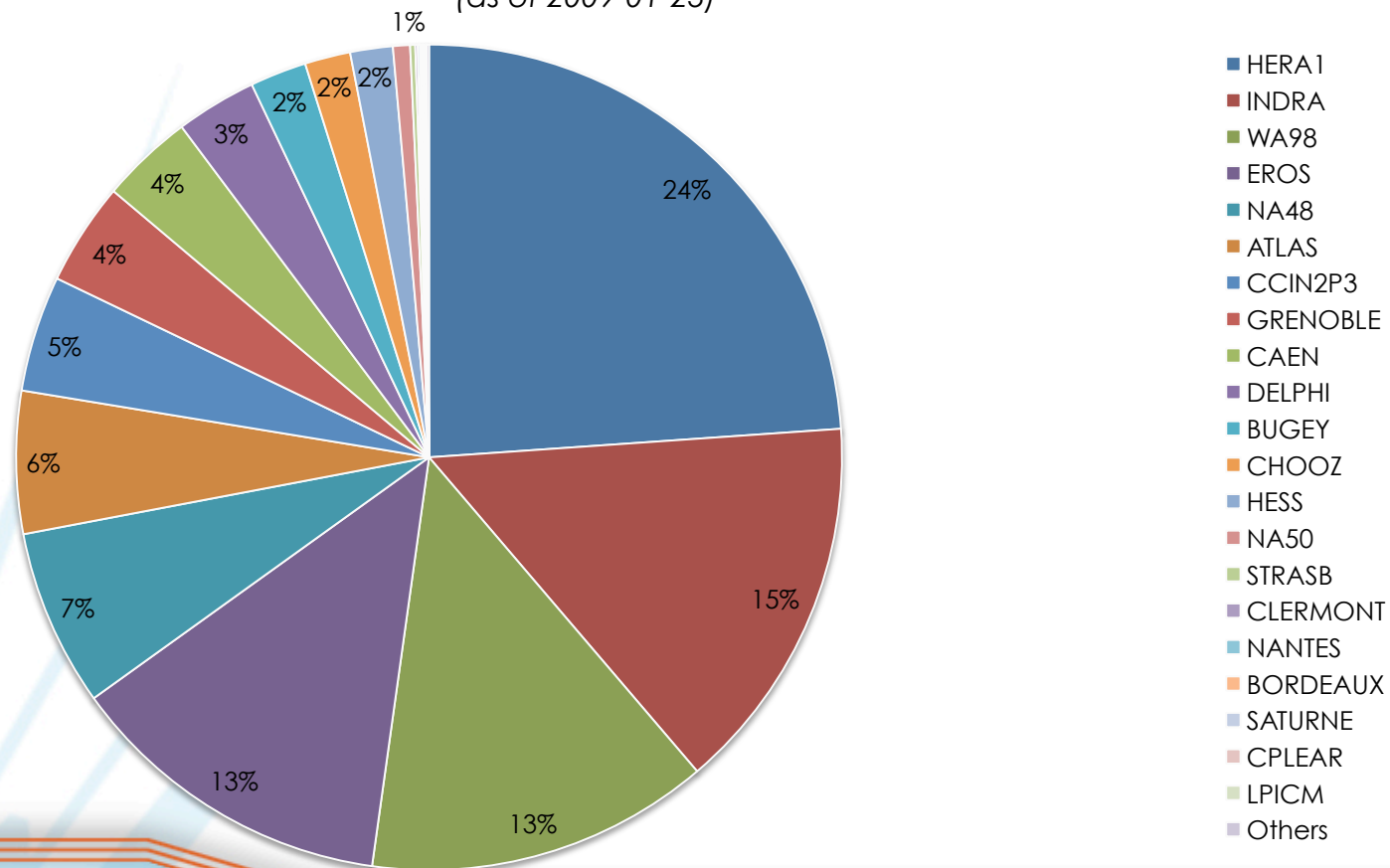


# ▶ What we do: the vault (cont.)



## CCIN2P3 – Experiments and labs using the cartridge vault

% of the number of cartridges  
(as of 2009-01-23)



# ▶ What we don't/can't do



- The following applies for data on tape cartridges in the vault
- We don't check the integrity of the data
  - In several cases, we don't have on site the technology to read those tapes
  - Although we keep and maintain updated a catalogue of cartridge metadata, we certainly don't know what the contents of those cartridges is
  - Experiments are aware of this, but they still require us to keep them
    - *Our liaison persons don't always feel entitled to decide on what to do with those data*
    - *For some experiments, some of the traditional liaison persons have already retired*
- We don't assign preservation metadata

## ▶ What we don't/can't do (cont.)



- We cannot ensure data is stored in file formats appropriate for long-term preservation
- We cannot ensure those data are still usable
  - The software for exploiting those data is under control of the experiments
- We are sure most of those data are not (easily) accessible !



# ▶ What we do with online/nearline data



- Migration of data on disk and under control of the MSS is under the site's responsibility
- We help experiments detect unused data
  - Both on disk and on tape (under control of the mass storage system)
  - This is done more or less regularly and often triggered by changes in technology
    - *For instance, introduction of more capacitive tapes*
- Experiments then decide what to do with the data
  - Remove them, archive them, keep them, do nothing

# ▶ What we do with online/nearline data (cont.)



- We help astro-particle physics experiments to make available their data to the community
  - Several ways: copying the data to an external repository, making the data available through a web site hosted by us, ...
- The experiment is responsible for preserving the data left by users whose account is disabled/closed
  - For instance, PhD students leaving the experiment/labs
  - Although the process is not sufficiently formalized

# ▶ Final Remarks



- Currently, we don't really preserve data, we preserve cartridges (!)
  - No strategy for long term data preservation is defined
  - A collective and coordinated effort from experiments, funding agencies and data centres seems essential for dealing with this issue
- The amount of data being collected by today's experiments, and the high degree of distribution (as a consequence of using grid technology) make data preservation of current experiments a big challenge
- We would be happy to contribute to whatever initiative the HEP community takes on research data preservation
- We would be happy to learn from other centres' good practices on this topic

# ▶ Questions/Comments



# ▶ Acknowledgements



Thanks to Suzanne Poulat and Philippe Olivero for providing material for this presentation.