

# Status Report of the DPHEP Collaboration: A Global Effort for Sustainable Data Preservation in High Energy Physics

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

## Abstract

Data from high energy physics (HEP) experiments are collected with significant financial and human effort and are mostly unique. An inter-experimental study group on HEP data preservation and long-term analysis was convened as a panel of the International Committee for Future Accelerators (ICFA). The group was formed by large collider-based experiments and investigated the technical and organisational aspects of HEP data preservation. An intermediate report was released in November 2009 addressing the general issues of data preservation in HEP and an extended blueprint paper was published<sup>1</sup> in 2012. In July 2014 the DPHEP collaboration was formed as a result of the signature of the Collaboration Agreement by seven large funding agencies (others have since joined or are in the process of acquisition) and in June 2015 the first DPHEP Collaboration Workshop<sup>2</sup> and Collaboration Board meeting took place.

This status report of the DPHEP collaboration details the progress over the past 3 years.



International Collaboration for Data Preservation  
and Long Term Analysis in High Energy Physics

---

<sup>1</sup> See <http://arxiv.org/pdf/1205.4667>.

<sup>2</sup> See <https://indico.cern.ch/event/377026/other-view?view=standard>.

DRAFT

EXECUTIVE SUMMARY .....	5
INTRODUCTION .....	6
THE DPHEP STUDY GROUP .....	7
THE DPHEP COLLABORATION AGREEMENT.....	8
THE DPHEP 2020 VISION .....	9
REQUIREMENTS FROM FUNDING AGENCIES.....	9
OPEN ACCESS POLICIES .....	10
DPHEP PORTAL .....	11
THE CERN GREY BOOK.....	11
THE DPHEP COLLABORATION AND IMPLEMENTATION BOARDS.....	14
USE CASES, COST MODELS AND BUSINESS CASES.....	14
BIT PRESERVATION AND STORAGE TECHNOLOGY OUTLOOK.....	19
STATUS AND ROADMAP OF CERNVM (CHEP ABSTRACT) .....	25
<b>DOCUMENTATION AND DIGITAL LIBRARY TECHNOLOGIES.....</b>	<b>25</b>
CERN PROGRAM LIBRARY DOCUMENTATION AND SOFTWARE .....	26
OPEN DATA AND DATA ANALYSIS PRESERVATION SERVICES FOR LHC EXPERIMENTS [ CHEP ABSTRACT ].....	27
HEP SOFTWARE FOUNDATION .....	28
RELATED PROJECTS, DISCIPLINES AND INITIATIVES .....	29
EU FP7 Projects .....	30
<b>CERN OPEN DATA PORTAL .....</b>	<b>33</b>
CERTIFICATION OF DIGITAL REPOSITORIES .....	33

<b>SITE / EXPERIMENT STATUS REPORTS (JUNE 2015)</b>	<b>36</b>
Belle I & II	36
BES III	36
HERA	37
LEP	38
Tevatron	38
BaBar	39
IPP	40
LHC	40
<b>TOWARDS A DATA PRESERVATION STRATEGY FOR CERN EXPERIMENTS</b>	<b>42</b>
<b>CHANGES WITH RESPECT TO THE BLUEPRINT</b>	<b>43</b>
<b>LESSONS FOR FUTURE CIRCULAR COLLIDERS / EXPERIMENTS</b>	<b>43</b>
<b>FUTURE ACTIVITIES</b>	<b>44</b>
<b>OUTLOOK AND CONCLUSIONS</b>	<b>44</b>
<b>APPENDIX A – THE DPHEP COLLABORATION</b>	<b>47</b>
<b>APPENDIX B – THE DPHEP IMPLEMENTATION BOARD</b>	<b>48</b>

## Executive Summary

- Significant progress has been made in the past years regarding our understanding of, and implementation of services and solutions for, long-term data preservation for future re-use;
- **However, continued investment in data preservation is needed: without this the data will soon become unusable or indeed lost (as history has told us all too many times);**
- Funding agencies – and indeed the general public – are now understanding the need for preservation and sharing of “data” (which typically includes significant metadata, software and “knowledge”) with requirements on data management plans, preservation of data, reproducibility of results and sharing of data and results becoming increasingly important and in some cases mandatory;
- The “business case” for data preservation in scientific, educational and cultural as well as financial terms is increasingly well understood: funding beyond (or outside) the standard lifetime of projects is required to ensure this preservation;
- A well-established model for data preservation exists – the Open Archival Information System (OAIS). Whilst developed primarily in the Space Data Community, it has since been adopted by all most all disciplines – ranging from Science to Humanities and Digital Cultural Heritage – and provides useful terminology and guidance that has proven applicable also to HEP;
- **The main message – from Past and Present Circular Colliders to Future ones – is that it is never early to consider data preservation: early planning is likely to result in cost savings that may be significant. Furthermore, resources (and budget) beyond the data-taking lifetime of the projects must be foreseen from the beginning.**

## Introduction

Shortly after the publication of the DPHEP Blueprint (see below), various inputs concerning the long-term preservation of HEP data were made to the group preparing the update to the European Strategy for Particle Physics. An updated strategy was adopted by a special session<sup>3</sup> of the CERN Council in May 2013 in Brussels, and this<sup>4</sup> includes the following statement:

The success of particle physics experiments, such as those required for the high-luminosity LHC, relies on innovative instrumentation, state-of-the-art infrastructures and large-scale data-intensive computing. *Detector R&D programmes should be supported strongly at CERN, national institutes, laboratories and universities. Infrastructure and engineering capabilities for the R&D programme and construction of large detectors, as well as infrastructures for data analysis, **data preservation** and distributed data-intensive computing should be maintained and further developed.*

As of 2013, with the appointment by CERN of a DPHEP Project Manager – one of the priorities identified in the Blueprint – the first steps towards transitioning to a Collaboration began. Seven institutes signed the Collaboration Agreement of May 2014, with additional (and often active) partners preparing to join.

After numerous workshops organized by and involving the Study Group, topical workshops on the “Full Costs of Curation” (January 2014)<sup>5</sup> and on “Common Projects and Shared Use Cases” (June 2015)<sup>6</sup> have been held.

The former has been instrumental in ensuring medium to long-term funding for the data preservation resources needed by the LHC experiments, whereas several CERN groups have committed support and services needed for the primary Use Cases agreed by these experiments (see below), which in many cases is matched by effort from the experiments and/or external institutes.

**The message that constant effort and investment is needed should not be lost. However this effort can be well justified by the measurable benefits. These include not only direct benefits to the (sometimes former) collaboration in terms of scientific papers and PhDs obtained, but also in terms of much needed publicity for HEP through educational outreach and “open access” activities.**

Future events where data preservation experiences and solutions can be shared will continue, as well as topical events as needs arise. (An event<sup>7</sup> is planned in conjunction with WLCG in Lisbon in February 2016, to prepare a detailed Data Preservation Plan following the OAIS and related standards.

---

<sup>3</sup> See <https://indico.cern.ch/event/244974/page/1>.

<sup>4</sup> See <https://indico.cern.ch/event/244974/page/7>.

<sup>5</sup> See <https://indico.cern.ch/event/276820/>.

<sup>6</sup> See <https://indico.cern.ch/event/377026/>.

<sup>7</sup> See <http://indico.cern.ch/event/433164/>.

## The DPHEP Study Group

The DPHEP study group was initiated in early 2009 and became a sub-group of the International Committee for Future Accelerators (ICFA) – emphasizing its global nature – later that year. Its goal was:

**High Energy Physics** experiments initiate with this **Study Group**<sup>8</sup> a common reflection on **data persistency and long-term analysis** in order to get a common vision on these issues and create a multi-experiment dynamics for further reference.

**The objectives of the Study Group are:**

- Review and document the physics objectives of the data persistency in HEP.
- Exchange information concerning the analysis model: abstraction, software, documentation etc. and identify coherence points.
- Address the hardware and software persistency status.
- Review possible funding programs and other related international initiatives.
- Converge to a common set of specifications in a document that will constitute the basis for future collaborations.

As well as running a series of workshops that rotated around all of the main HEP laboratories, it generated a Blueprint document that was well received by ICFA and was fed into the process for updating the European Strategy for Particle Physics.

The full Blueprint – which runs close to 100 pages – should be referred to for details regarding the motivation for and status of data preservation activities across all key laboratories (status in 2012).

It states:

*“Data from high-energy physics (HEP) experiments are collected with significant financial and human effort and are mostly unique. An inter-experimental study group on HEP data preservation and long-term analysis was convened as a panel of the International Committee for Future Accelerators (ICFA). The group was formed by large collider-based experiments and investigated the technical and organisational aspects of HEP data preservation. An intermediate report was released in November 2009 addressing the general issues of data preservation in HEP. This paper includes and extends the intermediate report. It provides an analysis of the research case for data preservation and a detailed description of the various projects at experiment, laboratory and international levels. In addition, the paper provides a concrete proposal for an international organisation in charge of the data management and policies in high-energy physics.”*

The DPHEP study group identified the following priorities, in order of urgency:

- **Priority 1: Experiment Level Projects in Data Preservation.** Large laboratories should define and establish data preservation projects in order to avoid catastrophic loss of data once major collaborations come to an end. The recent expertise gained during the last three years indicate that an extension of the computing effort within experiments with a person-power of the order of 2-3 FTEs leads to a significant improvement in the ability to move to a long-term data

---

<sup>8</sup> See <http://dphep.org> for further details.

*preservation phase. Such initiatives exist already or are being defined in the participating laboratories and are followed attentively by the study group.*

- **Priority 2: International Organisation DPHEP.** *The efforts are best exploited by a common organisation at the international level. The installation of this body, to be based on the existing ICFA study group, requires a Project Manager (1 FTE) to be employed as soon as possible. The effort is a joint request of the study group and could be assumed by rotation among the participating laboratories.*
- **Priority 3: Common R&D projects.** *Common requirements on data preservation are likely to evolve into inter-experimental R&D projects (three concrete examples are given above, each involving 1-2 dedicated FTE, across several laboratories). The projects will optimise the development effort and have the potential to improve the degree of standardisation in HEP computing in the longer term. Concrete requests will be formulated in common by the experiments to the funding agencies and the activity of these projects will be steered by the DPHEP organisation.*

*These priorities could be enacted with a funding model implying synergies from the three regions (Europe, America, Asia) and strong connections with laboratories hosting the data samples.*

## **The DPHEP Collaboration Agreement**

In order to implement priority 2 above (experiment-level data preservation is already under way in most cases and common “R&D” projects are already leading to services with a view to long-term support and sustainability), CERN has appointed a Project Manager (October 2012) and a Collaboration Agreement has been prepared. 9 institutes have now signed this agreement (CERN, DESY, HIP Finland, IHEP, IN2P3, KEK, MPP, IPP and STFC<sup>9</sup>) with several more in the pipeline.

The agreement, which largely reflects the recommendations of the Blueprint, includes the following goals:

*The Project, in coordination with the International Committee for Future Accelerators (ICFA), aims at:*

1. *Positioning itself as the natural forum for the entire discipline in order to foster discussion, achieve consensus and transfer knowledge in two main areas:*
  - a. *Technological challenges in data preservation in HEP,*
  - b. *Diverse governance at the collaboration and community level for preserved data,*
2. *Co-ordinate common R&D projects aiming to establish common, discipline-wide preservation tools,*

---

<sup>9</sup> Not yet formally ratified by a DPHEP Collaboration Board meeting.



3. *Harmonize preservation projects across the Partners and liaise with relevant initiatives from other fields,*
4. *Design the long-term organization of sustainable and economic preservation in HEP,*
5. *Outreach within the community and advocacy towards the main stakeholders for the case of preservation in HEP.*

All of these areas are currently being pursued actively and can be viewed in terms of a (slowly evolving) “2020 vision”.

## The DPHEP 2020 Vision

The “vision” for DPHEP – first presented to ICFA in February 2013 – consists of the following key points:

- By 2020, all **archived data** – e.g. that described in DPHEP Blueprint, including LHC data – should be easily **findable** and fully **usable** by the **designated communities** with clear (Open) access policies and possibilities to annotate further
- Best practices, tools and services should be well run-in, fully documented and sustainable; built in common with **other disciplines**, based on standards
- There should be a **DPHEP portal**, through which data / tools accessed
- **Clear targets & metrics** to measure the above should be agreed between **Funding Agencies, Service Providers** and the **Experiments (Collaborations)**.

Although there is clearly much work still to be done, this vision looks both achievable and the timescale for realizing it has been significantly reduced through interactions with other (non-HEP) projects and communities.

## Requirements from Funding Agencies

There have been numerous policy discussions and recommendations in recent years, some of which are reflected in the outputs of the (EU FP7) projects discussed below. A particularly clear statement can be found from the US Office of Science<sup>10</sup> that includes the following:

*All proposals submitted to the Office of Science (after 1 October 2014) for research funding must include a Data Management Plan (DMP) that addresses the following requirements:*

- *DMPs should describe whether and how data generated in the course of the proposed research will be shared and preserved.*

*If the plan is not to share and/or preserve certain data, then the plan must explain the basis of the decision*

<sup>10</sup> See <http://science.energy.gov/funding-opportunities/digital-data-management/>.

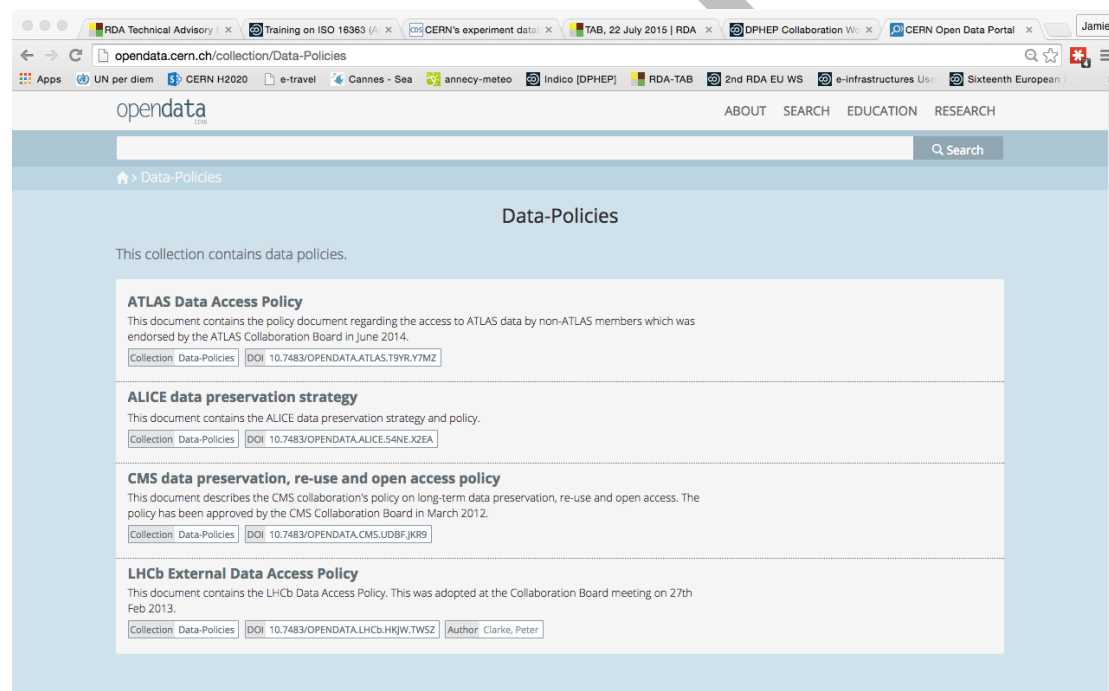
*At a minimum, DMPs must describe how data sharing and preservation will enable validation of results, or how results could be validated if data are not shared or preserved*

Similar requirements are coming (or have come) from other Funding Agencies and for International projects in particular it will be important to understand how to respond to these in a consistent manner. This is part of the debate that will continue, e.g. following the RECODE project recommendations covered below.

## Open Access Policies

The four main LHC experiments have approved Open Access policies<sup>11</sup> that, whilst they differ in detail, are broadly similar (and are being adopted by other experiments):

1. (Moving towards) Gold Open Access for Publications (DPHEP “level 1”);
2. Open Access to Specific Data Samples for Outreach (DPHEP “level 2”);
3. Open Access to (a fraction of the) Reconstructed data (after an embargo period) (DPHEP “level 3”);
4. Raw data<sup>12</sup> closed even to collaboration (today) (DPHEP “level 4”).



The “fractions” involved vary from 30-50% after a few years to (in some cases) 100% after ~10 years. (Just under 40TB of CMS data from 2010 have been released and 10TB of ALICE pp and PbPb data are expected to be released shortly. LHCb will release their first data only in 2018. For ATLAS, the plans are still unclear, but a volume similar to CMS or ALICE can be expected).

<sup>11</sup> See <http://opendata.cern.ch/collection/data-policies>.

<sup>12</sup> Most disciplines use a different notation, with “L0” corresponding to the raw data and L1/L2/L3 corresponding to calibrated and/or processed and/or derived data.

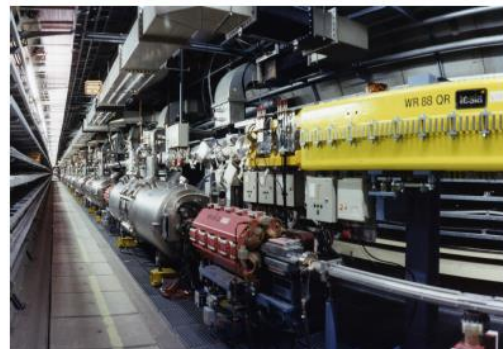
**Even though this applies to the reconstructed data, the volumes involved could end up being very significant and the technical and financial issues, particularly in the medium to long term (2020+) are not yet understood!**

## DPHEP Portal

First proposed in 2013, the initial idea was to federate the data preservation portals of the various laboratories and institutes involved, providing information on the experiments, data access and release policies, search capabilities and so forth. A much simplified and pragmatic approach is now being followed that can be embellished with additional capabilities as manpower allows – in particular for current and future experiments. A simple template is used to provide an overview of the experiment(s) and corresponding accelerator / collider and host laboratory, with drill-down to (largely existing) further detail as needed.

## HERA

- HERA was the largest particle accelerator at DESY
- It was the first internationally funded accelerator project and the joined effort of 11 countries
- Started in 1992, the storage ring served the international particle physics community for over 15 years
- The HERA experiments H1, ZEUS, and HERMES finished data taking in 2007 (Hera-B data taking ended in 2003)
- Up to now – and for the foreseeable future – no other electron-proton accelerator has explored electron-proton interaction at higher energies  
→ **Unique dataset**



Information on the data, documentation and software is provided with a standard look and feel, although details are expected to vary.

## The CERN Grey Book

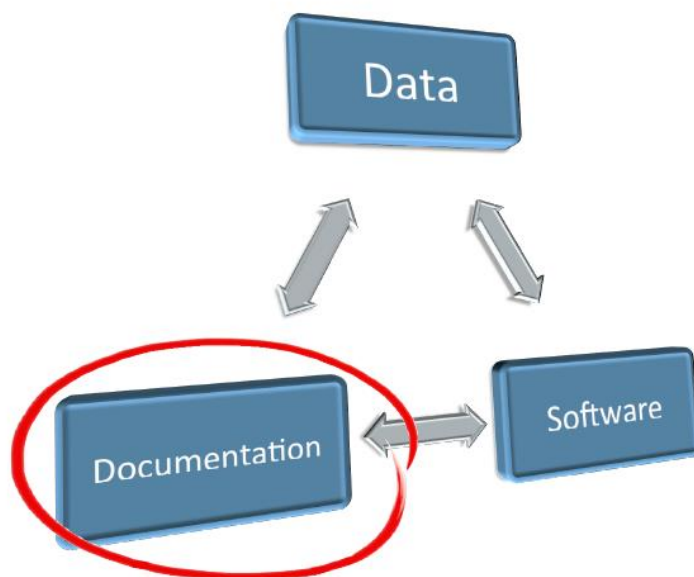
The list of experiments at CERN was published annually [from 1975](#) to 1999 in a printed version of the so-called Grey Book. Since the year 2000 CERN's experimental programme and projects are summarized electronically in the Grey

Book database. The information that the Grey Book contains for a given experiment includes:

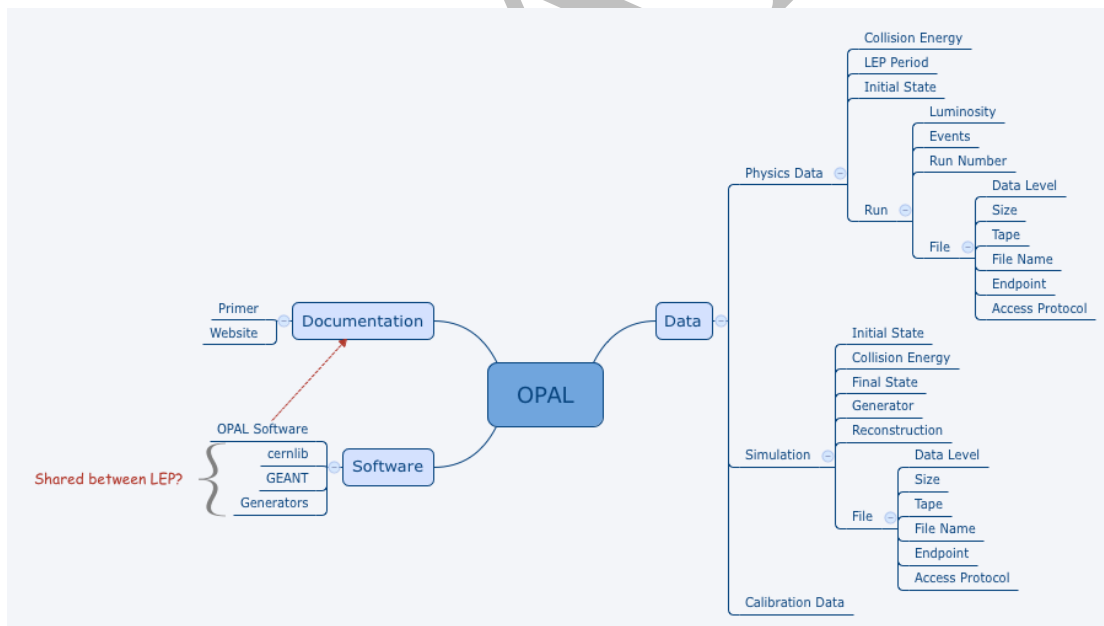
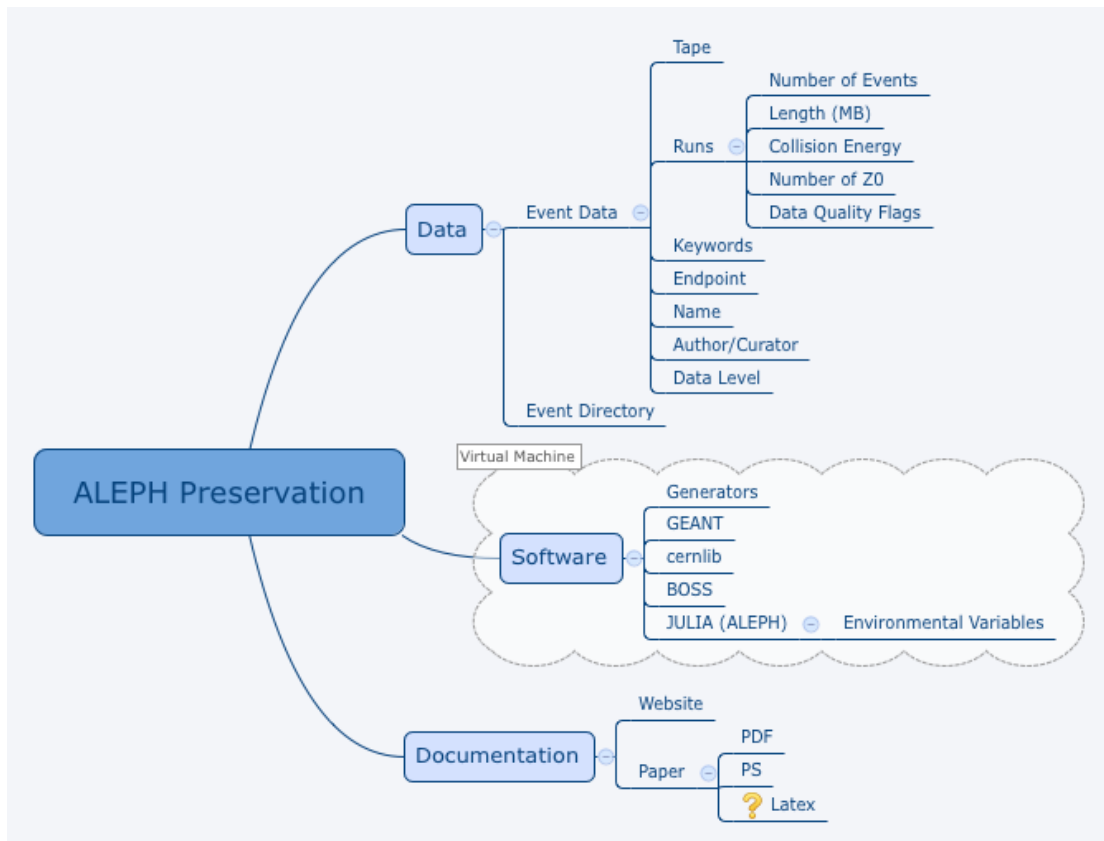
- A link to the experiments' Website;
- A pointer to the corresponding entry in the CERN Document Server (CDS) in the collection "Experiments at CERN";
- A similar pointer to the CDS collections "Committee Documents" and "Published Articles".

To link the Grey Book to the DPHEP portal (and vice-versa), an additional pointer will be added to point to the Data, Documentation and Software page in the DPHEP portal and a corresponding pointer (e.g. a "grey book" icon) between a given experiment's entry in the DPHEP portal to the Grey Book.

## Aspects of Data Preservation



Much of the information should be stable over time, with status reports (e.g. at DPHEP workshops, probably not more than annually) and updates to "HowTos" (updated for e.g. every new operating system release that is supported, changes in data access protocols etc. – hopefully less frequently but probably at least every 3-5 years) being the obvious exceptions.



For programmes such as those at the LHC, links to the analysis capture portal (for those authorised, i.e. collaboration members) and to the open data portal would additionally be provided. Links to external maintained sites – such as the active work on ALEPH data in INFN, that on OPAL data at the Max Planck Institute – would also fit naturally but not disturb the common look and feel.

## The DPHEP Collaboration and Implementation Boards

The DPHEP Collaboration Board (see Appendix A) consists of a representative of all of the institutes / bodies that have signed the DPHEP Collaboration Agreement. One meeting has been held so far, immediately after the DPHEP Collaboration Workshop of June 2015 and future meetings will be held approximately annually. The meetings are also open to future members of the Collaboration, as well as representatives from key projects such as DASPOS<sup>13</sup>.

The DPHEP Implementation Board (see Appendix B) meets more regularly and is composed of active participants in data preservation for HEP. The meeting frequency has dropped somewhat with time, given the relative maturity of a number of the data preservation projects as well as the occurrence of focused meetings on specific topics / technologies, such as CernVM, analysis capture and preservation and so forth.

The agendas of the meetings, as well as any material presented, can be found via the Indico category: <https://indico.cern.ch/category/4458/>. A web archive of the corresponding mailing list (requires authentication) can be found at <https://groups.cern.ch/group/DPHEP-IB/default.aspx>.

## Use Cases, Cost Models and Business Cases

Following numerous discussions, a set of common Use Cases has been agreed across the 4 main LHC experiments. With some small provisos, these are also valid for other experiments, including those reported on later in this document.

The basic Use Cases are as follows:

1. Bit preservation as a basic “service” on which higher level components can build;
  - Motivation: Data taken by the experiments should be preserved
2. Preserve data, software, and know-how<sup>14</sup> in the collaborations;
  - Foundation for long-term DP strategy
  - Analysis reproducibility: Data preservation alongside software evolution
3. Share data and associated software with (larger) scientific community
  - Additional requirements:
  - Storage, distributed computing
  - Accessibility issues, intellectual property
  - Formalising and simplifying data format and analysis procedure
  - Documentation
4. Open access to reduced data set to general public
  - Education and outreach
  - Continuous effort to provide meaningful examples and demonstrations

In general, Open Access is not currently considered for pre-LHC experiments that have well defined Open Access Policies. Furthermore, the “designated community”

---

<sup>13</sup> Data and Software Preservation for Open Science – see <https://daspos.crc.nd.edu/>.

<sup>14</sup> Additional Use Cases – not yet fully tested – help to define whether the “know-how” has been adequately captured. See the Analysis Capture section for further details.



(in OAIS terminology) is typically the (former) collaboration – although there is often considerable flexibility<sup>15</sup> in interpreting this restriction.

These Use Cases map well onto requirements now coming from Funding Agencies for data preservation, sharing and reproducibility. However, it is clear that we will have to work with them to understand and agree on what is technically possible, financially affordable and scientifically meaningful in this area.

A detailed cost model approximating<sup>16</sup> to that for LHC data shows that there is a significant upfront investment that drops rapidly with time. It is based on certain parameters, such as the use of Enterprise tape drives and media for the archive store, together with regular repacking to new, higher density media as this becomes available.

A very simple model that loosely matches the expected evolution in acquired LHC data volumes is shown below.

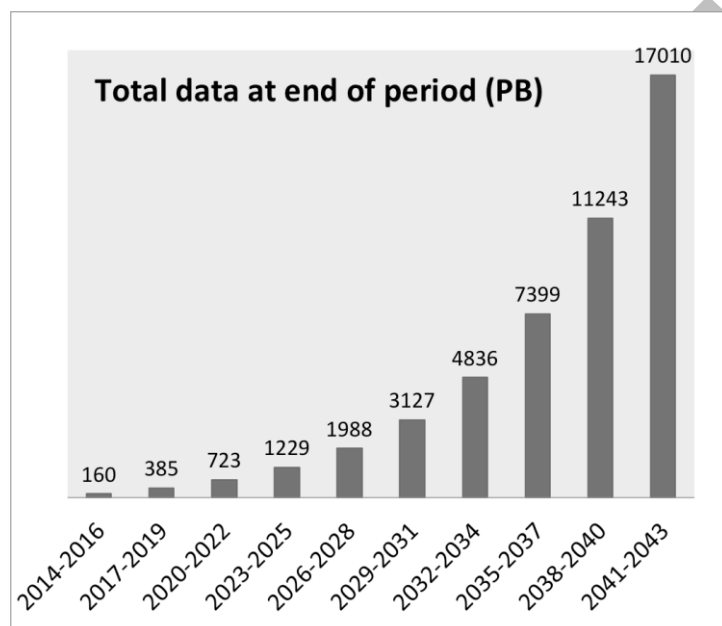


Figure 1 - Approximation to Evolution of LHC Storage for Cost Model

Based on publically available technology predictions and pricing information, we are able to calculate how much it would cost to store a single copy this information in a set of tape libraries (a 10% disk cache is included, as is a 3-year cycle for the media, after which all data are migrated forward to the next generation).

<sup>15</sup> In some cases it is sufficient to join the collaboration (typically by sending an e-mail to the Spokesperson); in others at least one former member of the collaboration must sign any papers and/or an appropriate disclaimer must be included).

<sup>16</sup> The cost model uses publically available pricing information and is thus suitable for sharing with other communities.

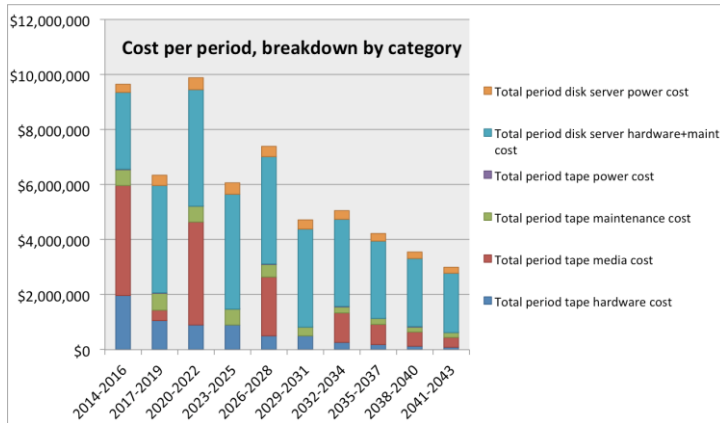


Figure 2 - Breakdown of Costs According to Storage Growth / Media Replacement

Not only does this – together with on-going data scrubbing – implement some of the key practice in OAIS and the associated certification procedures and hence allow us to offer “state-of-the-art” bit preservation, but we can also calculate the costs of a data store rising from several tens of PB initially to a few EB in the 2030s. Whilst much more detailed calculations are used in the LHC (WLCG) budget review and request process, this gives us at least a ballpark estimate for the costs involved and we can see that the cost over time averages to “just” \$2M / year (for such a vast and growing data store).

Comparatively, e.g. versus the cost of LHC computing, the cost of building and running the machine and its detectors, this is a “small number” – certainly much less than the cost of building a new machine in the future (at least with today’s technology)!

The “value” of the preserved data can be measured indirectly by the number of on-going analyses, publications and / or major conference presentations, as shown in the figures below for CDF, D0 and BaBar. These all show that there is significant activity that continues well after the end of data taking.



PRL + PRD + PLB + Eur (1988 2015-June-02)

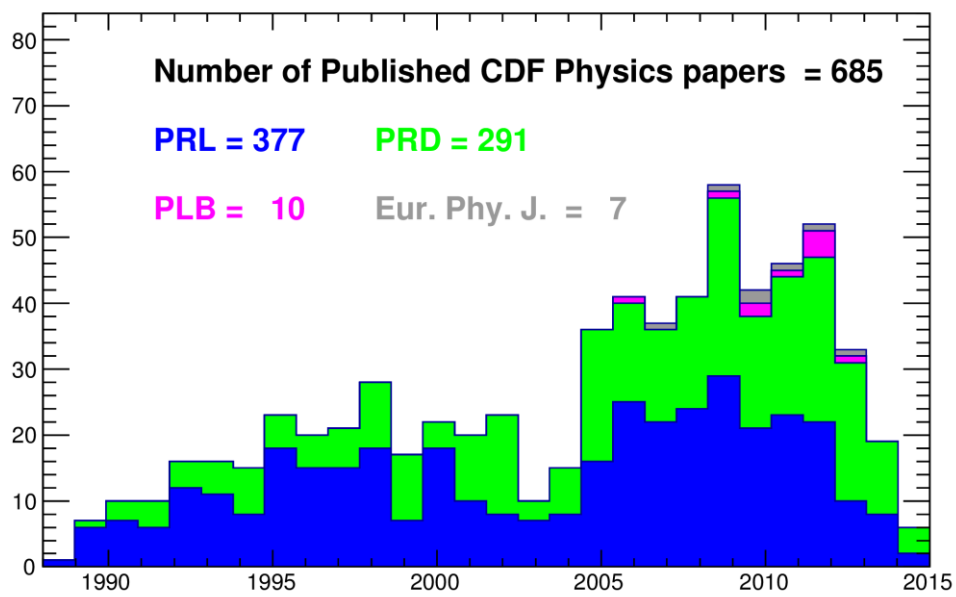


Figure 3 - Published Papers for the CDF Collaboration

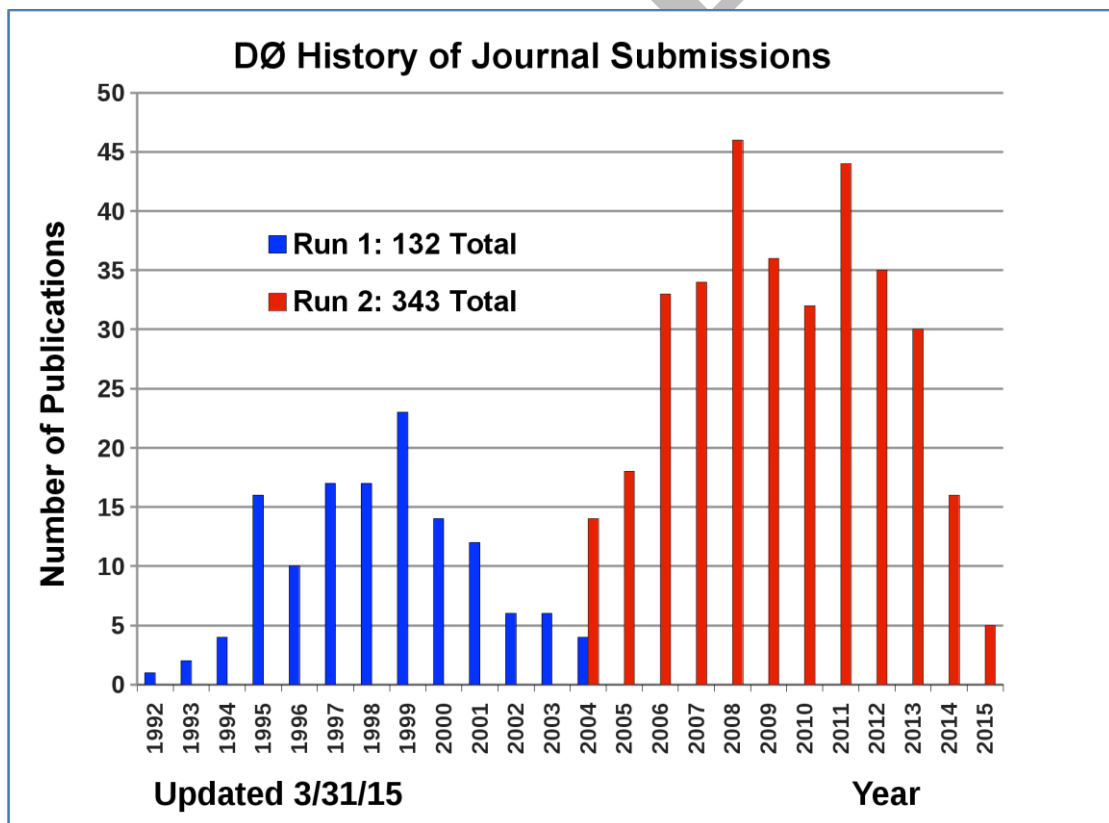


Figure 4 - Journal Submissions for the D0 Experiment

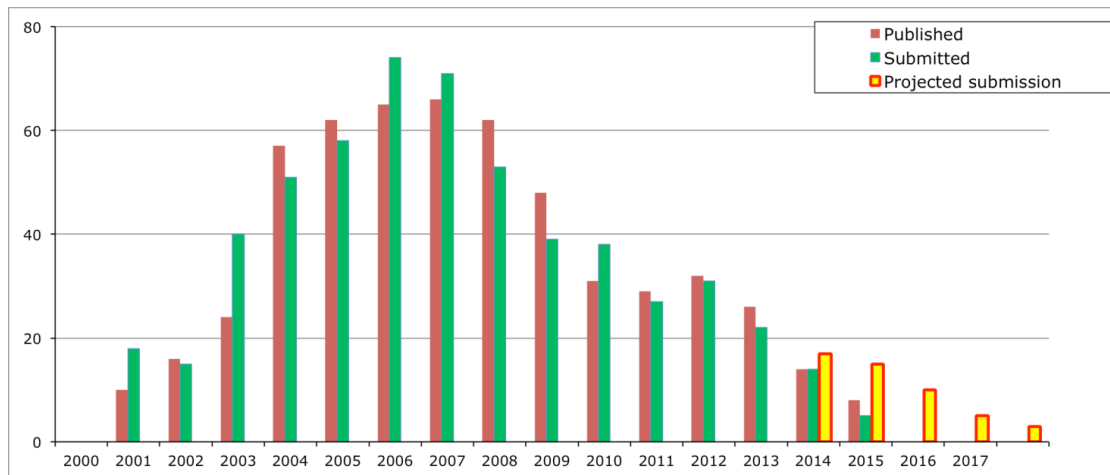


Figure 5 - Papers Submitted / Published for BaBar

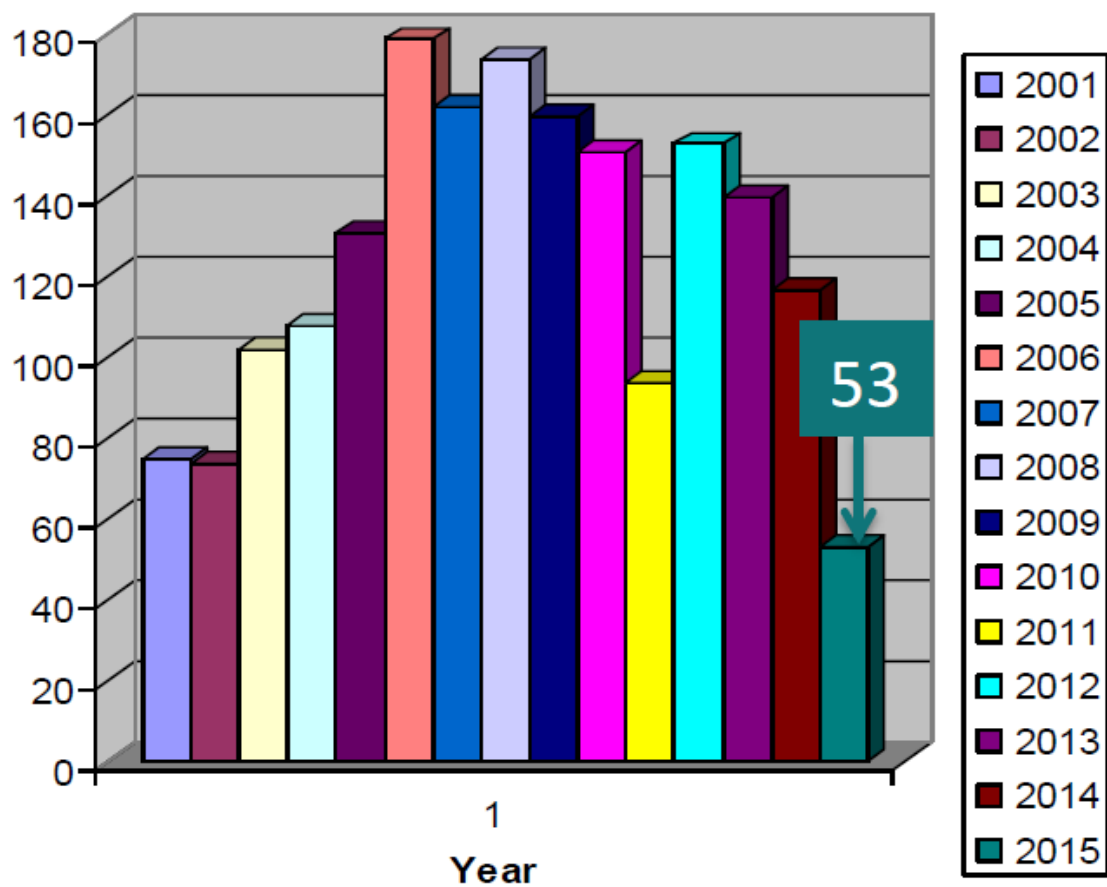


Figure 6 - Conference Talks by Year for BaBar

## Bit Preservation and Storage Technology Outlook

Bit preservation is an art in itself and – following the 4C project recommendations (see below) – is best performed at a limited number of “expert” sites, rather than across a multitude of smaller ones. This becomes even more important as densities increase – whereas user manipulation of individual tape volumes was common place in the LEP era, the latest generations of media requiring extreme clean-room conditions and prefer robots over humans!

The following graph shows the growth in data stored at CERN for the LHC and other experiments.

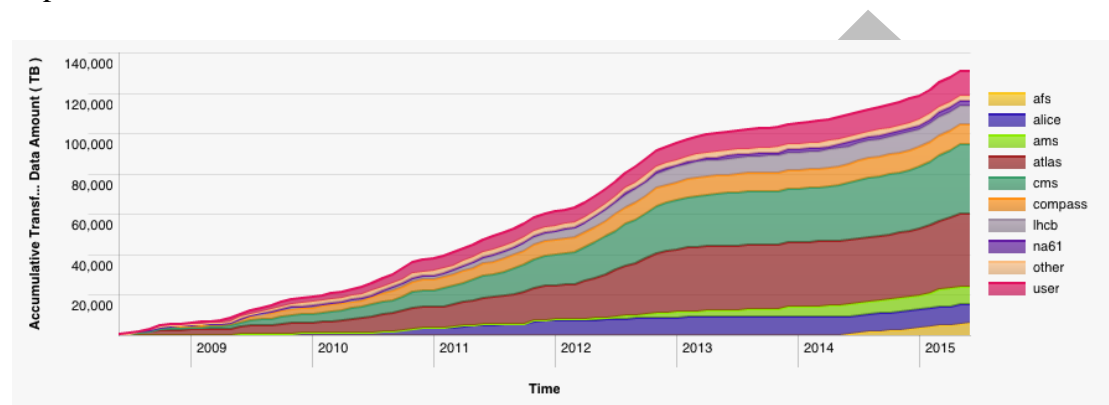


Figure 7 - Growth of Experiment (and AFS) data stored at CERN

Based on the anticipated data rates and volumes at LHC Run2 and future running periods, we predict a total data volume of a few EB (exabytes) in the 2030s.

Industry predictions (see below) suggest that cartridge capacity can continue to grow, at least over the next few years.

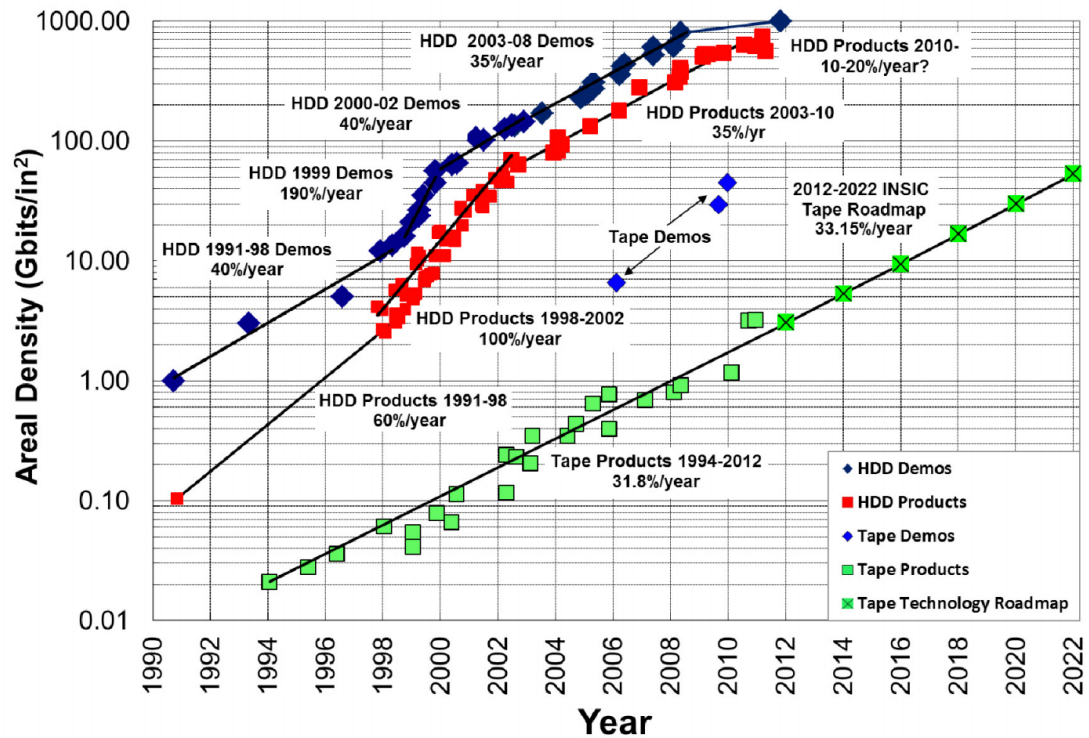


Figure 8 - Expected Evolution in Disk and Tape Technology

However, one has to inject a word of caution here – the tape market is shrinking, the source of enterprise drives has become a duopoly, with but a single supplier of high-density media.

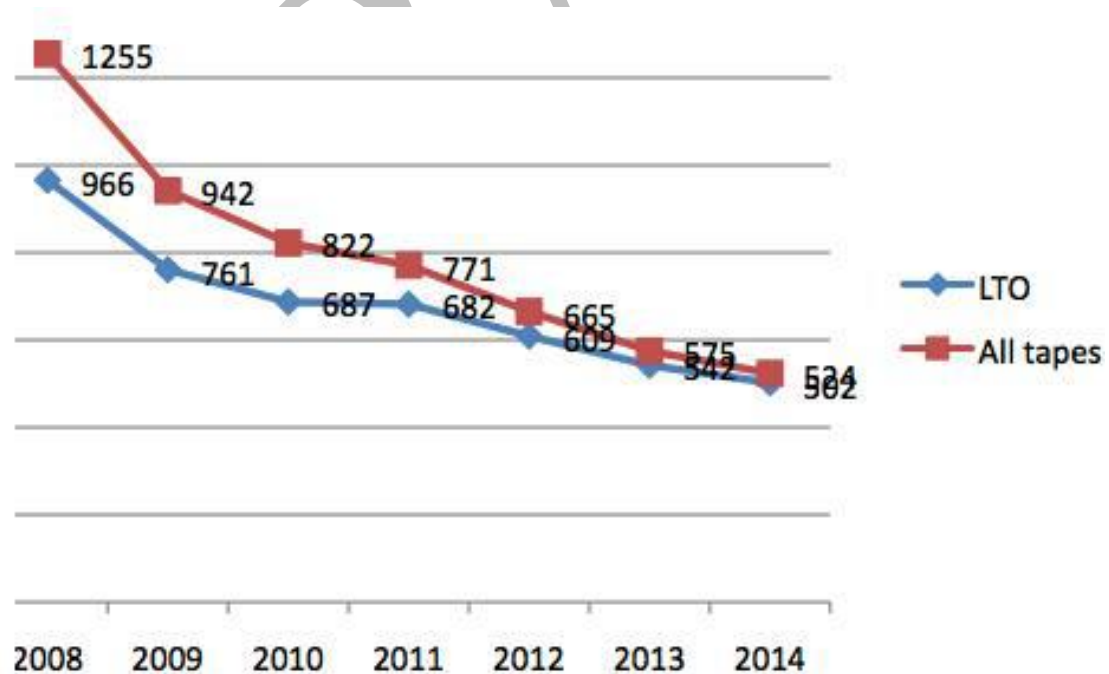


Figure 9 - Evolution of Tape Market

Some experts – such as David Rosenthal<sup>17</sup> - predict that Kryder's law (the “equivalent” of Moore's law for storage) will no longer hold true in the future. He warns that we should expect to pay more for storage. There have been many predictions of storage revolutions in the past – often involving optical or holographic storage. However, these have so far failed to materialise.

Looking back, we see a significant improvement in storage capacity with the number of cartridges required to storage all LEP data shrinking to an almost negligible number. So much so that now two tape copies are maintained at CERN, with a further read-only disk copy being setup. Given the additional copies maintained at a number of outside institutes for at least some of the LEP experiments, we have achieved a significant level of redundancy. Will this one day be true also for the LHC experiments?

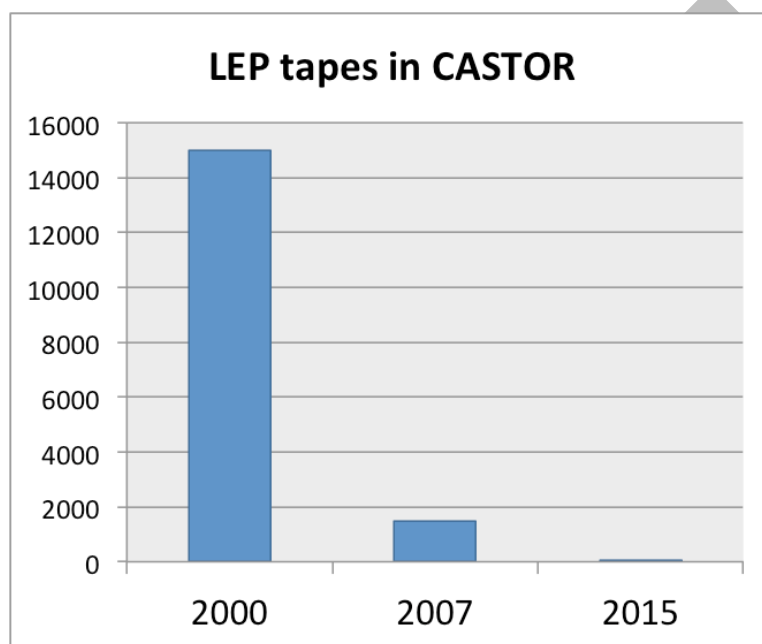
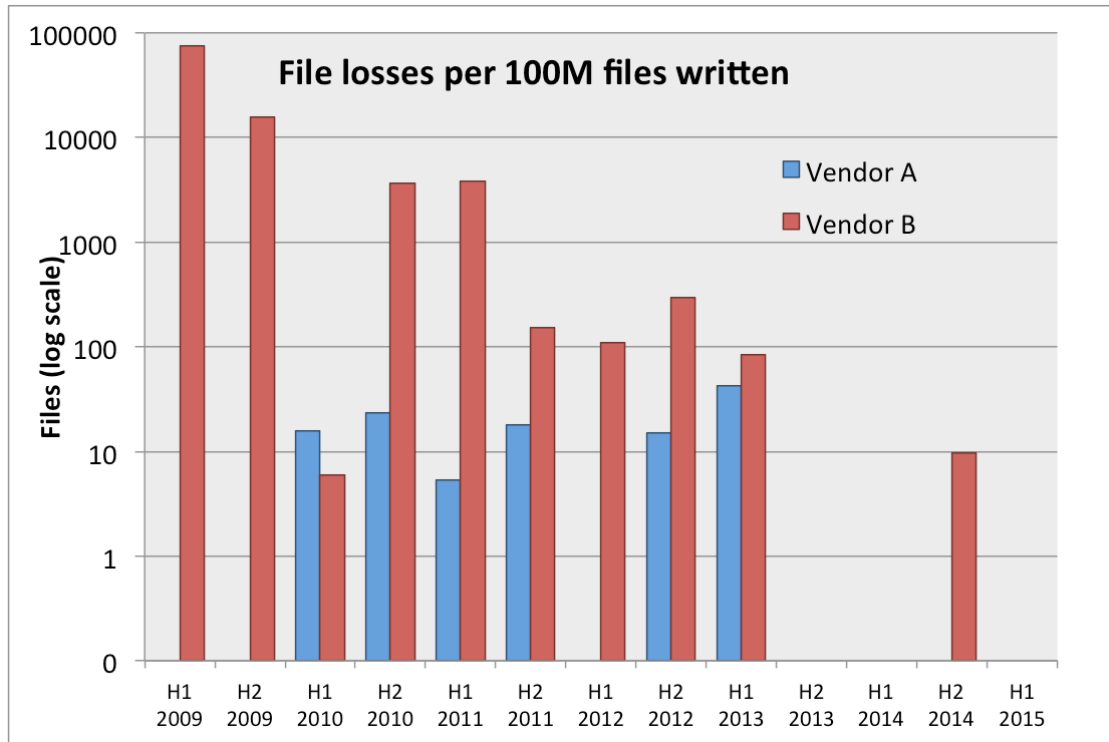


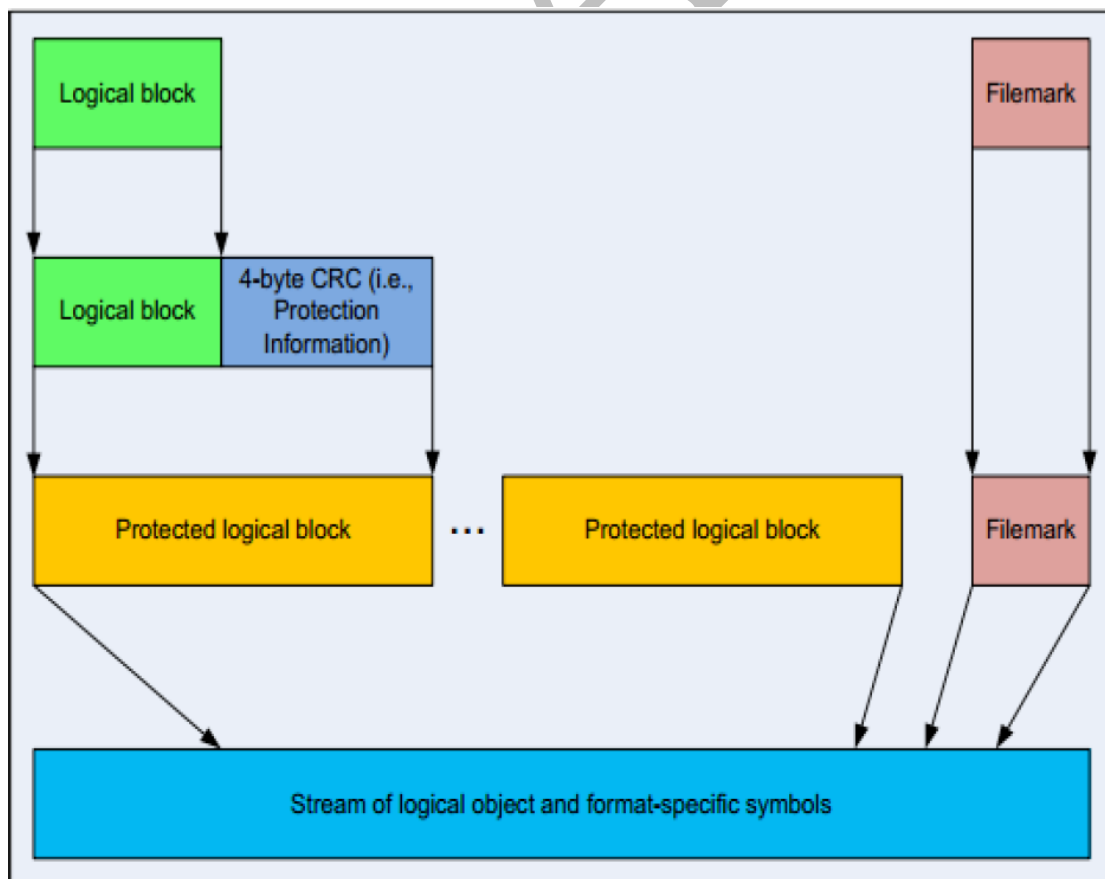
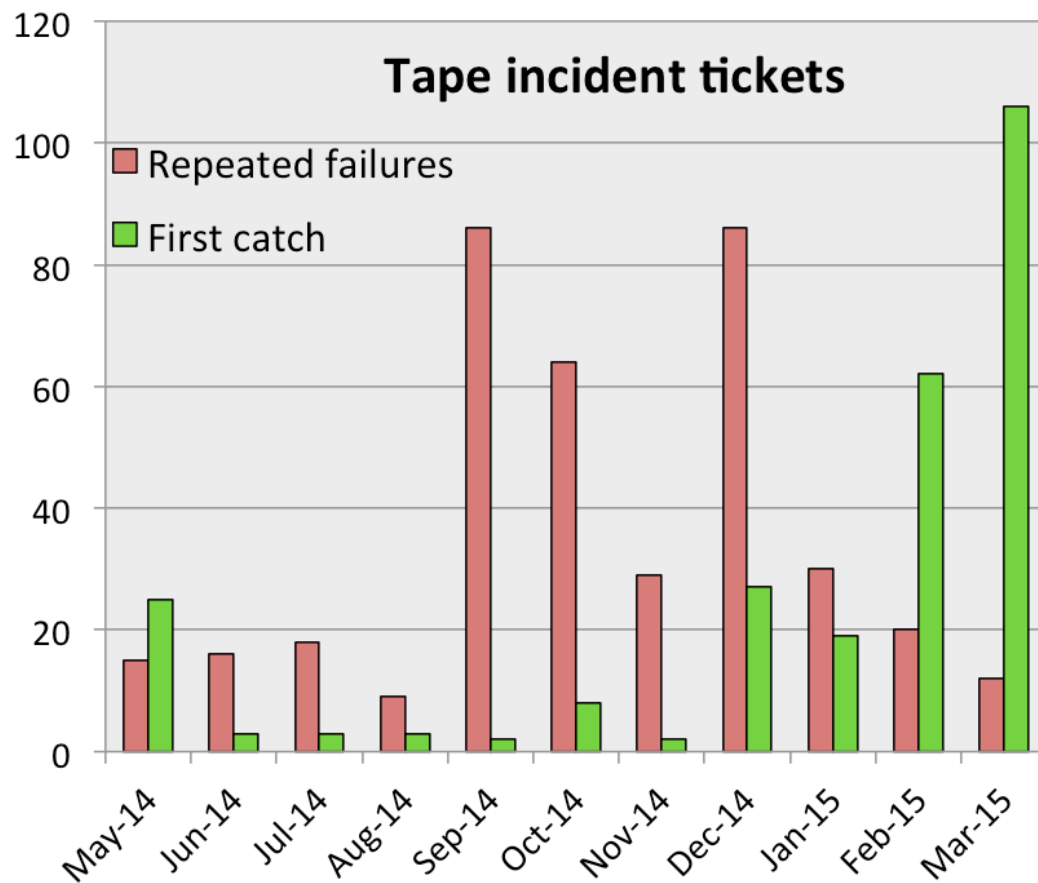
Figure 10 - Number of Tapes Required to Store LEP data

We can also see a measurable improvement in terms of reliability at the level of a single site. Additional levels of protection are foreseen, hopefully reducing data loss further still.

---

<sup>17</sup> See <http://blog.dshr.org/>.





This work can be compared to the recommendations of the National (i.e. US) Digital Stewardship Alliance (NDSA), shown in the table below.

	Level 1 (Protect your data)	Level 2 (Know your data)	Level 3 (Monitor your data)	Level 4 (Repair your data)
Storage and Geographic Location	<ul style="list-style-type: none"> <li>Two complete copies that are not collocated</li> <li>For data on heterogeneous media (optical discs, hard drives, etc.) get the content off the medium and into your storage system</li> </ul>	<ul style="list-style-type: none"> <li>At least three complete copies</li> <li>At least one copy in a different geographic location</li> <li>Document your storage system(s) and storage media and what you need to use them</li> </ul>	<ul style="list-style-type: none"> <li>At least one copy in a geographic location with a different disaster threat</li> <li>Obsolescence monitoring process for your storage system(s) and media</li> </ul>	<ul style="list-style-type: none"> <li>At least three copies in geographic locations with different disaster threats</li> <li>Have a comprehensive plan in place that will keep files and metadata on currently accessible media or systems</li> </ul>
File Fixity and Data Integrity	<ul style="list-style-type: none"> <li>Check file fixity on ingest if it has been provided with the content</li> <li>Create fixity info if it wasn't provided with the content</li> </ul>	<ul style="list-style-type: none"> <li>Check fixity on all ingests</li> <li>Use write-blockers when working with original media</li> <li>Virus-check high risk content</li> </ul>	<ul style="list-style-type: none"> <li>Check fixity of content at fixed intervals</li> <li>Maintain logs of fixity info; supply audit on demand</li> <li>Ability to detect corrupt data</li> <li>Virus-check all content</li> </ul>	<ul style="list-style-type: none"> <li>Check fixity of all content in response to specific events or activities</li> <li>Ability to replace/repair corrupted data</li> <li>Ensure no one person has write access to all copies</li> </ul>
Information Security	<ul style="list-style-type: none"> <li>Identify who has read, write, move and delete authorization to individual files</li> <li>Restrict who has those authorizations to individual files</li> </ul>	<ul style="list-style-type: none"> <li>Document access restrictions for content</li> </ul>	<ul style="list-style-type: none"> <li>Maintain logs of who performed what actions on files, including deletions and preservation actions</li> </ul>	<ul style="list-style-type: none"> <li>Perform audit of logs</li> </ul>
Metadata	<ul style="list-style-type: none"> <li>Inventory of content and its storage location</li> <li>Ensure backup and non-collocation of inventory</li> </ul>	<ul style="list-style-type: none"> <li>Store administrative metadata</li> <li>Store transformative metadata and log events</li> </ul>	<ul style="list-style-type: none"> <li>Store standard technical and descriptive metadata</li> </ul>	<ul style="list-style-type: none"> <li>Store standard preservation metadata</li> </ul>
File Formats	<ul style="list-style-type: none"> <li>When you can give input into the creation of digital files encourage use of a limited set of known open formats and codecs</li> </ul>	<ul style="list-style-type: none"> <li>Inventory of file formats in use</li> </ul>	<ul style="list-style-type: none"> <li>Monitor file format obsolescence issues</li> </ul>	<ul style="list-style-type: none"> <li>Perform format migrations, emulation and similar activities as needed</li> </ul>

Figure 11 - NDSA Levels of Digital Preservation



## Status and Roadmap of CernVM (CHEP abstract)

*Cloud resources nowadays contribute an essential share of resources for computing in high-energy physics. Such resources can be either provided by private or public IaaS clouds (e.g. OpenStack, Amazon EC2, Google Compute Engine) or by volunteers' computers (e.g. LHC@Home 2.0). In any case, experiments need to prepare a virtual machine image that provides the execution environment for the physics application at hand. The CernVM virtual machine since version 3 is a minimal and versatile virtual machine image capable of booting different operating systems. The virtual machine image is less than 20 megabyte in size. The actual operating system is delivered on demand by the CernVM File System. CernVM 3 has matured from a prototype to a production environment. It is used, for instance, to run LHC applications in the cloud, to tune event generators using a network of volunteer computers, and as a container for the historic Scientific Linux 5 and Scientific Linux 4 based software environments in the course of long-term data preservation efforts of the ALICE, CMS, and ALEPH experiments. We present experience and lessons learned from the use of CernVM at scale. We also provide an outlook on the upcoming developments. These developments include adding support for Scientific Linux 7, the use of container virtualization, such as provided by Docker, and the streamlining of virtual machine contextualization towards the cloud-init industry standard.*

## Documentation and Digital Library Technologies

Invenio is a [free software](#) suite enabling you to run your own [digital library](#) or document repository on the web. The technology offered by the software covers all aspects of digital library management from document ingestion through classification, indexing, and curation to dissemination. Invenio complies with standards such as the [Open Archives Initiative](#) metadata harvesting protocol (OAI-PMH) and uses [MARC 21](#) as its underlying bibliographic format. The flexibility and performance of Invenio make it a comprehensive solution for management of document repositories of moderate to large sizes (several millions of records).

Invenio has been originally developed at [CERN](#) to run the [CERN document server](#), managing over 1,000,000 bibliographic records in high-energy physics since 2002, covering articles, books, journals, photos, videos, and more. Invenio is being co-developed by an international collaboration comprising institutes such as [CERN](#), [DESY](#), [EPFL](#), [FNAL](#), [SLAC](#) and is being used by about thirty scientific institutions worldwide.

CERN, DESY, Fermilab and SLAC have built the next-generation High Energy Physics (HEP) information system, [INSPIRE](#). It combines the successful [SPIRES database content](#), curated at [DESY](#), [Fermilab](#) and [SLAC](#), with the [Invenio](#) digital library technology developed at [CERN](#). INSPIRE is run by a collaboration of CERN, DESY, Fermilab, IHEP, and SLAC, and interacts closely with HEP publishers, [arXiv.org](#), [NASA-ADS](#), [PDG](#), [HEPDATA](#) and other information resources.

INSPIRE represents a natural evolution of scholarly communication, built on successful community-based information systems, and provides a vision for information management in other fields of science.

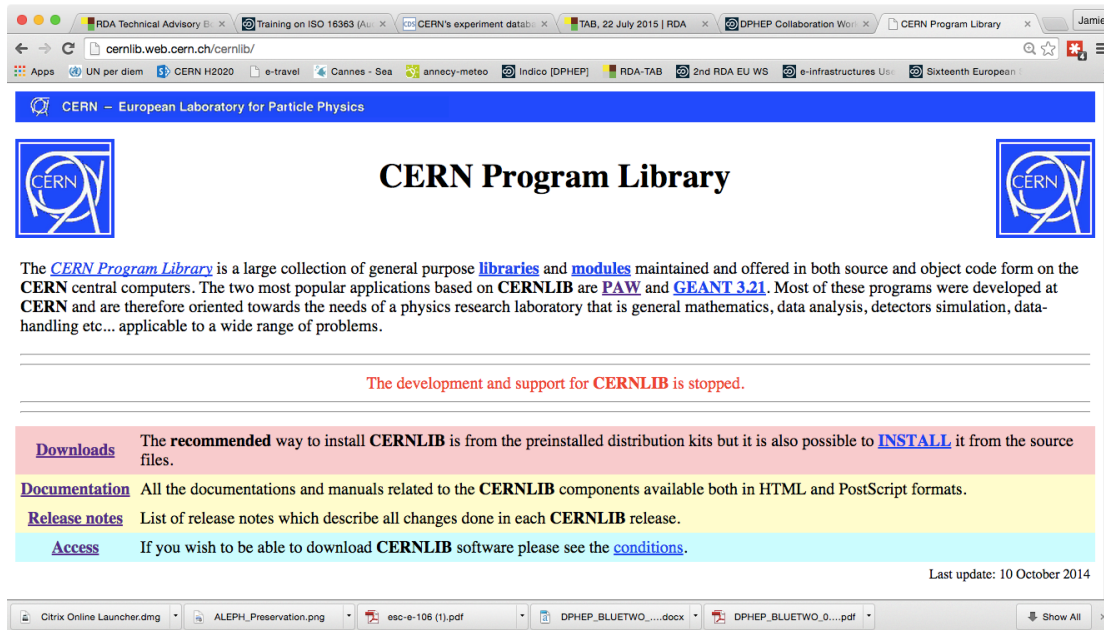
INSPIREHEP, the CERN Document Server, and other Invenio-based services make a good match for long-term preservation needs. However, during its active lifetime, experiments make use of many other systems, including Wikis, newsgroups, e-log books, Web pages and so forth. These are well suited to on-going production with relatively frequent updates but are not really aimed at long-term stability and preservation. The solution adopted is typically to “snap-shot” the information in this dynamic systems although more work – perhaps as a “common project” – would be very beneficial in this area. This is by no means a HEP-specific problem: it applies equally to all disciplines.

## **CERN Program Library Documentation and Software**

Many HEP experiments rely to a greater or lesser degree on the set of libraries known collectively as the “CERN Program Library” or simply CERNLIB. The documentation for these was last revised in the mid-1990s with the sources, marked up in LaTeX, stored at CERN in /afs.

In order to best preserve the documentation for the medium to long term, the following activities are currently underway:

1. Reformatting of the source files to produce PDF and/or PDF/A files with the latest fonts;
2. Capturing of the author and paper information, storing of the formatted files in the CERN Document Server using identifiers to refer to the authors and papers;
3. Addition of further meta-data to enable more powerful searches;
4. Agreement on a new “home” for the files - /afs has had a long life at CERN (some 20-25 years) but is targeted for replacement. Can we find an alternative that will be as long-lived?



Formal support for CERNLIB ceased over a decade ago – and development earlier still. However, it continues to be actively used in “data preservation” and re-use activities. Porting to future versions of Linux and an “official” version that the (past) experiments can trust is still desirable.

## Open Data and Data Analysis Preservation Services for LHC Experiments [ CHEP abstract ]

*In this paper we present newly launched services for open data and for long-term preservation and reuse of high-energy-physics data analyses. We follow the "data continuum" practices through several progressive data analysis phases up to the final publication. The aim is to capture all digital assets and associated knowledge inherent in the data analysis process for subsequent generations, and to make a subset available rapidly to the public.*

*A data analysis preservation pilot study was launched in order to assess the usual workflow practices in LHC collaborations. Leveraging on synergies between ALICE, ATLAS, CMS and LHCb experiments, the analysed data was followed through various "analysis train" steps, from the initial capture and pre-selection of primary data, through several intermediate selection steps yielding more greatly reduced datasets, up to the final selection of N-tuples used for producing high-level plots appearing in scientific journal publications. Most of the analysis chain is kept strictly restricted within a given collaboration; only the final plots, presentations and papers are usually made public. It is therefore essential to handle access rights and embargo periods as part of the data life cycle.*

*The study revealed many similarities between collaborations, even though the variety of different practices existing in different groups within the collaborations make it hard to reproduce an analysis at a later time in a uniform way. One recurring problem underlined by the study was to ensure an efficient "knowledge capture" related to user code when the principal author of an analysis (e.g. a PhD student) leaves the collaboration later.*

*The pilot solution has been prototyped using the Invenio digital library platform which was extended with several data-handling capabilities. The aim was to preserve information about datasets, the underlying OS platform and the user software used to study it. The configuration parameters, the high-level physics information such as physics object selection, and any necessary documentation and discussions are optionally being recorded alongside the process as well. The metadata representation of captured assets uses the MARC bibliographic standard which had to be customised and extended in relation to specific analysis-related fields. The captured digital assets are being minted with Digital Object Identifiers, ensuring later referencability and citability of preserved data and software. Connectors were built in the platform to large-scale data storage systems (CASTOR, EOS, Ceph). In addition, to facilitate information exchange among concerned services, further connectors were built to*

*the internal information management systems of LHC experiments (e.g. CMS CADI), to the discussion platforms (e.g. TWiki, SharePoint), and to the final publication servers (e.g. CDS, INSPIRE) used in the process. Finally, the platform draws inspiration from the Open Archival Information System (OAIS) recommended practices in order to ensure long-term preservation of captured assets.*

*The ultimate goal of the analysis preservation platform is to capture enough information about the process in order to facilitate reproduction of an analysis even many years after its initial publication, permitting to extend the impact of preserved analyses through future revalidation and recasting services.*

*A related "open data" service was launched for the benefit of the general public. The LHC experimental collaborations are committed to make their data open after a certain embargo period. Moreover, the collaborations also release simplified datasets for the general public within the framework of the international particle physics masterclass program. The primary and reduced datasets that the collaborations release for public use are being collected within the CERN Open Data portal service, allowing any physicist or general data scientist to access, explore, and further study the data on their own.*

*The CERN Open Data portal offers several high-level tools which help to visualise and work with the data, such as an interactive event display permitting to visualise CMS detector events on portal web pages, or a basic histogram plotting interface permitting to create live plots out of CMS reduced datasets. The platform guides high-school teachers and students to online masterclasses to further explore the data and improve their knowledge of particle physics. A considerable part of the CERN Open Data portal was therefore devoted to attractive presentation and ease-of-use of captured data and associated information.*

*The CERN Open Data platform not only offers datasets and live tools to explore them, but it also preserves the software tools used to analyse the data. It notably offers the download of Virtual Machine images permitting users to start their own working environment in order to further explore the data; for this the platform uses CernVM based images prepared by the collaborations. Moreover, the CERN Open Data platform preserves examples of user analysis code, illustrating how the general public could write their own code to perform further analyses.*

## **HEP Software Foundation**

The HEP Software Foundation (HSF)<sup>18</sup> facilitates coordination and common efforts in high energy physics (HEP) software and computing internationally. The objectives of the HSF as a community-wide organization include

- Sharing expertise;
- Raising awareness of existing software and solutions;
- Catalyzing new common projects;
- Promoting commonality and collaboration in new developments to make the most of limited resources;
- Aiding developers and users in creating, discovering, using and sustaining common software;
- Supporting career development for software and computing specialists;
- Provide a framework for setting goals and priorities, and attracting effort and support;
- Facilitate wider connections with other science fields.

Although not directly related to data preservation, it is clearly a forum where long-term sustainability of software can be discussed and contacts have been established with this in mind.

## **Related Projects, Disciplines and Initiatives**

The European Alliance for Permanent Access (APA) was set up as a non-profit organization, initiated as a Foundation under Dutch Law in September 2008. The goal of the Alliance is to align and enhance permanent information infrastructures in Europe across all disciplines. It is a networking organisation and a sustainable centre for advice and expertise on permanent access. The Alliance brings together seventeen major European research laboratories, research funders, and research support organisations such as national libraries and publishers. All its members are stakeholders in the European infrastructure for long-term preservation of and access to the digital records of science. Through the alliance, they are articulating a shared vision for a sustainable digital information infrastructure providing permanent access to scientific information.

CERN has been a member of the APA for many years, as well as having a seat on the Executive Board. Unfortunately, due to an unfavourable audit of an FP7 project in which the APA was involved, it is in the process of dissolution.

The APA held regular conferences where CERN and other partners presented. It also played a key role in FP7 projects including APARSEN, SCIDIP-ES, PRELIDA and others.

Through contacts with the APA, the activities of DPHEP have become much more widely known outside HEP, as well as specific activities such as peta- to exa-scale bit preservation and cost modeling.

---

<sup>18</sup> See <http://hepsoftwarefoundation.org/>.

A series of joint workshops have been held between the above projects and DPHEP at several meetings of the Research Data Alliance (RDA).

Material from these joint workshops can be found through the DPHEP Indico pages: <https://indico.cern.ch/category/4458/>.

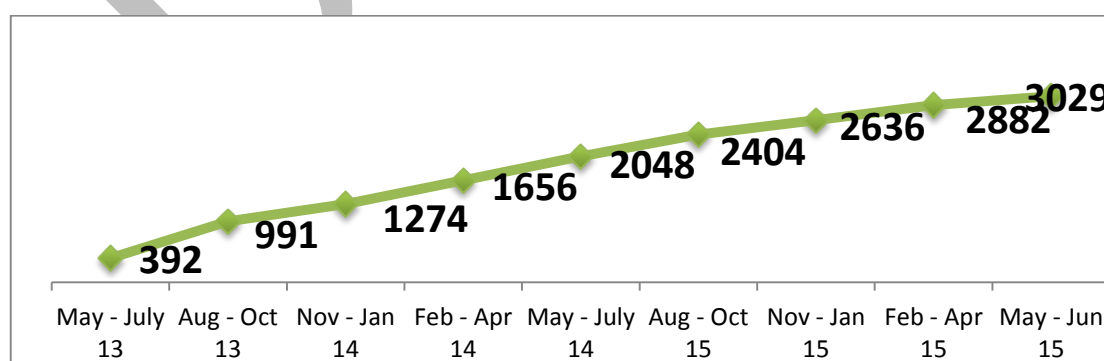
The RDA (<https://rd-alliance.org/node>) is now three years old, is supported by funding agencies in North America, Europe and Asia-Pacific and has a focus on data sharing and re-use. It holds two plenary meetings per year that include working meetings of numerous Working and Interest groups. (Working groups are supposed to deliver tangible outputs in around 18 months whereas interest groups are longer lived and are one mechanism by which working groups can be setup). Groups of particular interest to DPHEP include:

- The Preservation e-Infrastructure Interest Group;
- WG and IGs on (harmonization of) certification of digital repositories;
- Active Data Management;
- Reproducibility;
- Citation;
- And so forth.

As a networking event the RDA meetings can be particularly valuable and contacts made through the RDA as well as APA have helped establish and refine our vision, as well as providing channels whereby our activities can be more widely disseminated.

Furthermore, the RDA appears to be a central point for discussing “all things data” and clearly has the attention of the funding agencies. The RDA Europe arm of the “project” will likely be funded (in several stages) throughout the entire Horizon 2020 programme. CERN – on behalf of the EIROforum IT Working Group – is currently an Organizational Member and has a seat on the Technical Advisory Board (September 2013 – September 2015). It is also represented on the advisory board of the RDA Europe H2020 project.

The number of RDA members continues to grow roughly linearly, as the following graph shows.



## EU FP7 Projects

Two FP7 policy projects – 4C and RECODE – are also worthy of discussion, as summarized below.

## 4C and RECODE Policy Recommendations

4C was an FP7 project that terminated in January 2015 to help clarify the costs involved in data curation. Its goals were:

*“4C will help organisations across Europe to invest more effectively in digital curation and preservation. Research in digital preservation and curation has tended to emphasize the cost and complexity of the task in hand. 4C reminds us that the point of this investment is to realise a benefit, so our research must encompass related concepts such as ‘risk’, ‘value’, ‘quality’ and ‘sustainability’.”*

Its roadmap document<sup>19</sup> contains the following recommendations:

1. *Identify the value of digital assets and make choices;*
2. *Demand and choose more efficient systems;*
3. *Develop scalable services and infrastructure;*
4. *Design digital curation as a sustainable service;*
5. *Make funding dependent on costing digital assets across the whole lifecycle;*
6. *Be collaborative and transparent to drive down costs.*

With its leadership in providing scalable, sustainable services, HEP is well positioned to make key contributions in many of these areas. However, we must be aware of and plan for recommendation 5, which could have significant funding implications!

The Policy RECommendations for Open Access to Research Data in Europe (RECODE) project:

*“will leverage existing networks, communities and projects to address challenges within the open access and data dissemination and preservation sector and produce policy recommendations for open access to research data based on existing good practice.”*

As for 4C, this was also an FP7-funded project that recently terminated, again with a final set of policy recommendations.

As has happened with publications, the most likely course of events is that the Open Access to data movement will gain momentum. However, given the above-mentioned LHC policies and the volumes of data involved, we need to be prepared to answer the following questions:

1. Is it financially affordable?
2. Is it technically implementable?
3. Is it scientifically (or educationally, or culturally) meaningful?

---

<sup>19</sup> See <http://4cproject.eu/>.

The answers to these questions may well vary with time and also depend on the implementation(s) that we choose: Open Access is just one step in the progression towards Open Data and finally “Open Knowledge”.<sup>20</sup>

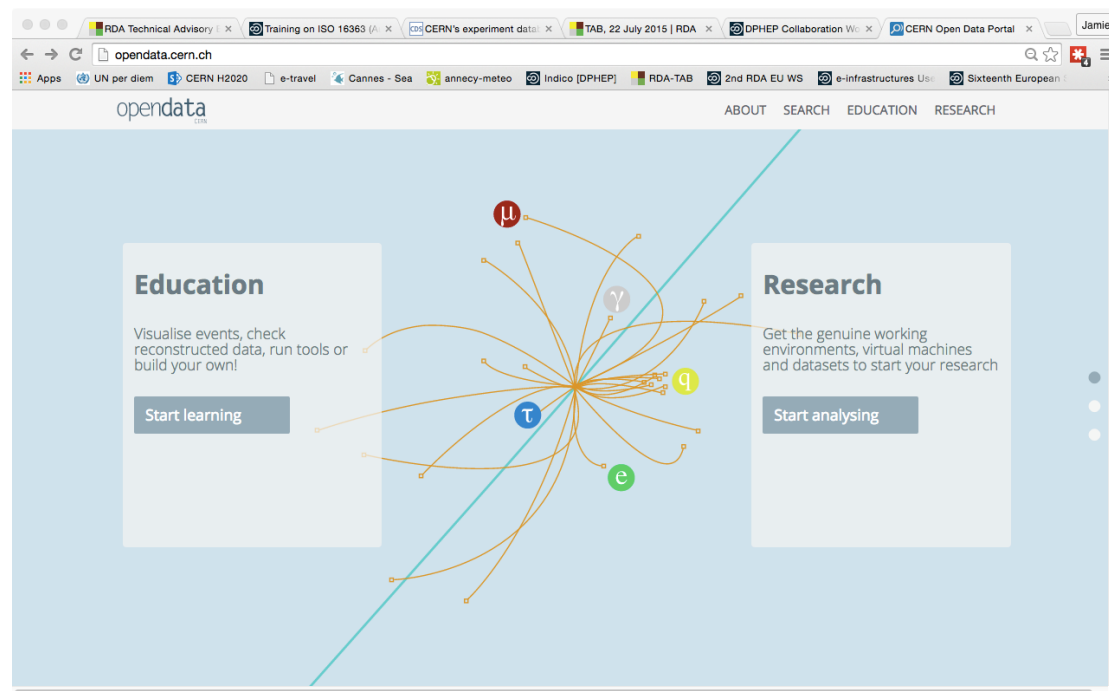
DRAFT

---

<sup>20</sup> An early but public draft of the Horizon 2020 2016-17 work programme states “*Research Infrastructures such as the ones on the ESFRI roadmap and others, are characterized by the very significant data volumes they generate and handle. These data are of interest to thousands of researchers across scientific disciplines and to other potential users via Open Access policies. Effective data preservation and open access for immediate and future sharing and re-use is a fundamental component of today’s research infrastructures and Horizon 2020 actions.*”



# CERN Open Data Portal



## Certification of Digital Repositories

Increasingly, the terms “trusted” or “certified” repositories are used: by data preservation projects, by communities requiring preservation services as well as by funding agencies in calls for project proposals. A number of methodologies exist – such as those from the Data Archiving and Networked Services (DANS) in the Netherlands, CODATA and finally a set of closely related ISO standards – that are in the process of being harmonised in the context of the RDA.

Following discussions at the WLCG Overview Board and following interest from the preservation community, a course was organised at CERN covering the ISO standards in this area, given by the authors of the standards involved.

There are three important ISO standards:

- **ISO 14721:2012** (OAIS – a reference model for what is required for an archive to provide long-term preservation of digital information)
- **ISO 16363:2013** (Audit and certification of trustworthy digital repositories – sets out comprehensive metrics for what an archive must do, based on OAIS)
- **ISO 16919:2014** (Requirements for bodies providing audit and certification of candidate trustworthy digital repositories – specifies the competencies and requirements on auditing bodies)

These three standards form a closely related family and an understanding of their principles and use will become increasingly important in establishing an internationally recognized set of trustworthy digital repositories.

Personnel followed this course from the WLCG Tier0 (CERN) and several WLCG Tier1 sites.

A checklist is available and it is foreseen – following further discussion in the WLCG and DPHEP communities – to proceed at least with a self-certification in 2016. This would help ensure that all of the necessary processes were in place, as well as identifying any gaps, for long-term preservation and re-use of HEP data. “Self-certification” discussions will form part of the DPHEP workshop that will be co-located with a WLCG workshop in Lisbon in February 2016.

Some of the metrics involved in obtaining certification are listed below.

<u>Metric</u>	<u>Supporting Text</u>	<u>Examples</u>
<b>3.1.1 THE REPOSITORY SHALL HAVE A MISSION STATEMENT THAT REFLECTS A COMMITMENT TO THE PRESERVATION OF, LONG TERM RETENTION OF, MANAGEMENT OF, AND ACCESS TO DIGITAL INFORMATION.</b>	This is necessary in order to ensure commitment to preservation and access at the repository’s highest administrative level.	Mission statement or charter of the repository or its parent organization that specifically addresses or implicitly calls for the preservation of information and/or other resources under its purview; a legal, statutory, or government regulatory mandate applicable to the repository that specifically addresses or implicitly requires the preservation of information and/or other resources under its purview.
<b>3.1.3 THE REPOSITORY SHALL HAVE A COLLECTION POLICY OR OTHER DOCUMENT THAT SPECIFIES THE TYPE OF INFORMATION IT WILL PRESERVE, RETAIN, MANAGE AND PROVIDE ACCESS TO.</b>	This is necessary in order that the repository has guidance on acquisition of digital content it will preserve, retain, manage and provide access to.	Collection policy and supporting documents, Preservation Policy, mission, goals and vision of the repository.
<b>3.2.1 THE REPOSITORY SHALL HAVE IDENTIFIED AND ESTABLISHED THE DUTIES THAT IT NEEDS TO PERFORM AND SHALL HAVE APPOINTED STAFF WITH ADEQUATE SKILLS AND EXPERIENCE TO FULFIL THESE DUTIES.</b>	Staffing of the repository should be by personnel with the required training and skills to carry out the activities of the repository. The repository should be able to document through development plans, organizational charts, job descriptions, and related policies and procedures that the repository is defining and maintaining the skills and roles that are required for the sustained operation of the repository.	Organizational charts; definitions of roles and responsibilities; comparison of staffing levels to industry benchmarks and standards.

<b>3.3.1</b>	<b>THE REPOSITORY SHALL HAVE DEFINED ITS DESIGNATED COMMUNITY AND ASSOCIATED KNOWLEDGE BASE(S) AND SHALL HAVE THESE DEFINITIONS APPROPRIATELY ACCESSIBLE.</b>	This is necessary in order that it is possible to test that the repository meets the needs of its Designated Community.	A written definition of the Designated Community.
<b>3.3.2</b>	<b>THE REPOSITORY SHALL HAVE PRESERVATION POLICIES IN PLACE TO ENSURE ITS PRESERVATION STRATEGIC PLAN WILL BE MET.</b>	This is necessary in order to ensure that the repository can fulfill that part of its mission related to preservation	Preservation Policies; Repository Mission Statement.
<b>4.1.1</b>	<b>THE REPOSITORY SHALL IDENTIFY THE CONTENT INFORMATION AND THE INFORMATION PROPERTIES THAT THE REPOSITORY WILL PRESERVE.</b>	This is necessary in order to make it clear to funders, depositors and users what responsibilities the repository is taking on and what aspects are excluded. It is also a necessary step in defining the information which is needed from the information producers or depositors.	Mission statement; submission agreements/deposit agreements/deeds of gift; workflow and Preservation Policy documents, including written definition of properties as agreed in the deposit agreement/deed of gift; written processing procedures; documentation of properties to be preserved.
<b>4.3.4</b>	<b>THE REPOSITORY SHALL PROVIDE EVIDENCE OF THE EFFECTIVENESS OF ITS PRESERVATION ACTIVITIES.</b>	This is necessary in order to assure the Designated Community that the repository will be able to make the information available and usable over the mid-to-long-term.	Collection of appropriate preservation metadata; proof of usability of randomly selected digital objects held within the system; demonstrable track record for retaining usable digital objects over time; Designated Community polls.
<b>5.1.1</b>	<b>THE REPOSITORY SHALL IDENTIFY AND MANAGE THE RISKS TO ITS PRESERVATION OPERATIONS AND GOALS ASSOCIATED WITH SYSTEM INFRASTRUCTURE.</b>	This is necessary to ensure a secure and trustworthy infrastructure.	Infrastructure inventory of system components; periodic technology assessments; estimates of system component lifetime; export of authentic records to an independent system; use of strongly community supported software .e.g., Apache, iRODS, Fedora); re-creation of archives from backups.

<b>5.2.1 THE REPOSITORY SHALL MAINTAIN A SYSTEMATIC ANALYSIS OF SECURITY RISK FACTORS ASSOCIATED WITH DATA, SYSTEMS, PERSONNEL, AND PHYSICAL PLANT.</b>	This is necessary to ensure ongoing and uninterrupted service to the designated community.	Repository employs the codes of practice found in the ISO 27000 series of standards system control list; risk, threat, or control analysis.
---	--	---

## Site / Experiment Status Reports (June 2015)

### Belle I & II

Preservation Aspect	Status
Bit Preservation	
Data	
Documentation	
Software	
Uses Case(s)	
Target Community(ies)	
Value	Quantitative measures (# papers, PhDs etc) exist
Uniqueness	
Resources	
Status	
Issues	
Outlook	

### BES III

Preservation Aspect	Status
Bit Preservation	A MD5 integrity check is done when data is copied from disk to tape Annual examination of tape library and LTO4 tapes (possibly moving to biennial due to risks to tapes)
Data	2750TB acquired 2009-2014 with annual growth of 450TB leading to 3450TB in 2020. Archive storage system based on CASTOR v1.8 with IBM3584 tape library, LTO 4 Current capacity for BESIII <ul style="list-style-type: none"> <li>• 2.7 PB, 2.2 PB used, 0.5 PB available</li> </ul> Remote replication of important raw data <ul style="list-style-type: none"> <li>• ~ 900 cartridges, 700 TB</li> </ul>
Documentation	<ul style="list-style-type: none"> <li>• DocDB: paper, technical notes, minutes...</li> <li>• Hypernews: notifications of software release, paper publishing ...</li> <li>• Indico: Conference slides,</li> <li>• Inspire: published paper</li> </ul>

<b>Software</b>	BOSS is an integrated software package that includes all the blocks required in BESIII data processing. For an old but stable version of BOSS, we preserve following items: <ul style="list-style-type: none"> <li>• A complete package of software,</li> <li>• A runnable virtual machine image</li> <li>• The puppet template and RPM repository from which a runnable OS is created,</li> <li>• Release documents, book-keeping parameters...</li> <li>• A functional validation is done according to the standard process of software release.</li> </ul>
<b>Uses Case(s)</b>	
<b>Target Community(ies)</b>	
<b>Value</b>	
<b>Uniqueness</b>	
<b>Resources</b>	Since the experiment is still working, budget and FTEs are shared with the operation of computing centre
<b>Status</b>	
<b>Issues</b>	
<b>Outlook</b>	The experiment is expected to stop data taking at 2022 and Lifespan of preserved data is expected to be about 15 years after then.

## HERA

<b>Preservation Aspect</b>	<b>Status</b>
<b>Bit Preservation</b>	
<b>Data</b>	Transferred to DPHEP area on DESY dCache. 2 tape copies (different media generations – 1.2 PiB) plus disk cache (700 TiB) for on-going analyses
<b>Documentation</b>	Non-digital documentation catalogued and stored in the DESY library archive; some digitized. Software notes in INSPIREHEP
<b>Software</b>	<i>“In the best of all worlds we would keep the software alive i.e. compilable on the latest Linux with the latest library versions”</i> We now follow a “freezing approach”, i.e. a VM with isolated storage and well defined set of external libs
<b>Uses Case(s)</b>	Continued analysis by former collaboration members
<b>Target Community(ies)</b>	Former collaboration
<b>Value</b>	Analyses, publications and PhDs continue to be produced
<b>Uniqueness</b>	Unique combination of initial state particles and energy
<b>Resources</b>	
<b>Status</b>	Transitioning (-ed) from experiment-specific to institutional solutions
<b>Issues</b>	Webservers: tension between production needs and long-term archiving.

	Do not underestimate the effort! Experiment expertise fades away quickly once funding stops. <b>Data preservation must be prepared whilst the collaboration exists and effort is available!</b>
<b>Outlook</b>	Continued ability to analyse data until 2020 (when support for SL6 stops); Migration to SL7 could extend this; Tape archive will life on.

## LEP

Preservation Aspect	Status
<b>Bit Preservation</b>	“State of the art” bit preservation with regular scrubbing and migration to new media
<b>Data</b>	2 copies on tape at CERN, an additional copy on disk (EOS) being setup. Additional copies exist outside CERN (ALEPH, OPAL and partial copy for DELPHI)
<b>Documentation</b>	Being revisited – to be “archived” in CERN Document Server for long-term preservation
<b>Software</b>	To be published into CernVMFS
<b>Uses Case(s)</b>	Continued analyses by former collaboration members
<b>Target Community(ies)</b>	Primarily former collaboration
<b>Value</b>	Analyses, publications and PhDs continue to be produced
<b>Uniqueness</b>	Unique – until and unless certain FCC options are implemented
<b>Resources</b>	Minimal resources for “bit preservation” and storage
<b>Status</b>	
<b>Issues</b>	Dependency on CERNLIB (no longer maintained)
<b>Outlook</b>	Expect to be able to analyse data (ALEPH, DELPHI, OPAL) until at least 2020. Until 2030 should be possible with < (<) 1FTE / experiment / yearß

## Tevatron

Preservation Aspect	Status
<b>Bit Preservation</b>	All data migrated to T10k technology (2 ½ years). Data integrity checks: After each copy during migration; Periodic reads from each tape. Long term future preservation of CDF data at INFN-CNAF, developed in collaboration with CDF and FNAL SCD.
<b>Data</b>	Two copies of raw data at FNAL, in different locations. In case of damage/loss analysis ntuples can be reproduced and/or eventually recovered from CNAF.
<b>Documentation</b>	All online webpages and code archived, still accessible from CDF webpages.

<b>Software</b>	<p>All online webpages and code archived, still accessible from CDF webpages.</p> <p>At the time of Tevatron shutdown</p> <ul style="list-style-type: none"> <li>• all code in frozen releases or in CVS repositories</li> <li>• based on 32-bit frameworks built on Scientific Linux 5 (but with compatibility libraries to older OSs)</li> </ul> <p>Long term future solution: build legacy release that contains no pre-SL6 libraries CVMFS for code distribution</p>
<b>Uses Case(s)</b>	Continued analyses by former collaboration members
<b>Target Community(ies)</b>	Primarily former collaboration
<b>Value</b>	Quantitative measures (# papers, PhDs etc) exist
<b>Uniqueness</b>	Unique initial state vs LHC; Multiple energy collisions (300, 900 and 1960 GeV)
<b>Resources</b>	FNAL R2DP project budgeted 4 (3) FTE in 2013, 3 (2.1) in 2014 and 0.3 (0.4) in 2015. (Expenditure)
<b>Status</b>	R2DP project complete
<b>Issues</b>	Both CDF and D0 use Oracle → licence cost is a long-term future challenge. Migration to open source db would require considerable human effort (need to rewrite the analysis software)
<b>Outlook</b>	Goal: Complete analysis capability (DPHEP “level 4”) through Nov 2020 (SL6 EOL) and beyond.

## BaBar

<b>Preservation Aspect</b>	<b>Status</b>
<b>Bit Preservation</b>	2.7PB of data of which 2PB (budget constraints) will be migrated to new media when supported by SLAC
<b>Data</b>	Data is stored on tape at SLAC and CC-IN2P3 (back-up only); Active data on disk accessed via xrootd.
<b>Documentation</b>	All the most used and fundamental information have been checked, updated and moved to a Media Wiki server, the BABAR WIKI
<b>Software</b>	
<b>Uses Case(s)</b>	Continued analyses by former collaboration members
<b>Target Community(ies)</b>	Primarily former collaboration
<b>Value</b>	Quantitative measures (# papers, PhDs etc) exist (>30 analyses on track for publication + ~20 with less clear future)
<b>Uniqueness</b>	Data will not be superseded by LHC – some by Belle II (not Y(3S))
<b>Resources</b>	0.35 FTE computing support for BaBar at SLAC by end 2015 + 1.55 FTE for data and user support
<b>Status</b>	
<b>Issues</b>	Much of the hardware is aging; Sun OS support will

	stop within 2 years and corresponding h/w be decommissioned
<b>Outlook</b>	Aim to preserve data for on-going analyses until 2018 with extension to 2020+ to match Belle II schedule. The technology at the base of the future operating model will be virtualization – all the services now running on physical hardware will soon run on virtual machines

## IPP

(not clear how to summarise in the following format – maybe just a text version of the talk?)

<b>Preservation Aspect</b>	<b>Status</b>
<b>Bit Preservation</b>	
<b>Data</b>	
<b>Documentation</b>	
<b>Software</b>	
<b>Uses Case(s)</b>	
<b>Target Community(ies)</b>	
<b>Value</b>	
<b>Uniqueness</b>	
<b>Resources</b>	
<b>Status</b>	
<b>Issues</b>	
<b>Outlook</b>	

## LHC

<b>Preservation Aspect</b>	<b>Status</b>
<b>Bit Preservation</b>	“State of the art” bit preservation with regular scrubbing and migration to new media
<b>Data</b>	Stored at WLCG Tier0 with additional copies across WLCG Tier1 sites
<b>Documentation</b>	
<b>Software</b>	“Published” into CernVMFS
<b>Uses Case(s)</b>	“Standard”
<b>Target Community(ies)</b>	Re-use of data within the collaboration(s), sharing with the wider scientific community, Open Access releases
<b>Value</b>	Landmark discoveries already made; significant potential for future “BSM” discoveries
<b>Uniqueness</b>	Unique data sets (both pp and HI) being acquired now - ~2035. Probably unique until “FCC” (2035-2050?)
<b>Resources</b>	Computing resources via Resource Review Board
<b>Status</b>	



<b>Issues</b>	Effort within the experiments is hard to find
<b>Outlook</b>	On-going activity on analysis capture and reproducibility. Regular public releases (according to individual experiment policies) and “master classes”

DRAFT

## Towards a Data Preservation Strategy for CERN Experiments

The updated Strategy for European Particle Physics<sup>21</sup>, approved by Council in May 2014, states that “*infrastructures for ... data preservation ... should be maintained and further developed.*”

In order to implement this strategy, the following proposals are currently under discussion. (The numbering reflects the draft proposal, where the paragraph above is point 1.):

2. Such infrastructures include *digital repositories*, where *copies* or *replicas* of the data are kept.
3. As host laboratory, it is expected that (from now on?) a copy of all data acquired by CERN experiments *and* targeted for long-term preservation be stored in the CERN digital repository. This will typically include all raw data and the final reprocessing pass and associated Monte Carlo datasets.
4. It is strongly recommended that one or more copies of the above data are maintained outside, at or spread over institutes that form part of the collaboration.
5. In order to ensure sufficient reliability and adherence to “best practices”, it is recommended that such repositories follow agreed guidelines / standards – this is currently being discussed in the context of WLCG for LHC data.
6. These guidelines not only include policies for the management of the repository itself, but also on access to data in the repository (adherence to agreed access policies and terms of use), as well as the *ingest* process, when data is “entered” into the repository. The latter is to ensure that appropriate and supported data formats are used, there is sufficient documentation, meta-data and other materials to permit use by the designated communities, and so forth.
7. The above recommendations could become part of a default strategy for CERN experiments, with implementation details – including variances on the above – provided in the Data Management Plan (DMP) for that experiment. DMPs are increasingly required by funding agencies for new and/or repeat funding and can be expected to be quasi-mandatory in the future.
8. As a minimum, the DMP of an experiment should detail the policy for storing replicas of data and the recovery mechanisms, both during and after the active lifetime of the associated collaboration.
9. These basic recommendations are expected to be supplemented by others – e.g. on “knowledge capture and preservation” – as we gain experience with preserved and open access data.

It is foreseen that this proposal will be discussed at CERN’s scientific committees – most likely starting with the LHCC, as an implementation based on the WLCG Tier0 and Tier1 sites could be a reality in the short to medium term.

---

<sup>21</sup> See <http://council.web.cern.ch/council/en/EuropeanStrategy/ESParticlePhysics.html>.

## Changes With Respect to the Blueprint

With respect to the DPHEP Blueprint, the following observations can be made:

- The pervasive use of INSPIREHEP and other Invenio-based solutions will come as no surprise;
- “Bit Preservation” (and loss) is more clearly defined, with extensive practical experience, albeit different implementations due to site preferences and requirements (hardware choices, funding schemes etc.);
- Virtualisation is more prominent with a better defined timeline (circa ten years);
- The use of CVMFS is a clear success story;
- Cost models and business cases are better understood, with quantitative measures across a variety of experiments;
- “Open Access” policies, embargo periods and the like are new but match well with the “*Zeitgeist*”;
- A variety of “end-of-life” scenarios have been realised: moving from experiment to site support, from host institution to former collaboration members and even porting to new systems and services, such as EUDAT.

These developments, as well as the concrete experience over the past three years, positions the DPHEP Collaboration well to make clear recommendations to future projects and experiments.

## Lessons for Future Circular Colliders / Experiments

**The main message – from Past and Present Circular Colliders to Future ones – is that it is never early to consider data preservation: early planning is likely to result in cost savings that may be significant. Furthermore, resources (and budget) beyond the data-taking lifetime of the projects must be foreseen from the beginning.**

Beyond that, the activities of numerous data preservation activities worldwide can be used as a guide to the type of activities, services and support that is required.

In other words, at least “observer status” from the FCC activities in the DPHEP Collaboration is to be strongly recommended.

For other future and / or current experiments the recommendations are similar:

- Align yourselves with the overall strategy and even implementation of other data preservation activities at your institute / laboratory or globally;
- Adopt mainstream and supported technologies where-ever possible;
- Understand the target communities for your data preservation activities, the Use Cases and the expected benefits and outcomes;
- Try to understand the costs – in particular those that are specific to your collaboration (and not “external” – e.g. host laboratory bit preservation services);

- Data preservation services and support for the LHC experiments can be expected to be provided for several decades: this may be a good place to start.

## Future Activities

Over the next period, one can expect progress to be made in the following areas:

- The establishment of a formal policy regarding data preservation for CERN experiments (perhaps linked to the approval process through the Research Board);
- At least a “self-audit” for the CERN Tier0 and WLCG Tier1 sites in the context of the WLCG project;
- Further developments in terms of Analysis Capture and Preservation;
- Further releases of Open Data through the CERN Open Data Portal;
- Harmonization of similar activities across various laboratories and projects;
- Extension of DPHEP’s activities to consider also those of potential FCCs;
- Clarifications regarding funding – of particular importance to past experiments where resources have already become sub-optimal;
- The continuation of regular meetings and workshops, aligning as much as possible with related events (WLCG, CHEP, HEP Software Foundation etc.);
- Further input to the next round of ESPP – building on concrete experience, results and remaining challenges.

The long-term management of the Collaboration also has to be considered – up to 2020 but also beyond.

## Outlook and Conclusions

There are clearly many similarities in the approaches being taken, the technologies deployed and the issues encountered. Regular reporting of results (possibly synchronised with major events such as CHEP) should be sufficient to ensure that coordinated approaches remain and that duplication is minimised.

The following quote<sup>22</sup> is traditionally attributed to Leslie Lamport – the initial author of LaTeX and an expert on distributed computing systems.

A distributed system is one in which the failure of a computer you didn't even know existed can render your own computer unusable.

This reminds us that data preservation is inherently unstable – with many components and dependencies, constant attention is required to ensure that the entire “system” remains usable. Some changes may be relatively minor, such as a name change in a webserver. Others can be much more disruptive, such as major change in operating system (think VAX/VMS to Unix) or programming language – even a standard-conforming language changes over time, with some constructs being first deprecated, then obsolete and finally unsupported.

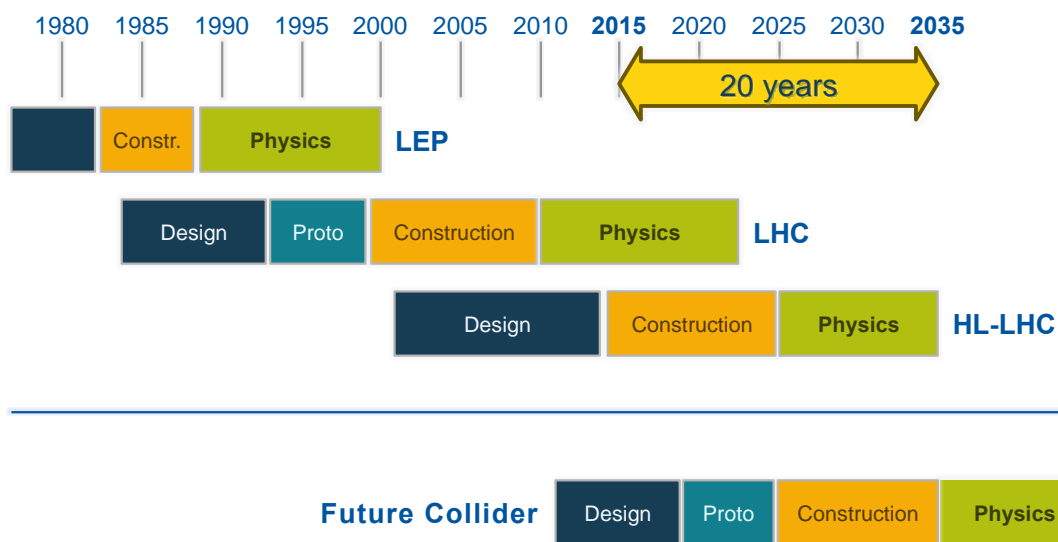
---

<sup>22</sup> See <http://research.microsoft.com/en-us/um/people/lamport/pubs/distributed-system.txt>.

Given the cost of today's storage and the likely evolution, there is no inherent cost why "data" cannot be stored more or less indefinitely. What is harder is to capture the necessary knowledge and validation procedures so that it can be used over long periods of time.

The "natural periodicity" of recent collider generations – some twenty years – is perhaps all one can hope for in terms of affordable data preservation. (Most LEP data – that of ALEPH, DELPHI and OPAL – may be usable somewhat longer, perhaps up to 25 / 30 years). Beyond that, re-use of the data will probably still be possible but may require a larger investment to "resuscitate", as has been done on rare (one?) occasion(s), notably for the JADE<sup>23</sup> experiment at the PETRA storage ring in DESY.

## CERN Circular Colliders + FCC



<sup>23</sup> See <https://wwwjade.mpp.mpg.de/> and the DPHEP Blueprint for further information.

# Content

1. Reports from labs and experiments
  - a. (we will start by collecting these, 4-6 pages from each contribution, also invited remotely, if necessary)
2. Physics case overview
  - a. the research we have gained, the added value
3. Review of DP models
  - a. the experience we have gained today's and tomorrow's
4. Review of DP-related technologies
  - a. and how they have been used in the past years
5. The resources and costs for DP, experience
6. A critical review of the DPHEP evolutions, the DPHEP Roadmap, The potential for common projects:

## Appendix A – The DPHEP Collaboration

<b>DPHEP Partner (May 2014 unless specified)</b>	<b>Location</b>	<b>Contact person</b>
European Organization for Nuclear Research, <b>CERN</b>	Switzerland	J. Shiers
Deutsches Elektronen-Synchrotron, <b>DESY</b>	Germany	D. South
Helsinki Institute of Physics, <b>HIP</b>	Finland	K. Lassila-Perini
Institute of High Energy Physics, <b>IHEP</b>	China	G. Chen
Institut national de physique nucléaire et de physique des particules, <b>IN2P3</b>	France	G. Lamanna
Institute of Particle and Nuclear Studies, High Energy Accelerator Research Organisation, <b>IPNS, KEK</b>	Japan	T. Hara
Max Planck Institut für Physik, <b>MPP</b>	Germany	S. Kluth
Institute of Particle Physics, <b>IPP</b> <b>(June 2015)</b>	Canada	R. Sobie
Science and Technology Facilities Council, <b>STFC</b> <b>(July 2015 – pending CB approval)</b>	UK	J. Bicarregui
Istituto Nazionale di Fisica Nucleare, <b>INFN</b> <b>(pending signature)</b>	Italy	M. Maggi

US labs might sign a “Letter of Intent” apparently? (Although they did sign the WLCG MoU).

## Appendix B – The DPHEP Implementation Board

(CERN e-group DPHEP-IB)

Alicia Calderon Tazon <Alicia.Calderon@cern.ch> Self added member

Andrew Branson <andrew.branson@cern.ch>

Andrii Verbytskyi <andrii.verbytskyi@cern.ch> Self added member

Benedikt Hegner <Benedikt.Hegner@cern.ch>

<boj@fnal.gov>

<cartaro@slac.stanford.edu>

<charles.f.vardeman.1@nd.edu>

David Colling <d.colling@imperial.ac.uk>

David Michael South <david.south@cern.ch>

<david.south@desy.de>

<denisov@to.infn.it>

Cristinel Diaconu <diaconu@cppm.in2p3.fr>

<dich@mail.desy.de>

<diesburg@fnal.gov>

Dirk Krucker <dirk.krucker@cern.ch> Self added member

<dirk.kruecker@desy.de>

Frank Berghaus <frank.berghaus@cern.ch>

<frank.berghaus@gmail.com>

<gang.chen@ihep.ac.cn>

<genevieve.romier@idgrilles.fr>

Gerardo Ganis <Gerardo.Ganis@cern.ch>

Gerhard Mallot <Gerhard.Mallot@cern.ch>

German Cancio Melia <German.Cancio.Melia@cern.ch>

<homer@slac.stanford.edu>

Jakob Blomer <Jakob.Blomer@cern.ch> Self added member

Jamie Shiers <Jamie.Shiers@cern.ch>

<jareknabrzyski@gmail.com>

Jetendr Shamdasani <Jetendr.Shamdasani@cern.ch> UWE

John Harvey <John.Harvey@cern.ch> Self added member

Kati Lassila-Perini <Katri.Lassila-Perini@cern.ch>

<kherner@fnal.gov>

<m.wing@ucl.ac.uk>

<marcello.maggi@ba.infn.it>

Marcello Maggi <Marcello.Maggi@cern.ch>

Marco Cattaneo <Marco.Cattaneo@cern.ch>

Maria Girone <Maria.Girone@cern.ch>

<matthew.viljoen@stfc.ac.uk>

Matthias Schroeder <Matthias.Schroeder@cern.ch>

<meenakshi\_narain@brown.edu>

<michael.d.hildreth.2@nd.edu>



Mihaela Gheata <Mihaela.Gheata@cern.ch>  
Miika Tuisku <miika.tuisku@iki.fi>  
Patricia Sigrid Herterich <patricia.herterich@cern.ch>  
<Pere.Mato@cern.ch>  
Peter Clarke <peter.clarke@ed.ac.uk>  
Predrag Buncic <Predrag.Buncic@cern.ch>  
Richard Mcclatchey <Richard.Mcclatchey@cern.ch>  
<Roger.Jones@cern.ch>  
Salvatore Mele <Salvatore.Mele@cern.ch>  
<silvia.amerio@pd.infn.it>  
<southd@mail.desy.de>  
Sunje Dallmeier-Tiessen <sunje.dallmeier-tiessen@cern.ch>  
<takanori.hara@kek.jp>  
Tibor Simko <Tibor.Simko@cern.ch>  
<Tim.Smith@cern.ch>  
<tpmccauley@gmail.com>  
Ulrich Schwickerath <Ulrich.Schwickerath@cern.ch>  
<wolbers@fnal.gov>  
<yves.kemp@desy.de>

DRAFT