

# 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
<b>THE DPHEP COLLABORATION AND IMPLEMENTATION BOARD.....</b>	<b>14</b>
<b>USE CASES, COST MODELS AND BUSINESS CASES.....</b>	<b>14</b>
<b>BIT PRESERVATION AND STORAGE TECHNOLOGY OUTLOOK.....</b>	<b>15</b>
<b>VIRTUALISATION AND VERSIONING FILESYSTEMS.....</b>	<b>15</b>
<b>DOCUMENTATION AND DIGITAL LIBRARY TECHNOLOGIES.....</b>	<b>15</b>
THE CERN PROGRAM LIBRARY: DOCUMENTATION AND SOFTWARE .....	15
<b>ANALYSIS CAPTURE AND REPRODUCIBILITY .....</b>	<b>16</b>
<b>RELATIONS WITH OTHER PROJECTS, DISCIPLINES AND INITIATIVES.....</b>	<b>16</b>
<b>CERN OPEN DATA PORTAL .....</b>	<b>18</b>
CERTIFICATION OF DIGITAL REPOSITORIES .....	18
<b>SITE / EXPERIMENT STATUS REPORTS (JUNE 2015) .....</b>	<b>21</b>
Belle I & II .....	21
BES III.....	21
HERA.....	22
LEP .....	23

Tevatron.....	23
BaBar.....	24
IPP.....	25
LHC.....	25
<b>TOWARDS A DATA PRESERVATION STRATEGY FOR CERN EXPERIMENTS.....</b>	<b>27</b>
<b>LESSONS FOR FUTURE CIRCULAR COLLIDERS.....</b>	<b>28</b>
<b>FUTURE ACTIVITIES.....</b>	<b>28</b>
<b>SUMMARY OF TECHNOLOGIES / “SERVICES” USED.....</b>	<b>28</b>
<b>OUTLOOK AND CONCLUSIONS.....</b>	<b>28</b>
<b>APPENDIX A – THE DPHEP COLLABORATION.....</b>	<b>30</b>
<b>APPENDIX B – THE DPHEP IMPLEMENTATION BOARD.....</b>	<b>31</b>

DRAFT

## 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/>.

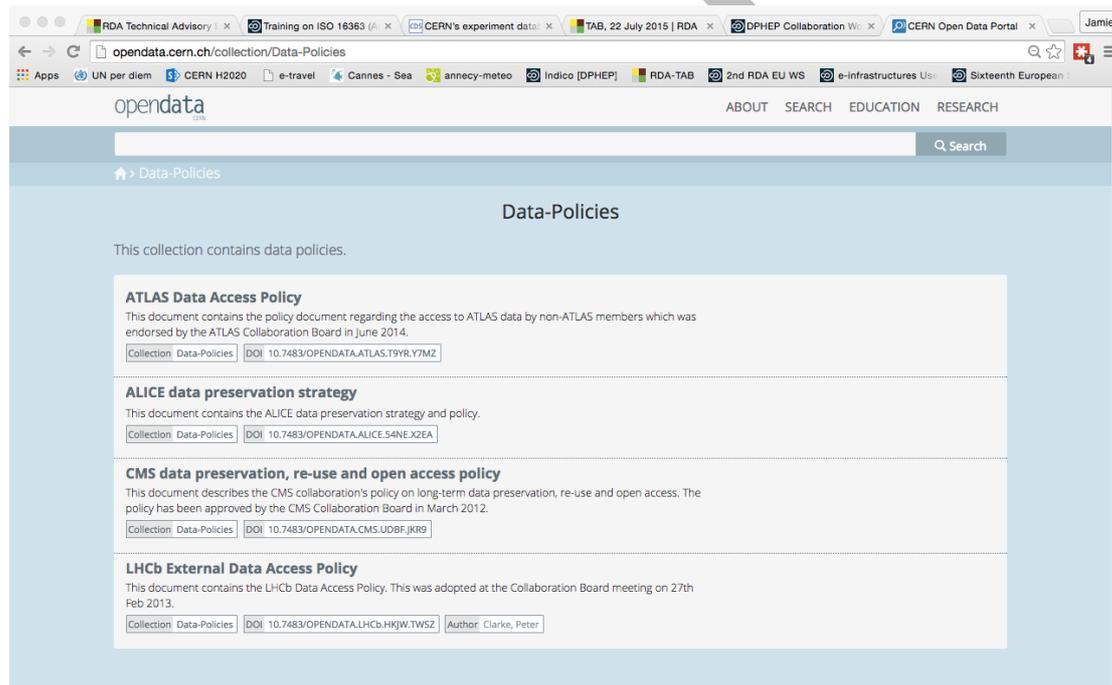
*At a minimum, DMPs must describe how data sharing and preservation will enable validation of results, or how results could be validated if data are not shared or preserved*

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 screenshot shows a web browser window with the URL [opendata.cern.ch/collection/Data-Policies](http://opendata.cern.ch/collection/Data-Policies). The page title is "Data-Policies" and it contains a list of four data policies:

- ATLAS Data Access Policy**: This document contains the policy document regarding the access to ATLAS data by non-ATLAS members which was endorsed by the ATLAS Collaboration Board in June 2014. DOI: 10.7483/OPENDATA.ATLAS.T9YR.V7MZ
- ALICE data preservation strategy**: This document contains the ALICE data preservation strategy and policy. DOI: 10.7483/OPENDATA.ALICE.S4NE.X2EA
- CMS data preservation, re-use and open access policy**: This document describes the CMS collaboration's policy on long-term data preservation, re-use and open access. The policy has been approved by the CMS Collaboration Board in March 2012. DOI: 10.7483/OPENDATA.CMS.UDBF.JKR9
- LHCb External Data Access Policy**: This document contains the LHCb Data Access Policy. This was adopted at the Collaboration Board meeting on 27th Feb 2013. DOI: 10.7483/OPENDATA.LHCb.HKJW.TWSZ. Author: Clarke, Peter

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.

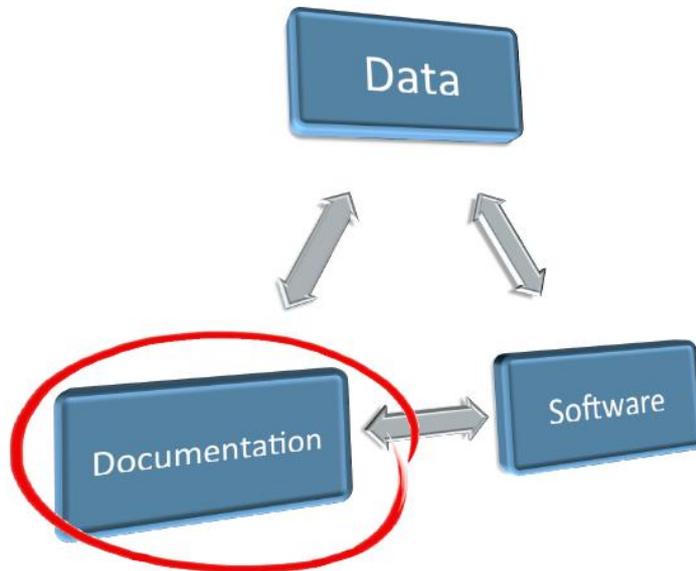
## HERA

- 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**



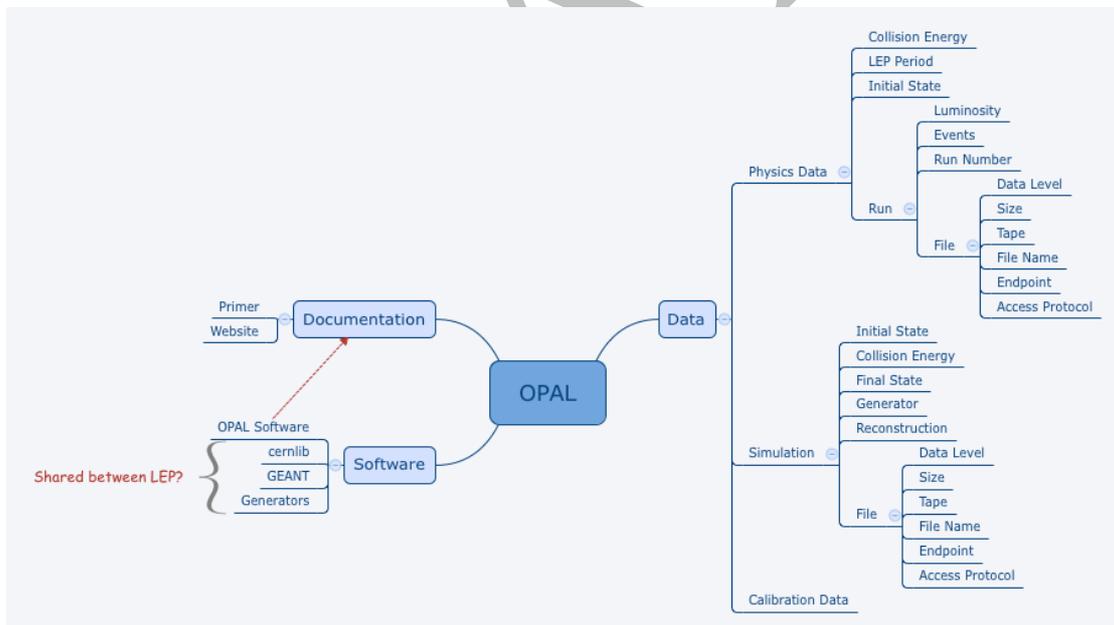
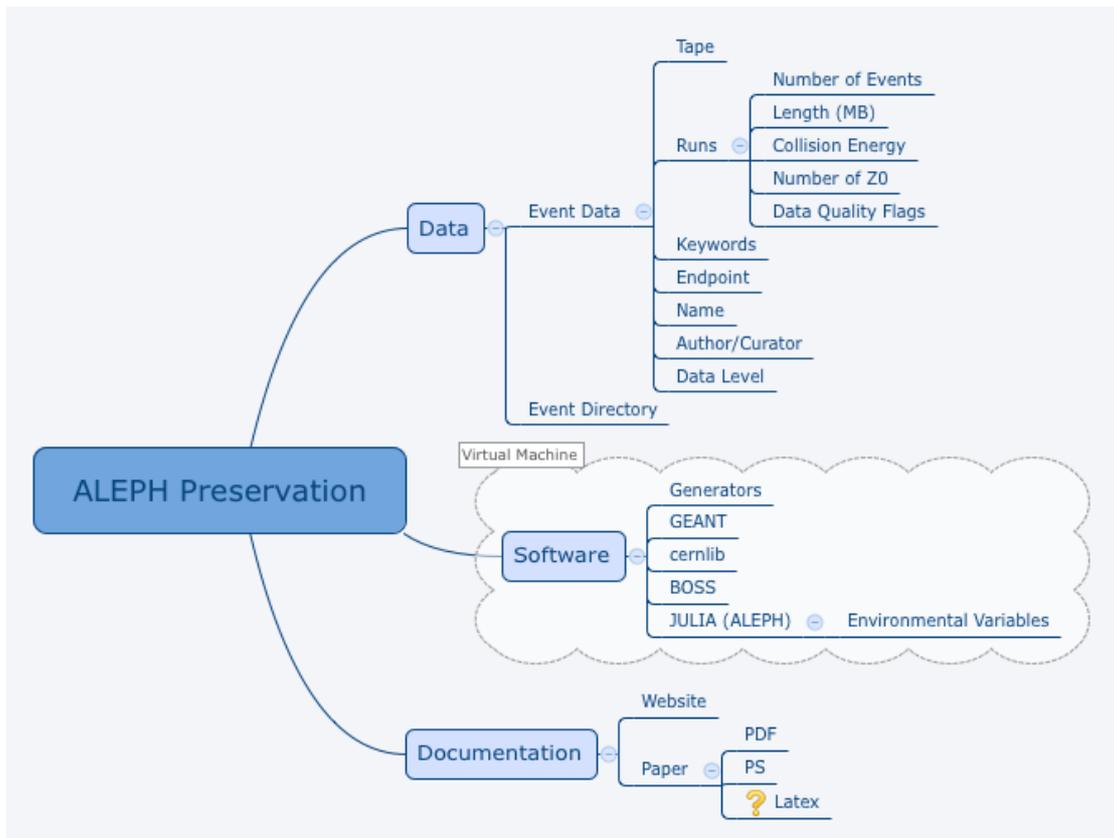
Information on the data, documentation and software is provided with a standard look and feel, although details are expected to vary.

# 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.

DRM



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 Board

## 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>13</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” (in OAIS terminology) is typically the (former) collaboration – although there is often considerable flexibility<sup>14</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>15</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.

[ more ]

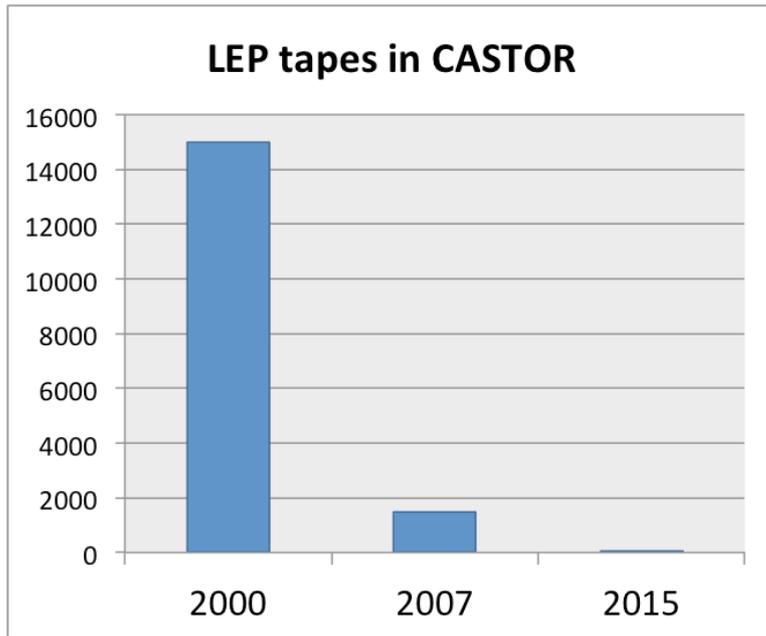
---

<sup>13</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.

<sup>14</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>15</sup> The cost model uses publically available pricing information and is thus suitable for sharing with other communities.

## Bit Preservation and Storage Technology Outlook



## Virtualisation and Versioning Filesystems

## Documentation and Digital Library Technologies

### 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.

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.

## Analysis Capture and Reproducibility

[ Suenje, Mike ]

## Relations with Other Projects, Disciplines and Initiatives

APA, APARSEN, SCIDIP-ES, RDA

...

## 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>16</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?

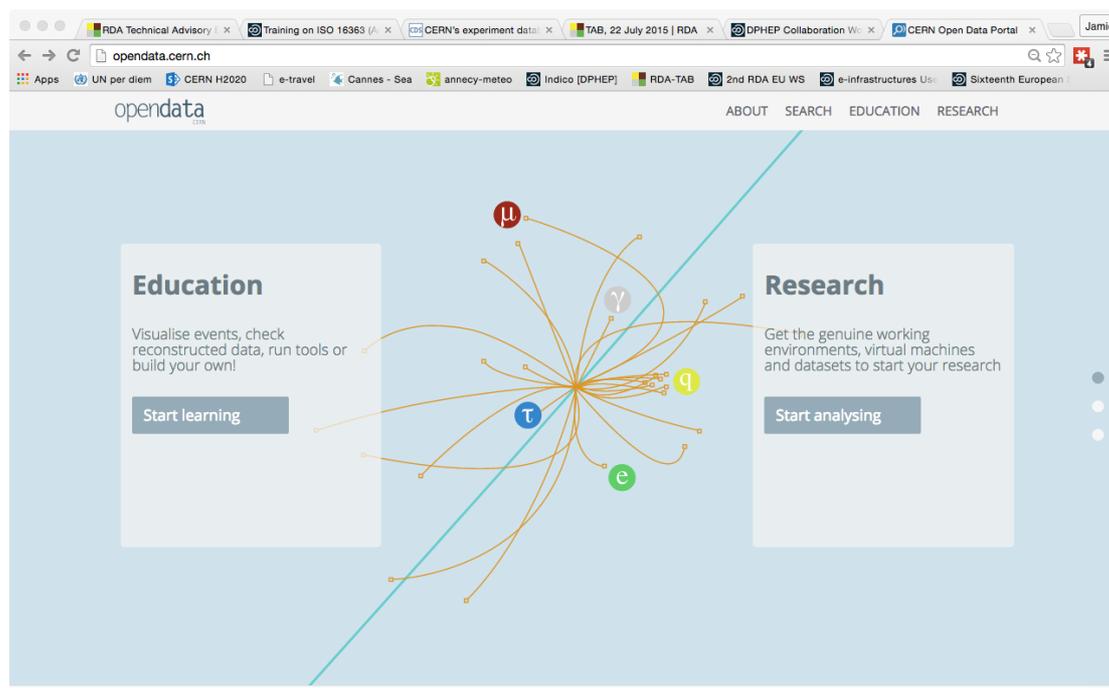
The answers to these questions may well vary with time and also depend on the implementation(s) that we choose: Open Access is just one step in the progression towards Open Data and finally “Open Knowledge”.<sup>17</sup>

---

<sup>16</sup> See <http://4cproject.eu/>.

<sup>17</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.

<p><b>3.3.1 THE REPOSITORY SHALL HAVE DEFINED ITS DESIGNATED COMMUNITY AND ASSOCIATED KNOWLEDGE BASE(S) AND SHALL HAVE THESE DEFINITIONS APPROPRIATELY ACCESSIBLE.</b></p>	<p>This is necessary in order that it is possible to test that the repository meets the needs of its Designated Community.</p>	<p>A written definition of the Designated Community.</p>
<p><b>3.3.2 THE REPOSITORY SHALL HAVE PRESERVATION POLICIES IN PLACE TO ENSURE ITS PRESERVATION STRATEGIC PLAN WILL BE MET.</b></p>	<p>This is necessary in order to ensure that the repository can fulfill that part of its mission related to preservation</p>	<p>Preservation Policies; Repository Mission Statement.</p>
<p><b>4.1.1 THE REPOSITORY SHALL IDENTIFY THE CONTENT INFORMATION AND THE PROPERTIES THAT THE REPOSITORY WILL PRESERVE.</b></p>	<p>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.</p>	<p>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.</p>
<p><b>4.3.4 THE REPOSITORY SHALL PROVIDE EVIDENCE OF THE EFFECTIVENESS OF ITS PRESERVATION ACTIVITIES.</b></p>	<p>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.</p>	<p>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.</p>
<p><b>5.1.1 THE REPOSITORY SHALL IDENTIFY AND MANAGE THE RISKS TO ITS PRESERVATION OPERATIONS AND GOALS ASSOCIATED WITH SYSTEM INFRASTRUCTURE.</b></p>	<p>This is necessary to ensure a secure and trustworthy infrastructure.</p>	<p>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.</p>

<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
<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>	Quantitative measures (# papers, PhDs etc) exist
<b>Uniqueness</b>	
<b>Resources</b>	
<b>Status</b>	
<b>Issues</b>	
<b>Outlook</b>	

### BES III

Preservation Aspect	Status
<b>Bit Preservation</b>	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)
<b>Data</b>	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>
<b>Documentation</b>	<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 / yearβ

## 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>	All online webpages and code archived, still accessible from CDF webpages. At the time of Tevatron shutdown <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> Long term future solution: build legacy release that contains no pre-SL6 libraries CVMFS for code distribution
<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?)

Preservation Aspect	Status
<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

Preservation Aspect	Status
<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>18</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>18</sup> See <http://council.web.cern.ch/council/en/EuropeanStrategy/ESParticlePhysics.html>.

## Lessons for Future Circular Colliders

### Future Activities

### Summary of Technologies / “Services” Used


### 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 in sync with major events such as CHEP) should be sufficient to ensure that coordinated approaches remain and that duplication is minimised.

DRAFT

# 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:

DRAFT

## 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