# Data Preservation in High Energy Physics: a collaborative perspective

Cristinel Diaconu, Ulrich Schwickerath
on behalf of the DPHEP collaboration

# Outline

- Apropos data preservation
  - Why preserve data
  - DPHEP collaboration history
  - Principles, data kinds and **risks**, preservation levels, status
  - Relationship to Open data/ open science
- Highlights
  - from our last workshop in autumn 2024
- Lessons learned and conclusions

# Quick intro about myself

- Until 2002, particle physicist at LEP, DELPHI experiment
- Working in CERN IT since 2005
- IT coordinator for scientific data preservation
- Member of the DPHEP collaboration
- About **24 years of history in DP**



http://dphep.org
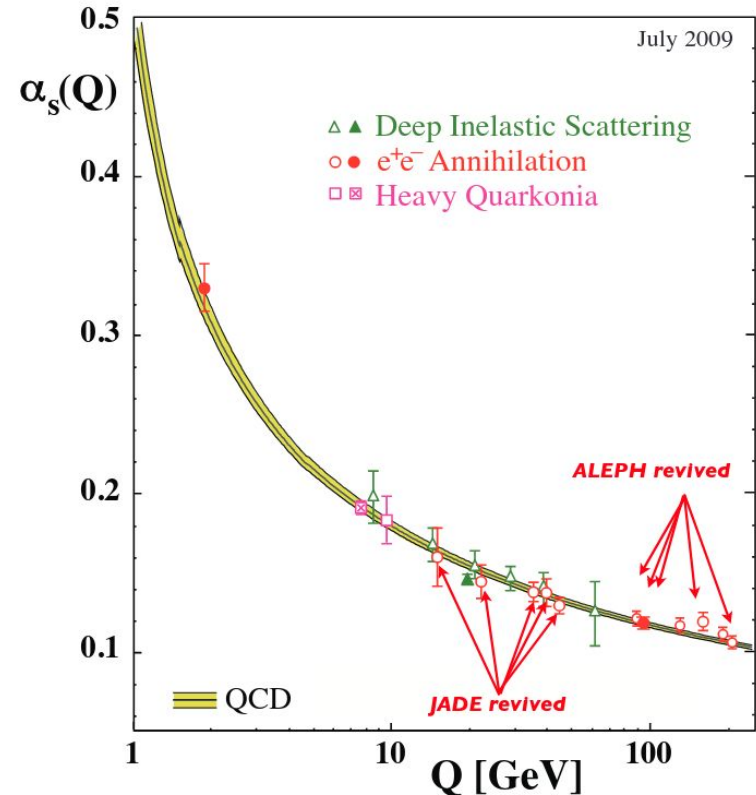
# Data preservation: why?

**Example:** the JADE experiment
- Experiment at Petra @ DESY
- **1978-1986**

Took about **8 years** to recover the data and the full software stack, **1995-2003**

[Full story](#) from Siggi Bethke (Talk at KEK)

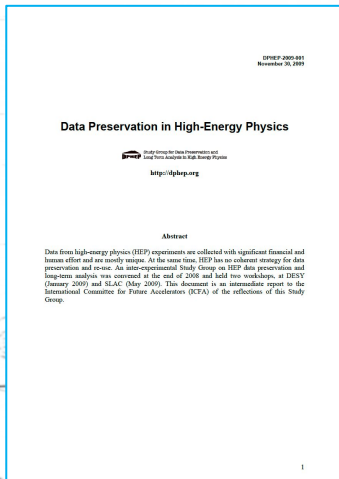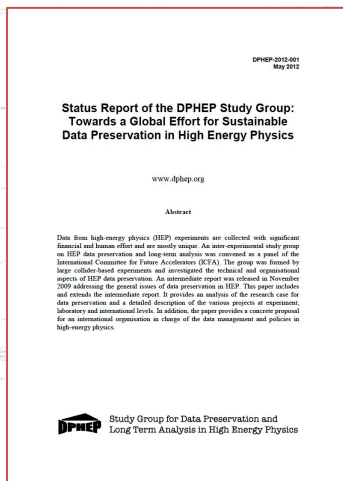**Proof of the energy dependence of the strong coupling constant**

# The DPHEP Collaboration

**2009**
Letter of intent

arXiv:0912.0255

**2012**
Blueprint

arXiv:1205.4667

**2015**
Collaboration MoU

arXiv: 1512.02019

**2023**
Decade report

arXiv: 2302.03583

**2025**
EPPS 2026 submission

2503.23619 [hep-ex]

# Workshops and activities

- DPHEP reports to the "International Committee for Future Accelerators" (ICFA)

- Regular workshops:
  - Every 2-3 years
  - 4th DPHEP workshop October 2024
  - Upcoming: 5th DPHEP workshop March 2026

- LEP and CERNLIB TF meetings

  https://indico.cern.ch/category/4458/

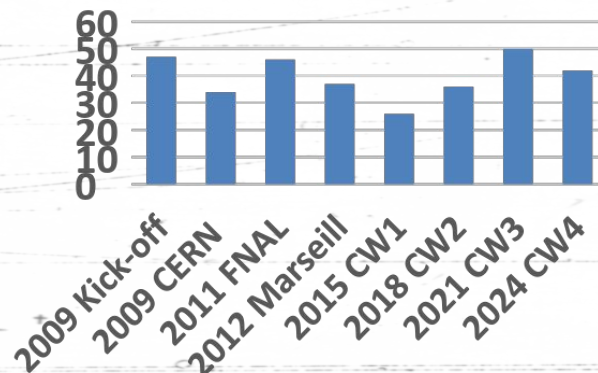| What is data ? | What is preservation ? |
|---|---|
| NOT only "files" | NOT a freezer, herbarium, museum, album, cellar … |
| Every digitally encoded information that was created as a result of planning, running and exploiting an experiment | The **process** of transforming a "high intensity/rapidly changing" computing system into a low intensity/slowly evolving computing system, while conserving the capability to extract new science from the data |
| | Requires a clear plan and long term organisation<br><br>● Within each collaboration<br>● and at international level (DPHEP) |

# Areas of data DP

## Bit - Preservation

Maintain and keep accessible:
- Raw data
- Reconstructed data
- Simulated data
- Meta data
- databases, conditions data

## Software preservation

Includes:
- Analysis framework
- Reconstruction software
- Simulation software
- Visualisation

Note: Keep the sources (alive). Binaries may stop working

## Documentation

- Published results
- Notes and internal documentation
- Technical documentation
- Manuals
- Web pages
- ….

## Analysis preservation

- What has been done ?
- How has it been done ?
- Which data sets were used ?
- Analysis software sources

**Ability to reproduce** published results long term

DPHEP

# Areas of DP: Risks

## Bit - Preservation

- physical media failure (e.g. tape damage)
- system errors
- human errors

Mitigation: dual-tape, external copies

## Software preservation

- Ever changing IT infrastructure
- Architecture and technology changes
- Storage systems access changes
- Security enhancements, e.g. deprecation of ciphers
- Compiler and computing language changes
- **Deprecation of required external dependencies**

## Documentation

- Quality of internal documentation
  - May not be intended for general public
  - Use of jargon
- Author agreements, copyright
- Paper only copies with bad quality
- Document format long term support

Mitigation: QC from the start of the experiment, digitization
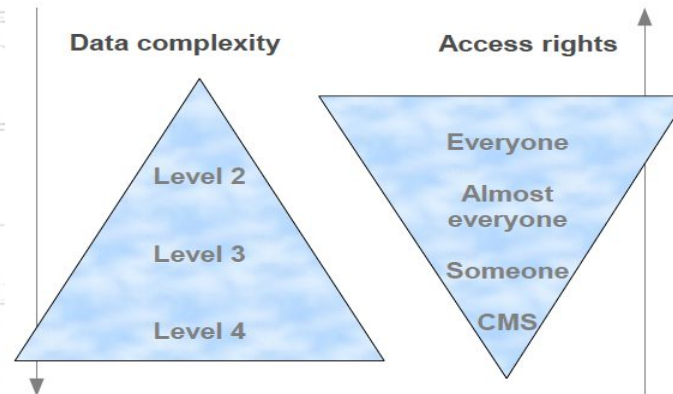
## Analysis preservation

- Badly written code
- Code not being shared and lost
- Code not maintained long term
- Proper documentation of the code and the analysis details

Mitigation: Guidelines for analyses, enforce these if possible

# Guidance into data complexity
## The "**DPHEP Preservation Levels**"

| Preservation Model | | Use Case | |
|---|---|---|---|
| 1 | Provide additional documentation | Publication related info search | **Documentation** |
| 2 | Preserve the data in a simplified format | Outreach, simple analyses | **Outreach, reanalysis** |
| 3 | Preserve the analysis level software and data format, frameworks | Full scientific analysis, based on the existing reconstruction | **Technical Preservation Projects** |
| 4 | Preserve the reconstruction and simulation software as well as the basic level data | Retain the full potential of the experimental data | |



Data complexity — Level 2, Level 3, Level 4

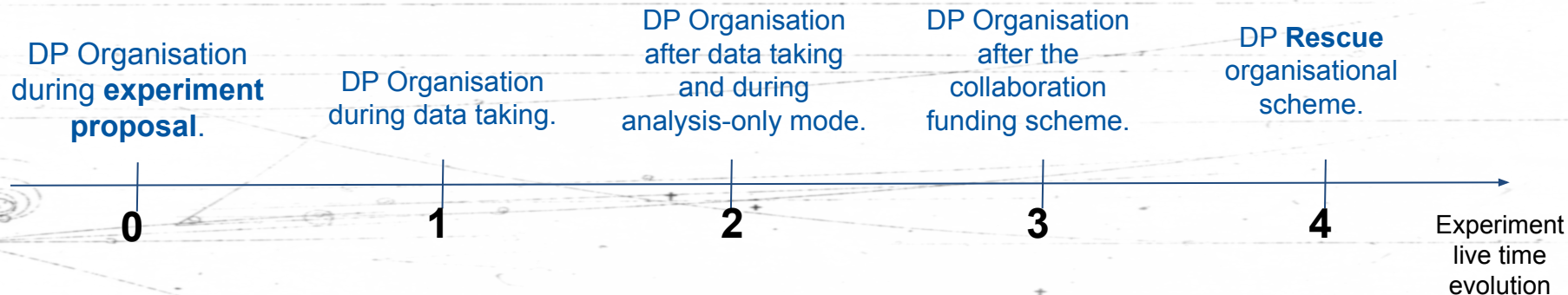Access rights — Everyone, Almost everyone, Someone, CMS

DPHEP

# A matter of collaboration as well (I)

The supervision and knowledge transfer/**capture** is essential at long term

Need to clarify the status and the rules

Various stages of organisation (CL) can be defined:

DP Organisation during **experiment proposal**.

DP Organisation during data taking.

DP Organisation after data taking and during analysis-only mode.

DP Organisation after the collaboration funding scheme.

DP **Rescue** organisational scheme.

**0**　　　　**1**　　　　**2**　　　　**3**　　　　**4**

Experiment live time evolution

DPHEP

# A matter of collaboration as well (II)

| Level 4: This organisation scheme is to be activated when: | |
|---|---|
| The host laboratory stops support and announce no long-term commitment. | The official collaboration/ data stewardship is stopped with no further plans (no step 3 is clearly defined). |

**Taking no action == decommissioning (deleting) the data.**

"Securely" storing/freezing the files and the latest version of the software is certainly not a substitute for a preservation project.

| Laboratory/ Collider | Experiment | Period | Preservation Level | Data Volume | Present status | CL |
|---|---|---|---|---|---|---|
| DESY/PETRA | JADE | 1979–1986 | 4 | 1 TB | Analysis running on preserved data; migrated from DESY to MPP | 4 |
| CERN/LEP | ALEPH, DELPHI, OPAL L3 | 1989-2000 | 4 2 | 0.5 PB | ADO: Analysis running on preserved data | 4 |
| DESY/HERA | H1 ZEUS | 1992 – 2007 | 4 3/4 | 0.5 PB 0.2 PB | Analysis running on preserved data | 3 |
| SLAC/PEP II | BABAR | 1999–2008 | 4 | 2 PB | Analysis running on preserved data; migrated from home lab to different centers | 4 |
| KEK/KEKB | Belle I | 1999-2010 | 4 | 4 PB | Analysis running on preserved data; Compatible with Belle II computing | 2 |
| FNAL/TeVatron | DØ CDF | 1983–2011 | 4 4 | 8.5 PB 9 PB | Archived on tapes | 4 |
| BNL/RHIC | PHENIX | 2000–2016 | 3 | 25 PB | Analysis running on preserved data | 3 |
| FNAL/ν-beam | Minerva | 2010–2019 | 3 | 10 TB | Analysis running | 2 |
| IHEP/BEPCII | BESIII | 2009–2030 | 4 | 6 PB | Collecting and analyzing data | 1 |
| CERN/LHC | ALICE, ATLAS, CMS, LHCb | 2010-2041 | 4 | O(1EB-10EB) | Collecting and analyzing data | 1 |

DPHEP

# Cost and Benefit

| | |
|---|---|
| **C1** | Host laboratories allocate **person power** and **computing resources** - specifically to DP. <br> (in % to the construction/operation costs) |
| **C2** | Collaborating laboratories participate in the effort: **replicate or take over** data and computing systems and provide technical assistance. |
| **C3** | Researchers and engineers participate outside their main research area. |
| **C4** | Innovative computing projects, including pluri-disciplinary open science initiatives, may offer attractive opportunities for data preservation and are therefore an indirect source of support. |
| **C5** | The proximity of a **follow-up experiment** clearly helps in structuring and supporting a data preservation project. |

| | |
|---|---|
| **B1** | **New publications** – counting here those executed with a strong involvement of the dedicated DP systems. |
| **B2** | **Derivative work**: Publications made by other groups/people using the new publications produced at B1. |
| **B3** | **Preserving** the **scientific expertise** and the leadership in the field of the experiment, possibly **boosting the transition** to a new experiment |
| **B4** | Technology expertise in robust data preservation. Improved ability to **plan for new experiments** and preserve their scientific potential at long term. |

## Figure of Merit:
FoM = B1/C1

DPHEP

# Conclusions after 10 years: the scientific output

DP is a **cost-effective way of doing fundamental research** by exploiting unique data sets in the light of the increasing theoretical understanding.

DP leads to
- a **significant increase in the scientific output** (>10% typically)
- for a minimal investment overhead (0.1%).
- …. As predicted in 2013

**Lesson:** When collisions are stopped, **~20% of the publications are still to come**, and half of them are unknown/unplanned!

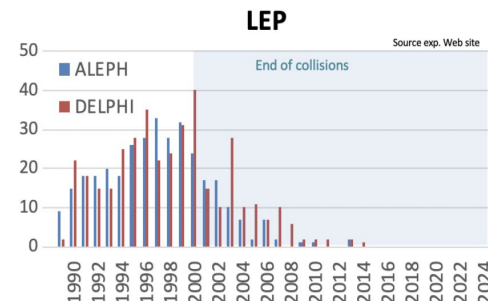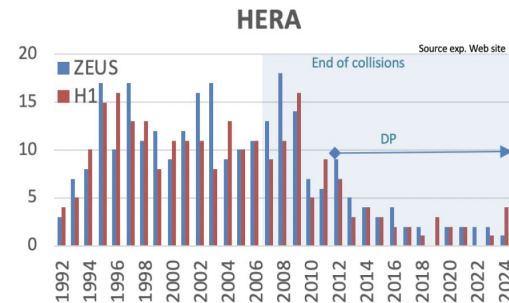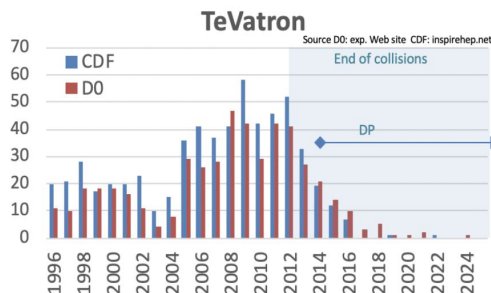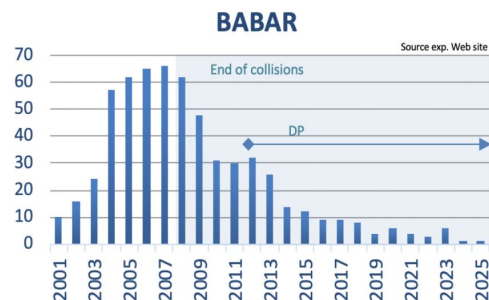LHC will have 3 decades ahead after the end of collisions !



| | Data taking stopped | Publications before 2012 | Publications 2012-2022 | Scientific return increase % |
|---|---|---|---|---|
| Babar | 2008 | 471 | 154 | 33% |
| H1+ZEUS | 2007 | 436 | 62 | 14% |

# … and ongoing

BaBar April 2023
The 600th paper

## News
News from the DESY research centre

ZEUS June 2023

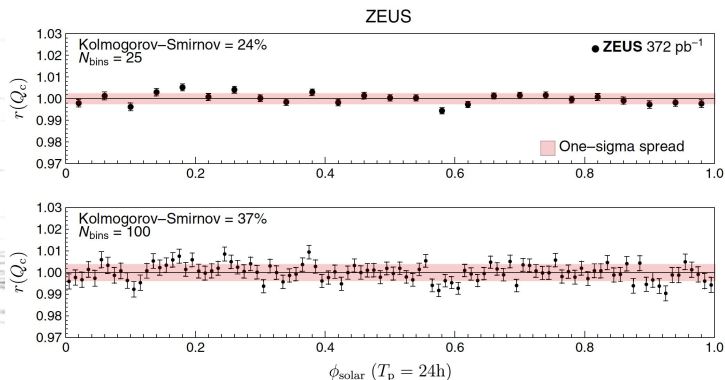2023/06/20

Back

### Do quarks interact with the cosmos?

HERA data places limits on the interactions between quarks and cosmic background fields

DESY's HERA collider, decommissioned in 2007, is still providing valuable results to scientists. A newly released paper shows that quarks, which were the main particles under investigation at the electron–proton collider, do not visibly interact with potential cosmic background fields. This means that they don't violate a fundamental symmetry of nature, the rotation and Lorentz invariance. HERA was specifically well-suited for studying quarks, so these results set important limits for other experiments and searches.

PHYSICAL REVIEW D **107**, 072001 (2023)

Study of the reactions $e^+e^- \to K^+K^-\pi^0\pi^0\pi^0$, $e^+e^- \to K^0_S K^\pm \pi^\mp \pi^0 \pi^0$, and $e^+e^- \to K^0_S K^\pm \pi^\mp \pi^+ \pi^-$ at center-of-mass energies from threshold to 4.5 GeV using initial-state radiation

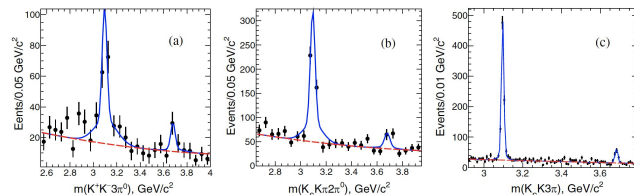J. P. LEES *et al.*　　　　　　　　　　　　PHYS. REV. D **107**, 072001 (2023)

FIG. 16. The $J/\psi$ invariant mass region for the (a) $K^+K^-3\pi^0$, (b) $K^0_S K\pi 2\pi^0$, and (c) $K^0_S K3\pi$ events. The curves show the fit functions described in the text.

## Unbinned Deep Learning Jet Substructure Measurement in High $Q^2$ ep collisions at HERA

**H1** Collaboration · V. Andreev (Lebedev Inst.)　Show All(148)

Mar 23, 2023

30 pages
e-Print: 2303.13620 [hep-ex]
Report number: DESY-23-034

H1 Mars 2023

# The JADE Experiment at the PETRA $e^+e^-$ collider -- history, achievements and revival

S. Bethke (Munich, Max Planck Inst.), A. Wagner (DESY)

Aug 23, 2022

58 pages

Published in: *Eur.Phys.J.H* 47 (2022) 16

e-Print: 2208.11076 [hep-ex]

DOI: 10.1140/epjh/s13129-022-00047-8 (publication)

Report number: MPP-2022-109

Experiments: DESY-PETRA-JADE

View in: ADS Abstract Service

📄 pdf    ⌕ cite    ⊟ claim

## Measurement of jet production in deep inelastic scattering and determination of the strong coupling at ZEUS

ZEUS Collaboration • I. Abt (Munich, Max Planck Inst.) Show All(80)

Sep 6, 2023

42 pages

Published in: *Eur.Phys.J.C* 83 (2023) 11, 1082

Published: Nov 27, 2023

e-Print: 2309.02889 [hep-ex]

DOI: 10.1140/epjc/s10052-023-12180-9 (publication)

Report number: DESY-23-129

Experiments: DESY-HERA-ZEUS

View in: ADS Abstract Service

📄 pdf    ⌕ cite    ⊟ claim    ⊟ datasets              ⌕ reference search    ↺ 14 citations

---

Preprint    PDF Available

## Analysis note: measurement of energy-energy correlator in $e^+e^-$ collisions at 91 GeV with archived ALEPH data

May 2025

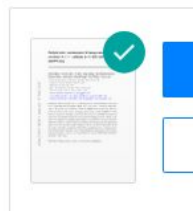DOI: 10.48550/arXiv.2505.11828

License · CC BY 4.0

Authors:

☐ **Hannah Bossi**        ⊙ **Yu-Chen Chen**

☐ **Yi Chen**            ☐

Show all 10 authors

---

Citations per year

## Search for baryogenesis and dark matter in $B^+ \to \Lambda_c^+ +$ invisible decays

BaBar Collaboration • J.P. Lees Show All(221)

Dec 9, 2024

7 pages

Published in: *Phys.Rev.D* 111 (2025) 3, L031101

Published: Feb 1, 2025

e-Print: 2412.06950 [hep-ex]

DOI: 10.1103/PhysRevD.111.L031101 (publication)

Report number: BABAR-PUB-24/001, SLAC-PUB-241025

Experiments: SLAC-PEP2-BABAR

View in: ADS Abstract Service

📄 pdf    ⌕ cite    ⊟ claim              ⌕ reference search    ↺ 0 citations

# Boosting the future experiments

- **HERA ⯈ EIC**
  - "Scientists today have a **renewed interest in HERA's particle experiments**, as they hope to use the data – and more precise computer simulations informed by tools like OmniFold – to aid in the analysis of results from future electron-proton experiments, such as at the Department of Energy's next-generation **Electron-Ion Collider (EIC)**. "

- **Possibly**
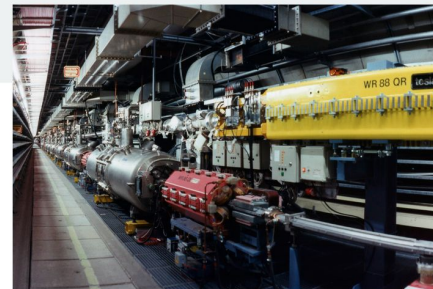  - **LHC ⯈ FCChh**
  - **LEP ⯈ FCCee**

**ARTICLE · MYSTERIES OF MATTER**

## How Do You Solve a Problem Like a Proton? You Smash It to Smithereens – Then Build It Back Together With Machine Learning

By **Theresa Duque**
October 25, 2022

New tool decodes proton snapshots captured by history-making particle detector in record time

CONTACT MEDIA@LBL.GOV →

Looking into the HERA tunnel: Berkeley Lab scientists have developed new machine learning algorithms to accelerate the analysis of data collected decades ago by HERA, the world's most powerful electron-proton collider that ran at the DESY national research center in Germany from 1992 to 2007. (Credit: DESY)

**Preserved data can be used to transfer knowledge, training/teaching, outreach or boosting new research programs**

# Relationship with Open Data

- Two different beasts, albeit related
- Not the same but complementary to each other
  - Open Data can be serve as a long term data preservation strategy
  - Can be made as easy as agreeing on a data access strategy, and then dropping access restrictions
- Open Data requires a certain level of data preservation
  - Can be restricted to a subset or an extract (ntuple) of the preserved data
- Open Data solves the long term data access problem

DPHEP

# 4th DPHEP workshop Oct 2024: Agenda

## Wednesday Oct 2nd

- 14:25 **ALEPH**; Jacopo Fanini
- 14:45 **CERNLIB;** Andrii Verbytskyi, Ulrich Schwickerath
- 15:00 **DELPHI** ; Dietrich Liko, Dr Ulrich Schwickerath
- 15:15 **OPAL** ; Matthias Schroeder
- 15:30 **DELPHI and OPAL event displays**; M.Schroeder

- 16:00 **ZEUS**; Achim Geiser
- 16:20 **H1** ; Speaker: Henry Klest
- 16:40 **JADE** Andrii Verbytskyi /Richard Hildebrandt
- 17:00 **PHENIX** ; Maxim Potekhin
- 17:20 **BaBar** : Marcus Ebert

## Thursday Oct 3rd

- 09:00 **KEK / Belle I & II** ; Takanori Hara
- 09:20 **BESIII** Gang Chen
- 09:40 **CERN Open Data portal** Pablo Saiz
- 10:00 **REANA** Marco Donadoni (CERN)
- 10:20 **CERN Analysis Preservation porta** P. Fokianos
- 11:00 **CERN Open Data: Policy/implementation**; J. Boyd
- 11:20 **ALICE** : David Dobrigkeit Chinellato
- **11:40 ATLAS**; Zach Marshall
- 12:00 **LHCb**; Dillon Fitzgerald
- 14:00 **Preserving ANTARES legacy data** ; Jutta Schnabel **14:20 PUNCH4NFDI** ; Achim Geiser
- 14: **40 CMS** ; Julie Hogan
- 15:00 **ICFA Data Lifecycle Panel** ; Kati Lassila-Perini
- 15:25 **DPHEP Collaboration**

# 4<sup>th</sup> DPHEP workshop

| Key points |
|---|
| **Significant progress** has been made since the last workshop, both on data preservation and opening of data. |
| While CMS has pioneered the publication of **open data**, the other LHC experiments are rapidly catching up now: ATLAS, LHCb, ALICE |
| **Open data policies** are increasingly applied also **beyond LHC**; useful synergy and data opening calendar |
| **LEP Data is** (mostly) **alive and active** ! DELPHI data has been recently released as open data |

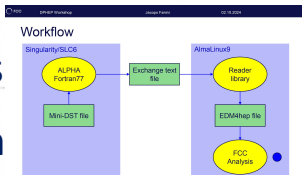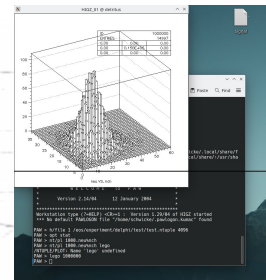| Observations |
|---|
| In many contributions continued funding of DP was mentioned as an issue. An example is the BaBar experiment, whose software is running on outdated hardware but there is no funding to replace them. |
| Transfer across generations is visible<br>• HERA, RHIC -> EIC<br>• LEP -> FCC |
| Bottom up approach starting from the people involved in the practical work |
| Packed agenda, with input from HEP and also astromomy |

# Highlights:LEP data is back !

- **ALEPH:**
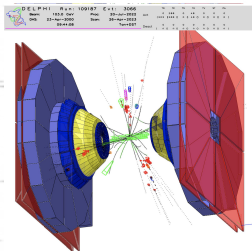  - Full stack on containers
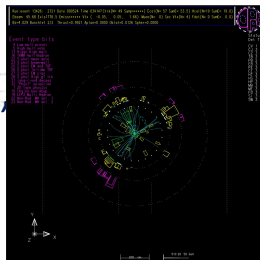  - EDM4HEP migration on

- **DELPHI:**
  - Full stack available, based on community CERNLIB
  - Available on Opendata.org

- **OPAL:**
  - Full stack resurrected, based on community CERNLIB
  - Plan to open the data as DELPHI did, work on this ongoing



## (Community) CERNLIB
  - 64bit support
  - Works on decent OS versions
  - Binaries soon available from /cvmfs/dphep.cern.ch
  - Support for xrootd

## Collaborative effort:

- Frequent LEP TF meetings
- CERNLIB and LEP collaboration
- Synergies e.g. for event displays of DELPHI and OPAL
- Interest in DELPHI about EDM4hep

# Best practices: Learn from the past

- ICFA panel on best practices being collected
  - Coming soon! Stay tuned!

- Learn from past experiments
  - Remember Murphies law

| Keep it simple | |
|---|---|
| Rely on open and free software | Remove license cost and avoid 3rd party closed source software |
| Reduce external dependencies as much as possible | Removed site specific dependencies as much as possible |
| Preserve documentation and knowledge | Ensure documentation quality |
| Use Automation | Continuous integration |
| | Automated testing |
| Plan for and enforce analysis preservation | Used to be a difficult task, and still is ! |

# Conclusions

| | |
|---|---|
| Significant/measurable impact of dedicated DP projects @expts./labs | – Production of high quality and unique scientific results at very low (non-zero) cost<br>   • **10%** output for less than **1%** investment: ✓<br>   • **Long term organisation proves to be productive**<br>– Signs of re-vigorating collaborations in the context of new projects, e.g. HERA-EIC; LEP-FCCee<br>– Case for longer term preservation: data sets parking<br>   • CDF, D0, Babar, LEP, JADE : carefully follow the usability in time |
| There is full coherence (but not total overlap) between **DP** and **Open Data/Science** | – LHC experiments consider both, looking forward to 2041+20/30 years<br>– Lesson: When collisions are stopped, 20% of the publications are still to come, and half of them are unknown/unplanned! |
| The (DP)HEP future is also considered | FCC, EIC : transfer of knowledge in DP from LHC/oldies |
| **And more is possible on education, training, outreach … via open data** | |

DPHEP