The past and future of Data Preservation in High Energy Physics

Collider data ageing: the bottle and the cellar



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



http://dphep.org



Cristinel DIACONU CPPM/CNRS/Aix-Marseille University



10 years ago.....

THE ADVENT OF PRESERVED DATA ANALYSIS

At the end of the first decade of the XXI century, many HEP experiments were close to the data taking.

The world faced the rise of the LHC experiments.

A big data deluge was threatening the HEP computing centers.

The DPHEP collaboration federates the data preservation initiatives, monitor and report the worldwide status of preserved data.

But the fight is not over, data deletion forces are still alive....

Data collection accellerates



Digital data are fragile

- Storage capacity is physically exceeded
- Unattended/orpahned data vanishes quickly



FIGURE 1.3: Information and Storage Source: J. Gantz January 2008 (revised). Used with permission.

In Blue Ribbon Task Force report

The six Vs of big data

Big data is a collection of data from various sources, often characterized by what's become known as the 3Vs: *volume, variety and velocity*. Over time, other Vs have been added to descriptions of big data:

VOLUME	VARIETY	VELOCITY	VERACITY	VALUE	VARIABILITY
The amount of data from myriad sources.	The types of data: structured, semi-structured, unstructured.	The speed at which big data is generated.	The degree to which big data can be trusted.	The business value of the data collected.	The ways in which the big data can be used and formatted.
	بېنې *	٢		O	

https://www.techtarget.com/searchdatamanagement/definition/big-data



Models of data preservation and acces

- Collaborations addressed this issue in a generic way
 - e.g. Blue Ribbon, APA, DPC, eSciDir, RDA ...



Source: Consultative Committee for Space Data Systems January 2002.

>

Where is your data?

MISSING DATA

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



Scientific Data

- Structured following a scientific plan
- > Diverse sources
- Large and expensive projects
 - Not easy to repeat
- Contain unique knowledge, objective/subjective rigidity in the long term
 - « time stamped »
 - « technology stamped »
 - « common knowledge stamped » ….
- > Data Observatories
 - SD usually contain more information than initially needed/intended
 - less than possible, more than intended

Exemple: HEP experiments in ± 10 ans



[not all programmes, dates are approximate, just to give the picture]

HEP Data



What is "data"?

> The persistent confusion:

 "data" = an operating system files, i.e. bits on a memory support such as disks, tapes etc.

This very superficial view, is useless for any running experiment and cannot exclusively apply to any useful thinking on the long term data preservation.

- Data = "every digitally encoded information that was created as a result of planning, running and exploiting an experiment"
 - digital data files: raw and processed, control/configuration, meta-data resuming environmental parameters, operational data etc.
 - software in all its forms (front-end, trigger, middleware, reconstruction, classification include machine learning setups, high-level analysis, visualisation etc.)
 - documentation files (internal/public notes, publications, manuals, workflows)
 - organisation and diffuse knowledge files (rules, minutes, meetings and slides, news, blogs, logbooks etc.)

Can HEP data "age"?

- The wine example (*)
 - Constituents
 - Bottle
 - Storage
 - Evolution (including "dumb phase")

(*) l'abus d'alcool est dangereux pour la santé à consommer avec modération <u>alcohol abuse is dangerous for your health and</u> <u>should be consumed in moderation</u>

- Not: a freezer, a herbarium, a museum, an album etc. ... and not a cellar
 And certainly not a "save it on tape, we will download it later if needed"
- Preservation: the process of transforming a "high intensity/ rapidly changing " computing system into a "low intensity / slowly evolving" computing system with conserving the capacity of extracting new science from the "data" (within its definition of above).
 - ensure physical existence of data from a digital point of view (see data definition above, all this has to be physically saved and secured at long term - and that includes software, of course) - note that this is the simplest and basically solved aspect of DP in HEP.
 - as an obvious (and relatively easy to solve) aspect of the previous item: identify and provide computing and storage resources.
 - ensure the functionality of the whole system, identify the potential risks and take appropriate measures as technology and community evolve. The level of complexity differ for the various aspects of the data. The simplest examples include the digital files, the documentation etc. that need only storage and access, i.e. rather standard operations independent of the experiment complexity in general. In contrast, specific experimental software and databases are much more difficult to keep functional across technological changes (hardware, operating systems etc.).
 - **define and identify the human resources** related to the research plan
 - oversight and manage the collaborative work and manage the preserved data analysis activity according to the DP design.
 - define and implement the policies for data usage, including opening the access to data to new collaborators and/or releasing the data to larger (not identified apriori) communities.
 - observe and update the physics case of the preserved data. It should be noted that the technical solutions and the necessary choices on the information to be dismissed while designing a long term preservation system should decouple as much as possible from the epoch-related physics case. Indeed, the door should remain open for unexpected analyses (see below discussion on preservation levels).

DPHEP Study Group (2009)

> arXiv:0912.0255



- An urgent and vigorous action is needed to ensure data preservation in HEP
 - Examples for the physics case explored
 - Data is rich and can be further exploited in most cases beyond the collaboration lifetime
- The preservation of the full analysis capability of experiments is recommended, including the preservation of reconstruction and simulation software
- An interface to the experiment know-how should be introduced: data archivist position in the computing centres
 - The preservation of HEP data requires a synergic action: collaborations, laboratories and funding agencies
- An International Data Preservation Forum is proposed as a reference organisation. The Forum should represent experimental collaborations, laboratories and computing centres

DPHEP Blueprint May 2012

- Full status report of the activities of the DPHEP study group, including:
 - An expanded description of the physics case
 - Defining and establishing data preservation principles
 - Updates from the experiments and joint projects
 - FTE estimates for these and future projects

Priority 1:

Next steps to establish fully DPHEP in the field

	DPHEP-2012-001 May 2012
Sta To Dat	tus Report of the DPHEP Study Group: owards a Global Effort for Sustainable a Preservation in High Energy Physics
	www.dphep.org
	Abstract
Data from financial an on HEP da Internationa large collid aspects of November 22 includes and case for dat experiment, p concrete pro management a	high-energy physics (HEP) experiments are collected with significant d human effort and are mostly unique. An inter-experimental Study Group ta preservation and long-term analysis was convened as a panel of the Committee for Future Accelerators (ICFA). The group was formed by ar-based experiments and investigated the technical and organisational by Badressing the general issues of data preservation in HEP. This paper extends the intermediate report was released in a preservation and a detailed description of the various projects at aboratory and international levels. In addition, the paper provides a model of an international organisation in charge with the data and policies in high-energy physics.
DPHEP	Study Group for Data Preservation and Long Term Analysis in High Energy Div

Local Action in experiments, laboratories	Data archivists: 0.5-1 FTE /lab					
Priority 2:	Project Manager: 1 FTE					
International	Technical support: 0.2 FTE					
organization	Contributions from Labs: 0.2/lab					
	(data archivists)					
Priority 3:	Project leaders: 1-2 FTE's/projects					
Transverse Projects	+ contributions from involved					
(examples considered)	experiments 0.2 FIES/expt.					

Data preparation:1-3 FTE/expt/2-3 years

arXiv:1205.4667

- October, 2012: CERN endorses the blueprint and appoints the DPHEP Project Manager (Jamie Shiers)
- Retain the basic structure of the Study Group, with links to the host experiments, labs, funding agencies, ICFA



Dear Dr. Diaconu,

Following the delivery of the final DPHEP blueprint, various inputs received into the European Strategy for Particle Physics symposium earlier this week and after consultation with my colleagues, I would like to inform you that CERN offers to provide the role of the initial DPHEP project manager.

We would propose to appoint Jamie Shiers in this role for an initial period of 3 years starting 1 January 2013, after which the role may be assumed by another laboratory, as suggested in the blueprint.

We would anticipate that during this period the DPHEP organization will be launched (year 1) and that the initial deliverables defined in the blueprint would be achieved.

CERN would also foresee participation in the other activities described in the document in areas such as R&D into the use of virtual machine technology for data preservation purposes (PH-SFT input to ESPP) and into the management of very large data stores.

Yours sincerely,

Sergio Bertolucci Director for Research and Computing

The DPHEP Collaboration

- Collaboration Agreement was signed in 2014
 - Give a clear sign of the will of labs to collaborate in this common challenge
- Members:
 - 2014: CERN, DESY, HIP, IHEP, IN2P3, KEK, MPP
 - 2015 IPP/Canada , 2017 UK/STFC
 - Active labs from US, Italy
 - have not formally joined, but are represented in the Collaboration Board.
- The DPHEP collaboration continue to act as an ICFA panel, as indicated in the Collaboration Agreement
 - About 60 contact persons FA, Labs, experiments
- DPHEP Activity
 - Global reports 2009(whitepaper), 2012 (blueprint), 2015, 2017 (global reports)
 - Collaboration meetings: 2015, 2017, 2021
 - Remote panel discussion March 2nd 2021
 - Reports to ICFA 2021 and 2022



Collaboration Agreement for the DPHEP Project

BETWEEN:

The Partners of the DPHEP Project (the "Partners") set out in Annex 1 to the Collaboration Agreement,

CONSIDERING THAT:

(1) Data from high-energy physics (HEP) experiments are collected with significant financial and human effort and are mostly unique;

(2) The Data Preservation and Long Term Analysis in High Energy Physics (DPHEP) project (the "Project"), an inter-experimental study group on HEP data preservation and long-term analysis, was initially formed by large collider-based experiments to investigate the technical and organizational aspects of HEP data preservation and convened by a Chair and a Project Manager as a panel of the International Committee for Future Accelerators (ICFA); Two reports were released, providing an analysis of the research case for data preservation and a detailed description of the various projects at experiment, laboratory and international levels;

(3) In its report of May 2012 (see Annex 2), the study group provided a concrete proposal for an international collaboration in charge of the Project and data management and policies in high-energy physics;

(4) The Partners have expressed their interest to take part in and contribute to the Project in order to implement the recommendations provided in the report referred to in Annex 2 and wish to formalize their collaboration through the present Collaboration Agreement;

(5) The mutual benefit of the Partners that shall result from collaboration between them;

HAVE AGREED AS FOLLOWS:

Organizational structure and decision mechanism

The organizational structure of the Project shall include the following entities:

- 1) International Advisory Committee (IAC)
- 2) Collaboration Board (CB)
- 3) Implementation Board (IB)
- 4) Project Manager
- 5) Chairperson

The DPHEP 2020 Vision

- The "vision" for DPHEP first presented to ICFA in February 2013 a 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.

DPHEP recent activities

- Remote discussion March 2021
- ICFA Mandate prolongued 2021-2024
- 3rd Collaboration meeting (remote) June 21-23, 2021
 - <u>https://indico.cern.ch/event/10431</u>
 <u>55/timetable/</u>
 - 22 contributions: experiments, dedicated projects

Panel remote discussion: March 2nd



- A decade perspective
 - The 2009 recommendations were crucial and are still valid:
 - address DP asap through dedicated projects
 - make it global via DPHEP
 - develop technologies

Scientific output from preserved data



HERA







Data Preservation projects labs: recent update

- **@DESY:** H1 (migration) and ZEUS (encapsulation) in great shape
 - successful transitions to the DP systems, publication plans continues and includes O(10) papers
 - objective: alive by 2030; New institutes joining (synergy with EIC)
- @CERN: strong LHC activity, LEP data/sw refreshed, OD/OS standards/technologies, DPHEP portal
 - Need for the continuation of the central management support
- @MPI: multi-experiment framework explored (JADE, HERA, OPAL)
 - JADE on a desktop
- **@KEK:** BELLE I data readable in Belle II framework ;
 - objective maintain Belle I data by 2023 (when the precision will be exceeded by the new data)
- @IHEP/BES3: The experiment is expected to stop data taking by 2022
 - Data to be preserved for 15 years
 - Strong support to DP national and international activities expressed
- @BNL/JLAB: DP activity ongoing (ATLAS, EIC), discussed with NPC
- **@Babar:** LTDA supported analysis since 2012. SLAC support ended in February. Data almost entirely copied to CERN/GridKa.
 - Data saved at CERN/GridKa: ~ 1.2 PB+ 0.5 PB (ongoing), Minimal user infrastructure for ongoing analyses and documentation hosted at U. of Victoria.
- @FNAL: (indirect news this time) transition to a DP system for both CDF (CDFDP) and D0 (R2DP)
 - Data stored/saved @FNAL+Italy, 500th paper from D0 in 2021

JADE

- JADE DP stack is based on open standards, does not rely on specific SW and is extremely portable. One can run it completely on desktop.
- "JADE collider experiment on your desktop".

Data Preservation model *circa* 1980-ies





2021

JADE software: recent developments

More portability, testing and documentation.

- GNU and IBM toolchains support extended with preliminary Intel^{NEW} and NAG^{NEW}. GNU is still the most stable one.
- More CI tests^{NEW}.
- Updated the site and documentation^{NEW}.
- Support for CentOS8^{NEW} and MacOSX10.15+ on x86_64^{NEW}



LEP

Date of paper

LTDP @LEP: Big Data Today - Peanuts Tomorrow New physics with Archeodata

17 January 2020

2021

2008

PHYSICAL REVIEW LETTERS 123, 212002 (2019)

Measurements of Two-Particle Correlations in e^+e^- Collisions at 91 GeV with ALEPH Archived Data

Anthony Badea,¹ Austin Baty,¹ Paoti Chang,² Gian Michele Innocenti,¹ Marcello Maggi,³ Christopher McGinn,¹ Michael Peters,¹ Tzu-An Sheng,² Jesse Thaler,¹ and Yen-Jie Lee,^{1,*} ¹Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA ²National Taiwan University, Taipei 10617, Taiwan ³INFN Sezione di Bari, Bari, Italy



On long-range pionic Bose-Einstein correlations – Including analyses of OPAL, L3 and CMS BECs –

Takuya Mizoguchi¹ and Minoru Biyajima² ¹National Institute of Technology, Toba College, Toba 517-8501, Japan ²Department of Physics, Shinshu University, Matsumoto 390-8621, Japan

February 23, 2021

Physical interpretation of the anomalous Cherenkov rings observed with the DELPHI detector

V. F. Perepelitsa ITEP, Moscow T. Ekelof Department of Physics and Astronomy, Uppsala University A. Ferrer IFIC, Valencia University B. R. French bernardfrench@blawein.ch



Babar today

T. Cartaro

Publications

- 595 papers published or submitted
 - 9 papers published in 2017, 0 8 in 2018, 4 in 2019, 6 in 2020
 - 3 in the pipeline so far in 2021, 0 few more expected later in 2021
- ~15 analyses active and on track for publication
 - Some are progressing slowly 0
 - 6 new analyses started last year and expect some 0 more this year
- 25 talks in 2021
 - 7 talks at EPS-HEP, and more already assigned 0
 - 26 talks given in 2020 (17 cancelled due to COVID-19) 0
 - Often shared talks (and collaborative analyses) with Belle 0
- Quality of physics results still excellent





But: SLAC LTDA decommissioned, moving to U. Victoria/CERN/CC-IN2P3/Grid-Ka **Open Data decided** 24

HERA: succesful DP, towards open

data

• H1: "Level 4" DPHEP strategy

- All data, full migration, including regular recompilation/validation
- Recent "technology jump" succesfull : in line with modern tools
 - "LHC"-like tools, ready for opendata

- ZEUS : "Level 3/4" DPHEP strategy
 - Root ntuples produced in the preparatory phase
 - easy to maintain/use/test/open

Synergy with future experiment: EIC

'H1Red' for simulated Pythia8.3 event





New topics/collaborators (EIC)



HERA C 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 nextgeneration <u>Electron-Ion Collider</u> (EIC).
- The EIC to be built at Brookhaven National Laboratory in partnership with Thomas Jefferson National Accelerator Facility – will be a powerful and versatile new machine capable of colliding high-energy beams of polarized electrons with a wide range of ions (or charged atoms) across many energies, including polarized protons and some polarized ions.

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 (\rightarrow)



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

https://newscenter.lbl.gov/2022/10/25/solving-the-proton-puzzle/

LHC Data Preservation

- Data Preservation and Open Access policies (already since 2012-2014)
 - DP is a « specification » included in the computing models and plans for upgrades
 - HEP Software Foundation Roadmap
- Strong initiative on Open Data and Open Science policy
- Concrete implementation and technology-oriented survey
 - Very active multi-experiment projects
 - data re-use, réanalysis, réinterpretation, outreach etc.
 - OpenData, Analysis Preservation, REANA...

A Roadmap for **HEP Software and Computing R&D** for the 2020s

HEP Software Foundation

arXiv:1712.06982

ime

tor

les.

ire.

ing

her.

for

CERN announces new open data policy in support of open science

A new open data policy for scientific experiments at the Large Hadron Collider (LHC) lers will make scientific research more reproducible, accessible, and collaborative

11 DECEMBER, 2020

nature physics

Explore Content Y Journal Information Y Publish With Us N

nature > nature physics > perspectives > article

https://www.nature.com/articles/s41567-018-0342-2

Perspective Open Access Published: 15 November 2018

Open is not enough



Other experiments expressed clear intention to join : LEP, JADE, H1/ZEUS, BaBar (HR is an issue)



Figure 2: CMS data release timeline.

Towards more standards

EDM4hep: the common language

- The Event Data Model describes the structure of the data
 - Challenge: can we have the same for all HEP experiments? LCIO shared by ILC and CLIC
- Heavily inspired by LCIO and FCC-edm



0

key4hep / EDM4hep and DPHEP?

- Key4hep / EDM4hep: framework with longer perspective than a single experiment
 - Not just another data format, but one that might become a standard
- Requires "migration", which may be a pain or not even possible
 - Workpower / Experts missing

CERNVM: the "freezer"

- Encapsulation may help here, both for migration and validation
- For LEP data, FCC-ee may provide a unique opportunity
 - Share to center-of-mass energies: 91.2 GeV, 160 GeV
 - Clear advantage in looking at what real data look like to understand bottle necks and limitations
 - Possible student projects
 - ALEPH: early investigations promising
 - ALPHA++ provides the relevant code for migration
 - Several ALEPH experts involved in FCC-ee studies

CERN Open Data portal: Status

configuration, event display

Otiborsimko

panying code for both education and research

Size: over 7600 records, over 900M files, over 2.4 PB

ATLAS released 13TeV educational samples

Non-LHC physics: OPERA neutrino physics data;

interest from PHENIX (RHIC/BNL); JADE tests

CMS completed 2010–11 proton-proton data; released

half of 2012 data; released 2010-11 heavy-ion data samples and corresponding pp reference datasets

Data science and Machine Learning (CMS, LHCb...)

Purpose: sharing event-level particle physics data and accom-**Content:** collision & simulated & derived datasets, software Explore more than two petabytes tools, analysis examples, VMs and containers, documentation, of open data from particle physics https://opendata.cern.ch ing file 1 of 11 -> File: ./5500/BuildFile.xm > Progress: 0/0 kiB (100%)
 > Verifying file BuildFile.xml expected size 305, found 305 cted checksum adler32:ff63668a, found adler32:ff63668a Command-line client to ease data download 2/3



CERN Open Data portal: Plans

December 2020: A common statement on the open data policy by CERN management and ATLAS, ALICE, CMS, LHCb and TOTEM experiments.

https://opendata.cern.ch/docs/cern-open-data-policy-for-lhc-experiments

- Prepare for forthcoming increase in open data publishing.
- Introduce flexible hot/cold disk/tape storage solution. Part of dataset files on disk, part on tapes.
- Simplify ingestion and exposure of experiment datasets (Rucio, Dirac).
- Automatise provenance testing and usage examples. (See the next presentation with REANA status overview.)

CERN announces new open data policy in support of open science

A new open data policy for scientific experiments at the Large Hadron Collider (LHC) will make scientific research more reproducible, accessible, and collaborativ 11 DECEMBER, 2020



2020. The four main LHC collabo ions (ALICE, ATLAS, CMS and LHCb) h en data policy for scientific experiments at the Large Hadron Collider (LHC), which wa ented to the CERN Council today. The policy commits to publicly releasing so-called level 3 scientific data, type required to make scientific studies, collected by the LHC experiments. Data will start to be released ion, and the aim is for the full dataset to be publicly available by the close ned. The policy addresses the growing movement of open s ducible accessible and collaborative

fic research in particle physics, as well as research in of scientific computing, for example to improve reconstruction or analysis methods based on machine lea echniques, an approach that requires rich data sets for training and validation

30

CERN Analysis Preservation and Reusable Analyses

nature physics

Explore Content V Journal Information V Publish With Us V

nature > nature physics > perspectives > article

Perspective | Open Access | Published: 15 November 2018

Open is not enough

Xiaoli Chen, Sünje Dallmeier-Tiessen 🗁, Robin Dasler, Sebastian Feger, Pamfilos Fokianos, Jose Benito Gonzalez, Harri Hirvonsalo, Dinos Kousidis, Artemis Lavasa, Salvatore Mele, Diego Rodriguez Rodriguez, Tibor Šimko 🗠, Tim Smith, Ana Trisovic 🗁, Anna Trzcinska, Ioannis Tsanaktsidis, Markus Zimmermann, Kyle Cranmer, Lukas Heinrich, Gordon Watts, Michael Hildreth, Lara Lloret Iglesias, Kati Lassila-Perini & Sebastian Neubert

- CAP : preserve analysis
 http://analysispreserva
 - tion.cern.ch/
- REANA : improve workflow
 - Run research data analyses on containerised compute clouds
 - <u>http://reana.io/</u>

CERN Analysis Preservation framework

Purpose: capture and preserve all elements needed to understand and reuse an analysis even several years later; take a consistent snapshot linking all the knowledge
Usage: describe analysis + deposit n-tuples, code etc via CLI and web UI + share with colleagues = preserve knowledge
Community: pilot with ALICE, ATLAS, CMS, LHCb

- content restricted to collaborations
- metadata interconnected with collaboration databases
- associated knowledge, e.g. CMS statistics questionnaire
- helps addressing increasing number of funding agencies asking for comprehensive data management policies
- run by CERN Scientific Information Service (P. Fokianos, K. Naim)

Øtihorsimko

REANA reproducible analysis platform

Purpose: run declarative computational workflows on containerised compute clouds

Usage: data + code + environment + workflow = computational reproducibility

Community: pilot examples with ALICE, ATLAS, CMS, FCC, LHCb; ATLAS search groups (SUSY, EXOT, HDBS) now require workflow preservation as mandatory for analysis approval

- promotes pre-producibility during active analysis phase to facilitate future preservation
- ▶ integration with GitLab; CI/CD mode
- verification of analysis examples and data provenance chain (CMS AOD reprocessing)
- support for hybrid compute workflows with multiple backends (HTCondor, Kubernetes, Slurm)

\frown		Apri and priCS the describe free of	grangt doddie dit navanen' i he do	Remethid cross sections in app collisions at 2 TeV $~\Lambda$ with dimensional transmission of the promptical dimensional transmission ρ^2 and identified splits (c) cross sections of the promptical ρ	iptent Therap
17		Open Daris valid Tris a fire validation	Han Ma MuManha	ar 2000 - S. 12 Mai and McNenter Agery detects. From to construct from an internet ONDER 4, 1, A series both 3.	- second Disease
744		Search for Myris			- updated
turbah 11		Search for W -	in /		
Passed 4		A specifier a sea	heavy gauge beaut it	Falcouring to an electron or mean and allow mean resulting is presented. The analysis was 2002 detection) weater age
		Search for Black	rates /	g collector of a contro of mass energy of F for its person had. The data sample corresponds to an integr.	e anato ap
	manufacture at 8 and 12 hor	,		des antifest?	
rational continuing of year's	n:) pearled to Antio mine	or decemp to the .	2 sents apr	Search for Verill the DLL tax, so indition from more	
CERN Analysis I	Preservation ^{III TA}	anti-tests. /	upticul 2 antici api	workflow#3 Open Sata wideleter Ph. McMeeting 2003	tests.
Less struct			ietend	dence workflow#5	testa.
CADI STATUS		par parter 1 d'de.			
- PAS	AUX	/	updated k seens age	workflowd2 Open Dan watalater MuManter 2000	trate
C (manual	100				
0	78	franci pe señel.	untered s untere a po		
Conglithulas	34				
PRE-APP	31				
PAS-PUB	24				
AdCounters 1	24	RECEIVE CHARGES	Name of the OWNER.	seemd to eveno Software. Conset: About: Status	
Their Assessed					
Cult-anded	17				
CHS WG		\$ cap-	client file	s listpid/-p <existing pid=""></existing>	
H	304				
Ц **	298				
U 100	221			n": "md5:f0428126e7cf7b0d4af7091c68ae2a	a9f",
9.5	211		filename'	": "file.ison".	
- HN	130			25	
90	130			EQ. FO L - (] 47-F DOTL 1/3 //01F/- 71	
0 104	104		10:200	52620-D600-4782-0410-1120101218C7	
C 105		1			
8					
8			checksun	n": "md5:926fb9c44251d70614ee42d34c53i	65Ъ6*.
<u> </u>			filename' filesize*:	': "Analysis_Notes_07112019.pdf", 160898,	
			id": "897	43c9b-106d-4235-8e96-23a164c7b1f4"	

https://analysispreservation.cern.ch

2/3



- MARE	
CATAR	

1910-1-1-2-5-5	

https://www.reana.io





REANA running on supercomputers (e.g. NERSC)

A word on FAIR

- The DPHEP objectives (2012) intrinsically comply with what has became to be known as FAIR principles (2016)
- Indeed, the data has to be

M. Wilkinson *et al.*, "The fair guiding principles for scientific data management and stewardship", *Scientific Data* Article No.160018 no. 3, (2016). 10.1038/sdata.2016.18.

- easy to find (F)
- accessible (A)
- and therefore -in a HEP collaborative context- (re)usable (R).
- The interoperability (I), identified as one of the long term goals ten years ago, is becoming a built-in specification of the recent computing systems as well.
 - Concrete steps have been achieved, with a few examples given, with a strong incentive originating from the open science policy or within structural projects such as WLCG.
- However, a clear strategy for a FAIR approach over the entire HEP field (including past, present and future experiments) is still to be defined.

- DPHEP can certainly contribute to such a global approach

Situation and trends

- Significant/measurable impact of dedicated DP projects @expts./labs
 - Production of high quality and unique scientific results at very low (non-zero) cost
 - 10% output for less than 1% investment: ✓
 - Long term organisation proves to be productive
 - Signs of re-vigorating collaborations in the context of new projects
 - HERÁ-EIC; LEP-FCCee
 - Case for longer term preservation: data sets parking
 - CDF, DO, Babar, LEP, Jade : carefully follow the usability in time
- LHC exps. very active in DP and **Open Data/Science**²
- 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)
- Global status report
 - Several remote editorial meetings
 - Quite advanced, to be released by Summer 2022



THE EUROPEAN PHYSICAL JOURNAL C

Regular Article - Experimental Physics



Impact of jet-production data on the next-to-next-to-leading-order determination of HERAPDF2.0 parton distributions



DESY 21-130, ISSN 0418-983

Measurement of lepton-jet correlation in deep-inelastic scattering with the H1 detector using machine learning for unfolding

> H1 Collaboration^{*} (To be submitted to Physical Review Letters) (Dated: August 30, 2021)



(Dated: August 30, 2021) The first measurement of lepton-jet momentum imbalance and azimuthal correlation in leptonproton scattering at high momentum transfer is presented. These data, taken with the H1 detector at HERA, are corrected for detector effects using an unbinned machine learning algorithm (OMNIFOLD), which considers eight observables simultaneously in this first application. The unfolded cross sections are compared to calculations performed within the context of collinear or transverse-momentum-dependent (TMD) factorization in Quantum Chromodynamics (QCD) as well

folded cross sections are compared to calculations performed within the context of collinear or transverse-momentum-dependent (TMD) factorization in Quantum Chromodynamics (QCD) as well as Monte Carlo event generators. The measurement probes a wide range of QCD phenomena, including (TMD parton distribution functions and their evolution with energy in so far unexplored kinematic regions.

> Data Preservation in High Energy Physics Status Report, Perspectives, and Plans DPHEP Global Report 2022

> > DPHEP Collaboration

Abstract

This document contemplates more than ten years of experience and global effort a pursue the preservation of data accumulated at large collider experiments.

Business as ... (not) usual



D. vom Bruch

Largest single internet exchange point: 14 Tbit/s





LHCb experiment @ CERN 40 Tbit/s



Next steps

- Objectives 2021-2024:
 - improve the awareness and stimulate improvements in Data Preservation
 - Scientific motivation, organisation, technologies, standards, outreach and education
 - Organise Workshops / issue Global Reports, link to other communities
 - reinforce and support the ongoing laboratory/experiment-based projects and their cooperation
 - keep alive data sets that (can) still produce science, keep track on parked data sets
 - support/develop the DP aspects for future experiments and encourage the transfer of knowledge
 - encourage open data and open science as a way to preserve data and knowledge
- Plans
 - Réinforce Laboratory and FA contacts
 - Release GR2022 and discuss it widely
 - DPHEP Workshop : September 2022
 - "Data Preservation for the future HEP+neutrinos"
- CERN support needed: focal point of ongoing major experiments and computing standards
- What would be the role of the future (EIC, FCC, etc.) experiments?



Big Scientific Data

Scientific research observes a dramatic increase in data and are questioning the long term future of this data



		Monday, 21 June		WEDNESDAY, 23 JUNE						
14:00 → 1	15:15 The Lar	dscape	09:00 → 11:30	Service /	Project Updates: Service/Project updates					
	14:00	Workshop Introduction and Goals Speakers: Cristinel Diaconu (CPPM, Abe-Marseille Université, CNRS/IN2P3 (FR)) , Dirk Duellmann (CERN)		09:00	CERN Open Data Portal Speaker: Tibor Simko (CERN)					
	14:10	An Open Data policy for LHC Speaker: Jamie Boyd (CERN)		09:15 CERN Analysis Preservation Speaker: Pamfilos Fokianos (CERN)						
	14:30	Welcome from CERN Management Speaker: Joachim Josef Mnich (CERN)		09:30 REANA Reproducible Analyses Speaker: Tibor Simko (CERN)						
	14:40	From LEP to fcc-ee. A bridge too far? Speaker: Jamie Shiers (CERN) DDUED dates and DDUED		09:45	Bit Preservation Speaker: Oliver Keeble (CERN)					
		Computing - t		10:00	Coffee					
15:15 → 1	15:35	Coffee		10:20	"Software Preservation" - virtualisation et al. Speaker: Jakob Blomer (CERN)					
15:35 → 1	18:20 Site / E	xperiment "Round table": Part 1		10:40	Managed service migration/retirement					
		ATLAS activities & plans Speakers: Lukas Alexander Heinrich (CERN), Marumi Kado (Sapienza Universita e INFN, Roma I (IT))		11:00 Discussion						
		CMS activities & plans	11:20 → 19:00	Site / Exp	periment "Round table": Part 2					
		Speaker: Kati Lassila-Perini (Helsinki Institute of Physics (FI))		11:20	BES III Speaker: Lu Wang (Computing Center,Institute of High Energy Physics, CAS)					
		LHCb activities & plans Speaker: Adam Morris (University of Bonn (DE))		11:40	KEK / Belle I & II Speaker: Takanori Hara (High Energy Accelerator Research Organization (JP))					
		ALICE activities & plans Speakers: Jochen Klein (CERN), Stefano Piano (INFN (IT))		12:00	Lunch					
	16:55 17:35	DESY / HERA H1 activities Speaker: Daniel Britzger (Max-Planck-Institut für Physik München) ZEUS activities Speaker: Achim Geiser (Deutsches Elektronen-Synchrotron (DE)) BaBar Speaker: Concetta Cartaro (SLAC)		14:00	LEP Session ALEPH, DELPHI and OPAL Status and Plans Speakers: Gerardo Ganis (CERN), Marcello Maggi (Universita e INFN, Barl (IT)), Matthias Schroeder (CERN), Ulrich Schwickera Opportunities offered by LEP data@edm4hep for future EW and Higgs factories Speaker: Gerardo Ganis (CERN) Discussion	th (CERN)				
	17:55	Discussion		14:45	CERNLIB					
			.	15:15	Coffee					
<u>3</u>	rd DP	HEP Collab. Meeting		15:35	FNAL / Tevatron					
<u>h</u> 1	ttps:/	//indico.cern.ch/event/1043155/		15:55	BNL / RHIC Speaker: Maxim Potekhin (Brookhaven National Laboratory (US))					
				16:15	Discussion on DPHEP futures	38				
				17:45	CB Meeting					

Scientific Data: what is it?



- Publications
- Documentation
- > Raw
- Processed data
- > Meta-data
- > Workflows
- Software (all branches)
- > Diffuse knowledge
-more...

C. Diaconu | DPHEP Status and perspectives

Technology, methodology Organisation

Ч С

0

ess

Set

DPHEP ressources for DP

• 2012 Blueprint

	Project	Goals and deliverables	Resources and timelines	Location, possible funding source, DPHEP allocation				
laboratory	Experimental Data Preservation Task Force	Install an experiment data preservation task force to define and implement data preservation goals.	1 FTE installed as soon as possible, and included in upgrade projects	Located within each computing team. Experiment funding agencies or host laboratories. DPHEP contact ensured, not necessarily as a displayed FTE.				
Experiment and Priority: 1	Facility or Laboratory Data Preservation Projects	Data archivist for facility, part of the R&D team or in charge with the running preservation system and designed as contact person for DPHEP.	1-2 FTE per laboratory, installed as a common resource.	Experiment common person-power, support by the host labs or by the funding agencies as a part of the on going experimental programme. A fraction 0.2 FTE allocated to DPHEP for technical support and overall organisation.				
	General validation framework	Provide a common framework for HEP software validation, leading to a common repository for experiments software. Deployment on grid and contingency with LHC computing also part of the goals.	1 FTE	Installed in DESY, as present host of the corresponding initiative. Funding from common projects. Cooperation with upgrades at LHC can be envisaged. Part of DPHEP.				
	Archival systems	Install secured data storage units able to maintain complex data in a functional form over long period of time without intensive usage.	0.5 FTE	Multi-lab project, cooperation with industry possible. Included in DPHEP person-power.				
	Virtual dedicated analysis farms	Provide a design for exporting regular analysis on farms to closed virtual farm able to ingest frozen analysis systems for a 5-10 years lifetime.	1 FTE	The host of this working group should be SLAC. Funding could come from central projects and can be considered as part of DPHEP.				
	RECAST contact	Ensure contact with projects aiming at defining interfaces between high-level data and theory.	0.5 FTE	Installed with proximity to the LHC, the main consumer of this initiative, with strong connections to the data preservation initiatives that may adopt the paradigms.				
	High level objects and INSPIRE	Extend INSPIRE service to documentation and high-level data object.	0.5-1.5 FTE	Installed at one of the INSPIRE partner laboratories.				
Multi-experiment Priority: 3	Outreach	Install a multi-experiment project on outreach using preserved data, define common formats for outreach and connect to the existing events.	1 FTE central + 0.2 FTE per experiment	A coordinating role can be played by DPHEP in connection with a large outreach project existing at CERN, DESY or FNAL. The outreach contributions from experiments and laboratories can be partially allocated to the common HEP data outreach project and steered by DPHEP.				
Global Priority: 2	DPHEP Organisation	DPHEP Project Manager	1 FTE	A position jointly funded by a combination of laboratories and agencies.				

Table 8: Resources required by projects of the DPHEP study group.



2018 status







BABAR needs Help!

- BABAR data actively being analyzed and high impact papers published (see slide 2). Expect this to continue to at least through 2021.
- SLAC management plans to stop hosting BABAR computing in February 2020 at which time the tapes with data will be ejected.
- DOE support ended in 2017, now running on international common funds (OCF).
- Looking for possibility of support and long term data preservation at
 - CERN,
 - GridKa (BABAR site for analysis and XRootD federated dataset main redirector),
 - University of Victoria (BABAR site for analysis, documentation, and tools support).
- BABAR lightweight VMs come with the latest software release and xrootd client included, running under the most common virtual machine players. Just add the data via the GridKa main XRootD redirector.

BABAR in Numbers

- 2PB of data on T10k-D tapes
 - raw, processed, Monte Carlo
 - Unique dataset at the Y(3S) resonance (no plan at the moment to run at the Y(3S) @ Belle II)
- Full environment enclosed in VMs (SL5,SL6)
- ~1TB of documentation, repositories, and dataset information (DBs, cvs, wiki, html)
 – Internal documents archived on INSPIRE

- 574 papers, ~10 papers/year past 3 years
- 231 members (semi-frozen author list)
 - Including PhD students in Canada, Germany, Israel, Italy, Russia, US
 - Associated theorists mine data to test new ideas
- ~20 analyses on track, ~10 more in the pipeline
 - Continue to have new analyses every year including joint *BABAR* -Belle analyses
- Students analyze BABAR data while working on Belle II and other experiments in construction/commissioning phase

HERA: succesful DP, towards open



- H1 has unique data and continues to produce physics publications, long after data taking ended
- A recent software modernisation program has been performed to allow this to continue using modern analysis tools, recent programming languages and on state-of-the-art platforms D. Britzger and D. South, H1 Data Preservation Status, DPHEP Preparatory Meeting, 2 March 2021



many EIC topics common with HERA



some EIC members have recently joined ZEUS to work on common analysis topics with real ZEUS data

Common Ntuple analysis model ZEUS

ZEUS Common Ntuple:

Motto: keep it simple!

flat (simple) ROOT-based ntuple (same format as PAW ntuple converted with h2root) containing high level objects (electrons, muons, jets, energy flow objects, ...) as well as low level objects (tracks, CAL cells, ...)

date:	4-06-2006 time: 00:06:30
E,=52.8 GeV	E _b =2.07 GeV
p,=0.583 GeV	p_=52.1 GeV
t,=-100 ns	t_=2.97 ns

Well tested |

almost all recent ZEUS papers (24 out of 25) based on Common Ntuples

Easy to maintain

transition sl5 -> sl6 -> sl7 completely transparent (just use newer ROOT version)

"Easy" to use

most recent ZEUS papers based on results produced by master students, PhD students or postdocs from remote institutes, e.g. related to EIC or Heavy Ion communities, using resources at DESY or MPP: analysis on DESY NAF/BIRD computing farm or at MPI/Garching



Low threshold for access to data by external groups

02.03.21

veriment: EIC

-

A. Geiser, DPHEP meeting

ZEUS physics papers





majority of papers produced in "data preservation mode" already since 2012 (25 papers)

since end of DESY funding 2014:

2015-20: 14 papers, 1 with > 500 citations 2021: expect 2-4 papers

long term: ~1-2 papers/year -> ~2030

expect ~10% of total ZEUS output

~80-90% of these would never exist without dedicated data preservation

ZEUS data preservation program is a success! some small official resources could double the output and/or allow Open Data

02.03.21

2

H1 DP



Figure 3: Left: Number of Monte Carlo events produced centrally by the H1 Collaboration. The years without MC production are related to a change of the computing environment, or no MC requests. Right: Number of H1 authors is increasing since 2019 due to retained analysis capabilities and new interest in *ep* physics. The colored areas indicate the data taking period (green), the period with active funding (yellow) and the period under the new collaboration agreement in *data preservation mode* (cyan). The number of corresponding publications is also indicated.



Figure 4: Number of ZEUS papers published and anticipated to be published per year (will be updated to 2022 version).

2018 status

DPHEP timelines

	2007	2008	2009	2010	2011	2012	2013	2014	2015	2016	2017	2018	2019	2020	2021	2022
НЕР	HERA tops	Babar stops	LHC starts	Belle I stops	Tevat ron stops				LHC Run 2							
DPHEP Groun			ICFA Panel		LHC exp. joine d	DPHEP Manger appoint ed at CERN		DPHEP Collaboration Agreements signed	1 st DPHEP Collab. Meeting		2 nd DPHEP Collab. Meeting				3 nd DPHEP Collab. Meeting	
DPHEP Docs			DPHEP White Paper			Bluepri nt Report			DPHEP Status Report 2020 Vision		DPHEP 2017 Status Report				ICFA mandate Renewed	DPHEP Status Report (in work)
DP Projects within expts.		Babar DP starts		HERA DP starts	BELLE DP starts	CMS DP Policy CDF/D0 DP starts Babar LTDAP operatio nal	ALICE, LHCb, DP Policies	ATLAS DP Policy H1/ZEUS DP systems operatio nal	CERN/LH C Open Data	CERN/ LHC Analys is Preser vation Tevatr on DP operat ional					CERN open science policy	

Start-up

Consolidation

DPHEP Collaboration

Preservation: where is the problem?



NATURE | NEWS

عربي

LHC plans for open data future

Researchers share results to keep them accessible.

Elizabeth Gibney

26 November 2013

"When the LHC programme comes to an end, it will probably be the last data at this frontier for many years. We can't afford to lose it." Storing the data is not a problem: hard drives are cheap and getting cheaper. The challenge is preserving knowledge that is less commonly stored — the software, algorithms and reference plots specific to each experiment. These often degrade or disappear with time, says Cristinel Diaconu of the Marseilles Centre for Particle Physics in France, who is chair of the international Data Preservation in Long Term Analysis in High Energy Physics (DPHEP) study group. He worries that if the data continue to be stored in their current state, physicists trying to decipher them in 10 years' time will be unable to reconstruct the discovery of the Higgs boson. "When the LHC programme comes to an end, it will probably be the last data at this frontier for many years," he says. "We can't afford to lose it."