

DEUTSCHES ELEKTRONEN-SYNCHROTRON
Ein Forschungszentrum der Helmholtz-Gemeinschaft



DESY 21-097
MPP-2021-87
arXiv:2106.11058
June 2021

**Preservation through Modernisation:
The Software of the H1 Experiment at HERA**

D. Britzger

Max-Planck-Institut für Physik, München

S. Levonian, S. Schmitt, D. South

Deutsches Elektronen-Synchrotron DESY, Hamburg

ISSN 0418-9833

NOTKESTRASSE 85 - 22607 HAMBURG

DESY behält sich alle Rechte für den Fall der Schutzrechtserteilung und für die wirtschaftliche Verwertung der in diesem Bericht enthaltenen Informationen vor.

DESY reserves all rights for commercial use of information included in this report, especially in case of filing application for or grant of patents.

To be sure that your reports and preprints are promptly included in the
HEP literature database
send them to (if possible by air mail):

DESY Zentralbibliothek Notkestraße 85 22607 Hamburg Germany	DESY Bibliothek Platanenallee 6 15738 Zeuthen Germany
---	---

Preservation through modernisation: The software of the H1 experiment at HERA

Daniel Britzger¹, Sergey Levonian², Stefan Schmitt², and David South²
for the H1 Collaboration

¹Max-Planck-Institut für Physik, Föhringer Ring 6, 80805 München, Germany

²Deutsches Elektronen Synchrotron (DESY), Notkestr. 85, 22607 Hamburg, Germany

Abstract. The lepton–proton collisions produced at the HERA collider represent a unique high energy physics data set. A number of years after the end of collisions, the data collected by the H1 experiment, as well as the simulated events and all software needed for reconstruction, simulation and data analysis, were migrated into a preserved operational mode at DESY. A recent modernisation of the H1 software architecture has been performed, which will not only facilitate on going and future data analysis efforts with the new inclusion of modern analysis tools, but also ensure the long-term availability of the H1 data and associated software. The present status of the H1 software stack, the data, simulations and the currently supported computing platforms for data analysis activities are discussed.

1 The H1 experiment at HERA

Operating during the years 1992 to 2007, HERA at DESY is so far the only high energy lepton–proton (ep) collider in the world to have been constructed, where 27.6 GeV electrons or positrons were brought into collision with 920 GeV protons, resulting in a centre-of-mass energy of 319 GeV. The collision of point-like leptons with hadrons made HERA a unique tool for precise measurements of the structure of the proton. Many other areas of particle physics were also accessible at HERA, including QCD and jets, heavy quark production, diffraction, electroweak physics, as well as the search for rare processes in ep collisions.

The H1 detector [1, 2] at HERA recorded the final state particles of ep collision events, and features tracking detectors closest to the beam pipe, surrounded by electromagnetic and hadronic calorimetry, a muon system, and several further subdetector components. Approximately 270,000 readout channels were employed by the H1 detector. A multi-level trigger system was employed to reduce the event-rate from the bunch crossing frequency of 96 ns (≈ 10 MHz), and selected events were then stored with a rate of 20–50 Hz in the RAW data format. The total volume of RAW ep collision data recorded by the H1 detector and suitable for analysis amounts to about 75 TB, and comprises approximately 1 billion events collected in the years 1996–2007.

Considering the planned Electron–Ion Collider in the US (EIC) [3], the proposed Large Hadron-electron Collider at CERN (LHeC) [4] and the proposed Electron-Ion Collider in China (EicC) [5], as well as many new related theoretical developments, the unique ep data from HERA retain their relevance for many years to come.

2 The H1 software stack

The H1 Collaboration developed and maintains a sizeable software stack, which was used during data taking, and continues to be used not only for subsequent data reconstruction and high-level data analysis, but also for the simulation of the detector response to high energy physics (HEP) processes. The software stack is also crucial to perform physics analyses of the H1 ep collision data. The relevant components of the software stack are briefly described in the following.

2.1 The core packages

The H1 core software is written almost entirely in FORTRAN 77, was designed in a machine-independent manner¹, and has served the experiment well since 1988. It is this software which creates the basic format, the Data Summary Tape (DST) from the RAW data, as well as providing the means to simulate ep collision events for Monte Carlo (MC) comparisons to data.

The FORTRAN software is based on two key components, BOS [6] and F-PACK [7], and is arranged in a series of strictly modular packages or modules, with self-contained sets of routines and clear input/output (I/O) interfaces. BOS is a dynamic memory management system for defined data blocks and their persistent I/O. The system supports a modular structure of the application program and portability for both the software and the data sets. F-PACK is a machine-independent general I/O package for data blocks. It performs automatic word format conversions for different machine representations (IBM, VAX, DEC, IEEE) and supports fast access to subsets of data through indexed files. An important feature of F-PACK is the support of keyed/ordered access files, which is essential for data-base-like applications. The modules of the H1 packages communicate with each other strictly only via BOS banks, and larger programs, such as the reconstruction program, consist of a plain series of module calls.

The H1 reconstruction package H1REC comprises approximately 50 different sub-modules for noise suppression, clustering, calibration, tracking, vertexing, tagging, for the combination of different sub-detector quantities, as well as electron identification and the reconstruction of deep inelastic scattering kinematic variables. The detector simulation package H1SIM is based on GEANT3 [8]. It has a modular structure and contains the geometry definition of all sub-detectors, the magnetic field map, and tools for fast simulation and the shower library. Taking the relevant run conditions from an Oracle database, MC events are produced in the same format as the data, with some additional information. The same reconstruction software as used on the data is then applied to the simulated MC events.

Further modules include the H1 global database H1NDB, the data management system DATMAN, the event filter modules H1L4 and H1ECLASS, trigger simulation modules H1TRIG and H1FTTEMU, as well as libraries for useful numerical utilities H1UTIL and utilities for MC production H1MCUTIL. Finally, the LOOK package is a general system for graphics applications in physics, and uses GKS [9] for basic graphical operations with a system of utility subprograms to prepare own display programs. As such, LOOK provides the core framework for the event display, H1ED.

Only a few mandatory external dependencies exist, namely CERNLIB [10], GEANT3, GKS, oracle-instant and the FORTRAN compiler, presently gfortran. The first three packages are no longer maintained, but H1 retains a copy in its internal software stack.

¹The H1 Collaboration employed many computing platforms and operating systems to analyse data, among them AIX, AILLANT, APOLLO, AXP, Ultrix, HPUX, IBM, IBMMVS, MAC, MIPS, OS-9, RTPC, VAC, VAXMVS, VMS, UNIX, SGI, SUN, LINUX, and most recently i686 and amd64 with RHEL or Scientific Linux distributions.

2.2 Analysis data and software: H1oo

In the year 2000, a few years ahead of the HERA II phase of data taking, H1 made the decision to develop a new object-oriented analysis core framework in C++ , H1oo [11–15]. The goal of the H1oo project was to improve the overall efficiency in H1 physics analysis and performance by providing a modern analysis environment, new standardised and fast event selection facilities and a new, common data format.

The design and development principles employed in H1oo were closely aligned to the ambitious developments of the emerging ROOT [16] framework at CERN. Concepts such as (multiple) inheritance, dynamic instantiation of objects via Names, consistent Setter and Getter functions for the use with the interactive C++ interpreter (CINT) and comprehensive source-code documentation are among the similarities in the design principles. H1oo also makes use of the generic collection of tools for high-performance I/O, as well as the interactive analysis environment and graphical opportunities provided by ROOT. The use of a global singleton class to provide convenient access to all relevant quantities is also very similar to the ROOT paradigm.

The H1oo framework is structured into about 50 sub-packages and consists of more than 600 C++ classes, inheriting from the core ROOT class TObject. The highly standardised data transfer methods between the sub-packages via BOS and F-PACK in the FORTRAN software described in the previous section also provided a stable access interface to read the DST event files from H1oo in C++. A more complete discussion of the structure of the H1oo software can be found in [14].

For H1oo, new data formats were designed to define new high-level analysis objects, such as jets, heavy-quark tagging information or (re-)calibrated particle candidates, with improved data I/O via TTrees stored in ROOT files. The H1oo data format typically used for analysis is in reality comprised of two persistent file formats, to be used in tandem: the "H1 Analysis Tag" (HAT) which contains simple, calibrated variables to perform a fast event selection, and the larger "Micro Object Data Store" (mODS) which contains additional information on identified particles. A third H1oo file format, the "Object Data Store" (ODS), may be accessed transiently during an analysis event loop, and provides an interface to the full information stored on the original DST file in ROOT format. This transient ODS access is the only part of an H1oo analysis which requires the FORTRAN software.

The H1oo framework provides many utilities for data analysis, wrappers for data access, templates for analysis codes, and standards for H1 physics analyses. An event display H1Red provides visualisation tools in the ROOT analysis environment, and features full backward compatibility with LOOK and H1ED.

Whilst at some level each H1 physics analysis retains its own specific high-level analysis code and workflows, the common H1oo framework greatly helped to standardise event selection, particle and event reconstruction, calibration, production of histograms and the assessment of measurement uncertainties and systematic errors. This has had additional benefits in terms of shared analysis code, expert knowledge and working environments and, perhaps most importantly, data handling where members of the collaboration all use the same file formats.

2.3 Software access, distribution and documentation

All H1 software and source code are available to the members of the collaboration through their DESY credentials. Pre-compiled executables, shared libraries and a working environment are provided centrally via the DESY-IT infrastructure for the supported operating system(s), which is currently CentOS7. The FORTRAN code is documented in about 60 internal

H1 software notes that were written in a standardised, journal-like format. Some packages were also presented and documented in journals or conference proceedings. Each of the approximately 10 000 subroutines is stored in a separate file and includes comprehensive, standardised in-code documentation as well as a version history.

The H1oo framework comes together with its own manual, tutorials, examples and a central, internal web-page. Online HTML code documentation is provided using ROOT's THtml class and features more than 2000 individual web-pages. A single, central web-server provided by DESY-IT hosts all web-resources from all parts of the collaboration, such as meeting notes, hardware documentation, run-dependent documentation, analysis notes, trigger details, papers and so on.

3 DPHEP and the data preservation model for H1

Data taking at HERA ended in June 2007, and was soon followed by other high energy physics experiments of the same generation, namely BaBar at the PEP-II e^+e^- collider at SLAC (April 2008) and the DØ and CDF experiments at the Tevatron $p\bar{p}$ collider (September 2011). In order to perform a detailed evaluation of how to preserve high energy physics data for long-term analysis, an inter-experimental study group "Data Preservation in High Energy Physics" (DPHEP) [17] was formed at the end of 2008 to systematically investigate all technical and organisational aspects of this subject. A series of six workshops took place between 2009 and 2012 and following a short interim report in 2009 [18], a full report was released in May 2012 [19].

The complete H1 RAW collision data comprises around 75 TB, the set of compressed DST data amounts to about 20 TB and the analysis level H1oo files are about 4 TB. Other data, such as random trigger streams, noise files, cosmic data, luminosity monitor and other calibration data amounts to a few TB. The total volume of preserved simulated MC sets is around the same size of the data, and hence the total volume of preserved data is about 0.5 PB. This is about the similar volume of RAW data written out in just three days by the ATLAS detector at the LHC.

A key issue in data preservation is that in addition to the data themselves the associated software also needs to be considered, such as the programs for data access, reconstruction, simulation and analysis programs. These programs also provide an important documentation of the data themselves. In contrast, without a well defined and understood software environment the scientific potential of the data is limited, e.g. the possibility to study new observables or to incorporate new reconstruction algorithms, detector simulations or event generators may not be possible. With this in mind, a model for data preservation was devised [20], based on a series of levels of increasing complexity into which projects could be assigned [21], as shown in Table 1.

Table 1. Data preservation levels defined by the DPHEP Study Group.

Level	Preservation Model	Use Case
1	Provide additional documentation	Publication related info search
2	Preserve the data in a simplified format	Outreach, simple training analyses and data exchange
3	Preserve the analysis level software and the data format	Full scientific analysis possible, based on the existing reconstruction
4	Preserve the reconstruction and simulation software as well as the basic level data	Retain the full flexibility and potential of the experimental data

As new experimental methods, new phenomenological models or new observables, are likely to be the prime reasons for analysing event data again, scenarios arise where only a

more comprehensive preservation model will provide the necessary ingredients. For example if a parameter in the reconstruction algorithm is kinematically limiting, or a new simulation written in a modern computing language provides only incompatible interfaces. This approach assumes, justifiably, that it would be impossible to rewrite experiment specific software such as detector simulation or reconstruction code from scratch, due to missing expert knowledge about real hardware components and complexity of these programs. With this in mind, H1 followed a DPHEP level 4 preservation model [22], and preserves the analysis level data formats, as well as the most basic level (RAW) data, and all software. It is clear that this level of preservation will necessarily include the full range of both experiment-specific and external software library dependencies. However, the benefit of such a model is that the full physics analysis chain is available and full flexibility is retained for future use, similar to a real running experiment.

Whilst beyond the scope of this paper, it is also worth noting that a significant DPHEP level 1 preservation effort has also been performed. This ensures all relevant digital and non-digital documentation concerning all aspects of H1 are safeguarded for future use [22]. A dedicated web server now hosts all of the documentation of the H1 experiment through static web pages, representing a single resource for knowledge transfer for the members of the collaboration. Furthermore, a new, simpler operational model was also officially adopted by the collaboration in July 2012 to ensure the efficient long term governance via the H1 Physics Board, which is made up of experts from all areas of the experiment.

4 Preserved operational mode: 2012–2020

The final version of the reconstructed data, DST 7, was produced in December 2010 with the inclusion of the HERA I data from 1996–2000. This reprocessing campaign utilised a dedicated "dataflow meta-computing framework" [23], which reached a stable performance of 60 million events per day. This represented a factor of 20 improvement with respect to the previous architecture employed during earlier HERA I reprocessing campaigns [24]. It was not until summer 2012 that the HERA I H100 data and MC to be used for analysis were fully prepared including all relevant alignments and calibrations. The H100 production release 4.0.25, which uses ROOT5.34, represented the end of major development, and was not only used to produce the data and MC for long term analysis but would also be the software release to be used for the remainder of the decade and beyond.

The H1 computing infrastructure included a dedicated 1200 core batch system, significant storage capabilities comprising both tape and disk, and several large working group servers. As part of the 2012 transition into the new "preserved operational mode", the majority of these resources were phased out and replaced by a computing infrastructure centrally managed by DESY-IT and largely shared with other experiments. The H1 data identified for preservation (see section 3) was relocated during 2014–2015 to a dedicated, dCache [25] based storage at DESY, featuring two copies of the RAW data on different tape media as well as an online pool for access to the most popular data.

As described in section 3, the H1 strategy is to retain the full flexibility of the data using a DPHEP level 4 data preservation model. This requires that the H1 software stack and all of its dependencies are continuously maintained and repeatedly updated and tested whenever a change in either the infrastructure or external dependencies is made. These continuous migrations were made possible due to the modular structure, the high quality and stability of the packages, and H1 further benefited from the choice of stable programming languages for the software, FORTRAN 77, C and C++98, and only a very moderate usage of shell and perl scripts.

The main OS employed by DESY-IT in 2012 was 32-bit Scientific Linux DESY 5 (SLD5), and whilst this became the baseline version for H1, work was begun that year to migrate to 64-bit, an important step to achieve given that next generation OS versions would only be available as 64-bit versions. SLD6 became the standard at DESY after 2015. Each OS migration revealed certain dependencies or required additional updates from DESY-IT, for example new versions of Oracle or dCache, which required only small changes to the H1 software to retain compatibility. For example, after the transition from SLD6 to CentOS7 in 2020, the library "libdcap.so" used for direct IO access to dCache systems was removed from the DESY computing environment, resulting in corresponding, minor adjustments to the H1 software.

The use of external software in the H1oo software was reduced as much as possible in 2012 and whilst those remaining, namely ROOT, fastjet [26] and neurobayes-expert [27], were migrated to the new operating systems, the last stable versions from 2012 were kept. Table 2 shows a breakdown of the different components of the H1 software stack, and the status of the various dependencies during the preserved operational mode (2012–2020).

Whilst the migrations required only moderate person power, due to necessary validations and the overall complexity of the software stack as such they had to be performed by experienced H1 software experts. To assist in this exercise a framework [28] was developed, consisting of a series of well defined validation tests, to check for consistency every time a part of the software stack or environment was updated. Whilst not fully deployed in an automated way, this project nevertheless showed its value during the first migrations to 64-bit operating systems.

Table 2. Breakdown of the H1 software stack and its dependencies for the preserved operational mode (2012–2020).

Component	Responsible	Maintained packages	Discontinued packages
H1 software	H1	H1 core software, H1oo	–
OS dependencies (continuous updates)	DESY-IT	Oracle, dCache, web-services, compilers, GNU utilities, gmake, system libraries	CVS
External dependencies (selected fixed releases)	H1	fastjet, neurobayes-expert, MC generators	CERNLIB, GKS, GEANT3, ROOT5, LHAPDF5, MC generators

5 The modernisation of H1 software and computing for the 2020s

5.1 Revisiting the status of the H1 software

In 2020 the status of the H1 software and the data analysis capabilities were again revisited. Although the entire software was successfully migrated to CentOS7 and the full data analysis capability is retained, the overall status of the H1 software had a few shortcomings. Some of these were introduced causally over the time period of the data preservation effort, and include the following:

- The programming languages and standards (C++98) are unattractive for young people to learn and slow down new developments and data analysis efforts
- Outdated dependencies, such as ROOT5, complicate the usage of modern data analysis techniques

- Modern tools have not in general been introduced
- No link to an externally maintained package repository is present, meaning that new packages had to be provided manually
- New dependencies may be incompatible, for example due to different compiler standards like C++20, or different interfaces such as MC event generators providing an event record in newer data formats
- The compilation and maintenance of the packages still requires a profound and detailed knowledge about the specific structure of the H1 software, as well as some insight into the historic development
- No concise build instructions of the entire H1 software stack are available
- The H1 core packages are bound to the DESY-IT infrastructure and cannot be relocated

Some of these issues are very specific to the H1 software or computing infrastructure as it developed organically over three decades, but others are of a more general nature and may be repeated in data preservation efforts by other experiments in the future. In particular for younger students, who want to perform a modern data analysis and also need to learn modern computing languages and techniques, the H1 software environment was not attractive.

5.2 Restructuring the H1 software infrastructure after 10 years

At DESY, the only supported and locally deployed platform for the H1 software is CentOS7, including all of the add-ons that it brings. However, in HEP the LCG package repository [29, 30] is commonly used as standard, and many recent dependencies, as well as modern compilers, are provided by LCG/AA [31] through `cvmfs` [32] for CentOS7. H1 has decided to adopt this dependence, update all external dependencies accordingly, and to provide its software in a similar scheme to the members of the collaboration.

In a first instance, the platform `x86_64-centos7-gcc9-opt` and the LCG version `LCG_97a` are selected for the next stable H1 software release. Many external packages are now provided through LCG, such as `ROOT`, `fastjet`, `neurobayes-expert`, thus reducing the maintenance cost to H1, as can be seen in Table 3. This may be contrasted with the previous situation shown in Table 2.

Table 3. Breakdown of the H1 software stack and its dependencies from 2020 onwards.

Component	Responsible	Maintained packages	Discontinued packages
H1 software	H1	H1 core software, H1oo	–
OS dependencies (continuous updates)	DESY-IT	Oracle, dCache, web-services, GNU utilities, git, gmake, system libraries	–
External dependencies (selected fixed releases)	H1	–	CERNLIB, GKS, GEANT3 (selected) MC generators
External dependencies (selected regular updates)	LCG	LHADPF6, ROOT6, compilers, fastjet, neurobayes-expert, MC generators, (and as back up option: Oracle, dCache, Git)	–

The transition to the GNU compiler collection 9.2.0, in contrast to the previously used version 4.8.5, required a rebuild of all of the core FORTRAN packages from scratch. Due to the high quality of those packages, the stability of the FORTRAN 77 standard, and the many previous migrations of these codes, this transition was accomplished without complications, and only one case of memory corruption and few incorrect bit-wise logical operators needed

to be fixed in about 950k lines of code (LOC). The discontinued packages CERNLIB, GKS² and GEANT3 could also be sourced. The important low-level packages BOS and FPACK, which are needed for the interface between the H1 core packages and the access to the data, have been tested successfully.

More important for data analysis is the migration of the H1 analysis framework H1oo, which is 300k LOC written in C++98 and where the last stable version was built in 2012 upon ROOT5.34³. The migration of H1oo to ROOT6 was non-trivial, since for efficient data access, H1oo makes extensive use of dynamic instantiation of named classes and dynamic scopes. This functionality was provided in ROOT5 through CINT, which was replaced by CLING in ROOT6. Consequently, H1oo had to strictly use the respective interfaces from ROOT6 and the direct access to CINT functions had to be omitted. A single class out of about 600 H1oo classes could not be re-factored and was dropped, but this did not limit the usability of the analysis framework. On the other hand, the new C++17 standard permits the implementation of convenient range-based for loops for H1oo array classes, thus simplifying the code. This feature facilitates the processing of particle lists in single events. The object-oriented file formats mODS and ODS remain compatible, but a single object-oriented database file had to be regenerated. Latest releases of fastjet v3.3.2 and neurobayes-expert v3.7.0 were integrated and successfully tested.

A complete release of all of the H1 software is provided to the collaboration in a single directory on afs, which is mirrored to cvmfs and the internal nfs shared file system at DESY. This directory contains the binaries, libraries and the include files for all H1 software packages, together with compatible database snapshots. A small setup script also initialises the compatible LCG release LCG_97a/x86_64-centos7-gcc9-opt. Compatible Oracle and dCache libraries are further provided from the LCG mirror, and serve as a backup option if the DESY-IT libraries fail.

5.3 Python-based and interactive data analysis with H1oo

One great benefit from the close integration of H1oo into the ROOT-ecosystem is obtained from the PyROOT developments in ROOT6 [33]. Through ROOT's automatised Python-C++ bindings, the full functionality of H1oo is available from Python. Together with the range-based operators in H1oo, this enables a fully pythonic analysis of the H1 data and makes use of the H1oo analysis framework. Interactive data analyses in Python, or using ROOT's C++ interpreter, are also possible, and may become a valuable option for future training or outreach activities.

Since H1oo is provided through a global shared file system (cvmfs or afs), similar to the LCG releases, it is fully relocatable and H1 data analysis or MC production can now be performed on any CentOS7 system anywhere in the world. Alternatively, standardised CentOS7 containers, for example from CernVM [34], can be employed on other platforms. For a quick start for students or newcomers, a few new code examples are provided to H1 Collaboration members, both in C++ and Python.

5.4 Future maintenance of the H1 software stack

Future migrations of the H1 software will be connected to given versions of LCG releases, featuring newer platforms, dependencies and compiler versions. In order to facilitate this dependence a few additional developments were done.

²GKS was obtained from <https://github.com/Starlink/starlink/tree/master/libraries/gks>

³H1oo was initially developed using ROOT version 2.

All code repositories were migrated to `Git`, which is the current standard for code management, an important transition given that most younger developers are unfamiliar with CVS (or going even further back, CMZ). The DESY-IT central `Git` repository hosting service `Bitbucket` [35], which provides authentication for H1 members through their DESY credentials, replaces the old CVS web tools, thus reducing the maintenance cost for H1.

A new set of build instructions has been developed that allow to rebuild the entire H1 software stack from scratch. Historic dependencies to central resources were removed and the build system was further simplified. The build system remains dependent mainly on `gmake`, which has proven to be very reliable over a long period of time. Standardised software tests are currently under development. About 50 examples of analysis code are also kept in the `Git` archive alongside the official H1 code, serving as a valuable reference for ongoing activities and for new data analysts.

5.5 Singularity and containerised workflows

The advent of containerisation and `Singularity` [36] provides great opportunities for the preservation of data and software in HEP. As binaries and libraries of the H1 software were retained in the current storage systems deployed by DESY-IT on together with the complete set up procedures, these can now conveniently be used with `Singularity` and compatible SLD5 or SLD6 containers. Both SLD5 and SLD6 binaries have been tested and full functionality has been proven with standard containers from DESY-IT. The use of `Singularity` represents the retrospective realisation of the DPHEP level 3 ‘encapsulation’ of the H1 software.

It is worth noting that whilst the use of containers and virtualisation offers potential solutions to data preservation, further development can only take place within the constraints of the original environment and technologies. The modernisation program described in this paper has not only increased the longevity of the H1 software but undoubtedly made it more accessible and attractive to the next generation of collaborators.

6 Summary and Outlook

The H1 experiment recorded a unique data set of lepton–proton collisions in the years 1992 to 2007 and developed a sizeable software stack for its processing and analysis. The H1 Collaboration maintains these data, the related software and the documentation. Several data analysis activities are still ongoing and new analysis projects are beginning due to the uniqueness of this scientific data set and the increasing interest of the HEP community in *ep* scattering. This interest is also reflected in the recent addition of new collaboration members. An additional 18 papers (8% of the total) have been published by the H1 Collaboration since 2012, which was already five years after the end of data taking.

All core software packages that were developed in the years 1988 to 2012 have been successfully migrated to a modern computing platform (`amd64` (x86-64), `CentOS7`, `GNU-compilers 9.2`). These software modules are required for data access, data processing, reconstruction, simulation, visualisation and, of course, high–level data analysis. External dependencies were updated to latest releases and are now obtained from the `LCG/AA` package repository. All other IT services are hosted by DESY-IT.

The common object–oriented data analysis framework `H1oo` is now based on `ROOT6` and supports the `C++17` standard. This framework facilitates the communication within the collaboration members, provides a common standard, and additionally provides highly valuable documentation. Some example programs and a few selected full high–level analysis codes

from previous publications are prepared for newcomers. The H1oo analysis framework is now fully accessible in Python and an online code documentation is available. H1oo now enables interactive analysis in C++ through CLING or more commonly in Python. Container solutions have been implemented for backward compatibility and software tests.

The programs and libraries are provided to the members of the H1 Collaboration on shared global file systems for convenience. All H1 software packages are now maintained in Git and new build instructions for a complete re-build of the entire software stack have been prepared. With these modifications to the H1 software architecture, H1 is confident to provide high quality data analysis capability of the unique HERA data in the future, using modern analysis tools, recent programming languages and on state-of-the-art platforms.

Acknowledgments

The authors would like to thank all collaboration members, past and present, as well as our colleagues in the DESY-IT department who helped make and continue to make the H1 experiment and DESY a uniquely enjoyable experience. Credit is also due to those involved in the early days of the DPHEP Study Group. We would particularly like to recognise and acknowledge the tremendous amount of work done by our late friend and colleague Vitaliy Dodonov, whose care, dedication and attention to detail made it possible to migrate the H1 software to 64-bit SLD6, laying the groundwork for the later migration to CentOS7 in 2020.

References

- [1] I. Abt *et al.*, The H1 detector at HERA, Nucl. Instrum. Meth. A **386** (1997) 310
- [2] I. Abt *et al.*, The tracking, calorimeter and muon detectors of the H1 experiment at HERA, Nucl. Instrum. Meth. A **386** (1997) 348
- [3] A. Accardi *et al.*, Electron Ion Collider: The next QCD frontier, Eur. Phys. J. A **52** (2016) 268
- [4] P. Agostini *et al.*, The Large Hadron-Electron Collider at the HL-LHC, arXiv:2007.14491
- [5] D. Anderle *et al.*, Electron-Ion Collider in China, arXiv:2102.09222
- [6] V. Blobel, BOS and related packages, Proc. 14th Workshop of the INFN ELOISATRON Project; Data Structures for Particle Physics Experiments: Evolution or Revolution, Erice, Italy (1990), p. 1–6
- [7] V. Blobel, The F-package for input/output, Proc. 6th Int. Conf. on Computing in High Energy Physics (CHEP 1992), Annecy, France (1992), p. 755–758
- [8] R. Brun *et al.*, GEANT 3, CERN DD/EE 84-1 (1984)
- [9] D. A. Duce, The graphical kernel system (GKS) ISO 7942, Computer Standards & Interfaces, Vol. 6, Issue 2 (1987) p. 235–237, DOI: 10.1016/0920-5489(87)90065-1
- [10] *CERN Program Library, CERNLIB*, URL: <https://cernlib.web.cern.ch/cernlib/> [accessed 2021-02-24]
- [11] U. Berthon *et al.*, New data analysis environment in H1, Proc. 11th Int. Conf. on Computing in High Energy and Nuclear Physics (CHEP 2000), Padua, Italy (2000), p. 700–703
- [12] M. Peez, The new object oriented analysis framework for H1, Proc. 13th Int. Conf. on Computing in High Energy and Nuclear Physics (CHEP 2003), La Jolla, USA (2003), arXiv:physics/0306124
- [13] J. Katzy, H1OO: An analysis framework for H1, Proc. 14th Int. Conf. on Computing in High Energy and Nuclear Physics (CHEP 2005), Interlaken, Switzerland (2004), DOI: 10.5170/CERN-2005-002.265
- [14] M. Steder, H1OO: A centralised analysis framework for the H1 experiment, Proc. 18th Int. Conf. on Computing in High Energy and Nuclear Physics (CHEP 2010), Taipei, Taiwan (2010), J. Phys.: Conf. Ser. **331** 032051 (2011), DOI: 10.1088/1742-6596/331/3/032051
- [15] P. Laycock, Ten years of object-oriented analysis on H1, Proc. 14th Int. Workshop on Advanced Computing and Analysis Techniques in Physics Research (ACAT 2011), Uxbridge, UK (2011), J. Phys.: Conf. Ser. **368** 012048 (2012), DOI: 10.1088/1742-6596/368/1/012048
- [16] *ROOT: Data Analysis Framework*, URL: <https://root.cern/> [accessed 2021-02-24]
- [17] *Data Preservation in High Energy Physics, DPHEP*, URL: <https://dphep.org/> [accessed 2021-02-24]
- [18] D. Asner *et al.* [DPHEP Study Group], Data preservation in high energy physics, arXiv:0912.0255
- [19] Z. Akopov *et al.* [DPHEP Study Group], Status report of the DPHEP Study Group: Towards a global effort for sustainable data preservation in high energy physics, arXiv:1205.4667

- [20] D. M. South, Data preservation in high energy physics, Proc. 18th Int. Conf. on Computing in High Energy and Nuclear Physics (CHEP 2010), Taipei, Taiwan (2010), J. Phys.: Conf. Ser. **331** 012005 (2011), DOI: 10.1088/1742-6596/331/1/012005
- [21] D. M. South, Data preservation and long term analysis in high energy physics, Proc. 19th Int. Conf. on Computing in High Energy and Nuclear Physics (CHEP 2012) New York, USA (2012), J. Phys.: Conf. Ser. **396** 062018 (2012), DOI: 10.1088/1742-6596/396/6/062018
- [22] D. M. South and M. Steder, The H1 data preservation project, Proc. 19th Int. Conf. on Computing in High Energy and Nuclear Physics (CHEP 2012) New York, USA (2012), J. Phys.: Conf. Ser. **396** 062019 (2012), DOI: 10.1088/1742-6596/396/6/062019
- [23] A. Campbell *et al.*, A dataflow meta-computing framework for event processing in the H1 experiment, Proc. 12th Int. Conf. on Computing in High Energy and Nuclear Physics (CHEP 2001), Beijing, China (2001), p. 651–655
- [24] S. Dagoret-Campagne *et al.*, Reprocessing H1 data on the IN2P3 computer farm in 1997, Proc. 9th Int. Conf. on Computing in High Energy and Nuclear Physics (CHEP 1997), Berlin, Germany (1997).
- [25] *dCache: Distributed Storage for Scientific Data*, URL: <https://dcache.org> [accessed 2021-02-24]
- [26] M. Cacciari and G. P. Salam, Dispelling the N^3 myth for the k_T jet-finder, Phys. Lett. B **641** 57 (2006), DOI: 10.1016/j.physletb.2006.08.037
- [27] M. Feindt and U. Kerzel, The NeuroBayes neural network package, Proc. 10th Int. Workshop on Advanced Computing and Analysis Techniques in Physics Research (ACAT 2005), Zeuthen, Germany (2005), Nucl. Instrum. Meth. A **559** 190 (2006), DOI: 10.1016/j.nima.2005.11.166
- [28] D. Ozerov and D. M. South, A validation framework for the long term preservation of high energy physics data, Proc. 20th Int. Conf. on Computing in High Energy and Nuclear Physics (CHEP 2013) Amsterdam, Netherlands (2013), J. Phys.: Conf. Ser. **513** 042043 (2014), DOI: 10.1088/1742-6596/513/4/042043
- [29] A. Pfeiffer, Overview of the LCG applications area software projects, Proc. IEEE Nuclear Science Symposium and Medical Imaging Conference (NSS/MIC 2004), Rome, Italy (2004), DOI: 10.1109/NSS-MIC.2004.1462660
- [30] S. Roiser, A. Gaspar, Y. Perrin and K. Kruszelecki, Servicing HEP experiments with a complete set of ready integrated and configured common software components, Proc. 17th Int. Conf. on Computing in High Energy and Nuclear Physics (CHEP 2009), Prague, Czech Republic (2009), J. Phys.: Conf. Ser. **219** 042022 (2010), DOI: 10.1088/1742-6596/219/4/042022
- [31] A. L. Hodgkins, V. Diez and B. Hegner, LCG/AA build infrastructure, J. Phys.: Conf. Ser. **396**, 052026 (2012), DOI: 10.1088/1742-6596/396/5/052026
- [32] P. Buncic *et al.*, CernVM: A virtual software appliance for LHC applications, Proc. 17th Int. Conf. on Computing in High Energy and Nuclear Physics (CHEP 2009), Prague, Czech Republic (2009), J. Phys.: Conf. Ser. **219** 042003 (2010), DOI: 10.1088/1742-6596/219/4/042003
- [33] M. Galli, E. Tejedor, and S. Wunsch, A new PyROOT: Modern, interoperable and more pythonic, Proc. 24th Int. Conf. on Computing in High Energy and Nuclear Physics (CHEP 2019), Adelaide, Australia (2019), EPJ Web Conf. **245** 06004 (2020), DOI: 10.1051/epjconf/202024506004
- [34] J. Blomer *et al.*, New directions in the CernVM file system, Proc. 22nd Int. Conf. on Computing in High Energy and Nuclear Physics (CHEP 2016), San Francisco, USA, (2016), J. Phys.: Conf. Ser. **898** (2017) 062031, DOI: 10.1088/1742-6596/898/6/062031
- [35] *DESY Bitbucket repository*, URL: <https://stash.desy.de/> [accessed 2021-02-24]
- [36] G. M. Kurtzer, V. Sochat and M. W. Bauer, Singularity: Scientific containers for mobility of compute, PLoS ONE **12(5)** e0177459 (2017), DOI: 10.1371/journal.pone.0177459