# HERA Data Preservation plans and activities

**J Szuba on behalf of the DESY Data Preservation Group**

Deutsches Elektronen Synchrotron, DESY, Notkestraße 85, 22607 Hamburg, Germany

E-mail: janusz.szuba@desy.de

**Abstract.** An international inter-experimental study group on data preservation and long-term analysis in HEP (DPHEP) was convened at the end of 2008 and held a series of workshops during 2009. The HERA experiments H1, ZEUS, HERMES as well as the IT division and the Library are well represented in DPHEP and efforts are now being made to form a coherent approach at DESY. Various options for preservation are explored, from permanent evolution (H1) to the use of virtualisation techniques (ZEUS). Both experiments have planned the computing and the associated resources until 2013 and now explore possibilities to ensure the maintenance of the data analysis capabilities beyond 2013. A common effort and additional resources may lead to longer viability of data analysis. Technical solutions have been investigated by DESY-IT and involve virtualisation systems tailored for long term software preservation as well as systems for self consistent data archiving and migration. The communication between experiments, DESY-IT and the Library have put forward the possibility for further common developments related to documentation scanning and storage as well as pilot projects within the HEP documentation system INSPIRE. The evaluation of such projects is ongoing and concrete proposals to ensure HERA data analysis after 2013 are expected in 2010.

## 1. Introduction

The issue of data preservation of HEP experiments is addressed in a systematic way within a study group on Data Preservation and Long Term Analysis in High Energy Physics (DPHEP). The group was formed at the end of 2008 and convened a series of workshops since then. The HERA experiments H1, ZEUS and HERMES as well as the IT division and the DESY Library (DESY Data Preservation Group) are well represented within DPHEP study group and now are making efforts to form a coherent approach on data preservation aspects.

In the following the plans and status of the individual experiments are presented, followed by the collaborative projects as well as documentation efforts.

## 2. Future Analysis Models

The generic scheme of HEP data analysis from raw data up to analysis ntuple-like objects is common for all experiments. The differences lie in a specific data format, abstraction level, software details etc. The present analysis models are not always well suited for the long term data preservation. Thus, all HERA experiments developed a strategy for their own data preservation.

### 2.1. ZEUS Analysis Model

The future analysis model adopted by ZEUS is based on the Common Ntuple project, started back in 2006. All data and wide range of MC samples is preserved in a flat ROOT ntuple format. Besides, maintaining the ability of simulation of new MC samples after the end of the

current analysis model is foreseen. The standalone MC simulation package, which includes a full chain from simulation to reconstruction to Common Ntuple production, will replace the current production system. With this package MC mass production can then be performed on real machines, given that necessary resources like GRID or local farm are available. Once the migration to a newer OS breaks the MC package can be frozen, a virtual image created and run on a virtual machine.

The necessary ingredient for any future analysis model is a validation tool which tests analysis results against changes in software or running environment. Validation can be performed on different levels of data: simulation, reconstruction and physics analysis. The simple tool used by ZEUS so far to compare different MC simulation software releases needs to be developed further to incorporate also the physics analysis level.

### 2.2. HERMES Analysis Model
In case of HERMES experiment, the full software analysis chain, based on MDST, is planned to be preserved. Also the ability of new MC mass production requires a fully functional data production chain. The complete data reproduction is however not foreseen, although would be in principle possible. Running the experimental software on virtual machines is not excluded and first tests were performed.

The validation procedures are in development phase.

### 2.3. H1 Analysis Model
H1 plans a rolling model of data preservation, with a production timescale of a few months interval. The rolling preservation model includes a regular recompilation of analysis level software and full production of analysis level data and MC. Several issues concerning external software dependencies still need to be resolved.

For this kind of preservation model good validation tools are essential. Such a scheme already exists to validate the data files content of the analysis level software between different releases. Expanding this scheme to include full analysis selections, as well as the simulation and reconstruction code is foreseen.
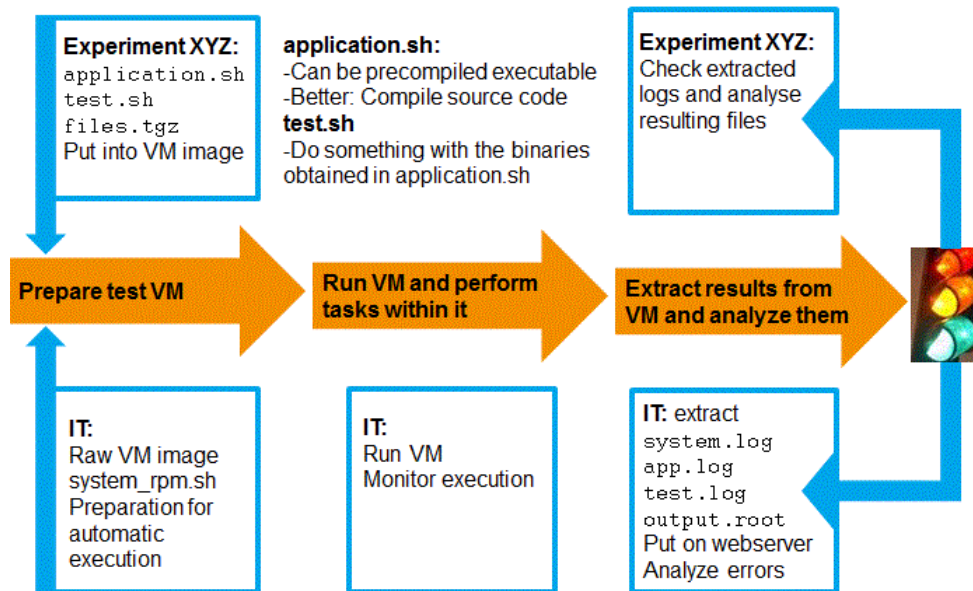
### 2.4. Common issues between experiments
Despite differences in computing models, all HERA experiments have many things in common which were identified in the process of planning future analyses. Among them are migration to newer operating systems, recompilation of experimental software with newer compilers or external software dependencies. All HERA experiments are currently moving to Scientific Linux 5. During this migration, several problems had to be overcome. The change in gcc compiler suite from g77 FORTRAN compiler to gfortran required the usage of non-default, but newer, gcc 4.4 compiler by H1 and ZEUS. Nevertheless, HERMES still uses gcc-3/g77 for compilation of their software on SL5, although efforts for full migration are still under consideration. Some of the problems remained, but they are relevant to obsolete external software like ADAMO or GKS. Also several external software dependencies are common for all experiments, like CERNLIB, ROOT, CLHEP, ORACLE, FastJet, Neurobayes, to name a few and can be dealt with on the same ground.

## 3. Analysis Software Validation Project
Giving the similarities mentioned above between all HERA experiments, like migration to a new operating system, recompilation of the experiment specific software, or furthermore validation of new software releases, testing access to data, running simulation, reconstruction and analysis programs in different environment and verifying the output, it is straightforward to think of a need for a unified validation suite.

One of the main proposed data preservation projects at DESY is the development of such a validation framework. A framework, which allows a rigorous test of experiment level software builds against changes in operating system, external software, running environment, is realized using virtualisation techniques and will prove invaluable in dealing with future migrations. The schematic workflow of the proposed framework is depicted in figure 1. The core of this



**Figure 1.** A workflow of the proposed experimental software validation scheme at DESY-IT.

framework is a virtual machine running clean installation of an operating system with additional dependency requirements defined by an experimental software. The experiments application is then downloaded, run and tested inside VM and finally the results are exported outside for checks of output and error analysis. There is a clear separation of tasks both for an experiment and IT side. The experiments provide application packages and validation tests, while IT maintains virtual images of different operating systems, runs virtual machines and provides resulting output.

A test version of such a framework was successfully installed, where the stability of a variety of software from the H1, ZEUS and HERA-B collaborations was tested against three different operating systems, showing the proof of principle of such a scheme. Figure 2 presents in a simplified form the results from the test run.

Data analysis on virtual machines has been tested by the HERMES collaboration, who could also participate in this project.

A full version of the validation suite may now be implemented at DESY-IT, to safeguard the HERA data for the long term.

## 4. HERA data formats for preservation

Data format and approximate sizes intended for preservation are briefly discussed here.

The final ZEUS data reprocessing to MDST format was completed in 2009 and no more reproduction are foreseen. The basic preserved data format are ROOT based Common Ntuples, which are produced iteratively improving and adding new content required by physics analyses. Ultimately RAW, MDST data and MC tapes will be removed from robots and stored in safe place. Thus, the total amount of space needed for preserving ZEUS data is reduced from the current 1 PB to approximately 100 TB.

**Figure 2.** Result of the first run of the test version of the validation framework.

Also the final H1 reprocessing of HERA II data was finished in 2009, currently the equivalent HERA I reprocessing is ongoing. Common analysis software H1OO, which started back in 2000, uses a ROOT based data format and is a basis of all H1 analyses. In addition, a current monthly MC production amounts up to 1/4 billion events, however it will decrease with time. H1 is planning to preserve RAW data, as well as at least one DST and analysis level versions for data and MC. The estimated total amount of data to be preserved for H1 is around 200-500 TB.

Main format for HERMES analyses is the MDST. The new production is planned before final freeze, using an improved calibrations for the last years of data taking. The total amount of data to be preserved on tapes is estimated as 150-200 TB.

Preservation of HERA-B data is under investigation within DESY-IT [1]. Total amount of data is currently around 250 TB, but it will decrease once a preservation model is established.

## 5. An Archival System for HERA Data
The scope of the validation framework described in the previous section does not foresee an examination of the condition of complete data sets, but rather the use of smaller samples to test software changes. The present dCache storage system at DESY-IT is not suitable for long term storage of HERA data. Additional complications arise due to the widespread use of HEP specific protocols. It is therefore proposed to develop a long term archive system, which would include: automatic migration to new media generation and technology, automatic data integrity checks, retrieval of the data themselves and metadata operations.

Such a system would serve not only for the HERA experiments data preservation purposes, but for all scientific groups at DESY, including rapidly growing Photon Science community.

## 6. Other Virtualisation Related Projects
Among main projects related to the long term data preservation, also a few smaller ideas emerged. One of them is a study of virtualisation techniques using H1OO within the CERN VM [2]. An example of running an analysis using a virtual image of the H1 environment was created although still access to the data remains an outstanding problem.

The second example is a test of possible usage of the Cloud Computing model, based on the Eucalyptus Private Cloud platform and the open source storage system CEPH [3].

## 7. Documentation projects

An effort has now begun to secure the status of the H1 non-digital documentation. This include cataloguing, organisation and digitalization where appropriate of H1 papers, notes, drawings, talks, etc. In terms of digital documentation, several initiatives are needed including streamlining the current H1 web content and working with the DESY Library and INSPIRE about future storage of electronic documentation and secondary data. On the more technical side, the database access and reliance of the H1 web server should be examined and a migration of the web server to the DESY-IT virtual environment should be foreseen, relieving the collaboration of future hardware requirements.

ZEUS non-digital documentation includes internal notes, transparencies from collaboration meetings before 2000, technical drawings etc. Consolidation, electronic cataloguing as well as digitalization of old theses and internal notes is planned. Finally, custody will be handed over to DESY library. ZEUS digital documentation resides mostly on the main web server. This include also specific technical documentation (detectors, trigger) and electronic log book. Migration of the main web server to newer hardware as well as content revision is foreseen.

HERMES made a large effort to move all important documentation, technical notes to a wikipedia structure. Revisions of the content will be performed by corresponding experts. For the long preservation only wiki pages are considered. Electronic logbook as a important source of information about data running conditions may be preserved in a virtual environment.

## 8. Conclusions

The current status of the ongoing activities and proposed projects were presented at the last DPHEP workshop at KEK, where other experiments showed interest in the further development of the validation project, including BaBar collaboration, who are also investigating such a validation system.

For all presented projects additional resources in terms of person power and financial support are needed. The Physics Research Commitee of DESY Laboratory fully endorsed the DESY data preservation group report status and requirements for the proposed common projects and the next steps were undertaken to fully ensure HERA data will be safeguarded and future analyses possible.

## References

[1] D.Ozerov 2010 Heritage Preservation of HERA-B Collaboration *these proceedings*
[2] CernVM http://cernvm.cern.ch/portal
[3] B.Lobodzinski 2010 Tests of cloud computing and storage system features for H1 *these proceedings*