

DPHEP@Max-Planck Institut für Physik

Andrii Verbytskyi

DESY DPHEP Group Meeting Hamburg, November 20, 2014 A unique data collected from past experiments can be re-analysed in the future with new methods and used to verify new theoretical predictions. MPP participates in the data preservation for the following experiments:

• H1

- ZEUS
- JADE
- OPAL

The main intention is to provide facilities for the physics analysis and to do physics analysis for in house experiments.

- Save all data...
- Save all software...
- Make (re-)analysis possible!

• We are not interested in data, software, codes, data bases, past computing facilities and outdated documentation;

• We are interested in physics analysis and physical results!

• Data, software, codes and data bases should be saved only for the reason of physics analysis.

Ideas about DPHEP@MPP and outside

- The size of the data is not a bottleneck;
- Save everything which could be useful with the estimated amount of preserved knowledge and manpower in the future.

- Save all the codes;
- Try to keep working things crucial for the future analysis;
- Any modification of the software should be done by experts or not done at all.

- Save electronic documentation in the most simple format;
- Make documentation (including as much technical notes as needed) available in InSpire or other data bases;
- Documentation should be about making an analysis in the future and not about making the analysis in the past;
- Provide documentation on Data Preservation. Examples are mandatory.

- BIG PROBLEM experiments had not made any decision on access policy. Maybe some options should be explained?;
- A future strategy should be developed as soon as possible;
- Note that approaches for DESY and RZG are very different.

- There should be enough CPU power to do analysis in a well defined, stable environment;
- The location of the resources should not be restricted to DESY or RZG;
- Some estimations of the required CPU power should be done.

- Validation of produced data (e.g. MC samples) in Data Preservation mode is needed;
- At least a toy analysis in Data Preservation mode is needed to test the DPHEP infrastructure.

Implementation of DPHEP at MPP (and outside)

- All OPAL data marked for preservation is copied to RZG. Available also at CERN;
- All JADE data marked for preservation is copied to RZG. Maybe also should be copied to DESY? That is DESY data!!!
- All ZEUS data marked for preservation is copied to RZG¹. Available also at DESY;
- 70% of H1 data probably marked for preservation is copied to RZG. No final list of data files yet. Available also at DESY.

 $^{^1 \}rm With$ an exception of PAW ntuples which were marked for preservation only recently and some new MC samples. To be done soon.

Source codes as is:

- OPAL software available on CERN AFS. To be maintained by CERN IT;
- JADE software available on RZG AFS. To be maintained by DPHEP@MPP;
- ZEUS software available on DESY AFS. To be maintained by DESY IT² Software also backed-up on a web-server;
- H1 software available on DESY AFS. To be maintained by DESY IT.
- The environment for the (analysis) software will be created inside virtual machines. See details below.

²Are we sure about it? Will be AFS there in 20 years?

- So far there is no intention to copy/duplicate the preserved documentation from any experiment to MPP or RZG storage or server. All of the experiments did a huge work to save their legacy;
- The future updates of the documentation for ZEUS and H1 (e.g. bookkeeping of new MC samples) should be clarified;
- Documentation on Data Preservation at MPP is in preparation.

• Current status:

- JADE data is provided as is a)on RZG AFS b) in ZEUS directory on RZG dCache (for ZEUS VO members only);
- OPAL data is available on RZG dCache for OPAL VO members;
- ZEUS data is available on RZG dCache for ZEUS VO members;
- (A lot of) H1 data is available on RZG dCache for H1 VO members;
- The data from RZG is available from any PC inside and outside MPP and RZG;
- So far a proper Grid certificate is needed. Unfortunatelly Grid certificates are not popular... Looking for better solutions – e.g. short living X509 credentials. However, that kind of service can be implemented only in big IT centre – DESY, RZG or CERN.

- Nowadays
 - For the read-the-data-only (e.g. ZEUS Common ntuple based analyses) it is possible to use virtually any existing system with ROOT or PAW;
 - There are about enough CPU resources to parasite on (local work-group servers in different institutes, Grid, NAF etc.);
 - Not clear how much of that resources are available for more complicated tasks like MC generation;
 - Not clear how long those resources will be available at all.
- For the future a set-up for virtual machines are prepared/in preparation. The final versions should be deployable on local machines, academical (CERN Openstack) and commercial cloud services (e.g. Amazon). This will assure virtually unlimited CPU capabilities for many years.

- In DESY a lot of work was done towards validation of ZEUS and H1 MC production in Data Preservation mode. Maybe some of those tests will be reused at MPP;
- Some toy analyses (use cases modelling possible future analyses) were done at MPP using RZG storage facilities (i.e. first analyses outside of DESY). The goal was to estimate how easy it will be to an analysis in the future. See description below.

Implementation: tests and validation, use cases for possible future analysis

- Possible future discovery of exotic particles (e.g. penthaquarks) will immediately rise a question about re-analysis of HERA data. Let's consider some possible cases:
 - With the preserved ZEUS data the reconstruction speed of penthaquark decay $P_{\bar{c}s}^- \rightarrow \pi^- K_S^0 \Lambda_0$ is 620 events/s/core for Xeon E5520@2.27GHz; Full processing time for HERAII data is 160 hours.
 - With the preserved ZEUS data the reconstruction speed of some analyses (e.g. leptoquark decays, Deeply Virtual Compton Scattering) can be as low as 10-100 events/s/core. Full processing time for HERAII data is more than 1000 hours;
- Generating MC with new models (see HERA symposium 2014 talks on MC) is even more challenging and CPU consuming procedure.

- RedHat-based systems SL5-i386, SL5-x86_64, SL6-i386, SL6-x86_64, SL7-x86_64 or corresponding CentOS with small installation;
- Kickstart-based installation from custom ISO image with all needed packages;
- H1/ZEUS/OPAL/JADE specific software,databases etc. included to the ISO;
- Timestamped ISOs should be available for data analysers in RZG dCache.

- SL5-x86_64;
- Kickstart-based installation from custom ISO image with gcc,gfortran,ROOT,PAW, cernlib, globus etc.;
- ZEUS event display (ZEVIS), file catalog (cninfo), MC files format converter (formoza), standalone MC production package (ZMSP)³;
- Set-up scripts for different data locations (e.g. DESY or RZG);
- Runs on VirtualBox (should be possible to run on other hypervisors) and uses data from DESY and RZG. Can be used outside of DESY/MPP/RZG;
- To be prepared for using on CERN Openstack.

³So far has problems, hopefully will be fixed in next days.

- OPAL: it looks the software is in a good shape, an expert is available. Perhaps will be done right after ZEUS;
- JADE: requires porting of the software from AIX, an expert is available;
- According to H1 plans, H1 software should be ready for recompilation on newer systems. Maybe start with SL7?

Conclusions and plans

- The Data preservation in MPP is moving forward;
- There is a clear goal and plan to implement it;
- Still, there are some decisions that should be made by experiments as soon as possible. It is impossible to do it from outside;
- Any help and commitment is welcome.