

Use case: CERN Analysis Preservation

Sünje Dallmeier-Tiessen¹, Artemis Lavasa¹, Tibor Šimko¹, Javier Delgado Fernández^{1,2},
 Pamfilos Fokianos¹, Robin Dasler¹, Anxhela Dani^{1,3}, Annemarie Mattmann^{1,4}, Ioannis
 Tsanaktsidis¹, Anna Trzcinska¹, Diego Rodriguez Rodriguez^{1,2}

Abstract— The CERN Analysis Preservation Framework is a central platform for the four LHC collaborations at CERN and it was developed to address the need for the long-term preservation of all the digital assets and associated knowledge in the data analysis process, in order to enable future reproducibility of research results. As the service continues to develop, new challenges arise, so in this paper we present the initial considerations for implementing a common data markup schema for all the information in the service, which is essential for enhanced search and discovery, and for supporting links with various other internal and external platforms.

LHC data is unique, precious and complex. Consequently, preserving this data, the software and the knowledge around a physics analysis in a comprehensive reusable manner is challenging. The CERN Analysis Preservation (CAP) service was started in partnership with the LHC collaborations (i.e. ALICE, ATLAS, CMS and LHCb) to try and overcome those barriers. It aims to enable individual physicists to preserve information as it is produced by making it easier to describe, preserve, find, exchange and hand over information in a fast paced environment of highly fluctuating personnel.

The CERN Analysis Preservation service (Figure 1) is a response to two parallel demands, internal and external. Within the Collaborations, the high and complex output of analyses and analysers results in significant challenges in terms of capturing and preserving the analysis and the knowledge around it, which is of particular relevance for future reuse and reproducibility of the research results. Externally, an increasing number of funding agencies have put in place data management policies, which demand the development of comprehensive data management frameworks for data and knowledge preservation, and for future reuse or even reproducibility of research outcomes. A similar development can be observed

¹ CERN, Geneva, Switzerland

² Universidad de Oviedo, Oviedo, Spain

³ Alexander Technological Educational Institute of Thessaloniki, Thessaloniki, Greece

⁴ Technische Universität Darmstadt, Darmstadt, Germany

among journals, which start making data and reproducibility measures mandatory for any paper to be published.

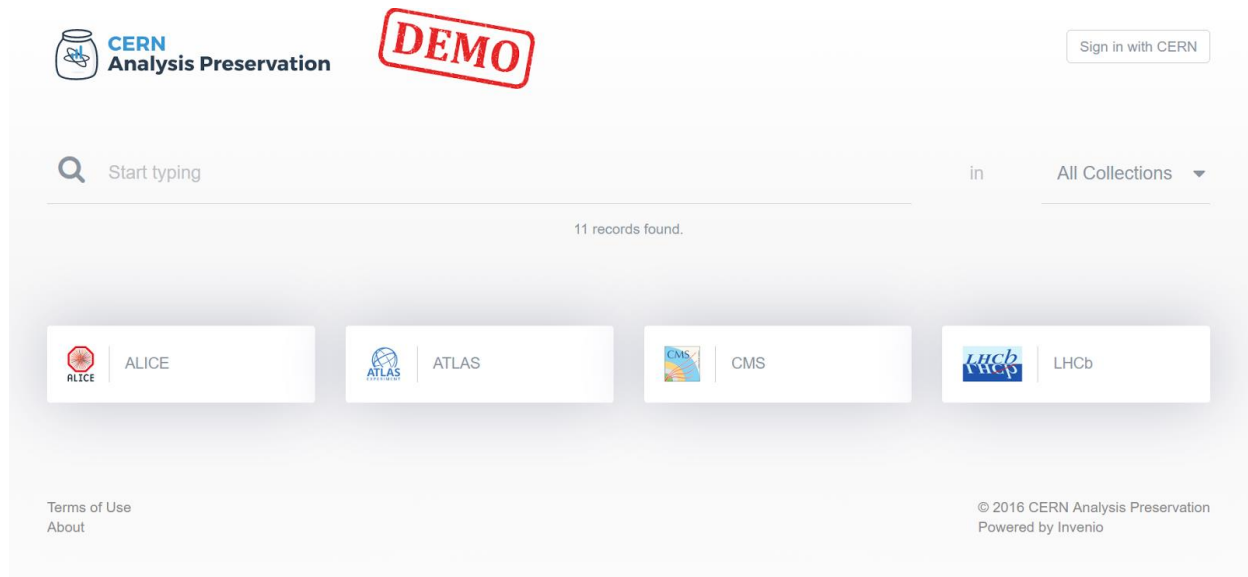


Figure 1: The CAP service

CAP is developed as an open source tool and built on the Invenio Open Source Digital library framework⁵, which also underpins services such as Zenodo⁶ or CERN Open Data⁷ and the data model is developed in JSON⁸, which is able to handle highly intricate metadata structures. CAP service functions are constantly being enriched and improved through CERN’s participation in the THOR project (Technical and Human Infrastructure for Open Research)⁹, a European Commission-funded project under the H2020 framework program; for example, through this project work is being done to implement ORCID authentication and to explore dynamic data citation.

As the tool gathers information throughout the research lifecycle, it is faced with complex materials, such as data, software, their dependencies and versions. To facilitate future reuse of such research objects, it is indispensable to enable the researcher to provide enough context information around the analysis (progress). A standard analysis “record” within the service contains detailed information about the processing steps, the datasets that are used, the software (version) used and, in addition, detailed physics information is made available (e.g. energy level, vetoes, cuts)¹⁰. It is expected that researchers/users of the

⁵ <http://invenio-software.org/>

⁶ <https://zenodo.org/>

⁷ <http://opendata.cern.ch/>

⁸ <http://www.json.org/>

⁹ <https://project-thor.eu/>

¹⁰ The schemas used: <https://github.com/cernanalysispreservation/analysis-preservation.cern.ch/tree/master/cap/jsonschemas>

tool will also use this information to search for interesting ongoing analysis (content).

Of course, to be able to benefit from the rich knowledge that could be made available in the tool, it would be best to standardise the information inside as much as possible. This would enable an unprecedented search functionality across these preserved content. Figure 2 below shows an example excerpt used by the CMS collaboration. Currently, the information provided in analysis records is tailored to each experiment, in order to address their particular needs, which makes standardisation even more challenging.

Did you produce n-tuples as an initial step to any measurement? ☒ no ☐ yes

Main Measurement Workflow

Please provide information about the main measurements of your analysis

Item #1

Measurement Description E.g. signal measurement of NNN in Z -> ee final state

Description Details If applicable, please provide a more detailed description for this measurement

Event Selection

Physics Objects

Item #1

Object Select Object

Relations + Add New Item

Vetos + Add New Item

Figure 2: Example excerpt from the schema developed for the CMS collaboration

First attempts have been made to look into existing vocabularies and in particular into using the Schema.org¹¹ vocabulary encoded in JSON-LD¹², which would provide much-needed support for Linked Data. From this initial research, it became apparent that a vocabulary that could provide wide enough coverage for our purposes without having to be extended significantly to reflect the community specific subject has not been developed yet. This is, of course, often the case when attempting uniformity within any given system with unique characteristics. More specifically, there are terms that would need to be standardised that may only find application in one experiment of the LHC. So, the most logical way to handle this situation would be to integrate a common data vocabulary for the more high-level information

¹¹ <https://schema.org/>

¹² <http://json-ld.org/>

across the platform and then proceed with extensions where needed, most likely on an experiment-by-experiment basis.

As the service continues to grow, more and more records will be introduced into the system, thus making the implementation of a common way of describing all this information more relevant than ever, not only for CERN Analysis Preservation, but for CERN's other platforms as well (e.g. the CERN Open Data Portal). This is particularly crucial to facilitate a meaningful search. As an example of a searchable field that would be facilitated by a common vocabulary and format is the particle reactions and observables for each data table (Figure 3). This has already been implemented in HEPData¹³, the High-Energy Physics open-access data repository, which is hosted by CERN.

Reactions	
P P --> P P	116
E+ E- --> HADRONS	108
PI- P --> PI- P	74
PI+ P --> PI+ P	64
PBAR P --> PBAR P	50
GAMMA P --> PI+ N	42
GAMMA P --> PI0 P	40
P NUCLEUS --> P X	39
P P --> X	37
PBAR P --> X	31
E+ E- --> Z0	29
E- P --> E- X	27
PI- P --> PI0 N	26
E+ E- --> MU+ MU-	24

Figure 3: HEPData reactions facet

Some work on this subject to address the need for formally representing and describing experimental results within the particle physics community has already been carried out (Carral et al., 2015). As we work to develop high-quality services for our users, it is essential that we build upon and expand these ideas, and the improvement and semantic harmonisation of our data models is an integral part of that.

¹³ <https://hepdata.net/>

REFERENCES

Carral, D., Cheatham, M., Dallmeier-Tiessen, S., Herterich, P., Hildreth, M. D., Hitzler, P., ... & Watts, G. (2015, October). An Ontology Design Pattern for Particle Physics Analysis. In WOP.