Data & analysis preservation at CERN

intro and IT aspects

Software & Computing Round Table 2 June 2020 Dirk Duellmann, CERN IT

Preserving the Deduction Chain & Preserving the Investment

- Preservation what and why?
 - data (measured or simulated) provides the opportunity to \bullet obtain scientific information via analysis
- Analysis steps: need to adhere to scientific quality standards
 - reproducible, repeatable
 - by other people, on different hardware, at a later time
- Scrutiny requires to preserve *used* data and analysis
 - benefits the originating project and science community as whole
- General data preservation
 - data is (re-)produced with a significant **public** investment
 - data/method reuse for a different purpose in the future \bullet
- benefit original producers but only in future
- benefit an external community Open Data







Community Context - DPHEP

- ICFA Study Group on Data Preservation and Long Term Analysis in High Energy Physics
- Includes lacksquare
 - many accelerator sites
 - experiments at different stages of their lifecycle
- Built foundation and structured discussion across the HEP community



Data Abstraction Hierarchy

 \sim KB/paper

 \sim GB/analysis

1. Provision of additional documentation for the published results

2. Simplified data formats for analysis in outreach and training exercises

 \sim TB/analysis

 \sim PB/year

 \sim GB/sec

3. Reconstructed data and simulations as well as the analysis level software to allow a full scientific analysis

4. Basic raw level data (if not yet covered as level 3 data) and their associated software which allows access to the full potential of the experimental data







Four data levels for capture, preservation and opening adapted from slide by Tibor Simko

Some milestones... (tombstones?)

- 2012: DPHEP Blueprint published
- **2013: CERN assumes role of DPHEP MGR**
 - Many events & w/s in <u>https://indico.cern.ch/category/4458/</u>
- FAs
- **2016: CERN services for LTDP@iPRES (Bern)** •
 - documentation
- 2017: DPHEP w/s on FAIR data, TDRs

 - audit as well
- protocols")
- 2019: Workshop on Sustainable Software: LEP+CERNLIB experience

2014: DPHEP Collaboration Agreement signed by CERN and other labs /

Covers all areas proposed to 2012/13 ESPPU: bit preservation, s/w preservation,

- TDR = Trustworthy Digital Repository (ISO 16363, CoreTrustSeal, DPC RAM, ...) – Self audit of CERN (ISO 16363) completed), discussion with T1s in ESCAPE to complete

• 2018: Science Europe Guidelines on Research Data Management ("domain



Focus on LHC experiments

- Driving current infrastructure and software developments
- Active collaborations are most exposed to the central resource conflict that needs to be addressed
- Assumption: a solution for large active collaborations
 - can be extended/adapted to other projects
 - prevents further unmanaged knowledge decay that may already have taken place in older projects

Data & Analysis Preservation Workshop

- LHC experiment focus (profiting from RHIC participation)
 - Experts and management from experiments, WLCG and IT
 - All experiment presentations and discussion notes at
 - https://indico.cern.ch/event/858039/
- Target: collect consolidated experiment input to policy and funding discussion

Four domains of best practices



- Analysis preservation can be adopted step-by-step
- Synergies between the domains

LHCb data and analysis preservation

Slide: Concezio Bozzi



REusable ANAlysis

- Plan to use REANA
 - Data: open data \bigcirc
 - Use input from CAP \bigcirc
 - Software: CVMFS \bigcirc
 - **Environment: LEGO trains** Ο
- Can be used for:
 - Rerun the train \bigcirc
 - Plot production with local macros \bigcirc
- **Procedure** is available on the REANA main page:
 - to run a lego train in a AliPhysics specific \bigcirc container
 - use as input the intermediate files of the lego \bigcirc system, not yet any pre-saved JSON configuration file
- In addition, simple <u>ALICE analysis demonstrator</u> submitted to the CERN Open Data portal works with ALICE Open Data VM, now available in REANA with docker container



Slide: Stefano Piano - Data & Analysis Preservation workshop









Technology Support

- Bit-preservation is performed routinely and independently by storage infrastructures
 - media replacement and disk failure tracking/prediction
 - eg @CERN: EOS with Castor/CTA as archive layer
- S/W and environment preservation
 - containers and GIT are used by all experiments

- Technology problems are largely solved in a common way
 - current industry standard -> migration path can be expected



Reproducible Analysis

- Reproducible analysis provides an immediate scientific benefit
 - as a principle: essential for quality science
 - as a tool (eg ReANA) helps
 - analyst productivity via automated re-execution
 - consolidated input to **publishing review** process
 - preservation by making analyst knowledge explicit early (input data, s/w, environment, meta-data)

REANA Reproducible Analyses

Status: pilot

Purpose: run containerised scientific workflows on diverse compute clouds (Kubernetes, Condor, Slurm)

Usage: data + code + environment + workflow = reproducibility

Community: pilot examples with ALICE, ATLAS, CMS, FCC, LHCb; synergies with astronomy, life sciences, machine learning

Notes: focus on "preproducibility" of analysis during its active phase; structure analysis in a reproducible way to facilitate its future "preservation"

Resources: shared Ceph storage and Kubernetes cluster; need for proper experiment accounting





http://www.reana.io

Slide: Tibor Simko



Computational workflows





Serial





CWL

Slide: Tibor Simko

REANA: NSF collaboration

Ongoing collaboration with **NSF** SCAILFIN project (NYU, Notre Dame University, etc) to deploy **REANA** using VC3 on **HPC** centers in the US





Slide: Jose Benito Gonzalez



CERN Analysis Preservation (CAP)

Responds to two parallel demands regarding data re-use and reproducibility:

- **Internal:** high complexity of experiments create major challenges in capturing and preserving analyses and related knowledge
- **External:** Funder policies that require comprehensive solutions for comprehensive data management and knowledge preservation

Accordingly the goals of this effort are:

- Capture all the elements needed to understand and rerun an analysis and link them together persistently
- Make analysis components easily discoverable, shareable and re-useable
- Flexible to respond to diverse needs of research teams







What is CAP? Get Started Integrations Documentation Log in

CERN Analysis Preservation

capture, preserve and reuse physics analyses



Capture



Collect and preserve elements needed to understand and rerun your analysis

Collaborate



Share vour analysis and components with other users your collaboration or group Reuse



Run containerized workflows and easily reuse analysis components

Slide: Kamran Naim

RCS-SIS | Open Science



CAP Features:

- Flexible data models (JSON-Schemas)
- FAIR practices [more info]
- Versioning of metadata and files
- Integration with related scientific services and universal identifiers (i.e. Github, Gitlab, Zenodo, ORCID, ROR, etc.)
- Ongoing integration with services that support remote execution and reuse (e.g. REANA) of computational workflows
- Piloted in collaboration with all the major LHC experiments at CERN





Slide: Kamran Naim

RCS-SIS | Open Science

Open data policy doi:10.7483/OPENDATA.LHCb.HKJW.TWSZ

- Adopted by the LHCb Collaboration Board on Feb 27th 2013
- Open data for Outreach and education:
 - Only a limited fraction of the complete LHCb data-set may be used. »
- Open data for Research:
 - may be varied for specific requests.»
- Big caveat

• »... LHCb already participates in outreach activities and will continue to do so. [...] The data are provided for educational purposes only, and are not considered suitable for publication.

• «... access will be granted to portions of the DST data five years after data is taken. The portion of the data which LHCb would normally make available is 50% after 5 years, rising to 100% after 10 years. All requests will be considered by the CB and the period and proportion

• «... LHCb is extremely resource limited at present. Therefore whilst this policy expresses a spirit of intent, we cannot commit to implementation of any capability on any specific timescale. Specifically in respect of open access we will not be able to undertake any significant development to support this without injection of additional resources. »

> LHCb data and analysis preservation 17

Slide: Concezio Bozzi



CERN Open Data http://opendata.cern.ch

opendata CERN			Ab
OPERA ×		Sort by: Best match 🗘 asc. 🗘 Display: detailed 🗘 20 results 🗘	
include on-demand datasets		Found 910 results.	
Filter by type			
✓ □ Dataset	904		
Derived	904	OPERA neutrino-induced charmed hadron event 10270021561	
 Documentation 	4	This OPERA detector event is a muon neutrino interaction with the lead target where a charmed	
About	1	hadron was reconstructed in the final state. The event data consist of Electronic Detector files	
Authors	2	(such as	
Guide	1	Dataset Derived OPERA	
□ News	2		
Filter by experiment		OPERA neutrino-induced charmed hadron event 222274169	
	26	This OPERA detector event is a muon neutrino interaction with the lead target where a charmed hadron was reconstructed in the final state. The event data consist of Electronic Detector files (such as	
□ ATLAS	125		
	3928		
LHCb	12		
OPERA	910	Dataset Derived OPERA	









CERN Open Data

Status: *production* (since November 2014)

Size: 7K records, 800K files, 2 PB size

Purpose: "big data" sharing of event-level particle physics data and accompanying code for both education and research purposes

Content: raw samples, collision & simulated & derived datasets, docs, configs, software tools, example analyses, VMs, event display

Community: ALICE, ATLAS, CMS, LHCb, OPERA (coming: JADE, Data Science)

Notes: independent expert curation; batch ingestion workflows with Collaborations



Slide: Tibor Simko

http://opendata.cern.ch



Integration Options with Existing Infrastructures

- Open Data access patterns seem well matched to the existing archive infrastructure at CERN or the distributed archive in WLCG
 - potentially: include also cloud based back-end storage eg funded via CPU purchases of external OD users
 - balance between cost distribution and long term availability
- Expectation: Open Data is findable and available, but not in all cases without access latency
 - caching techniques together with lower cost archive media should allow to steer media costs
 - closer integration with experiment data management should enable to further reduce duplication between open data and experiment resources

Common Technical Activities

- All experiments focus on similar model and same tools
- Production: CERN Open Data Portal (COP)
- Evaluation: ReANA, CERN Analysis Portal (CAP)
 - insure resource coverage to complete evaluation by significant part of the analysis user community
 - setup effective communication channels to facilitate technical exchange between experiments and service teams
- Use HSF/WLCG workshops to discuss and facilitate the adoption of common data preservation and reproducible analysis services and their integration with existing production

Data Management Plans

- All experiments anticipate the requirement of formal Data Management Plans (DMPs)

 - potentially without additional resources, unless a
- Proposal:
 - concrete experiment DMPs

integrated with existing services at hostlab and WLCG

demonstration of value for additional investment can be given

• Task experiment computing management, IT and WLCG to prepare an agreed DMP skeleton - to be used to derive

Towards a Common Policy

- A common Open Data Access Statement between CERN and LHC experiments could provide leadership
- Frame for common activities such as defining a resource reporting & review process for open / preservation data
 - profit from policy review several experiments are currently undertaking
- Towards FAs: clarification of potential hostlab responsibilities wrt open data and preservation
- Dedicated discussion via experiment spokespersons and directorate has been initiated
- WG has been actively working towards a CERN Open Data Policy draft

Open Science, Open Data & Data Preservation

- Open Science, Open Data and Data Preservation are closely related
 - obligation to return the full potential of a **public** investment to the public
 - on the scale of LHC experiments this is a complex goal requiring patience to develop a shared strategy resulting in quality science
- Open Data and Data Preservation require a policy/funding incentive
 - to resolve their likely resource conflict with more immediate priorities of scientific collaborations
- Goal: find a pragmatic balance for the collaborations, science community and public
 - between decaying analysis knowledge and decreasing storage prices
 - within external budget and policy constraints

Summary

- activities and data & analysis preservation.
- CERN management is currently being prepared
- impact of data privacy regulations locally and via DPHEP

• All LHC experiments made significant investments in open data

• Supporting services (COD, CAP, ReANA) exist and are evaluated by all LHC experiments. Their adoption by the analysis community at large should be discussed regularly eg at WLCG/HSF workshops.

 At this point Run1 Open Data releases have absorbed the initial adhoc resources. Without clarification on policy/funding incentives further releases would likely be suppressed by immediate priorities.

A common Open Data Statement between LHC experiment and

Expect to continue discussion on Data Management Plans and

Thank you.