



Towards Reproducible Research Data Analyses in LHC Particle Physics

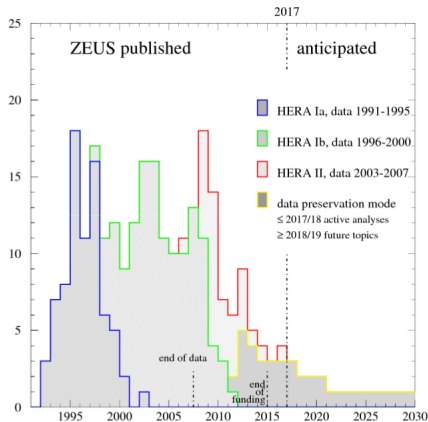
Tibor Šimko
CERN

ILIDE 2017 · Jasná, Slovakia · 3–5 April 2017

Preserving research data

Why? Scientific output timeline

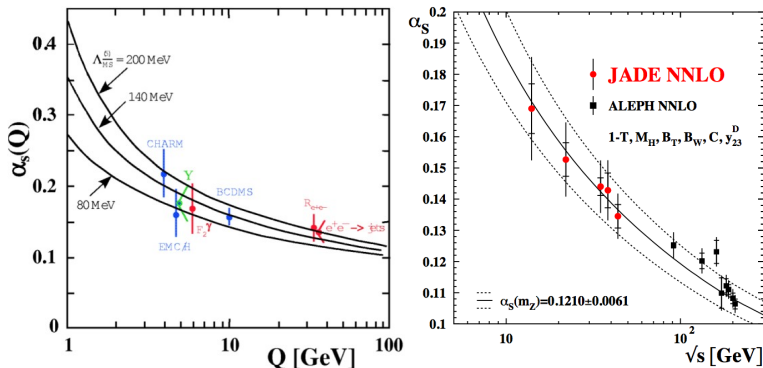
- significant number of publications after end of data taking
- example: ZEUS detector operating on HERA accelerator at DESY



Achim Geiser <https://indico.cern.ch/event/588219>

Why? Uniqueness of data

- JADE experiment (1979–1986) on PETRA accelerator at DESY
- JADE data still cover unique e^+e^- energy range in 2017
- JADE data being re-analysed even ~ 35 years later!



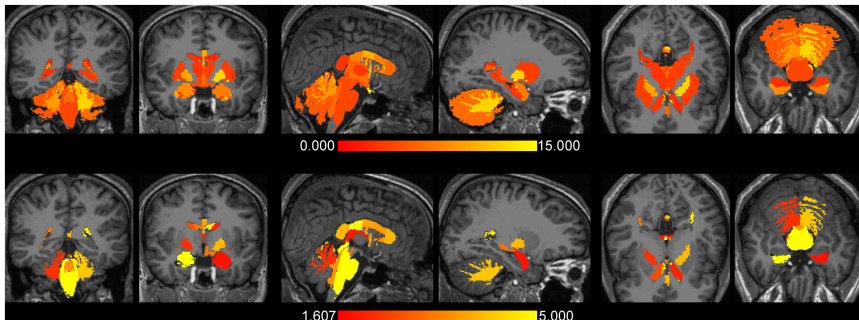
DPHEP <https://arxiv.org/abs/1205.4667>

Data alone is not enough...

The Effects of FreeSurfer Version, Workstation Type, and Macintosh Operating System Version on Anatomical Volume and Cortical Thickness Measurements

Ed H. B. M. Gronenschild , Petra Habets, Heidi I. L. Jacobs, Ron Mengelers, Nico Rozendaal, Jim van Os, Machteld Marcelis

Published: June 1, 2012 • DOI: 10.1371/journal.pone.0038234



$8.8 \pm 6.6\%$ (volume) and $2.8 \pm 1.3\%$ (cortical thickness)

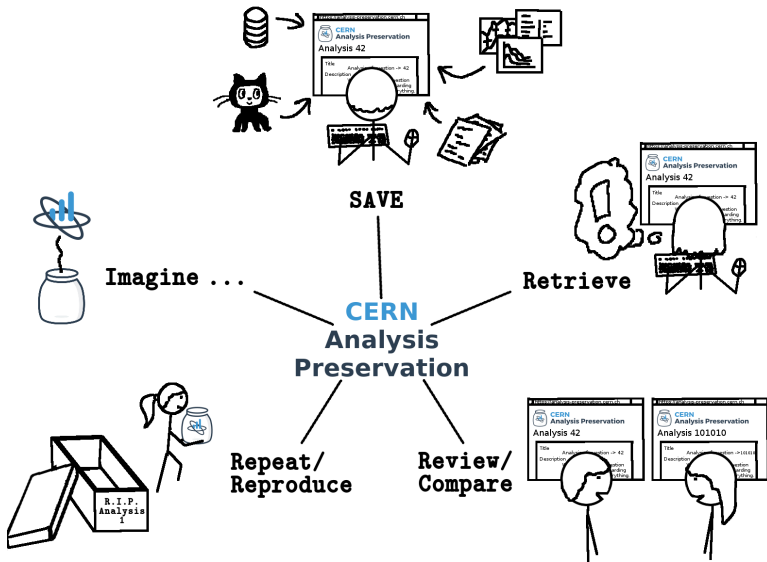
CERN Analysis Preservation

CERN Analysis Preservation

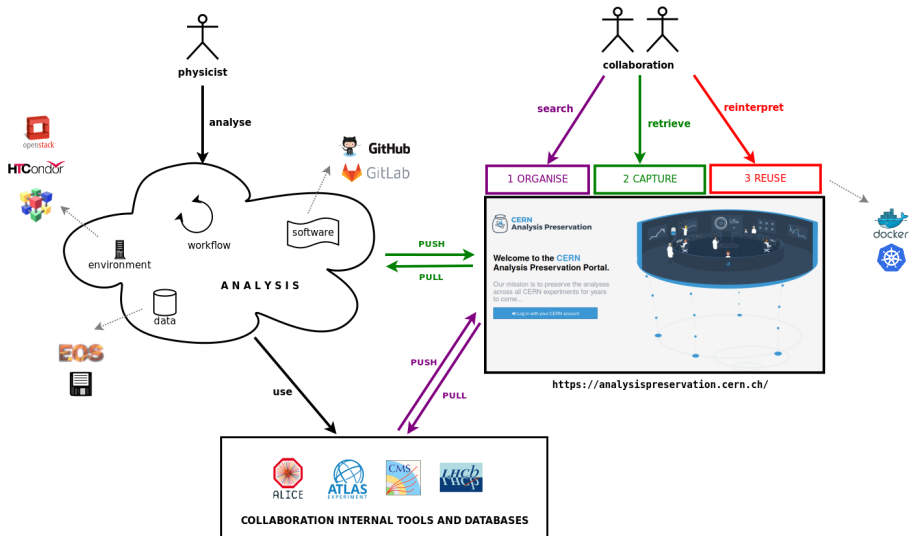
- A platform for **preserving knowledge** and **assets** of an individual physics analysis.
- Capturing the elements needed to **understand** and **rerun** an analysis even several years later:
 - ✓ data
 - ✓ software
 - ✓ environment
 - ✓ workflow
 - ✓ context
 - ✓ documentation
- Advanced **search** for high-level physics information
- Applying standard **collaboration access restrictions**

*Developed by CERN SIS and CERN IT in close collaboration
with LHC experiments*

Use cases



System overview



Three pillars

1 Describe

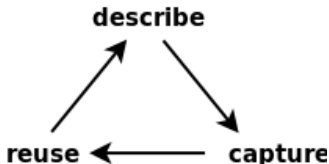
Knowledge modelling
Analysis description

2 Capture

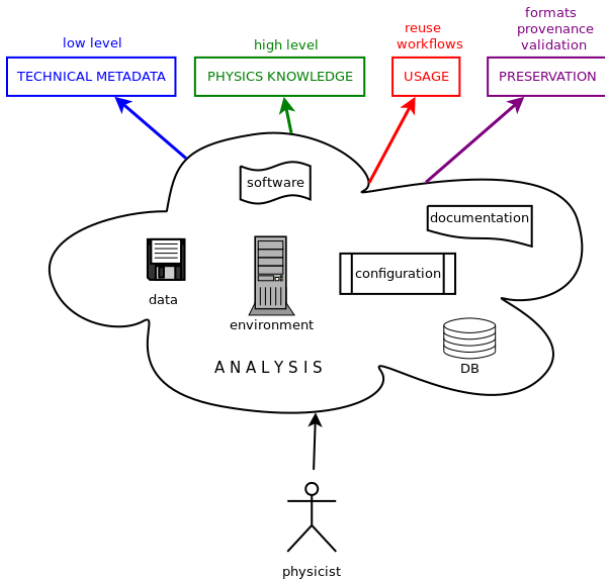
Push: deposit via API
Pull: ingest via grabbing

3 Reuse

Runnable components
Reinstantiate analyses on cloud



1. Describing an analysis

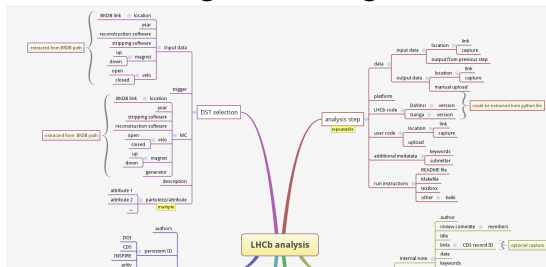


Knowledge representation


■ rare cross-discipline standards (W3C DCAT)

```
"primary_dataset": {  
  "@type": "dcat:Dataset",  
  "title": "/Mu/Run2010B-Apr21ReReco-v1/AOD",  
  "description": "Mu primary dataset in AOD format from RunB of 2010",  
  "licence": "CC0 waiver",  
  "issued": "2011-04-26 11:32:43",  
  "modified": "2011-05-02 21:22:30",  
  [...]
```


■ domain-specific knowledge modelling



Demo: rich physics objects info

**CERN**
Analysis Preservation

CMS ▾

Search 

Save as draft

Filter fields...

BASIC INFORMATION ▸ N/A

INPUT DATA ▸ N/A

N-TUPLE PRODUCTION ▸ N/A



MAIN MEASUREMENT WORKFLOW ▾

- Analysis Note Number
- User Code Base
- Description Details
- Event Selection
- Measurement Description
- Processing


AUXILIARY MEASUREMENT WORKFLOW ▸ N/A

Full Report... N/A


PHYSICS OBJECTS

Create  

Filter fields...
Physics Objects ✕
Add New


Object 

muon


Muon type 


GlobalMuon

NUMBER

<, >, =, <=, >= 

=

Number 

2 

Selection Criteria


☒ Tight

☐ Medium

☐ Loose

☐ Other

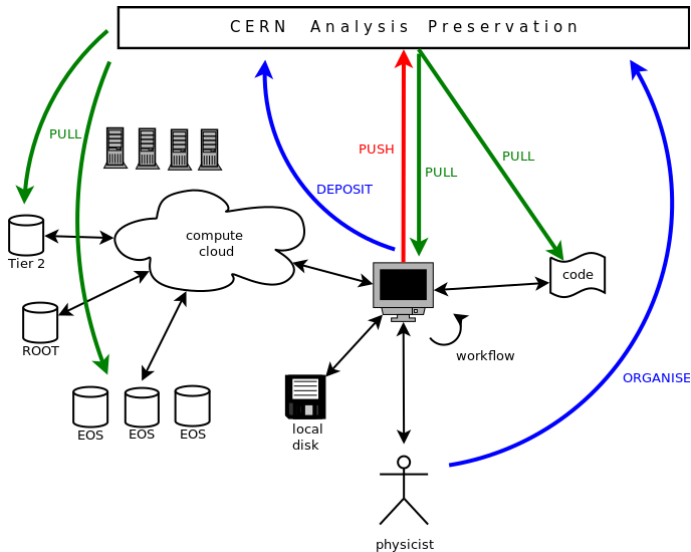
DISCRIMINATOR

Tag 

@tiborsimko

13 / 30

2. Capturing an analysis



Capturing analysis assets

- capturing **datafiles** from various sources:

- local storage
- institute network storage
- WLCG Tier 2 site

via various **protocols**:

- HTTP
- XRootD

- capturing code from various software **repositories**:


- Git
- SVN

- capturing **additional information** from various sources:

- collaboration information databases
- TWiki
- SharePoint

Taking consistent snapshot of information at a certain time

Demo: CMS n-tuple production

**CERN**
Analysis Preservation

CMS

Search

Create

Add New

Save as draft

Filter fields...

BASIC INFORMATION ▶ N/A

INPUT DATA ▶ N/A

N-TUPLE PRODUCTION ▼ N/A

- User Code Base
- Processing

MAIN MEASUREMENT WORKFLOW ▶ N/A

AUXILIARY MEASUREMENT WORKFLOW ▶ N/A

FINAL RESULTS ▶ N/A

ADDITIONAL RESOURCES ▶ N/A

N-TUPLE PRODUCTION

Please provide the n-tuples you used for your measurements

USER CODE BASE

Provide user code

URL

E.g. git@github.com:johnndoe/myrepo.git

Tag

E.g. v2.1

Revision Identifier

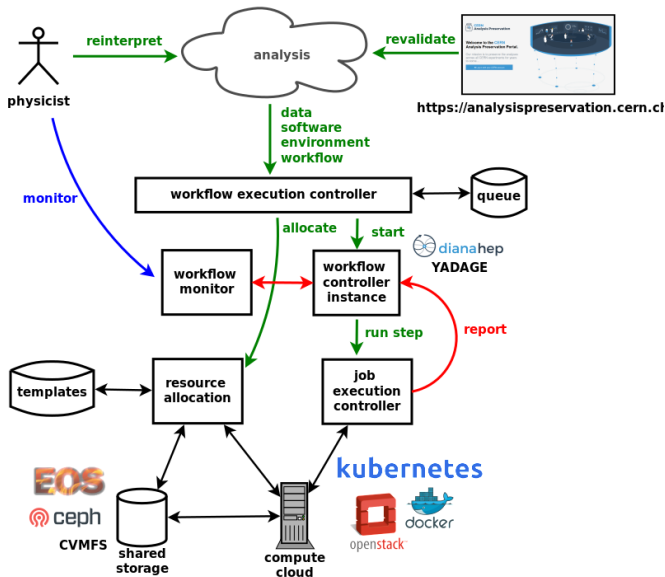
E.g. your git commit hash

PROCESSING

Configuration File

E.g. alt@github.com:johnndoe/./mv-confia-file.root

3. Reusing an analysis

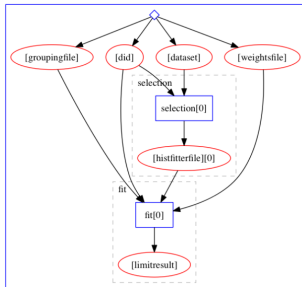


Demo: Reusable analysis pilot

- case study: ATLAS multi-B-jet analysis
- case study: LHCb Lb2LcD0K analysis



```
stages:
- name: selection
  dependencies: ['init']
  scheduler:
    scheduler_type: singlestep-stage
  parameters:
    dataset: {stages: init, output: dataset, unwrap: true}
    submitdir: '{workdir}/submitdir'
    outputprefix: '{workdir}/histfitter.root'
    did: {stages: init, output: did, unwrap: true}
    step: {$ref: 'selscript.yml#'}
- name: fit
  dependencies: ['selection']
  scheduler:
    scheduler_type: singlestep-stage
  parameters:
    bkgtree: 'root://eosuser.cern.ch///eos/project/r/recast/Bkg_2.4.15-2-0_merged.root'
    datatree: 'root://eosuser.cern.ch///eos/project/r/recast/Data_2.4.15-2-0.root'
    outputjson: '{workdir}/fitoutput.json'
    selectionoutput: {stages: selection, output: histfitterfile, unwrap: true}
    weightsfile: {stages: init, output: weightsfile, unwrap: true}
    did: {stages: init, output: did, unwrap: true}
    step: {$ref: 'fitscript.yml#'}
```



Lukas Heinrich <http://github.com/diana-hep/yadage>

Technology

Invenio digital library software



Integrated Library System

Manage MARC21 authority and bibliographic records. Curate records and run automated quality checks. Use circulation module with customisable borrower, item acquisition and interlibrary loan workflows.



Research Data

Capture and preserve research output. Harvest datasets, analysis code, virtual machine environment, configuration and knowledge information. Visualise data in the browser. Rerun preserved code on the cloud.




Multimedia Archive

Manage audio, photo and video material. Create thumbnails and derived formats. Customise portfolio search outputs. Create albums and playlists. Configure related material discovery. Tag multimedia content.

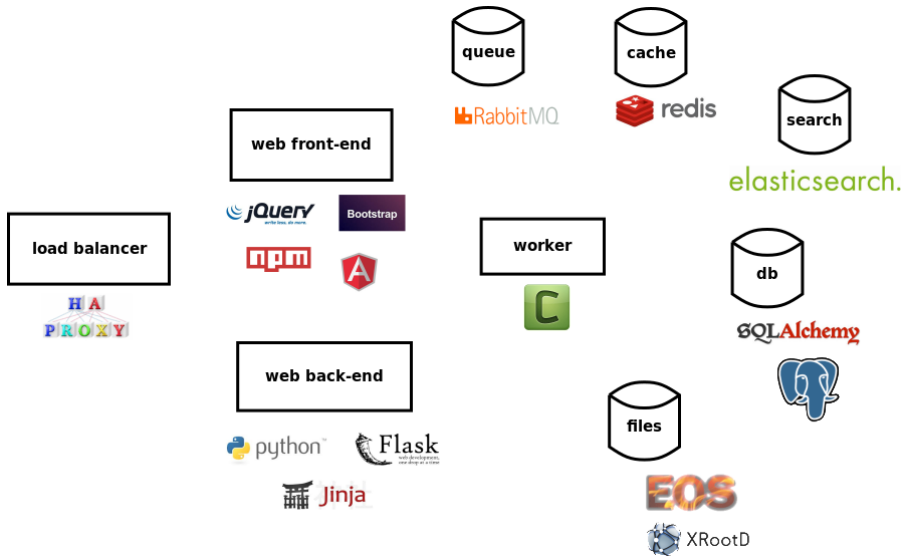


Institutional Repository

Publish articles, reports or theses of your institute. Organise content in collections. Configure ingestion workflows and approvals. Mint material with permanent identifiers. Disseminate material via OAI-PMH.

 <http://inveniosoftware.org>

Invenio technology stack

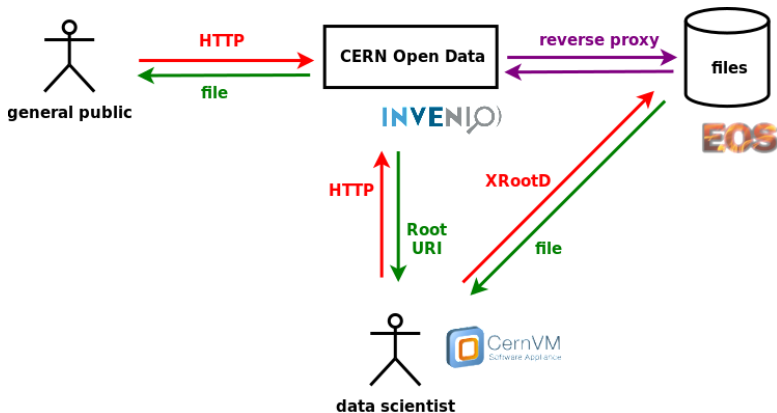


EOS

- **EOS** high-capacity low-latency disk-based storage system

🌐 <https://cern.ch/eos>

- **XRootD** protocol to access parts of data “on demand”



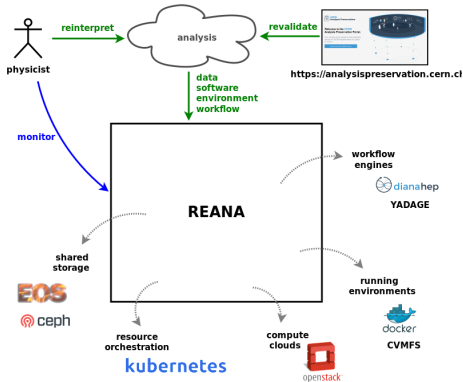
REANA = REusable ANAlyses

- building system to **instantiate** preserved analysis **on the cloud**

🔗 <https://reanahub.io>

- supporting **multiple scenarios**

- multiple computing clouds
→ CERN OpenStack
- multiple running environments
→ Docker with CVMFS
- multiple resource orchestration
→ Kubernetes
- multiple workflow engines
→ Yadage
- multiple shared storage systems
→ Ceph, EOS



- close **collaboration** with **DASPOS** and



Challenges

Social challenges

■ publish or perish culture

- time devoted to preservation = time taken away from the next paper?
- “preservation” platform \rightsquigarrow “live” platform

■ structured workflows

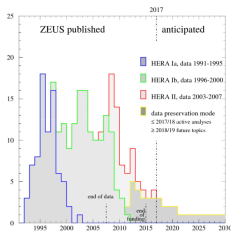
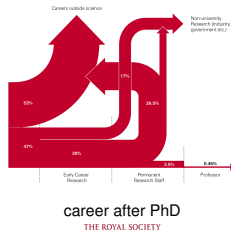
- “runnable READMEs”

■ funding agency requirements

- credible long-term data preservation plans
- standardisation of data formats

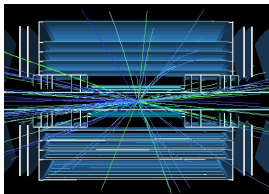
■ scientific **benefit vs cost** of preservation

- Achim Geiser’s study of ZEUS publishing history and long-term preservation efforts:
~10% more papers for <1% of total cost
(of which ~90% during active phase)

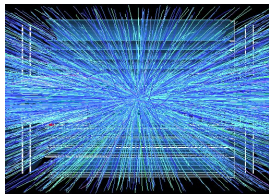


Technological challenges - data

- ever-increasing **data size**?
- example: CMS collision datasets released publicly
 - 2010: ~ 36 out of 40 pb^{-1} available publicly ($\sim 30 \text{ TB}$)
 - 2011: ~ 2.5 out of 5 fb^{-1} available publicly ($\sim 100 \text{ TB}$)
 - 2012: ~ 13 out of 22 fb^{-1} to be released this year ($\sim 0.8 \text{ PB}$)
- High Luminosity LHC upgrade proposal for ~ 2025
 - integrated data increase from ~ 300 to $\sim 3000 \text{ fb}^{-1}$



5 simultaneous collisions



400 simultaneous collisions

<http://atlas.web.cern.ch/Atlas/GROUPS/UPGRADES/>

Technological challenges - code

- ever-changing **computing technology**?
- **history lesson**: JADE dataset resurrection efforts
 - 1986: end of data taking
 - 2016: MPP ports original software from FORTRAN IV (1974), FORTRAN 77, Sheltran, Mortran, Assembler
 - big effort to port ~35 year old code; but data only 600 GB!
- **timeline** projections
 - suppose LHC data taking ends in 2035 — are we looking at being year-2065-compatible following JADE example?
- **encapsulated environments**
 - containers are great, but... Docker? Singularity? Shifter? others?
 - capture live remote service calls (databases)
- **pragmatic approach**
 - focusing at reusability in 1–5 year horizon is already helpful

Conclusions

CERN Analysis Preservation



Welcome to the **CERN Analysis Preservation Portal**.

Our mission is to preserve the analysis across all CERN experiments for years to come.

[Learn more about CERN Analysis Preservation](#)



CERN Analysis Preservation

<http://analysispreservation.cern.ch>

<http://github.com/cernanalysispreservation>

CERN Open Data

<http://opendata.cern.ch>

<http://github.com/cernopendata>

REANA

<http://reanahub.io>

<http://github.com/reanahub>

Invenio

<http://inveniosoftware.org>

<http://github.com/inveniosoftware>

CERN IT H. Hirvonsalo, D. Rodríguez, T. Šimko **CERN SIS** S. Dallmeier-Tiessen, R. Dasler, S. Feger, P. Fokianos, A. Lavasa, A. Mattmann, I. Tsanaksidis, A. Trzcinska **ALICE** M. Gheata, C. Grigoras, M. Zimmermann **ATLAS** K. Cranmer, L. Heinrich, A. Sanchez Pineda, D. Rousseau, F. Socher **CMS** A. Calderon, A. Geiser, A. Huffman, K. Lassila-Perini, T. McCauley, A. Rao, A. Rodriguez Marrero **LHCb** S. Amerio, B. Couturier, S. Neubert, A. Trisovic **CERN CernVM** J. Blomer **CERN Kubernetes** R. Rocha **CERN EOS** L. Mascetti **DASPOS** M. Hildreth, H. Meng, D. Thain, A. Vyushkov **DPHEP** J. Shiers