# Data & Analysis Preservation: Recent Experience and Implications for the EIC
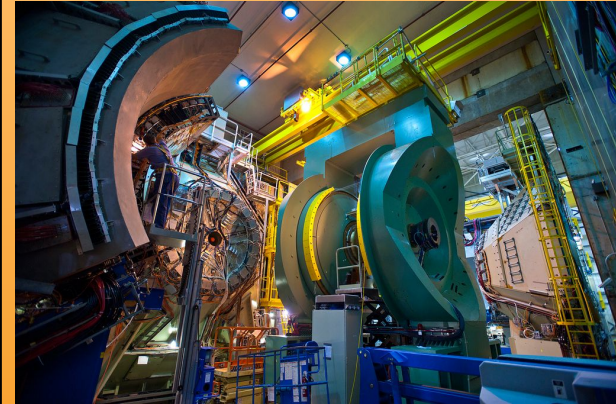
Maxim Potekhin
*Nuclear and Particle Physics Software Group (BNL)*

**Brookhaven**™
National Laboratory

***EIC CompSW Meeting***
07/13/2022

# Overview

- Data and Analysis Preservation (DAP) is by no means a new topic. Its importance in the past two decades was realized based on prior experience at LEP and later at the LHC – with large amounts of data recorded, and overall complexity of the software required to process it

- Both large RHIC Experiments – STAR and PHENIX – are involved in DAP activities

- This stack of slides is informed primarily by participation of the author in the DAP-related work for PHENIX and also in the DPHEP Collaboration led by CERN

- The goal today is to present some "lessons learned" and see if it can be useful for formulation of the initial EIC DAP policies and practices

Brookhaven
National Laboratory

# DAP Motivations (as they relate to the EIC)

- Traditionally, the goals of DAP are to provide new insights into previously analysed data, continue/finalize some analyses over an extended period of time, test new models, create educational and outreach opportunities etc

- However it is also true that implementing DAP helps ensure reproducibility of the results not only on the timescale typically associated with "preservation" – e.g. decades – but also on a much shorter timescale e.g. publication/conference cycle, for the simple reason that at its core it consists of a number of best practices of scientific computing and knowledge and data management

- cf. – even before an experiment starts, reproducibility of the simulation results needs to be guaranteed, documentation ready and accessible etc

# DAP: Challenges and observations

*If there is one lesson in this story it is the need to take a "holistic approach" – data without the software is often useless, as is software without build and verification systems and/or necessary additional data (alignment, calibration, magnetic field maps etc.) These are typically stored separately and involve distinct services that evolve on independent timescales and with lifetimes typically much shorter than the period for which the corresponding "data" needs to be preserved.*

https://doi.org/10.5281/zenodo.2653526 "Software Preservation and Legacy issues at LEP" (J.Shiers)

*No matter what preservation tools are developed that might enable reuse of software, analysis techniques, and data, if they are not conceived from the beginning as an integral part of the standard frameworks, retrofitting will be nearly impossible.*

https://arxiv.org/abs/1810.01191 "HSF White Paper: Data and Software Preservation to Enable Reuse"

Brookhaven
National Laboratory

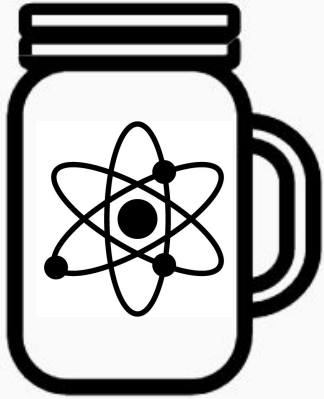# Knowledge Management

- Need to keep records of software provenance, dependencies, configuration, use etc

- Software preservation ≠ Analysis preservation

- Keep track of "data artifacts" such as conditions-type data which may be produced for the purposes of a particular analysis and depend on details known mostly to the people involved in this analysis (misc. cuts, maps, lists, numerical constants in macros etc)

- There is a legacy solution which is a requirement to record such info in a dedicated section of the "Analysis Note" which must accompany every paper, but in reality its efficacy is variable and often insufficient

- Hard to provide continuity of know-how as people move on

- Knowledge dissipates quicker than most people realize
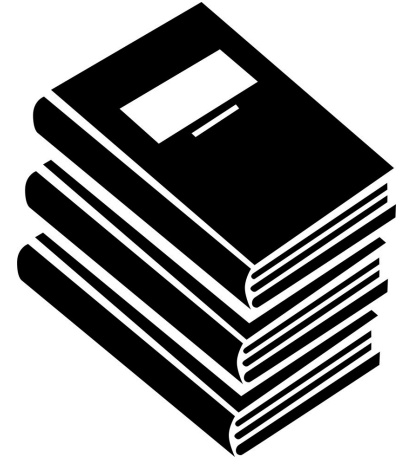
# DAP Components

Bit preservation

Analysis capture
(Knowledge
Management)

Software,
Images, CI,
Configuration

Web-based
documentation

Modern repositories
and portals for research
materials

Brookhaven
National Laboratory

# Tiers of Data Access (incl. HEPData and OpenData)

- **Level 1: Data Products used in publications.**
  - Such as data points and errors used in plots, in numeric format
  - cf. the "HEPData" portal: https://www.hepdata.net/

- **Level 2: Special Purpose Datasets (e.g. for education and outreach).**
  - Select datasets + virtualized or otherwise portable analysis software + documentation
  - cf. the "OpenData" portal: https://opendata.cern.ch/

- **Level 3: Reconstructed Open Data; may be released in future (e.g. based on policy)**
  - Implies a more complex analysis environment than in Level 2
  - Requires adequate software and computing infrastructure to be properly used

- **Level 4: Raw Data. Preserved, but not considered useful for release.**

# Experience with HEPData (Tier-1)



HEPData submissions are a standard practice at the LHC

Now mandated for all new publications in STAR and PHENIX

Revisiting older publication materials as time permits, using GitHub for material development in PHENIX

With established practice and protocols in place, the extra workload is not unreasonable

# The PHENIX OpenData entry - the first for a US-based experiment



- Tier-2/3

- The OpenData portal provides a point of synthesis for software, data and documentation

- Contents of this particular package:
  - Derived data (Ntuples)
  - ROOT macros
  - Detailed instructions (PDF)

- Subject area:
  - Analyses based on the EMcal data

- A promising tool to bring new collaborators up to speed, outreach etc

- Select data published by the LHC experiments

# OpenData capturing MC data provenance (CMS experiment)

## Example 1: Data provenance of simulated datasets

### Simulated dataset BulkGravTohhTohbbhbb_narrow_M-4500_13TeV-madgraph in MINIAODSIM format for 2016 collision data

/BulkGravTohhTohbbhbb_narrow_M-4500_13TeV-madgraph/RunIISummer16MiniAODv2-PUMoriond17_80X_mcRun2_asymptotic_2016_TrancheIV_v6-v1/MINIAODSIM, CMS Collaboration

Cite as: CMS Collaboration (2019). Simulated dataset BulkGravTohhTohbbhbb_narrow_M-4500_13TeV-madgraph in MINIAODSIM format for 2016 collision data. CERN Open Data Portal. DOI:10.7483/OPENDATA.CMS.7N4X.Z7FA

`Dataset` `Simulated` `Exotica` `Gravitons` `CMS` `13TeV` `CERN-LHC`

### How were these data generated?

These data were generated in several steps (see also CMS Monte Carlo production overview):

**Step LHE**
Release: CMSSW_7_1_16
Output dataset: /BulkGravTohhTohbbhbb_narrow_M-4500_13TeV-madgraph/RunIIWinter15wmLHE-MCRUN2_71_V1-v1/LHE
Note: To get the exact generator parameters, please see Finding the generator parameters.

**Step SIM**
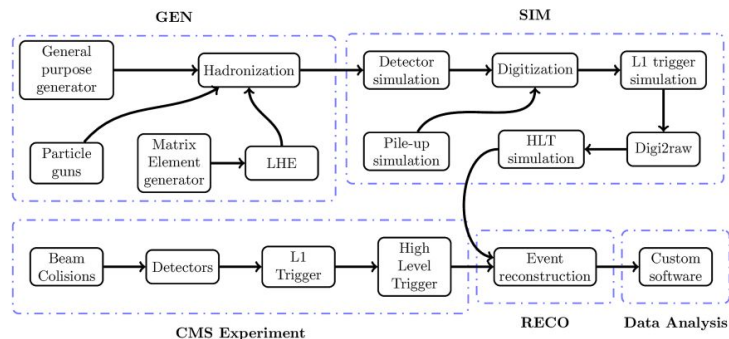Release: CMSSW_7_1_20
📄 Configuration file for SIM (link)
Output dataset: /BulkGravTohhTohbbhbb_narrow_M-4500_13TeV-madgraph/RunIISummer15GS-MCRUN2_71_V1-v1/GEN-SIM

**Step HLT RECO**
Release: CMSSW_8_0_21

GEN | SIM

General purpose generator → Hadronization
Particle guns
Matrix Element generator → LHE

Detector simulation → Digitization → L1 trigger simulation
Pile-up simulation → HLT simulation ← Digi2raw

Beam Colisions → Detectors → L1 Trigger → High Level Trigger → Event reconstruction → Custom software

CMS Experiment | RECO | Data Analysis

► full capture of data generation steps
► full capture of compute environments
► full capture of configuration files
► full capture of production scripts

Data records come with full provenance information

# Rivet (for completeness)

- RIVET: Robust Independent Validation of Experiment and Theory

- "Object-oriented C++ framework for analysis algorithms"

- Can import datasets (experimental data) from HEPData

- The interface to MC event generators is HepMC

- …not a priority for EIC at this point but good to keep in mind

Brookhaven
National Laboratory

# Zenodo@CERN – an example of a "community"/catalog

https://zenodo.org/communities/phenixcollaboration

- A modern digital repository

- >500 PHENIX items, uploads ongoing

- Branded, curated, discoverable, DOI

- Well-suited for long-term preservation
  - Also to support current activity: theses, analysis tutorials, presentations etc
  - Can be used to store data in almost any format and aggregation

- Good search capability
  - Keywords (especially if managed)
  - Elastic search

- Search within the "community" helps to improve results (i.e. cut on false positives)

# Zenodo – a few notes

- Zenodo has uses beyond storing just documents, images etc – e.g. it's used in astrophysics to store reference simulation data samples

- Citable software (DOI's issued for GitHub releases)

- A few tiers of access to documents, from public to restricted to private

- It does not handle the document workflow process e.g. comments and approvals, it was not built for that; however it does support versioning

- DOI capability is the key for long term preservation (existence of a migration mechanism is implied/suggested/committed to)

- In my opinion the most viable solution, with main limitation being the available effort for deployment and maintenance – can this be overcome by the EIC community in the long run, given the importance of such system? Can we get more help from CERN?

Brookhaven
National Laboratory

# Capturing the Software Environment



- We use images to capture a few PHENIX SW environments

- Part of the solution for the changing OS's and dependencies

- NB. interoperability between Docker and Singularity, and the CI angle

- In PHENIX, we are using GitHub to manage Dockerfiles, Docker Hub for image delivery and also a private Docker registry at BNL *to provision software to REANA*

- BNL SDCC provides crucial support in the creation of the "main" image

- CVMFS is used to offload large blocks of binary content thus keeping the image size reasonable, for the "maximal" image representing the complete environment

- NB. ideally, containers should be able to run *anywhere*

Brookhaven
National Laboratory

# REANA

- The PHENIX team successfully ran a few simple and one complex analysis cases

- The workflow description syntax in REANA is a clear improvement compared to a free form assembly of shell and other scripts, as it establishes a minimal structure (YAML)

- The learning curve is particularly easy for linear workflows; but DAGs may be complex

- Individual images can be set for steps in workflows i.e. quite a bit of flexibility

- Experience and REANA instances exist at CERN, U Manitoba, BNL etc

- Application for MC seems straightforward



reana    Home   Examples   Get Started   Documentation   News   Roadmap   Contact   Blog

reana

Reproducible research data analysis platform

| Flexible | Scalable | Reusable | Free |
|---|---|---|---|
| Run many computational workflow engines. | Support for remote compute clouds. | Containerise once, reuse elsewhere. Cloud-native. | Free Software. MIT licence. Made with ♥ at CERN. |

# REANA – a few notes

REANA allows the user to record crucial components of analyses:
- The software environment (by reference to images and libraries, environment etc)
- The workflow(s), in one of the available YAML formats
- Data components to be staged in and staged out, any other auxiliary files that are need

The "workspace" paradigm (essentially a sandbox) enforces completeness of the description and provision of well-defined dependencies.

Also of note is a good CLI, a full Python API and Jupyter integration (e.g. one can open a notebook inside a workspace). The workspace is persistent (if needed).

A variety of computational back-ends is supported (even simultaneously – via hybrid pipelines)

Active community, ongoing development and integration

# Sidenote: workflow diagrams – a low tech but effective tool

- One the right – the actual flowchart used in adopting the direct photon/EMCal analysis in PHENIX, to REANA

- Complementary to a good verbal description

- Now being considered as a potential requirement for future analysis notes in PHENIX, as an effective and relatively low-cost policy

- In combination with REANA, enhances knowledge sharing and transfer within and between working groups

# A note on "services" (from slide 4)

- There can be a variety of conditions-type data (including calibrations) which needs to be preserved and accessible for sensible DAP – so there must be a plan to keep them up and running over a long period of time, or be able to persist them in storage and spin them up on demand, other options…
  - There is experience at CERN with images containing "embedded" CVMFS

- While small "data artifacts" – custom channel maps etc – can be packaged with the analysis software in images or otherwise, it doesn't negate the need for a Conditions Database which is operational over a long period of time and is available "anywhere". Should we store conditions-type derived data used in individual analyses in the Conditions DB? This might help reproducibility since it takes care of version control

- File/data catalog is another service that must be built with durability in mind

# Lessons learned

- DAP: *plan and start early*
  - The effort will pay for itself by increasing overall productivity of the experiment
  - Will be hard or impossible to "catch up" later

- Avoid building in-house information systems, there are many tools available
  - State-of-the-art services such as GitHub, Zenodo, OpenData, HEPData, REANA, Rivet, Inspire (publication catalog) etc cover a vast majority of the experiments' needs
  - There must be no coupling to a particular MC/reco framework

- Containerization solves many of the challenges of capturing the software environment
  - Use it the right (portable) way, with services (DB) made accessible

- Create websites for the long haul (static site generation works well)
  - Avoid platforms that will require updates and maintenance in the long term e.g. Drupal
  - Any resource will become overgrown/obsolete in absence of a designated team of editors
  - Avoid resource fragmentation

# Access control on static websites (e.g. Jekyll-based)

- Sites are always deployed on web servers such as Apache or nginx

- Both of these platforms have auth/auth mechanisms that can be leveraged

- A Jekyll site deployed on nginx, with restricted access, has been tested

- Admittedly role-based access is harder to implement

# Effort profile?

- Example: a DB expert can create and maintain a database for an experiment etc

- DAP is different – it's a true team effort and commitment (one or two experts cannot achieve the goals e.g. knowledge management) – cf. it is impossible to know enough details of someone else's analysis or MC to effectively reproduce/preserve the study – it has been tried without success

- Of course subject experts will be helpful in any case

- Ideally, working group conveners would agree to a set of policies and practices (e.g. designate a few critical simulations for REANA and/or OpenData deployment)

- Historically, DAP has been an underfunded area at RHIC, while some LHC experiments allocated a healthy level of effort to DAP

- There is a fraction of FTE available to help now (yours truly)

# DAP@CERN

- CERN is the preeminent center of DAP expertise and a host to a number of world-class DAP facilities

- The DPHEP Collaboration is the hub of this activity https://dphep.web.cern.ch/

- Example: BNL is a member of the DPHEP Collaboration, both the facility and individual experiments (PHENIX) participated in the DPHEP meetings and other activities

- It would make a lot of sense for the EIC Project Detector collaboration to seek association with DPHEP, to leverage the considerable expertise and resources available through it
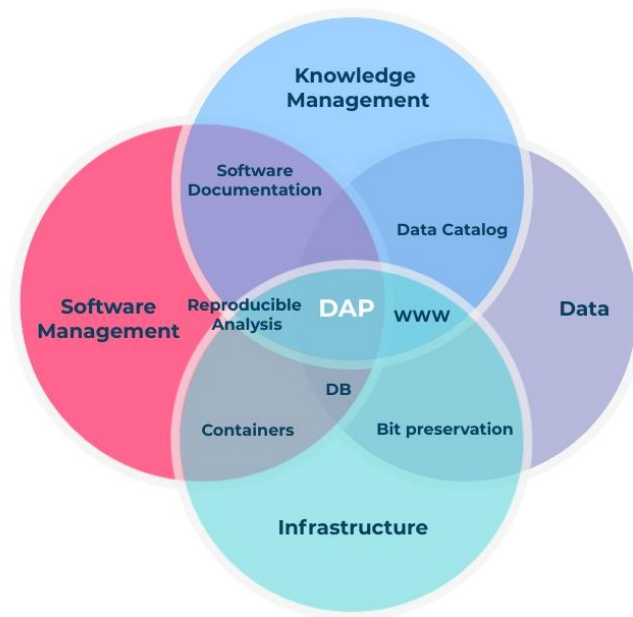
Brookhaven
National Laboratory

# Final Note

DAP practices have the potential to enhance quality of the science output in the near term by helping ensure reproducibility and robustness of the results

DAP focus on knowledge management is conducive to efficient knowledge transfer within the collaboration and across projects (cf. onboarding graduate students)

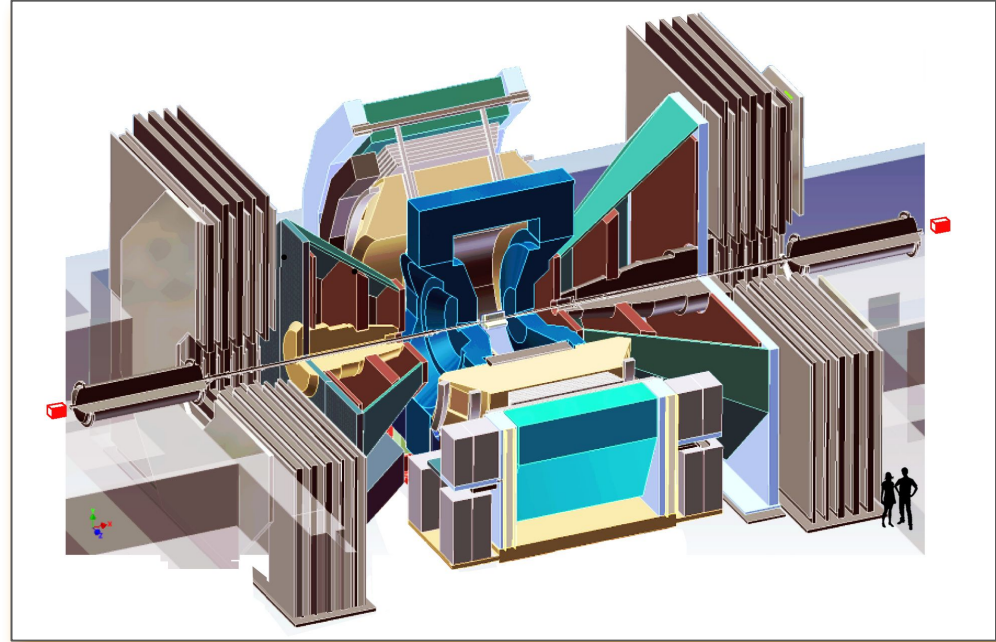Software management, packaging and containerization facilitates deployment

Modern digital repositories create efficient document management solutions on any time scale (cf. OpenData, Zenodo etc)

# Backup

Brookhaven
National Laboratory

# PHENIX

- "Pioneering High Energy Nuclear Interaction eXperiment"

- One of the two large RHIC experiments

- A large, complex general purpose detector with a considerable physics reach and *complex analyses*



- For more details please see the "PHENIX Collaboration Community" on Zenodo: https://zenodo.org/communities/phenixcollaboration/

Brookhaven
National Laboratory

# PHENIX today

- Data taking finished in 2016 with ~24PB of raw data accumulated
- Active analysis work underway (average ~10 articles a year in 2019-2022)
  - Total of >240 published papers + conference contributions (total ~1300 items), ~170 PhD theses and counting

- **All current analyses are done using preserved data**

# PHENIX: Legacy web infrastructure

- Information was spread across a few legacy web resources - the software, detector and subsystem information and other documentation

- Information was diluted with items once relevant for PHENIX but no longer current or aligned with its current and future needs

- PHP-based proprietary information systems e.g. document database (papers, talks, theses), numerical data archive etc became difficult to upgrade, maintain and keep secure - and experiences outages

Brookhaven
National Laboratory

# PHENIX: Software challenges

- Over the years, portability of the core software build procedures was largely lost
  - Build and configuration are specific and coupled to the computing site (BNL)
  - Situation not unique to PHENIX

- Due to compatibility issues most of the software is still built in the i686 environment
  - This does mean extra software packages that need to be installed on modern Scientific Linux
  - Getting hold of certain gcc/OS combinations etc may be a challenge sometimes

- ROOT5 still widely used (and is default) due to many legacy macros
  - Dependencies must also be addressed for the 32-bit build

- PHENIX leveraged Singularity to run production in a containerized environment
  - Motivated by reproducibility, keeping SL6
  - ...with the caveat that the software stack is in AFS - not suitable for REANA