# Data & Analysis Preservation: status update

Maxim Potekhin
*Nuclear and Particle Physics Software Group*

**BROOKHAVEN**
NATIONAL LABORATORY

***PHENIX DAP Meeting***
04/22/2021

**PH ENIX**

# Overview

- Zenodo, Website
- HEPData
- Docker
- REANA
- DAP@PHENIX School
- Open Data

# Zenodo + Website

- Added more conferences
  - Maxim and Gabor
  - GHP17&19, Moriond 19&21, QAT21
- Incremental updates on the website (keywords), minor fixes
  - ...work in progress, catching up with recent work done by Gabor
  - Healthy level of effort/activity
- EMCAL paper added/linked

# HEPData

- There is now a solution to the long standing problem of harmonizing the accuracy of values and errors (i.e. decimal places) in the HEPData submission packages - a Python script was shared by STAR for better control over notation
  - Further improved, committed to our repository
  - Currently in use by a few members
- For example PPG147 HEPData package has been updated (Takashi)
  - Ready to upload
  - Can Ron please be the official reviewer so the submission can be finalized?
- READMEs etc have been updated, spreadsheet up to date
  - Some items stalled, lack of available effort in groups
- PPG201 in development
- *PPG071 successfully revised (an old item where inaccuracies were discovered) - thanks Krista*
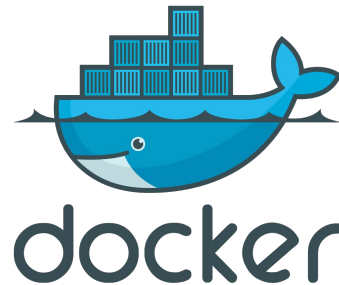
# HEPData - Christine's idea

- There are ~150 papers yet to be committed to HEPData (ballpark)
- This is important to ensure availability of PHENIX data in the long term
- Although the system is reasonable and user-friendly, preparation of HEPData material still takes non-trivial effort (some entries are genuinely complex)
- Due to volume this cannot be addressed by assigning this type of work to interns, summer students or other participants since it goes beyond educational purposes and is effectively a sizeable technical assignment
- Idea - hire undergrads at $12-15 an hour
- Estimated cost of processing ~150 papers is approx. $45k
- Students would work at U of Tennessee and perhaps other locations
- Not necessarily in the scope of BNL-funded work
- What are the possibilities for a proposal? Talk to the DOE Office of Science?
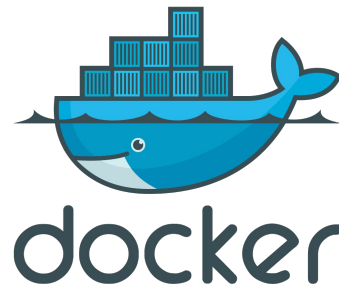
# REANA - quick recap (more to follow)

- REANA: captures the workflow, the software, the data
  - Workflows described as graphs in a YAML format
  - Software captured in containers and/or CVMFS (defined network storage)
  - Data can be uploaded to the sandbox or accessed via the network
- REANA - since it's run at BNL - can help solve the local database access issues (when needed by analyses)
- Containerization is on the critical path

# Docker: status and outlook

- Two approaches are being tested:
  - *1. Copying 32-bit binaries/dependencies to a Docker image*
    - Can be committed to Docker Hub
    - ...see next slide
  - *2. Use SDCC-provided images*
    - *The images have been made available recently, not tested yet - **on the to-do list***
    - *Can't be made public but we will provide them on request or from the Docker registry at BNL*

BROOKHAVEN
NATIONAL LABORATORY

# Docker "1": collection of 32-bit binaries

- *Goes way beyond copying a few libraries into the image*
- O(100) dependencies of both ROOT and PHENIX libraries
- ...not obvious (since libraries are loaded dynamically), dependencies on local installs of various packages - need to be collected, mostly by hand, labor-intensive
  - PHENIX libraries are 1.2GB, root5 is 240MB, package libs are tiny but numerous, with the base SL7 the total is ~2.5GB
  - Having a Docker-capable machine at BNL would have been of help
- **Created a working image that successfully runs PHPythia**
  - AFS mounted on a laptop to collect software
  - Need to test other use cases to validate dependencies
  - See next slide for screenshots
- If problems continue to pop up will probably shift focus to option "2"
- If successful, the image can go on Docker Hub

# Docker "1": 32-bit PHENIX container on a 64-bit laptop, Fun4All etc

# Docker



- ***Important - need simple yet realistic use cases/analyses to test these containers in operation***
  - Can the DAP team help? - install Docker, run a container with some macros?
  - Essentially need samples of simple macros, based on the PHENIX software stack
- Having a valid image is only 20% of work, need input from experts to package analyses for Docker (may be as simple as creating folder with all necessary components)
- Should pick those macros not requiring DB access at first

# REANA and Docker

- Progress with REANA will depend on successful creation of 32-bit images (PHENIX software including fun4all, root etc)
- Minimalistic ROOT macros (e.g. such as created by Gabor) can be used to demo REANA at the PHENIX School right away
- As mentioned earlier, a better estimate in ~2 weeks
- Please suggest analyses to be tested in REANA

BROOKHAVEN
NATIONAL LABORATORY

# School

- Leveraging the PHENIX website to support School activities
  - Makes it easier to stay organized, increases visibility of the site, encourages contributions
- Curated School material should probably be pre-uploaded to Zenodo
- Agendas can easily link to items committed to GitHub, Zenodo, the website
  - Materials no longer trapped in CDS and instead are easily discoverable
- Giving a simple REANA demo at the School makes a lot of sense
  - Currently there isn't a lot of REANA awareness, and it will be useful for students even beyond PHENIX
  - Accounts need to be created prior to School
- Content and complexity of containerized use cases is TBD
  - Need participation of the DAP group
- At least an intro to HEPData makes 100% sense - accounts need to be created prior to School

# Open Data

- Uploaded to CERN, admins notified, awaiting response
- No news for weeks now, will give them more time

**BROOKHAVEN**
NATIONAL LABORATORY

# Plans

- Docker images
- Ongoing HEPData work
- PHENIX School