# The MC Data Catalog - a possible starting point

Maxim Potekhin (BNL/NPPS)

May 12th 2021

**BROOKHAVEN**

# The problem

- *In the long term* a data catalog will likely become necessary

- Need to keep the data *discoverable and accessible*

- During the YR process data discoverability was identified as an issue - but this didn't have much on of impact due to scale
  - People relied on personal knowledge of data location and characteristics and a combination of Wikis, Google Drive folders and GitHub

- With detector simulations coming into the scope of the MC effort the data complexity will increase, bringing about
  - Issues of managing software configuration, both for MC generators and detector simulation, geometry management etc
  - Increased computational costs e.g. "losing" data becoming more expensive

BROOKHAVEN

# Modus Operandi

- If a *Workload Management System* (PanDA, Dirac etc) is implemented there is a need for integration and support for automated data management (including metadata)

- If simpler options are used i.e. direct submission to batch systems there is still a need to maintain records

- At the present stage in the development of the proto-collaborations and detector proposals for the EIC there is no immediately available solution to adopt without substantial effort upfront
  - However the need for such functionality may emerge in the next months - a problem of timescale

- Need a solution with a low cost of entry but extensible and forward-looking
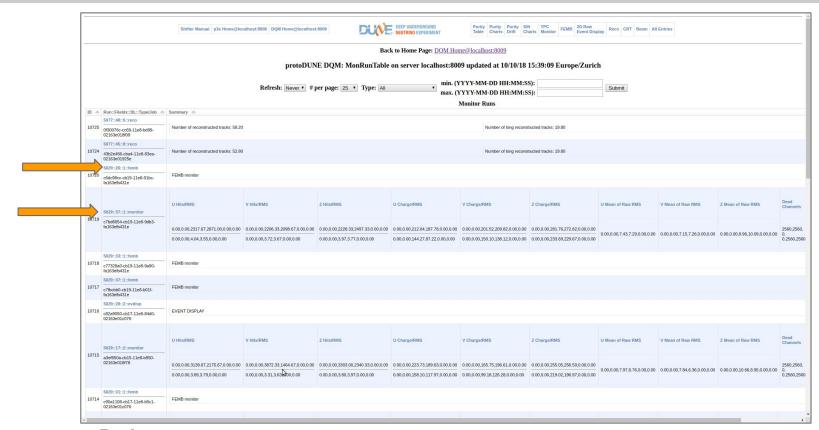
**BROOKHAVEN**

# Experience in protoDUNE-SP

- protoDUNE-SP was an experiment which ran at CERN in 2018-2020
  - A large-scale LAr TPC - a prototype of the Far Detector of the DUNE experiment

- Data Quality Monitoring: prompt processing i.e. several distinct types of jobs
  - Challenge: keep track of diverse data and their provenance, discover and navigate the data

- Solution
  - Save job configuration in a small JSON file
  - Establish a convention whereby each job (which can be of any variety) produces JSON files describing its outputs
  - Can be done in ROOT macros or in wrapper scripts
  - Result: automatic classification of the outputs and straightforward access in the UI

- Comment: in hindsight, using YAML instead of JSON would have been a functional equivalent but would provide much better readability.
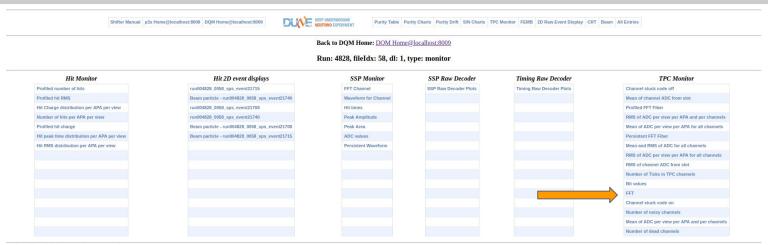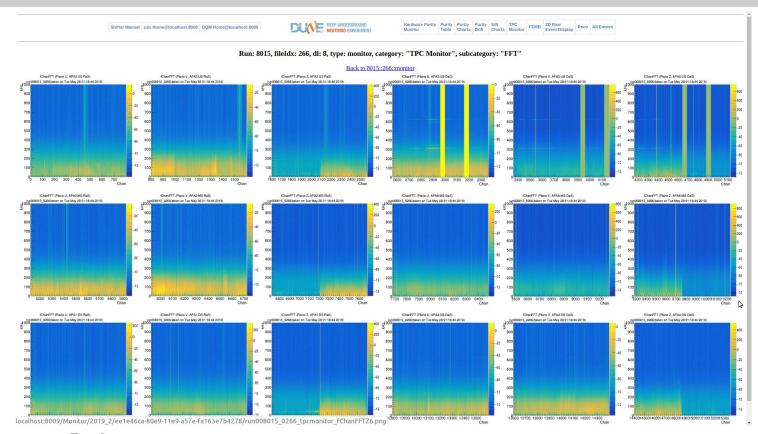
**BROOKHAVEN**

# protoDUNE-SP DQM screen: job summary

# protoDUNE-SP DQM screen: auto-generated menus

# protoDUNE-SP DQM screen: data products

# The descriptor

- A possible solution with a low cost of entry - yet hopefully future-proof - could be to establish a convention where files are accompanied with short descriptors/metadata records formatted in YAML (or JSON)
  - e.g. myFile.root comes with myFile.root.yaml (or, myFile.txt with myFile.txt)
  - The volume/cost of extra data is "small" (technically depends on implementation)
  - Far superior to using filenames and/or folder names as metadata - this is hard to extend and it may not scale well
  - Schema can be updated/augmented at a later time
  - The process is asynchronous i.e. no DB is updated in real time or at all
- If a folder is moved or copied, the descriptors remain with the data

# The contents

- The descriptor would keep information on the file provenance and various aspects of configuration of the software used (all TBD)
  - Type (e.g. MCEG vs G4 etc)
  - Version, references to configuration files (tags)
  - Number of events
  - md5, sha-1 or other hash
  - …
- Configuration files for MCGENs can be referred to by their SHA-1 hash in git/GitHub which guarantees non-ambiguity and audit trail, also is compact

# The catalog

- YAML provides a high degree of (organic) compatibility with the current technology used on the ECCE Software Documentation site (Liquid/Jekyll)

- A master catalog can be compiled and/or recompiled by skimming the data descriptors
  - Can be hosted as YAML or a simple Web app (e.g. Django-based) can be created
  - A few different search mechanisms are available

- **Not locking into any technology at this point**

- Since parsing YAML is "almost free" these metadata can be ingested in future databases chosen by the collaboration
  - Granted, metadata still needs thought and design
  - Future system can be RDBMS or noSQL-like in their properties

## BROOKHAVEN

# Questions, other considerations...

- Files vs datasets?
- Tags
- XRootD
- Capabilities of Rucio
  - Widely used in HEP - at scale
  - "Extensible Metadata" feature has become available in the past year
  - Deployment for EIC requires effort (TBD)

**BROOKHAVEN**