

Data & Analysis Preservation and Open Data: Experience in PHENIX

Maxim Potekhin

Nuclear and Particle Physics Software Group



Future Trends in Nuclear Physics Computing

09/30/2020



This presentation

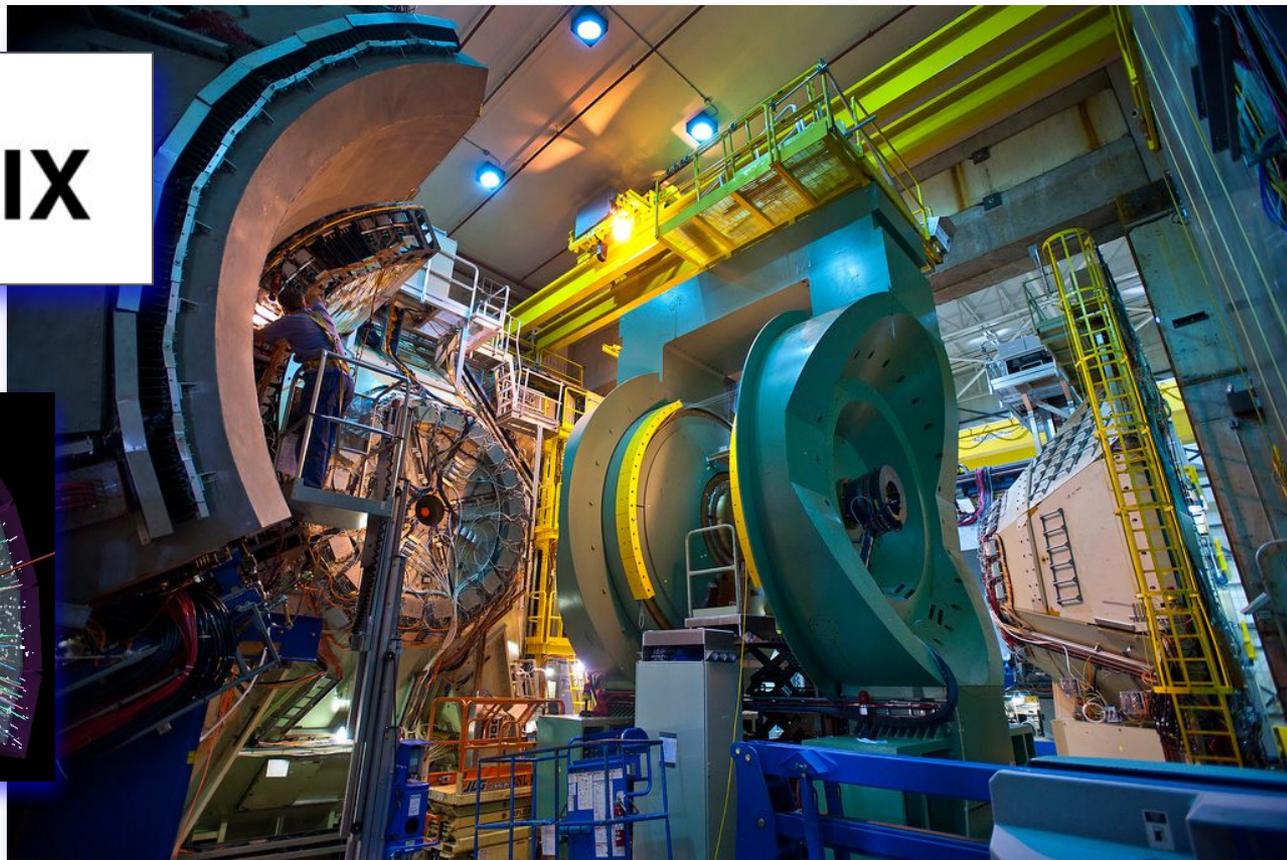
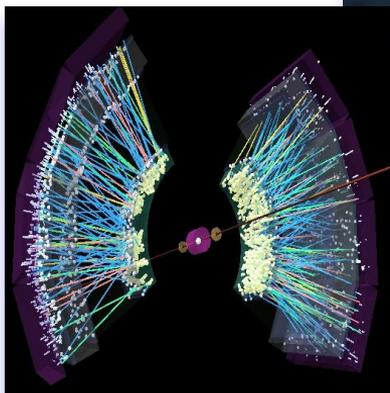
- The PHENIX Collaboration became actively engaged in Data and Analysis Preservation (DAP) in 2019 and is facing challenges in that area common to other experiments (e.g. at RHIC, JLab and elsewhere)
- The goal of today's presentation is to consider ways of leveraging the PHENIX DAP experience for the benefit of the community and identifying venues for collaboration in that area

An Overview

- PHENIX in a nutshell
- Data and Analysis Preservation (DAP)
- DAP: various aspects and challenges in PHENIX
- Technical solutions identified, leveraged or developed for the PHENIX DAP
- Lessons learned + potential for collaboration

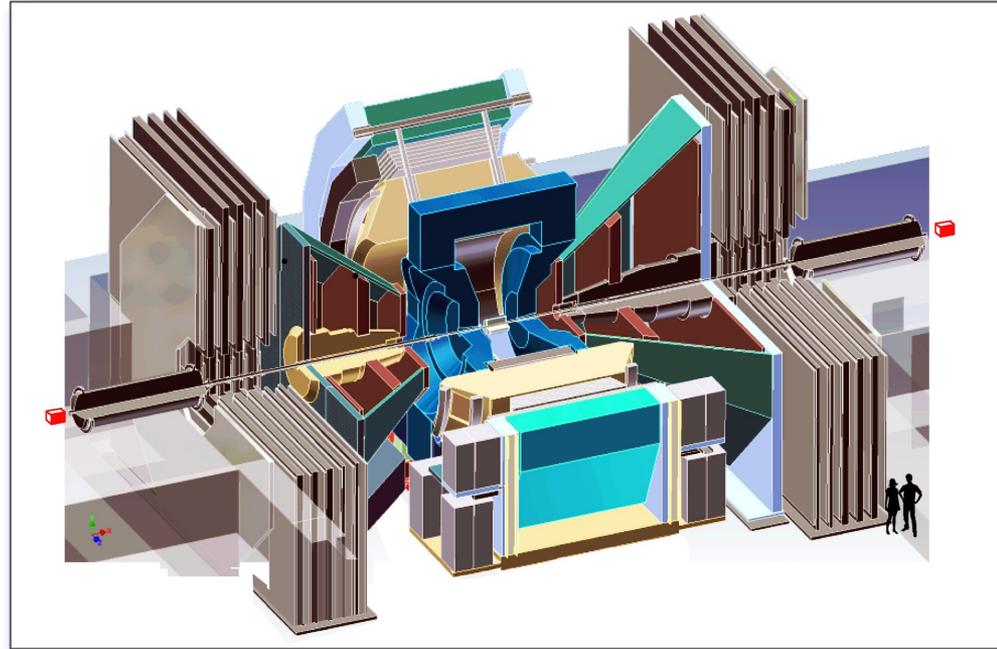
The logo for the PHENIX experiment, featuring the word "PHENIX" in a bold, black, sans-serif font. A red, stylized arc is positioned above the "H", and a white starburst symbol is placed between the "H" and the "E".

PHENIX



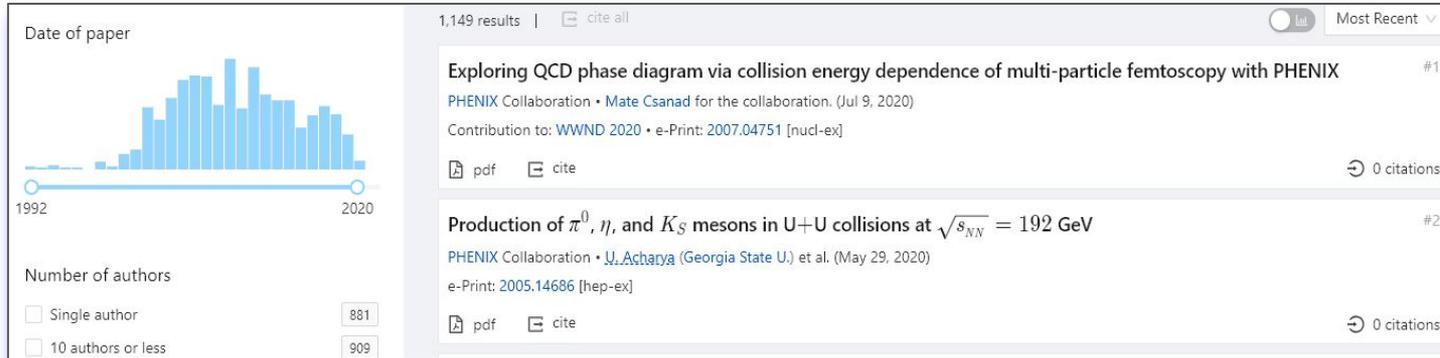
PHENIX in a nutshell

- “Pioneering High Energy Nuclear Interaction eXperiment”
- One of the two large RHIC experiments
- A large, complex general purpose detector with a considerable physics reach
 - Search for QGP and study of its properties, spin physics, other topics
 - Complex analyses
- Please see the “PHENIX Collaboration Community” on Zenodo, the CERN-based digital repository: <https://zenodo.org/communities/phenixcollaboration/>



PHENIX today

- Data taking finished in 2016 with ~24PB of raw data accumulated
- Active analysis work underway (average ~10 articles a year in 2019-2020)
 - Total of >240 published papers + many more conference contributions
 - A total of 165 PhD theses and counting
- Key active analyses considered a priority (effort limited, need to prioritize):
 - Heavy flavor in central and forward rapidity ranges
 - Low pT photons and thermal dileptons
 - High pT direct photon and photon-jets



Definition of the problem as per the “Software Preservation and Legacy issues at LEP” (J.Shiers)

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>

DAP

- DAP = Data and **Analysis** Preservation
 - A tremendous investment which needs to be preserved
 - Potential to provide new insights, validate new models, possibly perform cross-checks etc
 - Synergy with near and medium term goal of maintaining quality of analyses
- Recognized as an important work area by the HEP and NP communities
 - Solid coverage at CHEP 2019 and recent workshops (CERN and others)
 - cf. HSF Community White Paper <https://arxiv.org/abs/1810.01191> and other papers
 - DAP has become a standard item in the experiment and facility reviews
 - However, funding can be a problem - DAP requires additional organizational, software, documentation and other effort
 - While “cost justification” is generally acknowledged such justification does not always result in action and/or funding
- DAP implementation = *Infrastructure and facility + Knowledge Management*

Data preservation: the role of the facility



- DAP universally depends on continuity of services and expertise provided by the facility - this is especially true in PHENIX case where BNL SDCC is the only fully-capable computing site for that experiment.
- In addition to bit preservation (mass storage) the facility provides software building and provisioning capabilities (including containers, CVMFS etc), databases and much more.
- In terms of planning, the facility involvement should cover the time period relevant to a particular DAP effort. Cf. if there is a dependence on AFS it must be somehow satisfied over the years.

Facility contribution to the community infrastructure

- BNL SDCC is a collaborator in the InvenioRDM project
 - BNL lead: C. Gamboa
- A local instance of Zenodo at BNL has been deployed
- ...more details on Zenodo below

Nuclear and Particle Physics Software Group

- <https://npps.bnl.gov/>
- *“The Nuclear and Particle Physics Software (NPPS) Group in Brookhaven National Laboratory's Physics Department participates in a wide range of experiments across BNL's nuclear and particle physics programs. NPPS provides software and expertise across many technical areas, with a particular emphasis on common software solutions.”*
- NPPS provides DAP support for PHENIX and is a natural venue for cross-experiment and cross-site development and collaboration in this area.

PHENIX analysis Knowledge Management: software and data

- Software provenance, dependencies, APIs and configuration (e.g. which specific macros were run, with what arguments, in what sequence and under what conditions, linked to what external libraries etc)
- “Data artifacts” such as conditions- and calibrations-type data which may be produced for the purposes of a particular analysis and depend on details known only to the people involved in this analysis
 - Analysis-specific dead channel maps, recalculated efficiencies, “good run” lists etc
 - Fiducial cuts specific to the analysis (not always documented for reuse)
 - Numerical data in the code (unclear provenance)
- Lots of “moving” parts, not easy to capture and document
 - Even more difficult after the analysis is done, paper published and the team moves on

A legacy Knowledge Management tool: analysis notes

- By policy, each analysis progressing towards eventual publication is expected to be documented in an “analysis note” (similar to other experiments)
 - This is a core instrument to ensure the quality of analyses
- In addition to the physics content each note *should* contain enough information for a qualified physicist to reproduce the results
 - Description and location of (archived) macros and other software
 - Description and location of the data
 - Other data and software dependencies (additional configuration and conditions-type data etc)
- This is codified in the “analysis note template” and policies related to its use
- Analysis code is saved as a tarball (and in CVS for some components)
 - But not really designed for reusability
 - Most of the time no detailed instructions for software configuration and use are provided
- In summary, the analysis notes are not the ideal tool - unless there is an enforceable policy and they become part of the experiment’s culture early on

Knowledge Management: core software documentation

- “PISA” - the core simulation framework in PHENIX (based on Geant3)
- Fun4All - the principal reco framework
- In general, the existing PHENIX documentation for both components is obsolete and existing examples/tutorials are unreliable
- Information is spread over a few Web resources
- This is a work area in itself

Knowledge Management: the detector and its subsystems

- Detailed and evolving experiment configuration information is necessary to maintain the analysis capability
 - Hardware configuration/installation geometry, components upgraded, added or removed
 - Trigger configuration (e.g. EMCAL trigger masks), luminosity, various efficiencies, statistics etc
 - Hard-to-find information e.g. limits on mechanical precision of the detector alignment
- Much of this info spread over multiple Web resources e.g. run pages, Wikis, CAD pages, notes etc - information becomes diluted
- PhD theses are an important component of the documentation, and until recently they were stored in a proprietary DB on the PHENIX server
 - Now being moved to Zenodo (CERN) - complex queries are now possible
- Other detector-related references (NB. need to be mindful of copyright issues)

The HSF Community White Paper: “Data and Software Preservation to Enable Reuse”

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>

An example of a working directory used in an analysis.

Reverse engineering is difficult.

Future Trends

```
check_pdst_vertxfile.root
condor_bb.job
condor_cc.job
condor_hadrons.job
condor.job
condor_jpsi.job
condor_ktomu.job
condor_step3.job
corruptfiles
CorruptFiles_bb.txt
CorruptFiles_cc.txt
CorruptFiles_hadrons.txt
CorruptFiles_jpsi.txt
corruptfiles.txt
CorruptFiles.txt
dcar_bbbblack_jpsired_pdst.pdf
dcar_bbbblack_jpsired_pdst_ylog.pdf
dcar_bb_from_pythia_logy.pdf
dcar_bb_from_pythia.pdf
dcar_jpsi_from_pythia_logy.pdf
dcar_jpsi_from_pythia.pdf
DeadChannels.dat
diff_mydir_cesandir.txt
DST_embed_bb
DST_embed_bbbbar
DST_embed_cc
DST_embed_hadrons
DST_embed_hadrons_3gev
DST_embed_jpsi
DST_embed_ktomu
DST_pythia_hadrons
dst_pythia_hadrons-MB0-0000429896-0000.list
dstpythia.list
dstpythia.root
dstpythia.txt
DST_sim_bb
DST_sim_bbbbar
DST_sim_cc
DST_sim_hadrons
DST_sim_hadrons_3gev
DST_sim_hadrons_combined
DST_sim_hadrons_minexuanSample
DST_sim_jpsi
DST_sim_ktomu
embedpythia.pdf
env.log
env.txt
ERR
event_gen_bb.csh
event_gen_bb_step2.csh
event_gen_bb_step3.csh
event_gen_cc.csh
event_gen.csh
event_gen.csh~
event_gen_hadrons.csh
event_gen_hadrons_step1.csh
event_gen_hadrons_step2.csh
event_gen_hadrons_step3_reass.csh
event_gen_hadrons_step3_xuan.csh
event_gen_jpsi_combined.csh
event_gen_jpsi.csh
event_gen_ktomu.csh
geom.root
geom_run15_v2.root
get_nevents
get_nevents.C
hadd_files.sh
HV
jpsipdst.txt
landau_parameters.txt
libpicodst_object.so
listofcc.txt
listofpythia.txt
LOG
logall.txt
logfiles.txt
logfile.txt
loglog.txt
logout.txt
logs.txt
log.txt
make_picodstobj.C
mc_bjpsi.C
mc_bjpsi_C_ACLIC_dict_rdict.pcm
mc_bjpsi_C.d
mc_bjpsi_C.so
mc_bjpsi.h
merge_picodsts.C
model_bdecay.C
muid_tube_eff_north_Run15pp200.txt
muid_tube_eff_south_Run15pp200.txt
mut_disabledwipes.dat
mvcorruptfiles_bb.sh
mvcorruptfiles.sh
mylogfiles.txt
mypythia_cmu_pdst.pdf
mypythia.pdf
mypythia_pdst.pdf
mysimhad.txt
nokeys.txt
Nonexistingfiles.list
Nonexisting.txt
nooffiles.txt
notopen_bb_1.txt
notopen_bb_2.txt
notopen_bb_3.txt
notopen_bb_4.txt
notopen_bb_5.txt
notopen_bb.txt
NotOpenfiles.txt
notopen_hadrons_ajeeta.list
notopen_jpsi_1.txt
notopen_jpsi_2.txt
notopen_jpsi_3.txt
notopen_jpsi_4.txt
notopen_jpsi.txt
notopen.txt
old
old_pdst_bb
out.log
output_bb_step2.txt
output_bb_step3.txt
output_bb.txt
output_env.txt
PISA
PisaFullfiles.txt
pisa.txt
practice.list
pylist12.dat
#pythia_bb.C#
pythia_bb.C
pythia_bb.C~
pythia_bb_C_ACLIC_dict_rdict.pcm
pythia_bb_C.d
pythia_bb_C.so
pythia_bb.h
pythia_bb.h~
pythia_configuration
pythia_files_bb
pythia_files_bbbbar
pythia_files_bb_iter0
pythia_files_bb_list.txt
pythia_files_cc
pythia_files_hadrons
pythia_files_hadrons_3gev
pythia_files_hadrons_combined
pythia_files_hadrons_list.txt
pythia_files_hadrons_minexuanSample
pythia_files_jpsi
pythia_files_jpsi_combined
pythia_files_jpsi_list.txt
pythia_files_ktomu
pythiafiles.txt
pythia.txt
realdataBG-run15pp_file-forAjeeta.list
realdataBG-run15pp_file-forAjeeta_onefile.list
realdataBG-run15pp_file-forCesar.list
realdataBG-run15pp_file-forXuan.list
realdataBG-run15pp_file.list
realdataBG-run15pp_file_split-forAjeeta1.list
realdataBG-run15pp_file_split-forAjeeta2.list
realdataBG-run15pp_file_split-forAjeeta.list
realdataBG-run15pp_file_split-forCesar1.list
realdataBG-run15pp_file_split-forCesar2.list
realdataBG-run15pp_file_split-forCesar.list
realdataBG-run15pp_file_split-forXuan1.list
realdataBG-run15pp_file_split-forXuan2.list
realdataBG-run15pp_file_split-forXuan.list
realdataBG-run15pp_file_split.list
Run_reassociation.C
run_segments.list
Sim3D++.root
SimFileLists
singlemuon_embed_pdst.root
singlemuon_jpsi_embed_pdst-MB0-0000422070-0001-2.root
split.csh
split_dsts.C
splitfiles
svxPISA.par
temp.txt
totest.list
whencreated.txt
with_n_option.txt
Xuan_allFailedFiles.txt
xuanfilelist.txt
```

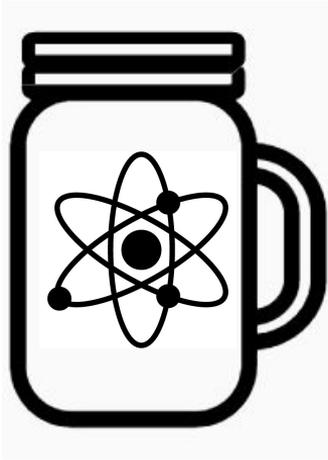
Challenges in PHENIX

- “Retrofitting is nearly impossible” (as it concerns various DAP tools)
 - So the next best thing is an accurate record of the analysis procedure and its components
- Diminishing manpower
 - Progressively smaller teams at each participating institution
- Dissipation of expertise, combined with insufficient maintenance of many types of documentation and its fragmentation
- Analyses still ongoing - pressure on PhD students, postdocs and other personnel to get it done with modest available resources
 - DAP requirements being a significant additional hurdle, hard to find reliable volunteers
 - Currently, making an extra effort is not rewarded in any meaningful way (see next slide)
- New students need training and solid information sources to learn from
- The PHENIX software is not easily portable

DAP and Software Sustainability

- A few points from the presentation by Daniel Katz at this workshop on 9/29/2020 are quite relevant for DAP, with regards to the analysis software
- Quote:
 - *Software development and maintenance is human-intensive*
 - *Much software developed specifically for research, by researchers*
 - *Researchers know their disciplines, but often not software best practices*
 - *Researchers are not rewarded for software development and maintenance in academia*
- The word “incentive” is used 3 times in Daniel’s talk
- Incentivizing could be the key to changing the culture - DAP is serious work and it needs recognition

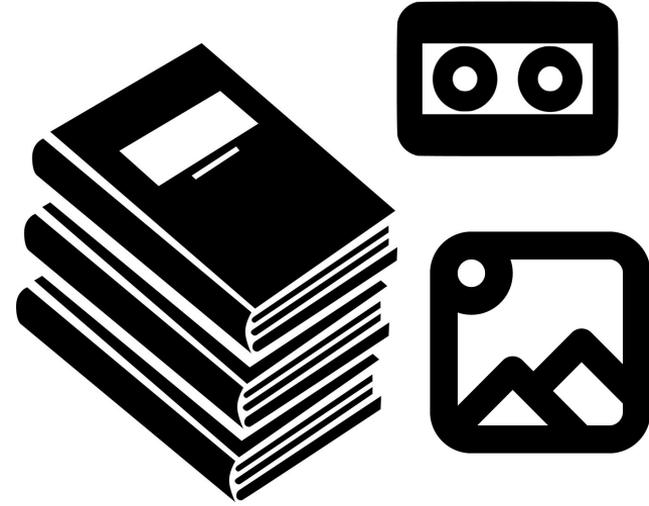
A practical approach to the components of the Knowledge Management in PHENIX



Analysis capture



Web-based documentation



A highly-functional repository for research materials

The Analysis Environment

- Core analyses take place within the SDCC facility at BNL
- AFS is used for software provisioning (regular builds etc)
- SL7
- A large part of the analysis code is managed in CVS (legacy)
- Some of the conditions data is kept as files accessed via scripts
 - Part of the reason for creation of custom data products managed in individual analyses
- The cost/benefit situation in PHENIX currently does not favor creation of containers for analysis purposes
 - Containers are used in production for separate reasons

Analysis Preservation

- In practice we can't capture all - or many - analyses for the following reasons
 - PHENIX started >20 years ago with no defined DAP policy, plan or framework
 - **this changed in 2019**
 - Once a postdoc or a graduate student leaves, there is little chance of knowledge recovery
 - Only current/recent analyses are viable for preservation
- There is an overlap between medium term analysis capability and longer term analysis preservation - an extra motivation for DAP right now
- *Creating up-to-date summaries of analysis techniques (e.g. PID etc) is universally helpful - progress is being made*

Web-based documentation: the new PHENIX website

- <https://www.phenix.bnl.gov/>
- The hub for all DAP information, designed to be a long-term resource
- Completely new layout and interface
- Content consolidated from a number of legacy web resources
- Eliminates in-house DB management, minimizes local storage and instead relies on external services for hosting content:
 - Zenodo
 - HEPData
 - GitHub
- Where possible, relies on DOIs for references, to achieve long-term integrity

The website technology

- Inspired by the websites of the HEP Software Foundation and NPPS
 - <https://hepsoftwarefoundation.org/>
 - <https://npps.bnl.gov/>
- The website is static (no DB) - using the Jekyll static website generator
 - HTML is generated, not written by hand
 - Easy to contribute (Markdown, YAML) and easy to maintain
 - Data content kept separately from the layout
 - Manipulation of structured data in Jekyll makes for a compact and efficient design of pages, and automated content generation - DB-like functionality at compile time
 - GitHub Pages for management and development + version control + Jekyll builds
 - Portability (easy to export/migrate the complete site as HTML)
 - Performance
 - Security
- The “Bootstrap” toolkit for the layouts/navigation (no custom JS/CSS needed)

Structured data used to generate the site

```
##### RUN 12
- run: run12
  title: Run 12
  period: 2011-2012
  coordinator: Xiaochun He, GSU.
  rhic:
  - {
    species: 'polarized p+p',
    energy: 100.2,
    lumi: '- /10 <i>pb</i><sup>-1</sup>',
    Nevents:
  }
  - {
    species: 'polarized p+p',
    energy: 254.9,
    lumi: '32/- <i>pb</i><sup>-1</sup>',
    Nevents:
  }
  - {
    species: '<sup>238</sup>U<sup>92+</sup><sup>238</sup>U<sup>92+</sup>',
    energy: 96.4,
    lumi: '0.2<i>nb</i><sup>-1</sup>',
    Nevents: 1.28/0.88
  }
  - {
    species: '<sup>63</sup>Cu<sup>29+</sup><sup>197</sup>Au<sup>79+</sup>',
    energy: 99.9+100.0,
    lumi: '5<i>nb</i><sup>-1</sup>',
    Nevents: 0.88/8.18
  }
  - {
    species: '<sup>197</sup>Au<sup>79+</sup><sup>197</sup>Au<sup>79+</sup>',
    energy: 2.5,
    lumi: '-',
    Nevents: Very short
  }
  ert_comment: Summary of thresholds (DAC values). Values in parentheses are for the PbG1.
  ert_thresholds:
  - '02/09/12, 358208, 30(29), 31(30), 29(29), 29(25), 920, Run12pp200 - Pedestal tuned - EMCal dynamic range ~25GeV'
  - '03/20/12, 364957, 30(29), 31(30), 29(29), 29(25), 920, Run12pp510 - EMCal dynamic range ~50GeV'
  - '04/23/12, 369200, 30(29), 31(30), 29(29), 29(25), 920, Run12UU193 - EMCal dynamic range ~25GeV'
  - '05/16/12, 372155, 31(30), 32(31), 30(29), 49(45), 920, Run12CuAu200 - EMCal dynamic range ~25GeV'
```

The new PHENIX website

<https://www.phenix.bnl.gov/>

The screenshot displays the PHENIX website interface. At the top, a navigation bar includes links for Experiment, Detectors, Software, Analysis, Results, Resources, and About. The main content area features the PHENIX logo and a large image of the detector interior. A text block describes PHENIX as the Pioneering High Energy Nuclear Ion Collider, noting that data-taking was finished in 2012 and large data samples are now being analyzed.

A secondary screenshot shows a dropdown menu for the 'Detectors' section, listing: Detectors Overview, PHENIX Photo Gallery, Run Configuration Gallery, Subsystems (highlighted), Central Arm Detectors, Muon Arm Detectors, Event Characterization Detectors, and Magnet.

On the right, a 'RHIC records + PHENIX run summary table' is shown, with a navigation bar for runs 01 through 16. The table below provides detailed run information:

Run	Species	Energy (GeV/nucleon)	Integrated Luminosity (Polarization L/T)	N _{events} [BBC _{30cm} /BBC _{narrow}]
01	¹⁹⁷ Au ⁹⁺ + ¹⁹⁷ Au ⁹⁺	65.2	1b ⁻¹	10M
02	¹⁹⁷ Au ⁹⁺ + ¹⁹⁷ Au ⁹⁺	100.0	24b ⁻¹	10M
	¹⁹⁷ Au ⁹⁺ + ¹⁹⁷ Au ⁹⁺	9.8	-	<1M
	polarized p+p	100.2	~ / 0.15 pb ⁻¹	3.78
	d + ¹⁹⁷ Au ⁹⁺	100.7+100	2.74nb ⁻¹	5.58
03	polarized p+p	100.2	0.35 ~ pb ⁻¹	6.68
04	¹⁹⁷ Au ⁹⁺ + ¹⁹⁷ Au ⁹⁺	100.0	241μb ⁻¹	1.58
	¹⁹⁷ Au ⁹⁺ + ¹⁹⁷ Au ⁹⁺	31.2	9μb ⁻¹	58M
	⁶³ Cu ²⁹⁺ + ⁶³ Cu ²⁹⁺	100.0	3nb ⁻¹	8.68
	⁶³ Cu ²⁹⁺ + ⁶³ Cu ²⁹⁺	31.2	0.19nb ⁻¹	0.58
	⁶³ Cu ²⁹⁺ + ⁶³ Cu ²⁹⁺	11.2	2.7nb ⁻¹	3.78
	polarized p+p	100.2	3.4/0.2 pb ⁻¹	858
06	polarized p+p	100.2	7.5/2.7 pb ⁻¹	2338
	polarized p+p	31.2	0.08/0.02 pb ⁻¹	288
07	¹⁹⁷ Au ⁹⁺ + ¹⁹⁷ Au ⁹⁺	100.0	813μb ⁻¹	5.18
08	d + ¹⁹⁷ Au ⁹⁺	100.7+100	80nb ⁻¹	1608
	polarized p+p	100.2	~ / 5.2 pb ⁻¹	1158

An example of a subsystem page

Experiment Detectors Software Analysis Results Resources

Electromagnetic Calorimeter

Write-ups

- DOI: 10.5281/zenodo.3833205 PHENIX Electromagnetic Calorimeter (EMCal) – Detector Basics (G.David)
- DOI: 10.5281/zenodo.3893972 Explanation of PHENIX triggers (A.Bazilevsky)

Theses

- DOI: 10.5281/zenodo.3885856 The Quark Gluon Plasma probed by Low Momentum Direct Photons in Au+Au Collisions at $\sqrt{s_{NN}}=62.4$ GeV and $\sqrt{s_{NN}}=39$ GeV beam energies (Vlad Khachatryan)
- DOI: 10.5281/zenodo.3885870 Inclusive jet production in proton-proton and copper-gold collisions at $\sqrt{s_{NN}} = 200$ GeV (Arbin Timilsina)

Publications

- PHENIX Calorimeter (NIM A 499, 2003, doi.org/10.1016/S0168-9002(02)01954-X)
- High Energy Beam Test of the PHENIX Lead-Scintillator EM Calorimeter High Energy Beam Test of the PHENIX Lead-Scintillator EM Calorimeter

Presentations

- DOI: 10.5281/zenodo.4007113 PHENIX Focus: Electromagnetic Calorimeter (Gabor David)

Variables and Accessors under PHCentralTrack Node (used for charged particle analyses)

Type	Name	Description
float	get_pemcx	x-component of the projection of the cgl track onto the EMC (cm)
float	get_pemcy	y-component of the projection of the cgl track onto the EMC (cm)
float	get_pemcz	z-component of the projection of the cgl track onto the EMC (cm)
float	get_lemc	path Length following particle trajectory from vertex to EMC
float	get_temc	time of the EMC hit. This time has been back-corrected inPHCentralTracks to be the physical time instead of the photon flash time. The reason is that the former is more useful for calculating properties of a charged track.
float	get_emcdphi	difference in phi (rads) between the track model projection and the hit in emc
float	get_emcdz	difference in Z (cms) between the track model projection and the hit in emc
float	get_emcdsphi	emcdphi variable normalized to SIGMAS (after calibrations)
float	get_emcdsdz	emcdz variable normalized to SIGMAS (after calibrations)
		position resolution of the EMCal depends upon the shower type, emcdphi variable in SIGMAS assuming the resolution appropriate for EM showers

CERN and community tools for DAP

DPHEP Data Preservation in High Energy Physics
Collaboration for Data Preservation and Long Term Analysis in High Energy Physics

Partners Accelerators Meetings ICFA Study Group About Us

FOLLOW THE LINKS BELOW TO FIND INFORMATION ON OUR PARTNER ORGANIZATIONS. EACH REPRESENT SOME EXPERIMENTS AND ACCELERATORS TO THE COLLABORATION FOR DATA PRESERVATION IN HIGH ENERGY PHYSICS.

BNL Home
CERN Home Data
CSC Home
DESY Home
Fermilab Home
IHEP Home
IN2P3 Home
INFN Home
IPP Home
KEK Home
SLAC Home
STFC Home

reana

Reproducible research data analysis platform

 HEPData

zenodo

providing
research data
services.

DPHEP Study
Group *et*
A common
reflection on data
persistence and
long term
analysis in High
Energy Physics.

opendata
CERN

iNSPIRE HEP

Leveraging CERN-based tools

- Using **InspireHEP** to keep the publication list and cross-reference HEPData
- **Zenodo**
 - Based on the Invenio framework created at CERN
 - Active migration of relevant documents from legacy services to Zenodo is under way
 - Using a curated list of keywords for - easy to find physics topics, detector elements, conference materials etc
 - Easy to cite using the DOIs (persistent identifiers)
 - BNL is a collaborator in the InvenioRDM project - the next generation repository
- **HEPData**
 - Based on the Invenio framework
 - Recently renewed effort to create submission packages and commit new data to the service
 - **As of this week, new policy by PHENIX IB in place: HEPData package as a prerequisite for publication**
- The CERN Open Data Portal - access under discussion, material in the works
- REANA, Rivet and other tools
 - Use of these tools is only possible with highly organized and documented active analyses
 - Retrofitting is quite hard, limited available effort likely means these tools won't be used
 - Instead, annotate/document analyses (in progress)

Zenodo@CERN - the PHENIX community

- Branded
- Curated
- Discoverable
- Indexed (keywords)
- Elastic search capability

The screenshot displays the Zenodo interface for the PHENIX Collaboration. The top navigation bar includes the Zenodo logo, a search bar, and links for 'Upload' and 'Communities'. The user profile 'phenix-dap-l@lists.bnl.gov' is visible in the top right. The main content area is titled 'PHENIX Collaboration' and features a 'Recent uploads' section with a search bar and a list of three items. Each item includes a date, a title, an author, a brief description, and an upload date. The first item is 'π0-hadron correlations in 200GeV Au+Au collisions' by Wong, Cheuk-Ping. The second is 'PHENIX measurement of system size dependence of low momentum photon production' by Esha, Roli. The third is 'Study of jet modifications at PHENIX using two-particle azimuthal correlations and high-pT hadrons' by Wong, Cheuk-Ping. A 'New upload' button is located at the top right of the content area. On the right side, there is a 'Community' section with the PHENIX logo, a description of the community's purpose, and details about its curation and creation. At the bottom right, there is a section titled 'Want your upload to appear in this community?' with a bullet point indicating that users should click the 'New upload' button.

An example of a PHENIX item on Zenodo

December 1, 2013 Thesis Open Access Edit

Low Momentum Direct Photons as a Probe of Heavy Ion Collisions

Petti, Richard
Thesis supervisor(s)
Drees, Axel

Essential to the study of heavy ion collisions are probes that are produced in the collision itself. Photons are a very useful probe of the collisions, since they escape the fireball virtually unmodified and carry with them information about the environment in which it was produced. Recent interest in low momentum direct photons has increased, due to the onset of the "thermal photon puzzle" and the apparent inability for typical models to explain both a large direct photon yield excess and large azimuthal production asymmetry (v_2) at low momentum measured by PHENIX.

Preview

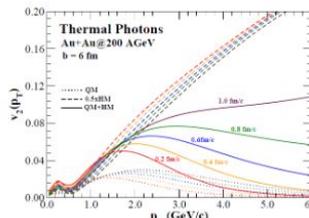


Figure 1.20: A calculation of the thermal photon v_2 from [23]. The dotted curves represent the v_2 of thermal photons emitted from the QGP, dashed curves represent the v_2 of thermal photons emitted from the hadron gas, and solid curves represent the time averaged thermal photon v_2 integrated over the entire evolution of the system. The various colors represent the calculation

Files (10.8 MB)

8 views 5 downloads [See more details...](#)

Indexed in **OpenAIRE**

Publication date: December 1, 2013
DOI: [10.5281/zenodo.3887326](https://doi.org/10.5281/zenodo.3887326)

Keyword(s): RHIC, direct photon, PID, emcal, PHENIX, hbd, zdc, run07, heavy ion

Awarding University: SUNYSB
Communities: PHENIX Collaboration

Keywords

Managing the keywords

Keywords

Listed on this page are *recommended* keywords used to tag materials placed on this site. Keywords are used in certain automated features e.g. aggregation and linking of materials pertaining to a particular topic, so their consistent use is encouraged. The keywords are case-sensitive.

The *same set of keywords* is used for materials uploaded to the Zenodo system under the umbrella of the [PHENIX Community on Zenodo](#). The list is compiled to provide better consistency of the subsequent queries. Zenodo is using a complex query mechanism which includes but is not limited to “elastic search” on the submission text (where applicable) so the effect of capitalization on queries is not always straightforward. In the following we adopt lowercase convention for all keywords. Note that a keyword on Zenodo can actually be a combination of words and white space (i.e. phrases). Multiple such combinations are allowed in a single query.

In the tables below, the keywords are grouped in categories. Each entry in the left column acts as a query link to *Zenodo*, for that specific keyword. Pages containing results of queries will open in a new tab/window.

General

Keyword	Description
decadal plan	Two documents describing the proposed PHENIX research program for two different time periods
phenix	Pioneering High Energy Nuclear Interaction Experiment (PHENIX)
rhic	Relativistic Heavy Ion Collider (RHIC)

Conferences

Keyword	Description
dnp19	DNP 2019
qm2019	Quark Matter 2019
sfjhf20	Santa Fe Jets and Heavy Flavor Workshop
wwnd2020	The 36th Winter Workshop on Nuclear Dynamics

Physics

Keyword	Description
au+au	Gold-on-gold collisions
Correlations	Various types of correlations

Working Zenodo links

HEPData: an example of a PHENIX entry

HEPData Search HEPData Search

About Submission Help Sign in

Browse all Adare, A. et al. Last updated on 2014-08-11 17:26 Accessed 799 times 99 Cite JSON

Hide Publication Information

Inclusive double-helicity asymmetries in neutral-pion and eta-meson production in $\bar{p} + \bar{p}$ collisions at $\sqrt{s} = 200$ GeV

The PHENIX collaboration

Adare, A., Aidala, C., Ajitanand, N.N., Akiba, Y., Akimoto, R., Al-Ta'ani, H., Alexander, J., Andrews, K.R., Angerami, A., Aoki, K.

Phys.Rev. D90 (2014) 012007, 2014.  <https://doi.org/10.17182/hepdata.64716>

Journal INSPIRE HEPData Resources

Abstract (download)  BNL-RHIC. Results are presented from data recorded in 2009 by the PHENIX experiment at the Relativistic Heavy Ion Collider (RHIC). We report the double-longitudinal spin asymmetry, A_{LL} , for π^0 and η production in $\sqrt{s} = 200$ GeV polarized \bar{p} - \bar{p} collisions. Comparison of the π^0 results with different theory expectations based on fits of other polarized data showed a preference for small positive values of gluon polarization, ΔG , in the proton in the probed Bjorken x , x_B , range. The effect of adding the new π^0 data to a recent global analysis of polarized scattering data is given.

Download All

Filter 9 data tables

Table 1
Data from Table 4
10.17182/hepdata.64716.v1/t1
 π^0 ASYM(LL) measurements from 2005.

Table 2
Data from Table 4
10.17182/hepdata.64716.v1/t2
 π^0 ASYM(LL) measurements from 2006.

Table 3
Data from Table 4
10.17182/hepdata.64716.v1/t3
 π^0 ASYM(LL) measurements from 2009.

Table 4
Data from Table 5
10.17182/hepdata.64716.v1/t4
 η ASYM(LL) measurements from 2005.

Table 5
Data from Table 5
10.17182/hepdata.64716.v1/t5
 η ASYM(LL) measurements from 2006.

Table 2 [10.17182/hepdata.64716.v1/t2](#)
Data from Table 4

π^0 ASYM(LL) measurements from 2006.

cmenergies 200.0

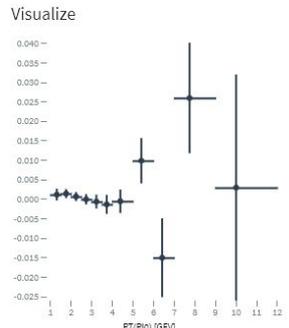
observables ASYM

phrases Inclusive Asymmetry Measurement Proton-Proton Scattering

reactions P P -> π^0 X

RE	P P -> π^0 < GAMMA GAMMA > X
SQRT(S)	200.0 GeV
PT(π^0) [GEV]	ASYM(LL)
1.3 (bin: 1.0 - 1.5)	0.0012 ± 0.0013 stat ± 0.00075 sys,rel,lumi. $\pm 8.3\%$ sys.pol.
1.5 - 2.0	0.00146 ± 0.00082 stat ± 0.00075 sys,rel,lumi. $\pm 8.3\%$ sys.pol.
2.23 (bin: 2.0 - 2.5)	0.0007 ± 0.00084 stat ± 0.00075 sys,rel,lumi. $\pm 8.3\%$ sys.pol.
2.72 (bin: 2.5 - 3.0)	0.0 ± 0.0011 stat ± 0.00075 sys,rel,lumi. $\pm 8.3\%$ sys.pol.
3.22 (bin: 3.0 - 3.5)	-0.0006 ± 0.0016 stat ± 0.00075 sys,rel,lumi. $\pm 8.3\%$ sys.pol.
3.72 (bin: 3.5 - 4.0)	-0.0013 ± 0.0023 stat ± 0.00075 sys,rel,lumi. $\pm 8.3\%$ sys.pol.

Visualize



Open Data

- “Canonical” data levels
 - Level 1 data provides more information on published results
 - Level 2 data includes simplified data formats for outreach and analysis training
 - Level 3 data comprises reconstructed collision data and simulated data together with analysis-level experiment-specific software
 - Level 4 data covers basic raw data
- Where does PHENIX stand?
 - Level 1 is covered by the current HEPData activity and auxiliary info committed to Zenodo
 - Level 2: work underway to create Ntuples illustrating analysis techniques (*next slide*)
 - Level 2: the Open Data Portal (*under discussion*)
 - Levels 3 and 4 are not practical for public access for much of the same reasons as exist in other experiments (e.g. access to calibrations and access to sites)

The Open Data effort: annotated Ntuples (G. David)

system ($x > 0$ is the west arm, negative z is South).

The *MBntup.root* file is produced from minimum bias data (no lower limit on single cluster p_T in *gnt* or pair p_T in *ggntuple*), whereas in *ERTntup.root* the threshold for single cluster p_T in *gnt* is 5 GeV, and the threshold for pair p_T in *ggntuple* is also 5 GeV. Note that here we restrict only the pair p_T , the energy of the individual clusters can be (and often is) significantly lower.

Variable name	Description
cent	Event centrality
vtxZ	z -vertex of the event
pt	Transverse momentum of the cluster
costheta	Polar angle of the cluster
phi	Azimuthal angle of the cluster
sec	EMCal sector of the cluster
ecore	“Core” energy of the cluster (γ -candidate)
ecent	Energy in the central tower of the cluster (γ -candidate)
tof	Time-of-flight in the central tower of the cluster (γ -candidate)
prob	Probability that the cluster is a photon (based on χ^2)
disp	Dispersion of the cluster (γ -candidate)
chisq	χ^2 from expected photon shape of the cluster (γ -candidate)
twrhit	Number of towers in the cluster (γ -candidate)
stoch	Combined variable to describe “photonness” of the cluster (γ -candidate)
x	x -position of impact point on the EMCal surface
y	y -position of impact point on the EMCal surface
z	z -position of impact point on the EMCal surface

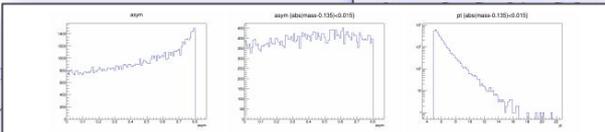


FIG. 2. ERT data, plots from the pair ntuple. Left: energy asymmetry distribution for all pairs.

```
ggntuple->Draw("mass", "mass<1.0");
ggntuple->Draw("mass", "mass<0.4&&pt>8.0");
ggntuple->Draw("mass">>htemp1, "mass<0.4");
ggntuple->Draw("mass">>htemp2, "mass<0.4&&chisq1<2.0&&chisq2<2.0");
htemp1->SetLineColor(1);
```

(2);

see Fig. 1.

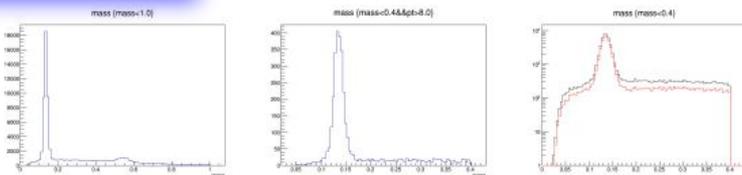


FIG. 1. ERT data, plots from the pair ntuple. Left: Invariant mass in the 0-1 GeV region. You can see a strong π^0 and a well-recognizable η peak. Middle: π^0 peak for pairs with p_T greater than 8 GeV/c. You can clearly see the combinatorial background outside the peak, which should

Ntuples: O(100MB) each
 Hosting options: Zenodo, Open Data

Status

- Infrastructure for the PHENIX DAP has been put in place
- Analysis capture is limited to annotation and documentation of a few use cases due to insufficient effort available and difficulty of retrofitting new structured solutions
- Within the scope described above, PHENIX relies on state-of-the-art services, platforms and tools
- The new durable DAP website has been commissioned and content is being added
- Ongoing creation and uploads of curated and tagged materials to Zenodo
- Systematized submissions to HEPData

Lessons learned

- DAP: plan and start early (should be a part of someone's job description)
- Avoid building in-house information systems wherever possible
 - State-of-the-art services such as Zenodo, HEPData, Inspire etc which cover a vast majority of the experiment's needs when properly used
- Create websites for the long haul: **the content needs to be curated**
 - Information gets fragmented, redundant and/or obsolete (many examples with Wikis etc)
 - Having misleading information is worse than none at all
- A flexible and easy-to-use Conditions DB is key to reducing the amount of ad-hoc data produced and utilized in analyses, hence facilitates DAP
- Prioritize analyses (even with good effort, we won't be able to preserve all)
- Once DAP-friendly frameworks and practices are established in the project the cost of DAP will likely be incremental due to economies of scale and available expertise within the research teams

Potential for Collaboration and Cooperation

- Community workshops and exchange of experience (Rivet/HEPdata workshops planned by the PHENIX community)
 - <https://indico.bnl.gov/event/8840/overview>
- PHENIX will be happy to share the DAP website development experience
 - Already taking place with shared components of the EICUG Software Group website
- Design of common analysis workflow templates that would be conducive to the analyses capture (cf. REANA)
- Conditions DB
- In future: developing shared expertise in advanced tools like REANA etc