# Importance of Data and Analysis Preservation

- The two decades of operation of RHIC have produced PetaBytes of invaluable data.

- Preserving and ensuring accessibility of these data and analyses is crucial to advancing scientific knowledge and fostering collaborative innovation.

- Data and Analysis Preservation:
  - Ensures long-term scientific impact, through potential future discoveries, verification, and reproducibility of results.
  - Supports the education and training of future scientists.
  - Enables reanalysis with new insights, ensuring data remains accessible as new theoretical advancements and discoveries emerge.
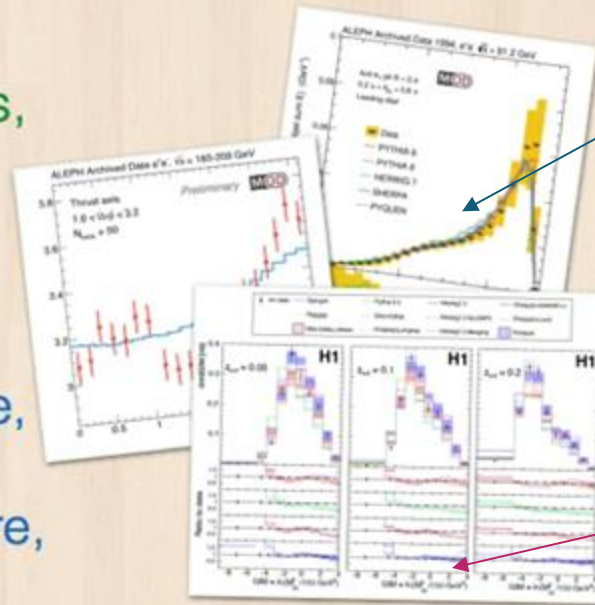
This effort aligns with DOE policy.
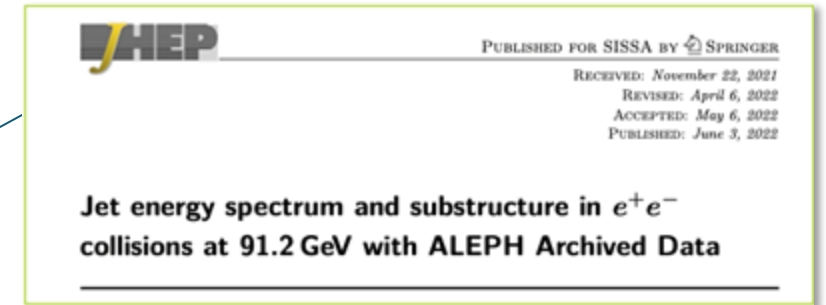
**Brookhaven**
National Laboratory

# Old Data is a gold mine!

The data from the ALEPH and H1 experiments were reanalyzed using algorithms developed long after the data collection had concluded.

ALEPH experiment at LEP ended data taking in 2000
Analysis published in 2022



ALEPH Archived Data 1994, e⁺e⁻, √s = 91.2 GeV

**Jet energy spectrum and substructure in $e^+e^-$ collisions at 91.2 GeV with ALEPH Archived Data**

H1 experiment at DESY ended data taking in 2007
Analysis published in 2024

**Measurement of groomed event shape observables in deep-inelastic electron-proton scattering at HERA**
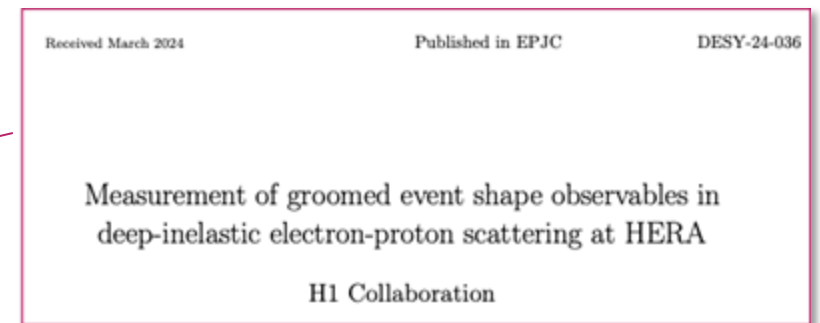
H1 Collaboration

## Summary

- Archived data is a gold mine with many exciting opportunities

  - QCD studies, new ideas, new algorithms, …

- Food for thought for ongoing experiments:
  preservation of knowledge,
  multiple MC samples,
  ability to rerun key software,
  low-level information, …
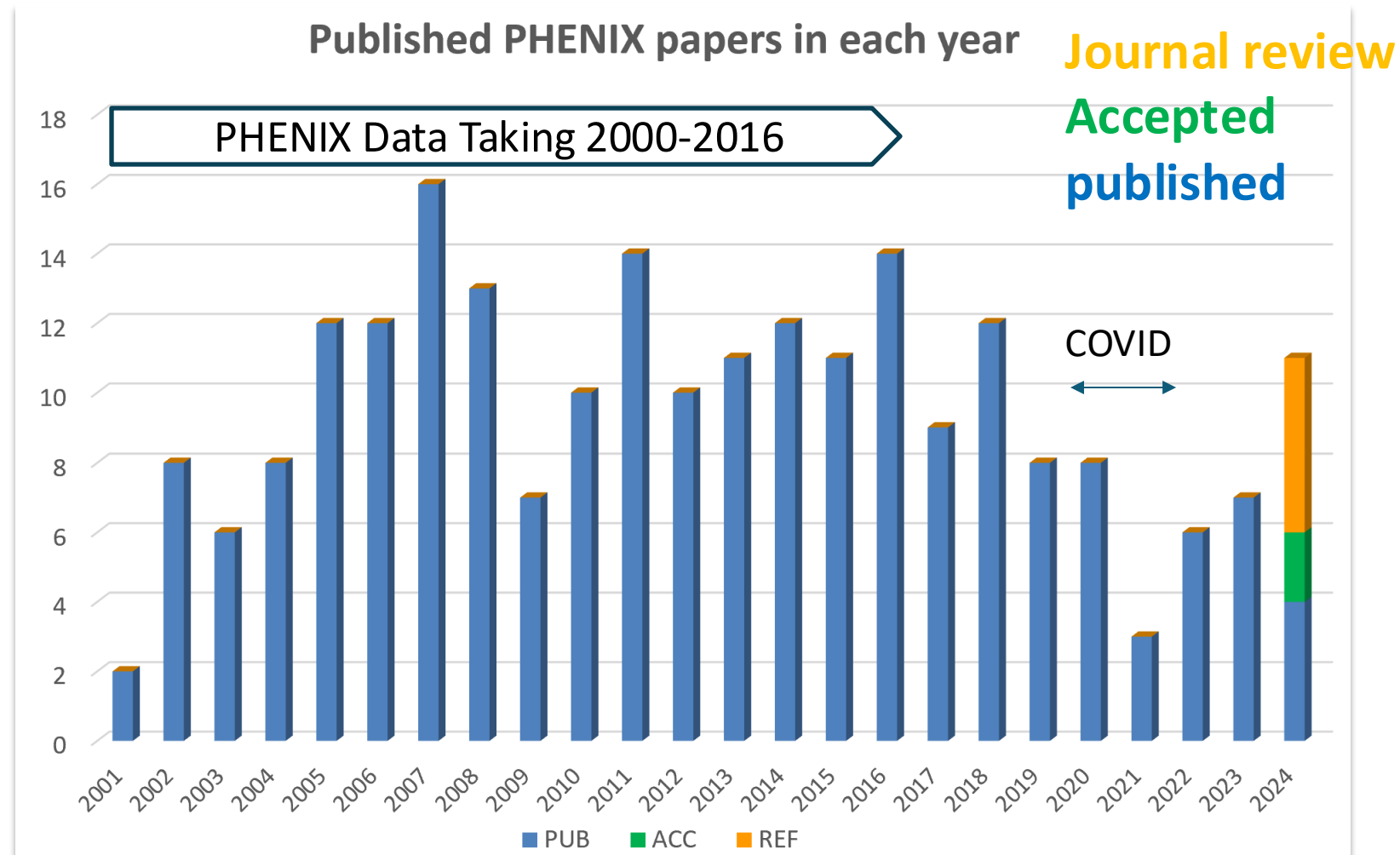
14

WORKSHOP ON
SOFTWARE INFRASTRUCTURE
FOR ADVANCED
NUCLEAR PHYSICS
COMPUTING

# Experiment publications extend well after the end of the data-taking: PHENIX example

**Published PHENIX papers in each year**

**Journal review**
**Accepted**
**published**

PHENIX Data Taking 2000-2016

COVID

Legend: PUB (blue), ACC (green), REF (orange)

# Experiment publications extend well after the end of the data-taking

33% of publications

15% of publications

EIC momentum is expected to drive increased re-analysis of HERA data.



Several years after the end of data taking, an experiment enters the Data and Analysis Preservation mode.

*Publications in Preservation mode may account for 10-20% of the total experiment publication set.*

# Publishing is not Preserving Analyses



- **PHENIX and STAR** publish their figures and tables in HEPData, along with differential distributions of data and predictions. This approach provides the main inputs for statistical interpretation.

- **Preserving analyses** requires more than just publishing data.
  - It involves saving raw data, simulated data, analysis code, conditions and comprehensive documentation.
  - This approach ensures that methods and context are fully understood, allowing others to reliably replicate the analyses.

Brookhaven National Laboratory

# Challenges for Reproduction of Analysis

Some of the challenges that must be overcome for *each* analysis to be preserved:

- Navigating complex analysis methods and workflows,

- Understanding the computing environment, detector and accelerator conditions,

- Preparing Data, analysis artifacts, and documentation for use by non-collaborators.

This requires early and ongoing effort from experts (analyzers, computing, and software) to capture knowledge before it fades.

The efforts in analysis preservation by the PHENIX (and STAR) collaborations have provided valuable insights for the RHIC community.

Brookhaven
National Laboratory

# Different levels of preservation

- Level 3 allows for re-analysis only.

- Most experiments (Zeus, H1, Babar, etc.) have adopted level 4, which allows for simulation and data reconstruction.

| | Preservation Model | Use Case | |
|---|---|---|---|
| 1 | Provide additional documentation | Publication related info search | **Documentation** |
| 2 | Preserve the data in a simplified format | Outreach, simple analyses | **Outreach, reanalysis** |
| 3 | Preserve the analysis level software and data format | Full scientific analysis, based on the existing reconstruction | **Technical Preservation Projects** |
| 4 | Preserve the reconstruction and simulation software as well as the basic level data | Retain the full potential of the experimental data | |

*Data Preservation Levels defined by the Data Preservation in HEP (DPHEP) Collaboration*

Aiming at level 4 for RHIC Data and Analysis Preservation

**Brookhaven**
National Laboratory

# Comprehensive RHIC Data and Analysis Preservation plan

- **Goal**: Ensure the integrity, accessibility, and longevity of RHIC data and analysis post-RHIC operation.

- **Comprehensive strategy**:
  - Engagement with experiments while they are actively analyzing data, integrating DAP into the publication process.
  - Implement standards and best practices as documented, for example, by the Data Preservation in High Energy Physics (DPHEP) collaboration
  - Future users of RHIC data may not be current collaborators, so the DAP strategy should be designed with this in mind.

- **Components of a DAP plan**:
  - Software Preservation: Maintain operability of older software using containerization.
  - Knowledge Management: Capture and preserve knowledge through detailed documentation.
  - FAIR Principles: Ensure data is Findable, Accessible, Interoperable, and Reusable.
  - Dedicated Portal: Provide comprehensive access to RHIC data and analyses

- **Continuous Improvement**:
  - DAP will evolve through collaboration with the RHIC user community.
  - Each experiment is different

- Note: bit and functional preservation (storage and infrastructure) is not covered (yet) by this project

A BNL Program Development for supporting RHIC Data and Analysis Preservation was approved for FY25.

Brookhaven National Laboratory

# Bit and functional preservation

- **Bit Preservation** ensures that digital files remain unaltered over time, preventing data loss and corruption.

- **Functional Preservation** ensures that data remains accessible and usable over time.
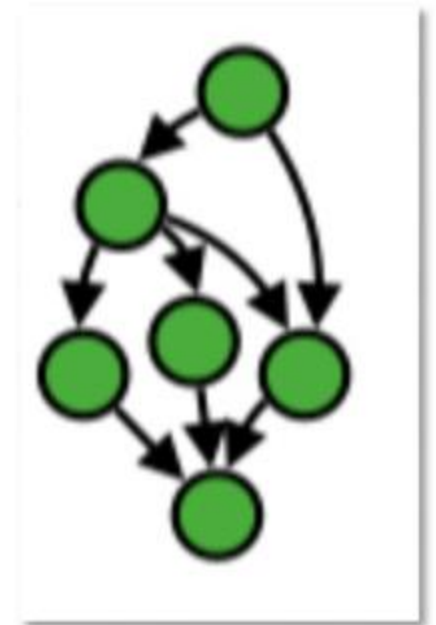
Both types of preservation are essential for a successful Data Analysis and Preservation strategy.

This means that data archives (tape systems, media, and software) and some computing infrastructure must be maintained and remain accessible.

**Brookhaven**
National Laboratory

# Software Preservation

- **Challenge**:
  - Maintaining older software across diverse operating systems and hardware architectures.
- **Solution**:
  - Use containerization technologies to encapsulate analysis software along with its dependencies.
  - Preserve workflow orchestration of analysis steps to ensure continued operability on future platforms.

While container technology ensures the portability of software environments, it remains confined to the original environment and technologies. Software must be modernized and adapted for successful long-term preservation and analysis.

**Brookhaven**
National Laboratory

# Knowledge Management

- **Importance**:
  - Preserve expertise through detailed accurate documentation before it fades over time.

- **Implementation**:
  - Document data formats, analysis pipelines, and scientific rationales behind specific analyses.
  - Record and document experimental conditions
  - Control and document software versions.
  - Use version control systems like Git to track changes and ensure traceability.

This is crucial and difficult

# Implementing FAIR Principles

- **Framework**:
  - Ensure data is Findable, Accessible, Interoperable, and Reusable (FAIR).
- **Implementation**:
  - Use standardized metadata and persistent identifiers (PIDs) to enhance discoverability.
  - Include comprehensive metadata for each dataset, describing its content, context, and provenance.
  - Assign unique IDs and metadata to each analysis for advanced search capabilities.

Brookhaven
National Laboratory

# A dedicated DAP Web Portal

Organize information for each analysis by providing links to  Data, Code, Workflow, Plots, Figures, and Documentation

- Single Discovery Entry:
  - Access various data and information about RHIC analyses.
  - Perform simple and complex searches for data.
- User-Facing:
  - Provides access to data analysis tools and pertinent information.
  - Links to repositories with detailed datasets and analyses.
  - Accessible through Federated Identity mechanisms initially for RHIC collaborators only
- Incorporating an AI assistant (Chatbot) leveraging Large Language Models:
  - To manage large volumes of information specific to RHIC.
  - Provide code snippets, facilitate data access, and answer complex questions.
  - Developed through collaborative learning with experts from the collaborations.
  - STAR has already a prototype.

**Brookhaven**
National Laboratory

# Summary

- Data and Analysis Preservation (DAP) at RHIC is vital for scientific integrity and accessibility.

- A comprehensive DAP at RHIC was proposed, establishing solid foundation for a preservation program which will require long-term support.

- Implementing successful RHIC Data and Analysis Preservation will require significant effort and commitment from the entire RHIC community, including collaborations and the Lab.

- A task force with all RHIC experiments would boost collaboration and efficiency.

**Brookhaven**
National Laboratory