

# USQCD Long-term Data Management Strategy and Plan

James N. Simone, Andreas S. Kronfeld  
Fermi National Accelerator Laboratory

Robert G. Edwards  
Thomas Jefferson National Accelerator Facility

January 12, 2021

## Abstract

This document is the data management plan for long-term storage of data generated by members of the USQCD Collaboration, with operations carried out by the LQCD and NPPLC computing sites.

## 1 Introduction

A data management strategy describes a coherent set of principles and objectives for managing data over all phases of its lifecycle: pre-planning, data production, distribution, active usage, and long-term preservation or eventual removal. A data management plan (DMP) describes the coordinated actions needed to accomplish the strategic goals. A DMP is an important tool to provide a roadmap for maintaining the provenance of scientific results. They are also increasingly required in proposals to funding agencies and as supplementary material to scientific publications.

Three distinct but related entities are involved in the DMP described in this document. USQCD [1] is a federation of U.S.-based researchers working on lattice QCD and other lattice field theories. It is administered by an Executive Committee (EC), which charges a Scientific Program Committee (SPC) with developing a proposal-driven program of research. USQCD is the scientific steward of two streams of funding from offices within the U.S. DOE Office of Science: the Nuclear and Particle Physics Lattice-QCD Computing (NPPLC) initiative, which is funded by the DOE Office of Nuclear Physics (DOE NP), and the LQCD extension III research program (LQCD for short), which is funded by the DOE Office of High Energy Physics (DOE HEP). LQCD procures computer time on institutional clusters at Brookhaven [2] and Fermilab [3]; NPPLC procures and operates dedicated cluster at Jefferson Lab [4]. LQCD, NPPLC, and the USQCD EC and SPC collaborate to maximize the scientific output of the USQCD membership, given the resources available.

This document focuses on a plan to manage computer data generated by USQCD members on the clusters or elsewhere, with a focus on long-term repositories for data of broad and lasting scientific value. These repositories will make use of the tape storage facilities located at the three sites. These long-term data will be managed separately from but consistently with the short-term storage that these facilities provide to scientific projects on a yearly basis.

This document is organized as follows. Section 2 highlights key aspects of the USQCD tape data management strategy, while Sec. 3 describes the role of tape storage. In Sec. 4, we describe important elements and objectives for the USQCD data management strategy. A template DMP for USQCD projects is provided in an Appendix.

## 2 Objectives of the USQCD DMP

This DMP for tape resources is designed to coordinate several objectives of USQCD collaboration:

- 2.1. Ensembles of lattice-QCD gauge configurations are a valuable community resource that have been instrumental in enabling many diverse physics projects. Gauge ensembles generated using USQCD resources are made available to the USQCD community and others worldwide, as discussed in the USQCD charter [5].
- 2.2. QCD gauge ensembles generated by USQCD members will remain an important scientific resource well beyond the two current funded programs, which are funded by the DOE HEP and DOE NP. Since the outset, USQCD has been committed to long-term data preservation program for QCD gauge ensembles, albeit in an *ad hoc* way.
- 2.3. The record of scientific achievements enabled through the availability of USQCD gauge ensembles is invaluable to present and future research. USQCD is working in partnership with inSPIRE-HEP developers and the DOE Office of Scientific and Technical Information (OSTI) to develop processes to link a publication to the gauge ensembles used in the publication via unique persistent Digital Object Identifiers (DOIs).
- 2.4. With modern lattice-QCD workflows, these procedures for sharing must be extended to further computationally expensive large-size data objects, such as Dirac-operator eigenvectors. Eigenvector sets, like gauge-field ensembles, are stored on magnetic tape and are now subject to the same sharing policies. DMPs will provide a useful means to document data compression methods, file formats, and metadata encoding for eigenvector sets.
- 2.5. USQCD is committed to ensuring research is reproducible and to providing open access to published data. Thus, this DMP provides expectations for maintaining long-term the provenance and integrity of published data and results.
- 2.6. A coherent data management strategy and policy will assist in effective allocation and utilization of distributed tape storage resources project-wide in order to maximize scientific productivity and impact.
- 2.7. The DMP described in this document is a recommended element of the DMP of the individual scientific collaborations within the USQCD federation. Such DMPs are a required part of grant proposals to the DOE, the NSF, and at nearly all computing facilities worldwide. Often, proposals omitting a DMP will be rejected without further consideration. This plan is available as a framework and a citable reference in developing an individual DMP for a scientific collaboration.
- 2.8. This DMP (and future revisions) will keep USQCD in compliance with U.S. government policies regarding access to data resulting from government-sponsored research. For example, in February 2013, the White House Office of Science and Technology Policy (OSTP) issued a memorandum directing Federal agencies that provide significant research funding to develop a plan to expand public access to research. Among other requirements, the plans must: “Ensure that all extramural researchers receiving Federal grants and contracts for scientific research and intramural researchers develop data management plans, as appropriate, describing how they will provide for long-term preservation of, and access to, scientific data in digital formats resulting from federally funded research, or explaining why long-term preservation and access cannot be justified.”

### 3 Roles of magnetic tape data storage

Magnetic tape storage is used for both short- and long-term data distinguished by their life cycles and data usage patterns. Hence, the desired features of the two storage categories will differ, and USQCD, NPPLC, and LQCD have distinct policies to allocate and manage each category of tape storage.

Short-term storage is appropriate for data that are used during the execution of a batch production campaign. Among these data are prestaged inputs to batch runs, data that are generated and shared among intermediate stages of the calculations, and data outputs from each batch job. These data

products reside on disk storage during the batch runs and most of such data are deleted from disk when no longer needed by the campaign. Short-term tape is used in conjunction with disk as a second storage tier that is historically less expensive than high-performance disk. Tape extends the total available capacity for short-term data storage but at the cost of higher access latencies. Workflows orchestrate the movement of data among storage tiers: moving data to tape to make room on disk for data needed immediately in jobs, then prestaging it back to disk as needed later in the campaign. The added complexity of managing a tiered storage system is especially cost effective when the intermediate data are costly to recompute. The short-term tape facility should have high bandwidth connectivity to the disk storage system. Short-term tape data are written and read frequently and often have several reads for each write. Short-term data are also frequently deleted. Short-term storage resources are awarded to each project by the USQCD SPC as part of the annual awards process. To this end, proposals must understand the cost of storage vs. recomputing data. Once a project has ended, short-term storage is guaranteed to remain in place for six months, at which time site managers may ask that it be cleared, particularly disk.

Long-term storage is needed for data of lasting high scientific value that are costly to produce, widely shared among scientific collaborations, and that enable a wide variety of computations in nuclear and high-energy physics. Ensembles of gauge configurations are the most familiar example. Long-term storage is also appropriate for data such as Dirac eigenvectors that are costly to compute and can be shared among many projects, thereby reducing the net cost of quark propagator solves.

Long-term storage is also appropriate for preserving a representative subset of the provenance (*e.g.*, code version, random number information, sample output logs, etc.) and all essential data used in scientific publications. As part of a commitment to open and reproducible research, such records should be retained for at least ten years following publication.

The predicted useful lifetime of long-term data such as gauge ensembles and published data records exceeds the planning and budgeted calendar of the current DOE SC program for lattice QCD computing. USQCD has collaborated with LQCD and NPPLC to develop a plan and budget specifically for the curation and preservation of USQCD data. The usage pattern for long-term data is that the data are written once and read infrequently, compared with short-term data. Reading a given dataset can be expected to be done in a sequence accessing a contiguous subset of data. Deletions are rare and may result once it is more cost effective recompute the data rather than to retain it on tape. Long-term tape storage will require access to temporary disk storage to cache data sets upon transfer to or from tape.

It can happen that a project is allocated computing at one site that requires data residing in long-term storage at another site. In these cases, the proposal's data management plan will be revised so that the SPC officially furnishes such projects with an allocation with the needed amount of staging disk space and a negligible amount of computing.

The USQCD SPC will allocate long-term tape storage by soliciting storage proposals from projects twice a year, with one solicitation coinciding with the yearly Call for Proposals for computing resources. These proposals must justify why USQCD resources (as opposed to, say, those at a home institution) are appropriate.

## 4 Elements of the tape data management strategy

The elements of the USQCD data-management strategy are organized into four sets of objectives and responsibilities: 1) those of the USQCD collaboration as a whole, 2) those of the LQCD research program and NPPLC initiative, 3) those of the cluster facilities operated (currently) at Brookhaven, Fermilab, and Jefferson Lab, and 4) those of the projects and researchers making use of USQCD resources.

### 4.1. USQCD Collaboration

- (4.1.1) USQCD acknowledges the need to preserve high-value data of lasting scientific interest. These data have useful lifetimes that may well extend beyond the horizon of current DOE funding for USQCD computing.

- (4.1.2) The yearly Call for Proposals (CfP) for scientific projects will announce the availability of long-term tape storage for valuable data. Scientific projects responding to this CfP will provide a justification designating their data for long-term storage, and they will provide statistics describing any existing data and projections for the rates they expect to produce long-term data in coming years.
- (4.1.3) The USQCD SPC will determine whether the stored long-term data are of sufficiently broad interest such that the cost need not be considered part of a project allocation but borne by the LQCD or NPPLC on behalf of the whole USQCD community.
- (4.1.4) Short-term storage will be granted to scientific projects as part of the SPC computing resource allocation process. Such tape storage is allocated for the duration of the project, plus the (guaranteed) six months following a project's end.
- (4.1.5) The USQCD EC, together with LQCD and NPPLC, will appoint and coordinate a committee of USQCD EC and SPC representatives, facilities staff, and USQCD users mandated to document use cases and requirements for a long-term storage facility.
- (4.1.6) USQCD will encourage and assist science projects to adopt best standards and practices for data organization, data storage formats, and metadata markup. Poor choices for any of these often impacts both the productivity of the project and the performance of facility data resources.
- (4.1.7) USQCD will provide a template DMP to encourage and assist research groups in creating a project specific DMP. Such plans are already a prerequisite in every proposal to the DOE and NSF. A thoughtful DMP will help researchers to understand the costs and their responsibilities with regard to the data they produce. A DMP is also useful to better understand and plan a balance compute vs storage requirements and costs. Such a template is provided in the Appendix.
- (4.1.8) The USQCD EC will devise a strategy for maintaining long-term data of broad utility beyond the lifetime of the hardware projects.

#### 4.2. LQCD research program and NPPLC initiative

- (4.2.1) LQCD and NPPLC will budget funding for disk and tape storage for both short- and long-term data, based on usage patterns of USQCD physics projects.
- (4.2.2) LQCD and NPPLC, together with the USQCD EC, will appoint and coordinate a committee of USQCD EC and SPC representatives, facilities staff, and USQCD users mandated to document use cases and requirements for a long-term storage facility.
- (4.2.3) LQCD will negotiate a memorandum of understanding (MOU) with each of its sites for hosting a long-term storage facility; NPPLC and JLab will maintain the understanding that data relevant to the JLab mission are kept in long-term storage.

#### 4.3. BNL, Fermilab, and JLab sites

- (4.3.1) The sites will deploy tape storage resources as spelled out in MOUs with LQCD and NPPLC, for USQCD's long- and short-term disk and tape needs.
- (4.3.2) The sites are responsible for maintaining integrity and access to data over the life cycle each data set, as informed by each USQCD project's DMP. This information will guide facilities in the management of project data.
- (4.3.3) Facilities hosting USQCD data are responsible for implementing the USQCD DM strategy. The details may vary among the sites. Each site has posted the policy details for users in the sites online documentation: BNL [6], Fermilab [7], and JLab [8].
- (4.3.4) LQCD and NPPLC will provision user accounts, provide documentation, and provide technical support for long-term storage, just as they do for projects' short-term storage and compute resources.

- (4.3.5) When tapes for long-term storage are scheduled to be migrated to newer media, the site managers will alert the USQCD EC & SPC—and also LQCD or NPPLC, whichever is pertinent. The SPC will then decide, in consultation with other stakeholders, which data need not be migrated.

#### 4.4. Researchers and USQCD science projects

- (4.4.1) Each scientific project must include its DMP in its proposal to the SPC and provide this (possibly revised) DMP to the site managers upon receiving an allocation.
- (4.4.2) The DMP must identify a data manager, who is responsible for making data management decisions on behalf of the project. The data manager and the project PI will be the primary points of contact regarding data. The data manager should be someone expected to remain a part of the USQCD Collaboration well after a project becomes inactive, although the duties can be reassigned to another project member. The project PI is the default data manager.
- (4.4.3) Projects seeking long-term data preservation should begin negotiating with suitable data hosting sites as soon as possible. Projects may also need to apply to the funding agencies or their home institutions for additional funds for data preservation.
- (4.4.4) Individual researchers and projects are ultimately responsible for the integrity of their data. They must develop and use backup procedures to prevent the catastrophic loss of critical data. Such data should be replicated among geographically separated locations.
- (4.4.5) Each project is expected to abide by the data management policies of the facility hosting their data. Under extraordinary circumstances a project may be unable to meet expectations. Projects should immediately inform the facilities and the SPC of the situation and provide justifications for seeking an exception.
- (4.4.6) Magnetic tape storage provides bulk data storage that is currently less expensive per GigaByte than either magnetic or solid state disk storage. Even so, total magnetic tape costs are a significant recurring cost for the project. Projects may request short-term tape storage for the active life of the project and its extensions. Projects are expected to remove short-term data, particularly disk, within six months after a project becomes inactive.
- (4.4.7) Projects may request long-term magnetic tape storage at the USQCD sites for critical results or community data, such as QCD gauge configurations. Typically, storage space is granted for the lifetime of the tape media, and a single migration to new tape technologies or media. Additional migrations may incur significant costs to the researchers and the USQCD collaboration. Projects must request long-term storage as part of their proposal and data management plan. The amount of storage and hosting institution will be negotiated by the SPC, EC, and the USQCD facilities.
- (4.4.8) Each project must have a plan for what is to happen to their data after it becomes inactive or ends. At the earliest stages of planning, a project should begin to develop a workable, cost-effective plan for data preservation. Data that are not of broad community value could, for example, be copied back to the researchers' home institutions.
- (4.4.9) Not all data created by a project need be shared or preserved. The costs and benefits for doing so has to be weighed during data management planning.
- (4.4.10) Projects are encouraged to preserve the provenance and integrity of data and results that appear in publications, as required, for example, by funding agencies or journals. Typically, it is impractical to store all lattice data from intermediate stages leading to publications. Researchers should preserve sufficient metadata, or a subset of log files, describing data as well as software versions at intermediate stages in order that it would be possible to statistically reproduce equivalent data. Data and analysis programs leading up to the final stages of a published analysis should be preserved and be made publicly available as supplementary materials. It is recommended that a machine readable format be used for published values,

errors, and correlations among data points. Researchers are encouraged to publish programs from the final stages of their analysis or a notebook of such programs as part of a publication's supplementary materials.

## References

- [1] USQCD. *USQCD Collaboration governance*. URL: <https://www.usqcd.org/collaboration.html> (visited on 11/01/2020).
- [2] USQCD and BNL. *BNL Computing Information for USQCD Users*. URL: <https://www.usqcd.org/bnl/> (visited on 11/01/2020).
- [3] USQCD and Fermilab. *Welcome to the Fermilab Lattice Gauge Theory Computational Facility*. URL: <https://www.usqcd.org/fnal/v2/> (visited on 11/01/2020).
- [4] USQCD and JLab. *Jefferson Lab LQCD*. URL: [https://lqcd.jlab.org/lqcd/lqcd\\_index.html#/](https://lqcd.jlab.org/lqcd/lqcd_index.html#/) (visited on 11/01/2020).
- [5] USQCD. *USQCD Collaboration Charter*. Aug. 2018. URL: <https://www.usqcd.org/documents/charter.pdf> (visited on 11/10/2020).
- [6] BNL RACF. *HEP Digital Data Management Policy*. June 8, 2019. URL: [https://www.racf.bnl.gov/about/policy/data-management/at\\_download/file](https://www.racf.bnl.gov/about/policy/data-management/at_download/file) (visited on 11/03/2020).
- [7] Fermilab SCD. *Fermilab LQCD Facility Data Management Guidelines and Policies*. Feb. 4, 2019. URL: <https://www.usqcd.org/fnal/v2/datamanagement.html> (visited on 02/04/2019).
- [8] Chip Watson. *Jefferson Lab Data Management Plan*. July 1, 2014. URL: <https://scicomp.jlab.org/DataManagementPlan.pdf> (visited on 11/03/2020).
- [9] University of California Digital Library Curation Center. *DMPTool: Build your Data Management Plan*. Feb. 2019. URL: <https://dmptool.org/> (visited on 11/03/2020).
- [10] University of California, Berkeley, Libraries. *UCB Research Data Management*. Feb. 2019. URL: <https://researchdata.berkeley.edu/services/data-management-plans> (visited on 11/03/2020).
- [11] MIT Libraries. *MIT Libraries data management resources*. Feb. 2019. URL: <https://libraries.mit.edu/data-management/> (visited on 11/03/2020).
- [12] U.S. DOE. *DOE Policy for Digital Research Data Management: Suggested Elements for a Data Management Plan*. Mar. 5, 2016. URL: <https://www.energy.gov/datamanagement/doe-policy-digital-research-data-management-suggested-elements-data-management-plan> (visited on 11/03/2020).
- [13] U.S. National Science Foundation. *NSF Data Management Plan Template*. July 1, 2014. URL: <https://www.mbl.edu/osp/files/2014/07/OSP-NSF-Data-Management-Plan-Template.pdf> (visited on 11/03/2020).

## Appendix: DMP template for projects

Online resources, including many university libraries, provide tools and guides to preparing a DMP [9, 10, 11]. The DOE Office of Science provides a guide to the suggested elements of a DMP [12]. Likewise, the NSF provides a DMP template as well [13]. The DOE and NSF templates are alike in organization and this template follows a similar outline. The on line templates [9, 10, 11] have the advantage that they can be used to create living documents and are therefore perhaps preferable to a static text document.

### App.1 DMP for USQCD project *project-name*

- *Briefly describe the project. You may cite the USQCD proposal for this project.*

### App.2 Plan revision history

**January 12, 2021** Plans updated to include configuration generation in 2020.

**June 1, 2020** Initial revision.

### App.3 Contact Information

- *Provide the name and affiliation of the projects PI or main point of contact. Also provide the name and affiliation of the project's data managers. The data managers and the POC are the project members having primary responsibility for making decisions regarding the data covered in this plan.*

Role	Name	Affiliation
PI		
Data manager		

### App.4 Data types

- *DMPs for science projects applying for long-term storage should tabulate the amount of data expected to be produced in each of the next three years going forward. Include in the table the year the data will be produced, a concise unique identifier, data types (e.g., gauge field, eigenvectors, etc.), expected data set size produced in that year, number of files, and file sizes. For existing datasets, add a single line to the table listing the years it was produced, e.g., 2016–2019 and the other statistics. Please add any additional comments following the table.*

Year(s)	Dataset ID	Data type	Total size [TB]	File size [GB]	File count
-2020		DWF configurations	200	1–90	~ 50k
-2020		DWF eigenvectors	800		~1500k
-2020		DWF $g - 2$ HLbL	400		
-2020		DWF propagators	200		
-2020		DWF distillation	?		
2020–2023		DWF configurations	100		
2020–2023		DWF eigenvectors	600		50k
2020–2023		DWF $g - 2$ HLbL, propagators	300		
		⋮			
		⋮			

## **App.5 Data and metadata standards**

- *For the data types tabulated in the previous section, briefly describe the file format and the codebases that produced the data. Describe software interfaces that are able to read the files. Describe the nature, formats and location of the metadata describing the data.*

## **App.6 Policies for access, sharing, and protection**

- *Gauge ensembles generated using USQCD resources are made available to the USQCD community. Describe your project's plan for providing access and sharing within USQCD. Describe your collaboration's policies governing allowable usage of these data by other USQCD projects. Also describe any mechanisms in place to protect these data from unauthorized access.*

## **App.7 Policies for re-use and redistribution**

- *Describe your collaboration's limitations, policies and plans on sharing these data outside of USQCD.*