

# Memorandum of Understanding regarding the PetaLibrary storage service

*University of Colorado Boulder, Research Computing*

## 1. Introduction

This is the memorandum of understanding (MOU) for the Research Computing (RC) PetaLibrary Storage Service. It includes both a Service Level Agreement (services provided to the customer) and an End User Agreement (terms agreed to by the customer).

## 2. Description of service

The PetaLibrary is a National Science Foundation-subsidized service for the storage and archival of research data. It is available to any US-based researcher affiliated with the University of Colorado Boulder.

### 2.1 Service classes

Two main classes of storage are available:

#### Active

- for data that is frequently written or read
- directly accessible (read/write) from computational resources within the CU RC environment
- accessible from outside the CU RC environment only via specific data transfer protocols (scp, sftp, gridftp/GlobusOnline) through designated gateway nodes
- mounted as /work/<projectname>
- not designed for highly I/O-intensive or massively parallel usage
- “snapshots” (on-disk incremental backups) can be enabled

#### Archive

- for data that is used infrequently
- accessible only via specific data transfer protocols (scp, sftp, gridftp/GlobusOnline) through designated gateway nodes
- mounted as /archive/<projectname>

### 2.2 Service Details

We offer five specific storage services. Each customer project will be associated with one of these services. Storage space must be pre-purchased, with a minimum purchase of 2 TB for 1 year.

#### Active - single copy on disk

- Appropriate for non-critical data
- No protection against accidental file deletion unless snapshots are enabled
- No protection against catastrophic failure of the primary storage system

### **Active with replication - as above but with files replicated at least once a day to a disk system in a different building**

- Appropriate for data that would be difficult to recreate
- Some protection against accidental file deletion
- Good protection against catastrophic failure of the primary storage system
- Good protection against loss of the primary data center
- Replication happens at least once per day

### **Active plus one copy on tape - a single copy of each file is saved to tape**

- Appropriate for data that would be difficult to recreate
- Good protection against catastrophic failure of the primary storage system
- Does not allow restoring all files to a “point in time”; rather, only the most recent version of each file is saved on tape
- Recovery of data from tape can only be done by RC technical staff
- Copy to tape is made once per day
- Does not work well with many (hundreds of thousands) of small (under 128 KB) files

### **Archive - single copy in hierarchical storage management (HSM)**

- Appropriate for non-critical data that is used infrequently
- Recently-used files are stored on disk
- Files that have not been accessed recently are automatically migrated to tape
- Files that have been migrated to tape are still visible in the HSM file system and can be accessed via the usual interface, but there will be a delay of at least several minutes while they are read from tape
- A maximum of 30% of the project’s pre-purchased data space can be resident on disk at any time
- The maximum size of an individual file is 2 TB. The average file size must be greater than 100 MB (enforced via a limit on the total number of files: 10,000 files per TB of space purchased.) It may be necessary to tar or zip your data to fit these requirements
- No protection against accidental file deletion
- No protection against mechanical failure of a tape cartridge

### **Archive plus one copy on tape**

- Appropriate for data that would be difficult to recreate but that is used infrequently
- Good protection against mechanical failure of a tape cartridge
- Does not allow restoring all files to a “point in time”; rather, only the most recent version of each file is saved on tape
- Recovery of data from tape can only be done by RC technical staff

## **3. END-USER EXPECTATIONS**

### **3.1 Usage**

See Appendix A for the End User Agreement, which must be signed and returned to RC by the PI before a project can be set up.

### **3.2 User/Group Administration**

Every research group using the PetaLibrary Service may provide RC with up to two official contacts: the Principal Investigator (PI) and an alternate “Point Of Contact” (POC) person. The PI must be a CU-Boulder faculty member. These two people are the only ones authorized to make changes to the set of users that comprise a group, to add or remove users on the access list, or to request other changes on behalf of the project.

Any change request that may impact a group will have to come through the PI or POC. Without such approval, RC and OIT will not act on group-level requests. Change and service requests must be made through the RC ticketing system.

The PI and POC should verify their list of users on a yearly basis. Note that while IdentiKey passwords may be terminated after a person leaves the university, his or her files will be retained in the PetaLibrary Service. This is done to ensure that shared files do not accidentally become unusable by a research group when one user account becomes inactive. It does place an additional burden on the PI and POC to purge old data periodically.

The PetaLibrary Service relies on CU-Boulder's IdentiKey and LDAP infrastructure. All users must have a valid IdentiKey and their LDAP record must accurately reflect their school/department affiliation. RC will work with OIT and your department administrators to clean up the LDAP records, if needed.

The IdentiKey requirement implies that all users are CU-Boulder faculty, students, staff, or affiliates. Note that any CU-Boulder faculty member can sponsor a non-CU person as an affiliate, and then add them to the service. RC cannot act as the sponsor.

### **3.4 Costs/Fees**

The cost to the PetaLibrary end users is heavily subsidized. End users pay actual media costs for disk/tape space plus a modest overhead (for datacenter rent and similar costs.) Please see Appendix B for the costs of the different services.

### **3.5 Ownership of Media**

The disks and tapes that your data resides on belong to RC and may not be removed from the PetaLibrary.

## 4. SERVICE EXPECTATIONS

### 4.1 Service levels

We define two levels of service:

**Sensitive**, where service is provided from Monday through Friday 8am to 5pm, including after-hours support on those weekdays with “best effort”. There is no guarantee of weekend or holiday support, though RC will try to respond to incidents impacting the operation of the service with “best effort.”

**Tolerant**, where service is provided during business hours only, Monday through Friday, 8am to 5pm.

RC will operate the PetaLibrary Service as a hybrid between the “Sensitive” and “Tolerant” service levels. The backend storage systems and gateway nodes will be Sensitive. Customer-initiated support requests (e.g., group membership or permission changes) will be Tolerant.

While every reasonable and good faith effort will be made to ensure the reliability and availability of the PetaLibrary and of the files stored on it, access to data in the PetaLibrary may be affected by circumstances outside of the control of Research Computing.

### 4.2 Accessing storage

- Active storage is directly accessible for read/write via NFS from computational resources within the CU RC environment, including login and compute nodes. It is also directly accessible from login and compute nodes in the Blanca (“condo”) cluster.
- Archive storage is directly accessible for read/write within the CU RC environment from the login and data-transfer nodes only. Access via the login nodes is provided to facilitate movement of data from other RC storage; it is emphatically not to be used for frequent reads/writes of archived data.
- To access Active or Archive storage from outside of RC, it is necessary to use scp, sftp, or rsync through login.rc.colorado.edu, or Globus through the “CU Boulder Research Computing” endpoint. Globus will almost always provide much higher throughput than scp, sftp, or rsync, and is the preferred protocol for any but the smallest transfers. All of these methods require authentication via RC’s One Time Password (OTP) system. If a customer requires automated transfers that would be awkward with OTP authentication, access through a data-transfer node using scp (with ssh keys) can be provided to hosts within the CU-Boulder campus network.

All Active or Archive projects are also eligible to use the Globus Sharing service, in which Globus Online can be used to provide read and/or write access to your data space to anyone with a Globus account. A CU Identikey or RC OTP is not necessary for Sharing access. RC will sponsor one free Globus Plus account (for Sharing administration) per PetaLibrary project.

### **4.3 Duration of Service**

The hardware infrastructure supporting the PetaLibrary Service is funded initially through a National Science Foundation grant. This support ends May 30, 2018. Operation of the PetaLibrary is expected to continue beyond that date thanks to ongoing University. In the event that RC ceases to provide the PetaLibrary or any comparable resource, RC will give at least 60 days advance notice. It will be the responsibility of the customer to transfer their data to other storage resources within that time window.

### **4.4 Amount of storage available**

The existing PetaLibrary infrastructure can support about 950 TB of Active storage and about 1100 TB of Archive storage. Customers needing more than 75 TB of Active storage or 100 TB of Archive storage should consult with RC about availability.

### **4.5 Reporting**

A report detailing the project's current storage usage will be emailed within 10 days after the end of each month to the PI and the POC.

At least once per year, RC will provide a list of users to all PIs and POCs showing the names and IdentiKeys for every user who has access to the project's data.

### **4.6 Maintenance**

Planned maintenance of the RC infrastructure, including the PetaLibrary hardware, will take place on the first Wednesday each month. RC will broadcast the announcement on the rc-announce email list (which all users are strongly encouraged to join) and the RC web site.

### **4.7 Change management**

RC will try to announce to all PetaLibrary customers any major changes to the system a minimum of 30 days in advance. Exceptions may be critical security updates and bug fixes that improve the stability of the system significantly.

### **4.8 Performance**

RC will provide suggested throughput performance targets for each service type; however, since the PetaLibrary is a shared infrastructure actual performance may vary depending on workload from other customers. Please note that PetaLibrary is not an ideal storage target for active databases.

### **4.9 Refunds**

We are not able to process refunds or pro-rate the PetaLibrary fee for any time lost due to repairs or maintenance events (planned or otherwise), nor for any datacenter-related down time.

### **4.10 User training**

Training on file transfer to and from the PetaLibrary is part of our RC workshop series. RC can schedule a 30 minute consultation session to help customer groups learn to use the PetaLibrary most efficiently.

## 5. HELP/SUPPORT REQUESTS

PetaLibrary Service users may make help and support requests through the RC ticket system. While RC will make every effort to respond to and resolve support requests, those support requests which require domain-specific knowledge or expertise may not be able to be handled by RC alone. In these cases, the support request may be forwarded to the POC for assistance.

## APPENDIX A:End User Agreement

University of Colorado Boulder PetaLibrary use agreement

I, the undersigned, agree that my use of and access to the digital storage facility known as the University of Colorado Boulder “PetaLibrary” shall be in accordance with all of the following stipulations:

\_\_\_\_\_ I will not store on the PetaLibrary any files that are US government classified, or subject to the US federal Health Insurance Portability and Accountability Act (HIPAA), the US federal Family Educational Rights and Privacy Act (FERPA), the International Traffic in Arms Regulations (ITAR), or other restrictions described at <http://www.colorado.edu/avcit/policy>

- IT Policies: <http://www.colorado.edu/avcit/policy>
  - University of Colorado System Administrative Policies
    - <https://www.cusys.edu/policies/>
  - CU-Boulder Policies and Guidelines
    - <http://www.colorado.edu/about/policies>
  - OIT Campus-wide IT Policies
    - <http://www.colorado.edu/avcit/campus-policies>
- Health Insurance Portability and Accountability Act (HIPAA):
  - <http://www.hhs.gov/ocr/privacy/hipaa/understanding/index.html>
- Family Educational Rights and Privacy Act (FERPA):
  - [http://registrar.colorado.edu/regulations/ferpa\\_guide.html](http://registrar.colorado.edu/regulations/ferpa_guide.html)
- International Traffic in Arms Regulations (ITAR):
  - [http://pmdt.c.state.gov/regulations\\_laws/itar\\_official.html](http://pmdt.c.state.gov/regulations_laws/itar_official.html)
- CU-Boulder Private and Restricted Data Security Requirements:
  - <http://www.colorado.edu/avcit/sites/default/files/attached-files/data-definitions.pdf>

\_\_\_\_\_ I will not store on, or distribute from, the PetaLibrary any files that are copyrighted, except with the express permission of the copyright owner.

\_\_\_\_\_ If any of the files that I store on the PetaLibrary are subject to one or more agreements with any Institutional Review Board (IRB), including but not limited to the IRB of the University, then I will take full responsibility for ensuring full compliance with such agreement(s). \_\_\_\_\_

\_\_\_\_\_ If I am collaborating with colleagues who are at institutions outside of the United States of America (that is, outside of both US states and US territories), then I will take full responsibility for ensuring that those colleagues do not access the PetaLibrary themselves, but rather I and/or other members of my team who are at US institutions will access the PetaLibrary on behalf of the entire team.

\_\_\_\_\_ I understand that, if and when I cease to be employed by and/or a student at an institution in the United States of America, then access to my files on the PetaLibrary will be available only to those of my collaborators who are employed by and/or students at US institutions.

\_\_\_\_\_ I will take full responsibility for ensuring that my use of the PetaLibrary is in full compliance with the most current version of the University's Acceptable Use Policy, currently accessible at

<http://www.colorado.edu/avcit/sites/default/files/attached-files/resources.pdf>\_\_\_\_\_ If I am one of the Principal/Co-Principal investigators of a team, then I will take full responsibility for ensuring that any other members of the team are likewise in full compliance. \_\_\_\_\_ I will provide the RC with a yearly report of research activities supported by the PetaLibrary and I will acknowledge the PetaLibrary in publications as following: "This work was supported in part by the Scalable Petascale Storage Infrastructure funded by NSF under grant OCI-1126839"

\_\_\_\_\_ I understand that the ability of the University to provide the PetaLibrary is contingent on continued University funding and cooperation; that the University provides the PetaLibrary on an as-is basis, and while every reasonable and good faith effort will be made to ensure the reliability and availability of the PetaLibrary and of the files stored on it, the Research Computing at CU Boulder makes no guarantees with respect to its reliability or continued availability.

\_\_\_\_\_ In the event that Research Computing ceases providing the PetaLibrary or any comparable resource, then I will take full responsibility for transferring any and all relevant files to other storage resources, and in a timely manner.

\_\_\_\_\_ I will take full responsibility for ensuring that I keep abreast of and comply with changes to any of the relevant laws, policies and circumstances described above.

(Please initial to demonstrate agreement with each of the above statements.)

Signature: \_\_\_\_\_ Date: \_\_\_\_\_

Print Name: \_\_\_\_\_

Institution: \_\_\_\_\_

Contact Information:

Email: \_\_\_\_\_ Phone: \_\_\_\_\_

APPENDIX B:

Order Form Please fill out and email to [rc-help@colorado.edu](mailto:rc-help@colorado.edu)

Project name: \_\_\_\_\_

(should be specific to group or project; will be used in directory path; should not change over life of project; cannot include spaces. Good example: smith\_mri. Bad example: labdata.)

Service	Cost/TB/yr	TB (min 2)	Years (min 1) *	Total Cost
Active	\$65			
Active plus replication	\$130			

Active plus one copy on tape	\$80			
Archive (HSM)	\$35			
Archive plus one copy on tape	\$45			

Requested start date \_\_\_\_\_

Speedtype to charge \_\_\_\_\_

Speedtype owner \_\_\_\_\_ Project PI

\_\_\_\_\_  
 IdentiKey \_\_\_\_\_ Email \_\_\_\_\_ Dept

\_\_\_\_\_  
 Project POC \_\_\_\_\_ IdentiKey \_\_\_\_\_ Email

\_\_\_\_\_  
 Dept \_\_\_\_\_

Do you want snapshots enabled (Active options only)? Y / N

I have read and agree to the PetaLibrary Storage Service MoU.

\_\_\_\_\_  
 Signature of PI

Date

\* Cannot extend beyond May 30, 2023.