# Large File Ingest and Retrieval

- Ingesting Large Files via the REST API
- Direct Comparison of Different Transfer Methods
  - Comparison of Upload and Download Times for Different Transfer Methods
    - Copying Files Between Federated Filesystem and Repository Storage
- Range Retrieval

Large files can be uploaded via the REST API. Transfer times for uploading to the repository via the REST API are about the same as copying using NFS, and moderately faster than using SCP.

Java command-line options and System properties can be used

- -Xmx2048m maximum memory Java can use
- -Dfcrepo.home=/path/to/data set the directory for permanent data
- -Djava.io.tmpdir=/path/to/tmpdir set the directory for temp files. Data uploaded via the REST API is written to a temp file in this directory, so this directory should have enough free space for the largest files you will upload. (1) if running under tomcat, tomcat replaces your java.io.tmpdir with the value of CATALINA\_TMPDIR, so you should set that environment variable instead)

#### Ingesting Large Files via the REST API

Based on the tests below, we believe arbitrarily-large files can be ingested and downloaded via the REST API (tested up to 1TB). The only apparent limitations are disk space available to store the files, and a sufficiently large Java heap size.

- Platform: lib-devsandbox1.ucsd.edu (all data on NAS to handle large files)
- Repository Profile: Minimal
- Workflow Profile: Upload/Download Roundtrip

| File Size | Upload                  | Download               |
|-----------|-------------------------|------------------------|
| 256GB     | 3h51m34s (18.87MB/sec)  | 43m09s (101.25MB/sec)  |
| 512GB     | 7h49m43s (18.60MB/sec)  | 1h29m15s (97.90MB/sec) |
| 1TB       | 15h41m21s (18.57MB/sec) | 3h21m44s (86.63MB/sec) |

### **Direct Comparison of Different Transfer Methods**

Based on the tests below, we believe arbitrarily-large files can be uploaded and downloaded via the REST API, using either repository storage or a federated filesystem (tested up to 1TB). The only apparent limitations are disk space available to store the files, temp directory capacity, and a sufficiently large Java heap size.

- Platform: lib-devsandbox1.ucsd.edu (all data on NAS to handle large files)
- Repository Profile: Federation
- Workflow Profile: Repository/Federation/NFS/SCP Comparison

#### **Comparison of Upload and Download Times for Different Transfer Methods**

| Transfer Method       | File Size | Upload               | Download             |
|-----------------------|-----------|----------------------|----------------------|
| REST API (Federated)  | 1TB       | 732 min (84 GB/sec)  | 246 min (250 GB/sec) |
| REST API (Repository) | 1TB       | 339 min (181 GB/sec) | 250 min (246 GB/sec) |
| SCP                   | 1TB       | 383 min (160 GB/sec) |                      |
| NFS                   | 1TB       | 336 min (183 GB/sec) |                      |

## Copying Files Between Federated Filesystem and Repository Storage

| Source               | Destination          | File Size | Copy Time            |
|----------------------|----------------------|-----------|----------------------|
| Repository storage   | Federated filesystem | 1TB       | 402 min (153 GB/sec) |
| Federated filesystem | Repository storage   | 1TB       | 345 min (178 GB/sec) |

## Range Retrieval

Retrieving a byte range is supported and has been tested with 1TB files for both repository storage and federated filesystem. There is an integration test in the standard test suite for verifying that range retrieval works. By default, this test uses a small binary size to avoid slowing down the test suite, but the size is configurable so it is easy for a developer to test files as large as local disk space allows.