

# A Tape RSE for Extremely Large Data Collection Backups

CHEP 2024  
Krakow, Poland  
October 21 - 25, 2024

---

Andrew Hanushevsky, SLAC

<http://xrootd.org>

# The Instigator: Vera C. Rubin Observatory

- # Simonyi Survey Telescope and LSST Camera (8.4 m primary, 3.2 Gpix, 6 filter bands)
- # Legacy Survey of Space and Time (10 years, >18K deg<sup>2</sup> sky coverage, >825 visits per point)
- # Finishing Construction, Pre-Operations, and Commissioning; Operations to start in mid 2025



Render



Reality (Dec 2021)

# The LSST Backup Challenge

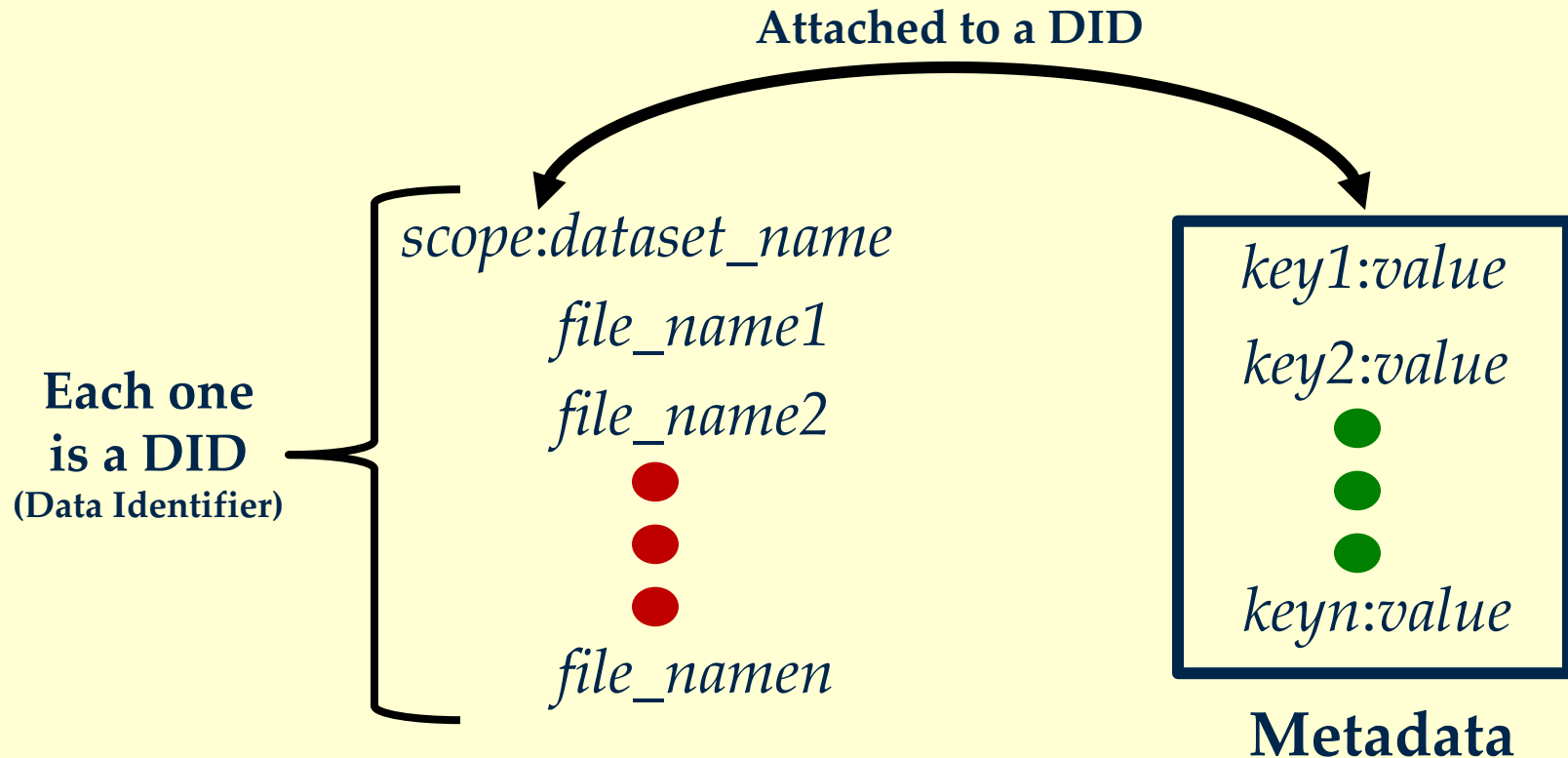
- # ~15 Rucio datasets per night (4K-8K/year)
  - Up to 20K files per dataset
    - Size range from a few KB to ~1 GB
    - Average size of dataset is approximately 100 GB
  - All of which need to be backed up
    - About 8 PB/year
- # Up to 12M additional dataset products/year
  - Each approximately 150GB
    - Undetermined number need backups

# Backup/Restore Requirements

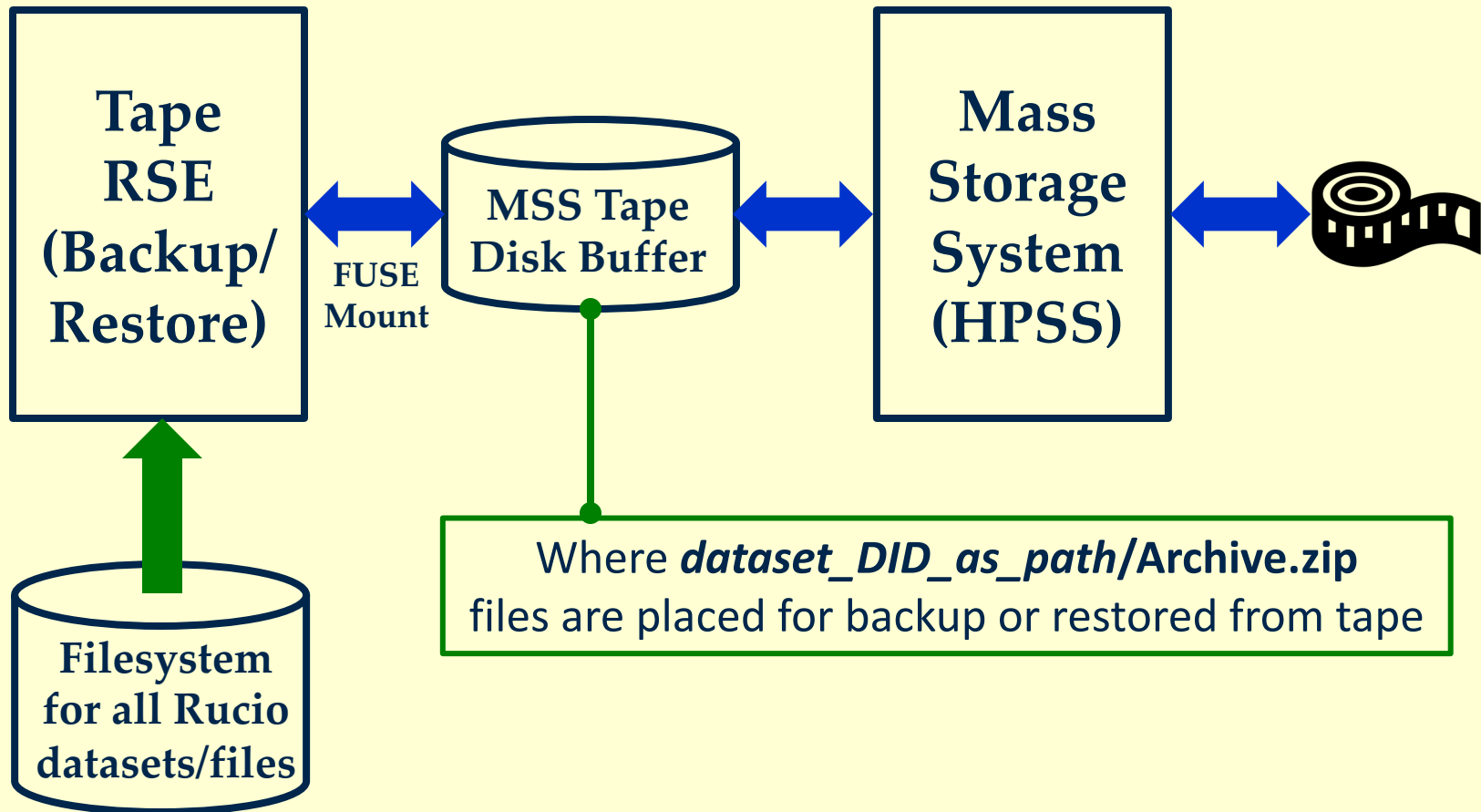
---

- # Simple way to create a backup
  - No complicated submission interfaces
- # 99+% assurance that a backup created
  - Ability to discover datasets not backed up
- # Ability to restore
  - Full datasets
  - Individual files from a dataset

# Rucio Conventions



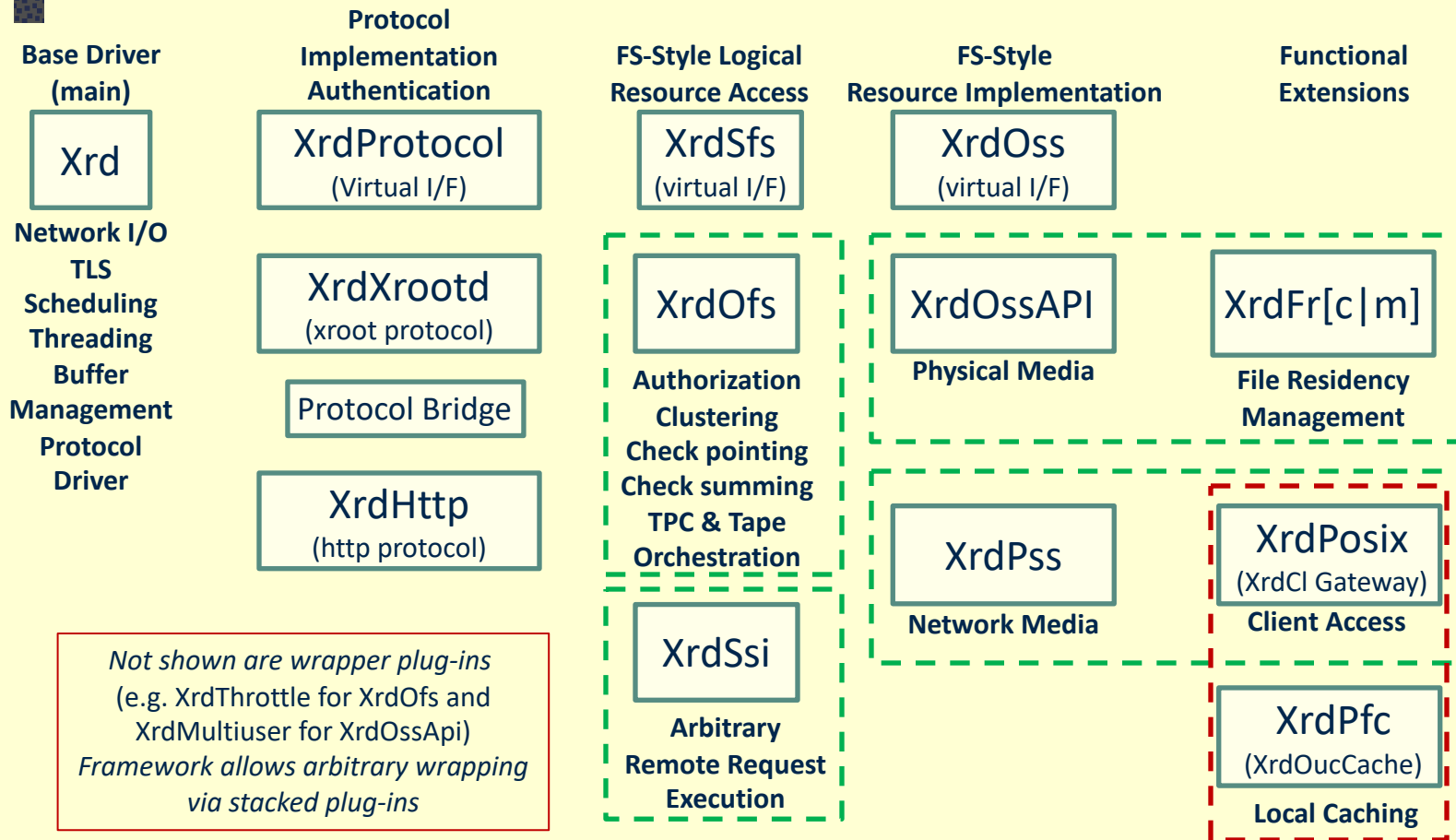
# Backup/Restore Environment



# The Backup Approach

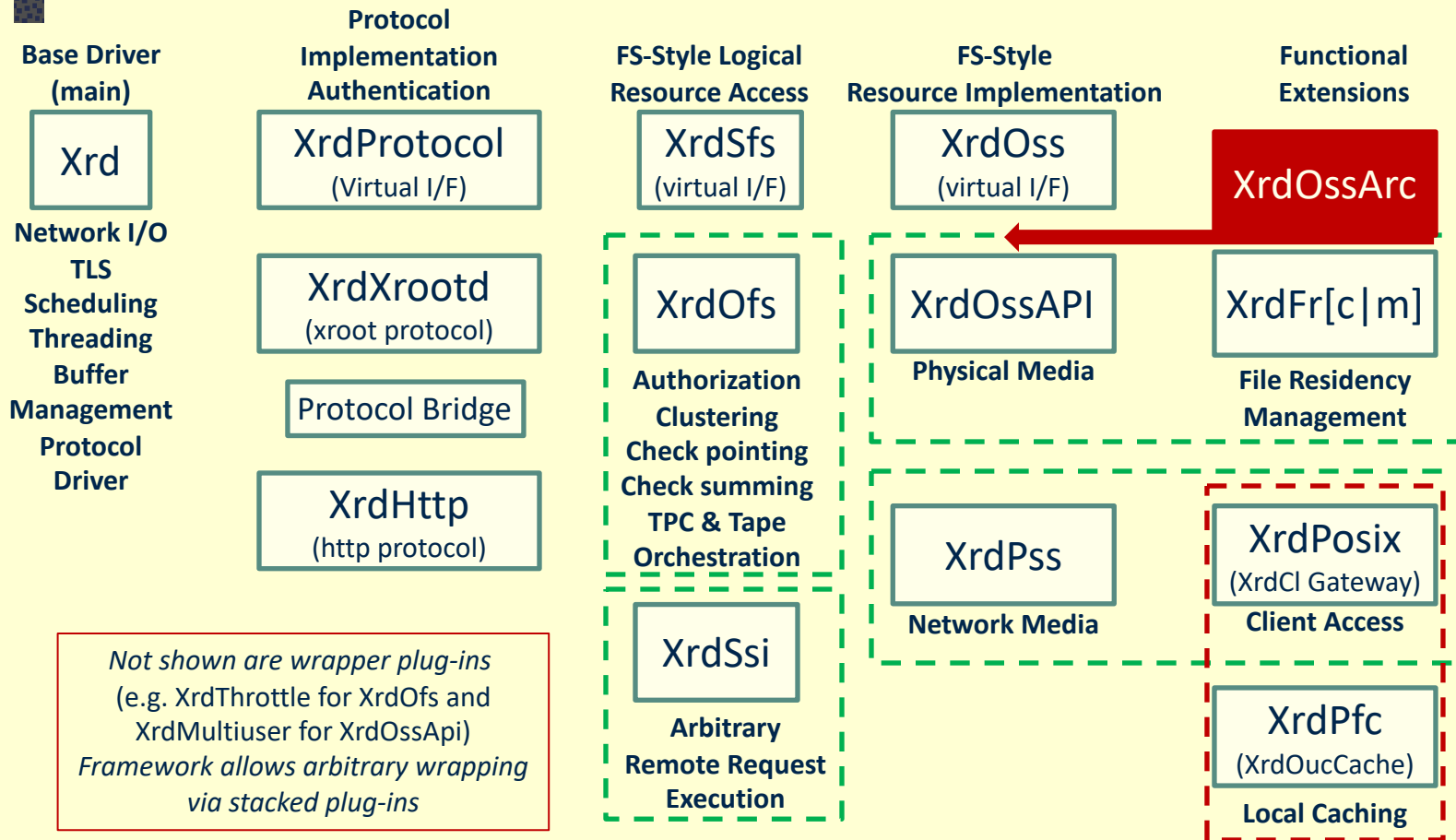
- # Dataset associated w/ backup metadata
  - **arcBackup** is the metadata key
    - *taperse:need* - *taperse* needs to backup ds
    - *taperse:done* - *taperse* completed the backup
- # The tape RSE is a specialized xroot server
  - Standard server with extra plug-in
    - libXrdOssArc.so (configurable)
      - Orchestrates the backup

# XRootD Plug-In Architecture





# The Stacked OssArc Plug-In



# Backup Orchestration

## # Coordinates 3 Python scripts

- XrdOssArc\_Archiver

- Creates zip file of dataset contents

- Stages archive to the Mass Store System for backup

- XrdOssArc\_BkpUtils

- Various Rucio dependent functions

- Used to setup backup and finalize the backup

- XrdOssArc\_MssCom

- Various MSS dependent functions

# The Backup Steps

## # Do Forever

- Create list of datasets to be backed up
- For each dataset using config scheduling
  - Setup logical dataset contents via symlinks
    - Default assumption all DID's are accessible to RSE
  - Invoke archiver to create and stage-out backup
    - Create *dataset\_DID/Archive.zip*
      - Zip member names are the file DIDs
    - Move archive to MSS stage-out file system buffer
    - Cleanup by removing logical dataset
  - Finalize the backup by updating DS metadata

# The Restore Approach I

## # Individual files

- Copy out file from the archive RSE

- `xrdcp xroot://rse//backup/dataset_DID~file_DID dest`

- Can also use your favorite HTTP copy program

- Currently, using a tilde (~) to separate *ds* from *fn*

- Only unassigned special character without UTF req

- Done by Rucio lfn to pfn plugin for the RSE

- Considering other alternatives like CGI “arc.fn=”

- Restore typically requires a stage-in from tape

- Does not require unpacking the zip file

- OssArc plug-in knows how to read/extract zip files

# The Restore Approach II

## # Full dataset

- Copy out the full archive zip file (fast restore)
  - `xrdcp xroot://rse//backup/dataset_DID~Archive.zip`
    - Can also use your favorite HTTP copy program
  - Unzip archive in-place implies special processing
- Potentially use Rucio replicate (by file restore)
  - Requires registering every file as replica
    - Needs to reopen the previously closed dataset or
      - Use a separate shadow dataset which is another issue
    - Doesn't solve the shared file restore problem
      - Note that backup includes all shared files

# Conclusions

- # The Tape RSE provides a needed service
  - Passive Rucio dataset backups
  - Active restores
    - Plan to make process more transparent
- # Can be used in many environments
  - Active components are Python scripts
    - Can be easily modified or replaced
      - Accommodate different storage environments

**That's All!**

---

**Questions?**