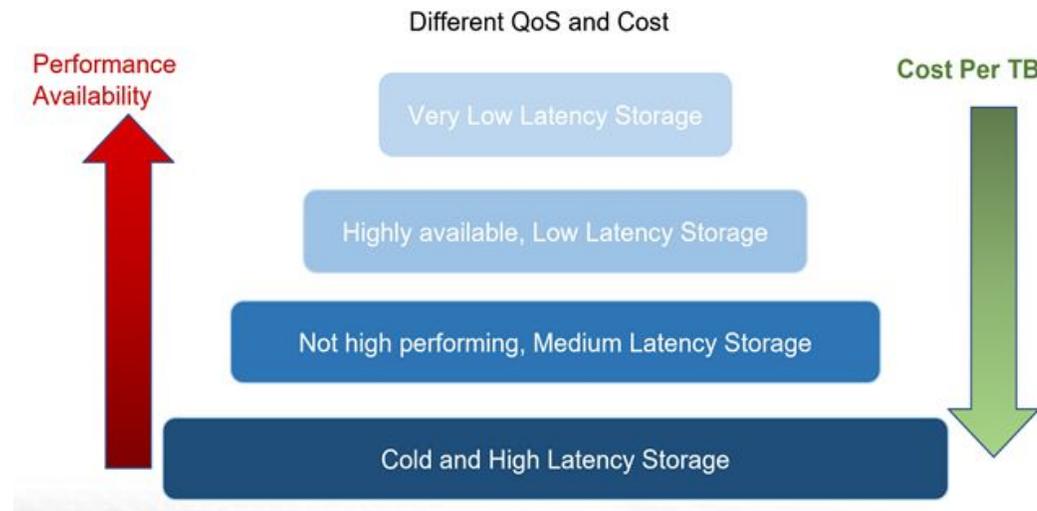# Rucio QoS Current State

- OSG All-Hands Meeting

- Matt Snyder, Hironori Ito, Martin Barisits, Doug Benjamin

18 Mar 2022

@BrookhavenLab

# Quality of Service

- QoS vital in the realization of Storage Cost reductions
- Tagging data with QoS allows Tier 1 sites to place data in the most cost effective system that satisfies user/client requirements
  - SSD, spinning disk or tape?  (ie move data to lower cost storage as needed)
  - How many replicas on disks? 1 copy or 5 copies (data deduplication)
- Enables sites to deploy more targeted and cost efficient storages systems
- Allowing users to specify their required performance might lead to more efficient use of the storage.
  - Storage cost is proportional to performance
  - Less cold data on the high performance disk.
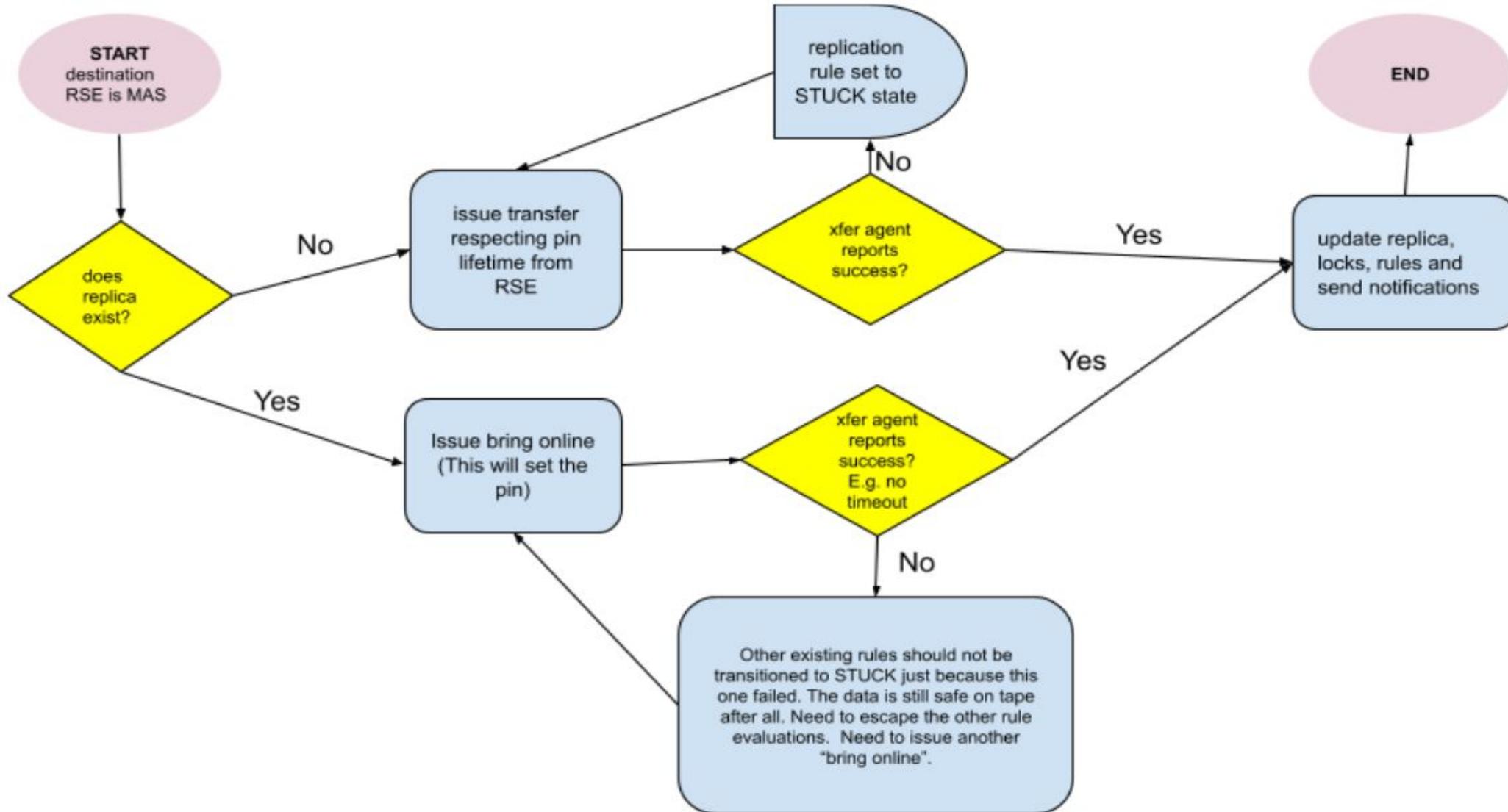


Slide courtesy Hironori Ito

# Phase I Current State

- We want to support something beyond binary data classification (DISK vs. TAPE)
  - allow clients to set their data demands
- Support MAS-type RSEs such as tape storages with very large disk buffers
  - in general, the quicker the data is moved to tape the more money is saved
- Real-world issue:
  - BNL has noticed in looking at its dCACHE disk-only pools that ATLAS has "cold" data on disk. We would be better served parking some of this data on tape temporarily in an automated manner. This way disused data is moved to least expensive storage (ie Tape) quicker.
  - There are issues with changing access latency to the data. How will PanDA behave when it "thinks" it's accessing low latency storage on HPSS?
- BNL contributing code to Rucio to achieve QoS policy

Brookhaven
National Laboratory

# Phase I Current State

- Our Phase I implementation of this so far:
  - Handle 2 attributes of the RSE
    - staging_required: True
    - maximum_pin_lifetime: {some value in seconds}
  - For rule_grouping if replica does not exist need to perform initial, full transfer
    - Rucio will submit a STAGING request even if data is already on storage
    - Panda cannot be notified immediately so rules are set to REPLICATING state waiting for the STAGING to be successful
  - We have functionally tested behavior of replicas when the transfer is required (no replica on MAS destination)
    - when replica already exists we need to be careful to respect

**Brookhaven**
National Laboratory
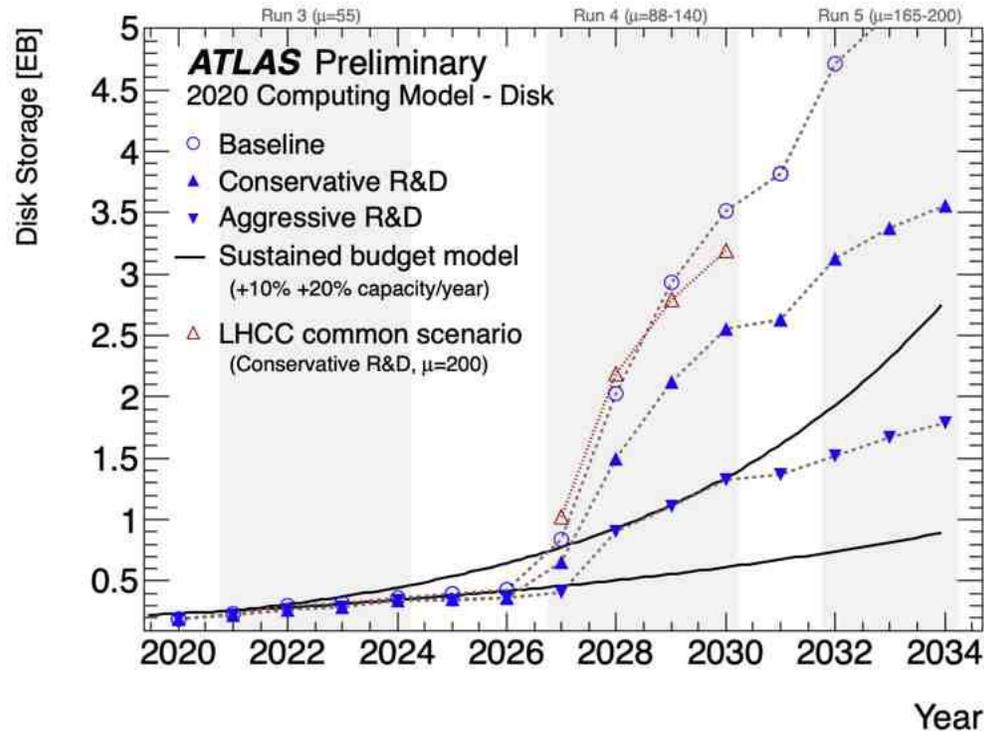
# Flowchart for Phase I Changes

# Phase II

- Add more QoS workflows/classes (e.g., SSD, SPINNING)
  - these could take advantage of changes being made in STAGINGAREAS
- Perhaps include support for location-based policies
  - RSE pin lifetimes could inherit from a site-wide strategy
- Some questions that have arose:
  - What happens when site deletes replicas from DISK?
  - How do we notify Rucio?
    - What would this message look like?
  - How do we query for storage replica state for a replication rule?

# Timeline (2022)

- Our aim is to finish the code changes with initial BNL-site functional tests by end of March
- We have to tread lightly and carefully as these changes would touch current workflows (such as lock and rule transitions)
  - Test cases include:
    - PanDA behavior as data slides down latency tiers
    - Pinning and expiring datasets on HPC scratch disk
    - Lock transitions: OK on success and STUCK on failure
    - Notifications properly sent on SUCCESS and FAILED transfers


- Once merged we can test with ATLAS workloads

# Rucio QoS - Summary

- This project's aim is to enable ATLAS collaboration to use storage smarter and more efficiently by moving unneeded files to cheaper storage.  Perhaps other experiments (DUNE) could benefit.

- It could also be used to speed up high-IOP workloads by moving data up to more performant storage.

# ...and thanks!

Contributors: Martin Barisits, Hironori Ito, Doug Benjamin

**Brookhaven** National Laboratory