# Data management and LHC

Eric Lancon

# The LHC



Large Hadron Collider

Lake of

CMS

LHCb

ALICE

ATLA

CERN

This is ATLAS

# Computing for LHC: WLCG



**Worldwide LHC Computing Grid**

**170** Data centres

**40** Countries

**800'000** Cores

**500** PB Disk

**750** PB Tape

**3** Tbps Network

*Tiered Structure*

| | |
|---|---|
| **Tier-0** | CERN |
| **Tier-1** | Large data centres |
| **Tier-2+3** | Universities and Laboratories |

*Heterogeneous Computing*

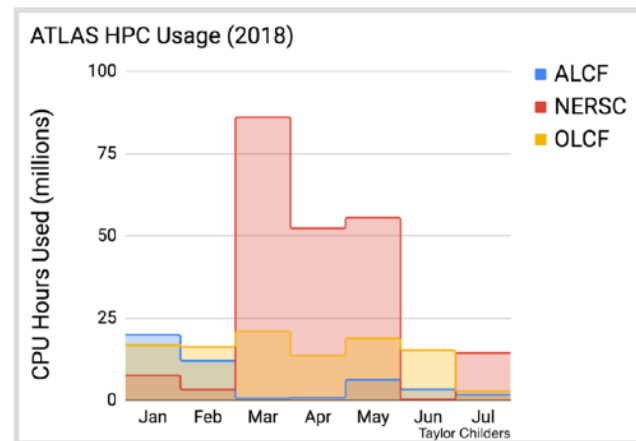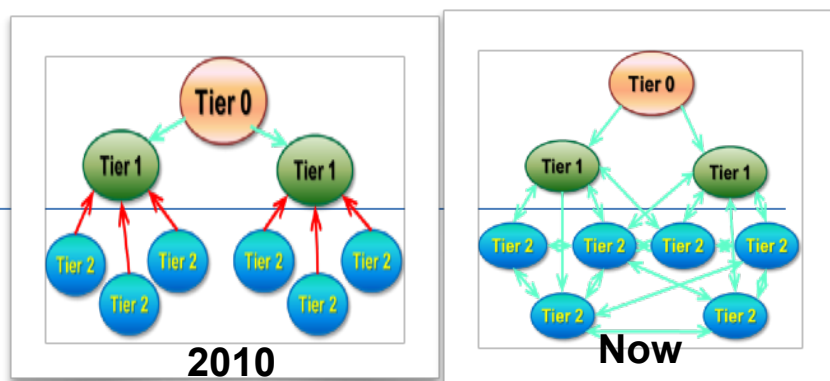Data centres (partly) supported by national Funding Agencies

Centres may host and support other projects

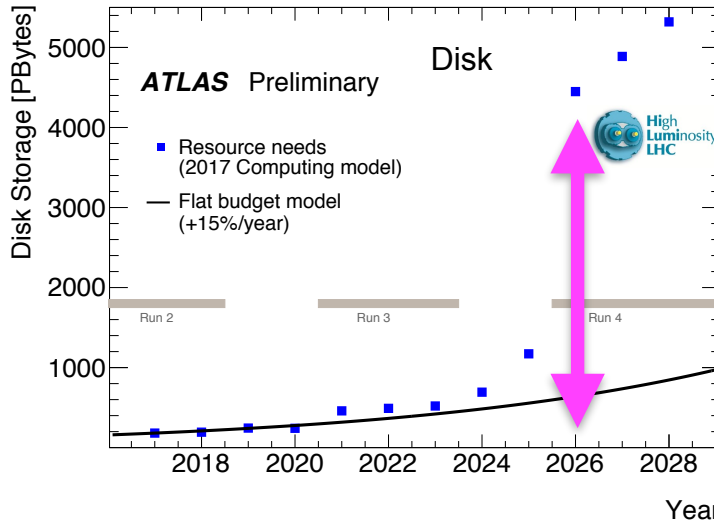Pledged storage and compute based on an MoU

# WLCG Tiers



2010 — Now

- From hierarchical to 'democratic' structure, data exchange between any tier, thanks to network capabilities!

- Sites are still providing resources with defined CPU/Storage ratio (No CPU only tier)

- BUT decoupling of storage and compute resources with increased usage of non-WLCG CPU provider like HPCs

BROOKHAVEN
NATIONAL LABORATORY | Scientific Data and Computing Center

# Upcoming (2026) : High Luminosity-LHC

ATLAS resource requirements



Required increased storage capacity ~10x today
(under current computing models)

# Time to re-think data distribution

# Elements of ATLAS experiment data movement

- Data Management Layer : RUCIO
- File Transfer Layer : FTS
- WAN: ESnet (+ GEANT, Signet,...)
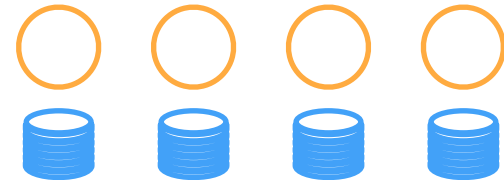- Storage and interfaces
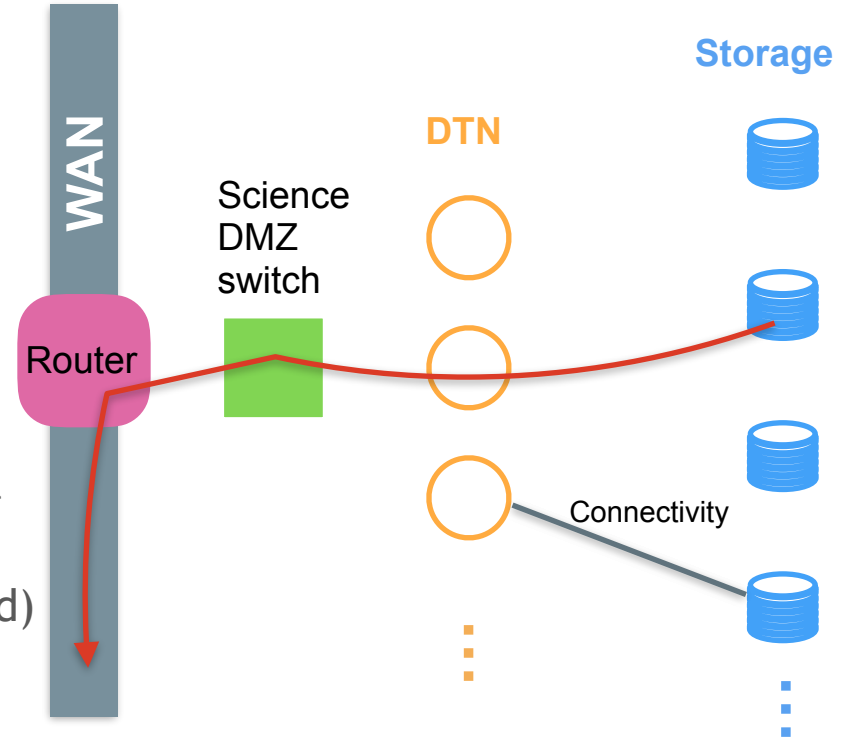
Data Management

File Transfer

WAN

# (Some) performance parameters for data movement

- Storage technology & hardware
- Connectivity
- DTNs
- Switch & Router
- WAN
- Protocols
  - Difference between filling the bandwidth and efficiently reading data
  - Currently GridFTP for transfers and
  - Xrootd for reading (caching/read-ahead)
- … (file size)

**WAN**

**Router**

Science DMZ switch

**DTN**

**Storage**

Connectivity

# WAN: perfSONAR

**288 Active** perfSONAR instances
- Tier1/Tier2 coverage
- Continuously testing 5000+ links
- Testing coordinated and managed from central place
- Dedicated latency and bandwidth nodes at each site
- **Analytic platform to analyses data and send alarms and warnings**

BROOKHAVEN NATIONAL LABORATORY | Scientific Data and Computing Center
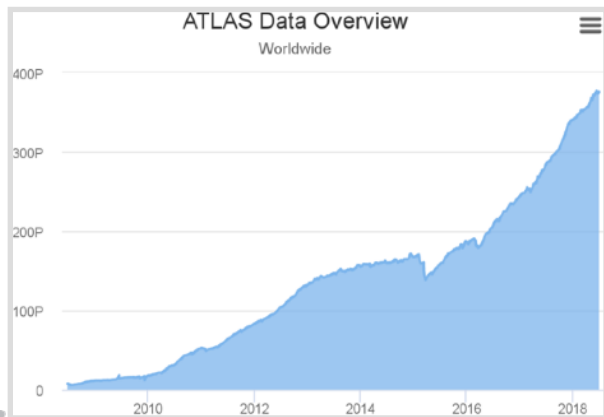
# WAN: perfSONAR

# ATLAS Experiment & RUCIO

Data volume approaching 400PB

10M containers, 20M datasets, 1B files

5K accounts

1-2PB transfered/day, 3PB deletions/day

130 sites, 600 storage endpoints



## Data Management

Global namespace to federate across different storage systems

Control & accounting of data and users

Declarative data management with policies and rules

Transfer orchestration with priorities, shares and activities

Popularity-based replication, caching and deletion

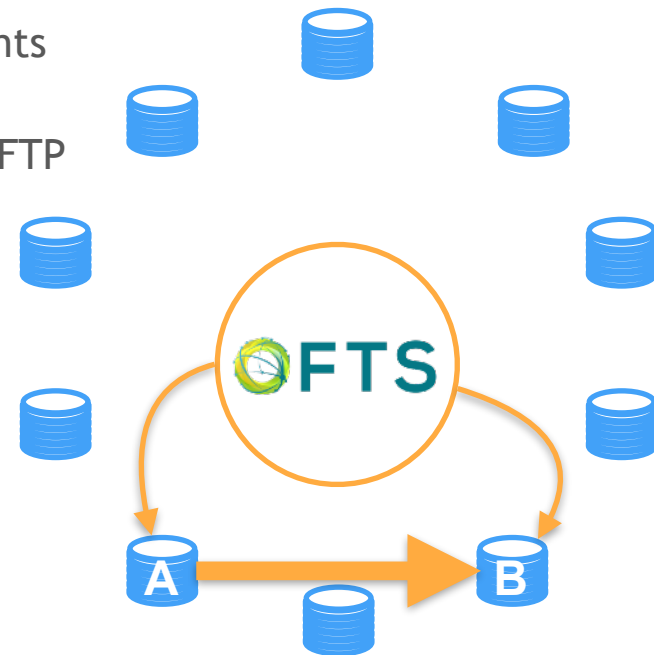Events & messages for synchronisation with other tools

Consistency & repair of broken and missing data

and much more …

# File Transfer System
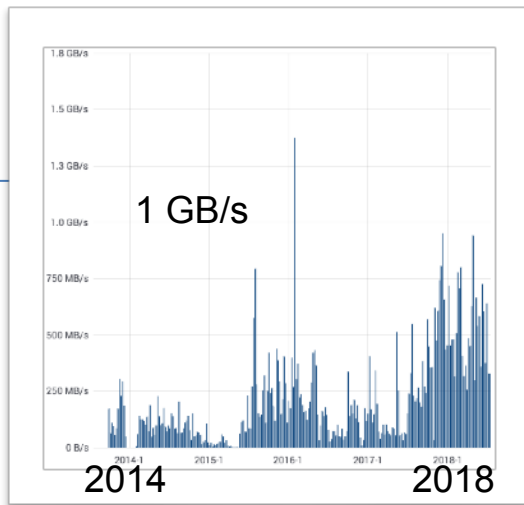
**File Transfer**

- Basic principle: 3rd party transfer between 2 end points
  - FTS should be authorized to talk to A & B
  - A & B should talk same protocol, currently GridFTP
  - Testing HTTP & XRootd

- Accept bulk requests
- Scheduler, shares
- Parallel file transfers
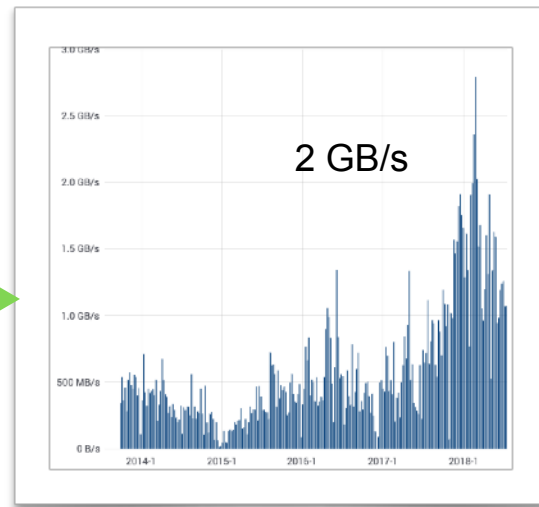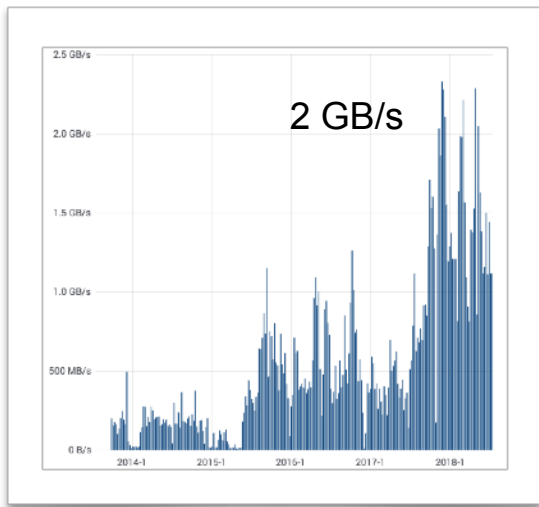- Adaptive auto-tuning
- Multihop
- Session re-use
- ...

Server at BNL

# FTS transfers



1 GB/s

2014     2018

Average transfer rates
Since 2014

2 GB/s

**Within US**

To
US

From US

2 GB/s

# Heterogeneity of storage

- Age
- Latency
- Resiliency
- Size (US: 38 usable PB at 5 locations)
- Can be geographically distributed
- Overhead (raw->usable space)
  - 2-3 actual copies
  - RAID
  - Erasure

- Technology :
  - dCache
  - XRootd
  - Ceph
  - GPFS
  - …

Different funding, Different locations, History….
MoU specify availability to receive data only

BROOKHAVEN
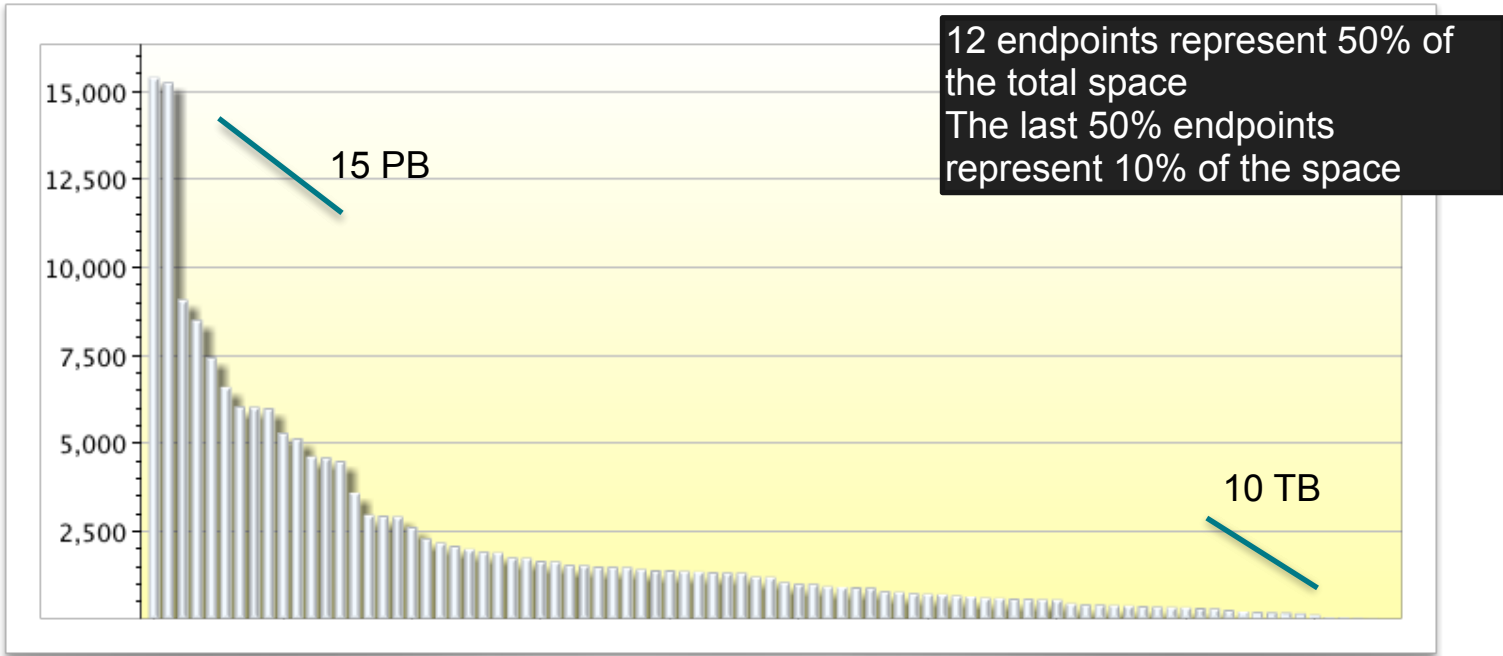NATIONAL LABORATORY | Scientific Data and Computing Center

# Heterogeneity of storage

- Age
- Latency
- Resiliency
- Size (US: 38 usable PB at 5 locations)
- Can be geographically distributed
- Overhead (raw->usable space)
  - 2-3 actual copies
  - RAID
  - Erasure

- Technology :
  - dCache
  - XRootd
  - Ceph
  - GPFS
  - …

Differences are (almost) not taken into account by current data placement policies

# ATLAS: Storage size / end point
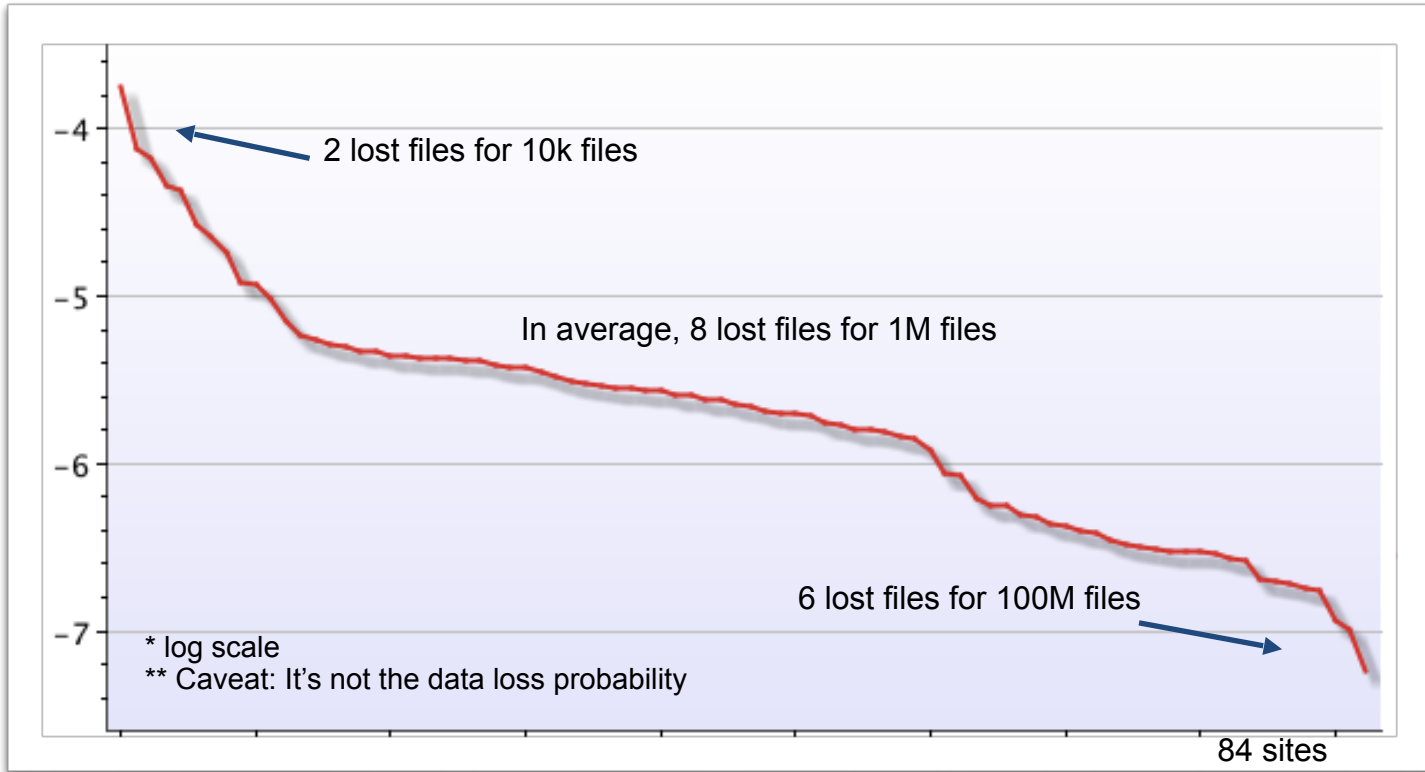


15 PB

12 endpoints represent 50% of the total space
The last 50% endpoints represent 10% of the space

10 TB

**Fragmentation of disk space + operational effort**

BROOKHAVEN
NATIONAL LABORATORY | Scientific Data and Computing Center

# ATLAS: Lost files frequency / site



2 lost files for 10k files

In average, 8 lost files for 1M files

6 lost files for 100M files

* log scale
** Caveat: It's not the data loss probability

84 sites

# ATLAS: Actual transfer rates between end points



Disk to Disk

Network matrix

Transfer at

0.1 KB/s
1 KB/s
10 KB/s
100 KB/s
1 MB/s
10 MB/s
100 MB/s
1 GB/s
10 GB/s

* One month statistic: Maximum throughput during one hour
** Throughput < 100 KB/s  can be due to less transfers and statistics

The real figure of merit!

# And tapes ???

- Available at a few locations (Tier-1s)
- Decoupled from disk storage
    - Used as archive
    - Scheduled access
- Underutilized
- Reliable, cheaper than disk
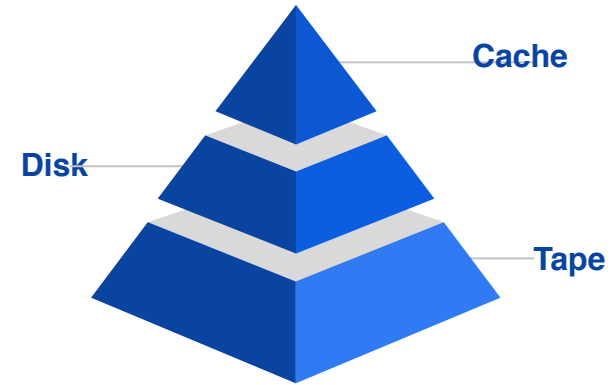- Ongoing tests of optimized tape access via 'tape carrousel'

# Storage distribution evolution

- Reduce number of end points
  - Larger storage entities
- Increase tiered hierarchy
- Introduce QoS
  - Reliability
  - Availability
  - Throughput
  - Redundancy
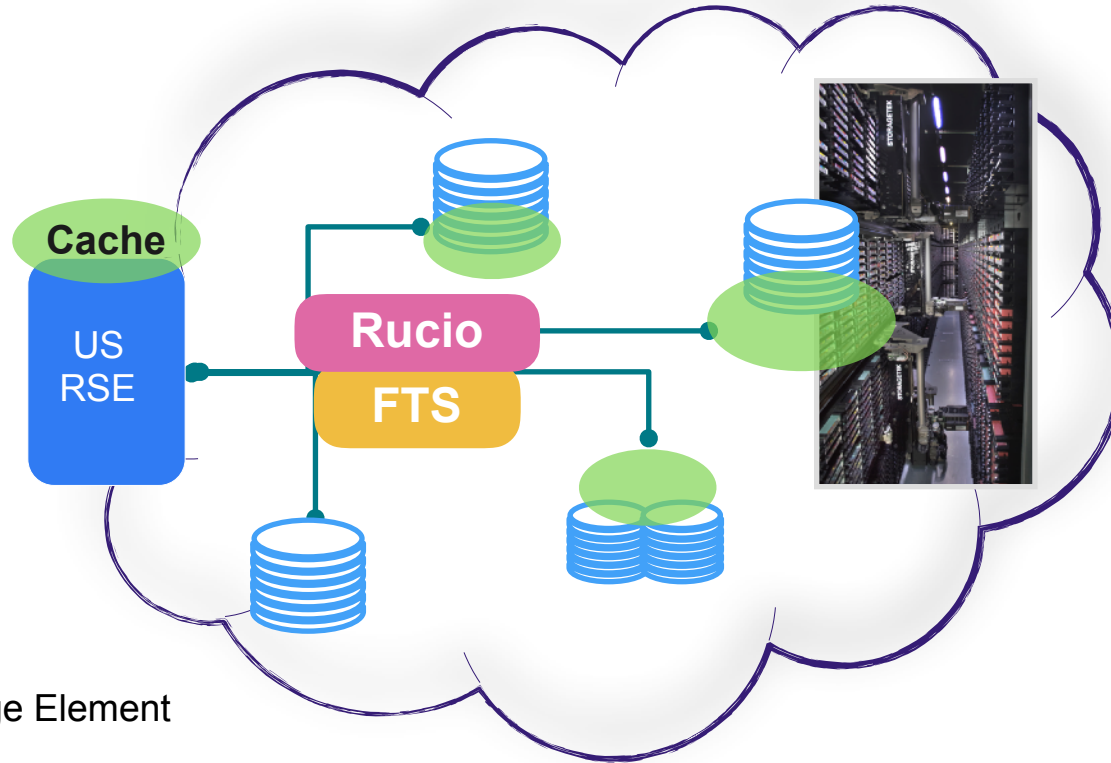  - …



Cache

Disk

Tape

BROOKHAVEN
NATIONAL LABORATORY | Scientific Data and Computing Center

# Proposed prototype for US ATLAS

- Expose unique storage entry point to WLCG
- Internally
  - Different QoS for various components
  - Increasing usage of tape as foundation
- Caching when needed at storage (and CPU) locations
- Rucio (and FTS)
  - To provide unified name space
  - To manage storage hierarchy and data placement and replication

**BROOKHAVEN** NATIONAL LABORATORY | Scientific Data and Computing Center

# Internally : Redirect, Move & Cache



RSE: RUCIO Storage Element

BROOKHAVEN NATIONAL LABORATORY | Scientific Data and Computing Center

# Thank you...