# Tape Evolution pre-GDB report

Alastair Dewhurst

# Pre-GDB on Tape Evolution

- Well attended, in excess of 70 people connected.

- Focus was on the underlying tape systems.
  - Accessing Tape systems is well covered by DOMA-TPC.

- To provide some context I requested some input from VO about their workflows.

- First session was about Tape optimizations.

- Second session covered longer term goals.



Alastair Dewhurst, 10th March 2021

# Data Carousel

- ATLAS Data Carousel has resulted in significant performance improvements since 2018.

- During recalls tape drive throughput is ~30% maximum.
- We want to maximize use of hardware.
  - Can we increase this to 70%?
- Writing larger files can help (~20%)
- The way files are stored on tapes can also help (~20%).

| Sites | 2018 Phase I Test (MB/s) | 2020 Reprocessing (MB/s) |
|---|---|---|
| CERN (CTA Test) | 2000 | 4300 |
| BNL | 866 | 3400 |
| FZK | 300 | 1600 |
| INFN | 300 | 1100 |
| PIC | 380 | 540 |
| TRIUMF | 1000 | 1600 |
| CC-IN2P3 | 3000 | 3000 |
| SARA-NIKHEF | 640 | 1100 |
| RAL | 2000 | 2000 |
| NDGF | 500 | 600 |

Table 1: Stable Rucio tape throughput for the ATLAS Tier-1 sites and CERN, measured from the 2020 reprocessing campaign, with comparison to the Phase I results

# LHCb planned usage

- High level requirement to be able to complete data re-processing in 4 months.

Data Taking

Re-processing

# KIT - TSM to HPSS Migration

- Migrating their data from an Oracle SL8500 to Spectra Logic TFinity Library.

- Migrating from TSM to HPSS as the underlying tape system.

- Improved their scheduling software.

- Increased number of concurrent requests from 2k to 30k.

- Eliminated bottlenecks with new hardware and network.

- Recall rate is improved by a factor of 3 per drive.



Alastair Dewhurst, 10th March 2021

# RAL - CTA and Tape Migration

- RAL provided an update on progress to migrate to CTA.
- RAL has also been migrating data from SL8500 to Tfinity Library.





Remaining tapes to repack

**~20 tapes/day**

**~25 tapes/day**

**Allocate more drives for repacking**

Alastair Dewhurst, 10th March 2021

# CNAF - Dynamic resources allocation

- CNAF noticed that some tape drives were unused even when large scale activity was ongoing.
    - Due to static partitioning of drives between VOs and activities.
- Monitor tape drive usage in InfluxDB and use that information to allocate drives where needed automatically.
    - Increased average throughput by 85%.



Alastair Dewhurst, 10th March 2021

# BNL - Tape Recalls and MAS

- Hiro gave a very detailed talk about directory based sorted writes to tape to maximize recall speed.

- Conclusions were, to improve efficiency:
  - Larger files
  - Fewer gaps between files on Tape
  - Read requests come in bulk to HPSS.

- MAS archives to tape unused files.

- Created 8PB additional space on disk.

**Volume per datatype**

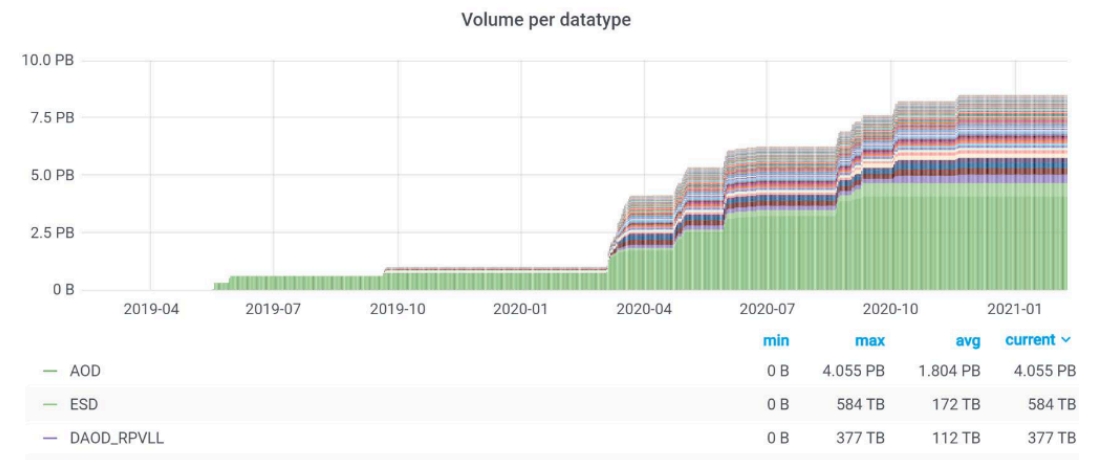

| | min | max | avg | current ⌄ |
|---|---|---|---|---|
| — AOD | | 0 B | 4.055 PB | 1.804 PB | 4.055 PB |
| — ESD | | 0 B | 584 TB | 172 TB | 584 TB |
| — DAOD_RPVLL | | 0 B | 377 TB | 112 TB | 377 TB |

Alastair Dewhurst, 10th March 2021

# RAO (1)

- RAO stands for Recommended Access Ordering.
  - It orders the way data is read from tape and can make a significant difference in performance.
- A Tape is:
  - 960m long
  - Is made up of 4 bands, only one of which can be read at a time.
  - Each band if made up of 52 wraps, which the read head needs to read in the correct direction.
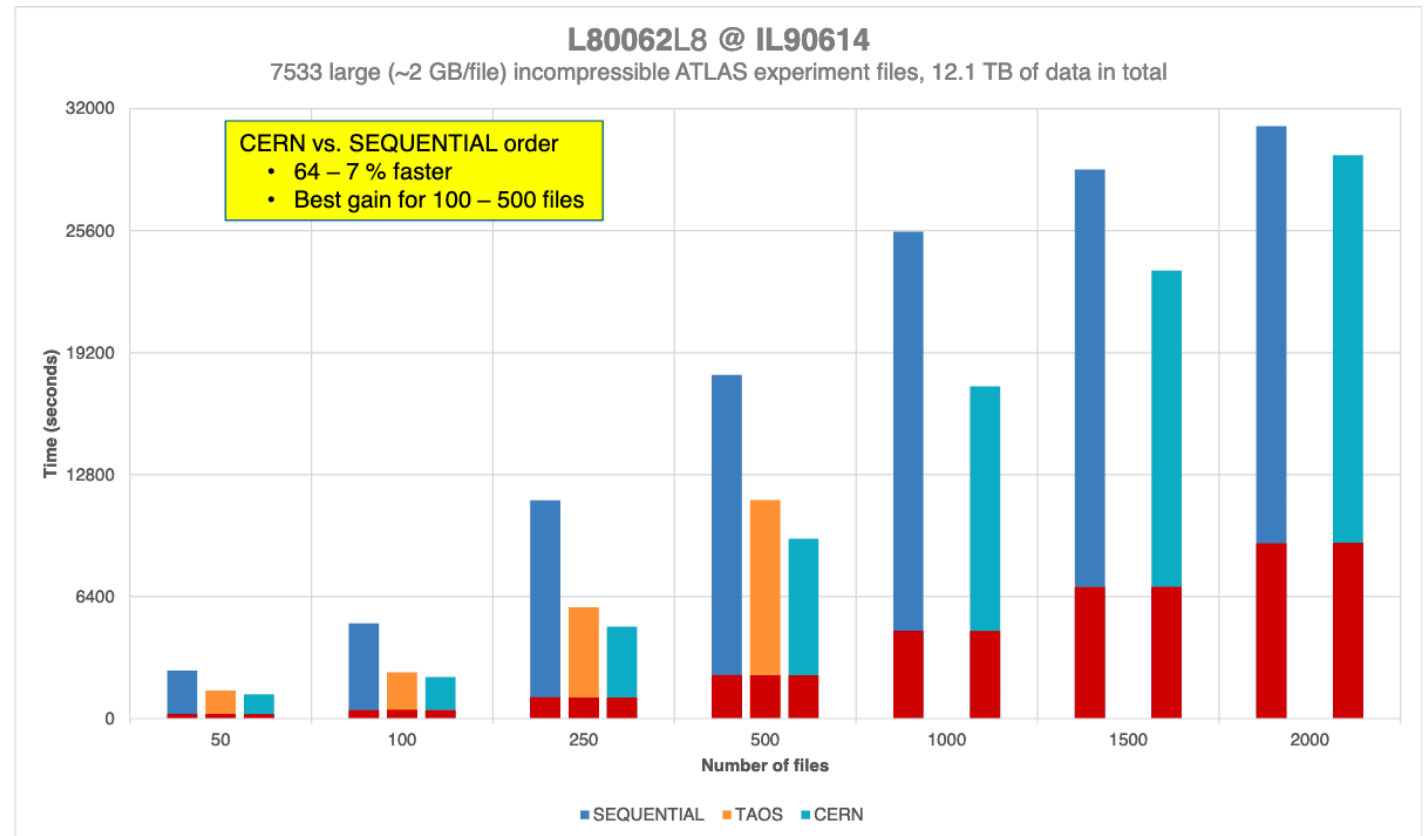  - Each wrap holds 32 tracks which actually contain the 1 and 0s.

# RAO (2)

- Enterprise drives have RAO.

- CERN have implemented RAO for LTO drives.

- Produces up to 64% improvement against sequential reads.

- Lower gains when there are many files (or very few) because sequential works ok.



**L80062L8 @ IL90614**
7533 large (~2 GB/file) incompressible ATLAS experiment files, 12.1 TB of data in total

CERN vs. SEQUENTIAL order
- 64 – 7 % faster
- Best gain for 100 – 500 files

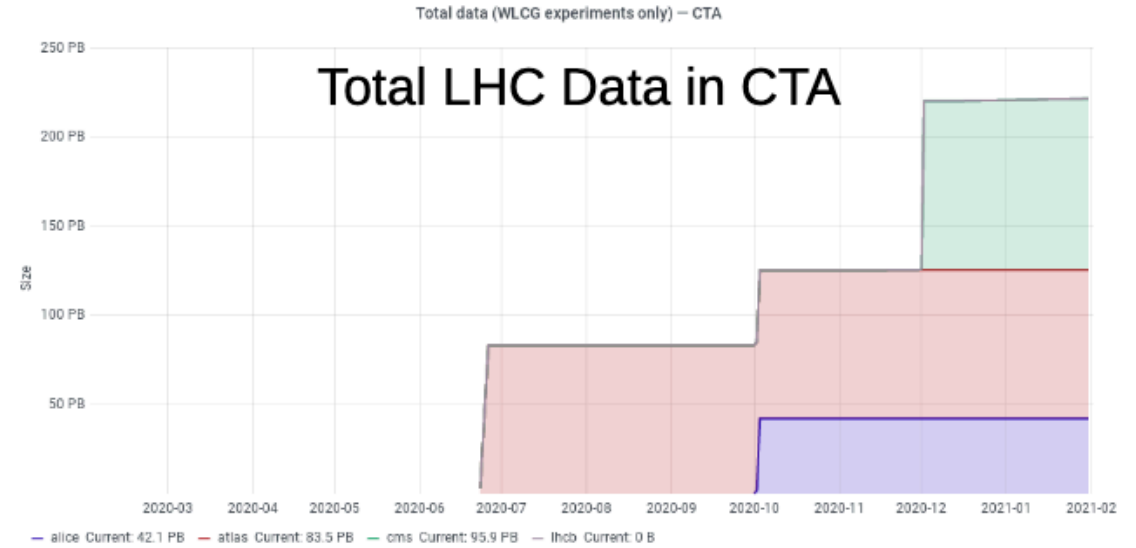Alastair Dewhurst, 10th March 2021

# Cost modelling

- Shigeki Misawa gave a detailed presentation on long term cost modelling.
    - This was a requested follow up from the WLCG Workshop.
- Inefficient use of tape hardware drives up costs!



Alastair Dewhurst, 10th March 2021

# CTA Outlook

- CTA has deployed production instances for Alice, Atlas and CMS.

- Lots still to do. Focusing on ensuring critical features are ready for start of Run 3.

  - CTA is feature-frozen in preparation for Run 3.

- CTA is designed to exploit tape hardware to the max.

  - SSD buffers are needed for throughput.

- CTA is in the early stages of establishing a community.

  - Welcomes contributions / collaborators

Total data (WLCG experiments only) — CTA

Total LHC Data in CTA

250 PB

200 PB

150 PB

100 PB

50 PB

2020-03  2020-04  2020-05  2020-06  2020-07  2020-08  2020-09  2020-10  2020-11  2020-12  2021-01  2021-02

— alice Current: 42.1 PB    — atlas Current: 83.5 PB    — cms Current: 95.9 PB    — lhcb Current: 0 B

Alastair Dewhurst, 10th March 2021

# dCache / DESY Plans

- Tigran gave an update of DESY tape setup and use cases.

- Also gave an update of some of the new dCache developments.

- DESY are evaluating HSM (Tape) Software to replace OSM.
  - Maximize Tape hardware efficiency

- Possible candidates:
  - Plan A (Open Source): CTA or Enstore
  - Plan B (Proprietary): HPSS or TSM

- DESY are looking to make a decision in the next 6 months.

- It should be noted that FNAL are also evaluating the long term status of Enstore.
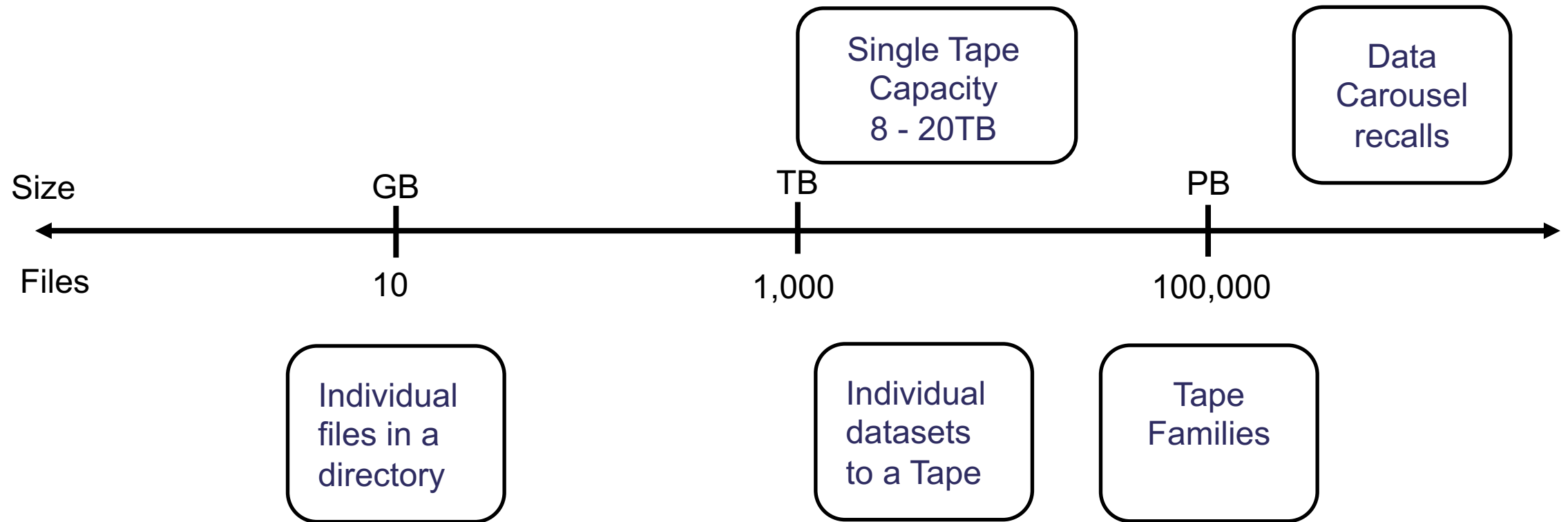
Alastair Dewhurst, 10th March 2021

# Motivation

- I wanted to organize the meeting for the following reasons:
  - Other sites have been interested in RAL's CTA evaluation and the progress we are making with deploying it.
  - I am aware that other sites have been making big improvements in their tape performance and 2020 was not the easiest year to collaborate and share ideas.

- I think Tape can play an incredibly important role in solving HL-LHC data challenges because it is so much cheaper than Disk.
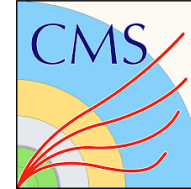  - Tape systems are complicated and I feel some opportunities are not capitalized on.

Alastair Dewhurst, 10th March 2021

# Tape Recalls

- There is a lot of discussion about how to optimize writes.
- Use cases vary, but they are all probably good enough.

Single Tape Capacity 8 - 20TB

Data Carousel recalls

Size

GB        TB        PB

Files

10        1,000        100,000

Individual files in a directory

Individual datasets to a Tape

Tape Families

Alastair Dewhurst, 10th March 2021

# Accessing Tape Endpoints

Communication between VOs and Tape frontends is done really well by various DOMA groups.

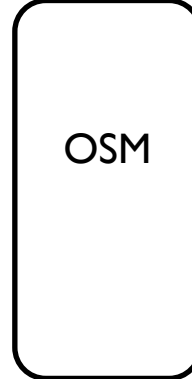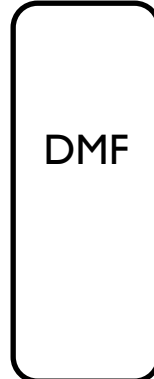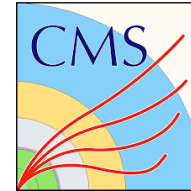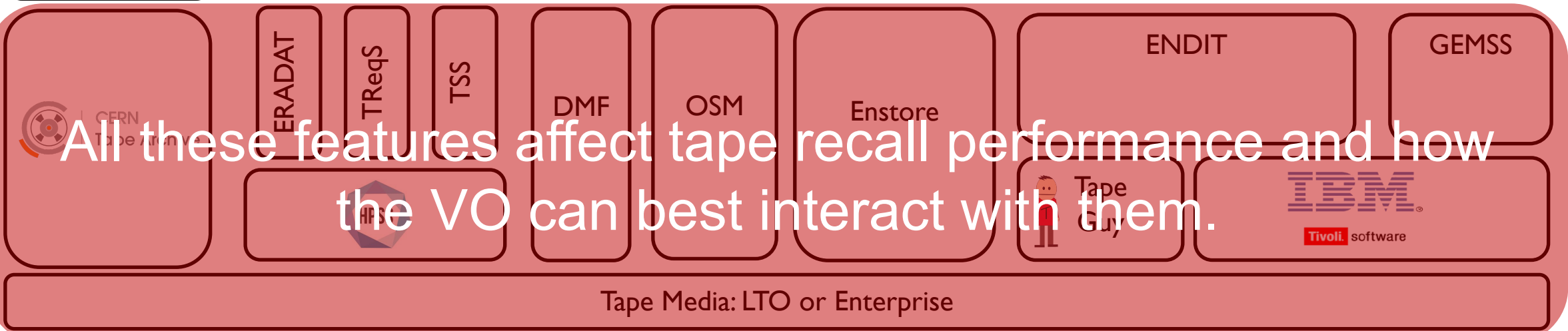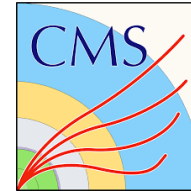| | | | |
|---|---|---|---|
| CASTOR CERN Advanced STORage manager | CERN Tape Archive | dCache | StoRM Storage Resource Manager |
| RAL | CERN | BNL, CNAF, FNAL, JINR, KISTI, KIT, IN2P3, NDGF, PIC, SARA, Triumf | CNAF |

# Optimizing Tape Endpoints

# Optimizing Tape Endpoints



All these features affect tape recall performance and how the VO can best interact with them.

| CERN, RAL | BNL | IN2P3 | FZK | SARA | DESY | FNAL, PIC, JINR | Triumf | KISIT, NDGF | CNAF |

Alastair Dewhurst, 10th March 2021

# A more consolidated future?

Optimizations between the frontend and the tape backend will necessarily be site specific. Sites do collaborate, maybe more could be done?

DMF

Tape Guy

With Recommended Access Ordering the performance difference between Enterprise and LTO should greatly reduce.

Alastair Dewhurst, 10th March 2021

# Conclusions

- An awful lot of work has gone into improving and optimizing tape systems in the last few years.

- Tape drives are expensive, therefore a lot of optimization focus on making the best use of this.

- There is currently a fantastic opportunity for greater collaboration between Tape providers.

- I am sure there will be many smaller discussions but I would like to arrange another pre-GDB (or similar) in ~9 months.

Alastair Dewhurst, 10th March 2021