

# Updates to the ATLAS Data Carousel Project

M.Borodin, D. Cameron, A. Klimentov, T. Korchuganova, M. Lassnig,  
T.Maeno, H. Musheghyan, D. South and *X.Zhao*

On behalf of the ATLAS Computing Activity

CHEP, 2023

Norfolk, VA, USA

## *Team Effort ---*

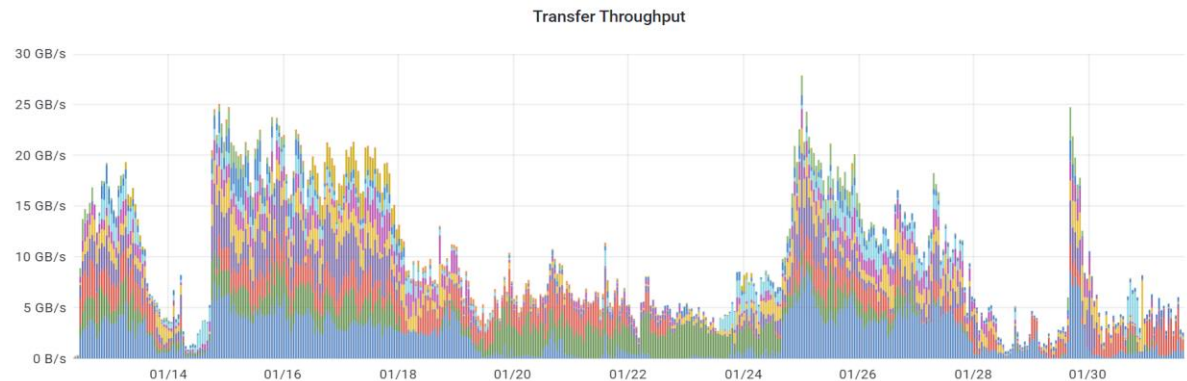
- *workflow management team(WFM)*
- *distributed data management team (DDM/Rucio)*
- *distributed production and analysis team (DPAs)*
- *operations team(Ops)*
- *monitoring team*
- *ATLAS distributed computing(ADC) coordinators and experts*
- *CERN T0 and all T1s storage and tape experts.*

# Outline

- Overview of Data Carousel
- Data on demand R&D
- Tape smart writing R&D
- Conclusion

# Data Carousel in Production

- Facing the data storage challenge of HL-LHC, ATLAS started the Data Carousel project in June 2018, which is a tape-driven workflow, jobs get inputs from tape directly.
- Since 2020, Data Carousel has been in production for managed production campaigns like RAW data reprocessing, derivation and Monte Carlo simulation.
  - Peak overall rate 20~25GB/s
  - 64PB data staged from tape in 2022.
  - Significant disk space savings (less than half of AOD on disk)



Tape staging from T0/T1s in 2023-01 (ATLAS DDM dashboard)

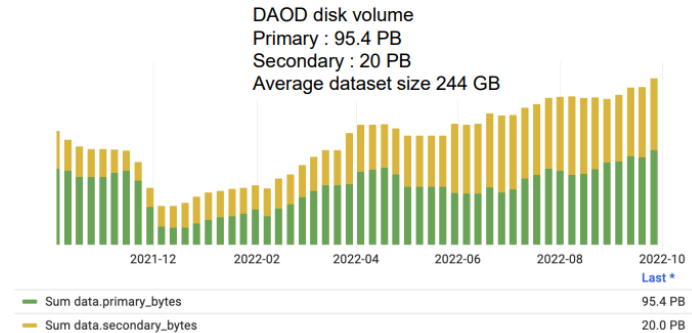
# Current R&D in Data Carousel

- DAOD-on-demand
  - Remove disk replicas for data that has not been accessed for a predefined period.
  - When users request them, they will be either staged from tape or recreated by following the original production steps.
- Tape smart writing
  - Intelligent algorithms for file placement on tape in order to retrieve data back more efficiently
  - Long term strategy to achieve optimal tape usage in Data Carousel
- Both are in the list of ATLAS HL-LHC demonstrator projects, as preparation for the TDR.

# DAOD on Demand ...

# Objectives and General Considerations

- Saving disk space
  - DAOD are inputs to user analysis jobs
  - large DAOD data volume are kept on disk but many of them lie unused
    - ATLAS DDM lifetime model : a procedure applied three or four times per year that identifies and deletes old and unused data, based on a set of policies.
    - Lifetime model exception list: users/physics groups can ask for exceptions for certain DAOD datasets, whose lifetime then extended.
    - Studies show that only ~20% of the DAODs in the exception list are accessed in the next 12 months
- Application and approval of lifetime model exceptions is a very labor-intensive procedure, for both the physics groups and ATLAS Distributed Computing operations.



# Possible Scenarios and Pro/Cons

- DAOD recreation -- delete old and unused DAODs, reproduce them on demand if needed
- DAOD on tape -- copy (archive) datasets to tapes and release disk space for secondary copies of other popular data

## “DAOD recreation”

Pro

Technically feasible

Cons

- How quickly we can reproduce DAODs
  - Including AOD staging from tape
- We need to prove the reproduced data are the same as the deleted ones
- Software development is required (workflow, dataflow, etc)

## “DAOD on tape”

Pro

Use Data Carousel (DAOD just another data format)

Cons

- How quickly we can stage DAOD from tape
- Extra load on tape system
  - DAOD datasets are in general smaller than other data types (like AOD).



# Demonstrators

- We plan to do demonstrator on both options.
- Demonstrators will be done in steps and metrics will be assessed after each step.
  - Step 1 : Proof of Concept (ongoing ...)
    - Choose small DAOD sample from the lifetime model exception list, make sure their parent AOD datasets are available.
    - Let the workflow management systems (ProdSys2/PanDA) to reproduce them measure TTC (time-to-completion)
    - Write the DAOD datasets to tape (preferably the same tape site where the parent AOD datasets reside). Stage the DAODs back to disk, measure the staging time.
    - Measure the time to stage the parent AODs from tape.
    - Initial results (on a 5 DAOD dataset sample) :
      - datasets recreated by PanDA (results to the right)
      - DAOD dataset staging time 5~6 hours (FZK)
      - ProdSys2 PoC for dataset recreation developed
  - Step 2 : ‘user scale’
    - Choose 50~100TB DAOD sample from exception list, again for more than 1 tape sites
    - (Re-)run derivation production, reproduce DAOD from AOD, measure TTC
    - Stage the whole DAOD data sample from tape, and stage the whole AOD data sample from tape.
  - Steps 3 : ‘production scale’ – TBD after ‘user scale’ step
- Demonstrator will only be successfully with Tier-1s and CERN participation and expertise.
  - To mitigate the impact on the results from differences among different site tape systems and site operational models, we try to run these tests with more than 1 tape sites. Right now BNL, RAL and FZK Tier-1s take part in the demo.

New task	Original task	Nucleus	Input size (GB)	New completion time (hrs)	Old completion time (hrs)
31920802	28174352	BNL-ATLAS	5537	6	90
31937895	28174668	TRIUMF-LCG2	2573	5	146
31938004	27562319	IN2P3-CC	9446	10	23
31941344	27726219	FZK-LCG2	2321	4	59
31943808	28174463	CERN-PROD	1929	4	96

# Metrics and Deliverables

- Time to completion (TTC)
- Validation of the new DAOD datasets in the DAOD recreation option
- Doing a spreadsheet/python notebook exercise to calculate how much space a no-exception lifetime model would involve, how much CPU we'd need to burn to recreate deleted datasets and how much data will be staged in. This relied on knowing the fraction of datasets on exceptions lists that are actually used.
- Load on tape systems
  - How much extra tape bandwidth will be needed for the two options respectively, extrapolated to HL-LHC conditions.
  - How efficient can DAOD be staged from tape ? If the efficiency is lower than the other data types, and DAOD staging has a lower latency requirement than other data types, it will lead to a bigger resource purchases at tape sites.

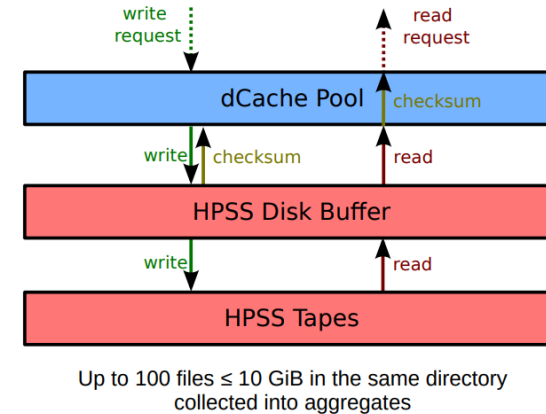
# Tape Smart Writing ...

# Overview and the Plan

- Tape exercises by ATLAS and other WLCG experiments have shown that the key to optimal tape usage is to better co-locate files on tape, so called “smart writing”.
  - A “catch all” phrase that encompasses the multitude of techniques that are possible to intelligently lay out data on tape to enhance read performance.
  - Techniques to improve read access must match how the data is written to tape with how data is read back.
- ATLAS tape sites are at different stages w.r.t. implementing site level smart writing solutions. To facilitate the process, we will do a demonstrator with sites :
  - Run exercise with selected sites, to demonstrate and assess existing smart writing solutions;
  - Proposal from several sites to develop tape system simulator to evaluate possible optimization techniques. The simulation can be compared with real case exercise results, to provide guidance for further improvements, not only for sites doing the simulation, but also to help other sites finding and improving their own solutions.
- Throughout the process, other deliverables are expected, including :
  - Mechanism to pass data grouping metadata from ATLAS DDM to site tape system
  - Tape system monitoring at both site and VO level, e.g. installed bandwidth capacities vs delivered throughput.

# Smart Writing Exercise with KIT (1/2)

- KIT is migrating to HPSS tape system, with smart writing feature
  - Dataset is grouped together using tape file families, which are random numbers and not associated with specific namespaces.
  - Files are written to tapes in file aggregates --- when recalled, the entire aggregate is staged (Full Aggregation Recall mechanism)
  - For more details please refer to the [KIT poster: Efficient interface to the GridKa tape storage system](#)



The smart writing exercise will assess and demonstrate performance of the new tape system setup, with the following metrics :

- Overall tape write/read throughput
- Tape write/read throughput per tape drive

# Smart Writing Exercise with KIT (2/2)

- Two phases test plan
  - Functional test (ongoing)
    - Validating the full chain of tape smart writing, all components work as expected. HPSS operations and monitoring handled by site experts internally.
  - Real use case in production
    - Once KIT HPSS is put in production for ATLAS, we will conduct tape write/read tests using O(100TB) data samples, using ATLAS Data Carousel and DDM services, monitoring at both VO and site level.
- Success of this exercise is contingent on the availability of metadata, which tells the tape system the name and size of the dataset each file belongs to, during tape writing.
  - ATLAS dataset size varies a lot, site needs the dataset size information to dynamically control how many tape drives will be used to migration a dataset to tape, to meet the throughput requirements while keeping files from the same dataset co-located on tape.
  - The metadata will be injected into FTS transfer requests by DDM/Rucio, passed to KIT tape system via FTS and dCache storage services.
    - An initial implementation using URL parameters in the https transfer request has been tested successfully
    - Long term solution is for Rucio to include metainfo as json dictionary in the FTS transfer job, FTS then transforms it into HTTP headers in the WebDav requests to the tape endpoint.

# Tape system simulation

- Several sites (Tier-0, Tier-1) have plan for a tape system simulator
  - Instrument and analyze ATLAS write/read access patterns by accessing I/O “transaction” logs at various levels (ADC, FTS, site storage), simulate real world I/O on the tape system, identify opportunities for optimization.
  - Implement an evaluation framework, e.g. “re-playing” read and write transactions on tape system simulator, to assess the effectiveness of different data placement strategies in achieving desired data layout on tape.
- Collaboration between ATLAS and sites
  - Better understanding of ATLAS tape access patterns
  - Identify data grouping opportunities
  - Collect information about which parameters, metrics and measures much be defined and agreed to improve tape usage and mechanisms through which the identified mechanisms might be achieved and monitored.
  - Access to the historical log information from various ATLAS Distributed Computing services.

# Conclusion

- Data Carousel has been running successfully in ATLAS production since 2020, resulting in significant disk space savings.
- New R&D projects are underway to
  - Explore new opportunities of disk space savings by moving unpopular data to tape and recreate them on-the-fly if needed
  - Address long term tape usage optimization strategies, starting from tape writing, to prepare for the HL-LHC scale.
  - Collaboration between experiment and sites is crucial for the success of these projects.