

CERN Tape Archive: a distributed, reliable and scalable scheduling system

Eric Cano¹ Vladímir Bahyl¹ Cédric Caffy¹ Germán Cancio¹ Michael Davis¹ Oliver Keeble¹
Viktor Kotlyar² Julien Leduc¹ Steven Murray¹

¹CERN, Geneva

²Institute for High Energy Physics, Protvino, Russia

CERN Tape Archive and EOS

- CERN Tape Archive (CTA) is a tape backend for EOS and possibly dCache
- EOS CTA is a replacement for CASTOR instances
- CTA manages all aspects of the tape system
 - Physical location on tape: file catalogue
 - Data transfer: tape server
 - Transfer queuing: scheduling system
- All main 4 LHC experiments migrated from CASTOR to CTA
- Other experiments have or are being migrated

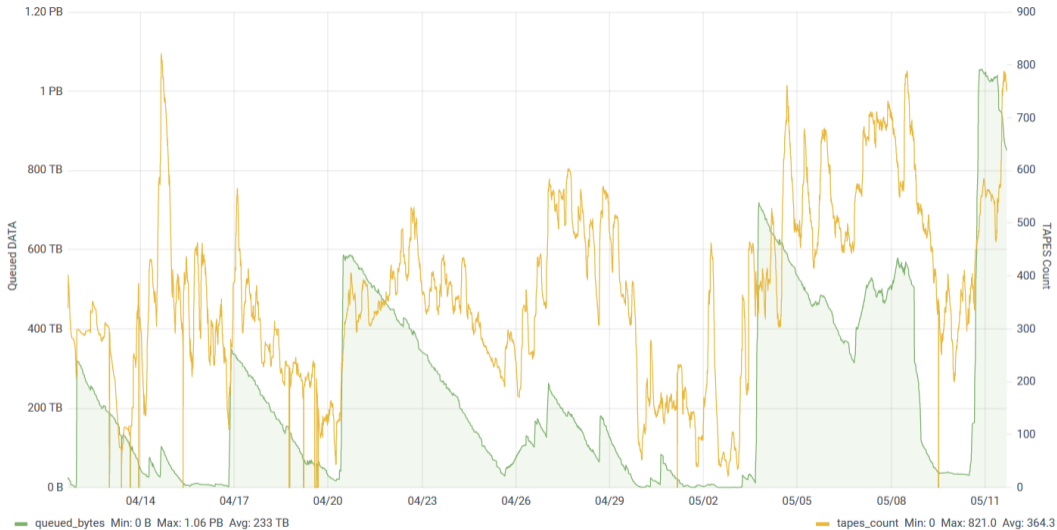
Issues and opportunities with CASTOR + EOS

- EOS and CASTOR disk systems are redundant services
- CASTOR dependent on Oracle's PL/SQL at its core and hence vendor locked-in
- Design constrained by past database limitations that no longer exist
 - Limited number of connections to a database server
- Scattered logic over multiple databases and daemons in CASTOR
 - Multiple databases hold parts of the data
 - Multiple daemons take decisions in incremental steps lacking global view
 - Ex: Commit to write to a tape located in a library with no free drive
 - Communication channels between daemons costly to maintain and evolve

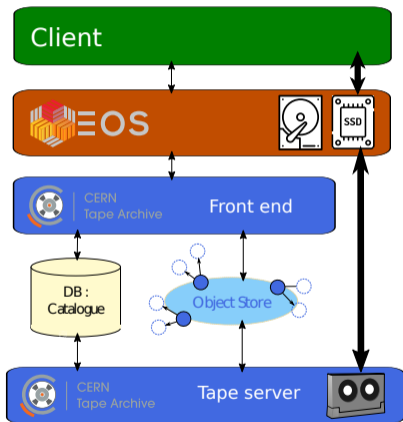
Tape queueing: more than FIFOs

- Tapes is the cheapest storage, with good bandwidth
 - 360 MB/s per drive
- But very high latency (1 min)
- Efficiency highly dependent on scheduling
- Transfers are accumulated until enough work warrants a mount
- Up to 10^6 s of files are queued on 100s to 1000s of queues
- Mount scheduling decides which queue to process
- Then, within a mount, which files to process

RETRIEVE requests by TAPE POOL: queued DATA (left) vs. TAPES count (right) over a month

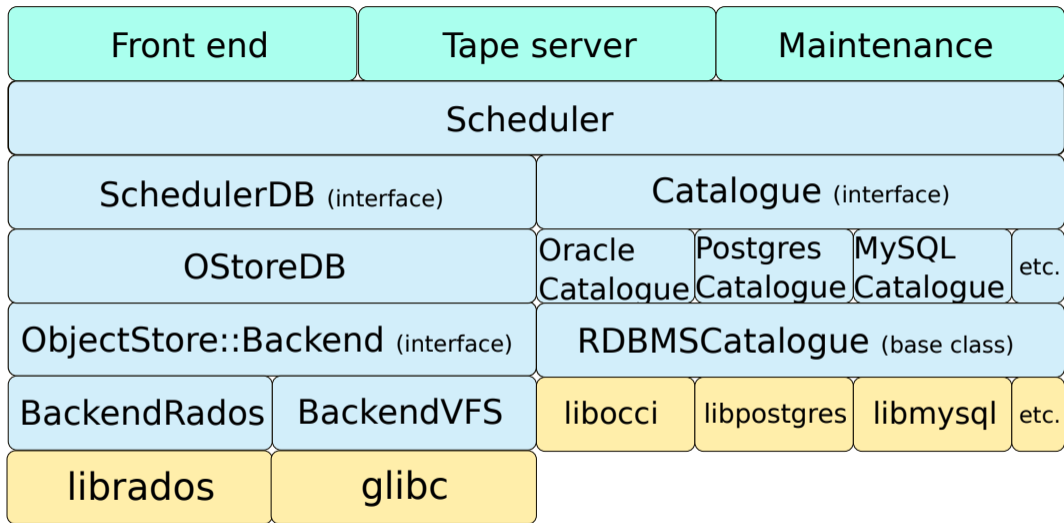


CTA architecture



HDD icon: <https://commons.wikimedia.org/wiki/File:Hard-drive.svg>
SSD icon: <https://commons.wikimedia.org/wiki/File:Ssd.svg>
Tape icon: https://commons.wikimedia.org/wiki/File:Tape_cmta_cassette_backup.svg

- Direct access to shared metadata for all processes
- Two metadata storage systems
 - Database for persistent data: files and tapes catalogue
 - Object store for transient data: queueing information and drive statuses
 - Multiple backends for both
- Scheduler (C++ class): a single entry point shared by all processes of CTA
- Self-scheduling by drives: examines all queues, pick one to process
 - Filter possible mounts, prioritize, pick first



Scheduler DB, aka Objectstore

- Tree of queues, requests, statuses for drives, tapes and processes
- Serialization based on Protocol Buffers
 - Only hard dependency in CTA
- Backends provide object level access
 - Lockable, atomically updated
 - Asynchronous access
 - Reliable
- Need multi-object coherence
 - With potential process crashes
- Scaling needed (few 100 Hz per queue)
 - 100s of drives
 - 1000s of queues
 - 10^6 files per queue

Reliability and scalability

- Algorithm achieving multi-objects operations safety
- Heartbeat based process crash recovery
- Scalability achieved with appropriate contention domains
 - Cutting structures in smaller parts where needed
 - Sharding queues for limited size of updates
- Lockless reads whenever appropriate
- Rate achieved by batching and asynchronicity
- Contention limited in practice to write access to queues
- Latency hiding mechanism added in front end (user visible)

Conclusion

- CTA successfully replaces the tape part of CASTOR after 20 years of service
- Provides a new scalable and efficient queueing and scheduling mechanism
- Flexible design allowing interfacing with other front-ends
- Designed for long-term with no vendor or technology lock-in