# Scheduling Evolution in CERN Tape Systems

**~1990 SHIFT** · manual scheduling

**~1998 STK libraries** · robotic scheduling,  namespace → CERN IT
repack introduced !

**next ~22 years CASTOR** · centrally deployed scheduling

**~2020 CTA** · distributed multi-threaded scheduling
using ObjectStore technology

**~ 2024** · CTA adopted by wider community !
**preparation for Run 4** · new challenges ahead

# Scheduling DB Stores Transient Metadata

**Disk buffer**
XrootD SSI / gRPC
WorkFlow Engine

**request**

**CTA Frontend**
XrootD SSI / gRPC
Scheduler thread

**data**

**CTA Tape Server**
tape drive daemons
Scheduler thread

**Scheduler DB**
Ceph ObjectStore

**transient metadata**

**permanent metadata**
tape file namespace

**Catalogue DB**
Oracle
PostgreSQL

# Scheduler DB Implemented as ObjectStore

## Architecture

- **motivated by performance for archival and retrieval queueing operations**

- **multi-threaded interface to Ceph**

- **protobuf serialised objects in key/value store**
  **(archive, retrieve requests and queues)**

- **code design ensures**

  ▸ **high performance**

  ▸ **scaling**

  ▸ **reliability**

    **(despite > storage round trips than DB)**
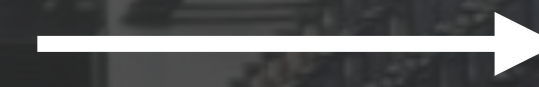
- **delivered very well for Run 3**

## Development Cost

- **complex distributed transaction management system**

- **high maintainability cost**
  - ▸ additional dependency
  - ▸ requires extensive learning effort

- **scheduler logic is tightly coupled to ObjectStore implementation**
  - ▸ lack of indexes is a serious constraint on implementation of scheduler algorithms (e.g. cancelling requests)
  - ▸ as we scaled up to ~200 tape drives, global locking caused scheduler contention

- **workflows beyond original design proven difficult**

**but** →

design allows multiple backends
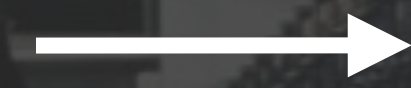
and

code complexity reduction possible

and

we have off-the-shelf solution = Relational DB

## Operational Cost

- **ObjectStore requires additional technology expertise**      ⟶

  ▸ sites which do not run Ceph

  ▸ FTS scheduling using Postgres

- **complex backend object structure**

  ▸ **object introspection, forensics and cleanup is difficult**

  ▸ **"schema" updates difficult to manage**

- **high priority fixes still required several times a year**

  (e.g. for object deletion, empty shard handling, infinite
  loops, global locking issues, repack exhausting resources,
  object size handling, etc.)

> **we can consolidate on common technologies**
>
> and
>
> **have Relational DB rows for operators**
>
> and
>
> **simpler operational tools**

# Advantages of Relational DB as the Solution

## New Scheduler DB
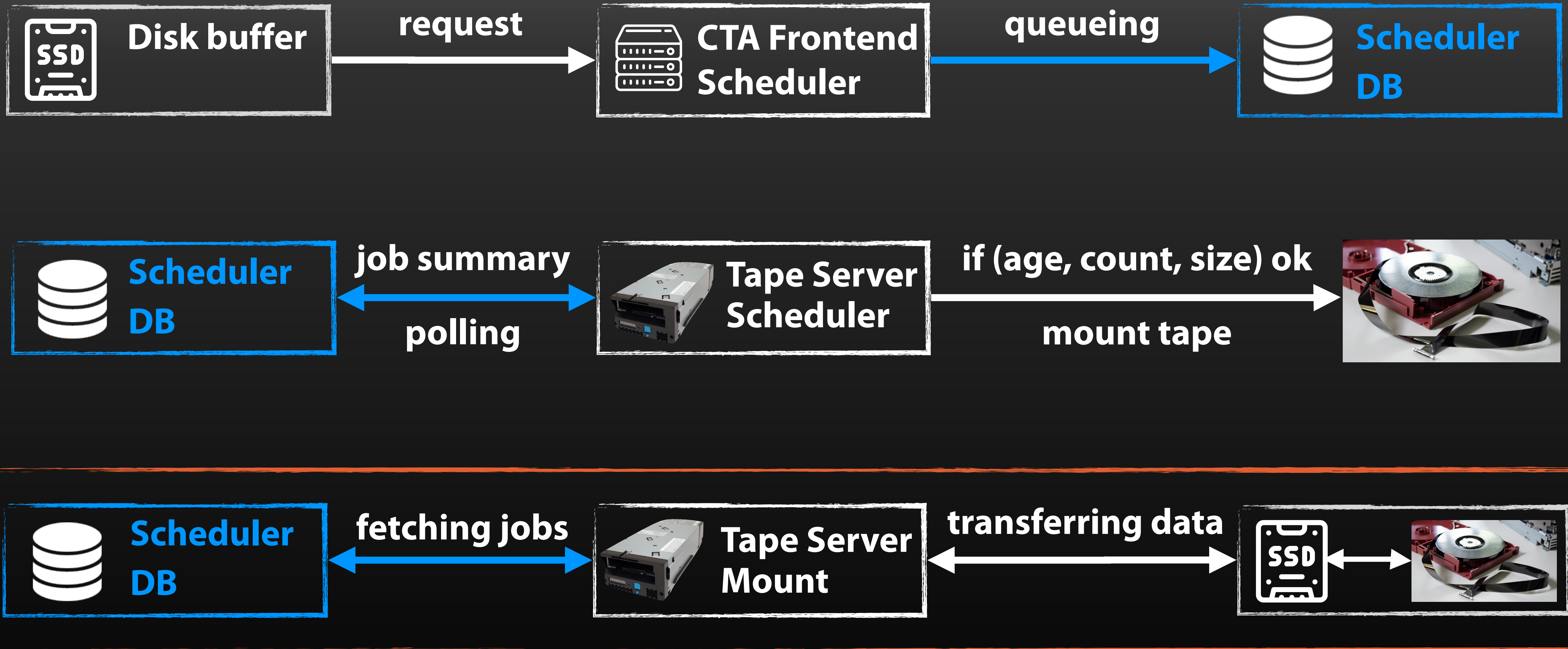
- **lower code complexity**

  ▸ **no overhead of transactional management code**

    MVCC, indexing (+sync) "for free" from DB

- **straightforward schema updates**

- **multi-index queues → more flexibility to improve and extend Scheduling algorithms**

- **extensible to multiple database backends**

  (PostgreSQL DB, Oracle DB)

- **no global scheduling lock**

time

**Disk buffer** → request → **CTA Frontend Scheduler** → queueing → **Scheduler DB**

**Scheduler DB** ← job summary / polling → **Tape Server Scheduler** → if (age, count, size) ok / mount tape →

**Scheduler DB** ← fetching jobs → **Tape Server Mount** ← transferring data →

**blue = scheduler DB tasks**

**time**

**transfer finished**

**Scheduler DB** ← **Tape Server Mount** → **update catalogue** → **Catalogue**

**status update**

**fetching transfers**

**Scheduler DB** ↔ **Tape Server DiskReporter** → **reporting** → **Disk buffer**

**to report**

SSD

**blue = scheduler DB tasks**

# Scheduler DB Implementation Architecture

**DB logic**

**Interface to DB**

**Scheduler logic**

**Interface to DB logic**

```
Scheduler
(cta/scheduler)
```
→
```
SchedulerDatabase
(cta/scheduler)
```
→
```
OStoreDB
(cta/scheduler/OStoreDB)
```
→
```
objectstore
(cta/objectstore)
```

```
RelationalDB
(cta/scheduler/rdbms)
```
→
```
rdbms
(cta/rdbms)
```

## Design Allows Multiple Scheduler DB Implementations

- **use CTA generic RDBMS interface for Postgres DB implementation**

- **workflow oriented tables (Archive/Retrieve/Repack), views, sequences, etc.**

  **(file transfer job = row in a table)**

# Finished Implementation of Archival Workflow



**Disk buffer**



## Archival

✓ **queueing**

✓ **job summaries per mount decision**

✓ **job fetching for transfer**

✓ **management of failures and retries**

✓ **transfer status reporting**

✓ **improvements in CTA rdbms layer**

## Lower Granularity Locking

✓ **drive deciding to mount a tape**

- **lock per logical library**

  (prevents empty mounts)

✓ **fetching the jobs from the queue**

- **lock per tape pool per workflow**

  (avoids interwoven row sets per drive)

✓ **no global lock on scheduling anymore !**

# Functional Testing of Archival Workflow

## Setup

- **Full System deployment in Minikube** (CTA+EOS)

  ‣ **1 disk for buffer (EOS MGM & FST)**

  ‣ **1 disk for tape drive**

- **External  Catalogue: Oracle DB**
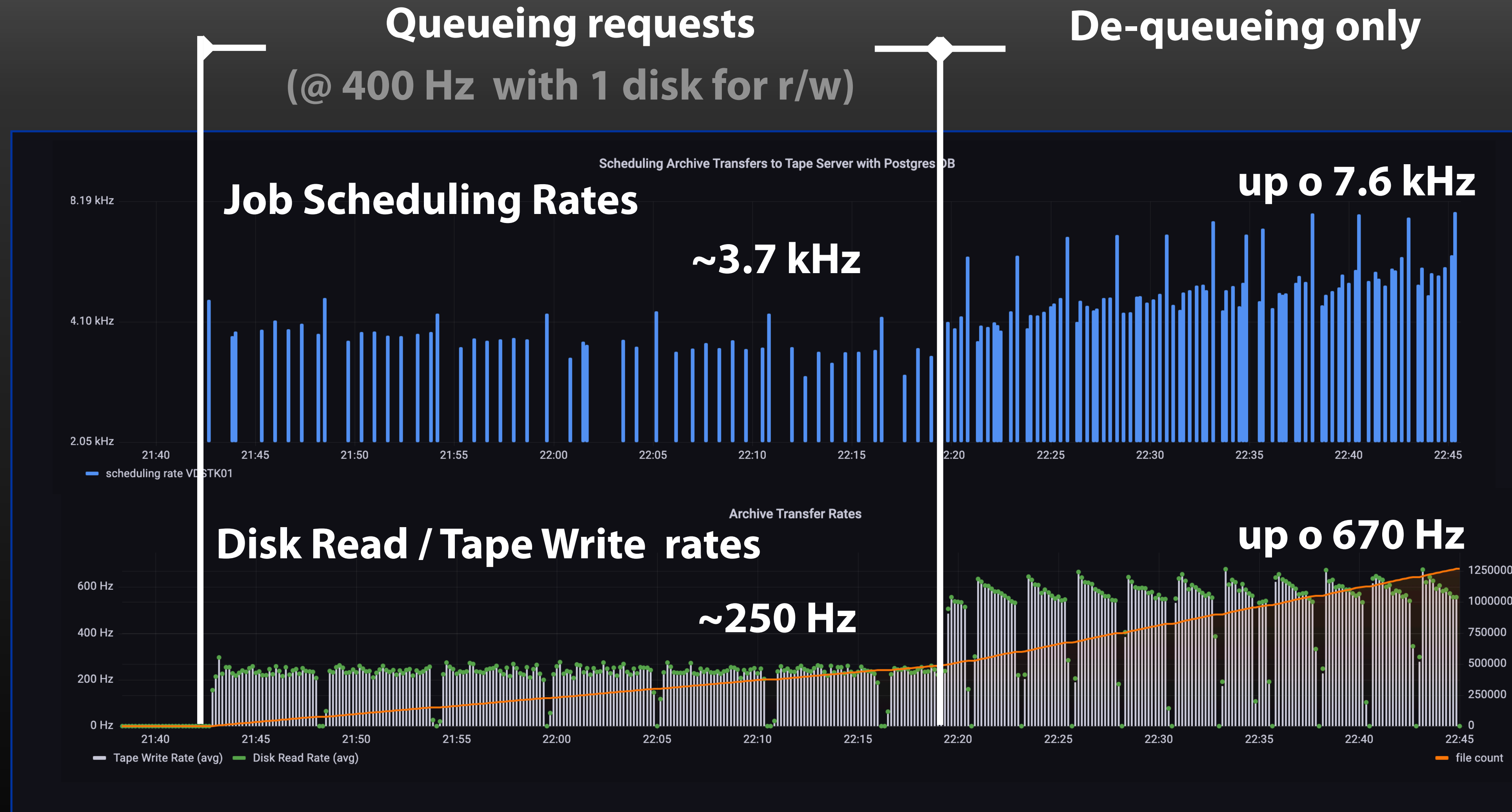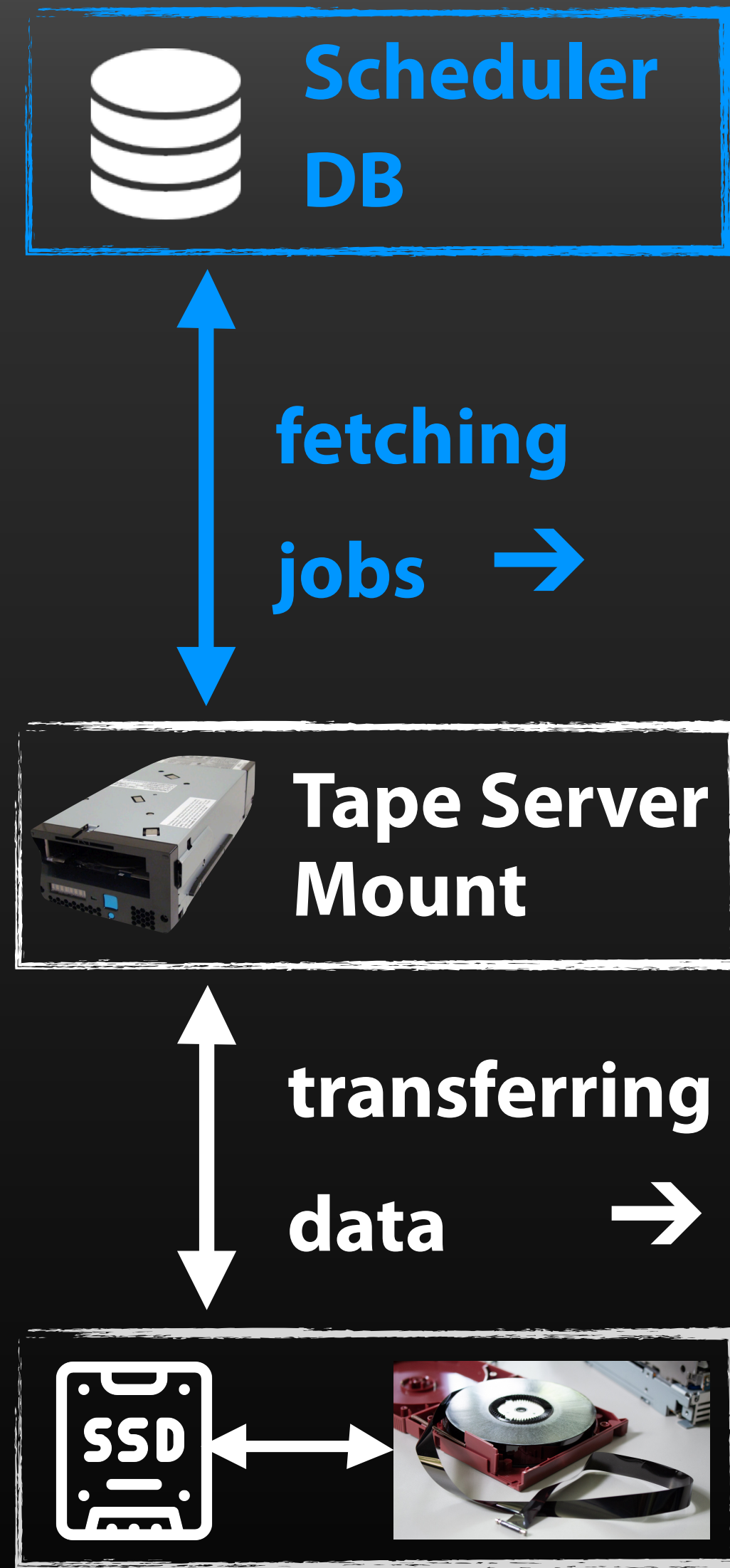
- **External  Scheduler: PostgresDB**
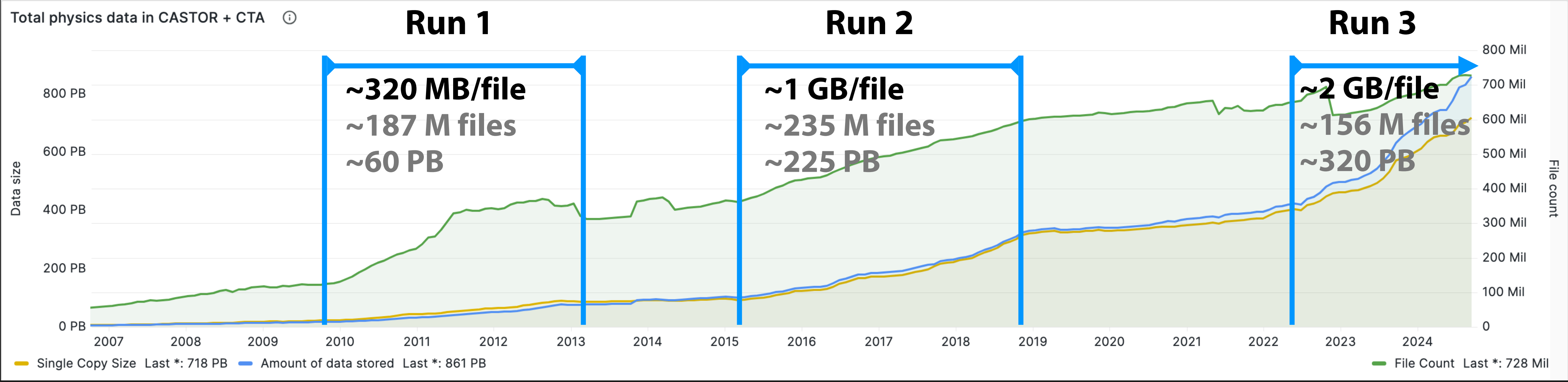
## Load

- **1M files, 128 B each**

- **1 tape pool** (30 MB tapes)

**Kubernetes deployment for performance and scalability tests is in works !**

# Functional Test of Archive Scheduling

Scheduler DB

fetching jobs →

Tape Server Mount

transferring data →

SSD

**Queueing requests**

(@ 400 Hz with 1 disk for r/w)

**De-queueing only**

Scheduling Archive Transfers to Tape Server with Postgres DB

**Job Scheduling Rates**

up o 7.6 kHz

~3.7 kHz

8.19 kHz

4.10 kHz

2.05 kHz

21:40   21:45   21:50   21:55   22:00   22:05   22:10   22:15   22:20   22:25   22:30   22:35   22:40   22:45

— scheduling rate VDSTK01

Archive Transfer Rates

**Disk Read / Tape Write rates**

up o 670 Hz

~250 Hz

600 Hz

400 Hz

200 Hz

0 Hz

21:40   21:45   21:50   21:55   22:00   22:05   22:10   22:15   22:20   22:25   22:30   22:35   22:40   22:45

1250000
1000000
750000
500000
250000
0

— Tape Write Rate (avg)   — Disk Read Rate (avg)   — file count

# Minimal File Rate Requirements

**Total physics data in CASTOR + CTA** ⓘ

Run 1 | Run 2 | Run 3

~320 MB/file
~187 M files
~60 PB

~1 GB/file
~235 M files
~225 PB

~2 GB/file
~156 M files
~320 PB

— Single Copy Size  Last *: 718 PB   — Amount of data stored  Last *: 861 PB   — File Count  Last *: 728 Mil

## Conservative estimates

- file sizes grow
  → keeps minimal rate down

|  | RUN 3 | RUN 4 |
|---|---|---|
| CTA SLA | 50 GB / s | 125 GB / s |
| Avg. File Size | 2 GB | 2 GB |
| Avg. Throughput | 25 Hz | **63 Hz** |

**Postgres DB Scheduling**

**several times faster than ~ 100 Hz = peak rate of Run 3**

# Implementation and Deployment Roadmap

**Implementation:**

**Code Optimisation:**

Retrieve and Repack

Performance & Scaling Tests

**Deployment:**

Kubernetes Test Cluster

Repack Instance

Production

Q4 2024 — Q1 2025 — Q2 2025 — Q3 2025 — Q4 2025 — Beyond Q1 2026

**Not planning to replace ObjectStore on a timescale of Run 3**

# Summary

## Next Scheduler DB Evolution

✓ **meets with Run 4 challenges**

✓ **satisfies the needs of larger CTA Community**

✓ **decreases future development and operational costs**

## Implementation Status

✓ **PostgreSQL as CTA Scheduler DB for Archival Workflow is functional**

✓ **Retrieve and Repack are coming next**

✓ **improved locking granularity**

## Long Term Prospects

✓ **tape supply logic part of scheduling**

✓ **refactoring and cleaner handling of several workflows:**
   *(leveraging DB features)*

   ‣ **request/file deletion**

   ‣ **multiple retrieves per file**

   ‣ **priority queues**

# Thank you !

Thank you for your attention

... special Thanks to few of my CERN colleagues for their help
(David S., Elvin S., Joao A., Julien L., Michael D., Steven M., Pablo O. C., Vlado B.)

I welcome any questions or comments !

# Backup

## CERN Tape Archive

- Archive of the physics data
- Provisioned capacity: ~1.18 EB

- Libraries:
  - 4 x IBM TS4500
  - 2 x Spectra Logic TFinity
- Drives:
  - 40 x TS1170, 46 x IBM1160
  - 88 x LTO-9, 10 x LTO-8
- Media:
  - 150 PB on 3592JF, 150 PB on 3592JE, 227 PB on 3592JD
  - 551 PB on LTO-9, 17 PB on LTO-8, 59 PB on LTO-7M



- Backup of the business data
- Licensed capacity: ~15 PB

- Libraries:
  - 1 x IBM TS4500 (partitioned)
  - 1 x Spectra Logic TFinity (partitioned)
- Drives:
  - 10 x LTO-9
  - 10 x LTO-8
- Media:
  - 12 PB on LTO-8
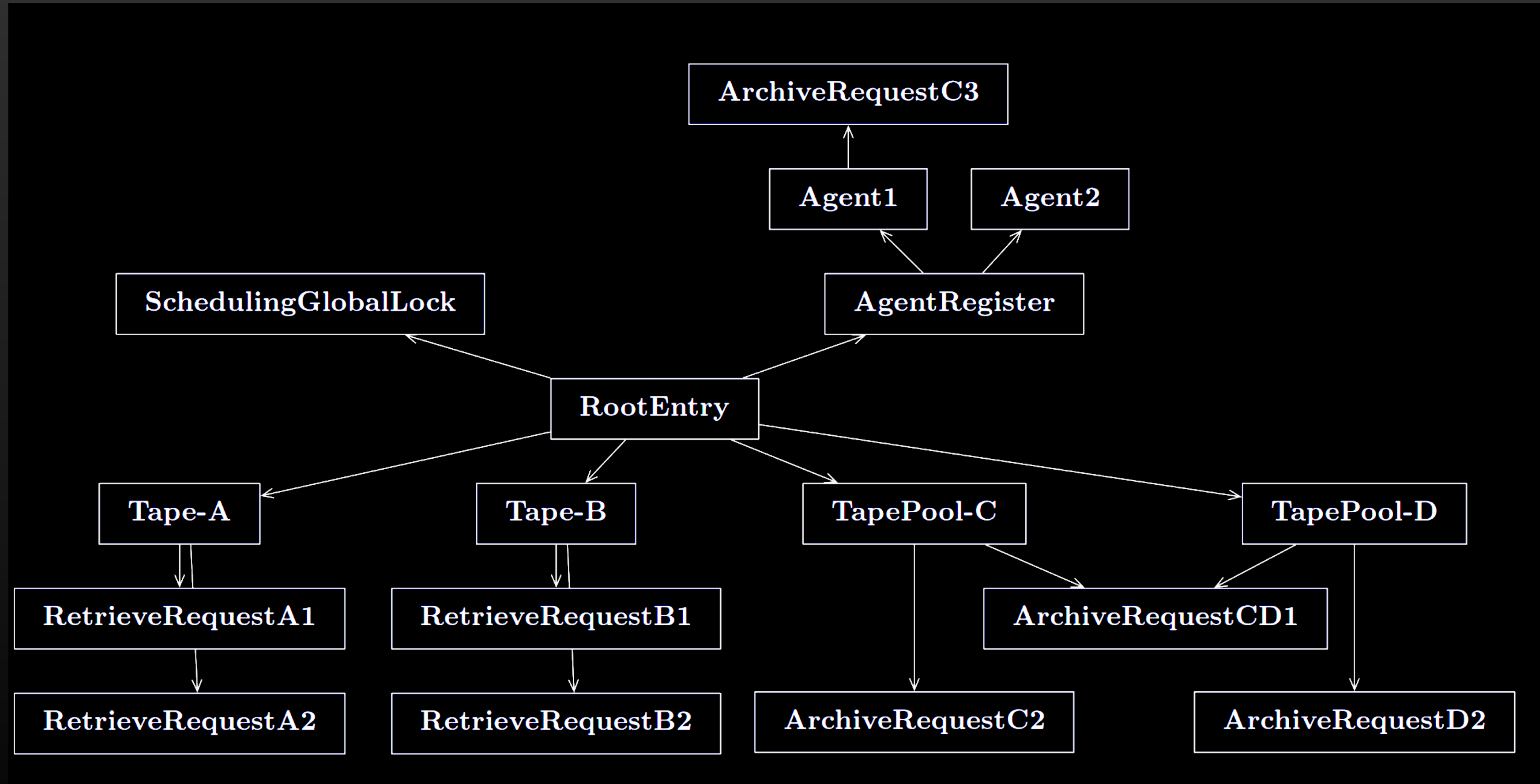  - 11 PB on LTO-7M

# New Scheduler DB: Relational DB

## Challenges

- **independent solutions for ObjectStore-coupled scheduler logic**

- **high performance, reliability and scalability IF**

  - ▸ **DB features exploited smartly !**

    **(e.g. not counting all rows for every query)**

  - ▸ **requires optimisation efforts per use-case**

  - ▸ **relies on diligence of developer with DB queries and DB configuration**

- **LHC Run 4 !**

# Scheduler DB Implemented as ObjectStore

## Architecture

# IO Limited Sanity Check

1 thread queueing

2 threads writing to tape