

I/O Load Scheduler for Grid Mass Storage

Christos Tziortzios

Christos.Tziortzios@os3.nl

Introduction

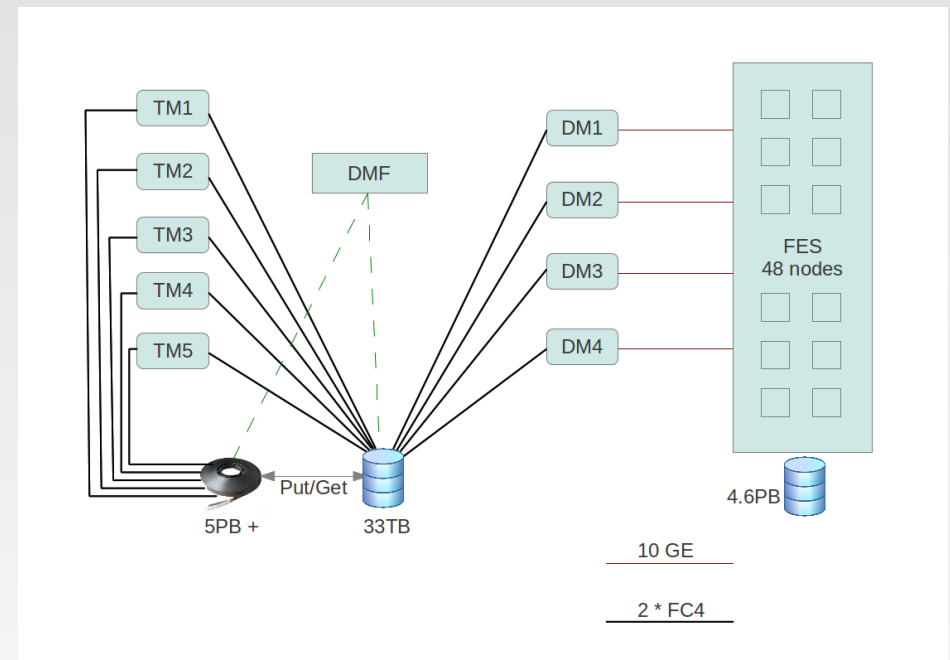
- SARA manages enormous amounts of data produced by CERN (LHC), LOFAR and more
- More than 5 PB stored on tapes at the moment
- Hierarchical Storage Management
 - Disk front end
 - Tape back end

Research Question

Is it possible to use an intelligent scheduling mechanism in order to control the data flow between the Front End Storage and Grid Mass Storage more efficiently?

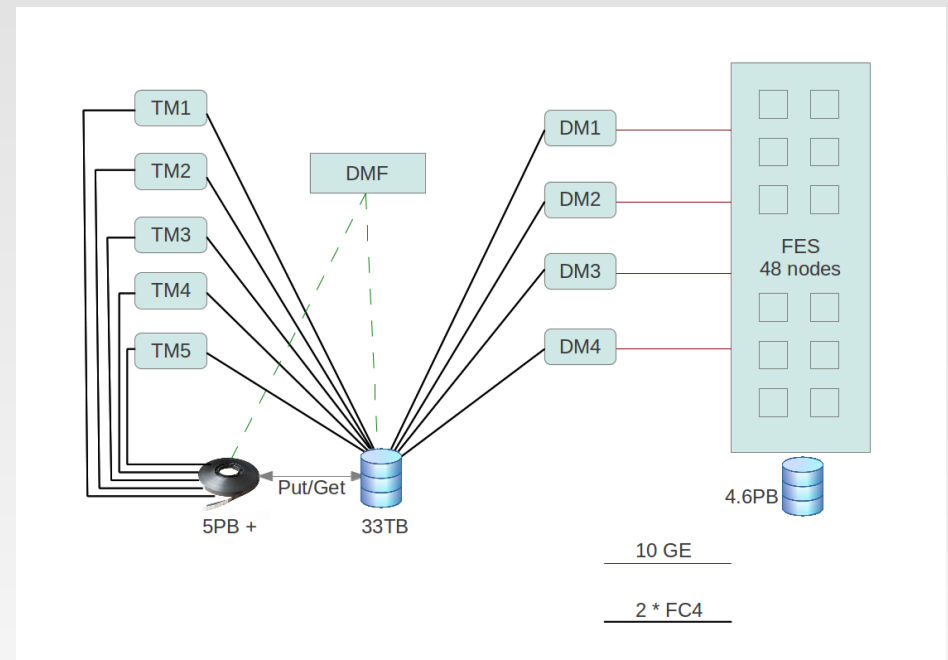
Infrastructure

- Front End Storage
 - 48 Nodes
- GridMS
 - 4 Data Movers (DM)
 - 5 Tape Movers (TM)
 - 20 Tape Drives
 - 33 TB disk
 - Data Migration Facility (DMF) takes care of put and get operations

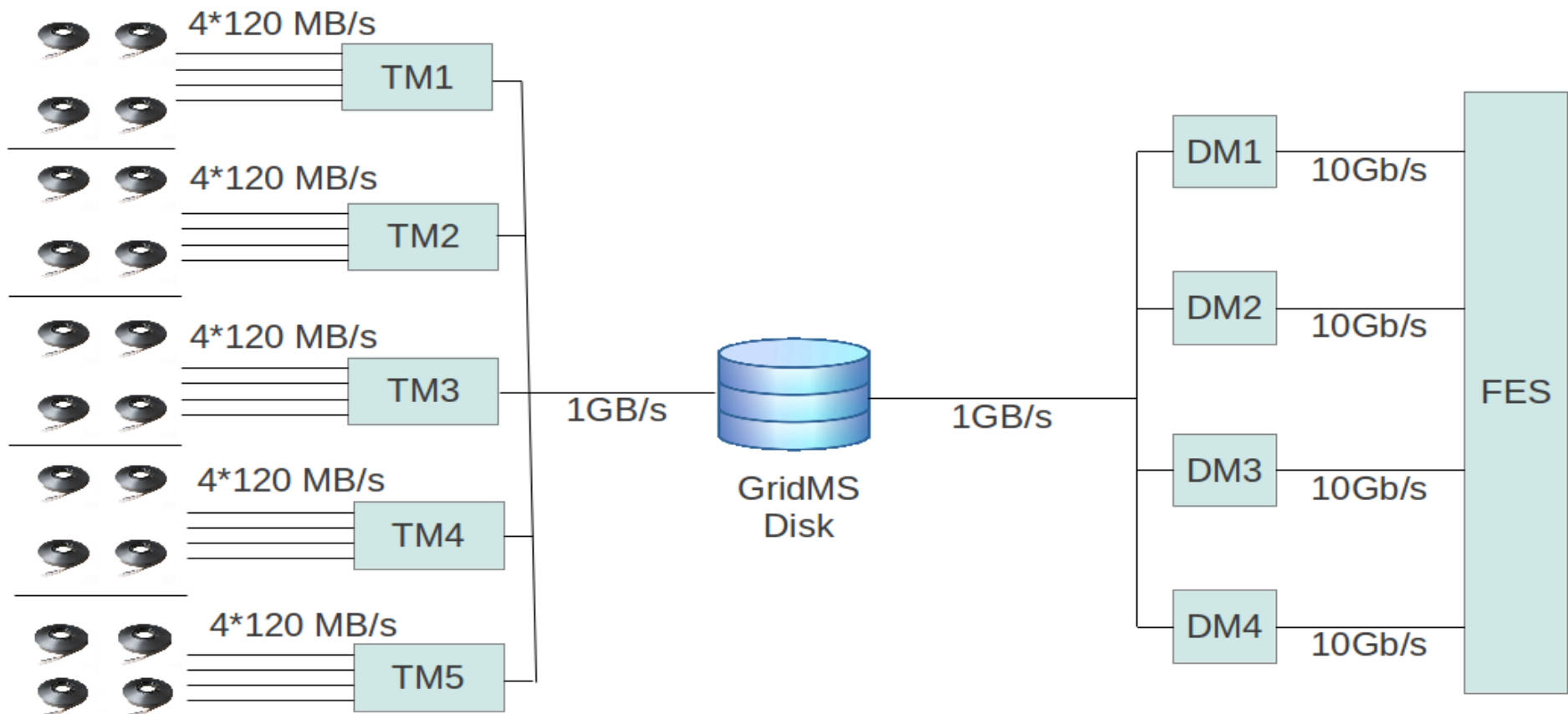


Performance Issues

- Random I/O leads to drop in performance.
- No job scheduling on groups of FES Nodes or User level.
- Only one transfer per FES node at a time, may lead to idle bandwidth
- Limited disk bandwidth



Disk Bandwidth Problem



Tape Drives

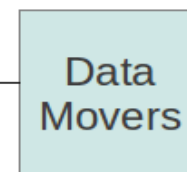
19 Gbps



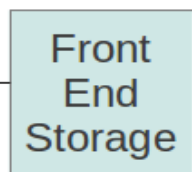
8 Gbps



8 Gbps



40 Gbps



Operations

- Operations between FES and GridMS (handled by our scheduler)
 - Store
 - Restore
 - Checksums (Both in FES and GridMS disk)
- Operations between GridMS disk and Tape (handled by Data Migration Facility)
 - Put
 - Get

Software Used

- TORQUE resource manager
 - Normally gives processes access to CPU time or memory
 - We are interested in disk I/O and bandwidth
- Maui Cluster Scheduler
 - Scheduling and Fairshare options

Tests and Results (1)

- No test environment
- Store and Restore jobs first submitted to the queue
- Successfully checked Priority and Fairshare Components
 - Priority depending on User
 - Fairshare based on short term historical data
- Maui overrides TORQUE priorities
- Different Maui and TORQUE configurations tested
 - Node allocation

Tests and Results (2)

- Requesting resources
 - Walltime: predicted by user.
 - Disk space: only works for one filesystem, SARA plans to have multiple, one filesystem for each project

Tradeoff: Accurate requests for resources increase efficiency - underestimating resources may lead to killing jobs

Conclusions

- Implemented a prototype solution for store and restore operations.
 - Advanced Scheduling.
 - Idle bandwidth would no longer be a problem.
- Disk space resource would work with the current infrastructure but not with multiple file systems.
- Current scheme works reliably. Changes in the working environment may introduce bugs.
- Reliability: Testing environment needed.

Questions