# Observations made while running a multi-petabyte storage system

## IEEE MSST 2010
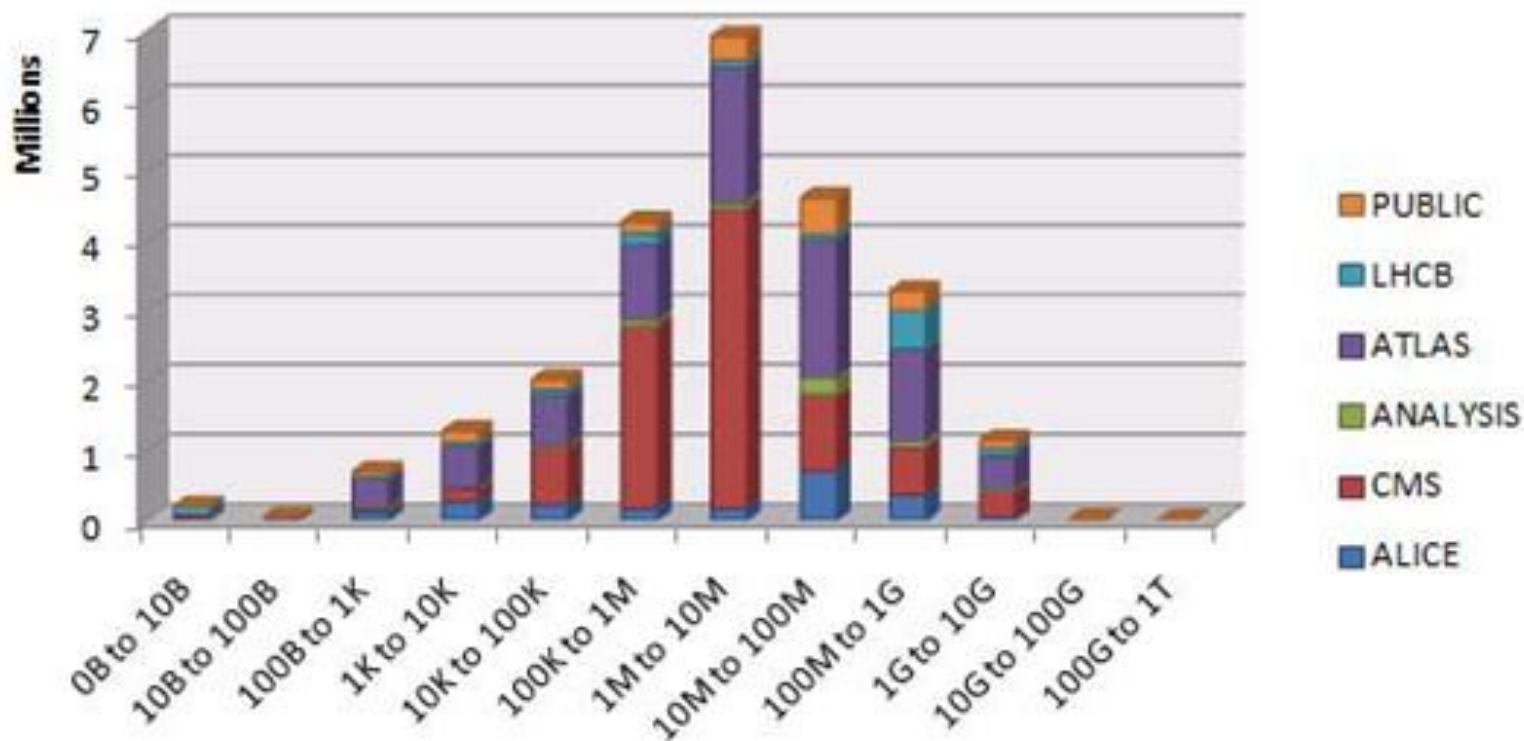## Lake Tahoe

*Miguel Coelho dos Santos,*

*Dennis Waldron*

*CERN*

CERN IT Department
CH-1211 Genève 23
Switzerland
**www.cern.ch/it**

DSS

- 25PB of data (10PB on disk, 20PB on tape)
- 145M files
- 1.300 disk servers
- 28.000 disk drives
- 130 tape drives
- 7 tape libraries

*Observations made while running a multi-petabyte storage system - 2*

DSS

CERN IT Department

| Disk Server Downtime | 3% |
|---|---|
| Disk Drive Failure | 87 failures/month ($2.9\times10^{-6}$ fail / hour op.) |
| Disk Array Failure | 0.9 failures/month ($1.8\times10^{-7}$ fail / hour op.) |
| Tape Failure | 1.7 failures/month ($5.4\times10^{-8}$ fail / hour op.) |

Large distributed storage systems need to handle component failures automatically, particularly regarding data availability and metadata consistency.

- The environment:
  - Continuous need to move data (and associated metadata) so hardware can be serviced or replaced (end of warranty)
  - High media access times both on commodity disks (random access) and tapes in general
  - Media capacity increases not matched with similar bandwidth increases
  - No significant change in file sizes
- Metadata and data operations growing with capacity increase
- Operations per media device increasing

DSS

- Measuring and monitoring is CRITICAL

- Log messages generate high rates of messages/second. Joins are problematic.

- Monitoring based on summaries collected by different components.

- Allow information to be plugged into standard tools for visualization and long term storage.

- Must be part of the product from the start

# Questions?

*Observations made while running a multi-petabyte storage system - 7*