Data Challenges and Fabric Architecture

General Fabric Layout



10/22/2002

Bernd Panzer-Steindel, CERN/IT

Benchmark, performance and testbed clusters (LCG prototype resources)

- → computing data challenges, technology challenges, online tests, EDG testbeds, preparations for the LCG-1 production system, complexity tests 'surplus' resource (Lxshare) for running experiments and physics production (integration into Lxbatch) Requests are mainly in number of nodes, not Si2000
- \rightarrow 400 CPU server, 100 disk server, ~250000 Si2000, ~47 TB

Main fabric cluster (Lxbatch/Lxplus resources)

- → physics production for all experiments Requests are made in units of Si2000
- → 650 CPU server, 130 disk server, ~ 350000 Si2000, ~ 82 TB

Certification and Service nodes

 \rightarrow ~60 CPU server, ~5 disk server

Level of complexity

Physical and logical coupling



Current CERN Fabrics architecture is based on :

- In general on commodity components
- Dual Intel processor PC hardware for CPU, disk and tape Server
- Hierarchical Ethernet (100, 1000, 10000) network topology
- NAS disk server with EIDE disk arrays
- RedHat Linux Operating system
- Medium end tape drive (linear) technology
- OpenSource software for storage (CASTOR, OpenAFS)



Status check of the components - CPU server + Linux-

Nodes in centre : ~700 nodes running batch jobs at ~ 65% cpu utilization during the last 6 month

Stability :7 reboots per day0.7 Hardware interventions per day (mostly IBM disk problems)

→ Average job length ~ 2.3 h, 3 jobs per nodes == Loss rate is 0.3 %

Problems : Level of automatization (configuration, etc.)

Status check of the components - Network -

Network in the computer center :

- 3COM and Enterasys equipment
- 14 routers
- 147 switches (Fast Ethernet and Gigabit)
- 3268 ports
- 2116 connections

Stability : 29 interventions in 6 month (resets, hardware failure, software bugs,etc.)

Traffic : constant load of several 100 MB/s, no overload

Future : tests with 10 GB routers and switches have started, still some stability problems

Problems : load balancing over several Gb lines is not efficient (<2/3), only matters for current computing Data Challenges

LCG Testbed Structure used for e.g. Benchmarks

GigaBit Gigabit Ethernet Fast Ethernet

100 cpu servers on GE, 300 on FE, 100 disk servers on GE (~50TB), 20 tape server on GE



Bernd Panzer-Steindel, CERN/IT

Aggregate disk server Network traffic



10/22/2002

Bernd Panzer-Steindel, CERN/IT

Disk stress tests :

30 servers with 232 disks running for 30 days I/O tests (multiple streams per disk, random+sequential) ~ 3 PB
→ 4 disk server crashes and one disk problem (~> 160000 disk MTBF)
(IBM disk problems last year)

Stability :

About 1 reboot per week (out of ~200 disk servers in production)and ~one disk error per week (out of ~3000 disks in production)

Disk server tests : 66 CPU server (Gigabit) , 600 concurrent read/write streams into 33 disk server for several days → 500 MB/s write + 500 MB/s read Limited by network setup + load balancing

Disk server : dual PIII 1 GHz, 1 GB memory, Gigabit, ~ 500 GB mirrored Performance : ~ 45 MB/s read/write

aggregate, multiple streams, network to/from disks

STORAGE BANDWIDTH - NUMBER OF STREAMS - BLOCK SIZE



 \rightarrow improved memory bandwidth

Problems : still high performance fluctuations between Linux kernel versions, Linux I/O still room for improvements **Block Size[bytes]**

Status check of the components - Tape system -

Installed in the center :

- Main workhorse = 28 STK 9940A drives (60 GB Cassettes)
- ~12 MB/s uncompressed, on average 8 MB/s (overhead)
- Mounting rate = ~45000 tapes per week

Stability :

- About one intervention per week on one drive
- About 1 tape with recoverable problems per 2 weeks(to be send to STK HQ)

Future :

- New drives successfully tested (200 GB, 30 MB/s) 9940B
- 20 drives end of October for the LCG prototype
- Upgrade of the 28 9940A to model B beginning of next year

Problems : 'coupling' of disk and tape server to achieve max. performance of tape drives

Utilization of Tape drives



Data Challenges with different Experiments



Mixture of hardware (disk server) and software (CASTOR,OBJ,POOL) optimization

Bernd Panzer-Steindel, CERN/IT

Status check of the components - HSM system -

Castor HSM system : currently \sim 7 million files with \sim 1.3 PB of data

Tests:

- ALICE MDC III with 85 MB/s (120 MB/s peak) into CASTOR for one week
- Lots of small tests to test scalability issues with file access/creation, nameserver Access, etc.
- ALICE MDC IV 50 CPU server and 20 disk server at 350 MB/s onto disk (no tapes yet)

Future :

Scheduled Mock Data Challenges for ALICE to stress the CASTOR system

Nov 2002	200 MB/s (peak 300 MB/s)
2003	300 MB/s
2004	450 MB/s
2005	750 MB/s

Conclusions

- Architecture verification okay so far more work needed
- Stability and performance of commodity equipment satisfactory
- Analysis model of the LHC experiments is crucial
- Major 'stress' (I/O) on the systems is coming from Computing DCs and currently running experiments, not the LHC physics productions

Remark : Things are driven by the market, not the pure technology → paradigm changes