

LCE Reports

French Atomic Energy Commission (CEA)

Stéphane Thiell stephane.thiell@cea.fr



CEA Computing Center

CEA/DAM Computing Center



LUG 2007

April 2007

Current supercomputers at CEA/DIF



- o Linux Kernel 2.6
- o Lustre 1.4.7 and Lustre 1.6 Beta 7
- o Quadrics and Infiniband interconnects

- o Linux Kernel 2.6
- o Lustre 1.4.8
- o Infiniband interconnect

TERA-10 Computing Center Architecture



Two Lustre architectures





- o Use of device-mapper multipath
- o~60 clients

HPC Cluster File System

- Lustre 1.4.7 and 1.4.8 with vendor support (Bull)
- Quadrics Elan4 or InfiniBand interconnect
 - Native Quadrics or OpenFabrics LND (qsnet/o2ib)
 - Network is dedicated to a cluster
- Dual attached DDN 9550 with fibre channel disks (WT, 8+1+1)
 - 🖛 16 LUN of 1 TB
 - 2.3 GB/s per DDN couplet
 - OSS are 16 or 8-cores Itanium servers
- Performance oriented
 - 100 GB/s on checkpoint/restart like benchmark
 - Single client performance: 2.2 GB/s Write, 1.4 GB/s Read
- TERA-10 day production is 30 TB



HPC Lustre Scalability





- Shared File System
 - Lustre 1.6 Beta 7
 - InfiniBand interconnect
 - Native OpenFabrics LND (o2ib)
 - Network is shared by several clusters
 - Dual attached DDN 9550 with SATA disks (WT, 8+2)
 - 🖛 48 LUN of 8 TB
 - ▶ 1.5 GB/s per DDN couplet
 - OSS are 4-cores Xeon servers
 - 1 FS shared by several clusters
 - Capacity oriented
 - Multi Petabytes FS (target is 2+ PB)
 - Single client performance: 472 MB/s Write, 371 MB/s Read

Initially created with 2 DDN couplets = 2 * 330 TB

Few weeks ago extended in half a day with 2 DDN

- No need to reformat
- Now at 1.3 PB in one FS

💙 root@cupidon7:~ - cupidon -	Konsole <2>					
cprot00-0ST00b4_UUID	7.2T	215.8G	7.OT	2% /cea/cache_pr	ot[0ST:180]	*
cprot00-0ST00b5_UUID	7.2T	227.3G	6.9T	3% /cea/cache_pr	ot[0ST:181]	
cprot00-0ST00b6_UUID	7.2T	228.8G	6.9T	3% /cea/cache_pr	ot[0ST:182]	
cprot00-0ST00b7_UUID	7.2T	228.3G	6.9T	3% /cea/cache_pr	ot[0ST:183]	
filesystem summary:	1.3P	184.9T	1.1P	14% /cea/cache_pr	ot	
[root@cupidon7 ~]#						
[root@cupidon7 ~]#						
[root@cupidon7 ~]# df /cea/cache_prot						
Filesystem 1K-blocks Used Available Use% Mounted on						
@o2ib: @o2ib:/cprot00						
1415155704608 198563072040 1216592369544 15% /cea/cache_prot						
[root@cupidon7 ~]#						
[root@cupidon7 ~]#		-				
[root@cupidon7 ~]# df -h /cea/cache_prot						
Filesystem Size Used Avail Use% Mounted on						
@o2ib: @o2ib:/cprot00						
	1.3P 185T	1.2P 15	% /cea/cac	he_prot		2
[root@cupidon7 ~]#						•

Shared InfiniBand Network Topology





2. Device-mapper multipath I/O failover (Linux DM MPIO)

- EMC path prioritizer on the MGS/MDS nodes
- **clariion** path checker for EMC
- Home made DDN path prioritizer for OSS nodes
- tur path checker for DDN





Lustre HSM Project Update

Goal: Add ILM features in Lustre (target 2008)

- ILM stands for Information Lifecycle Management
- HSM stands for Hierarchical Storage Management
- Migration inside Lustre or between Lustre and an external backend storage
- Designed to allow interoperability with multiple external backends, including HPSS
 - Use of userspace commands to access external backend system
- All files are always visible in the file system, but file data can reside:
 - On the primary storage (Lustre file system)
 - On the backend storage (internal or external storage system)
 - On both
- Metadata (size, ...) are always up-to-date

Involve the migration of file system objects

Migration enables multiple Lustre features (HSM, caches for Lustre proxy services, space rebalancing, LAID rebuild, etc.)

• Working at a FID granularity level

- MDT FID (full file)
- OST FID (file object)
- File access by FID feature
 - FID is used as the reference key in the backend storage system
 - Lustre namespace is independent from backend namespace

• File system space management will be either:

- Automatic
 - 📂 at OST level
 - ► at FS level (MDT)
- On-demand
 - eg. based on a provided list of files to purge (e2scan)

Purge method

- Keep start/end of FID data using the punch() call
- At OST level (ext3/ldiskfs)
- At FS level (lustre)



Questions ?