
Experiences with 10 Months HP SFS / Lustre in HPC Production

Roland Laifer

**Computing Centre (SSCK)
University of Karlsruhe**

Laifer@rz.uni-karlsruhe.de

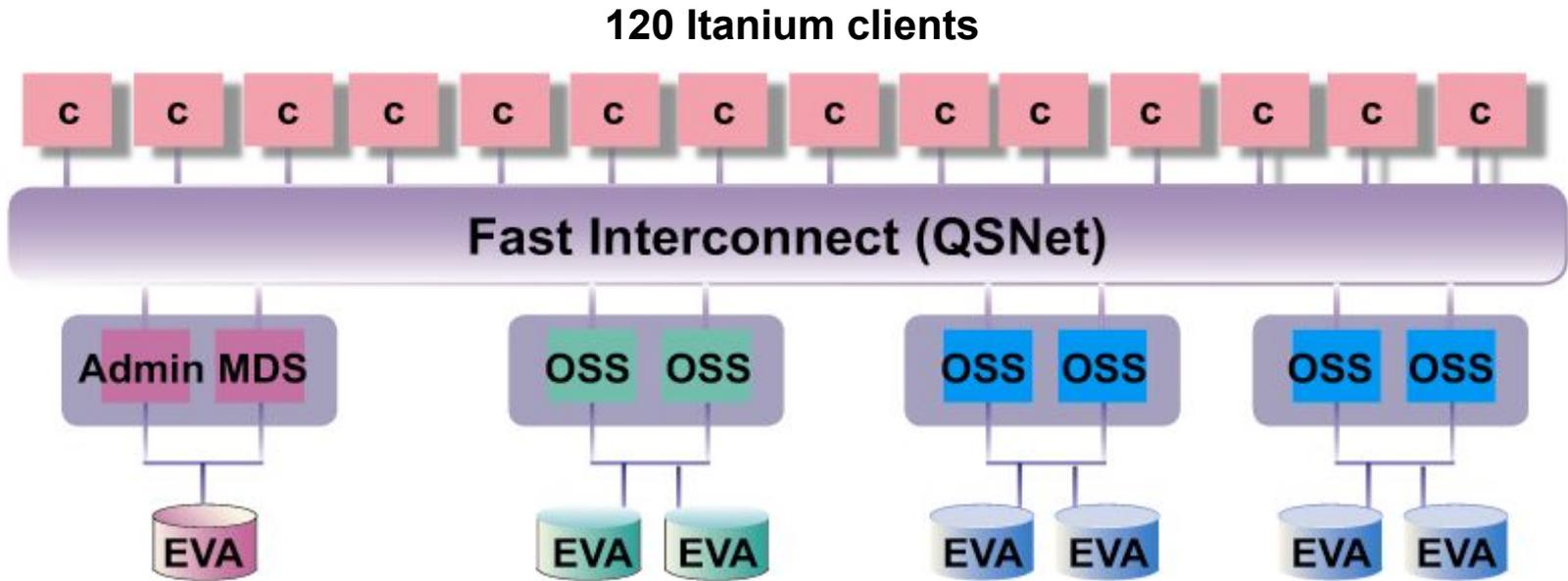


Outline

- » **Status**
- » **Problems**
- » **Positive aspects**
- » **Troubleshooting**
- » **Performance of OSS components**



HP SFS on SSCK's HP XC6000



MDS and Admin for
\$HOME and \$WORK

- Allows > 50 million files

\$HOME

- 3.8 TB storage

\$WORK

- 7.6 TB storage

Legend

Admin: Administration Server
MDS: Metadata Server
OSS: Object Storage Server
EVA: EVA5000 storage array
C: Client



Status of HP SFS at SSCK (1)

- » **Production started in January 2005**
 - **HP SFS was used for home directories from the beginning**
 - This was a good decision!

- » **HP SFS software version is still 1.1-1**
 - **Over time, 5 patches were installed**
 - see following slides
 - **SFS version 2.0 was not installed because of XC version 2.1 dependency**
 - Waiting for XC version 3 in order to do migration install from XC version 2.0
 - Different SFS versions on client and server are currently not supported

- » **A lot of storage hardware had to be replaced**
 - **15 FC disks and 3 EVA controllers**
 - This is much more than expected by the MTBF



Status of HP SFS at SSCK (2)

- » **Interconnect adapter on servers is single point of failure**
 - **Heartbeat between server pairs uses shared storage**
 - Room for improvement in future SFS version

- » **Both file systems were heavily used**
 - **One user group was doing most of the IO**
 - Parts of their applications now use local disks
 - This means more bandwidth is available for other users

- » **Problem management system at www.itrc.hp.com is bad**
 - **Updates are *not* sorted by date**
 - **Information is manually copied between 2 different systems**
 - QuIX and WFM
 - This causes delay or information to be lost

- » **HP SFS is pretty stable and stability has improved**



Problems with HP SFS version 1.1-1 (1)

- » **MDS stopped working and had to be manually rebooted**
 - **Bugs in LDLM locking caused all ll_mdt threads to become hung**
 - Fixed with patch SFSV1.1-1-153-26-70

- » **OSS dumped frequently**
 - **LBUG in filter_grant_sanity_check()**
 - Fixed with patch SFSV1.1-1-1000233568

- » **Filesystem hangs because server lost Quadrics connection**
 - **Installed new Quadrics driver to further investigate problem**
 - Problem happened 3 times

- » **File system did not start after reboot of server pair**
 - **Bugs in hpls-clumanager and hpls-db**
 - Fixed with patch SFSV1.1-1-1000240369



Problems with HP SFS version 1.1-1 (2)

» HP SFS logging did not work

- Lots of harmless error messages filled up eventlog
 - Fixed with patch SFSV1.1-1-1000243131
- 2 GB limit of eventlog

» Sometimes EVA controller failover did not work

- This caused IO errors in applications
- New FC drivers were installed
 - Fixed with patch SFSV1.1-1-1000279279-290319
- After storage problems OST luns might be mounted read-only

» Patch installation decreased HP SFS throughput by 30 %

- Problem is still under investigation



Positive aspects

- » **Pretty long periods without major problems**
 - **Until summer one month, now up to 3 months**
 - **OSS failover usually after few weeks**
 - **Failover usually worked**
 - **Only 3 of the 7 problems were really critical**
 - **MDS hangs, Quadrics adapter failure, and failing EVA controller failover**

- » **Most problems simply caused file system and applications to hang**
 - **Time limit of some batch jobs had to be extended**
 - **Most times only \$WORK file system was affected**

- » **High level support was good**
 - **Was able to fix all main problems**



Troubleshooting HP SFS (1)

» Check for SFS server errors during the last 4 days

- `evlview -m -f 'age < 4' | grep -i LustreError`
 - If no errors are reported everything works fine
 - Otherwise error messages can be critical or insignificant

» Check if SFS services are balanced

- `echo "show server" | sfsmgr`

» Check status of services

- `echo "show filesystem" | sfsmgr`

» Check if enough local disk space is available on south4

- `ssh south4 df /var/log/dump`
 - Filesystem `/local` on OSS is hidden

» Check for Lustre error messages on a XC client

- `grep -i LustreError /var/log/`hostname``



Troubleshooting HP SFS (2)

» Check on client if file system throughput is normal

- `time dd if=/dev/zero of=testfile bs=1M count=10000`

» Explanation of some Lustre error messages:

- Sep 14 01:35:54 xc1-ls3-adm kernel: LustreError: 4147:0:
(filter_io_24.c:172:filter_direct_io()) **short write?** expected 524288, wrote -5
 - **Writes are failing, i.e. you might have a critical problem with your storage**
- Sep 14 01:35:57 xc1-ls3-adm kernel: Lustre: 4144:0:
(debug.c:244:portals_run_upcall()) Invoked portals upcall /
usr/opt/hpls/bin/hpls_upcall **LBUG**,filter_io_24.c,filter_commitrw_write,343
 - **An LBUG triggers a failover, i.e. the server will probably go down**
- Oct 27 08:40:51 xc1-ls7-adm kernel: LustreError: 3348:0:
(ldlm_lockd.c:318:ldlm_failed_ast()) **### blocking AST failed (-107):** evicting
client c7774_data_e03790ea4e@NET_0x166_UUID NID 0x166 (0:358) ...
 - **Client with Quadrics node id 358 fails to relinquish a lock in time**
 - **Possible reason is client shutdown without unmounting cleanly**



Performance of OSS components

» Quadrics Elan4

- Internally about 1300 MB/s
- Only PCI-X adapters exist

» PCI-X bus on servers

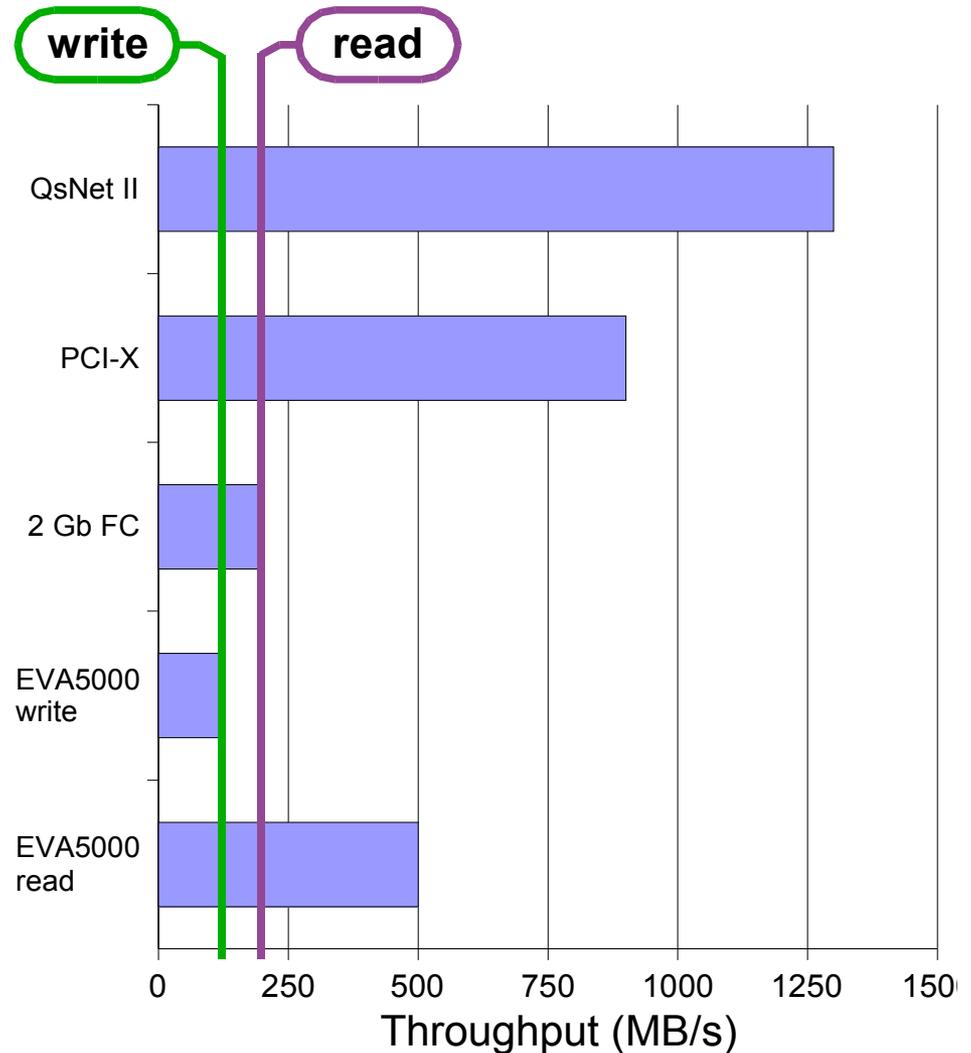
- About 900 MB/s

» Dual-ported FC adapter

- About 195 MB/s
- Actually only 1 port is used

» EVA5000 storage array

- About 120 MB/s for writes
- Nearly 500 MB/s for reads



Summary

- » **It's not easy to understand Lustre messages**
 - **Better documentation of error messages is necessary**

- » **Advantages of HP SFS compared to free Lustre:**
 - **Additional software for failover and ease of administration**
 - **Good support**

- » **Performance is sufficient for our current applications**
 - **Performance monitoring can help to improve applications**
 - **It can also help to identify system problems**

- » **Stability of HP SFS was good and seems to be improving**
 - **We had expected to run into more problems**

- » **We can recommend HP SFS !**

