# Latest Production Experiences with HP SFS

**Roland Laifer**

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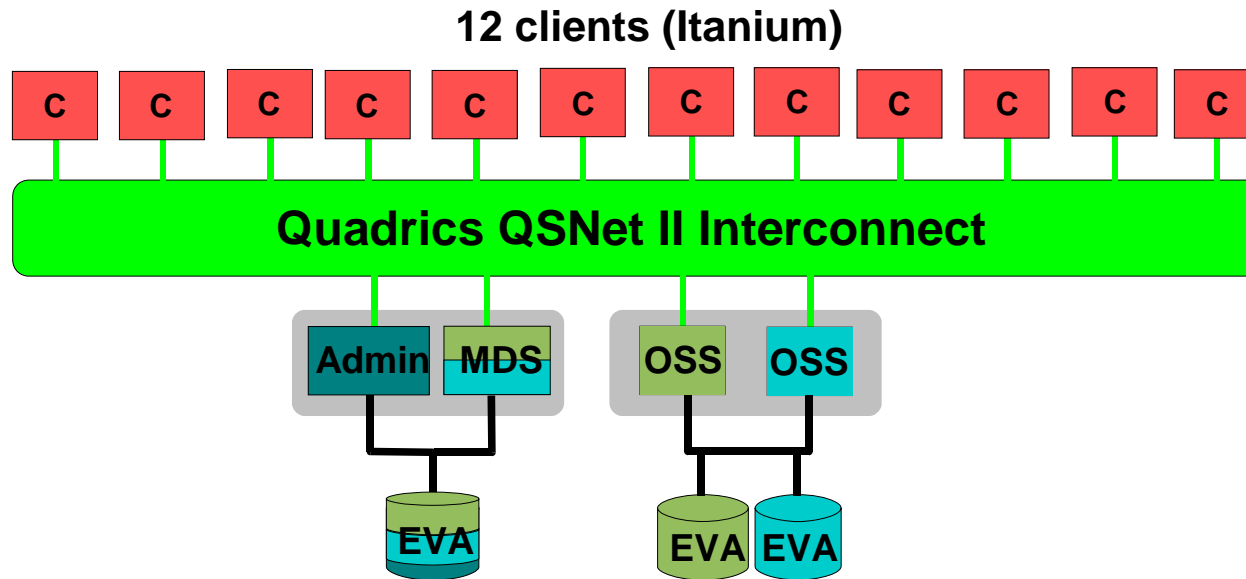**University of Karlsruhe**
**Germany**

**Laifer@rz.uni-karlsruhe.de**

# Outline

» **Description of SSCK's 4 HP SFS systems**

» **Performance graphs**

» **HP SFS versus open source Lustre**

» **Configuration decisions for our new SFS system**

» **Some not fully solved problems**

» **Operational experiences**
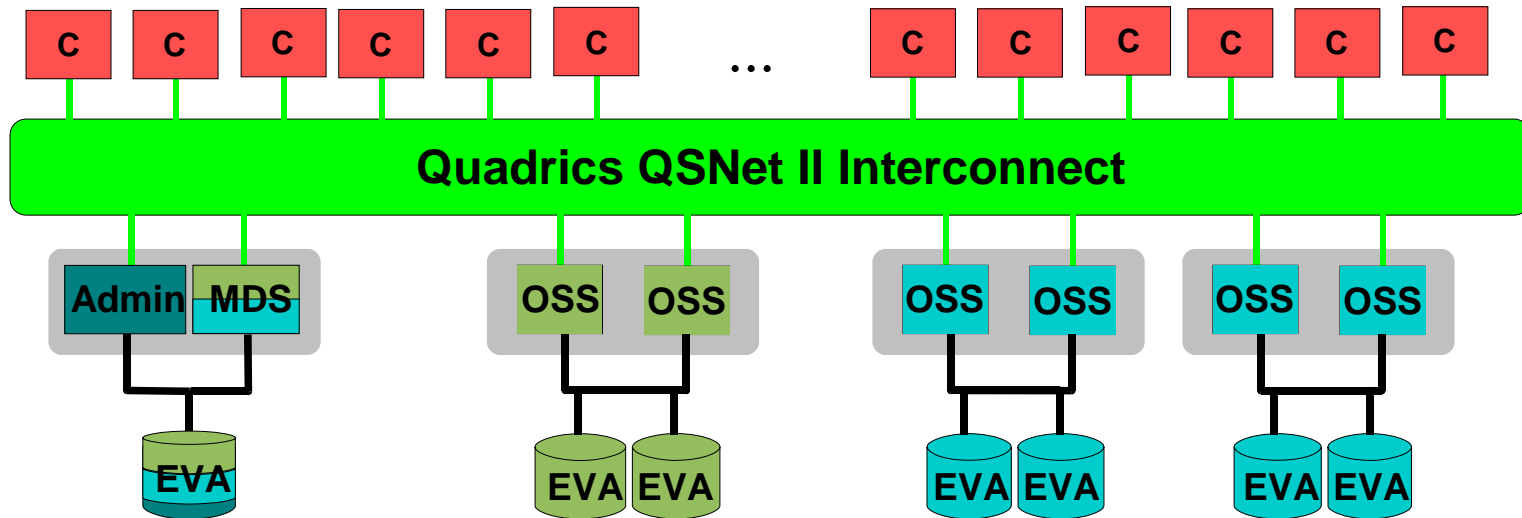
» **Future plans**

# Itanium test system (xc0)

**12 clients (Itanium)**

| C | C | C | C | C | C | C | C | C | C | C | C |

**Quadrics QSNet II Interconnect**

**Admin** **MDS** **OSS** **OSS**

**EVA** **EVA** **EVA**

|  | **$HOME** | **$WORK** |
|---|---|---|
| **Capacity** | 0.5 TB | 0.5 TB |
| **Write performance** | 120 MB/s | 120 MB/s |
| **Read performance** | 190 MB/s | 190 MB/s |

Universität Karlsruhe (TH)
**Rechenzentrum**

Roland Laifer

# Itanium production system (xc1)

**120 clients (Itanium)**

| C | C | C | C | C | C | … | C | C | C | C | C | C |

**Quadrics QSNet II Interconnect**

**Admin** **MDS**     **OSS** **OSS**     **OSS** **OSS**     **OSS** **OSS**

**EVA**     **EVA** **EVA**     **EVA** **EVA**     **EVA** **EVA**

|  | $HOME | $WORK |
|---|---|---|
| **Capacity** | 3.8 TB | 7.6 TB |
| **Write performance** | 220 MB/s | 380 MB/s |
| **Read performance** | 340 MB/s | 580 MB/s |

**Notes:**

- **Performance is reduced by fragmentation**
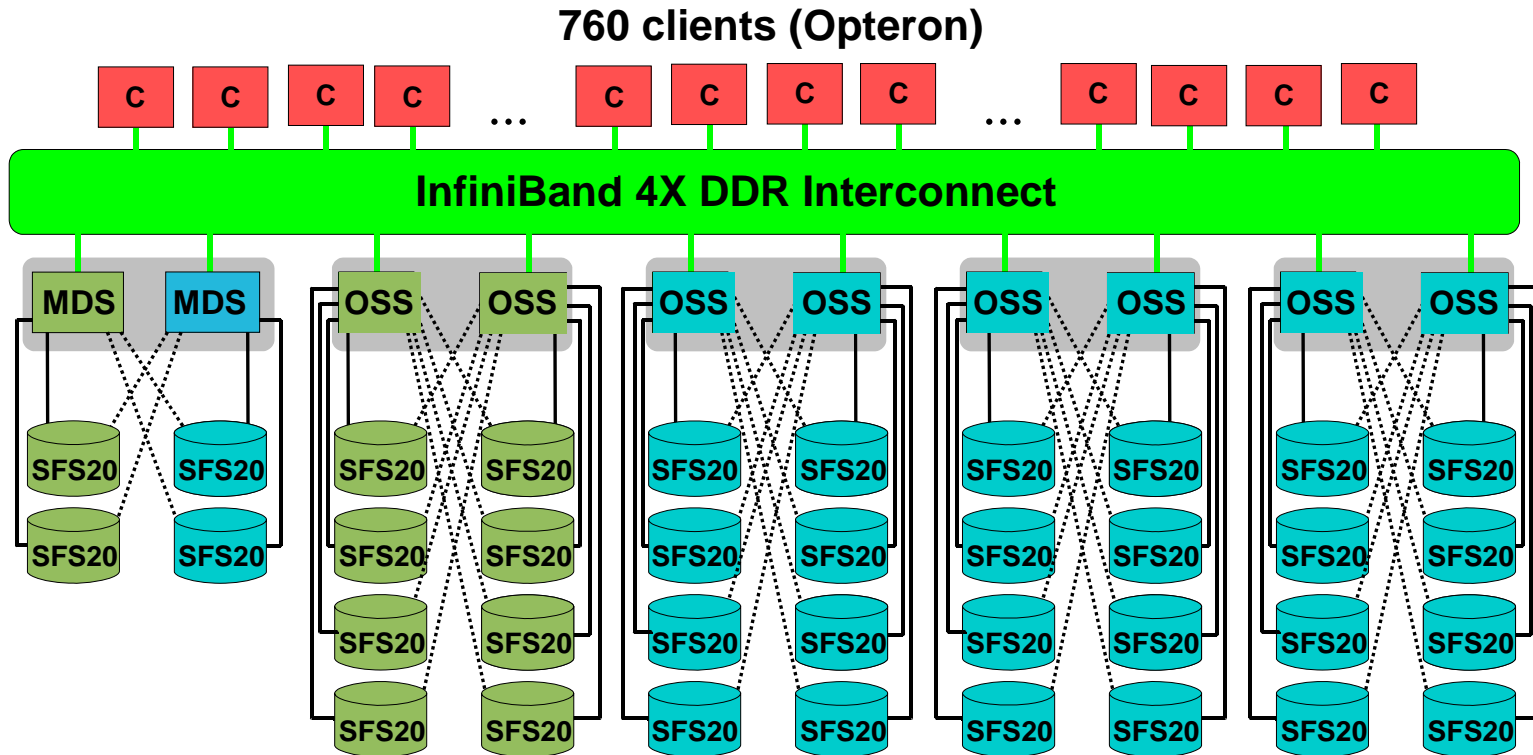- **Higher fragmentation of $WORK**

# Opteron test system (xc3)

**12 clients (Opteron)**

| | C | C | C | C | C | C | C | C | C | C | C | C |

**InfiniBand 4X DDR Interconnect**

**MDS**  **MDS**    **OSS**  **OSS**

**SFS20**  **SFS20**    **SFS20**  **SFS20**

**SFS20**  **SFS20**

| | $HOME | $WORK |
|---|---|---|
| **Capacity** | 2 TB | 4 TB |
| **Write performance** | 90 MB/s | 180 MB/s |
| **Read performance** | 150 MB/s | 300 MB/s |

**Notes:**

- **$HOME file system uses mirrored OST luns**
- **SFS20s use RAID ADG**

# Opteron production system (xc2)

**760 clients (Opteron)**

| C | C | C | C | … | C | C | C | C | … | C | C | C | C |

**InfiniBand 4X DDR Interconnect**

| MDS | MDS | | OSS | OSS | | OSS | OSS | | OSS | OSS | | OSS | OSS |

SFS20 SFS20 SFS20 SFS20 SFS20 SFS20 SFS20 SFS20 SFS20 SFS20

SFS20 SFS20 SFS20 SFS20 SFS20 SFS20 SFS20 SFS20 SFS20 SFS20

SFS20 SFS20 SFS20 SFS20 SFS20 SFS20 SFS20 SFS20

SFS20 SFS20 SFS20 SFS20 SFS20 SFS20 SFS20 SFS20
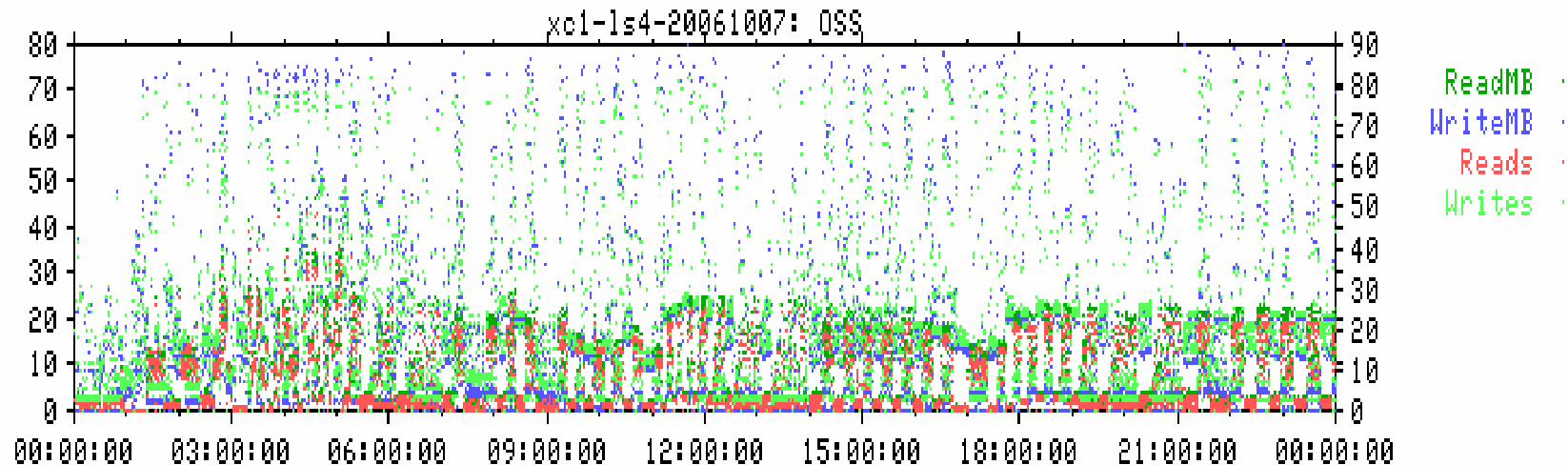
|                   | $HOME     | $WORK      |
|-------------------|-----------|------------|
| Capacity          | 8 TB      | 48 TB      |
| Write performance | 360 MB/s  | 1850 MB/s  |
| Read performance  | 600 MB/s  | 3000 MB/s  |

**Notes:**

- **$HOME file system uses mirrored OST luns**
- **Preliminary results for $WORK: was only tested once**

Universität Karlsruhe (TH)
**Rechenzentrum**

Roland Laifer

# Performance graph for one OSS of xc1



xc1-ls4-20061007: OSS

» **Applications with high I/O load:**

- **Computer algebra application**
  - **Could create output files in TB range**

- **Applications doing scratch I/O on each task**
  - **Capacity of local disk is not sufficient**

- **ABAQUS**

# HP SFS versus open source Lustre

» **HP SFS**

    – **Easy installation, configuration and upgrade**

    – **Additional software for failover, management and client build**

    – **Additional tests and patches to supply hardened Lustre version**

    – **Very good support**

    – **System health check, SFS log database and email alerts**

    – **Performance monitoring**

    – **Good documentation**

» **Open source Lustre**

    – **Flexibility in choice of server and storage hardware**

        • **Hard job to find appropriate storage, good drivers and firmware levels**

    – **Flexibility to use newest software versions**

        • **Possible impact on stability**

    – **No license costs**

# Configuration decisions for our SFS system on xc2

» **Default stripe size of 4**

- **Wanted to have very good performance from a single node**
  - **I/O is often done from a single task of a large parallel job**

- **Offers best load distribution on $HOME (4 OSTs)**

- **Metadata performance with stripe size 1 is not much better**

» **Use RAID ADG (RAID6)**

- **With huge storage capacity high risk to loose data with RAID5**

- **Moderate performance reduction (10% for writes)**

- **No capacity reduction with 250 GB disks and fully populated SFS20s**

» **On SFS20 use rebuild_priority=*medium***

- **Performance is much better during rebuild than with default**
  - **26 MB/s versus 4 MB/s when using rebuild_priority=*high***

- **Rebuild time is not extensively higher than with default**
  - **12 hours versus 5.5 hours when using rebuild_priority=*high***

# Configuration decisions for SFS on xc2 (continued)

» **Use OST lun mirroring for file system $HOME**

- **Broken SFS20 controller would normally not hang up the file system**
  - **This is not true if service lun is located on the broken SFS20**

- **Possibly break the mirror if the capacity is no longer sufficient**
  - **Solution without restoring the data is theoretically possible**

» **Distribute the MDS services of the 2 file systems**

- **Load distribution to Admin and MDS node**

- **Makes the file systems independent of each other**

# Some not fully solved problems

» **Fragmentation reduced performance by 10 to 30%**

  – **Fix needs recreation of file systems**

  – **Risk is reduced on newer systems because of ext3 extents**

» **Many broken FC disks**

  – **Rate is much higher if I/O load on system is high**

  – **Number of broken disks was lower during last months**

» **SFS20 with service lun is single point of failure**

  – **Creates extreme load on Admin node and stops complete system**
    • **This problem is under investigation**

  – **Mirroring service luns would be a good enhancement**

# Operational experiences

» **Only one complete outage during last 10 months**

- **Both OSS crashed permanently**
    - **Started after broken EVA controller was repaired**
    - **Reason: LAST_ID was not incremented while objects were created**
    - **Fix needed file system check**
    - **Delete dumps if hidden file system /local is full**
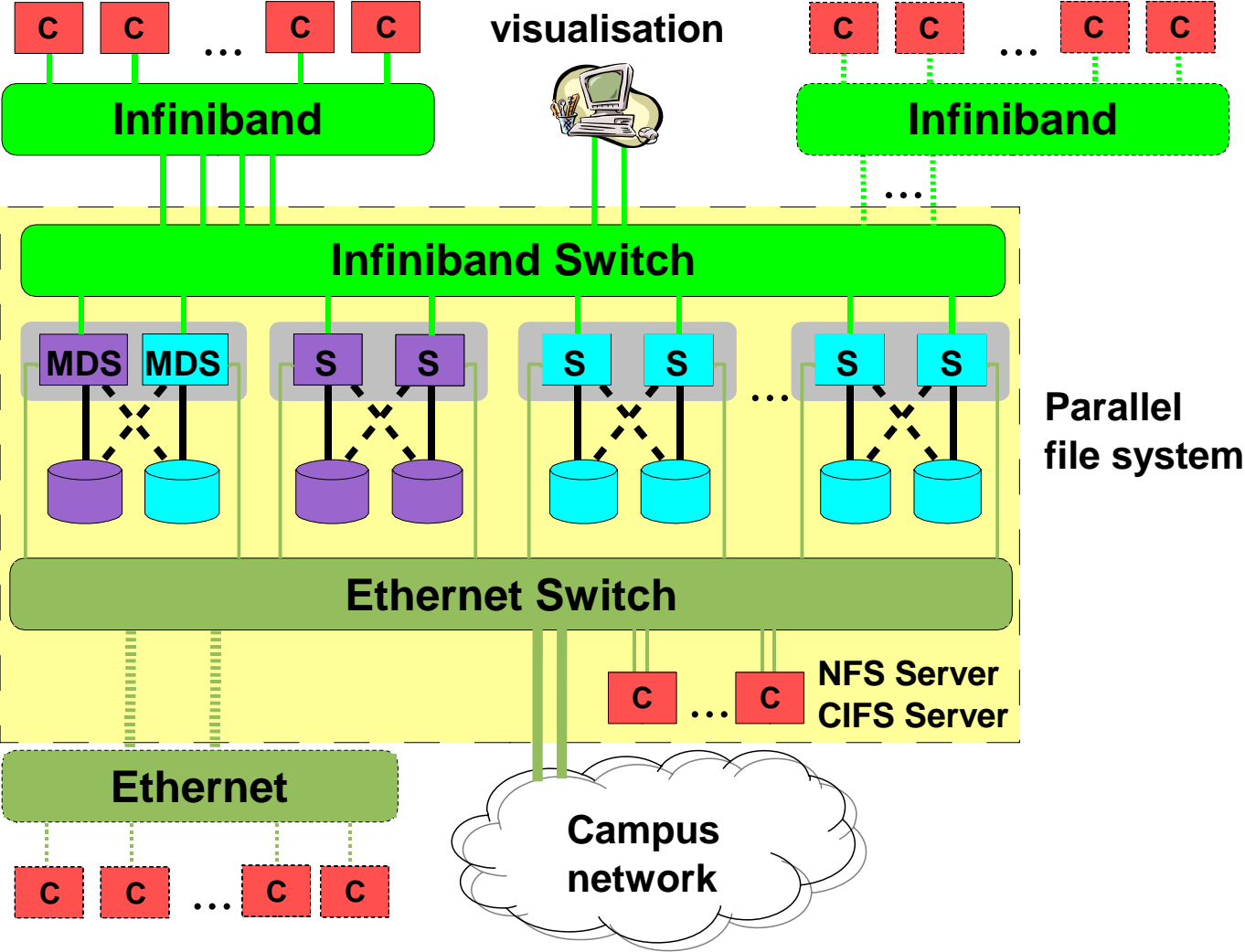
» **Administrative challenge to identify critical errors**

- **LustreError on client and server might indicate a critical issue**
    - **Lots of error messages which are not really critical**
- **Use syscheck to check the system's health**

» **New applications sometimes create new errors**

- **E.g. MPI-IO test program causes lots of errors on clients**
- **Some error messages appear when high load is created**

» **Collectl performance monitoring on client to identify critical users**

# Future plan for a central parallel file system

# Additional requirements for central parallel file system

» **Version compatibility**

  – **Upgrade of all clients together with servers is not reasonable**

» **Reduced kernel and distribution dependency**

  – **Support for more kernels and distributions is required**

  – **Patchless client might help**

» **User level security**

  – **Need to export file systems with high performance to untrusted clients**

  – **Kerberos security should provide this feature**
    • **Was unfortunately delayed several times**

» **Server system upgrade while file systems are online**

  – **File systems should have no downtime**

  – **This could be possible by upgrading servers in failover mode**

# Summary

» **Lustre provides a scalable and stable parallel file system**

» **HP SFS supplies additional features**

    – **which make it a real product**

» **Some non-default configuration settings could be useful**

» **Further experiences with HP SFS:**

    – **http://www.rz.uni-karlsruhe.de/dienste/lustretalks**