# Lustre@GSI
# - a Petabyte Filesystem

## Walter Schön, GSI

# Topics

- **Architecture and Hardware**
- **Cross site Lustre connection**
- **Managing 1 Pbyte of Data**
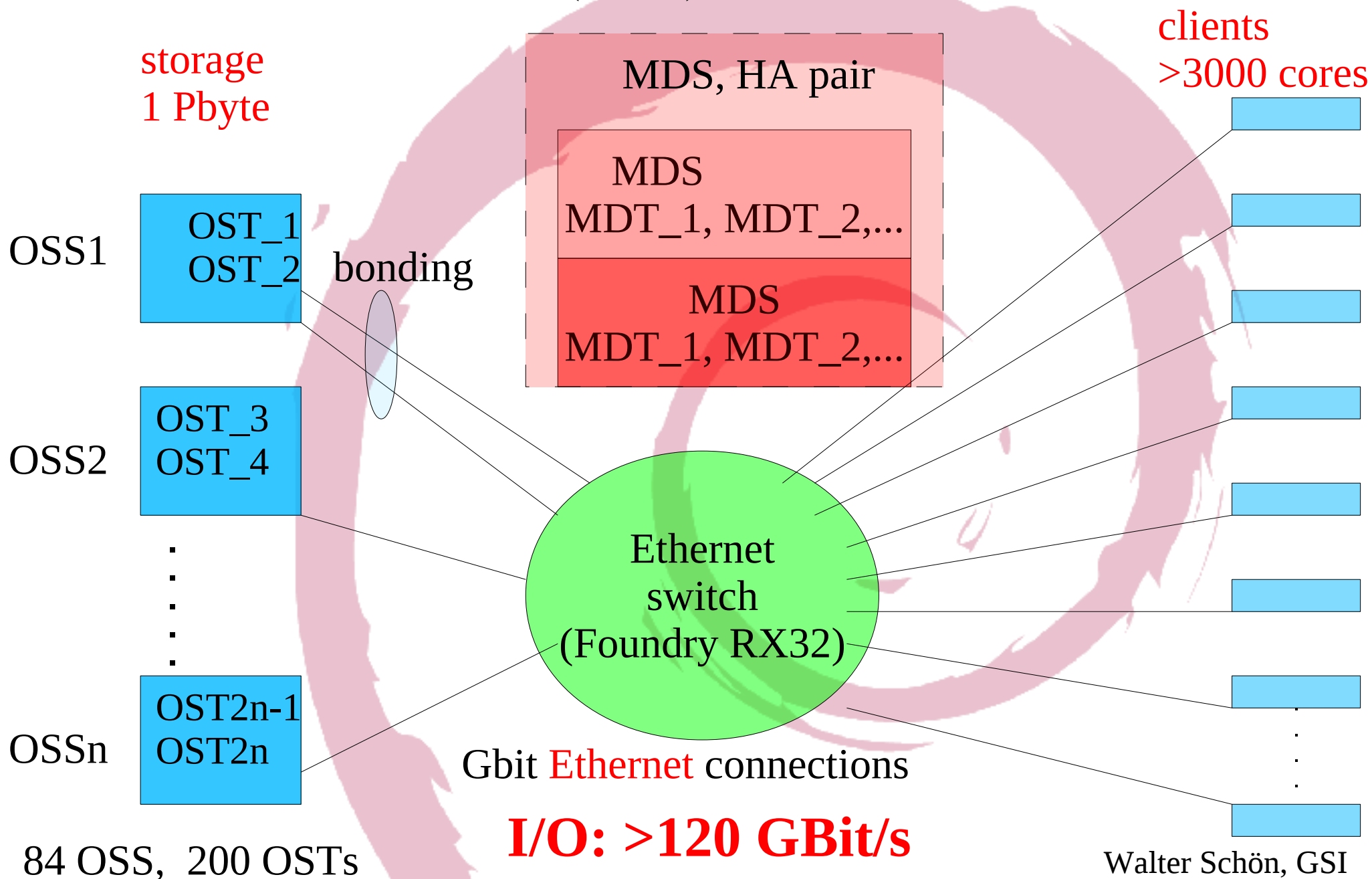- **The Dark Side of  Lustre**
- **Outlook**

# lustre@GSI:

## Online storage for

- Alice Analysis (Tier2)
- GSI Experiments
- Theory Groups
- FAIR Simulations

Walter Schön, GSI

# Lustre Cluster Architecture

Lustre 1.6.7.2 debian,  2.6.22 (server), 2.6.28 clients

storage
1 Pbyte

clients
>3000 cores

MDS, HA pair

MDS
MDT_1, MDT_2,...

MDS
MDT_1, MDT_2,...

OSS1

OST_1
OST_2

bonding

OSS2

OST_3
OST_4

OSSn

OST2n-1
OST2n

Ethernet
switch
(Foundry RX32)

Gbit Ethernet connections

**I/O: >120 GBit/s**

84 OSS,  200 OSTs

Walter Schön, GSI

# Lustre@GSI: Online Storage for Experiments and Theory Groups

Number of MDS: --------------- 3 – HA Pair, one standby
Number of **OSSs**: -------------- 84
Number of **OSTs**: ------------- 200
Number of disks: -------------- 1600 ( including RAID5+6, spare)
Size: ---------------- 1 Pbyte
Number of clients (cores): -- > 3000
Max. number of files: --------- $2.5 * 10^8$
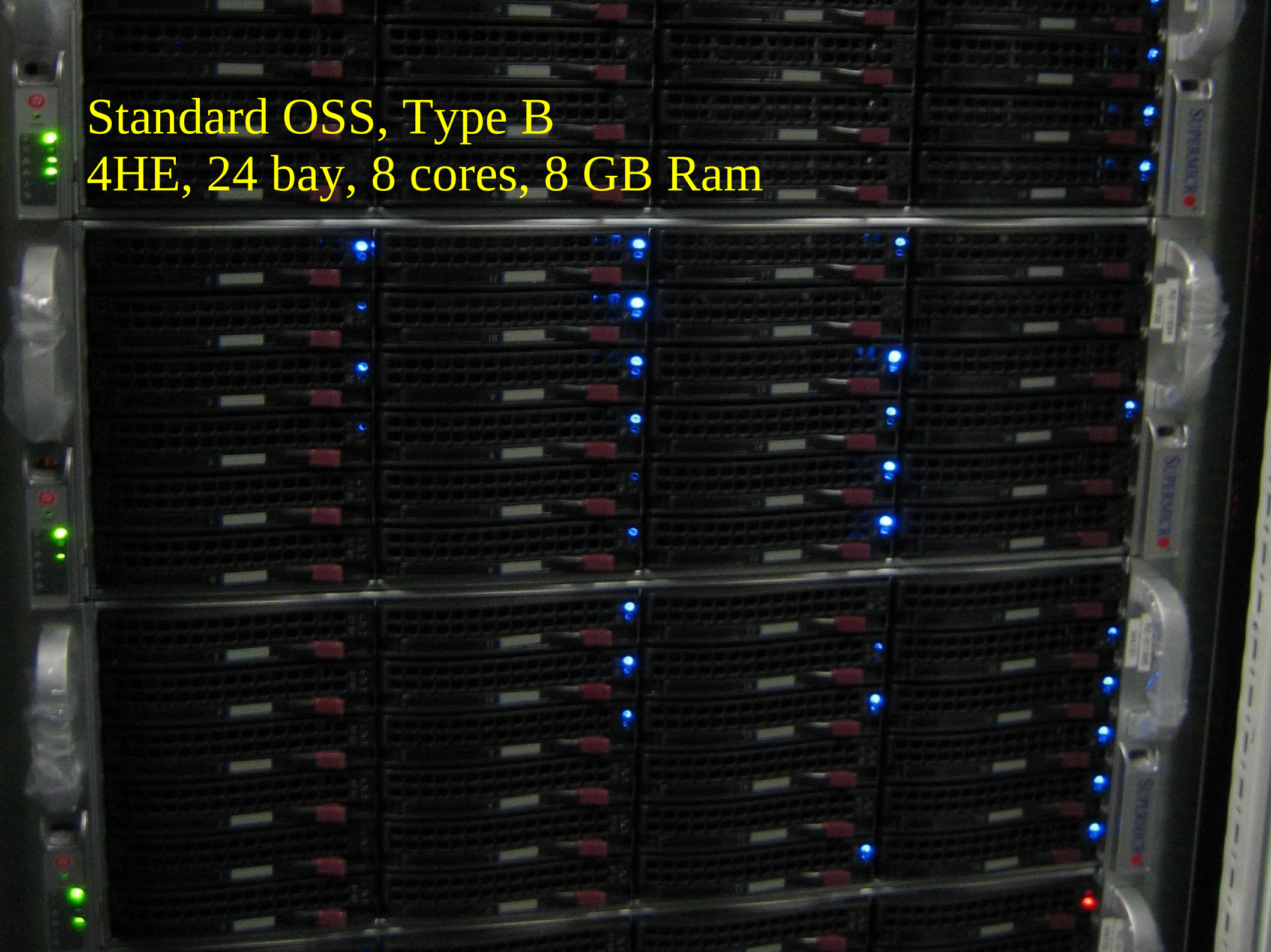Files in lustre: --------- $5 * 10^7$

OSTs : B: 1 TB, WD green line, RAID 6 + spare
A: 0.5 TB, WD RAID 5
MDT : WD "raptor" 150 GB, 14 disks in RAID 10

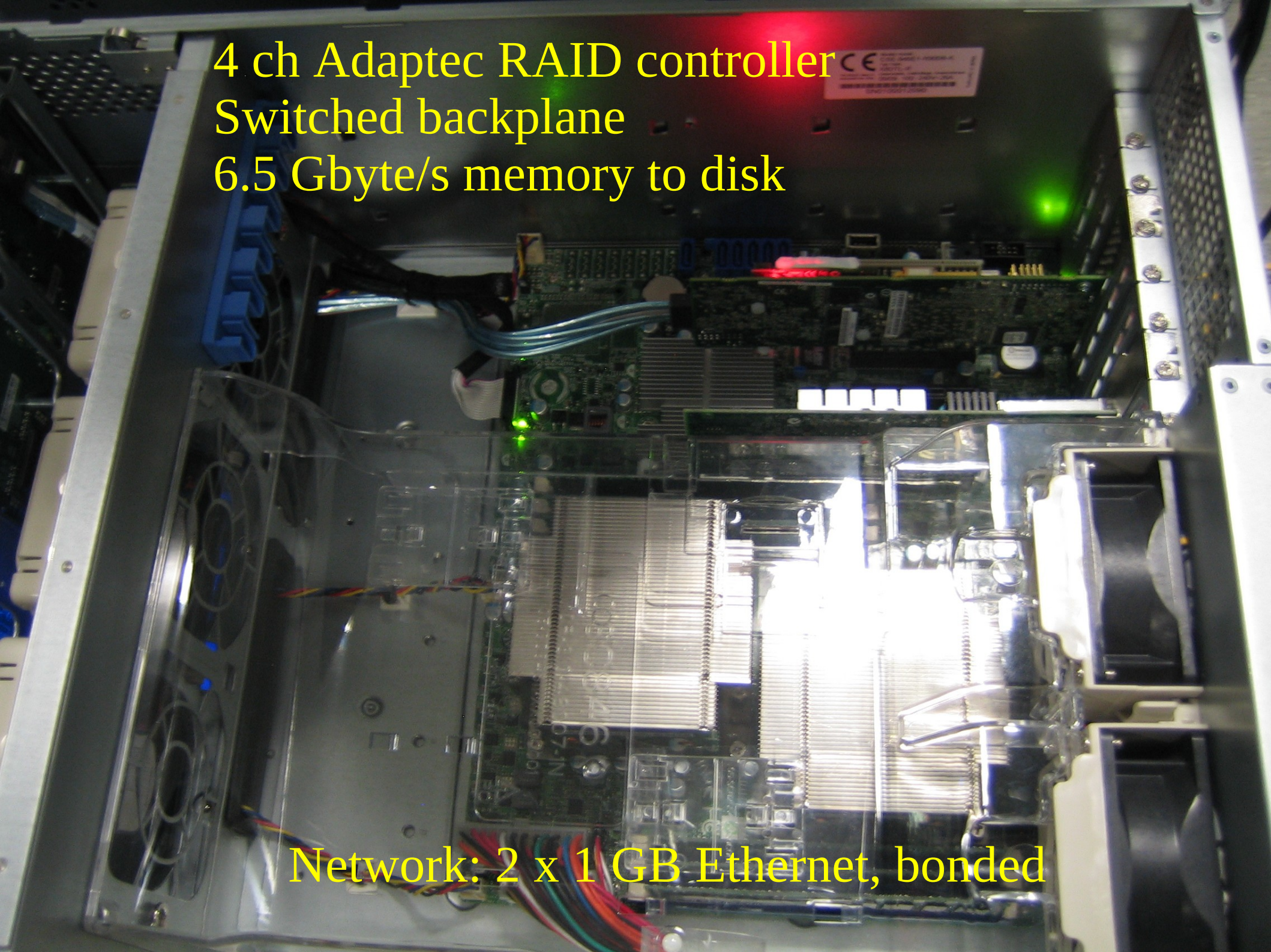MDS: 8 core Xeon 3 GHz, 32 GB Ram

Walter Schön, GSI

Standard OSS, Type B
4HE, 24 bay, 8 cores, 8 GB Ram

4 ch Adaptec RAID controller
Switched backplane
6.5 Gbyte/s memory to disk

Network: 2 x 1 GB Ethernet, bonded

# I/O Throughput of a real Analysis - ALICE

**Aggregate network connections: 120 GBit/s**

**Measurement of the Alice Analysis train I/O trouput:** **>50 Gbit/s** **using 2000 cores.**

**Total number of cores now: 3000 cores**
**Dezember 2009: 4000 cores**
**=> aggregate I/O limit can be reached**

Walter Schön, GSI

# Hardware Failures in 1 Year of Operation

· Disks: about 12 disks of 1600 damaged during operation

·      No double hit => no data loss, no break

· File server: 1 OSS (out of 90) damaged ( in the first hours of operation)

·      Operation downtime 1h

· RAID-Controller (200)

    1 battery status: failed, no break

·   1 damaged (MDS) => HA fail over, 40 min break

·   1 "bad stripe" event on one OST out of 200, few files are affected

=> bug under "rare conditions" in Adaptec firmware, patched now

·

**=> 3 hardware related unscheduled downtimes (1h in total)**
**in one year operation**

New Hardware with redundant RAID controllers for the MDS in testing mode

Walter Schön, GSI

# Software Faillures in one Year of Operation

- A: "Flat Network Situation"
  => many incidents, dying servers, confused lustre, dying linux
  => lustre down each 10 days, restart of MDS necessary
- B: after transfer in dedicated "HPC" network in summer
  => no more incidents ........... ( 120 days of operation  now ..... )

A:
- "strange"  packets on the mds interface,
- broadcast storms on the mds interface
B:
- HPC clients und lustre servers are in an own network segment (VLAN):
=> **clean packets, no more broadcast storms,**
   **no more problems in operation**
   **120 days of operation – 1 scheduled downtime 2h ( demonstration)**

Walter Schön, GSI

# The Dark Side of Lustre

...**lustre is more a formula 1 racing car than a "Volkswagen"....**
- **Complex system**
- **Vulnerable to (network) communication problems**

**Annoying lustre bug in our setup:**
- **Quota is not working after upgrade**
 **=> Bug report opened, not solved yet**

Walter Schön, GSI

# Lustre Cross Site Connection

Since September  cross site connection to 100 TFlop "scout cluster"
=> testing a cross site model for future data analysis
Remote compute power in the region with direct access to
GSI online mass storage

**First Results:**
**Testing with single and multiple lustre clients:**
**" just working", no problems yet.**
**Present connection: 1 Gbit now, 4 Gbit soon.**

Walter Schön, GSI

# Lustre  Management

**Dynamic expansion of the FS space in a production system**
· Expansion from 0.7 Pbyte to 1 Pbyte : successful, no break necessary


**Audit of a large FS**
In a Pbyte fs with $10^8$ entries simple questions like:
- List of  top users?
- List of top files?
- List of top groups?  etc.
- File space used by group "xyz" ?
 ...can be very boring/time consuming using traditional unix tools...
…  **and performed by 2000 users can be DoS Problem for the MDS !**

Walter Schön, GSI

# Lustre Management

**Audit Alternatives?**

**=> Robinhood Filesystem Monitor:**
 Audit and purge tool for large file Systems, advanced Capabilities for lustre
http://sourceforge.net/projects/robinhood, developed by CEA
- Parallel threats on clients reporting results to a central mySQL DB
  => no (small) stress for the MDS!
- Lustre capabilities only for lustre 1.8 and 2.x
- operated at GSI on lustre 1.6.x without special lustre features .... testing

**First Results:**
- **Fast ( parallel threats)**
- **Low noise on the MDS**

Walter Schön, GSI

# Outlook

**December**
- Increasing production lustre size to 1.4 Petabyte in December
- Increasing I/O Bandwidth to 150 Gbit/s
- Increasing number of lustre clients to 4000 (cores)
- Increasing lustre cross site connection to 4 Gbit/s

**2010:**
- Increasing lustre to 2.5 Petabyte
- Increasing I/O Bandwidth to 300 ? Gbit/s
- Increasing number of lustre clients (cores) >> 4000 .......
- Upgrade production cluster to 1.8.x series
- Introducing 10GB Ethernet and/or Infiniband

Walter Schön, GSI