



# Experiences with Lustre in Bern

Gianfranco Sciacca

Laboratory for High Energy Physics, University of Bern, Switzerland

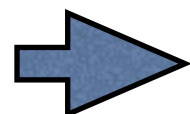
Parallel file systems for HPC

Community HPC-CH

Thursday 28 October 2010, ETH Zurich



# Lustre systems in Bern

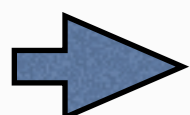


## UBELIX cluster - Central Informatik Diensten

~1100 cores

SunFire X2200 and DELL PowerEdge 1950 dual-quad core with Gentoo

~44TB on Lustre (home areas and scratch)



## LHEP ATLAS T3 cluster - LHEP

~200 cores

SunFire X2200 dual-quad core with ROCKS and CentOS

~11TB on Lustre\* (scratch only)

\* experimental setup



Both clusters expected to grow in size in the near future



# Lustre systems in Bern

Both clusters are embedded in National and International infrastructures

<http://www.nordugrid.org/monitor/> (>55k CPUs)

<http://giis.smsg.ch> (~7.3k CPUs)

Country	Site	CPU	Local	Grid	Queueing
Denmark	LSCF (NBI)	4	0+0	0+0	
Denmark	Morpheus (NBI)	5	0+0	0+0	
Denmark	Steno (DCSC/KU)	3184	392+1899	700+451	
Denmark	Steno Tier 3 (DCSC/KU)	3184	168+2138	738+0	
Finland	Ametisti (M-grid)	260	1+52 (queue inactive)	0+68	
Finland	Jade	768	35+482	0+4148	
Finland	Kiniini (CSC)	56	0+0		
Finland	Korundi (M-grid, HIP)	400	48+61		
Finland	Kvartsi (M-grid)	192	0+64		
Finland	Liuske (CSC test)	8	0+0 (queue inactive)		
Finland	Murska	2176	0+1689		
Finland	Opaali (M-grid)	88	0+28 (queue inactive)		
Finland	Pythia (M-grid)	80	0+3		
Finland	Topaasi (M-grid)	82	0+17 (queue inactive)		
Germany	Uni Lübeck - INB	16	0+6 (queue inactive)		
Iceland	Jotunn (Uol)	164	0+4		
Norway	EPF (UiO/FI)	20	0+5		
Norway	fimm (BCCS/UiB)	832	0+58		
Norway	hexagon (BCCS/UiB)	5552	0+5896		
Norway	stallo (HPC/UiT)	5600	0+3345		
Norway	Titan A (UiO/USIT)	4624	549+2619		
Norway	Titan B (UiO/USIT)	4624	0+3167		
Slovenia	Arnes	252	252+0		
Slovenia	SiGNET	1140	0+1142		
Slovenia	SiGNET	1140	1098+44		
Sweden	Ada (C3SE)	1000	0+664		
Sweden	Beda (C3SE)	1744	0+1632		
Sweden	Grad (SweGrid, Uppmax)	512	466+31		
Sweden	ISV	3	0+1		
Sweden	Neolith	6440	0+6888 (queue inactive)		
Sweden	Ritsem (SweGrid, HPC2>	536	486+1		
Sweden	Ruth (SweGrid, PDC)	134	134+0 (queue inactive)		
Sweden	Siri (SweGrid, Lunarc)	328	316+19		
Sweden	Smokerings (NSC)	496	215+278		
Sweden	Svea (SweGrid, C3SE)	512	390+67		
Switzerland	Bern ATLAS T3	212	136+0		
Switzerland	Bern UBELIX T3 Cluster	1072	296+258		
Switzerland	Geneva ATLAS T3	222	109+15		
Switzerland	Manno PHOENIX T2	1520	264+438		
Switzerland	Manno PHOENIX T2	1520	270+867		
UK	UKI-SCOTGRID-GLASGOW	1984	0+1407		
Ukraine	BITP Cluster	96	0+0		
Ukraine	IAP Cluster	24	0+26		
Ukraine	ICBGE Cluster	32	0+191		
Ukraine	ICYB SCIT-3	828	0+35		
Ukraine	ILTPE Cluster	88			
Ukraine	IMBG Cluster	24			
Ukraine	IMMSP Cluster	24			
Ukraine	IMP Cluster	72			
Ukraine	Inparcom Cluster	192			
Ukraine	IOP Cluster	80			

Country	Site	CPU	Local	Grid	Queueing
	arcXWCH at HEPIA, Gen>	612	0+0		0+0
	Bern ATLAS T3	212	136+0		175+0
	Bern UBELIX T3 Cluster	1072	296+258		38+31
	GC3 Grid Cluster	312	0+0		0+0
	Geneva ATLAS T3	222	105+17		76+0
Switzerland	Manno PHOENIX T2	1520	258+868		104+274
Switzerland	Manno PHOENIX T2	1520	250+417		155+262
	OCI Grid Cluster	59	0+0		0+0
	SMSCG - Vital-IT	1064	0+421		0+279
	USI-ICS Cluster	328	0+175		0+0
	WSL Grid Cluster	384	0+295		0+448
<b>TOTAL</b>	<b>11 sites</b>	<b>7305</b>	<b>1044 + 2427</b>		<b>548 + 1294</b>

Country	Site	CPU	Local	Grid	Queueing
Switzerland	Bern ATLAS T3	212	136+0		
Switzerland	Bern UBELIX T3 Cluster	1072	296+258		
Switzerland	Geneva ATLAS T3	222	109+15		
Switzerland	Manno PHOENIX T2	1520	264+438		
Switzerland	Manno PHOENIX T2	1520	270+867		



# Lustre at UBELIX



David Gurtner - Central Informatik Diensten - UniBe

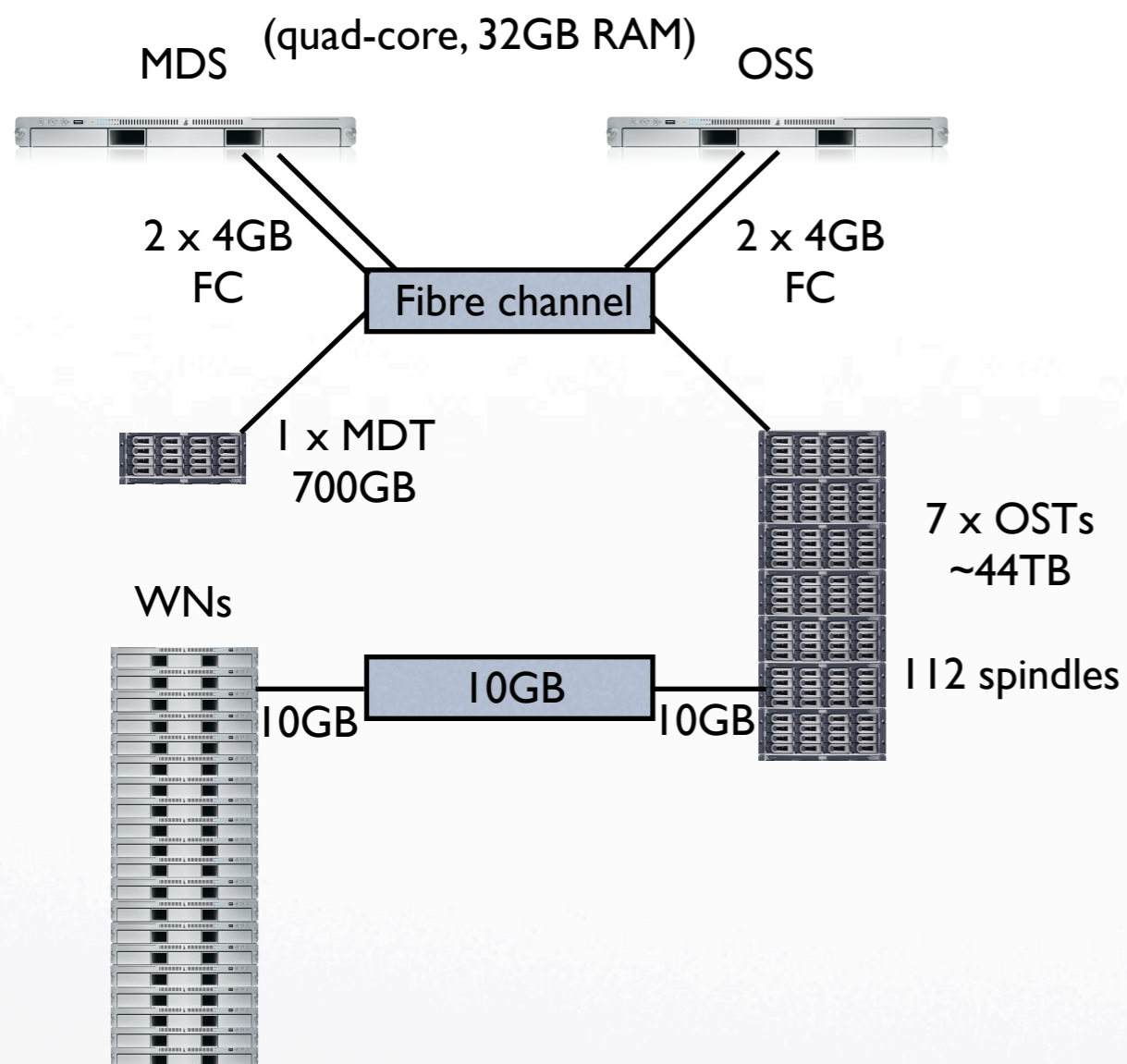




# Lustre at UBELIX



David Gurtner - Central Informatik Diensten - UniBe

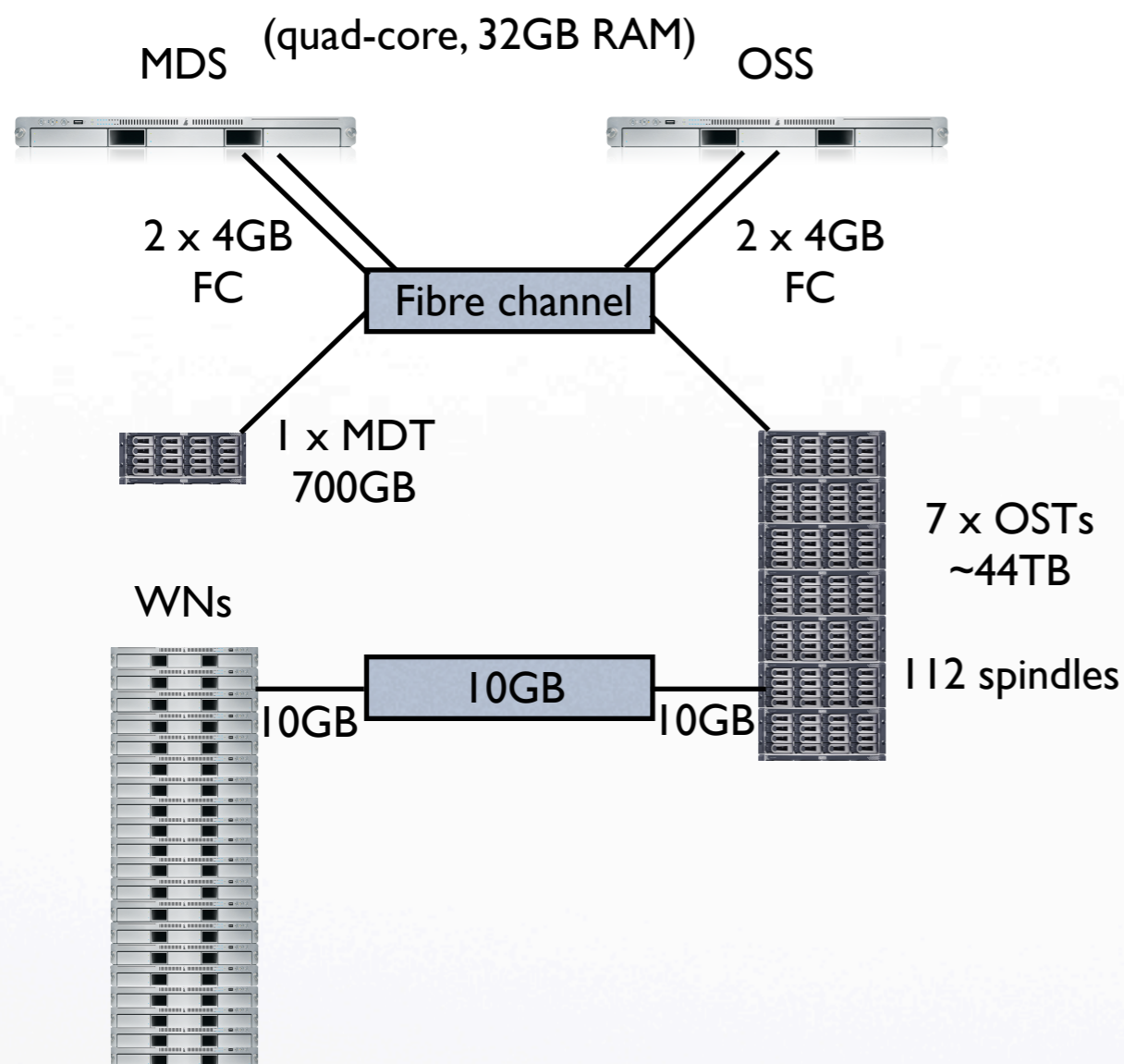




# Lustre at UBELIX



David Gurtner - Central Informatik Diensten - UniBe



## Why Lustre?

- Current system was designed and installed by former admins (> 2 years ago)
- Full design brief somehow “lost”
- Previously using NFS, not scalable/performant for larger cluster sizes
- Open source and free solution
- ...
- Also performance figures not available

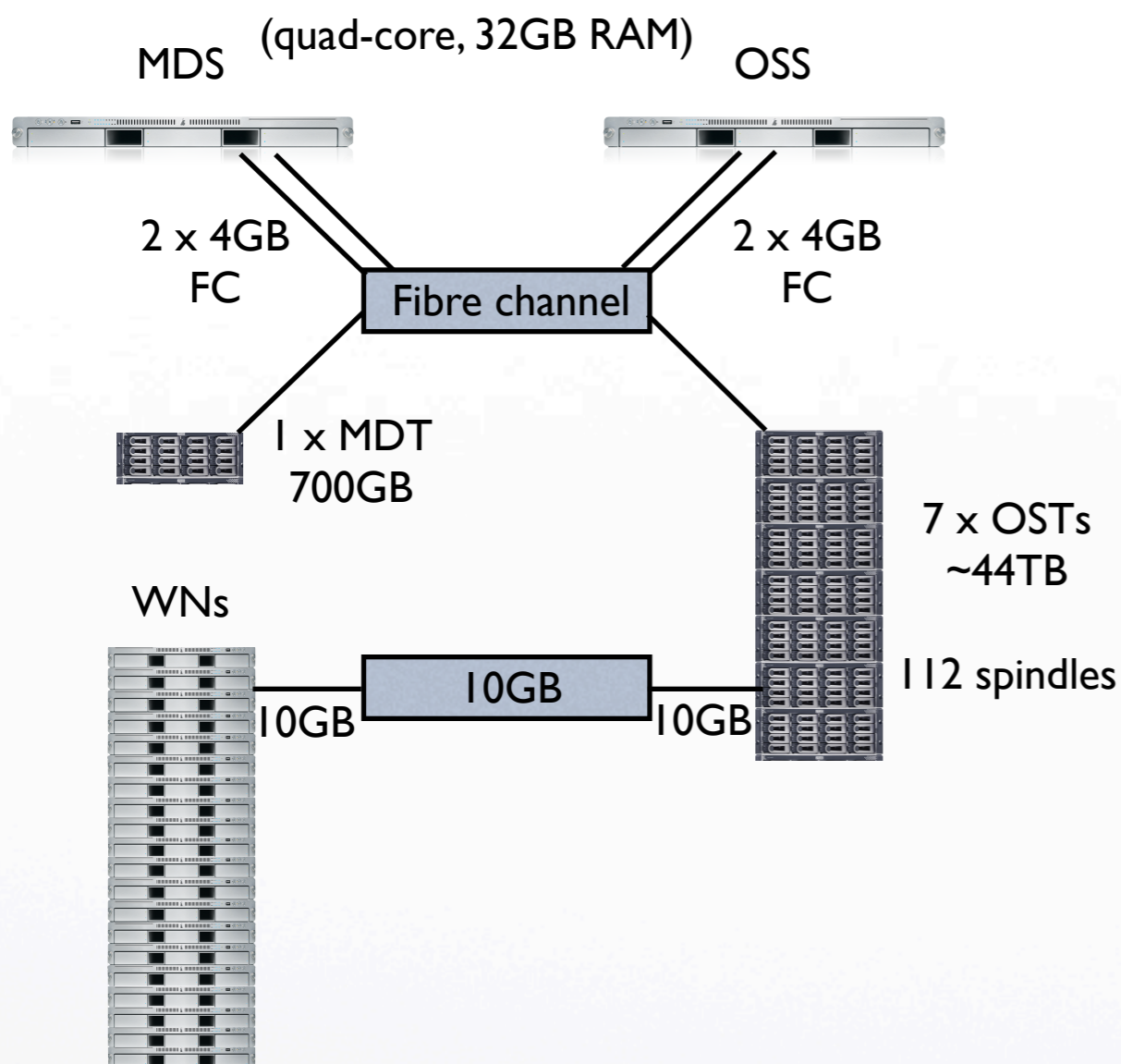




# Lustre at UBELIX



David Gurtner - Central Informatik Diensten - UniBe



## Experience/Issues

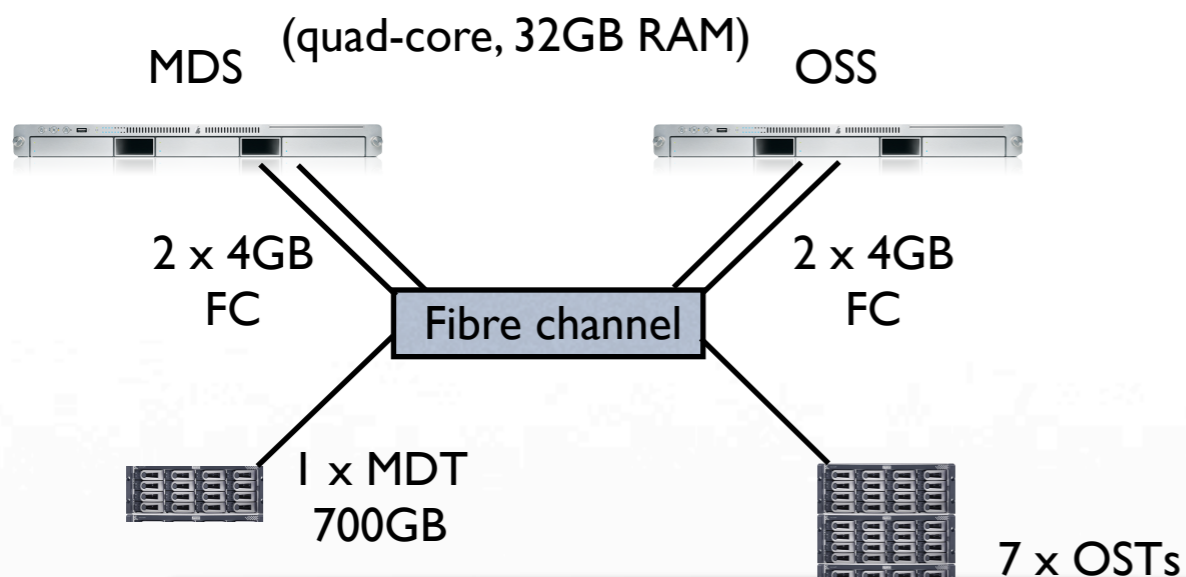
- Instabilities: failures/crashes due to:
  - Network glitch/outage
  - Server failure
  - High I/O usage



# Lustre at UBELIX



David Gurtner - Central Informatik Diensten - UniBe



- Inadequate choice of HW:
  - 700GB for MDT, only ~20GB used
  - Only 2 servers, 1 OSS serving 7 OSTs
- Unreliable HW: OST RAIDs crash "regularly"

## Experience/Issues

- Instabilities: failures/crashes due to:
  - Network glitch/outage
  - Server failure
  - High I/O usage

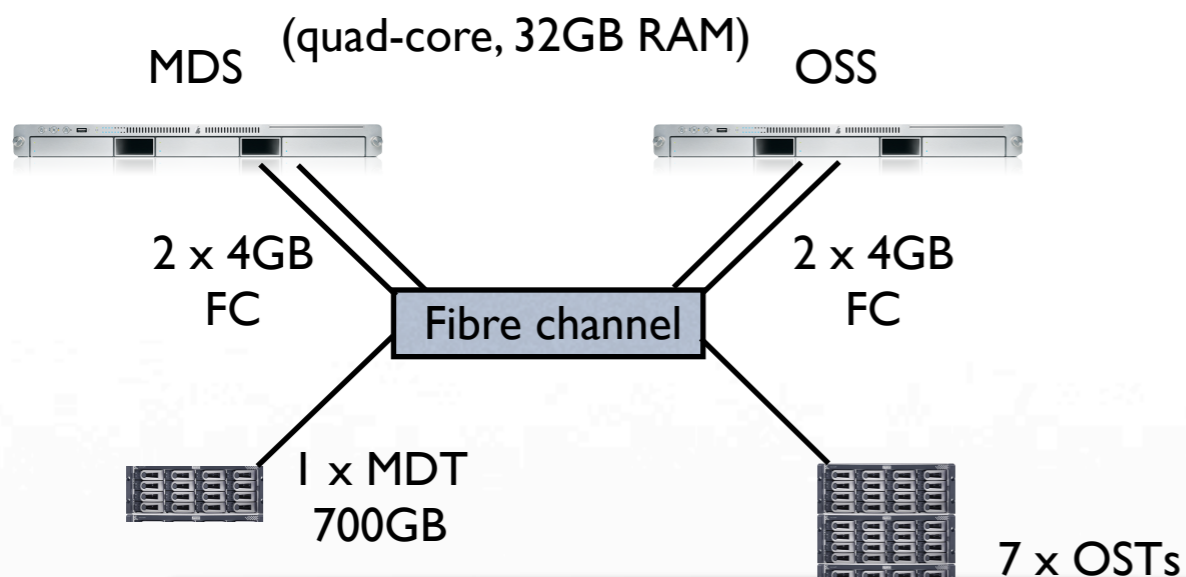




# Lustre at UBELIX



David Gurtner - Central Informatik Diensten - UniBe



- **Inadequate choice of HW:**
  - 700GB for MDT, only ~20GB used
  - Only 2 servers, 1 OSS serving 7 OSTs
- **Unreliable HW:** OST RAID's crash "regularly"

## Experience/Issues

- **Instabilities:** failures/crashes due to:
  - Network glitch/outage
  - Server failure
  - High I/O usage
- **Static setup:** difficult/impossible to add more OSSs to balance rising load due to higher usage (more compute nodes)
- **No quota:** at the time of installation there was no quota support in Lustre
- **Hard to manage/upgrade:** running different versions on server/client leads to segfaults and memory leaks / write own kernel patches (client/server) etc...



## Outlook

- Unhappy with current system
- Old, no sense in expanding it
- Can now fit much larger amount of storage in same rack space and at similar power consumption level
- Development/support roadmap uncertain (SUN/Oracle... ??? ...)
- Lustre on way to be phased out
- Going for different Open Source solution (proprietary still needs sysadmin support!)
- Reliability is most important, but must keep an eye on price
- Considering Ceph for new system: <http://ceph.newdream.net>





# Lustre at LHEP



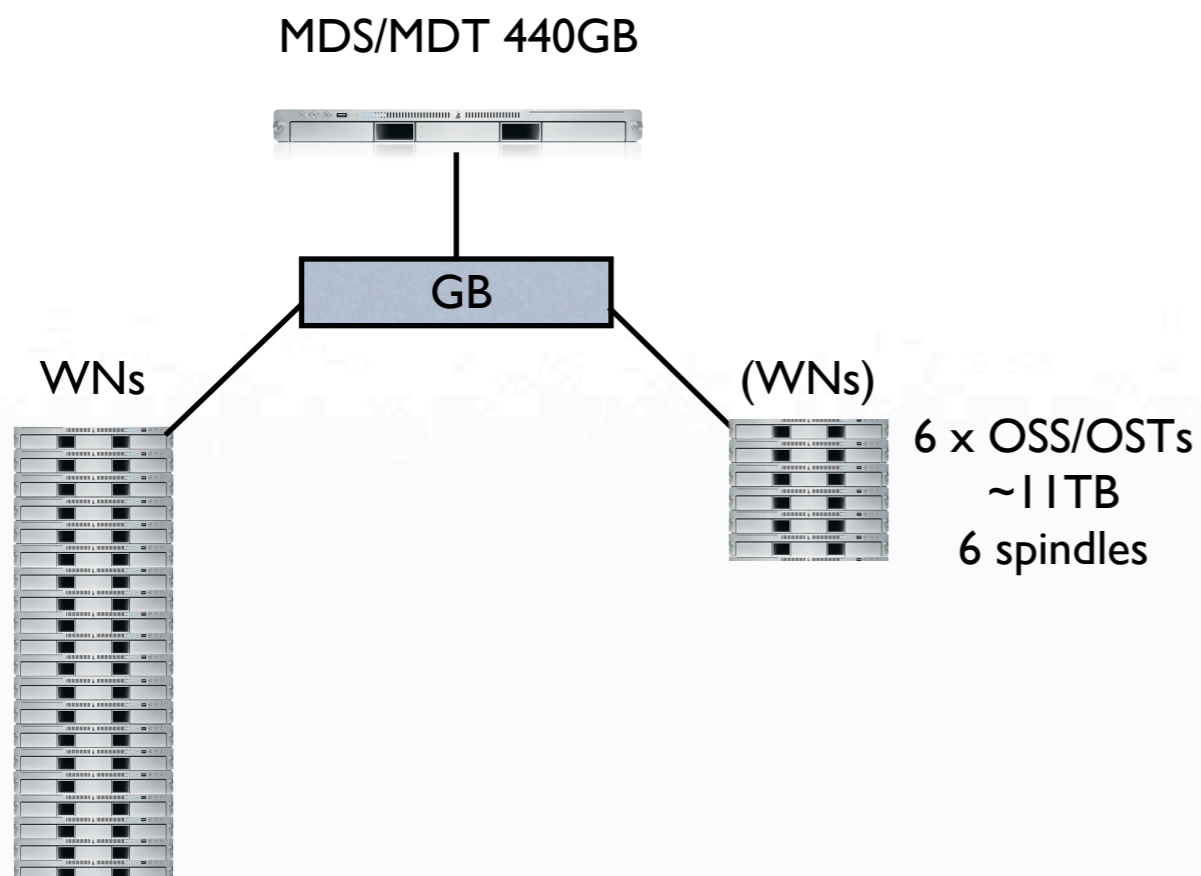
Sigve Haug, Gianfranco Sciacca - LHEP UniBe



# Lustre at LHEP



Sigve Haug, Gianfranco Sciacca - LHEP UniBe

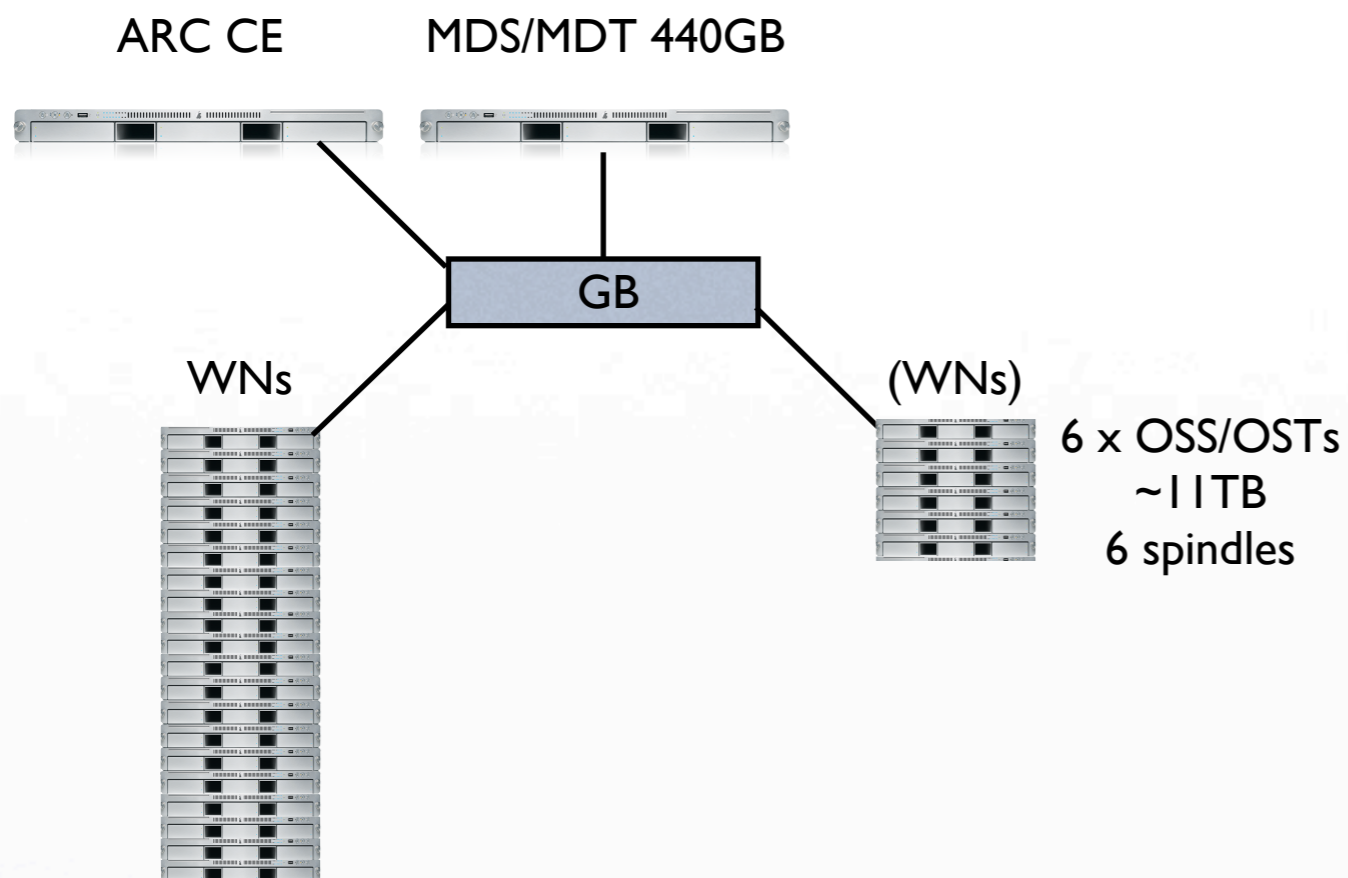




# Lustre at LHEP



Sigve Haug, Gianfranco Sciacca - LHEP UniBe



## Why Lustre?

- Cluster runs the ARC middleware (interface to nordugrid - [www.nordugrid.org](http://www.nordugrid.org))
  - needs shared “session” FS (for job execution)
  - performs much better with shared cache
  - replace NFS, ... experimenting with Lustre...

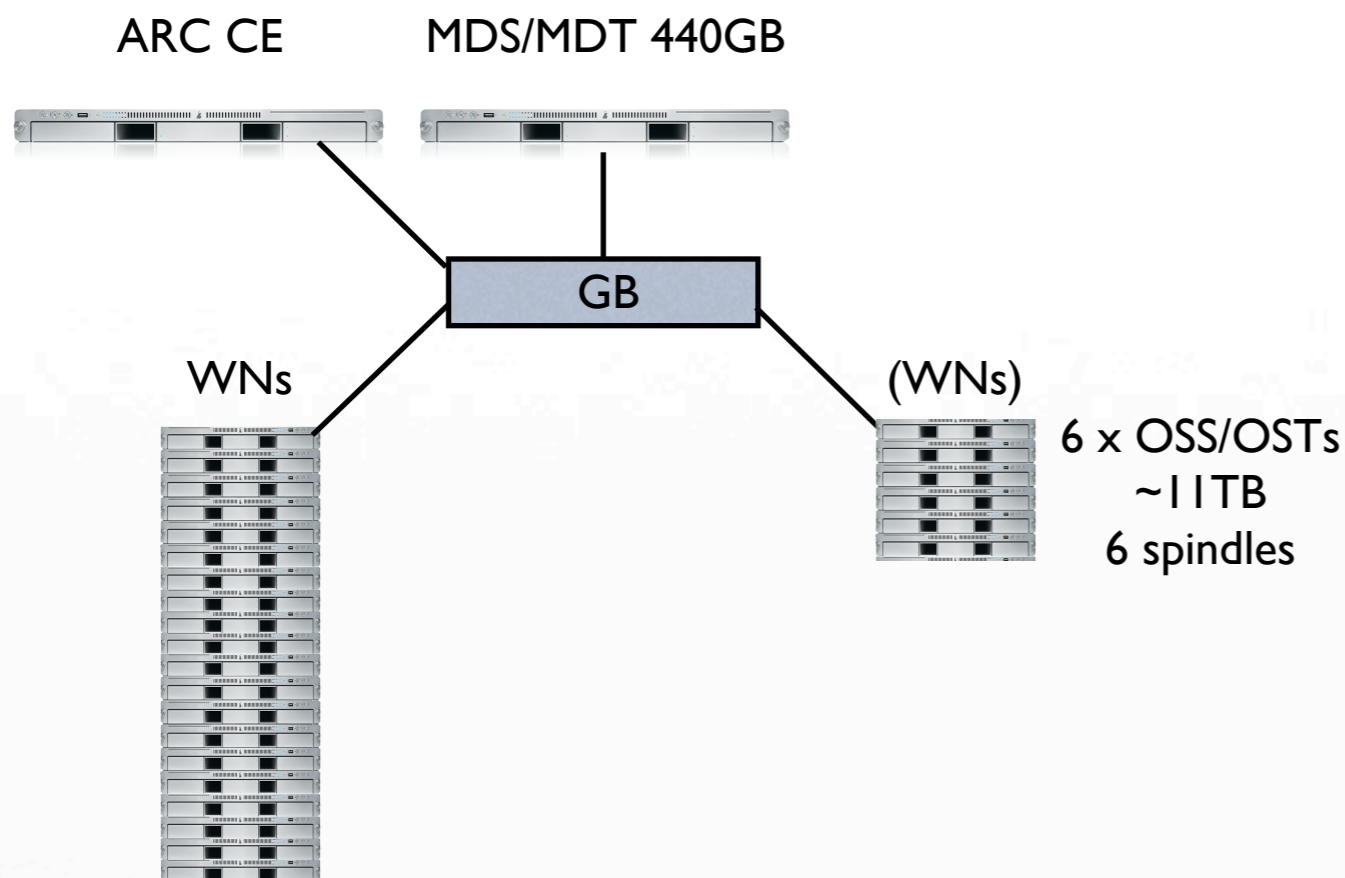




# Lustre at LHEP



Sigve Haug, Gianfranco Sciacca - LHEP UniBe



## Why Lustre?

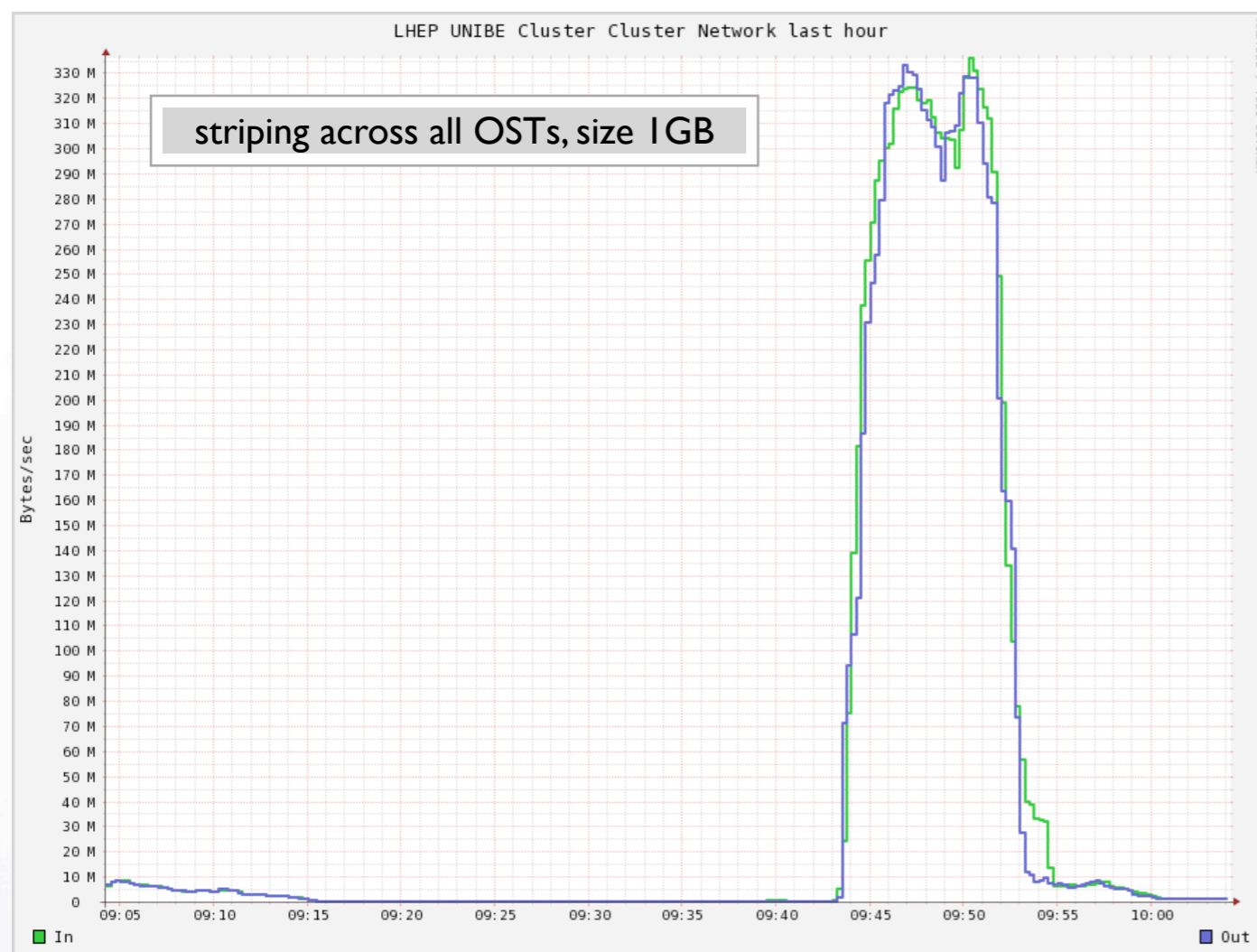
- Cluster runs the ARC middleware (interface to nordugrid - [www.nordugrid.org](http://www.nordugrid.org))
  - needs shared “session” FS (for job execution)
  - performs much better with shared cache
  - replace NFS, ... experimenting with Lustre...
- FS scalable to ~1000 job slots
- Better performant than NFS
- Very low cost: use spare HD slot on WNs. (almost) no extra HW/power
- Reliability? well...
- Did not benchmark the system (it's been a bumpy ride!)
- However...



# Lustre at LHEP



Sigve Haug, Gianfranco Sciacca - LHEP UniBe



live (batch) system going from idle to ~75% occupied

>3 times better than a 24-disk RAID (xfs, NFS over GB)

can probably do better than this

## Why Lustre?

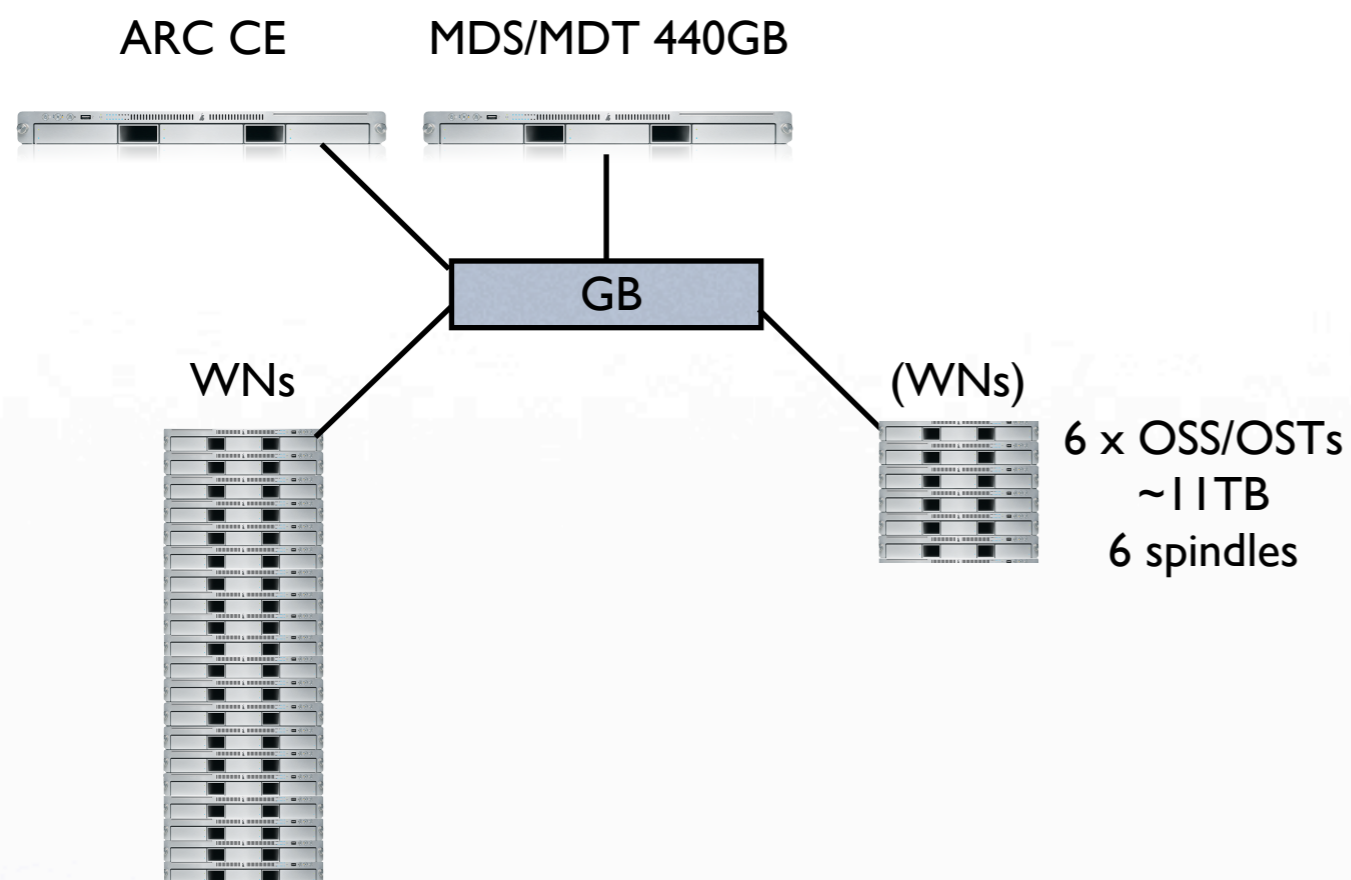
- Cluster runs the ARC middleware (interface to nordugrid - [www.nordugrid.org](http://www.nordugrid.org))
  - needs shared “session” FS (for job execution)
  - performs much better with shared cache
  - replace NFS, ... experimenting with Lustre...
- FS scalable to ~1000 job slots
- Better performant than NFS
- Very low cost: use spare HD slot on WNs. (almost) no extra HW/power
- Reliability? well...
- Did not benchmark the system (it's been a bumpy ride!)
- However...



# Lustre at LHEP



Sigve Haug, Gianfranco Sciacca - LHEP UniBe



striping across all OSTs, size 1GB

A (small scale) success story? Almost...

- Teething troubles with MDS/MDT
- Spontaneous kernel panics (SW? HW? not know yet...)
- Replaced expensive-ish machine with recycled WN: all ~well since!
- Prone to suffer of side effects of other (un-related) problems elsewhere in the system ( /var 100% full on CE...)
- Recent security patching: must re-build kernel and lustre modules:
  - server kernel (lustre patched): way too involved procedure for quick patching
  - decided to go for patchless kernel on clients (disable job execution on WNs/OSSs :-)
- Maintenance effort considerably >0





# Lustre at LHEP



Sigve Haug, Gianfranco Sciacca - LHEP UniBe

## Outlook

- Still in highly “experimental” phase
- Driven by the needs of our current application (and gain valuable experience in the process)
- At the present stage, we don’t need “stellar performance” or “unprecedented scalability”: currently this basic Lustre implementation fulfils our needs (at a very low cost)
- Can possibly close one eye on reliability (no-one likes downtimes, but it’s a scratch area after all)
- Patch Lustre kernel against CVE-2010-3081 on OSSs and (try to) resume batch job execution on these nodes. Should that work: gauge performance impact (if any)
- Cluster will more than double in size very soon (expect up to >500 jobs on it at any given time):
  - Re-deploy higher spec MDS
  - If our current Lustre hits the limit, tweak implementation accordingly (more OSS/OSTs)
  - Re-think our approach to reliability when cluster will have a wider variety of users



# Experiences with Lustre in Bern



David Gurtner, Sigve Haug, Gianfranco Sciacca - ID/LHEP UniBe

THANKS!