

GridKa Site Report

HEPiX Spring 2009

Manfred Alef, Jos van Wezel
Steinbuch Centre for Computing



Forschungszentrum Karlsruhe
in der Helmholtz-Gemeinschaft



Universität Karlsruhe (TH)
Research University · founded 1825

■ CPU:

- +350 nodes (175 Supermicro twins)
 - Intel Xeon L5420
 - 16 GB RAM
 - 2 SATA disks (operating system, local user data)
 - SL5 x86_64 (see below)

■ Disk storage:

- Installation completed of 3.5 PB DDN storage
 - RAID-6 with 1 TB drives, 60 drives per shelf
 - Connected to 12 quad core servers with Intel 10GE cards
 - GPFS filesystems shared between 4 nodes
 - Meta-data mirrored on dedicated SAS Luns
 - Storage for dCache, xrootd, NFS

- **Tape storage:**
 - Preparing installation of STK8500 library
 - install all 10000 slots (2000 accessible via key)
 - Prevents HW expansion downtime in the future
 - Starting with 6 LTO4 drives

■ Tape management:

- TSM is reaching administrative and stability limits
- Added eRMM for library failover and drive pooling
- Using HiStore to monitor drive and media failures/errors

■ (For detailed questions see [talk of] Artem Trunov)

■ Migration to SL5:

- Till April 2009:
 - SL4 i386
- +350 worker nodes in production since April
 - SL5 x86_64
- Problems with the SL5 subcluster detected so far:
 - Alice: ✓
 - Atlas:
 - Missing packages, e.g.
 - compatibility packages,
 - ghostscript, X11 libraries, and X11 devel packages
 - (SELinux is disabled by default)
 - CMS: ✓ (minor problems)
 - LHCb: ?
 - Other VOs: no problems reported so far
 - Open issue: CPU time reporting

CPU Accounting Issue

■ Weird reporting of CPU consumption on some WNs running SL5

■ Output from 'top' command:

```
c01-024-134 > top -bin 1
top - 16:32:31 up 6:09, 1 user, load average: 8.06, 7.92, 7.85
Tasks: 193 total, 9 running, 184 sleeping, 0 stopped, 0 zombie
Cpu(s): 0.0%us, 0.3%sy, 41.2%ni, 58.4%id, 0.1%wa, 0.0%hi, 0.0%si, 0.0%st
Mem: 16443372k total, 8152636k used, 8290736k free, 200808k buffers
Swap: 32764556k total, 0k used, 32764556k free, 6684332k cached
```

PID	USER	PR	NI	VIRT	RES	SHR	S	%CPU	%MEM	TIME+	COMMAND
8096	alef	29	4	112m	87m	1408	R	9999.0	0.5	1339:27	soplex_base.x32
8219	alef	29	4	111m	87m	1408	R	9999.0	0.5	1185:09	soplex_base.x32
8284	alef	29	4	111m	87m	1408	R	9999.0	0.5	1110:53	soplex_base.x32
8415	alef	29	4	111m	87m	1408	R	9999.0	0.5	1044:21	soplex_base.x32
8453	alef	29	4	111m	87m	1408	R	9999.0	0.5	988:01.08	soplex_base.x32
8475	alef	29	4	111m	86m	1408	R	9999.0	0.5	907:43.59	soplex_base.x32
8471	alef	29	4	111m	86m	1408	R	9999.0	0.5	934:43.96	soplex_base.x32
8463	alef	29	4	111m	87m	1408	R	9999.0	0.5	959:53.51	soplex_base.x32
8550	alef	26	10	12736	1036	712	R	208.5	0.0	0:03.09	top

```
c01-024-134 >
```

- **Weird reporting of CPU consumption on some WNs running SL5**
 - Occurs on about 3% of certain classes of WNs
 - Disappears after reboot ...
 - ... but reappears after n reboots:

machine 1:	n >150
machine 2:	n = 21
machine 3:	n = 12
 - No impact on HEP-SPEC06 scores

- **Differing results caused by flavor of installed glibc**
 - Differences in some results between GridKa and other sites
 - Tracked down to libm issue
 - Packages in the SL4 i386 distro are of i386 flavor. However, a few packages are provided as i386 and i686, e.g. kernel and glibc.
 - At GridKa, the i386 glibc had been selected by mistake.

- (Which of the 2 results is the right one?)

Thanks to Rod Walker for the investigations!

Batchsystem (PBS Pro) Issues

- **All worker nodes are managed by only 1 batch system**
 - **Fairshare Algorithm based on required CPU performance (kSI2K, HEP-SPEC06)**
 - **2 serious problems encountered (scaling issues?):**

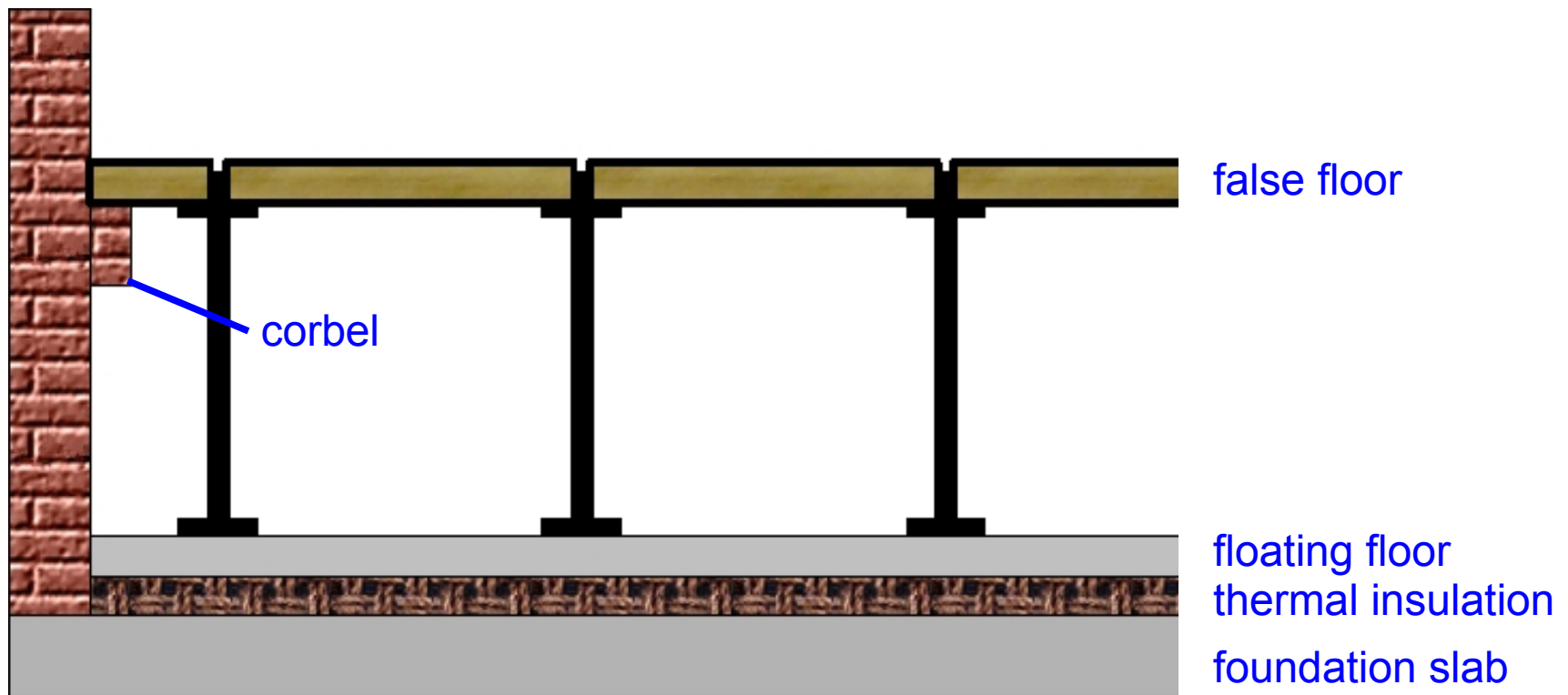
- **Scheduler launches jobs on WN which is temporarily unavailable**
 - What (repeatedly) happens:
 - WN temporarily unavailable (timeout in server to mom communication)
 - PBS reports "node down, communication closed" to the logfiles but ...
 - ... it tries to start hundreds of jobs on this particular node even though
 - Of course, they cannot run on that node!
 - Command 'qstat -f' reports:
 - state=Q (queued)
 - exec_host already set to that unreachable node
 - Now hundreds of jobs are waiting for being started on that particular node :-)
 - **No bug fix available so far!**
 - Statement from PBS support: 'Bug is already fixed in the latest release' ...
 - ... but the same problem reappeared a few days after the update :-)
 - Waiting for the next PBS release ...
 - No easy way to get the waiting jobs running

Batchsystem (PBS Pro) Issues

- **PBS server not responding**
 - Probably caused by network issues (timeouts)?
 - Couldn't restart the PBS server processes
 - Temporary fix (Xmas holidays!):
Reboot of PBS server from cron job :-)

Infrastructure: Stability of False Floor

- Load carrying capacity of the raised floor: **1.5 t / m²**
 - Floor board size: 75x75 cm²

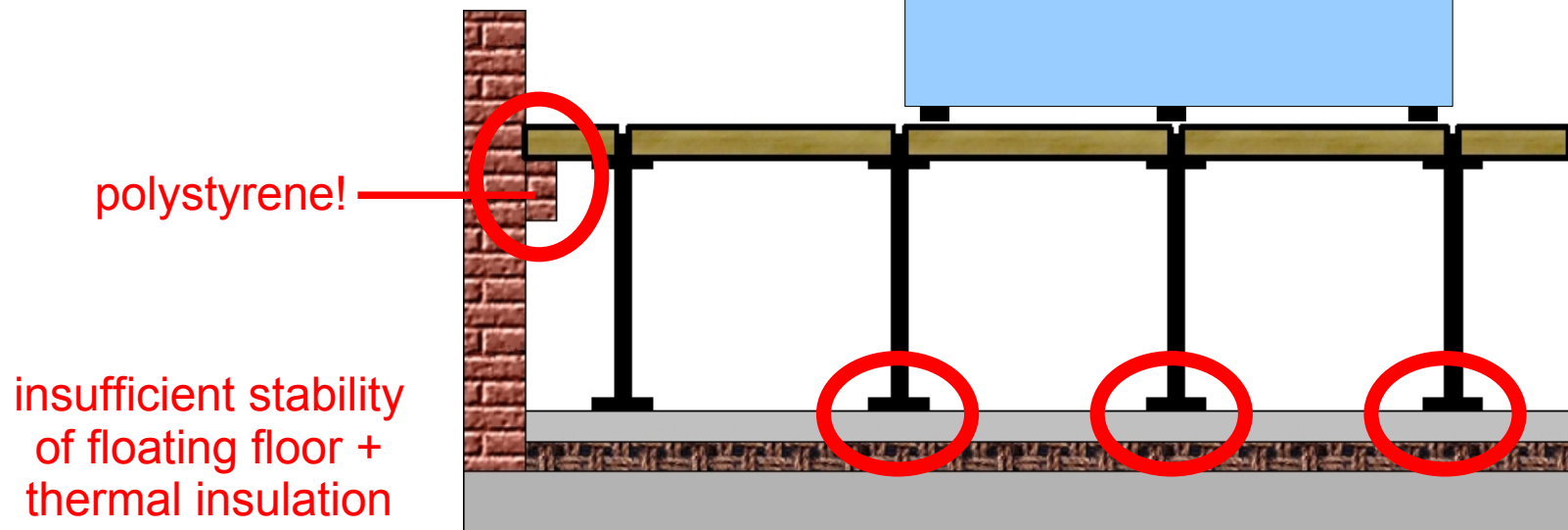


Infrastructure: Stability of False Floor

■ Newly installed high-density storage systems (60 hard disks per 4U):

> 1 t / rack

- Rack footprint: 80x150 cm²
- Floor board size: 75x75 cm²



■ Proposed solution: second subconstruction to improve stability

Questions, comments?