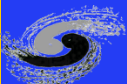


IHEP Site Report

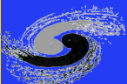


Hepix Fall 2013 Workshop
Shi, Jingyan
Computing center, IHEP



Outlet

- **BEIJING-LCG2 Tier2 EGI Site**
- Local Farm
- SDN Network



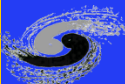
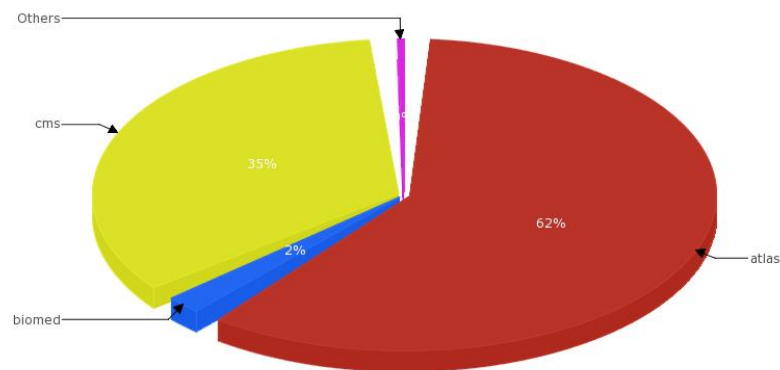
Resource of T2 site

- **Mainly Support:**
 - CMS & ATLAS
- **1500 job slots**
- **400TB storage**
 - 320TB for dCache
 - 80TB for dpm

Elapsed time per vo

Developed by CESGA 'EGI View': / normelap-HEPSPEC06 / 2013:1-2013:10 / SITE-VO / all (x) / GRBAR-LIN / 1

BEIJING-LCG2 Normalised Elapsed time (HEPSPEC06) per VO

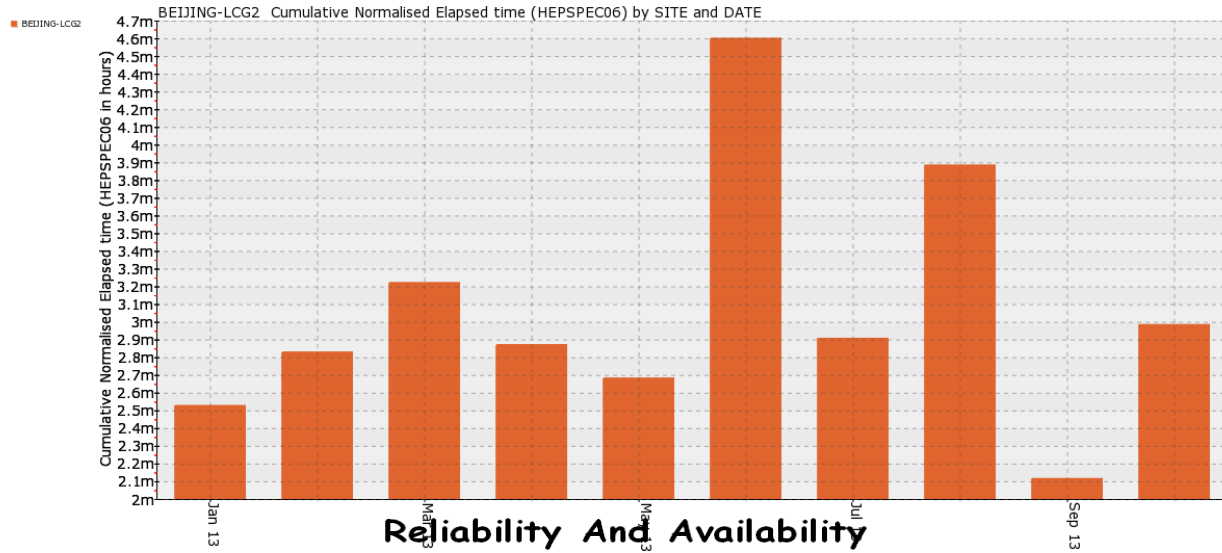


BEIJING-LCG2 Site report

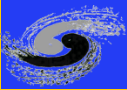
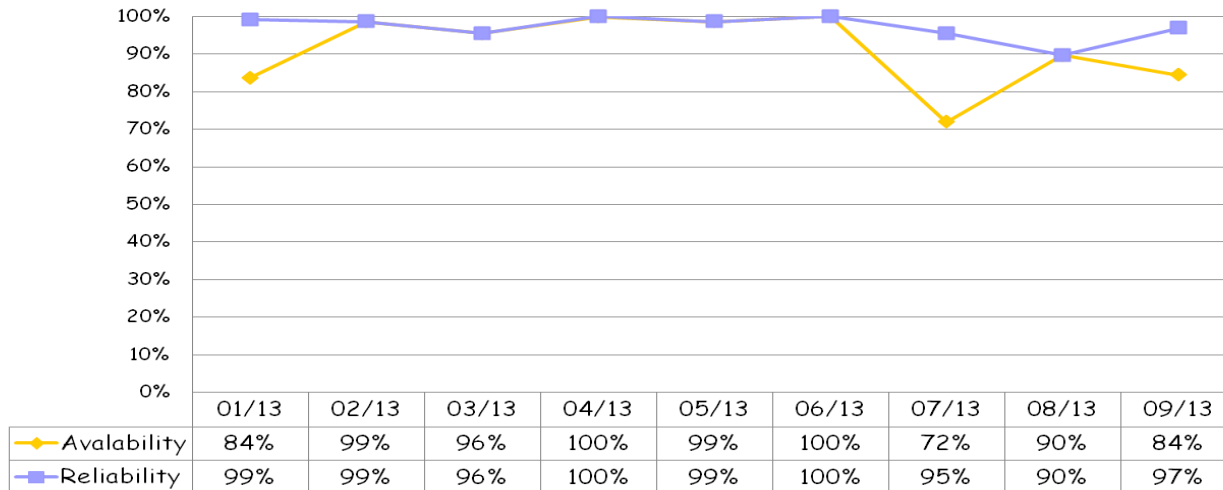
Elapsed cpu time by date

Developed by CESGA EGI View: / nomlap+HEPSPEC06 / 20131-201310 / SITE-DATE / all (x) / GRBAR-LIN / i

2013-10-24 08:44

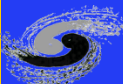


Reliability And Availability



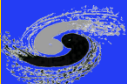
Site Operation

- All WNs Upgraded to SL 6.4 x86_64
- Atlas disk arrays are replaced by the new disk arrays
 - Capacity reduced from 320TB to 80TB due to the limit budget
 - Only production jobs
- DPM and CreamCE have been upgraded to SHA2 compatible
- dCache servers will be upgraded to SHA2 compatible before 1th Nov.



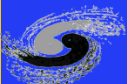
New CA Server

- **Old CA Server**
 - version 1.0
 - SHA2 unsupported
 - Cert's DN with email address caused error with some grid jobs
- **New CA server**
 - Version 3.0
 - New CA Root Cert's DN without email has been released
 - SHA2 supported
 - Signs new SHA2 Cert for users



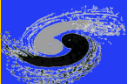
Outlet

- BEIJING-LCG2 Tier2 EGI Site
- Local Farm
- SDN Network



Resource of Local Farm

- Support several HEP experiments and some bio-med experiments
 - BES, DYB, Nano etc.
- 7500 job slots
- 3PB+ disk storage
- 5PB+ tape storage



Scheduler

- **Torque+Maui**

- **Heavy job schedule tasks**

- Generally, 100,000 jobs in the 50 queues each day

- **Version: Torque- 2.5.5 Maui-3.2.6**

- No way to schedule gpu task: extra scripts are provided

- **Performance tuning**

- Dedicated short queue created to increase the resources efficiency

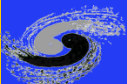
- More statistic on the jobs running added

- Integrated with monitor system

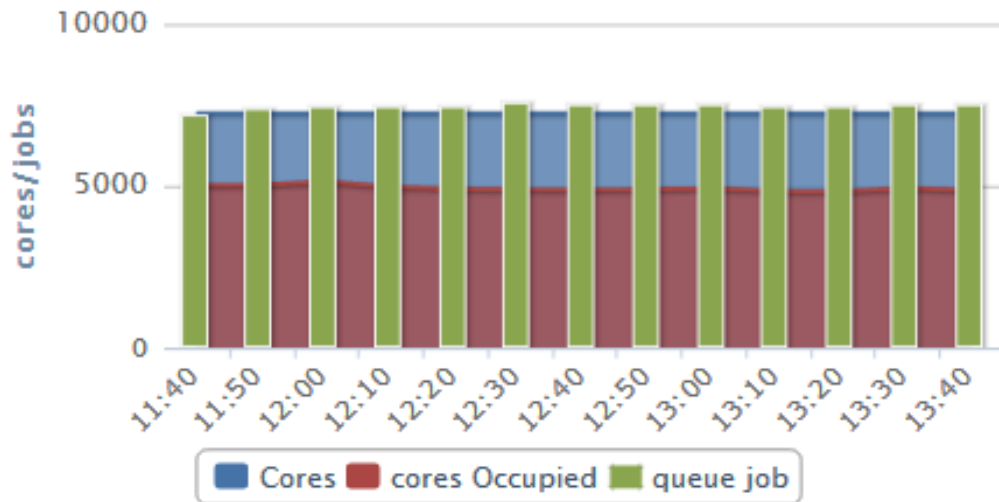
- Unexpected crashed nodes detected by the monitor system could be excluded automatically

- More scheduler policy are under discussion

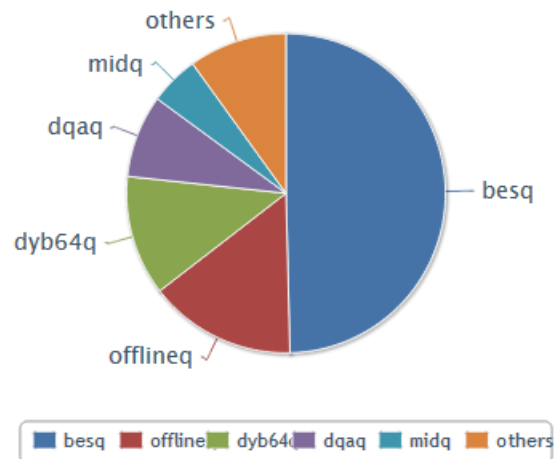
- Start thinking about the new scheduler plan



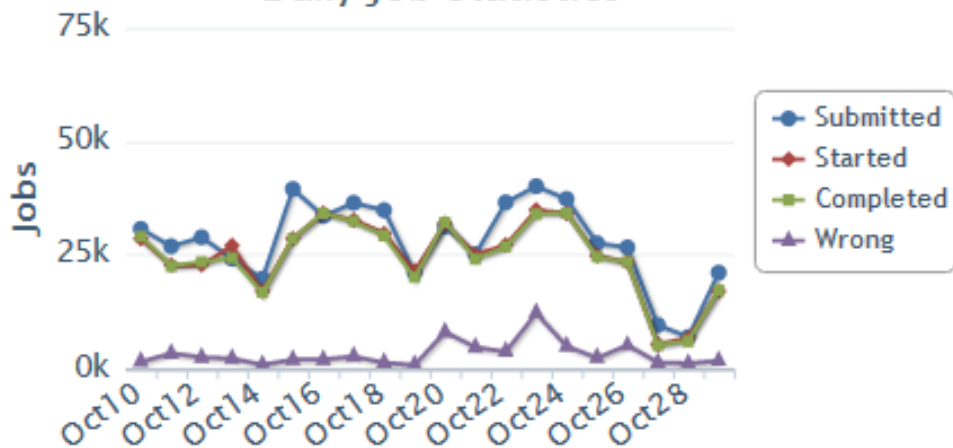
Computing Resource Utility ALL Applications



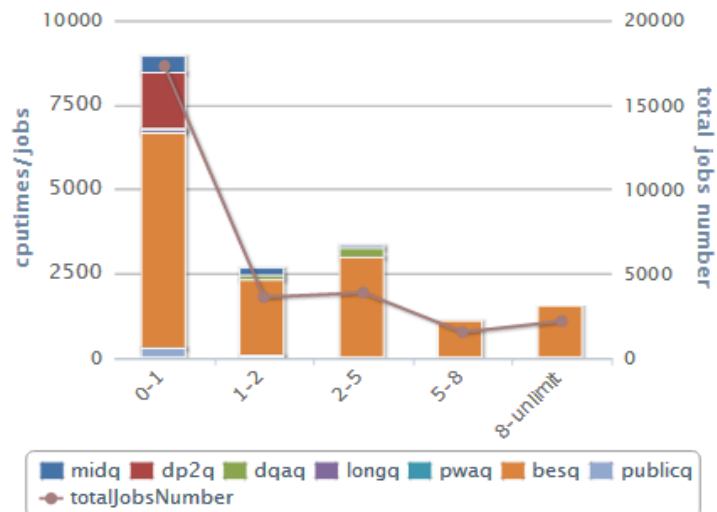
Started jobs



Daily Job Statistics



Computing Cputimes Distribution bes3-farm-s15



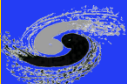
Migrate Quattor to Puppet

- **Migration plan**

- Puppet in production environment -- 1st, Oct
- Basic groups and modules definition --15th Oct
- New coming nodes are managed by puppet --31st, Oct
- Working nodes are managed by puppet --30th, Nov
- All machines are managed by puppet --31st, Dec

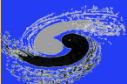
- **Current Status**

- Puppet has been setup in the production environment
 - Foreman 1.3 + Puppet 3.3.1
- Modules are coded and groups are divided
- More than 10 nodes are managed by puppet
 - Gpu nodes
 - Working nodes



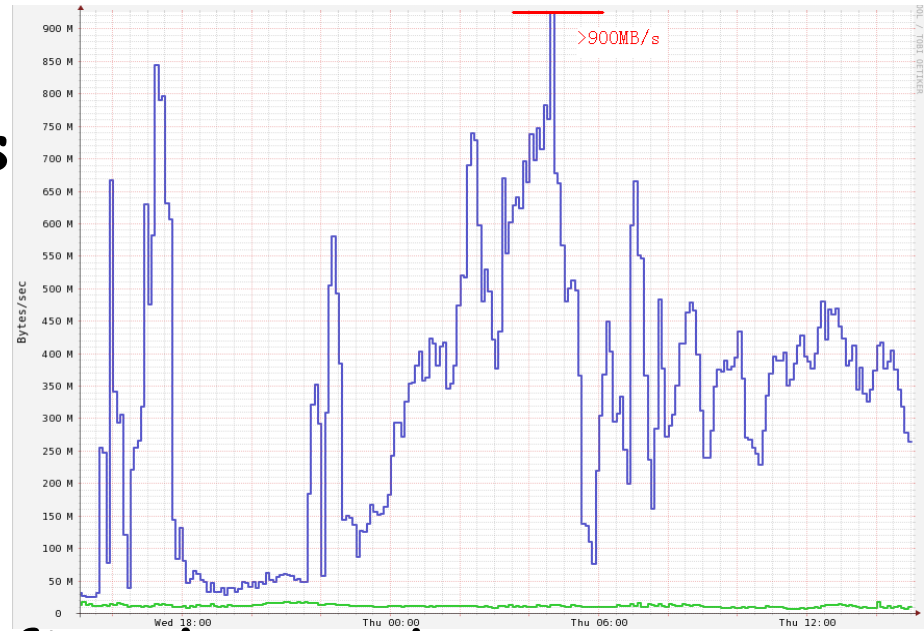
Storage - lustre

- **3+PB lustre storage for most of users**
 - 43 file servers, 5 mount points
- **1 PB new lustre storage will be added**
 - Hardware burning testing
- **Detection jobs actions**
 - Define typical jobs I/O actions
 - Detecting the job's I/O action running on worknode
 - Integrated with scheduler's log to give an over view of the cluster status
 - What kind of jobs are running
 - Performance tuning



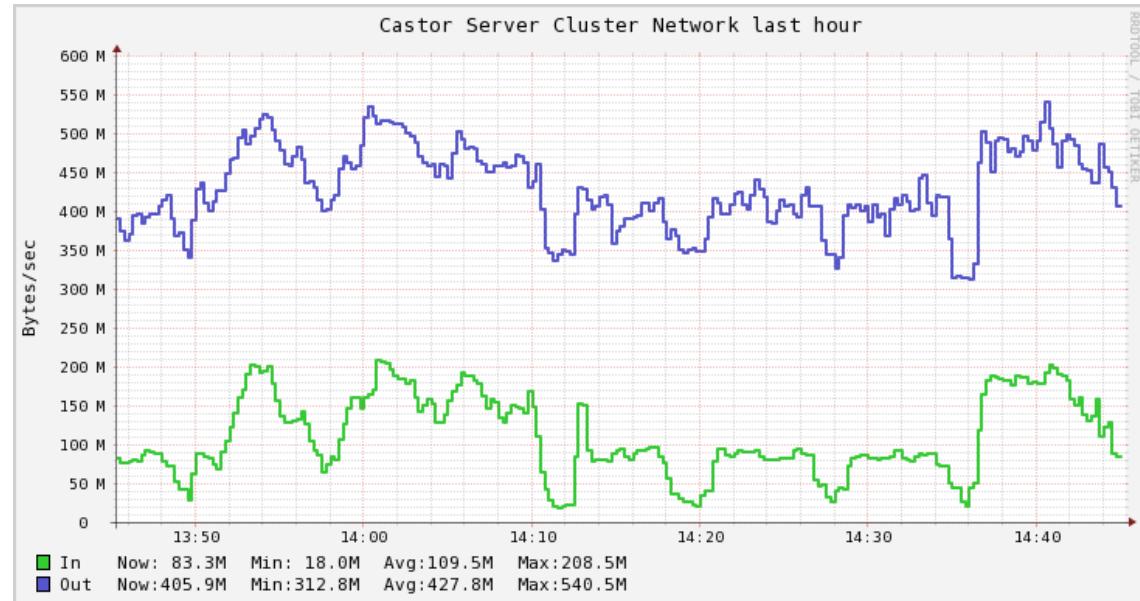
Storage - gLuster

- 186TB storage provided for cosmic-ray experiments
- Performance are quite satisfied
 - 4 file servers provides 3.2GB/s peak performance
 - Not stable at the beginning - complained by the user
- gLuster Code modified to fits the user's request
 - Unified layout
 - Double-lookup instead of all bricks lookup
 - Lost directories in filled on demand automatically
 - File access can't be blocked when some brick is down
- Metadata service for gLuster are under development



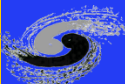
Tape Lib

- **Modified Castor-I to manage the Tape LibS**
 - 10 disk servers + 190TB cache
 - 10 tape servers + 26 drives
 - 2 tape libraries + 2 robots
- **Experiment On-line data storage and backup storage**
- **Performance testing**
 - Copy files from tape lib to Luster disk
 - Files are copied to Luster directly skipping the disk cache
 - Files are sequenced by its location on the tape before the process of copy
 - The peak performance could reach 500MB/s



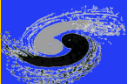
Difficulty we are facing

- **2/3 hardware have expired their guarantee periods**
 - **Some of them have been running over 5 years**
 - **Hardware error happened more and more frequently**
 - **For example: Disk Error 10 disk errors/week**
 - Disk rebuild process is quite slow and unbearable
 - Under the risk of data lost
 - **Memory stick, power module, disks are maintained by ourselves**
 - **Manpower consumed**



Outlet

- BEIJING-LCG2 Tier2 EGI Site
- Local Farm
- **SDN Network**



SDN Status@IHEP

Goal

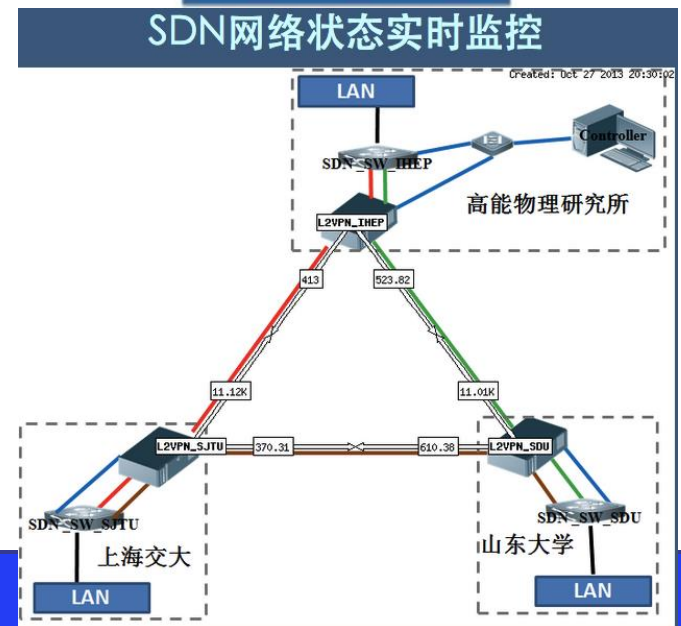
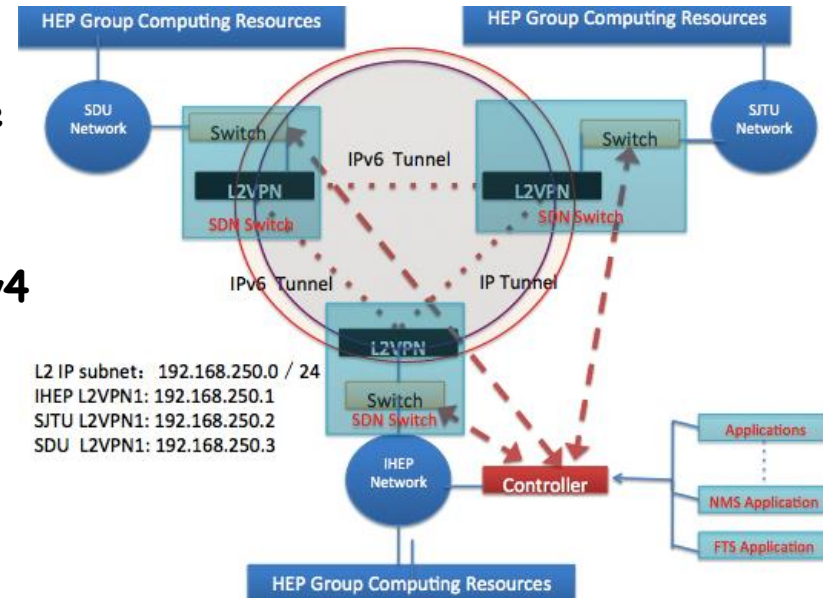
- A flexible, reliable and high performance HEP data transfer network in China
- Both IPv4 and IPv6 supported
- The traffic can be switched between IPv4 and IPv6

SDN@IHEP

- End user network
- Backbone network (IPv6 & IPv4)
- SDN Switch (L2VPN gateway & Openflow supported)
- Control center (API to Application)
- Applications (FTS/NMS/.....)

Status

- The SDN test bed is running well
- The traffic will choose the path based on the available bandwidth between source and destination sites
- We are trying to optimize the IPv6 performance (COS policy)



Thank you!

