# CERN IT Department

# HEPiX Report

Helge Meinhard, Steve Murray,
Miguel Coelho dos Santos / CERN-IT

Computing Seminar / After-C5

21 November 2008

# Outline

- Meeting organisation, site reports, data centre track (Helge Meinhard)

- Storage (including Castor), benchmarking (Miguel Coelho dos Santos)

- Operating systems and applications, virtualisation, network and security, miscellaneous (Steve Murray)

# HEPiX

- Global organisation of service managers and support staff providing computing facilities for HEP
- Covering all platforms of interest (Unix/Linux, Windows, Grid, …)
- Aim: Present recent work and future plans, share experience; provide technical advice to IHEPCCC (now suspended)
- Meetings ~ 2 / y (spring in Europe, autumn traditionally in North America)

- Held 20 – 24 October at Academia Sinica in Taipei, Taiwan
  - Second meeting in Asia ever…
  - First one was in Tsukuba (JP) in 1991!
  - Intention: Offload US and CA, attract participants from Asia
  - Excellent organisation (Simon Lin, Vicky Huang, Stella Shen, Jill Lin et al.) set very high standards
    - Difficult to meet for other sites…

- Format: Pre-defined tracks with conveners and invited speakers per track
  - Still room for spontaneous talks – either fit into one of the tracks, or classified as 'miscellaneous'
  - Again proved to be the right approach
  - Judging by number of submitted abstracts, storage remains a very hot topic
- Full details, slides and tons of photographs: http://indico.twgrid.org/internalPage.py?pageId=1&confId=471
  - Alan Silverman's trip report available, too

- 86 registered participants, of which 12 from CERN
  - Andreeva, Cass, Coelho dos Santos, Duellmann, Field, Lopienski, Lo Presti, Meinhard, Murray, Otto, Ponce, Silverman
  - Other sites: **ASGC,** CASPUR, CEA, DESY Hamburg, DESY Zeuthen, FNAL, GSI, **HAII Bangkok,** IN2P3, INFN, JSI Ljubljana, **KEK, KISTI (South Corea),** KIT, LAL, LBNL, **Manila, U Melbourne,** U Michigan, NIKHEF, NDGF, NSC Sweden, **Punjab (Pakistan), Putra (Malaysia),** Prague, RAL, SLAC, **U Taipei, Tata, Tokyo,** TRIUMF, U Victoria
  - Compare with St Louis (autumn 2007): 62 participants, of which 11 from CERN

- 64 talks, of which 16 from CERN
  - Compare with St Louis: 59 talks, of which 15 from CERN

- Next meetings:
  - Spring 2009: Umeå (Sweden, 25 – 29 May 2009)
  - Autumn 2009: Not yet fixed (LBNL interested)
  - Further application: GSI

**CERN IT Department**

- Cooling, power, space
  - Shortfalls mentioned by all large and most small centres
  - Solutions include new data centres (RAL), temporary structures (like Blackbox – not without issues), refurbishing rooms/halls built for other purposes, using remote rooms or commercial hosting, upgrading structures in place
  - Too little focus (IMHO) on power-efficient systems
  - More problems due to harmonics reported
  - Increasing number of sites retire old stuff aggressively

# Site Report Highlights (2)

- ## CPU worker nodes
  - Mixture of blades, twins and traditional 1U pizza boxes (sometimes even with redundant PSUs)
  - Mostly Xeons (Harpertowns), some Opterons (driver issues for chipset)

- ## Disk storage
  - Mixture of storage-in-a-box (either "white boxes" or Thumpers), Infortrend (or similar) arrays with FC (2400 TB in INFN SAN), larger units (e.g. DataDirectNet in Karlsruhe – 16 PB over 3 years)
    - Separating storage from serving nodes appears to be architectural criterion at some sites

- ## Tape storage: SL8500, LTO drives

# Site Report Highlights (3)

- ## File systems
  - Increasing interest in cluster file systems
  - A number of sites mentioned GPFS
    - Mostly qualitative statements, no performance data
  - Lustre
    - Detailed reports by DESY and GSI

- ## Windows: Still some (a lot of?) reluctance concerning Vista (in particular 64-bit) and Office 2007

# Site Report Highlights (4)

- Communications
  - More and more sites considering 10 GigE to storage nodes
  - VoIP and unified messaging at FNAL
  - IPv6: Dedicated talks, but little activity
- Software
  - Quattor, Lemon mentioned a few times
  - InDiCo used at most major sites (including ASGC)
  - DPM
  - Castor
  - Subversion taking off (and over from cvs)
  - MAUI: Scaling problems; LSF (and GPFS): Issues with core-based licencing

CERN
**IT**
Department

- Miscellaneous
  - Number of name changes, structure changes, personnel changes, …
  - ITIL mentioned by a number of sites, FNAL going for ITIL V2 and ISO 20000 certification within two years
  - Some sites reported to start (only) now deploying IPMI
  - Operational coverage: Some Tier1 sites working with 7 hours/working day operator coverage only

- Storage issues
- Castor
- **Data centres**
- Operating systems and applications
- Virtualisation
- Network and security
- Benchmarking
- Miscellaneous

# Data Centre Track (1)

- Presentations:
  - New CERN computer centre (Tony Cass)
  - Push-pull registry for discovery of OWL-S Web services (Hafiz Farooq Ahmad)
  - Data centre thermodynamic CFD (Enrico Mazzoni)
  - Roma Green computing centre (Giovanni Organtini)
  - High-density displays in CluMan project (Miguel Coelho dos Santos, for Sebastian Lopienski)

# Data Centre Track (2)

- Data centre thermodynamic CFD (Enrico Mazzoni)
- Report about work in progress
- Motivation: higher energy densities require more detailed knowledge
- Input: machine room design parameters (CAD model), server properties from data sheets
- CFD results have correct tendency
    - But are a few degrees to high
    - Will iterate: change input parameters to make the results match better with room measurements
- CFD is a valid tool to obtain useful information about cooling

- Roma Green computing centre (Giovanni Organtini)
  – Rather small layout with 14 racks (Knürr 17 kW)
    - Battery UPS with 120 kVA feeding racks and water pumps
    - 3000 l water at 12 deg C
  – Currently at 45 kW for servers, 100 kW total
  – Estimate for full load: 90 kW for servers, 59 kW for services
    - Total / servers = 1.65
    - Estimate for air cooling: ~ 70 kW more

HEPiX Fall 2008
ASGC (Taiwan)

# Storage issues, Castor and Benchmarking report

**Miguel Coelho dos Santos**
**FIO/FS**
**November 2008**

# Storage track – day one

| | |
|---|---|
| 12:00 | |
| 13:00 | **[40] Data Integrity and Data Security**<br>by Dr. Ted PANG (Infortrend)<br>(2nd Conference Room: 13:30 - 14:00) |
| 14:00 | **[41] Computing System Performance Bottlenecks for IOs and DTS Solution**<br>by Dr. Hiro TAKAHASHI (DTS)<br>(2nd Conference Room: 14:00 - 14:30)<br><br>**[43] HEPiX Storage WG Update**<br>by Dr. Andrei MASLENNIKOV (CASPUR)<br>(2nd Conference Room: 14:30 - 15:00) |
| 15:00 | **[119] Lustre File System news and roadmap**<br>by Peter BOJANIC<br>(2nd Conference Room: 15:30 - 16:15) |
| 16:00 | **[121] CERN storage update**<br>by Dirk DUELLMANN<br>(2nd Conference Room: 16:15 - 16:45)<br><br>**[122] Lustre at GSI**<br>by Thomas ROTH<br>(2nd Conference Room: 16:45 - 17:15) |
| 17:00 | |

## Agenda

### End-to-End Data Protection

### Undetected Error in Storage System

### What is DIF

| | | 512 | 514 | 516 | 519 |
|---|---|---|---|---|---|
| 512 bytes of data | | GRD | APP | REF | |

16-bit guard tag (CRC of 512-byte data portion) —
16-bit application tag —
32-bit reference tag —

- **8 bytes** DIF for each standard 512 block of data
  Logical Block Guard 2bytes – CRC (Check bit error)
  Logical Block Application Tag 2bytes – User Defined
  Logical Block Reference Tag 4 bytes – LBA (Check displacement)

- **SATA T13/EPP(External Path Protection) uses same format**
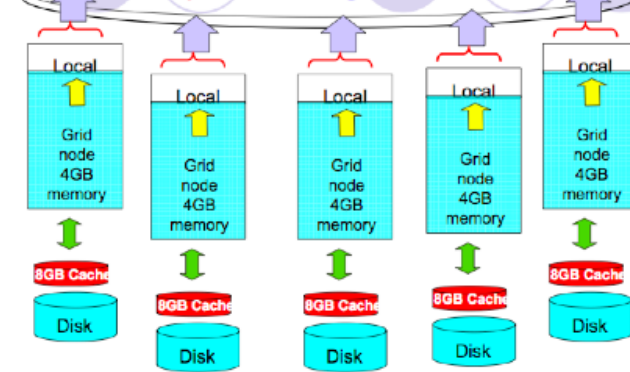
Infortrend Reliable Networked Storage Solutions

## Proposed concept:
## Balanced Memory Architecture

- Unbalanced performance shape creates the bottleneck of computer utilization.
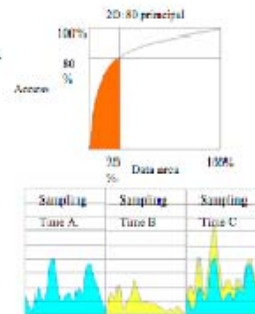- DTS basis memory management maintains balanced performance hierarchy shape.

DTS memory basis memory hierarchy.

Traditional Computer memory hierarchy.

18

## Proposing *DTS L4 cache memory system* configurations

Reduce a nuclear destruction semiconductor memory error on the local memory. 17

## DTS Intelligent caching algorithm

- Decide the most busy block address area.
- As 80 /20 principal, DTS intelligent cache keep these data on DTS cache table memory effectively.
- The DTS system enable most effective IO transfer performance by own cache technology and DTS working table's physical speed.

| Sampling Time A. | Sampling Time B. | Sampling Time C. |
|---|---|---|

22

## IO performance test result: DTS VS traditional computer

We researched write IOPS under DTS memory configuration and standard computer configuration as well. The result was very positive. For instance, 512 100%Write0%Read Random result is 211.51 IOPS but DTS based is 33354.32 IOPS under eight multiple accesses situation. Also 1K 80%Write20%Read Random result is 237.14 IOPS but DTS based is 32020.75 IOPS under eight multiple accesses situation. (Refer figure 6)
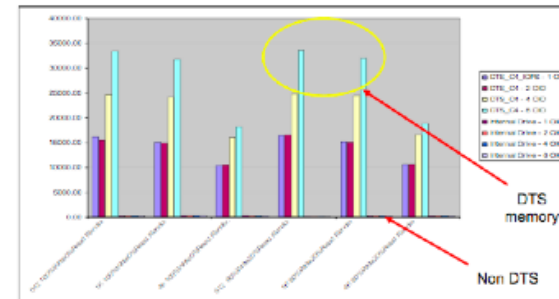
Figure 6: DTS type vs Traditional computer IOPS test result.

(Intel Xeon 2.4 GHz, 4GB System memory, 4GB Working Table, Target File Size = 2GB 25 Write Back Policy, Windows2000 Advance SP4 Server Target Disk size 4GB)
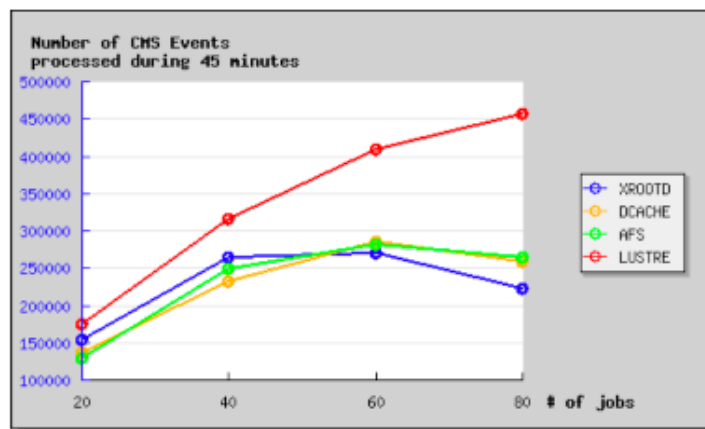
◆Sun

## World's Fastest and Most Scalable File System

- Parallel, scalable shared POSIX file system
- Key benefits
  - Petabytes of storage – one name space
  - Tens of thousands of clients
  - High-performance heterogeneous networking and routing
  - High availability
  - Open source, multi-platform and multi-vendor
  - Object-based architecture
  - Supported by Cray, Sun, and Others

**Lustre Releases**

◆Sun

### Lustre 1.8
(9/18 Code Freeze, 11/21 GA)
- Clients interoperate with 2.0
- Adaptive Timeouts, VBR
- OST Pools
- OSS Read Cache
- Client SMP Scalability
- Service Tags

### Lustre 2.0
(11/1 Code Freeze, 4/1/09 GA)
- Security (Kerberos)
- Server Change Logs
- Replication
- Commit on Share
- Client IO Subsystem Enhancements

### Lustre 3.0
(Schedule being worked now)
- Size on MDS
- Improved MDS concurrency
- OSS Write Cache

- Portability Improvements
- Routed Network Improvements
- HSM Infrastructure + HPSS
- OST Space Management
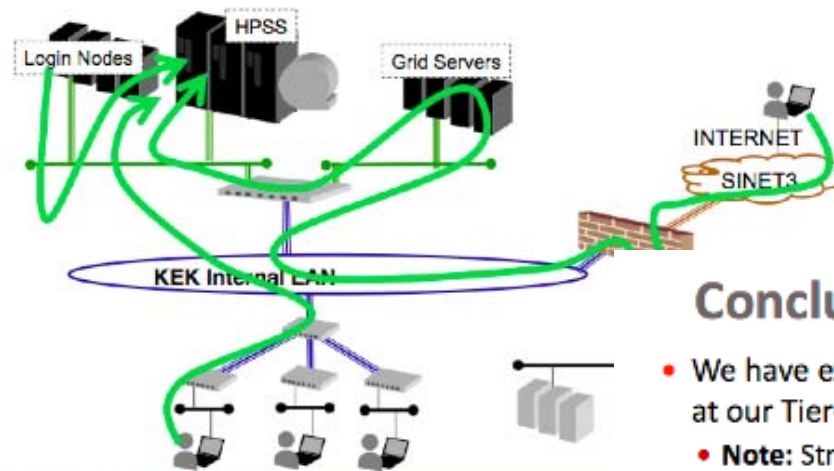- Windows Native Client

## Outline

① Permanently Remove OSS

② Expansion of the Cluster

③ Lustre Bug: OSS Freeze

④ Lustre Malpractice

### First results (16/10/2008)

Number of CMS Events processed during 45 minutes

XROOTD
DCACHE
AFS
LUSTRE

# of jobs

In this summary graph, Lustre seems to be almost twice as efficient compared to the other methods. As AFS, dCache and Xrootd seem to be very close to each other, they may have had a common blocking factor such as the local file system (and NOT the data access protocol).

CERN - IT Department
CH-1211 Genève 23
Switzerland
**www.cern.ch/it**

GSI

CERN

**22**

## Summary

CERN**IT** Department

- Additional requirements for analysis at CERN have been collected and integrated into the development plan for DM components at CERN
  - a initial set-up has been opened rapidly to experiment users based on CASTOR/XROOTD

- Shared long term vision between analysis and data management teams
  - more de-coupled management of tape and disk to prepare for future disk based set-ups
  - moving as close as possible to POSIX semantics and a simple file system view for the end-user
  - plan to resolve remaining questions on MSS integration and end-user files systems in the analysis set-up at CERN and with other storage providers during 2009

- Short term development projects on track to remove remaining constraints - planned deployment in the 2009 LHC run
  - tape aggregation, low-latency I/O scheduling, end-to-end monitoring, meta-data performance & consistency, GC algorithms

# Storage track – day two

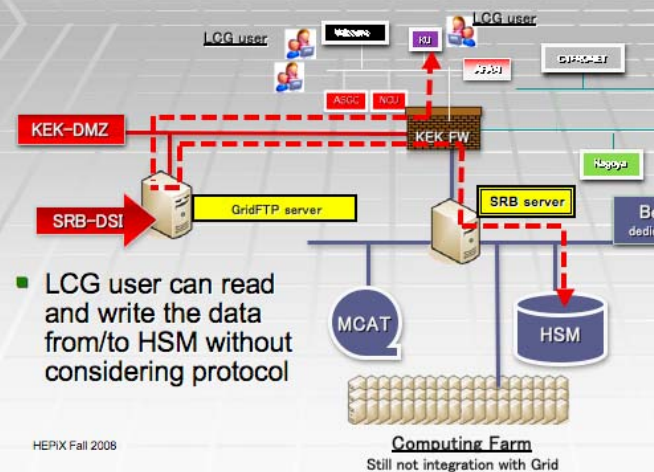| 11:00 | [123] **Optimizing Common Storage Systems**<br>by Benjeman MEEKHOF<br>(2nd Conference Room: 11:00 - 11:15) |
| --- | --- |
| | [182] **Finding a Practical Distributed Filesystem**<br>by Benjeman MEEKHOF<br>(2nd Conference Room: 11:15 - 11:30) |
| | [124] **A data grid environment with HPSS and GridFTP at KEK**<br>by Satomi YAMAMOTO<br>(2nd Conference Room: 11:30 - 12:00) |
| 12:00 | [125] **SRB system for Belle experiment**<br>by Yoshimi IIDA<br>(2nd Conference Room: 12:00 - 12:30) |

Uni Michigan tried AFS, NFS, dCache, Lustre and xrootd. Lustre may be revisited after important changes in 1.8 and 2.0

## Conclusions

- We have examined a few options for our grid-storage needs at our Tier-2
  - **Note:** Strengths and weaknesses are based on our own Tier-2 situation
- dCache was selected given the built-in SRM support and broad use with ATLAS.
- Very interested in NFSv4 and Xrootd options in the future and we may end up with two systems for meeting different needs...
- Other interesting options not explored: GlusterFS, GPFS, ...
- **We welcome comments about where perhaps we are doing things wrong, or if we as a site have not explored software that could better meet our needs.**
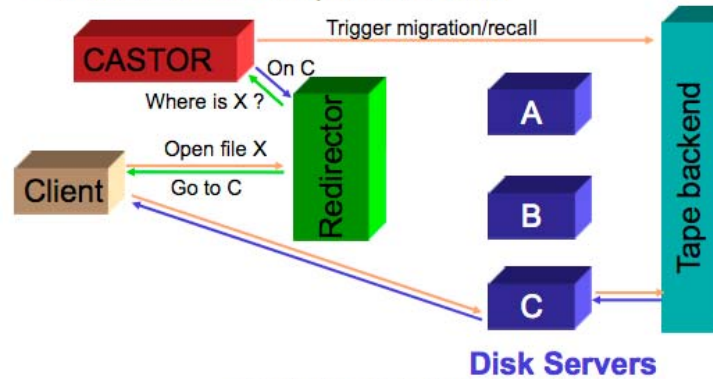
Finding a Practical Distributed Filesystem

10/21/2008

# CASTOR track – day three

11:00

[133] **CASTOR status and plans**
by Sebastien PONCE
(2nd Conference Room: 11:15 - 11:35)

[171] **SRM2 and monitoring projects in CASTOR**
by Dr. Giuseppe LO PRESTI (CERN)
(2nd Conference Room: 11:35 - 11:55)

[170] **CASTOR Operational Experiences**
by Miguel COELHO DOS SANTOS
(2nd Conference Room: 11:55 - 12:15)

12:00

[181] **Increasing Tape Efficiency**
by Dr. Steven MURRAY (CERN)
(2nd Conference Room: 12:15 - 12:35)

# CASTOR

## XROOT in CASTOR

- Client connects to a redirector node
- The redirector asks CASTOR where the file is
- Client then connects directly to the node holding the
- CASTOR handles tapes in the back



## Release 2.7 Key Features

- Support for SLC4
- Support for multiple frontends and backends
- Support for the SRM methods agreed in the WLCG MoU addendum (May '08)
  - **srmPurgeFromSpace** to allow better handling of Disk1 storage classes
  - **srmCopy** fully functional including support of source space token
  - **srmChangeSpaceForFiles** definitely deprecated
  - improved **srmAbortRequest**
- Improved management of space tokens
  - Admin tools for the service manager
  - Space allocation remains static
- Improved logging

## Problem Areas

- **Work done** • Write more data per tape mount

- **Current work** • Use a more efficient tape format

  The current tape format does not deal efficiently with small files

- **More ideas** • Improve read efficiency

  Require modifications from disk to tape

## Recent activities

- CCRC-08 was a successful test
- Introducing data services piquet and debugging alarm lists
- Ongoing cosmic data taking
- Created REPACK instance
  - Needed a separate instance to exercise repack at larger scale
  - Problems are still being found and fixed
- Created CERNT3 instance (local analysis)
  - New monitoring (daemons and service)
  - Running 2.1.8 and new xroot redirector
  - Improved HA deployment

8

27

# Storage Summary

- CASTOR
  - 5 Talks in total. Replay next week.
  - Analysis, xroot, scalability, SRM, tape performance, monitoring
- Lustre
  - 3 Talks
  - Used by 6 of top10 supercomputer sites, 40% of top100
  - Client cache and kerberos on the way
  - Showing good results, some important issues to be addressed (MDS scalability for example)
- Industry
  - 2+1 Talks (Data security & Performance + Lustre/SUN)
- Other
  - 4 Talks
  - Performance
  - Grid enabling (attaching gridftp servers to) existent solutions

**FIO**

**CERN IT Department**

| | |
|---|---|
| 10:00 | **[139] CPU Benchmarking at GridKa - Update**<br>by Manfred ALEF<br>(2nd Conference Room: 10:30 - 11:00) |
| 11:00 | **[140] Report from the HEPiX Benchmarking Working Group**<br>by Helge MEINHARD<br>(2nd Conference Room: 11:00 - 11:30) |
| | **[141] Power Efficiency of Servers**<br>by Helge MEINHARD<br>(2nd Conference Room: 11:30 - 12:00) |

**FIO**

**CERN IT Department**

## CPU Performance Measurements

KIT — Karlsruhe Institute of Technology

- ■ **New CPU types:**
  - ■ Intel E5430 (2.66 GHz Harpertown)
    - ■ Worker nodes at GridKa
    - ■ 2 batches from different vendors
      - ■ Lot #1:
        - ■ System: IBM x3550
        - ■ RAM: 8x 2GB DDR2-667 FB-DIMM
      - ■ Lot #2:
        - ■ System: Supermicro X7DCL-i
        - ■ RAM: 4x 4GB DDR2 reg. ECC
  - ■ AMD Opteron 2356 (2.3 GHz Barcelona)
    - ■ Test system at GridKa
      - ■ Hardware details:
        - ■ System: Supermicro H8DMU+
        - ■ RAM: 8x 2GB DDR2 reg. ECC

## Power Consumption of Cluster Nodes

KIT — Karlsruhe Institute of Technology

- ■ **New CPU types:**
  - ■ Intel E5430 (2.66 GHz Harpertown, **80W TDP**)
    - ■ Worker nodes at GridKa
    - ■ 2 batches from different vendors
      - ■ Lot #1:
        - ■ System: IBM x3550
        - ■ RAM: 8x 2GB DDR2-667 FB-DIMM
      - ■ Lot #2:
        - ■ System: Supermicro X7DCL-i
        - ■ RAM: 4x 4GB DDR2 reg. ECC
  - ■ AMD Opteron 2356 (2.3 GHz Barcelona, **75W wattage**)
    - ■ Test system at GridKa
      - ■ Hardware details:
        - ■ System: Supermicro H8DMU+
        - ■ RAM: 8x 2GB DDR2 reg. ECC

| Year in | Description | Power efficiency |
|---|---|---|
| | GHz, 2 GB, 1 disk, tower | 6.95 |
| | GHz, 2 GB, 1 disk, 7520 chipset, 4U | 6.76 |
| | GHz, 2 GB, 1 disk, 7320 chipset, 4U | 7.01 |
| | GHz, 4 GB, 1 disk, 7320 chipset, 4U | 6.63 |
| | GHz, 2 GB, 1 disk, 7320 chipset, 1U | 6.39 |
| | GHz, 2 GB, 1 disk, 7320 chipset, 1U | 6.40 |
| | .66 GHz, 8 x 1 GB, 1 disk, 5000P, 1U | 13.98 |
| | .00 GHz, 4 x 2 GB, 1 disk, 5000P, 1U | 15.02 |
| | .00 GHz, 8 x 1 GB, 1 disk, 5000P, 1U | 15.17 |
| | .00 GHz, 8 x 1 GB, 1 disk, 5000P, 1U | 22.30 |
| | 33 GHz, 8 x 2 GB, 2 disks, 5000P, twin | 29.02 |
| | 33 GHz, 8 x 2 GB, 2 disks, 5000P, twin | 28.60 |
| | 33 GHz, 8 x 2 GB, 2 disks, 5000P, twin | 29.50 |
| | 2.33 GHz, 8 x 2 GB, 2 disks, 5000P, blade | 28.22 |
| | 2.33 GHz, 8 x 2 GB, 2 disks, 5000P, blade | 35.60 |
| | 2.33 GHz, 8 x 2 GB, 2 disks, 5000X, blade | 37.81 |
| | 5420, 4 x 4 GB, 2 disks, 5100, twin | 54.71 |
| | 5420, 4 x 4 GB, 2 disks, 5100, twin | 56.97 |
| | E5410, 4 x 4 GB, 2 disks, 5100, twin | 50.48 |

# Benchmarking WG

## Agreements

CERN IT Department

- Focus on benchmarking of processing power for worker nodes
- Representative sample of machines needed; centres who can spare a machine (at least temporarily) announce this to the list
- Environment to be fixed
- Standard set of benchmarks to
- Experiments to be invited (i.e. code
  - Check how well experiments' co standard benchmarks

## Choice of Benchmark (

- SPECall_cpp2006 adopted
- Environment: WG proposed
  1. SL4 x86_64
  2. System compiler gcc 3.4.6
  3. Flags as defined by LCG-SPI architects' forum)
     `-O2 -fPIC -pthread -m32`
  4. Multiple independent parallel runs

## Conclusions

CERN IT Department

- WG has reached its primary goal – propose a new standard benchmark for HEP
- Good participation of LHC experiments – thanks!
- Thanks to all collaborators for their active, constructive participation
- Special thanks to Ian Gable for acting as secretary
- Gonzalo's working group looking at WLCG-specific aspects
- Next steps for HEPiX WG:
  - Presentation at CHEP in Prague (abstract: Michele Michelotto)
  - Statistical analysis of results
  - Writeup

# Report on
# HEPiX Fall 2008 Taipei

Operating Systems and Applications

Virtualization

Network and Security

Miscellaneous

Steven Murray, 21 November 2008                    Slide 33

- **Subversion(SVN) 2009**

- **Computation interface to NVidia graphics card**

- **Live CD and NTFS(read/write)**

- **Smooth transition to Exchange 2007**

- **Used by Atlas since their first data challenge**

CERN IT Department
CH-1211 Genève 23
Switzerland
**www.cern.ch/it**

Steven Murray, 21 November 2008

Slide 34

**CERN IT Department**

## DES CERN Central CVS Service

**CERN IT Department**

- Hosts over 330 Software Projects
  - 29 for Atlas
  - 46 for CMS
  - 8 for LHCb,.....
- Over 3000 developers registered
- Over 90 GBytes of source code
- Creates 250 Remedy tickets per year
- Over 100,000 commits per month

CERN - IT Department
CH-1211 Genève 23
Switzerland
**www.cern.ch/it**

4

CERN IT Department
CH-1211 Genève 23
Switzerland
**www.cern.ch/it**

Steven Murray, 21 November 2008

Slide 35

## DES SVN vs. CVS

| Feature | SVN | CVS |
|---|---|---|
| Speed | Faster | Slower |
| Permission | Full | Limited |
| File types | All | Limited |
| Off line operations | Yes | No |
| Repository format | Database | File system |
| Locks | No | Yes |
| Atomic commits | Yes | No |

CERN - IT Department
CH-1211 Genève 23
Switzerland
www.cern.ch/it

8

CERN IT Department
CH-1211 Genève 23
Switzerland
**www.cern.ch/it**

Steven Murray, 21 November 2008

Slide 37

DM

CERN IT Department

## Scientific Linux
### What we see in the future

- When RHEL 5 Update 3 comes out, releasing S.L. 5.3
  - Estimate - March 2009
- When RHEL 4 Update 8 comes out, releasing S.L. 4.8
  - Estimate - June 2009
- When RHEL 6 comes out, releasing S.L. 6.0
  - Estimate - January 2010

# Virtualization

- **Self service virtual machine in 10 minutes**

- **chroot-based VM still going strong since 2004**

- **Platform Enterprise Grid Orchestrator integrates virtualisation and the Grid**

- **Hyper-V and XenCenter can work together**

- **CISCO and F5 solutions for network virtual network services**

Current Server Self Service

Select an OS, type in a budget code and click Request. 10 minutes later, the user will receive an email notifying that his server is available.

CERN**IT**
Department

## Virtual Server 2005 vs. Hyper-V

CERN**IT**
Department

- Functionality:

| | Virtual Server 2005 | Hyper-V |
|---|---|---|
| Hypervisor Type | Hosted | Native |
| 64-bit Virtual Machines | No | Yes |
| Multi Processor Virtual Machines | No | Yes, up to 4 core |
| Virtual Machine Memory Support | 3.6 GB per VM | 64 GB per VM |
| Scriptable / Extensible | COM | WMI |

- Performance!
  - Live example: Windows Media Server based on Hyper-V can host more clients with less resources – proved during the LHC First Beam event
  - Enlightened I/O for storage, networking and graphics
    - Makes use of VMBus (instead of using the emulation layer) while using some high level protocols like SCSI
    - Supported for guests running: Windows Vista, 2008 and SUSE Linux

CERN IT Department
CH-1211 Genève 23
Switzerland
www.cern.ch/it

## Comparison to conventional VMs

- Pros:
  - Easy to implement for the current grid middleware
  - Negligible resources overhead
  - Works on all hardware
  - 4 years of production reliability
  - Can be used for interactive work as well (Atlas)
  - Can support most recent hardware
- Cons:
  - Shared kernel
  - No resources isolation/virtualization apart from the underlying filesystem
  - Only works for linux
  - Shared grid accounts with the basesystem
  - Security concerns for possible chroot-jail breaks

2008/10/23      Grid Infrastructure at SiGNET and Chrooted Systems      16

- **2 CAs, 1 RA and smart cards**

- **DBPowder will become open source**

- **Codian MCU will in production February 2009**

- **Good progress during last year**

- **Only one solution is long-term – the transition**

- **Cyber space isn't getting any safer**

CERN IT Department
CH-1211 Genève 23
Switzerland
**www.cern.ch/it**

Steven Murray, 21 November 2008

Slide 43

# Miscellaneous

- **VO communities detect problems first**

- **RT is good for out-of-the-box usage**

- **Tier-2 praguelcg2 will use cfengine everywhere**

- **GPFS usage recommendations made**

- **Standards are the only long-term solution**

CERN IT Department
CH-1211 Genève 23
Switzerland
**www.cern.ch/it**

Steven Murray, 21 November 2008
Slide 44

**DM**

CERN IT Department

## Conclusions

- With GPFS, there are a variety of methods available to connect to a GPFS filesystem

- We have evaluated a number of ways of providing high performance access to GPFS filesystems, ranging from direct access to GPFS filesystems (F/C) to the gateway model (using 10GE and Infiniband)

- The method chosen depends on considerations of
  - Matching the back-end B/W of the serving cluster
  - Cost considerations
  - Interconnect availability on the client cluster

- Using Infiniband networks (already present for other purposes in the cluster) to transport I/O allows us to
  - Fully utilize the IB fabric
  - Cost-effectively scale the interconnect

DM

CERN IT Department

## Conclusions

**ERSC**
NATIONAL ENERGY RESEARCH
SCIENTIFIC COMPUTING CENTER

- For small/medium size clusters, the gateway model provides a scalable, low-initial-investment way of providing high-performance access.
  - Additional connectivity can be provided with small incremental costs

- Access over 10GE networks has been shown to be efficient – however, the large scale deployment of 10GE to clients may be limited by cost considerations at this time

- Access over F/C is also highly efficient, and is useful for small clusters – however, costs of connecting to F/C fabric have to be considered.

Office of Science
U.S. DEPARTMENT OF ENERGY

BERKELEY LAB

CERN IT Department
CH-1211 Genève 23
Switzerland
www.cern.ch/it

Steven Murray, 21 November 2008

Slide 46