# *Design, Status, and Experience with BaBar LTDA*

Concetta Cartaro
SLAC

On behalf of

*BaBar* Computing Group

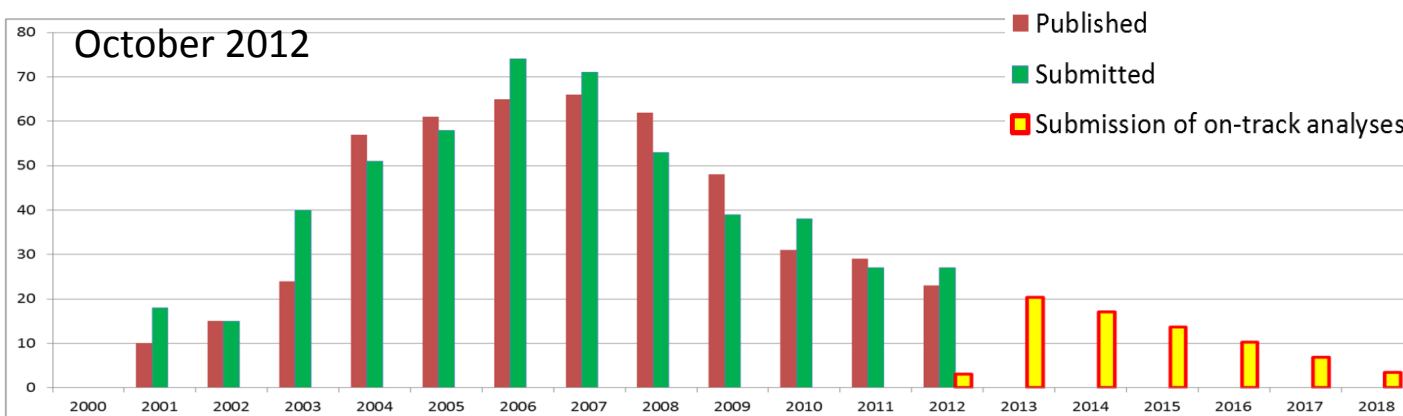DPHEP Workshop – Marseille, November 19th, 2012

# *OUTLINE*

- BaBar data and choices for the future
- The Long Term Data Access project
- Update since May 2012
  - DPHEP hosted @ CHEP 2012
  - `https://indico.cern.ch/getFile.py/access?contribId=12&resId=0&materialId=slides&confId=171962`
- Managing the cluster
- BaBar long term planning
- Conclusions

# BABAR DATA

- BaBar has collected data from Oct 22nd 1999 to Apr 7th 2008
  - 800TB of raw data, 1.2 PB from the last data reprocessing
  - 500th paper accepted on October 16th
    - 508 papers published/accepted/submitted
      - 31 published/accepted in 2012 – 2 more than in 2011
  - 74 on track analyses
    - Plus ~30 analyses progressing slower (generally lacking manpower)
    - Possibilities for new previously unforeseen analyses including discovery analyses
- BaBar (and Belle) data will not be superseded by LHC data
  - Belle II and SuperB will do it in 5-10years
  - Some datasets expected to remain unique for longer (Y(3S) dataset)



October 2012

Legend:
- Published
- Submitted
- Submission of on-track analyses

SLAC NATIONAL ACCELERATOR LABORATORY

# *LONG TERM DATA ACCESS*

- Insure the ability to support 50 to 80 analysis of the BaBar data until at least 2018 preserving:
  - Data, conditions and calibrations, releases and tools, databases, capability of running production and user jobs
    - This means that in 5 years from now it will be possible to add a new decay mode, produce the MC events and the relevant skims, and perform a completely new analysis developing new selection code, fitting procedures, etc.
  - Documentation
- Providing a stable environment
  - Last validated OS enclosed in a virtualization layer running the BaBar Framework minimizing the effort needed to maintain the system
    - → Need to address: hardware support, security risks, keep know-how on OS, Framework, etc…
- Open formats
  - Data format is based on ROOT which is open and will be part of the system
  - Databases will move away from Oracle and will be stabilized on MySql
  - Code is written in open formats: C/C++, Tcl, Perl, Python.
- Data Storage
  - 2PB (including raw data will be stored on tape in two Tier A sites (SLAC, CC-IN2P3)
  - Most used data will sit on disk

SLAC NATIONAL ACCELERATOR LABORATORY

# *THE LTDA CLUSTER FACTS (I)*

- Cisco 6506 network switch with 2x10Gb link card and 192Gb ports

- 9 infrastructure servers (Dell R410/R510)
  - 3 front end machines (`bbrltda` load balanced pool), 1 cron server, 1 test server, 2 infrastructure servers (network and identification services), 2 database servers (mirrored)

- 54 batch and storage servers
  - Dell R510: dual 6-core Intel Xeon X5675, 3.07GHz, 48GB RAM, 12x2TB disks
  - 4 were the prototype (dual 6-core Intel Xeon  X5670, 2.93GHz, 48GB RAM)
  - 11x2TB  disks (no raid) used to stage data through XROOTD
  - 1x2TB used as local scratch
  - 12 physical cores, 24 cores with hyper threading
    - 1 physical core  used for the host and the XROOTD services
    - 11 cores (22 w/ hyper-threading) dedicated to batch with one VM per core

C. Cartaro @ DPHEP

SLAC NATIONAL ACCELERATOR LABORATORY

# *The LTDA Cluster Facts (II)*

- 20 additional batch servers (new!)
  - Dell R410: dual 6-core Intel Xeon X5675, 3.07GHz, 48GB RAM, 2x2TB disks mirrored (for OS + local scratch)
  - 12/24 cores used to run batch jobs (VMs)

- 2 NFS servers (1 new!)
  - Sun X4540 Thor server: 12 cores, 32 GB memory and 32TB of effective storage
  - One for local home directories and code repositories and one for user data

- The LTDA cluster is in production mode since March 21st 2012
  - On time and on budget
  - 1.33 PB of disk space for data and users
  - 1668 job slots
  - SL4, SL5, SL6 platforms available

- All active BaBarians have a 1GB home directory on the LTDA

SLAC NATIONAL ACCELERATOR LABORATORY

20 Batch Servers
(no XROOTD)

4 Prototype
Servers
(batch+XROOTD)

50 Batch and
XROOTD Servers

Back

Switch

9 Infrastructure
and Login
Servers

2 NFS Servers

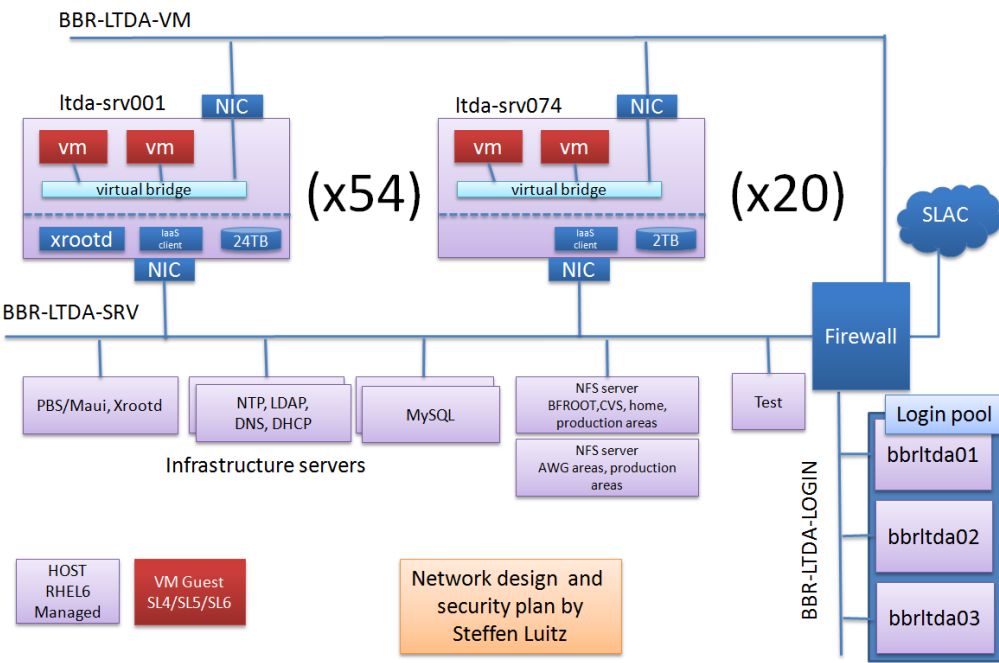C. Cartaro @ DPHEP

**SLAC** NATIONAL ACCELERATOR LABORATORY
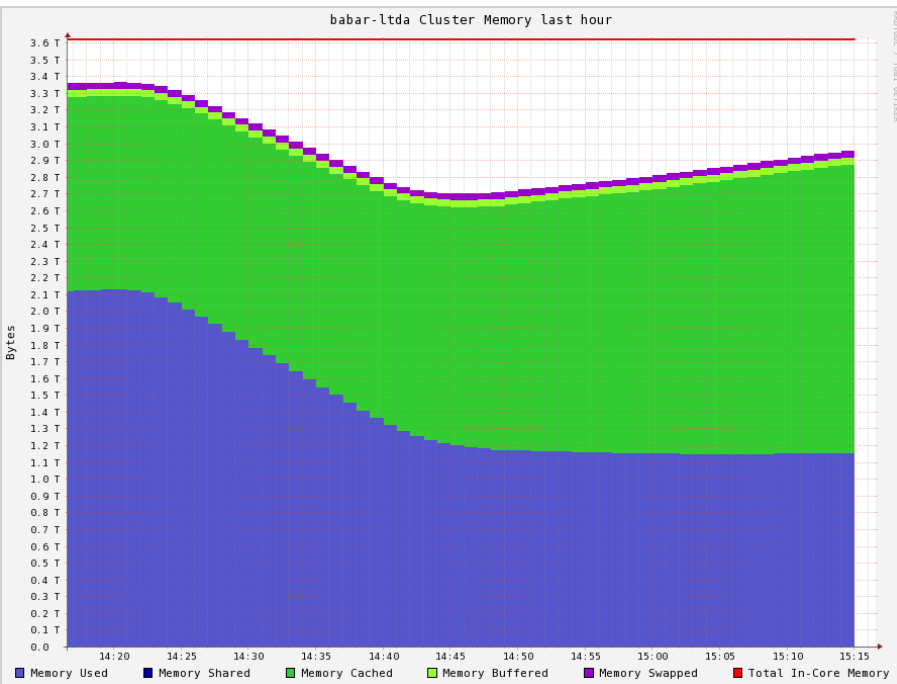
# VIRTUALIZATION & NETWORK

- Security threat associated to a VM connected to a network running old OS
  - Images are read-only, qcow2 produces a temporary file with changes to OS and scratch area and it is deleted when the VM's shut down

→ **Risk based approach assuming that the VMs are compromised**
  - Isolation of back versioned components with firewall rules
  - Physical hosts centrally managed by SLAC CD



- VMs are not allowed to connect to SLAC network or the world
- The Login network is protected from the VM network
  - Allow one way ssh from Login to VM network
  - VMs are not allowed to write over the Login network
- Well defined services between VM network and SRV network
  - Infrastructure (DNS, LDAP, NTP), file service (Xrootd, nfs), batch scheduling
  - LDAP is a subset of the SLAC Kerberos list mapped on /nfs internal home directories
- Allow SRV and Login networks use SLAC infrastructure

# *MEMORY USAGE*



On a sample analysis with about 1500 parallel jobs the memory usage was reduced from 2.1TB to 1.2TB!!!
Freed memory is used for caching files, but can also be used for memory intensive jobs (merging skim output needs at least 4GB RAM).

- 48GB of RAM/server
  - 24 VMs with 2GB RAM
  - 22 VMs on machines with storage space
    - One physical core left for xrootd
- RAM is also needed for the system itself
- Deduplication for identical blocks already used on filesystems
- "Kernel Samepage Merging" (KSM) introduced in kernel 2.6.32
  - Same memory pages are merged together into a single one among different processes!
  - most effective for a lot of identical processes
    - that's VMs!

Marcus Ebert

C. Cartaro @ DPHEP

SLAC NATIONAL ACCELERATOR LABORATORY

# *NFS, ZFS AND BACKUPS*

- NFS servers
  - 40TB zpools, 2 hot spares for 32 TB usable space
  - Compression enabled on /home (factor ~2 gain)
- ZFS snapshots implemented for /home and /BFROOT (releases, packages and cvs root directory) for user error recovery
  - snapshots are read-only, so it's protected against user error
  - frequent snapshots every 15min, overwritten every hour; the frequent full hour snapshot becomes the hourly snapshot after the next hour; the midnight hourly snapshot becomes a daily snapshot at next midnight; daily snapshots are overwritten every month; daily backups of the last 30 days snapshots are recursively created for all zfs under the given top zfs in one single operation
    - BFROOT is more slowly changing and no 15 min snapshot is implemented for it
- Tape backup for catastrophic events
  - All areas are backed up to tape every day and kept for 30 days
    - Root files are omitted because they are considered reproducible

Marcus Ebert

C. Cartaro @ DPHEP

SLAC NATIONAL ACCELERATOR LABORATORY

# *JOB SUBMISSION*

- PBS/Torque is used to manage the batch resources and Maui is the batch scheduler

- The virtualization layer uses qemu with kvm support directly
  - Moved away from libvirt due to instability

- Need to create the network interface for the VMs
  - 24 MAC addresses per host and usage status stored in local db

- PBS Prologues and Epilogues scripts are used to create and destroy the VM's and the needed network environment



Marcus Ebert, Kyle Fransham – LTDA developers

SLAC NATIONAL ACCELERATOR LABORATORY

# A LITTLE MORE DETAIL ...

User submits a job from one of the login machines

`sleep 60`

Submit Filter

**PBS-Server** collects the job, creates a script that will restore the user environment and then executes the user job.

Scheduler

**Maui** Scheduler of the system, knows where and when to allocate jobs

PBS –Client Ltda-srv002

PBS-Client Ltda-srv003

PBS clients on the ltda-srv0xx batch servers

**EXECUTE HOST**

PBS-Client on the ltda-srv001

**TARGET VM**

`sleep 60`

**VM** runs the user job like a normal batch machine with the environment from the submit host

**Prologue script** creates network interface, looks up free MAC in local database, creates VM with specified parameters and network, sends job script to the VM via ssh

**Epilogue script** Destroys the VM and the network interface , copies the log to the final location and frees the MAC address in the local database

SLAC NATIONAL ACCELERATOR LABORATORY

# MONITORING



Daily

System check

Jobs running

Queues check

Marcus Ebert

SLAC NATIONAL ACCELERATOR LABORATORY

# *IMPROVING USER EXPERIENCE*

- Usually submit a job and wait for the result

- Limitations are due to the restrictions on the VMs

- Interactive VMs for SL4,5,6 platforms are always available to the users in a natural way
  - `ssh sl4` will redirect passwordlessly (shosts) to an interactive VM running SL4
    - No waiting, no special commands

- Skimming tools adapted for user case
  - Create individual skim cycles
  - Easy to define and handle large amounts of jobs
  - All jobs info stored in the database for reproducibility and traceability

Marcus Ebert – LTDA developer
Douglas Smith – BaBar tools and db expert

SLAC NATIONAL ACCELERATOR LABORATORY

# *WHAT GOES ON THE LTDA*

- LTDA is dedicated to both users and production
  - Both simulation production and skimming run on the system at a constant level (200 slots) when needed
  - One server is dedicated to our OPR (raw data reconstruction application) processing
- The LTDA resources will help BaBar face both the loss of resources within SLAC
  - Dedicated production resources will disappear soon
    - Oracle did not renew the maintenance contract for >5years old SUN equipment
    - Machines turned off as newer hardware arrives (to recover LSF licenses)
- …and the loss of the TierA sites
  - Right now all TierA sites are still supporting BaBar analysis but this may already not be true in the near future
    - INFN Padova will shut down the tape library with raw data backup at the end of the year, and other two sites will run out of BaBar funds at middle/end of 2013

# PERFORMANCE TESTS

- VM's vs SLAC batch queue
  - CPU X5670 @ 2.93GHz vs (x5355 @ 2.66GHz and x5570 @2.93GHz mixed)
  - RH5, HT off → 2.7% in favor of LTDA

- Hyper-threading on/off
  - 40% slower with HT on, but SL6 faster than SL5 by 35% CPU time, 15% wall time

- I/O performance tests
  - XROOTD tested for heavy load and scalability
  - The capability of the cluster exceeded any possible demand of BaBar applications

SLAC NATIONAL ACCELERATOR LABORATORY

# PROBLEMS & SOLUTIONS

- Intrinsic dependency between services running on different machines forced the boot order of the servers
  - Often caused delays during outages
  - **Solution:** Remove the dependencies and use automount on all servers to make the cluster independent of the boot order

- Red Hat updates delivered to all hosts automatically
  - This caused long outages in in the past
    - Kernel update with network bug: VMs not reachable
    - Automount bug caused crashes when used with LDAP
  - **Solution:** Develop a validation system in order to test the updates before delivering them to all hosts
    - Ltda-test server available for testing, validating and releasing updates to the whole cluster
    - Remove not essential software packages to reduce the list of updates

Marcus Ebert

SLAC NATIONAL ACCELERATOR LABORATORY

# DOCUMENTATION

- Strong push toward documentation clean up, ease of access, and clarity

- All most used and fundamental info are being checked, updated and moved to a Media Wiki server, the *BABAR WIKI*
  - Old pages clearly marked but kept online for archival purposes
  - Detector pages and other pages that will supposedly never change again will be left in their original location

Wiki main page

# DOCUMENTATION WORKING GROUP

- The effort needed is not trivial

- The Documentation Working Group is coordinating the migration effort aided by an advisory committee

  – Many new students joined the effort but the input from senior members of the Collaboration is fundamental

  – There are 10 official members (plus some less official...) in the DWG but we promote the migration to the wiki as a Collaboration effort

  – Experts sign-off on the content of migrated pages

Matt Bellis, Alessandra Filippi – DWG coordinators

C. Cartaro @ DPHEP

SLAC NATIONAL ACCELERATOR LABORATORY

# *BABAR STILL ROCKS!*



http://www-public.slac.stanford.edu/babar/

BaBar public page: Abi Soffer, former PAC & senior DWG Member

arXiv:1207.5832

SLAC NATIONAL ACCELERATOR LABORATORY

# *MANPOWER, EXPERTISE & BUDGET*

- Designing and maintaining something like the LTDA through the years requires many talents and careful planning
  - *BABAR* experts
    - Releases, databases, data management and documentation
      - Plus virtualization support
    - The Collaboration will have to provide such expertise
      - Difficult to disentangle general BaBar support and LTDA related support, but at the moment 1-1.5 FTE is a good estimate for LTDA only right now
  - Computing experts
    - Network, security, system and networks administration, …
    - 0.5 FTE/year foreseen after 2012
- Costs (not FTEs)
  - Hardware and refreshment program
  - Recharge (somewhat unknown) costs
  - Red Hat entitlements for virtualization
    - Use SL virtual machines to avoid extra cost (SL is not supported at SLAC)

SLAC NATIONAL ACCELERATOR LABORATORY

# *LTDA BEYOND THE CLUSTER*

- Future of the Collaboration and the data
  - Time to start thinking about what we will do in the future and what will happen to our data and internal documents
    - Go public? If yes, when and how?
    - Use of Inspire
  - Very strong opinions within the Collaboration
  - The challenge will start at our next Collaboration Meeting at the end of January

C. Cartaro @ DPHEP

SLAC NATIONAL ACCELERATOR LABORATORY

# ACKNOWLEDGMENT PAGE

- Thanks to the LTDA developers
  - Coordinator:  Tina Cartaro
  - BaBar software expert: Homer Neal
  - Development and system administration of LTDA: Marcus Ebert
  - Network design: Steffen Luitz
  - Virtualization expert:  Kyle Fransham
  - System performance and CDB: Igor Gaponenko
  - Databases, tools and production: Douglas Smith and Tim Adye
  - Computing Division experts
    - System setup and adminstration: Booker Bense, Lance Nakata, Randall Radmer and all the Unix-Admin team
    - Xrootd expert: Wilko Kroeger
    - Network setup: Antonio Ceseracciu
    - BaBar-SLAC Computing Division liaison: Len Moss
  - Thanks to the PPA Management and the DOE for the strong support

SLAC NATIONAL ACCELERATOR LABORATORY

# CONCLUSION

- System went into production on March 21st
  - On time and within budget
  - Since then it has been already extended by 20 batch nodes and one more nfs server
- All BaBarians have a LTDA home
  - About 10-15 users very active
- In use for production
  - Behaves like a Tier site
- Improving the setup and fine tuning
  - Unexpected problems from updates took the system down for many days in two occasions
  - Improving user experience
  - Simplify things wherever possible

SLAC NATIONAL ACCELERATOR LABORATORY

# *Q & A*

- Why LTDA and not the GRID or the commercial clouds?
  - Code and data is directly accessible from the common NFS area and don't need to be included in the VM image thus allowing memory and startup efficient VM's and easy use of new releases and access to a global CVS repository.
  - Files are directly accessible from the NFS areas instead of needing to copy them to the VM space. A job has the same local read access privileges as those of the user.
  - Many research computing clusters using virtualization allow only certain directories to be mounted for importing files thus resulting in extra effort from the user to create a setup appropriate for batch submission which is often different from the development setup.
  - One can log into the VM that the job is running on and diagnosis resource usage problems.
  - GRID systems are frequently very difficult to debug because one can not directly observe the job as it is processing.
  - One can start interactive VM sessions for development/debugging work.
  - Users will not have to invest time and effort to acquire funds for using a commercial cloud system.
  - A platform equivalent to that for the batch jobs will always be available for development and debugging purposes.
  - The setup allows users to migrate seemlessly to the LTDA system with only a few minor restrictions on what one can do.
- Why doesn't one just use the standard batch system but with VM's?
  - The whole SLAC batch system would have to be adjusted to be behind a firewall to protect against use of insecure platforms and other projects running behind the firewall could be affected.
  - New tools and technologies may allow it in the future
- The lifetimes of the latest RHEL releases have been extended, doesn't it remove the need for the LTDA?
  - No because the expertise and person power will cease to exist to do full release validations and possible update of the hundreds of packages of code every time a new security patch (typically coming along with a new version of glibc) is released.
- What about dependencies on the virtualization system?
  - Look for alternatives (xen, …) and ultimately use hardware emulation
- What were/are the instabilities of libvirt?   ie, why is qemu called directly?
  - There have been at least two things which made problems: 1) libvirt crashed often with segfault with no obvious reason;  2) jobs just hang in the queue after the real job within the VM was already finished (VM was for unknown reasons not destroyed and the job not cleared from the queue system after finishing). Sometimes only the VM was killed, but not due to the prologue script but by a direct kill to the process which didn't freed up the MAC address and let the ssh processes running. Since we changed to call qemu directly, this never happened again. Also the structure is as simple as with libvirt, maybe even simpler since we don't need an additional layer and the use case within the LTDA is very limited. Also using qemu makes it easy to make changes to the base image using the same prologue/epilogue script structure. (it's easy with libvirt tools somehow too, but it's not needed at all and one more reason not to use it in the LTDA case)

C. Cartaro @ DPHEP

# PAST PROBLEMS & SOLUTIONS

- **qemu couldn't be started**
  - error message: kvm_create_vm: Interrupted system call
  - for about 0.5% of all jobs
  - jobs listed in the queue until wall time is over
  - known bug of kvm/qemu
  - should be solved in new versions
  - expect this problem to be gone on RHEL6.2

- **maui often shuts down without any hint in the log file**
  - seems to happen always with a high load on the scheduler and the network and more than 600 jobs

→ solution:
  - don't allow torque to push jobs to maui
  - only Maui looks every 10s for new jobs
  - for one schedule cycle only 10jobs/user are considered
  - if there are free nodes and waiting jobs, then let maui wait 4s between sending jobs to the nodes
  - users could also put in their scripts a delay of 1s between submission of jobs

- **very high network usage on some server**
  - due to the loading of one condition file in cond24boot09
  - not seen for cond24boot11

→ solution:
  - reduplicate conditions files on more servers

- **input collections or conditions couldn't be found**
  - to many open connections in xrootd
  - Network problems on the xrootd client hosts
  - connection couldn't be established
  - wrong mounted hard disk

→ solution:
  - correctly mount the hard disk on ltda-srv005
  - reduplicate the conditions on many servers to reduce the load on a single one
  - tune the tcp parameters on all ltda-srv0xx
  - use a timeout in xrood for the connections to the clients

- **NFS server stopped to give new nfs exports out**
  - after some runs with more than 1000VM in parallel no new nfs mounts have been possible
    - this includes the home mount using automounter on the login machines for new logins
    - all existing mounts still worked
  - seems like a limit in the nfs server, maybe in open network ports, was reached

→ solution:
  - unix-admin changed some nfs related settings, we will see in the future if it's enough

SLAC NATIONAL ACCELERATOR LABORATORY

# *PERFORMANCE TESTS*

- VM's vs SLAC batch queue

- Hyper-threading on/off

- I/O performance tests

C. Cartaro @ DPHEP

# *VM vs Bare Metal*

| Name of the output file | USER CPU time on the LTDA | VM name on the LTDA | USER CPU time on the normal system | Host name on the normal system(CPU type) | Differences in used time | percentages of cpu time used more on LTDA |
|---|---|---|---|---|---|---|
| LambdaC-Run1-OnPeak-R24c-1.out | 5521 | bbr-ltda-vm049 | 4212 | hequ0167 (Intel(R) Xeon(R) CPU X5570 @ 2.93GHz) | 1309 | 31.07 |
| LambdaC-Run1-OnPeak-R24c-10.out | 5572 | bbr-ltda-vm050 | 6765 | fell0171 (Intel(R) Xeon(R) CPU X5355 @ 2.66GHz) | -1193 | -17.63 |

~1500 jobs later …

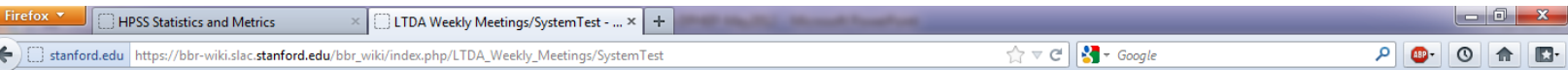| | | | | | | |
|---|---|---|---|---|---|---|
| LambdaC-Run6-OnPeak-R24c-99.out | 5274 | bbr-ltda-vm037 | 6862 | fell0249 (Intel(R) Xeon(R) CPU X5355 @ 2.66GHz) | -1588 | -23.14 |
| all | 9019444 | | 9272133 | | -252689 | -2.72 |

LTDA host machines:
Intel(R) Xeon(R) CPU X5670 @ 2.93GHz (prototype machines)
Rhel5 / Hyper-threading off

SLAC NATIONAL ACCELERATOR LABORATORY

# *HYPER-THREADING TESTS*



Marcus Ebert

C. Cartaro @ DPHEP

# COMPARISON TABLES

- CPU intensive jobs show no difference when single job or 11 jobs run on a machine
    - but variation for repeated tests is up to 30%
- I/O intensive jobs can be up to 300% slower when 11 jobs run in parallel on a single machine compared to a single job/machine
- BetaMiniApp is about 20% slower when running 11 jobs in parallel on a single machine compared to a single job/machine
- CPU time used for all Run1-6 jobs when running 11jobs/LTDA machine in parallel is comparable to running same jobs on the general SLAC queue
    - LTDA used about 2% less CPU time for all jobs

- CPU time used for a BetaMiniApp job is comparable to HT off when running same number of jobs in parallel
- single BetaMiniApp can use up to 50% more CPU time when running 22jobs in parallel instead of 11
- CPU intensive again independent of the number of parallel jobs
    - <10% slower than HT off
- I/O intensive jobs can be up to 900% slower when running 22jobs/machine compared to a single job

**HT off vs HT on**

|  | 11 jobs SL5 NFS image | 22 jobs SL5 NFS image | difference |
|---|---|---|---|
| CPU time | 7849934 | 11243061 | +40% |
| Wall time | 8697739 | 11996994 | +38% |

|  | 11 jobs SL6 NFS image | 22 jobs SL6 NFS image | difference |
|---|---|---|---|
| CPU time | 5152806 | 7286484 | +41% |
| Wall time | 7245687 | 10267120 | +42% |

**SL5 vs. SL6**

|  | 11 jobs SL5 NFS image | 11 jobs SL6 NFS image | difference |
|---|---|---|---|
| CPU time | 7829719 | 5120721 | -35% |
| Wall time | 8371660 | 7200235 | -14% |

|  | 22 jobs SL5 NFS image | 22 jobs SL6 NFS image | difference |
|---|---|---|---|
| CPU time | 11233998 | 7249142 | -35% |
| Wall time | 11987646 | 10214567 | -15% |

C. Cartaro @ DPHEP

SLAC NATIONAL ACCELERATOR LABORATORY

# RESULTS



True in general, not LTDA specific

## conclusion

- CPU intensive jobs show no dependency on the number of parallel jobs
- I/O intensive jobs depend heavily on the number of parallel jobs
  - that's expected
- for BetaMiniApp jobs:
  - running the same number of jobs in parallel shows no difference between HT off and HT on
  - running 11jobs/machine shows no difference to the general SLAC queue
  - running 22 jobs/machine slows down single jobs up to 50% compared to 11jobs/machine
  - the difference for CPU/Wall time between usage of NFS or local images is only very small
    - for using NFS images the network load on wain062 is higher
    - but no problem so far since it has a 4G etherchannel and peak value was around 3.2G
    - could become a problem when adding more servers
  - running the same binary on SL6 instead of SL5 reduces CPU time by about 35%
  - using all cores for VM's slows down single jobs compared to 11(HT off) or 22(HT on) jobs/machine
  - time to finish all Run1-6 jobs using 22jobs/machine is about 15% shorter than for 11jobs/machine
    - this number depends heavy on the number of jobs
    - for no HT ~1500 jobs mean ~2times full load while for HT on it's only 1x full load (+ running only some jobs/machine for both)
    - time difference will be much smaller if one uses only 500 jobs but with much more events processing/job
    - time difference will be much larger if one uses 3000 jobs but with much less events processing/job
    - difference can be between 0% and ~30%

## proposal for the final system

- use HT on
- allow 22jobs/machine
- use NFS images
  - switch to local images could easily be done if we see problems with more servers
- switch to a release which can be build and run on SL6
- reduce the wall time for the general queue again to let not run jobs with too many events processing
- repeat some tests once RHEL6.x is installed on all machines

C. Cartaro @ DPHEP
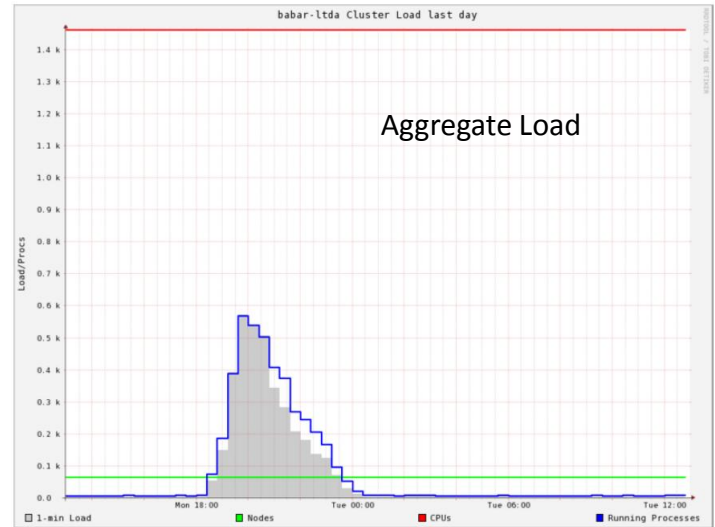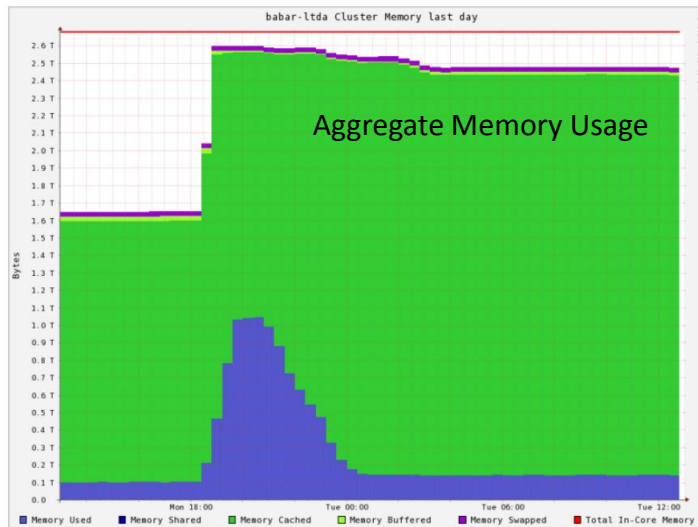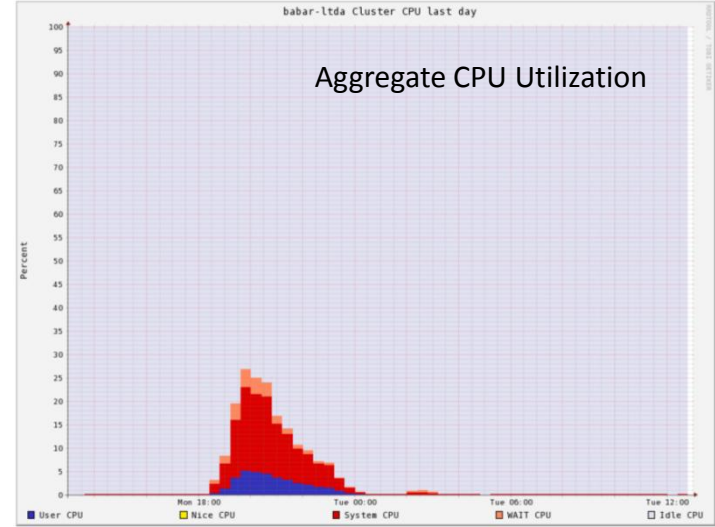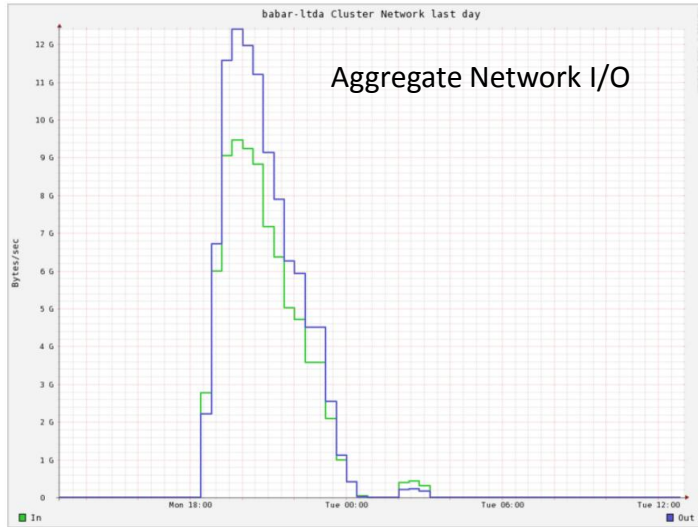
SLAC NATIONAL ACCELERATOR LABORATORY

# I/O Performance test For Xrootd

- The main goal of the test was to see how much data can be delivered under extreme load by LTDA XrootD installation to clients processes. Resource (memory and CPU) utilization was also being monitored during the tests using Ganglia. Scalability of the XrootD installation has been tested as well.

Igor Gaponenko

SLAC NATIONAL ACCELERATOR LABORATORY

# MONITORING THE TEST


Aggregate Network I/O


Aggregate CPU Utilization


Aggregate Memory Usage


Aggregate Load

SLAC NATIONAL ACCELERATOR LABORATORY

# RESULTS



Aggregate Memory Usage

Aggregate Load

## Preliminary conclusions

- the architecture of the cluster is capable of delivering up to 12 GB/s of data to client applications
- moderate resource (CPU and memory) utilization under extreme I/O load leaves more than enough room for client job to perform any useful processing on the data
- the I/O capability of the cluster far exceeded any possible demand from known BABAR applications run (or to be run) on the Cluster
- hence, the cluster can be easily expanded with storage-less (pure) compute nodes should this be needed by the BABAR experiment
- actual limits of the expansion can be easily drawn by cross-correlating performance numbers of this test against I/O requirements of the applications. For instance, based on prior tests run on the cluster it's probably safe to estimate at least 5 times greater number of compute nodes in the Cluster as compared with its current configuration.

This page was last modified on 28 February 2012, at 14:08.

This page has been accessed 212 times.

Privacy policy   About Bbr_wiki   Disclaimers