



Service Challenge Meeting
"Service Challenge 2 update"

James Casey, IT-GD, CERN
IN2P3, Lyon, 15 March 2005





Overview



- Reminder of targets for the Service Challenge
- Plan and timescales
- What is in place at CERN Tier-0
- What has been tested at Tier-1 sites
- What is still to be tested at Tier-1 sites
- Issues outstanding from CERN Meeting
 - Communication
 - Monitoring
- Problems seen so far



Service Challenge Summary



- "Service Challenge 2" Started 14th May (yesterday)
- Set up Infrastructure to 6 Sites
 - NL, IN2P3, FNAL, FZK, INFN, RAL
- 100MB/s to each site
 - 500MB/s combined to all sites at same time
 - 500MB/s to a few sites individually
- Use RADIANT transfer software
 - Still with dummy data
 - 1GB file size
- Goal : by end March, sustained 500 MB/s at CERN



Plans and timescale - Proposals



- Monday 14th- Friday 18th
 - Tests site individually. 6 sites, so started early with NL !
 - Monday IN2P3
 - Tuesday FNAL
 - Wednesday FZK
 - Thursday INFN
 - Friday RAL
- Monday 21st – Friday 1st April
 - Test all sites together
- Monday 4th – Friday 8th April
 - 500MB/s single site tests
- Monday 11th - ...
 - Tape tests ?
- Overlaps with Easter and Storage Management workshop at CERN



What is in place – CERN (1/2)



- Openlab Cluster extended to 20 machines
 - oplapro[55-64,71-80].cern.ch
 - All IA64 with local fast disk
 - Allows for concurrent testing with
 - a 4 node setup (oplapro55-59) for tier-1 sites
 - 2 x 2nodes for gLite SC code (oplapro60-63)
 - *radiantservice.cern.ch* is now 9 nodes (oplapro71-79)
 - Radiant Control node (oplapro80)
- Raidant software improvements
 - Updated schema and state changes to match gLite software
 - Will make future changes easier
 - Cleaned up some of the logic
 - Channel management improved
 - E.g. Can change number of concurrent transfers dynamically
 - Better logging
 - Load-generator improved to be more pluggable by site admins



What is in place – CERN (1/2)



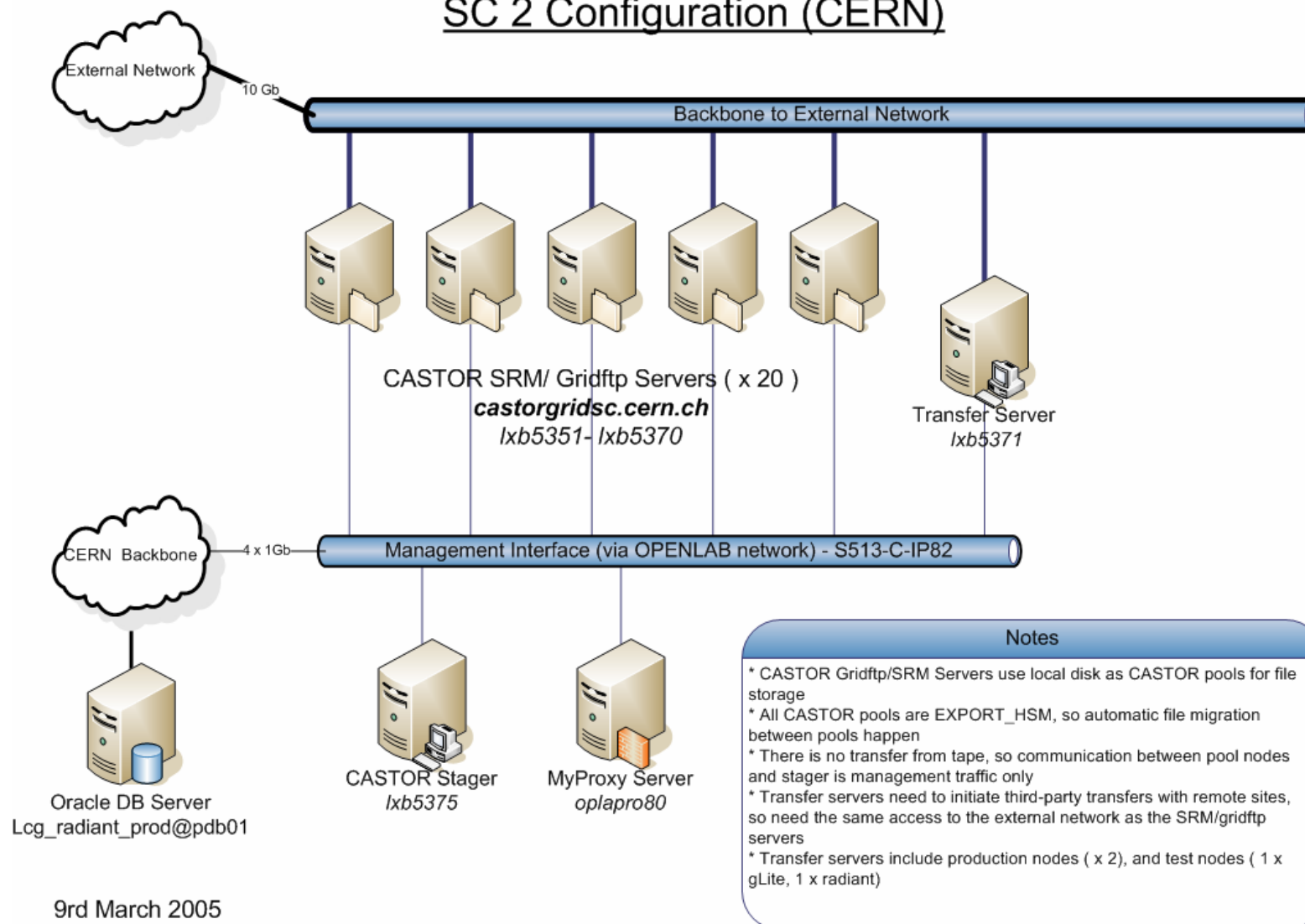
- New IA32 cluster for SC2
 - Managed as a production service by Castor Operations Team
 - Alarms and monitoring in place
 - 20 node lxf(s[51-70]) load-balanced gridftp/SRM
 - Uses local disks for storage
 - *castorgridsc.cern.ch*
 - *Still Waiting for connection to external 10Gb connection*
 - *Expected this week*
- 4 New transfer control nodes
 - 1 for SC2 production
 - 1 for SC2 spare
 - 1 for gLite production
 - 1 for testing (iperf, etc...)



SC CERN Configuration



SC 2 Configuration (CERN)



9rd March 2005



What is outstanding - CERN



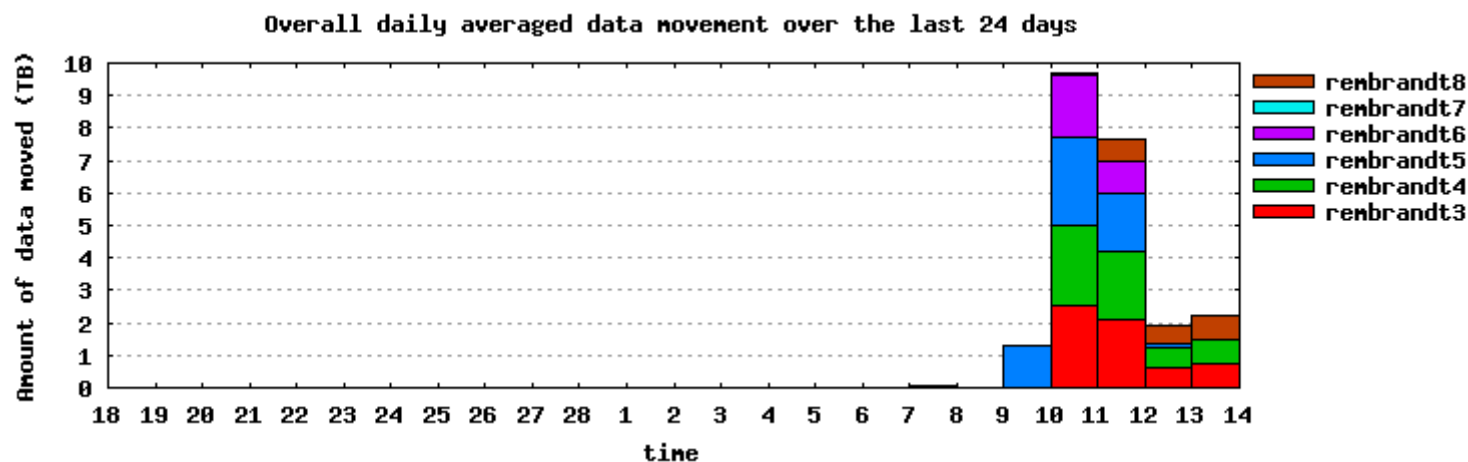
- External network connectivity for new SC2 cluster
 - We will use openlab machines until this is ready
- Radiant Software
 - Curently only pure gridftp has been tested
 - SRM code was written a long time ago, but not tested with a remote SRM
 - RAL will use this code
 - More agents needed
 - to set 'Canceling' jobs to 'Cancelled'
 - To debug error messages of jobs in 'Waiting' and send back to queue when ready
 - Some more CLI tools for admins
 - Summary of what transfers are in progress
 - Last hour Failure/Success rate per site



What is tested – SARA



- NL
 - Nodes in place for a few weeks
 - Testing with iperf and globus-url-copy done
 - Up to 80 MB/s peak single stream disk to disk
 - Sustained 40MB/s
 - Network currently over 3 x 1Gb etherchannel
 - 150MB/s sustained disk to disk





What is tested – IN2P3



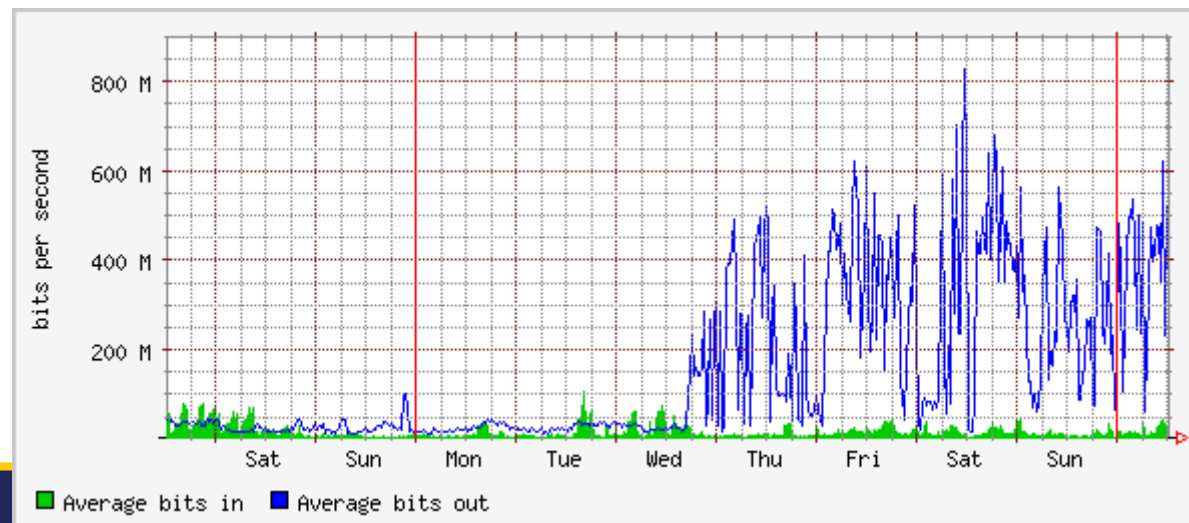
- IN2P3
 - 2 nodes that were tested in January
 - iperf memory-to-memory tests
 - ~300Mb/s to a single host
 - Single and multiple
 - Can get 600Mb/s to both hosts combined
 - Better gridftp single stream performance now seen
 - 60MB/s single stream disk to disk peak
 - 35-40MB/s sustained
 - Problem with one RAID controller stopped full tests to RAID disk
 - using local disk for now in `ccxfer04.in2p3.fr`
 - Fixed?



What is Tested - FNAL



- SRM to SRM from CASTOR SRM to dCache SRM
 - Also tested gridftp to SRM
 - Data pulled from Fermi
- Problems found and resolved in DNS handling
 - Java DNS name caching interacted poorly with the round robin DNS load-balancing (Fixed quickly by Timur)
 - Showed up with srmCp
 - All 40 transfers went to one host!





What is tested – FNAL/ BNL



- Transfers to tape have been tested
 - CERN SRM to FNAL tape at 40-50MB/s via FNAL SRM
- Transfers will be controlled by PHEDEX
 - Will be put in place early next week
- Will send data to 5 USCMS tier-2 sites
 - Florida, Caltech, San Diego, Wisconsin and Purdue
- BNL
 - Tried SRM to SRM transfers
 - Had problems with SURL format causing failures
 - Dual attached node at BNL caused debugging difficulties
 - Routing asymmetries



What is outstanding – Tier-1 sites



- RAL
 - dCache in place – good performance seen (500MB/s ?)
 - Waiting on UKLIGHT network line final provisioning
 - As of yesterday, Netherlight had provisioned the fibre to CERN
- INFN
 - Cluster now in place (8 nodes)
 - Network tests done
 - Can fill the 1Gb pipe
- FZK
 - Cluster in place – will use gridftp to start
 - Network tests with GEANT done
 - Still some problems to fill whole provisioned 10Gb pipe
 - Will test after SC2



Communication



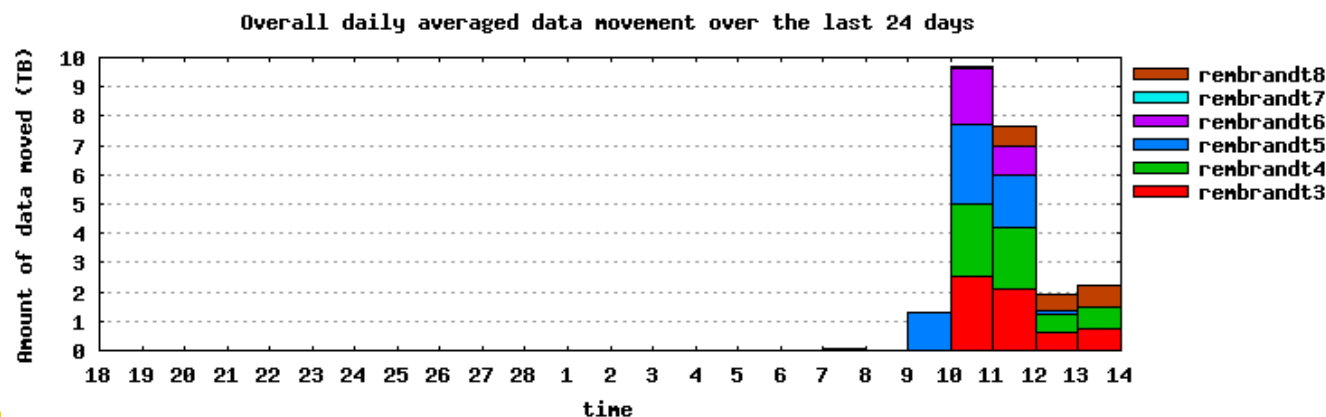
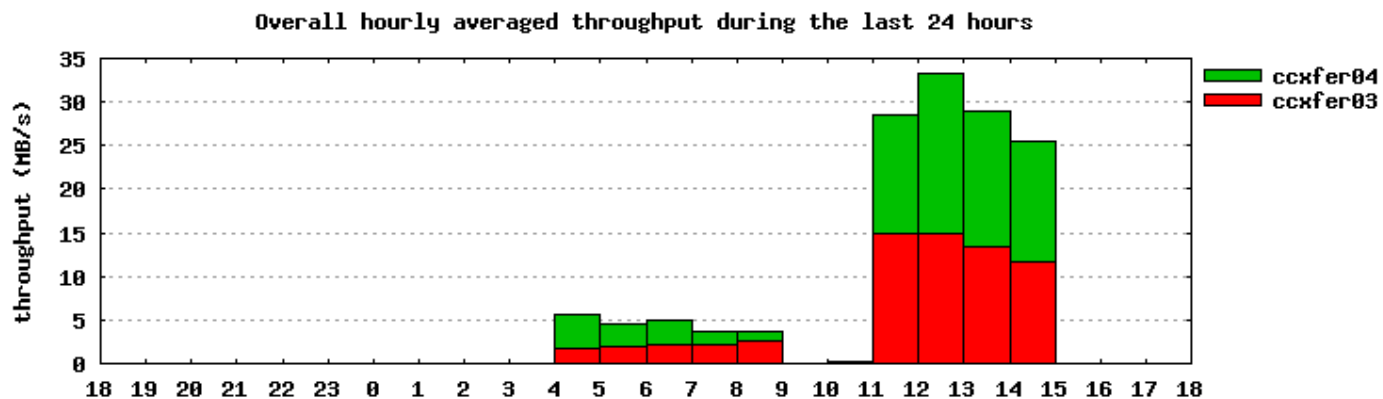
- Using **service-challenge-tech** list to distribute information
- Each site has now specified a mail alias for the local service challenge team
 - **Still waiting on one from FZK**
- Service Challenge meetings
 - **What Schedule ?**
 - Weekly (daily ?) with all sites?
 - What day?
- **Still work to be done on documentation**



Monitoring



- Ron Trompert wrote gridftp analysis tool for SARA
 - Extended it to analyze CERN logfiles to all sites
 - <http://challenge.matrix.sara.nl/SC/>





Monitoring



- For SC1 network statistics gathered from routers at CERN
 - Provided by Danny Davids from CERN IT-CS Group
- For SC2 extended to other sites
 - Uses SNMP from the router
 - Can also be from a switch (any virtual interface)
 - Gathered by hosts at CERN
 - Bytes in/out, packets in/out gathered
 - Stored in central database
- **Visualization tool outstanding**
 - Have identified effort from Indian LCG Collaborators
 - Will be coming to CERN this month



Problems seen so far



- Long timeouts/transfers (~ 2500 secs) in gridftp transfers
 - Intermittent failure (105 times to FNAL in total)
 - Seen by FNAL and NL
 - Now have some log messages from FNAL to aid the debugging
- Quite sensitive to changes in conditions
 - Extra traffic on disk servers seems to alter performance a lot
 - Changing number of concurrent transfers also
- Performance is not constant with radiant control software
 - Some parts needs to be understood better (next slide)



Radiant load-balancing



- Some performance loss due to way Radiant is scheduling transfers
 - When one transfer finishes, next scheduled transfer is not to that host
 - Gives asymmetric loading of destination hosts
 - A property of “load-balanced” gridftp servers
 - Not really load-balanced with round-robin – not awareness of behaviour of other nodes
 - Can end up with all transfers going to one host, with all others idle
 - One host ends up heavily loaded and very slow
 - Loss of all throughput



Solutions



- The longer term solutions are
 - move to SRM which does the load-balanacing (and with srm-cp this might go away)
 - move to gridftp 3.95 which could be adapted to have a load-balanced front-end
 - modify radiant (or the gLite software) to do better scheduling
 - Network and destination-SURL aware ?
 - put in a load-balanced alias for the destination cluster which does N of M load-balancing (i.e. M hosts where the N most-lightly loaded are in the alias.)
 - Used currently for castorgrid, and will be used for castorgridsc
- Move to SRM for all sites is coming, so putting in an appropriate amount of effort for this is required
 - How much is appropriate?



Summary



- Good progress since last meeting
 - Especially the gridftp monitoring tool from Ron
- Sites are all nearly ready
- Still a lot of work to be done
 - But now seems to be debugging underlying software components rather than configuration, network and hardware issues