

# Table of Contents

Introduction.....	2
Challenge environment.....	3
Infrastructure at IN2P3-CC.....	3
Network.....	3
Hardware/software.....	3
Configuration.....	4
Transfers details.....	4
Overview results.....	5
Disk-disk transfers.....	5
Disk-tape transfers.....	5
Disk-disk transfers detailed results.....	6
Daily average throughput.....	6
Quality of transfers.....	6
Disk-disk transfers comments.....	7
Legend.....	7
Analysis.....	7
Disk-tape transfers detailed results.....	8
Daily average throughput.....	8
Quality of GridFTP transfers.....	8
Quality of migration transfers.....	8
Disk-tape transfers comments.....	9
Legend.....	9
Analysis.....	9
Disk-disk transfers BIS.....	11
Results.....	11
Conclusion.....	12
Conclusions.....	13
Annexes.....	15
A1. Head node stats (disk-disk).....	15
A2. ccxfer13 stats (disk-disk).....	16
A3. ccxfer15 stats (disk-disk).....	17
A4. ccxfer16 stats (disk-disk).....	18
A5. Load on head node during DB backup.....	19
A6. FTS feeding issues.....	20
A7. Head node stats (disk-tape): Wed 19/04- Wed 26/04.....	21
A8. ccxfer09 stats (disk-tape): Wed 19/04- Wed 26/04.....	23
A9. ccxfer10 stats (disk-tape): Wed 19/04- Wed 26/04.....	24

# Introduction

This document describes the LCG Service Challenge 4 throughput phase at IN2P3-CC

**Author:** Lionel Schwarz

**Date:** 19. Apr. 2006

## Goals:

- Sustain 200MB/s from T0 disk to IN2P3-CC disk for 2 weeks with understanding of rate drops and ability to run above the nominal rate.
- Sustain 75MB/s from T0 disk to IN2P3-CC tapes for a week with understanding of rate drops and ability to run above the nominal rate.

## Planning:

- Mar 30-Apr 2: Preparation
- Apr 3-Apr13: Disk-disk transfers
- Apr 14-Apr 17: Unattended disk-disk transfers
- Apr 18-Apr 24: Disk-tape transfers

# Challenge environment

## *Infrastructure at IN2P3-CC*

### Network

- Direct link to CERN at 10gbs (dedicated)
- Switch to Renater/Geant 1gbs link in case of problems (shared)

### Hardware/software

- Dcache/SRM Head node
  - Pentium III 1.4 GHz, 1GB RAM
  - SL3.0.3 (kernel 2.4.21-20.ELsmp)
  - Globus 2.4, Java 1.5.0\_01, PostgreSQL 8.1.0 server for SRM-dCache, dcache-server-1.6.6-5 with SRM door and GSIDCAP door, pnfs-postgresql-3.1.10-3
  - All data (SRM database, pnfs, dCache log files, dCache billing) are stored on a 70GB ext3 mirrored FS
- dCache disk servers (pool nodes)
  - Disk tests: *ccxfer13+ccxfer15+ccxfer16* (the theoretical maximum aggregated rate is 300MB/s)
  - tape tests: *ccxfer09+ccxfer10* (the theoretical maximum aggregated rate is 200MB/s, but HPSS resources were allocated for 100MB/s maximum)
  - *ccxfer13, ccxfer09, ccxfer10*: V40Z Quaq-pro AMD opteron 2GHz, 8GB RAM, 4 TB in RAID5 attached to IBM DS8300 unit formatted in XFS
  - *ccxfer15, ccxfer16*: Transtec Bi-pro Xeon 3GHz, 4GB RAM, 4 TB locally attached formatted in XFS
  - SL3.0.5 (kernel 2.4.21-32.0.1.EL.XFSsmp)
  - Globus 2.4, Java 1.5.0\_01, dcache-server-1.6.6-5 with GridFTP door, RFIO client (rfcp) for the HPSS connexion
  - 2\*1Gb interface (1 for GridFTP traffic, 1 for rfcp traffic)
- 10 GridFTP servers (run on all pool nodes)
- HPSS
  - buffer disk: 1TB (dedicated) with immediate migration

- Up to 6 9940b drives available (shared with other activity)

## Configuration

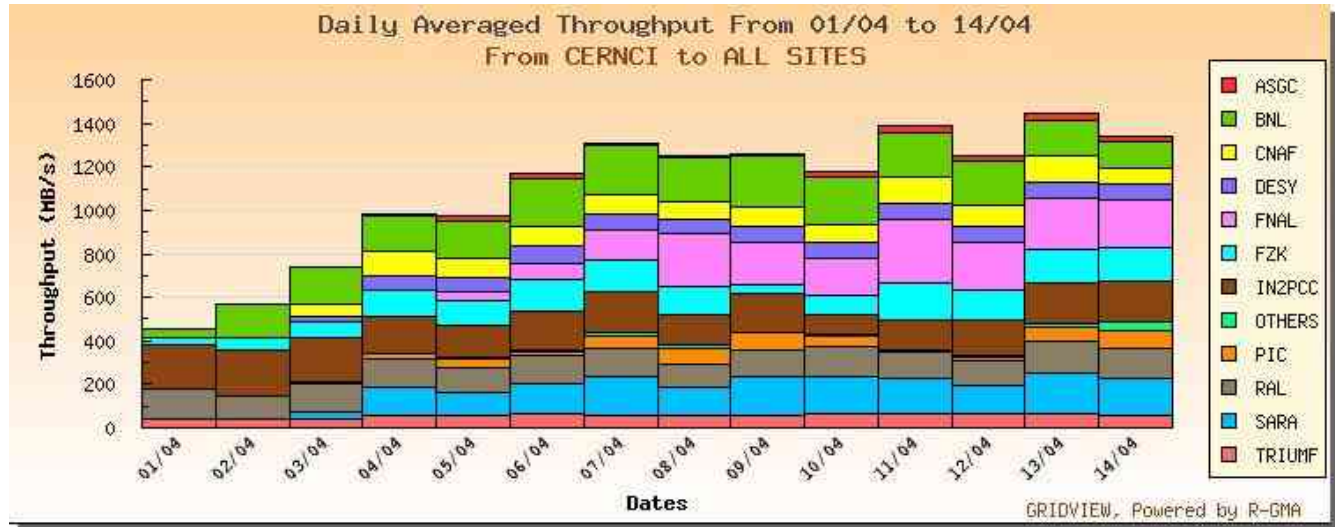
- Linux kernel parameters on disk servers
  - /proc/sys/vm/bdflush: 10 500 0 0 500 3000 0 20 0
  - /proc/sys/vm/min-readahead: 2048
  - /proc/sys/vm/max-readahead: 2048
- Buffer size
  - GridFTP: 1MB
  - RFIO: 1MB

## *Transfers details*

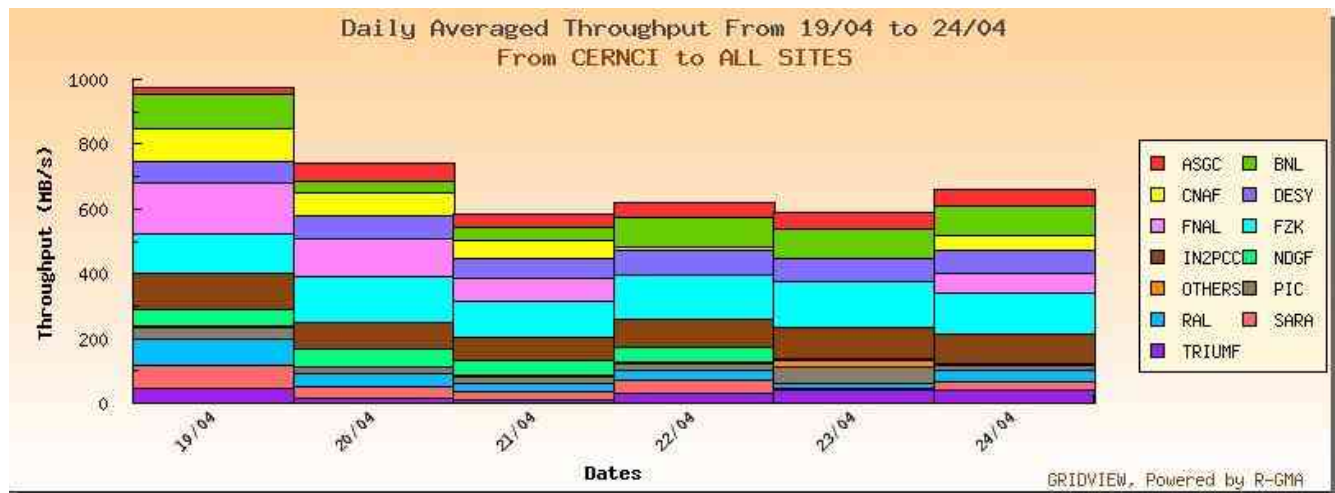
- FTS Channel CERN-IN2P3 managed by CERN and IN2P3
- Transfers initiated by CERN with “dteam” proxy
- FTS in “3rd party” mode
- DISK endpoint: “srm://ccsrm.in2p3.fr:8443//pnfs/in2p3.fr/data/dteam/disk” (Automatic deletion (every hour) of files older than 3 hours)
- TAPE endpoint: “srm://ccsrm.in2p3.fr:8443//pnfs/in2p3.fr/data/dteam/hpss/tape-sc4” (Manual deletion to re-use cartridges)

# Overview results

## Disk-disk transfers



## Disk-tape transfers



## Disk-disk transfers detailed results

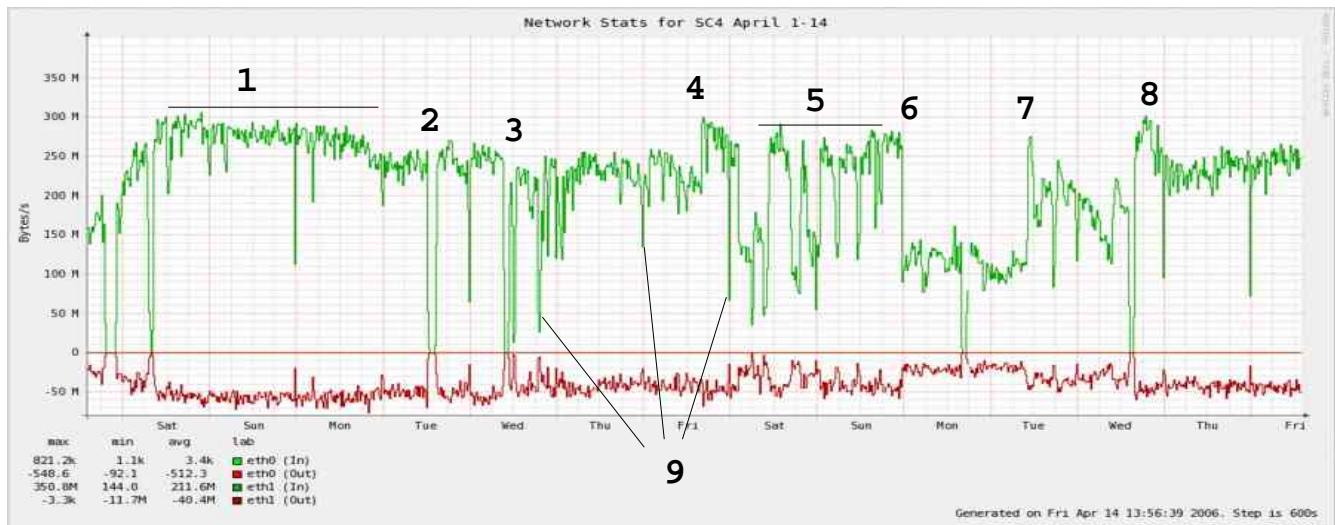
### Daily average throughput



### Quality of transfers

371 errors / 201273 transfers => error rate = **0.18%**

## Disk-disk transfers comments



### Legend

- eth1 out: outgoing traffic from dCache pools (Side effect of FTS 3rd party mode)
- eth1 in: incoming traffic to dCache pools (GridFTP from T0)

### Analysis

1. 01-03/04: sustained transfers at > 200MB/s
2. 04/04: not understood
3. 05/04: load generator problem @ CERN, RFIO problem @ CERN
4. 07/04: number of files increased to 30: rate > 200MB/s
5. 08-09/04: very irregular rate: FTS@CERN not able to feed sites on a regular basis (see annexe A6)
6. 10/04: From midnight, rate dropped at 100MB/s, no errors, no network problems. Not understood although there were too few concurrent transfers (see Annexe A6)
7. 11/04: number of files increased to 45: rate doubled (200MB/s) but then dropped slowly back to 100M/s: it seems that sites have strong concurrence on the FTS server
8. 12/04: after FTS DB intervention @ CERN, rate > 200MB/s
9. Each day at 23:40: rate drops due to high load on head node during PostgreSQL DB backup

# Disk-tape transfers detailed results

## Daily average throughput



Note: tape transfers started 19/04 at 15:00 GMT

## Quality of GridFTP transfers

180 errors / 40393 transfers => error rate = **0.44%**

## Quality of migration transfers

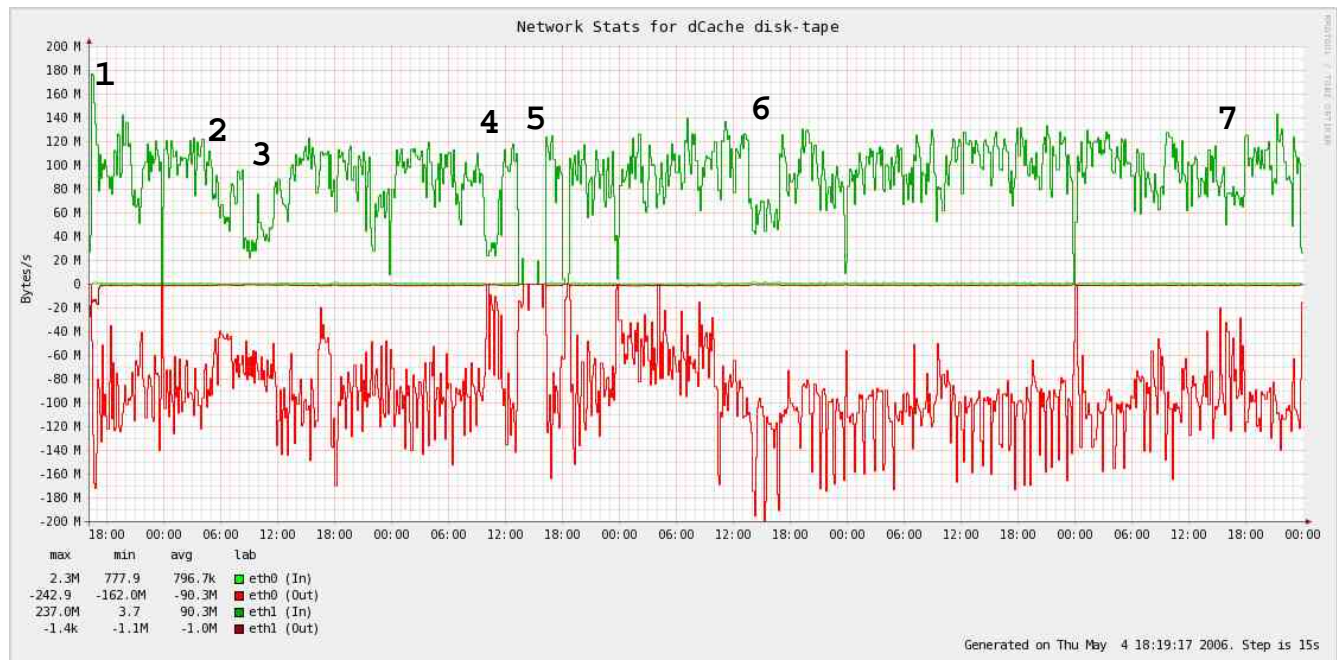
306 errors / 40123 transfers => error rate = **0.76%**

Note: when there is a problem in HPSS (for example disk full on 19/04) there is 1 error/file migrated although the error is the same. If we remove 19/04 from the statistics:

8 errors / 37633 transfers => error rate = **0.02%**



## Disk-tape transfers comments



### Legend

- eth1 in: incoming traffic to dCache pools (GridFTP from T0)
- eth0 out: outgoing traffic to HPSS

### Analysis

- 19/04: Start of disk-tape transfers
  - 10 files in the channel
  - 5 pools on 2 servers
  - 1 concurrent transfers on each pool
  - 1 HPSS stream on each pool
- 20/04: Due to a bad transfer, a mover hanged last night at 19:00 GMT and when the other pools on the machine ("ccxfer09") got full (at 4:00 GMT), no transfers did happen anymore on this machine because of dCache load balancing. The rate fell because all transfers went a single server ("ccxfer10"). When the stuck mover was killed transfers could resume on ccxfer09
3. The number of files received on the channel started to be very low, even 0 sometimes. The situation got better at 8:30 GMT without any intervention at IN2P3-CC. All channels affected.
- 21/04: CASTOR errors due to an unavailable file system and at the same time tape recall being stuck after an intervention on the IBM robotics two days ago
5. FTS out of service

6. 22/04: The night before, the migration process was stuck on both pools on one of the disk servers ("ccxfer10"), so the buffer was half full in the morning. To allow the disk to be emptied, I have disabled access to "ccxfer10" for incoming transfers. For 3 hours, all transfers went to one single host ("ccxfer09") that is why the rate dropped. However, why the migration process stopped is not understood.
7. LHCb transfers are using the same pools as dteam and their files are smaller (700MB), that could explain the drop. Increased the number of files from 10 to 20 and increased the number of movers from 5 to 10.
8. 25/04: end of disk-tape transfers

## Disk-disk transfers BIS

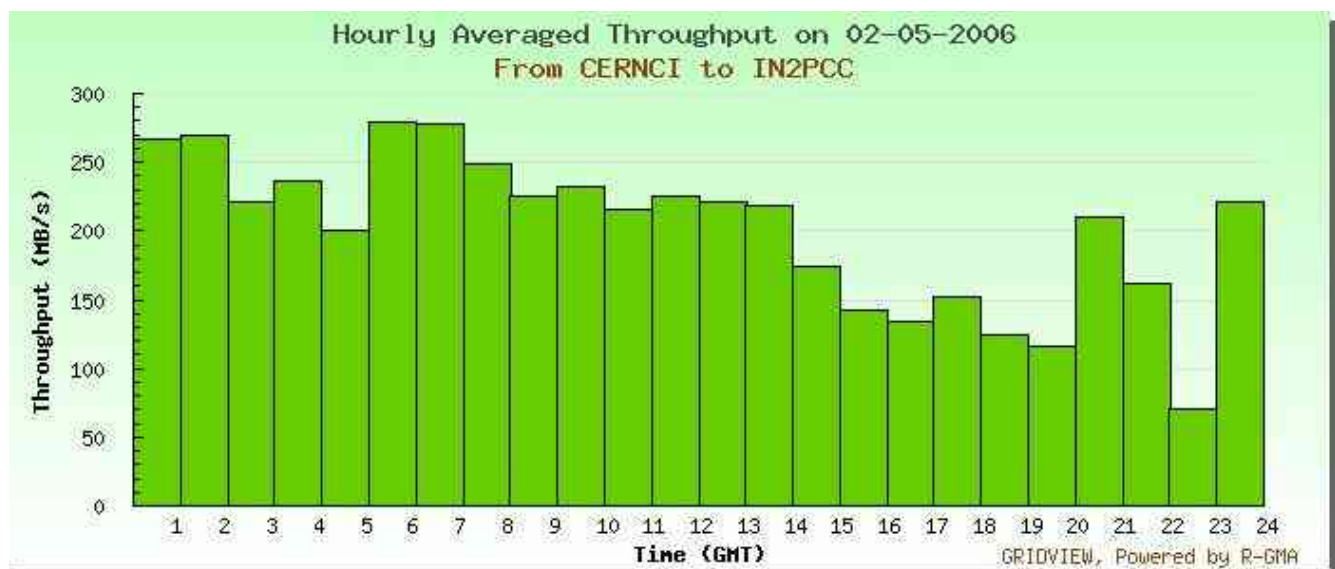
After the successful disk-tape transfers, it was decided to switch back to disk-disk transfers to understand why the nominal rate (200MB/s) could not be sustained during the first period.

We used the same configuration as first period with 1 more disk server.

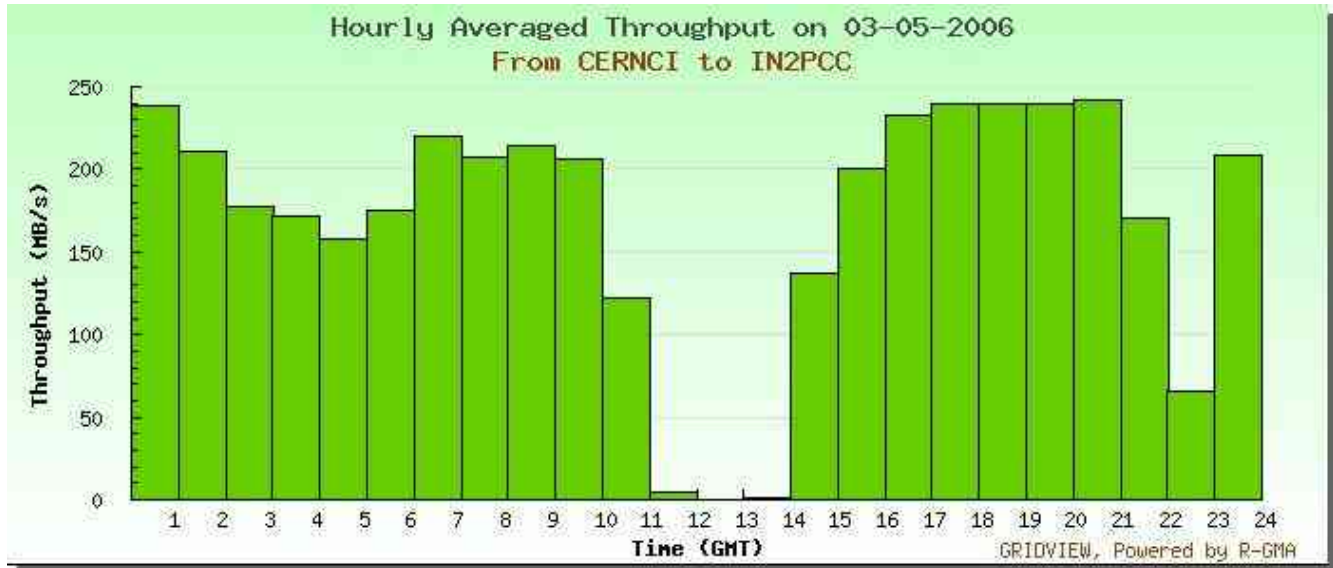
### Results



It was proved that IN2P3 could sustain nominal rate and above (+25%) for a few days using the same infrastructure as first period. Adding the 4<sup>th</sup> server allowed to run 25% more than the nominal rate. It was switched off on early 02/05 and the nominal rate was sustained that day despite a CASTOR problem in the afternoon (see SC4 blog).



On 3<sup>rd</sup> of May, the average was below the nominal rate due to SRM get failures (due to modifications in the CASTOR SC4 cluster), transfer failures due to a stuck tape recall at CERN and single stuck transfer blocking four channels in the FTS DB (see SC4 blog).



## Conclusion

During the second period, the overall rate from CERN was lower than during the first period (most sites were still doing tape transfers) and could explain why FTS could send average number of file close to the limit (40). This was not the case during the first period, where there was a strong concurrence between sites (FNAL and BNL had a very high number of files).

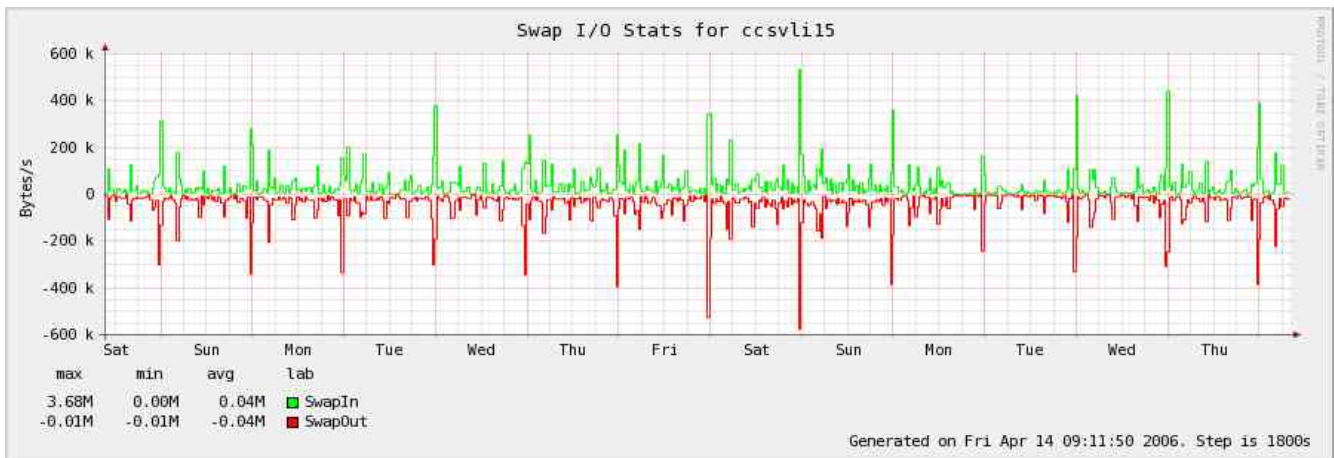
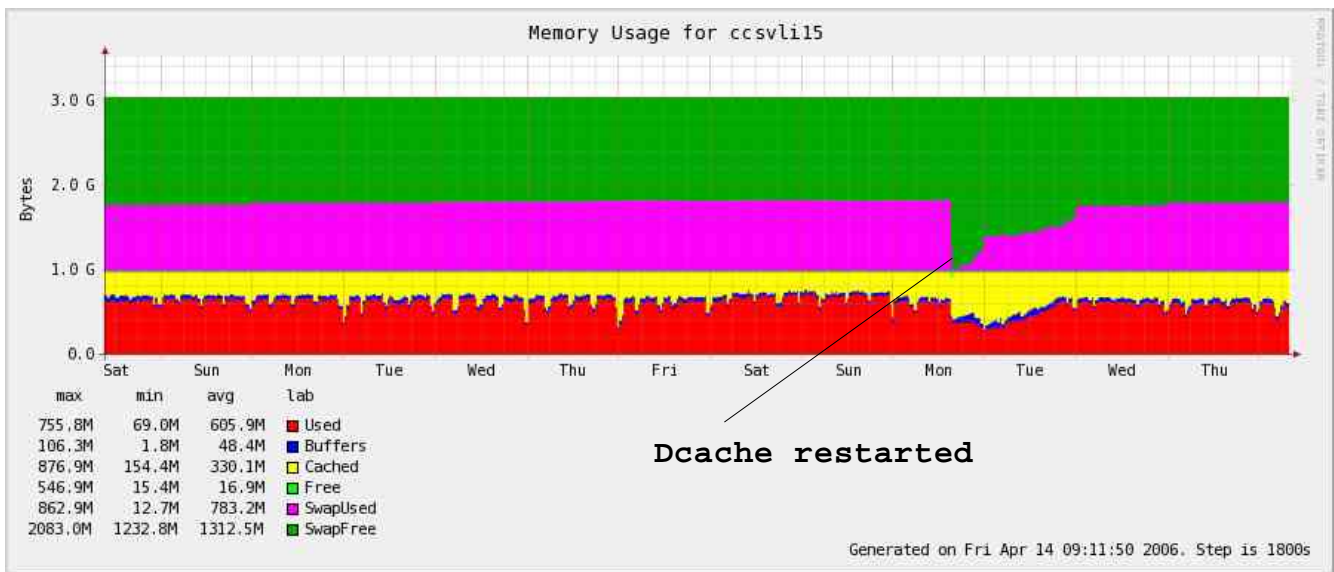
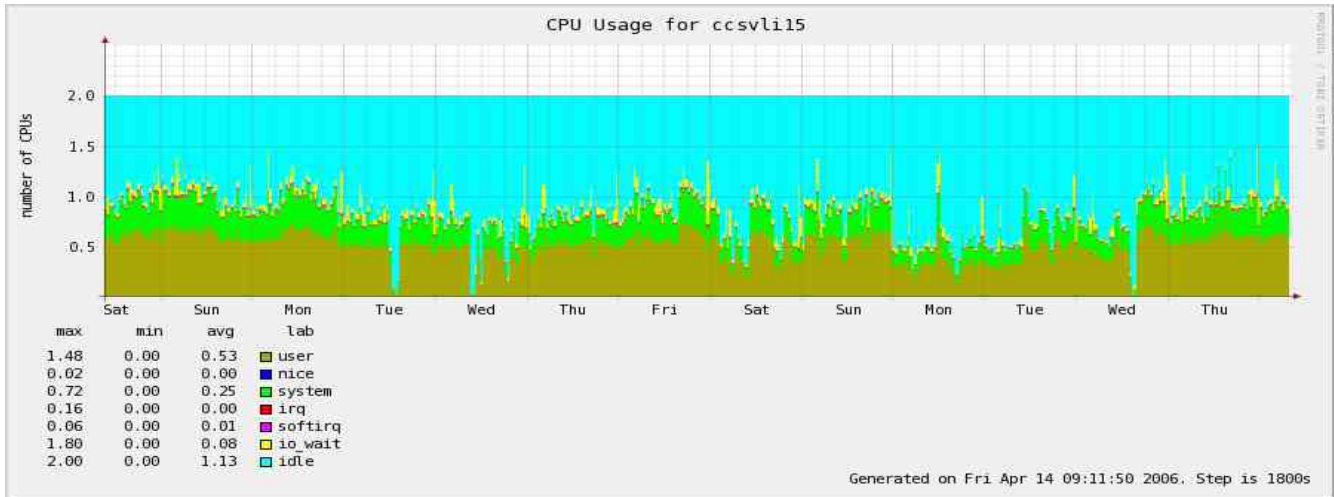
## Conclusions

<i>General remarks</i>	<i>Status</i>
All services (network, disk, dCache, HPSS) were stable except one disk server crashed (ccxfer13)	Problem understood already happened on this type of software. Setup allowed to ensure stability of service and transfer throughput with 1 server less.
Disk throughput target meet at the beginning but not sustained over the first period	It seems there was a strong concurrence between sites on the FTS server. The number of files in the channel had to be tuned in dependence with the activity of the other sites (eg number of files for FNAL, BNL)
The head node memory resources were not sufficient particularly during the daily backup	Memory was upgraded to 2GB on 25/04. <b>ACTION:</b> Split core services on 2 machines (in May).
The pool nodes resources were not fully used	
Disk throughput target sustained over the second period (26/04-03/05) at nominal rate+25%	Achieved with 4 disk servers and 40 files in the channel. Switching off 1 server did not make the rate decrease.
<i>FTS Issues</i>	<i>Status</i>
Irregular FTS feeding lead to an average number of transfers half the maximum allowed	
Extra inter-nodes traffic due to “3rd party” FTS mode	<b>ACTION:</b> SRM copy tests need to be scheduled (May?)
<i>dCache Issues</i>	<i>Status</i>
Stuck movers after failed transfers never cleaned	Submitted to dCache support (#596). <b>ACTION:</b> Timeout set to 600 seconds
GridFTP cell not removed when transfer failed	Submitted to dCache support (#600) <b>ACTION:</b> wait for next dCache release (1.6.7)
HSM Migration stopped	Submitted to dCache support (#602)
Deletion of precious files	Submitted to dCache support (#603) <b>ACTION:</b> Install patch provided by dCache (May)
<i>HSM Issues</i>	<i>Status</i>

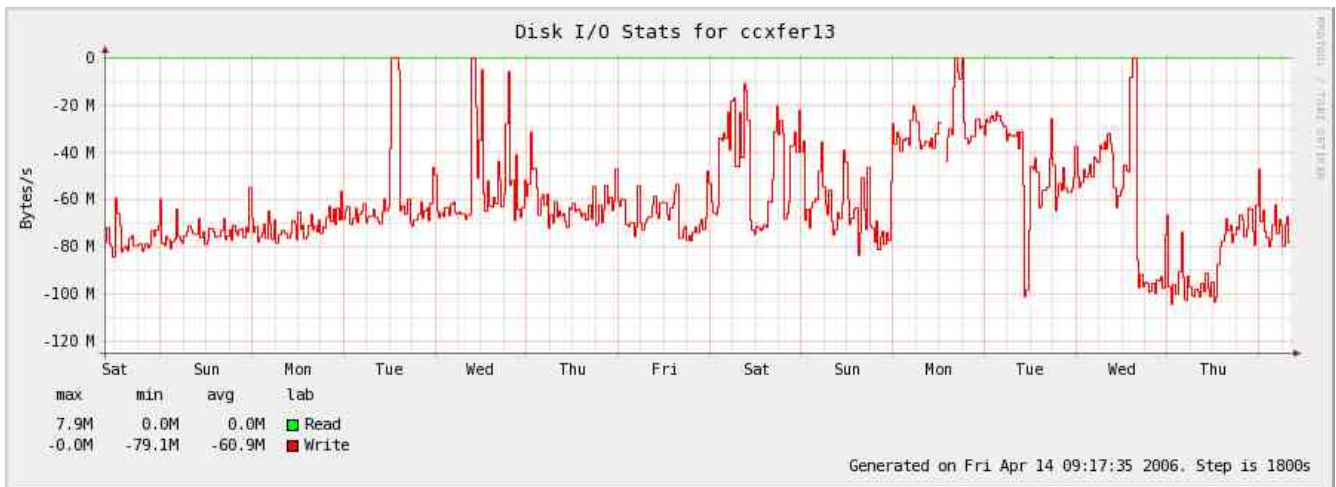
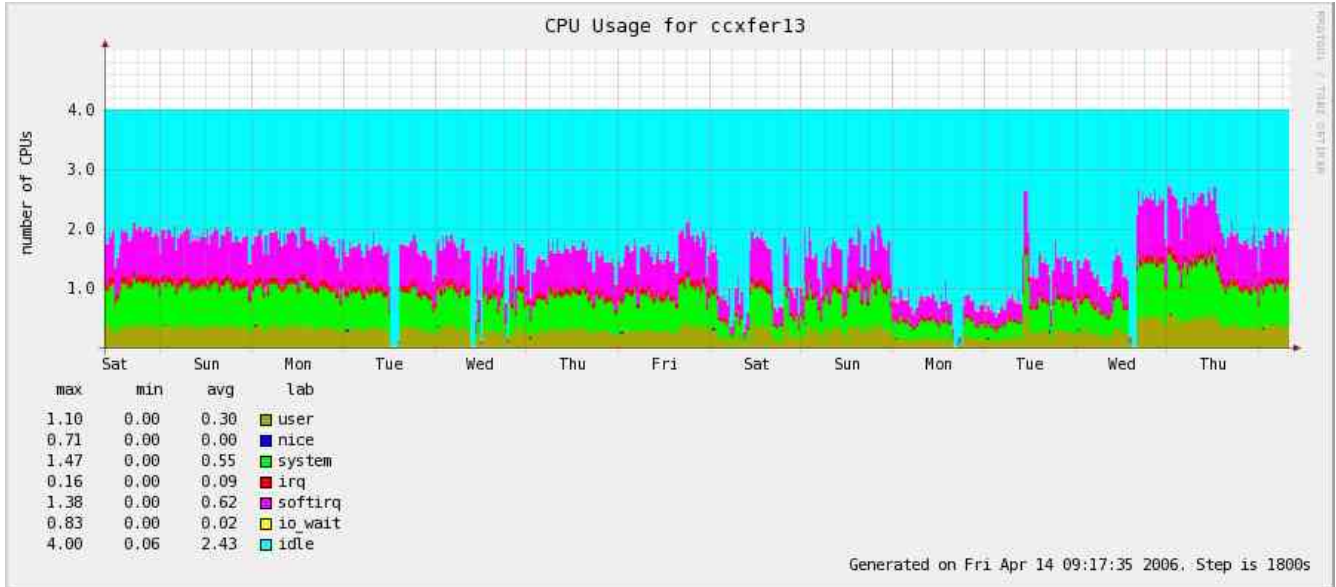
HSM Migration stopped	Submitted to dCache support (#602)
When HPSS disk is full, dCache should temporarily stop flushing precious files	Need to improve errors handling between HPSS and dCache

# Annexes

## A1. Head node stats (disk-disk)

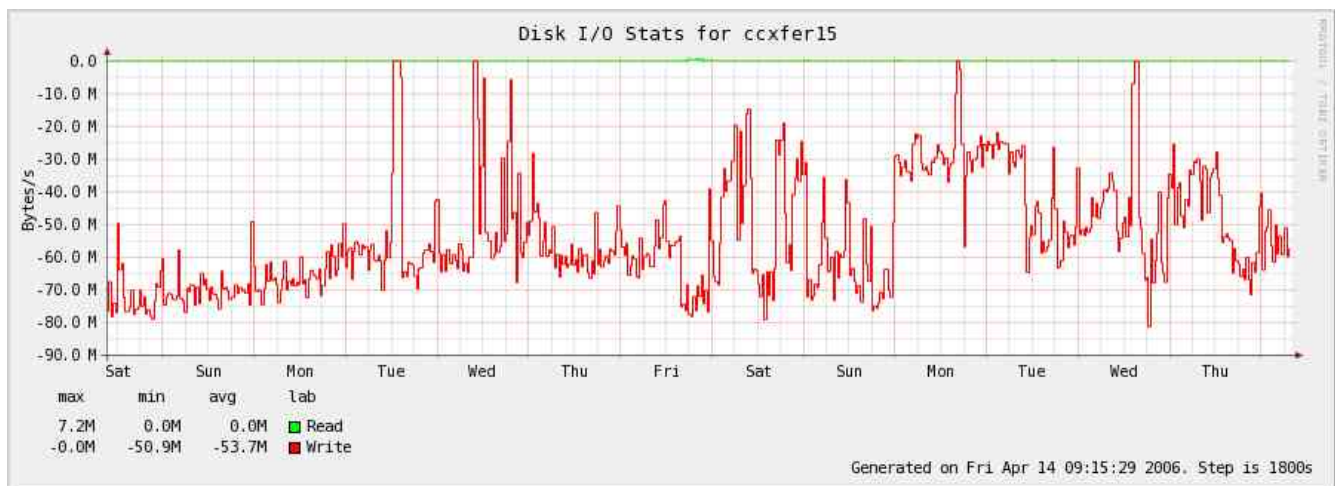
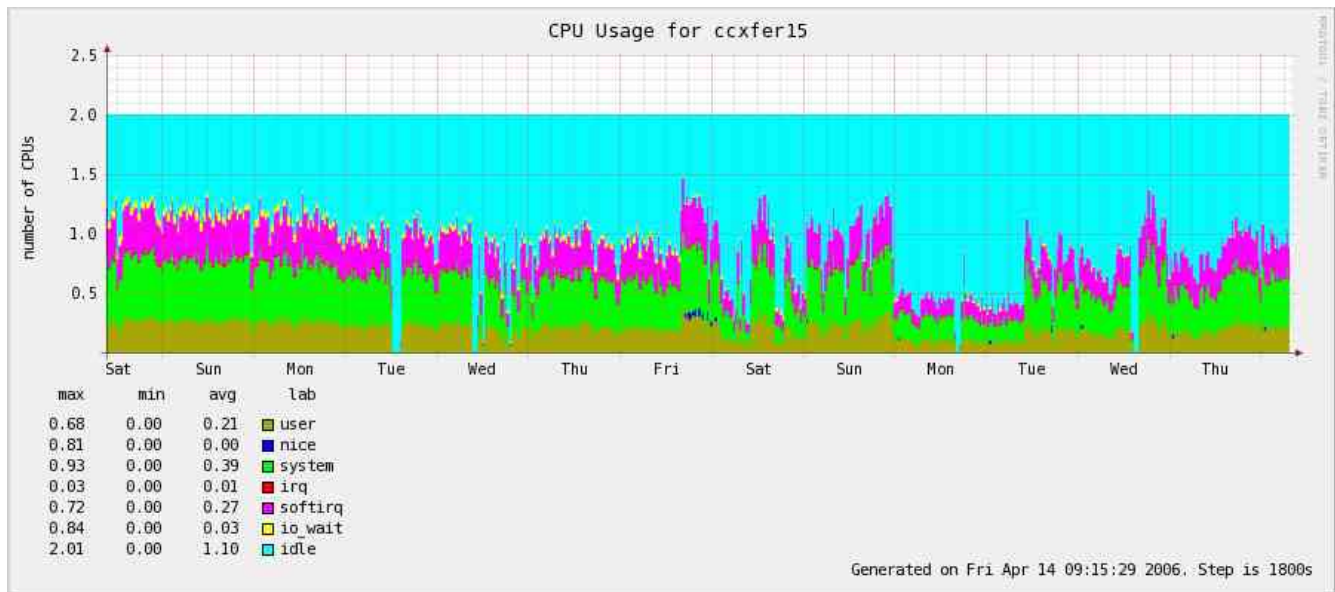


## A2. ccxfer13 stats (disk-disk)

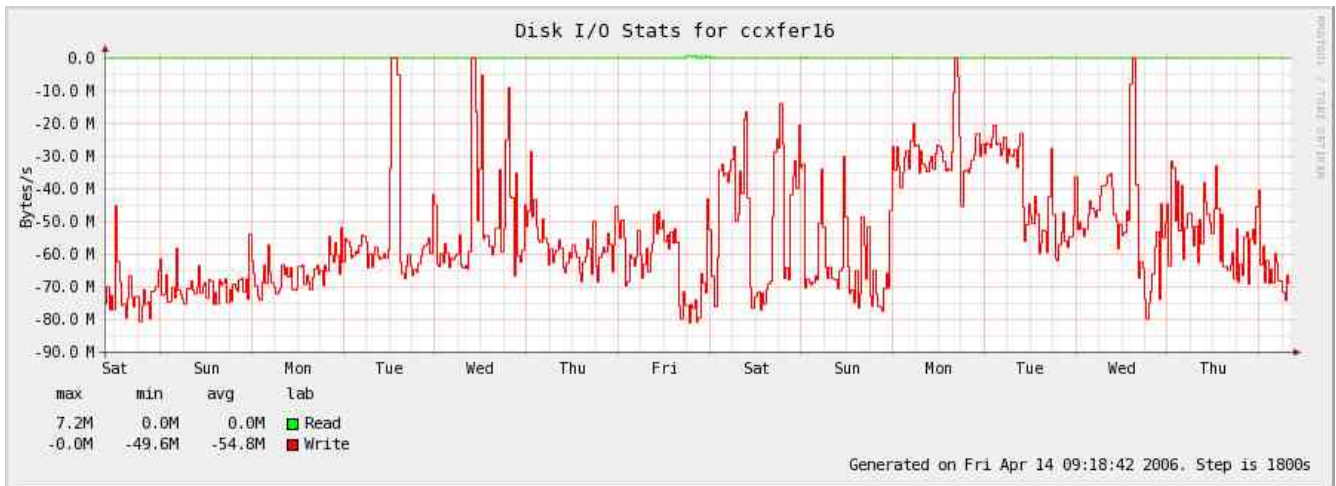
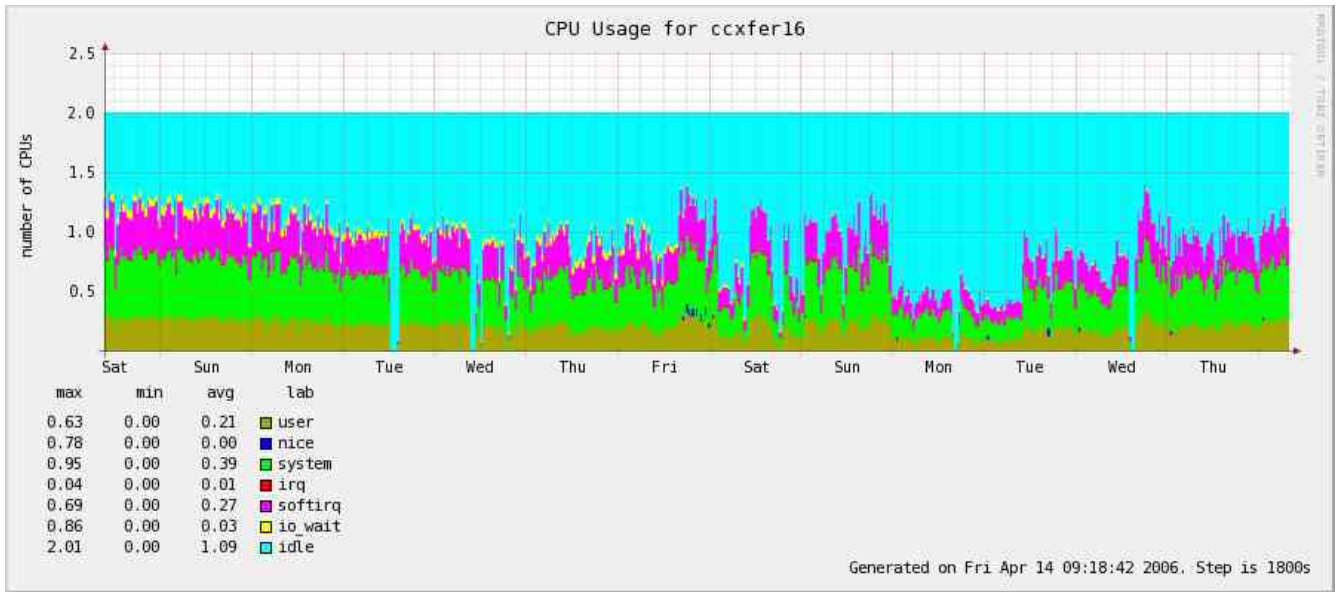




### A3. ccxfer15 stats (disk-disk)

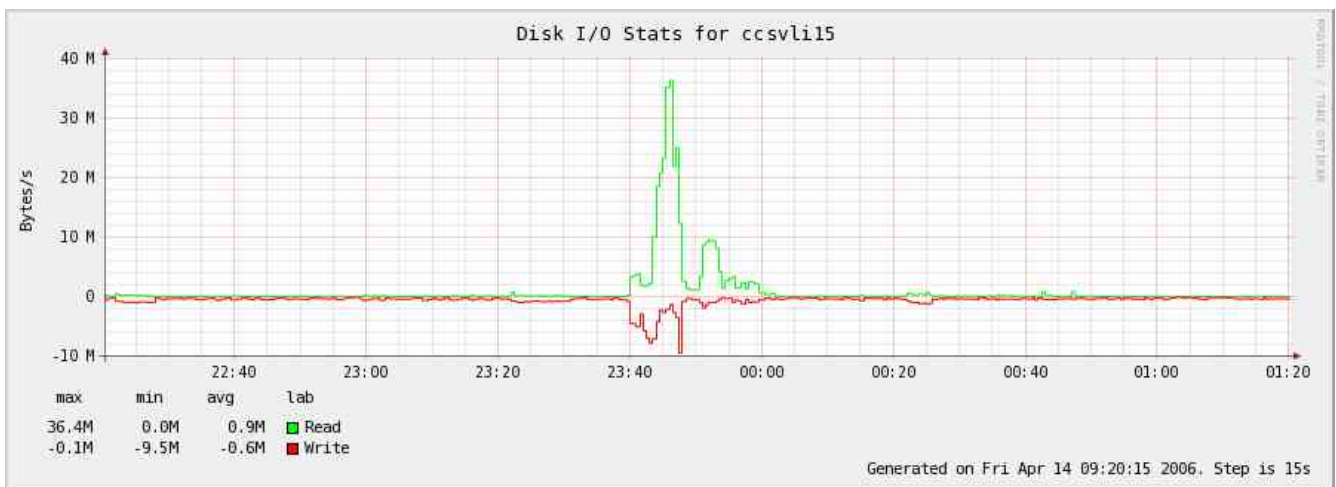
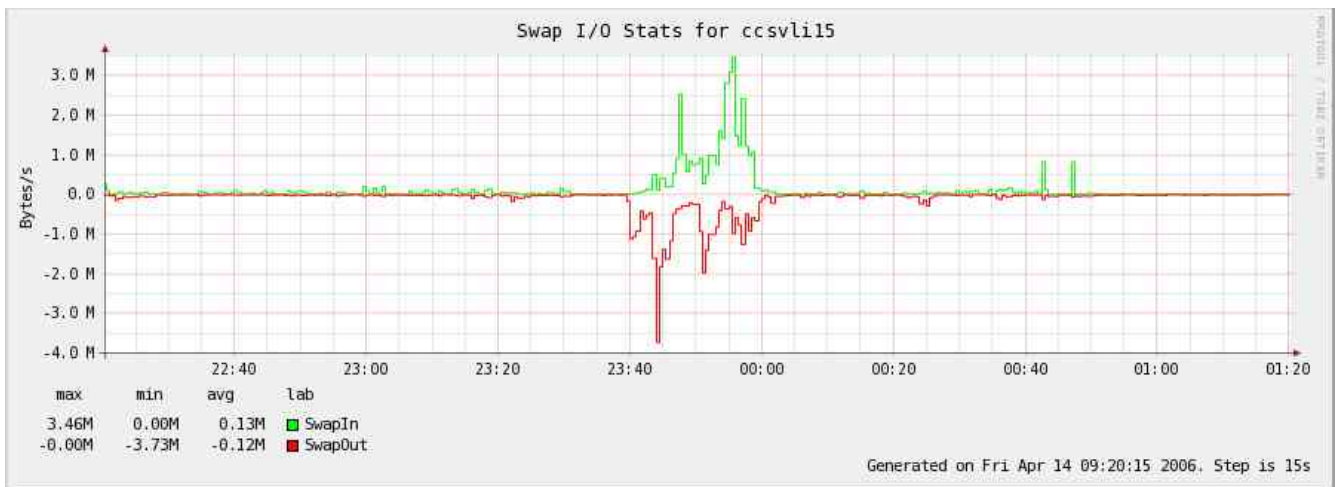
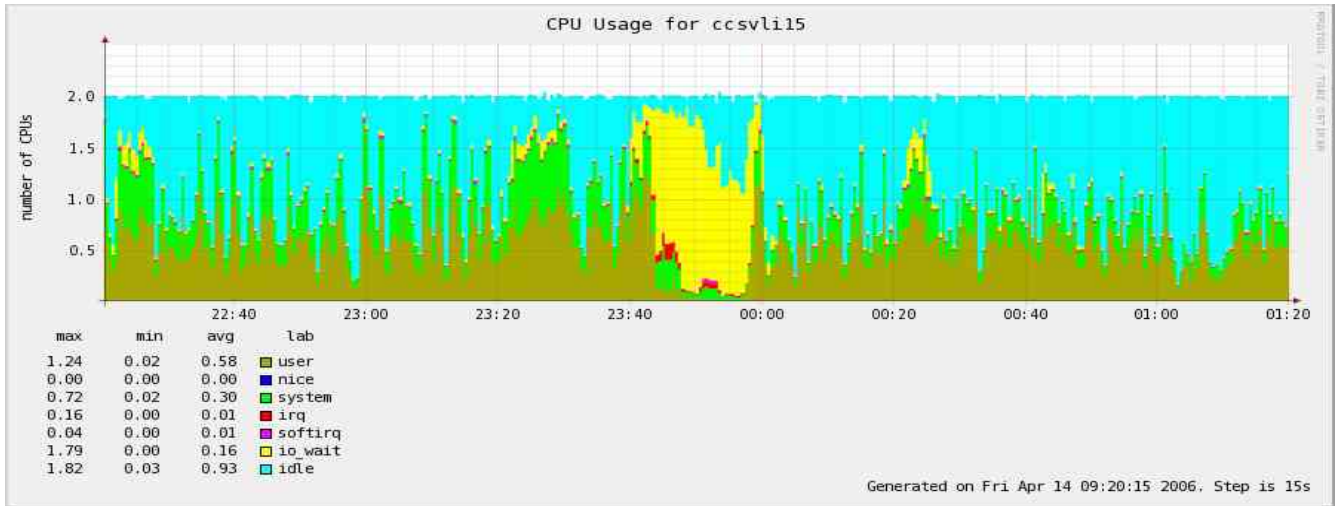


## A4. ccxfer16 stats (disk-disk)

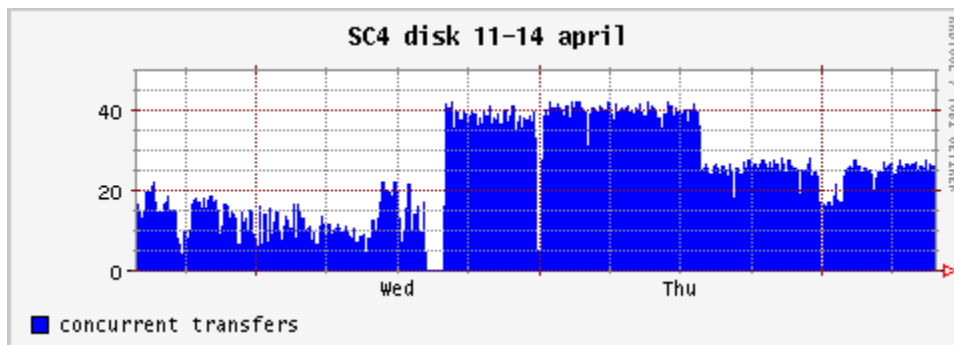
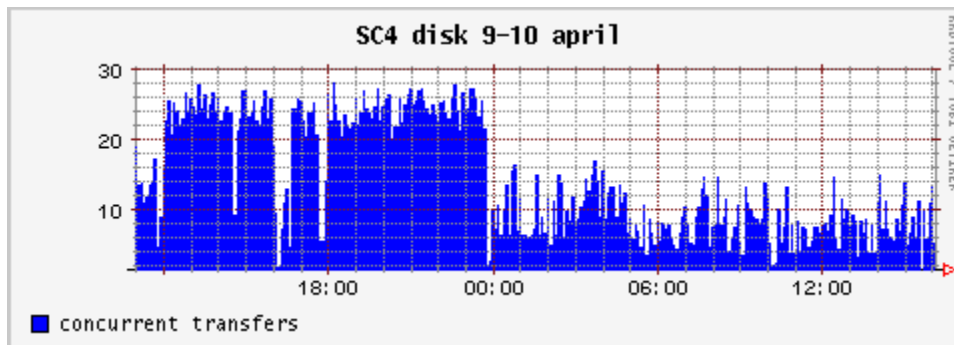


## A5. Load on head node during DB backup

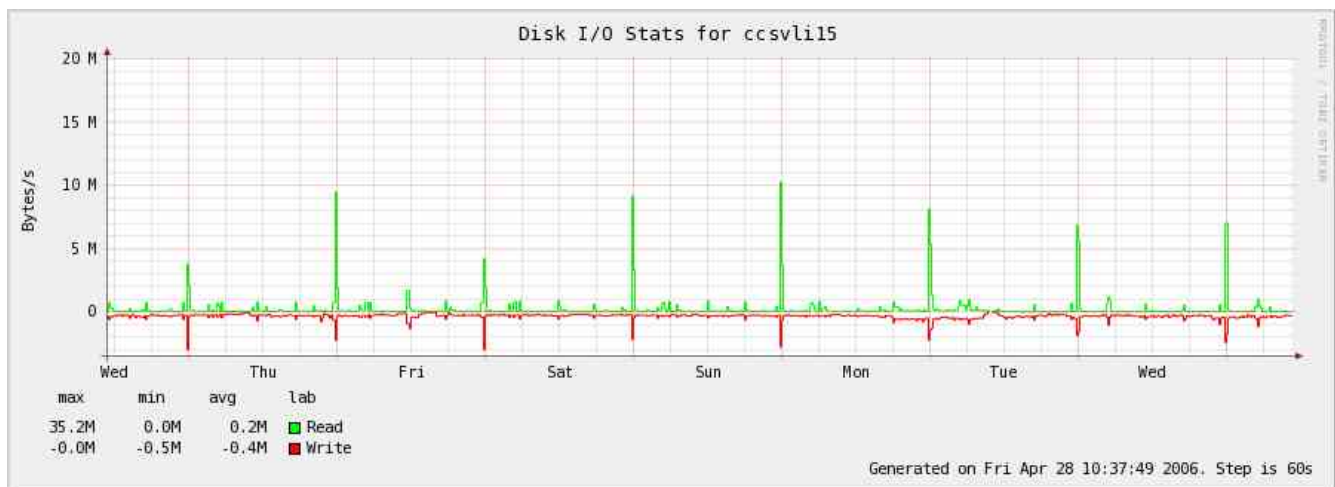
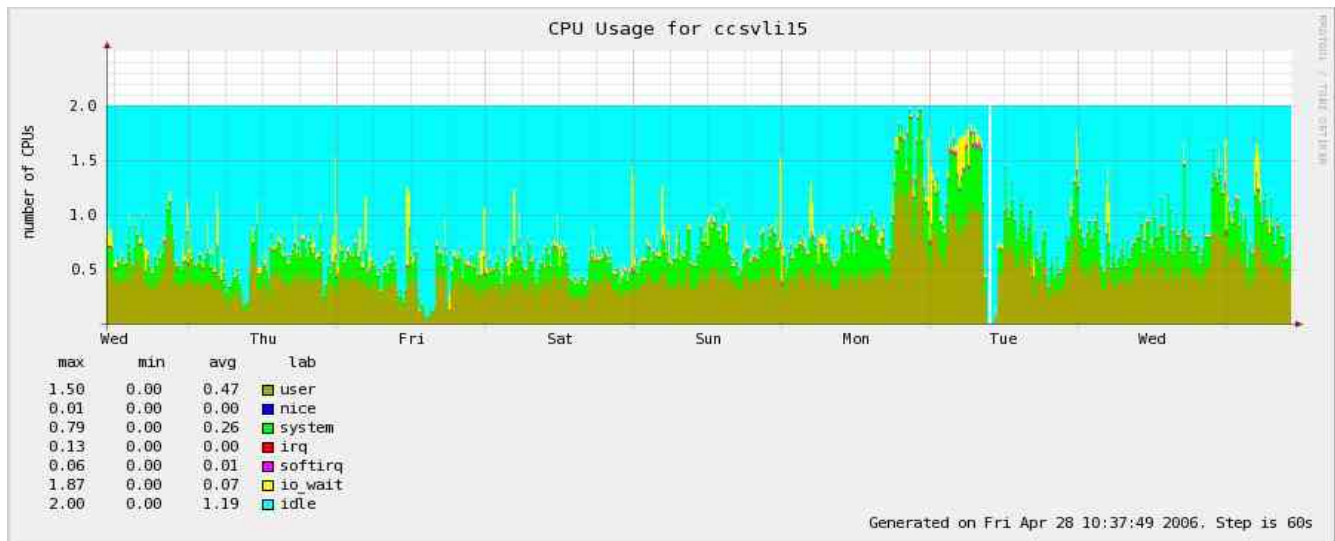
The DB backup is done at 11:40 PM every night. All DB are dumped and transferred to TSM.



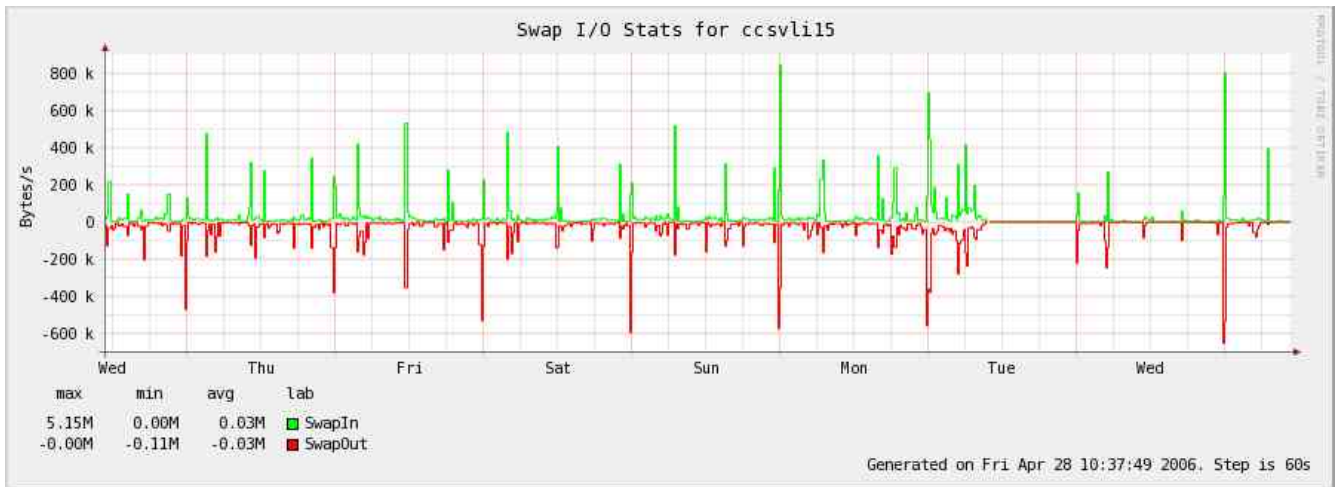
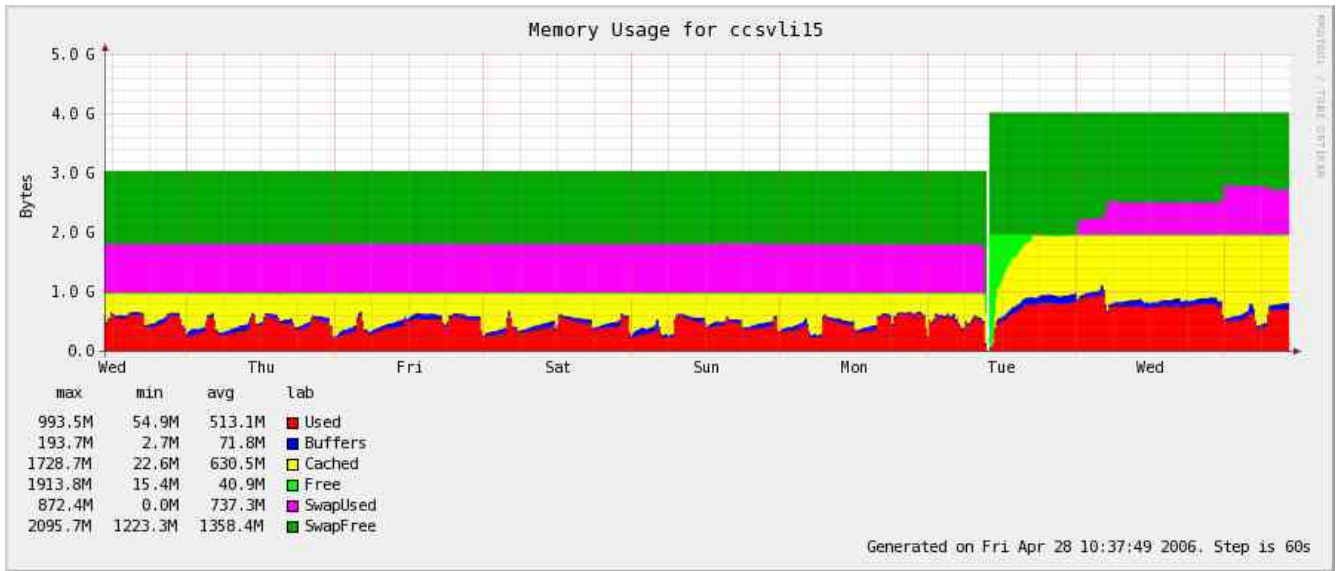
## A6. FTS feeding issues



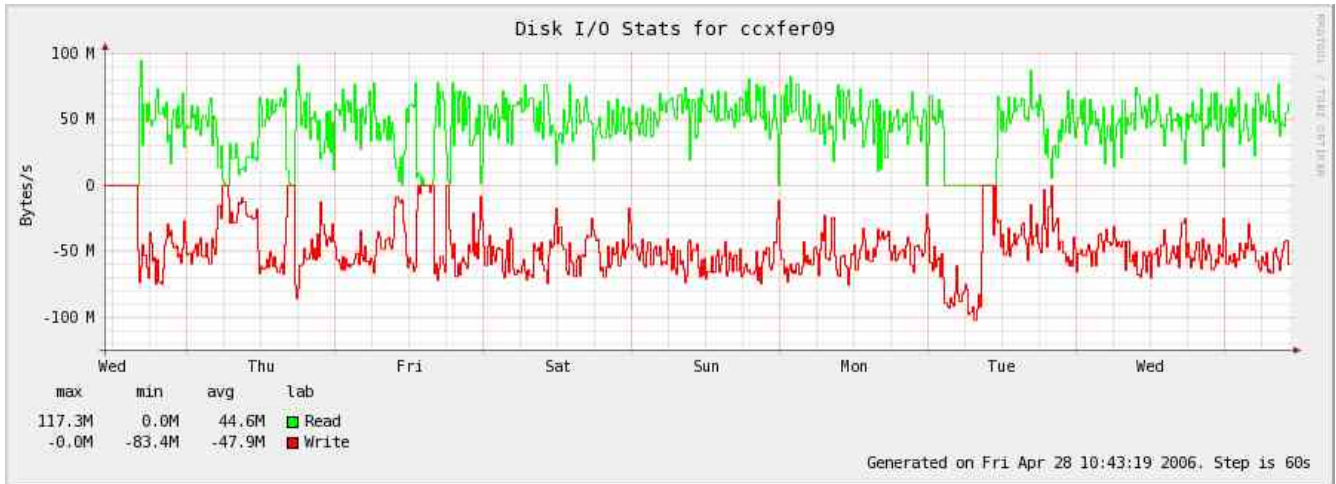
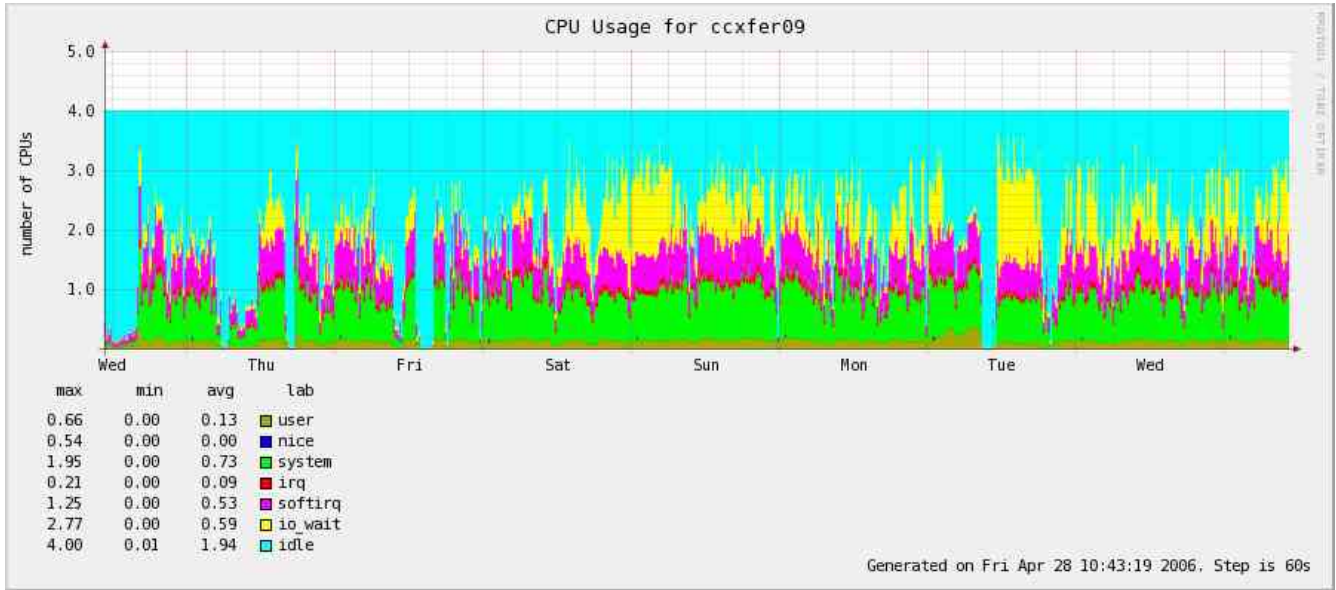
## A7. Head node stats (disk-tape): Wed 19/04- Wed 26/04



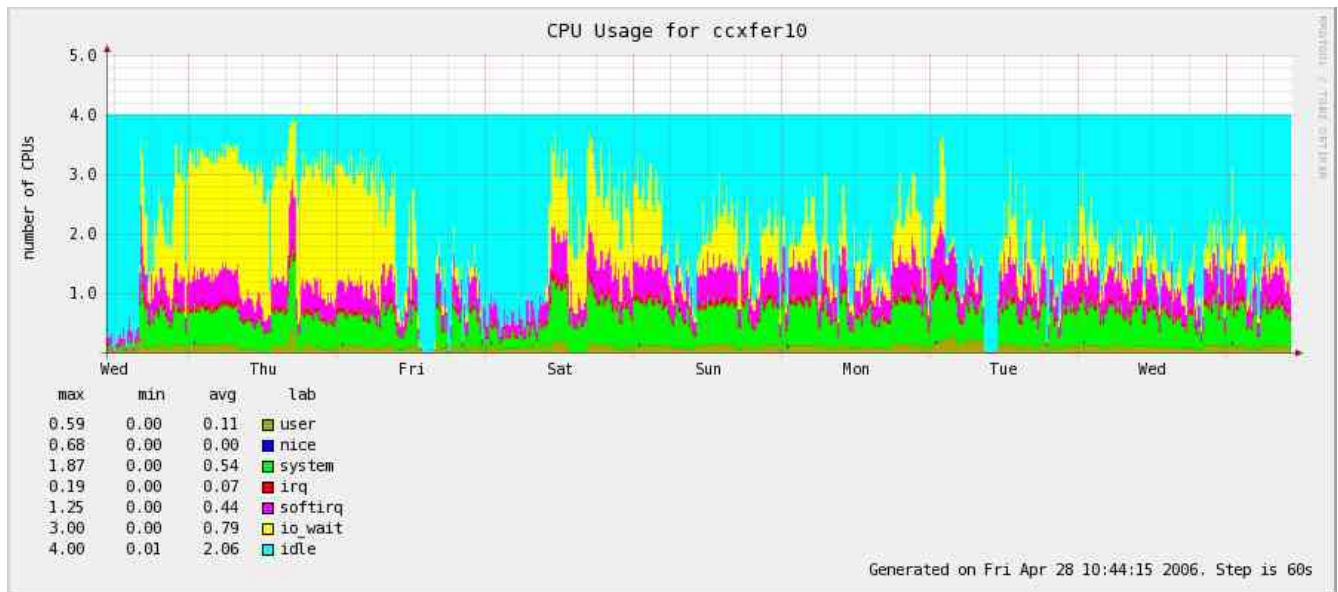
- All nights: dump databases and backup
- Every 4 hours: gridmap-file-update



## A8. ccxfer09 stats (disk-tape): Wed 19/04- Wed 26/04



## A9. ccxfer10 stats (disk-tape): Wed 19/04- Wed 26/04



Friday 21/04: change kernel parameters (

