# ALICE update - use of CASTOR2@CERN

CASTOR2 Delta Review

# ALICE Data Challenges - DAQ and Offline

- **DAQ**
  - Periodic, going to continuous, collaboration ALICE/IT to validate state-of-the-art HW and SW
    - Objectives
      - Top throughput: 3 GB/s event-building (met and exceeded last year)
      - Sustained rate: 1 GB/s to tape during 1 week (in progress)
      - Storage of TPC muon test data
- **Offline**
  - Distributed data production and analysis on the GRID
  - Objectives: Test of the Offline computing model
    - Production of MC data for software and detector performance studies
    - Tests of the LCG Grid services - Workload management, Storage management, network
    - Incorporation of computing centres in the ALICE Grid
    - Tests of ALICE application (ROOT, AliRoot, GEANT4, FLUKA) and Grid (AliEn) software
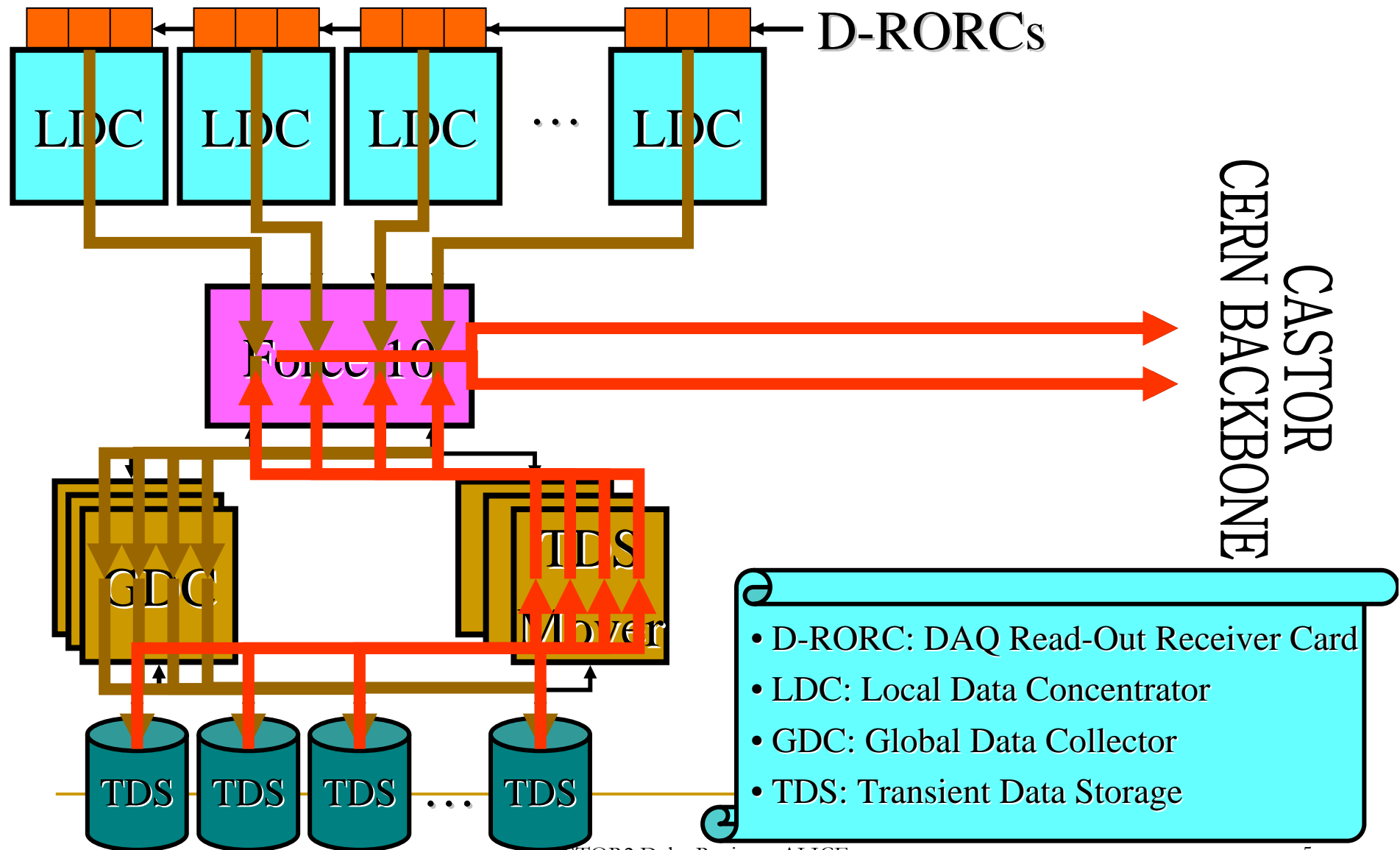    - End user data analysis

# ALICE Data Challenges - DAQ and Offline

- **Data challenges setup details**
  - DAQ **ADCVII** is running in the 'production LHC startup' conditions
  - ALICE **PDC'06** is running with conditions very close to the setup at LHC startup
- DAQ data challenge tests
  - Hardware, data acquisition, event building, storage
  - Networking
  - CASTOR client/server
  - Data registration and access in ALICE Offline Grid framework
- Offline data challenge tests
  - gLite WMS, SRM, FTS
  - Site storage (elsewhere) and CASTOR2@CERN

# DAQ ADC VII setup

- HW infrastructure:
  - 3 racks each with:
    - 16 TB of disk pool
    - 1 FiberChannel Switch
    - 9 hosts equipped with HostBusAdapters:
    - ADIC StorNext cluster file system allows the sharing of the disk arrays by the hosts
  - Force 10 Ethernet switch (same used for the CERN backbone)
  - All hosts have one GB Ethernet NIC used to move data
- SW infrastructure
  - LDCs, GDCs and TDS Movers run ALICE DAQ system DATE
  - TDS Movers run ALICE TDS Manager software and CASTOR client software
    - rfcp with MD5 checksum
- CASTOR2 - single instance 'castoralice'
  - 28 disk servers
  - Tape units for 1GB/sec tape throughput tests

# Architecture



D-RORCs

LDC    LDC    LDC    ...    LDC

Force 10

CASTOR
CERN BACKBONE

GDC

TDS Mover

TDS    TDS    TDS    ...    TDS

- D-RORC: DAQ Read-Out Receiver Card
- LDC: Local Data Concentrator
- GDC: Global Data Collector
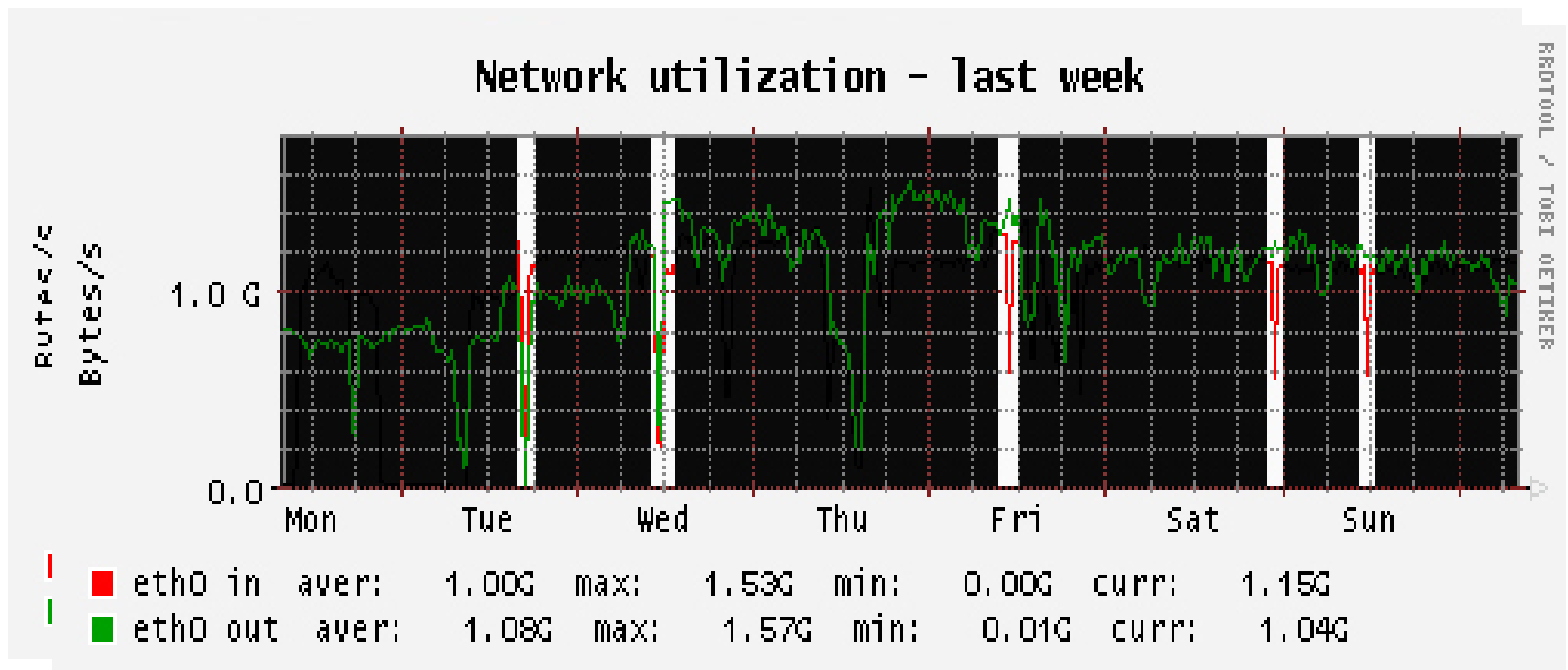- TDS: Transient Data Storage

# Pending items from DAQ ADCVII

- Random freezes
- Periodic slowdowns
- Low throughput/client
- Stuck clients

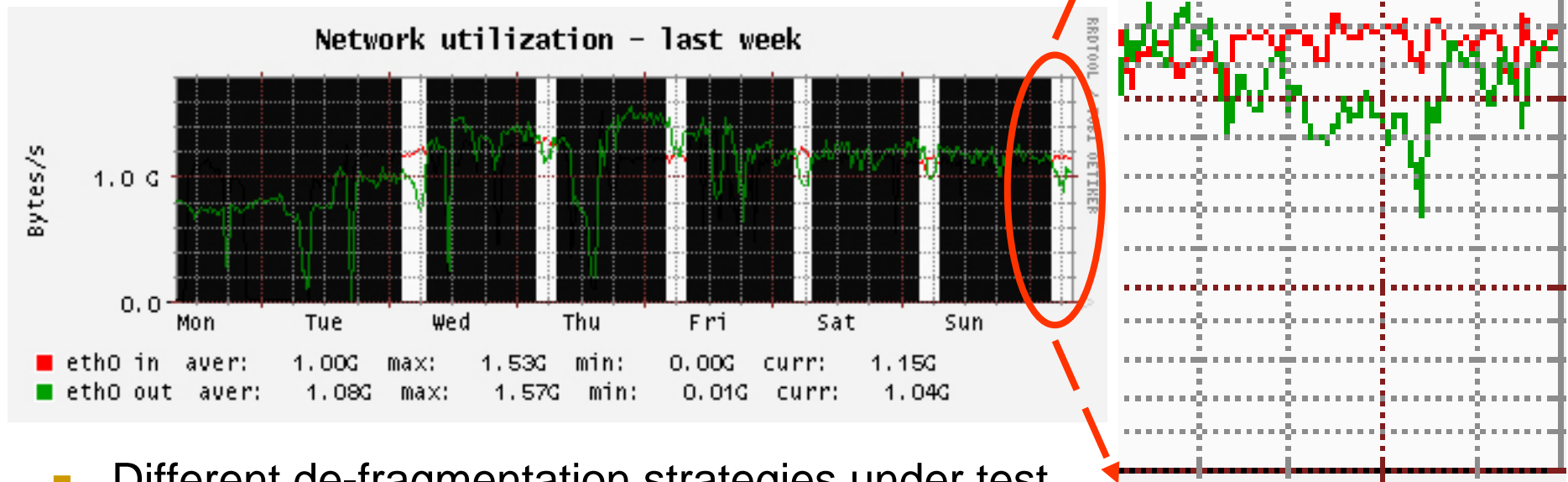# Freezes

- Effect discovered in ALICE tests
  - All rfcps blocked
  - Software in CASTOR runs a periodic check-and-restart
  - Investigation in progress



Network utilization – last week

# Slowdowns

- Almost every night around 04:00
- Possibly due to CASTOR2 disk servers de-fragmentation
- No effects on the DAQ (only on CASTOR)



Network utilization – last week

| | | | | |
|---|---|---|---|---|
| ■ eth0 in  aver: | 1.00G | max: 1.53G | min: 0.00G | curr: 1.15G |
| ■ eth0 out aver: | 1.08G | max: 1.57G | min: 0.01G | curr: 1.04G |

- Different de-fragmentation strategies under test
- Using a different file system on the servers @ IT may help
    - impossible for the moment
- Other slowdowns during the day
    - for example between 18:00 and 22:00 (but not every day)

# Low throughput and stuck clients

- ALICE uses "standard" rfcp to transfer files into CASTOR

- One single rfcp uses ~ ½ of the capacity of the outgoing NIC link, non-linear increase of transfer speed with multiple rfcps

  - Full use of the outgoing link with 9 simultaneous rfcps, still dips in the throughput

  - With/without MD5 checksum: same results

- rfcps gets stuck - ~ 10 times/day (out of 90,000 transfers)

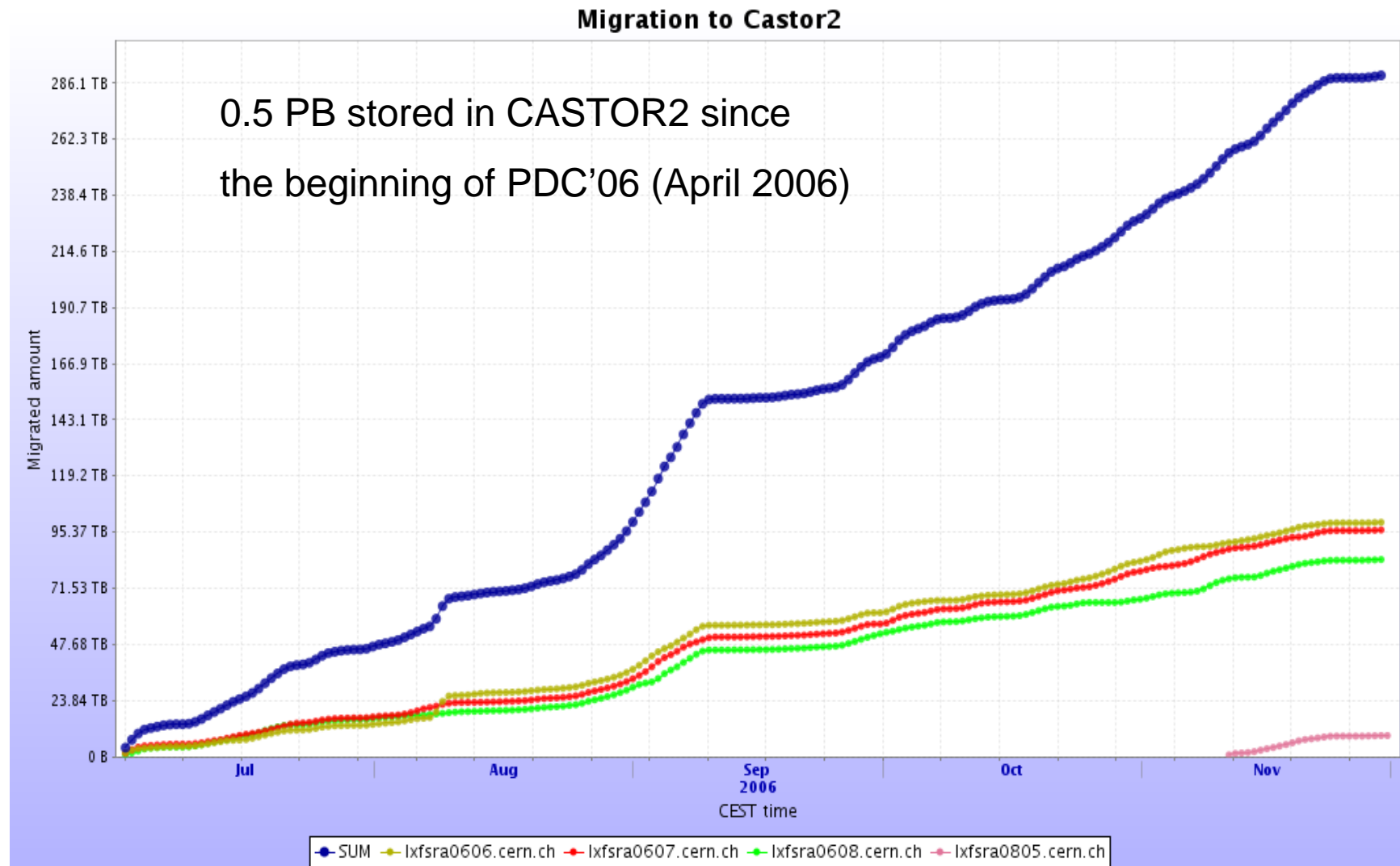  - Timeout and trap for slow transfers added on ALICE side

# Counters

- **Statistics of typical 'DAQ daily transfer activity'**
  - File size: slightly below 1 GB
  - ~ 1 GB/s/day (sustained)
  - ~ 90,000 files/day moved OK
  - ~ 320 transfers/day fail with various error messages
  - ~ 240 transfers/day trigger a timeout

- **Undetected, any of these failures (random freezes, slowdowns, low throughput/client, stuck clients) has the power to block the challenges**
  - Workarounds and protections are implemented in the Offline/DAQ transfer clients
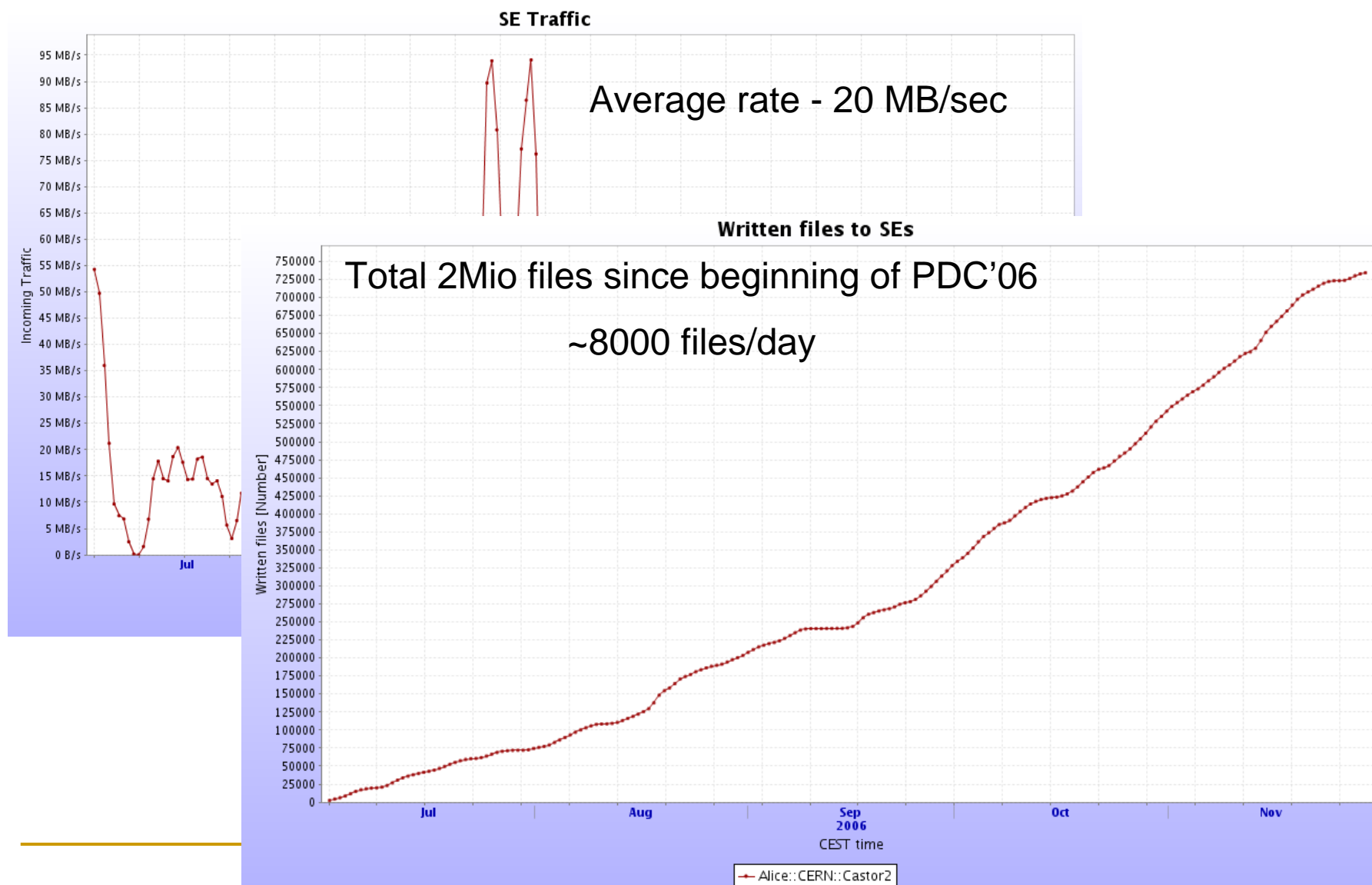  - Need to be fixed

# Offline PDC'06 structure

- Currently 50 active computing centres on 4 continents:
  - Average of 2000 CPUs running continuously
  - MC simulaton/reconstruction jobs (~9 hours/job) produce 1.2 GB output, spread over several storage elements
    - 1 GB (simulated RAW data + ESDs) - stored at at CASTOR2@CERN

- CASTOR2 - single instance 'castoralice'
  - 4 tactical disk buffers (4x5.5TB) running xrootd
  - Act as a buffer between CASTOR2 and clients running at the computing centres
  - CASTOR2 stagein/stageout - through a stager package running on the xrootd buffers
- FTS data transfer
  - Direct access to CASTOR2 (through SRM/gridftp)
  - Scheduled file transfer T0->T1 (five T1s involved)
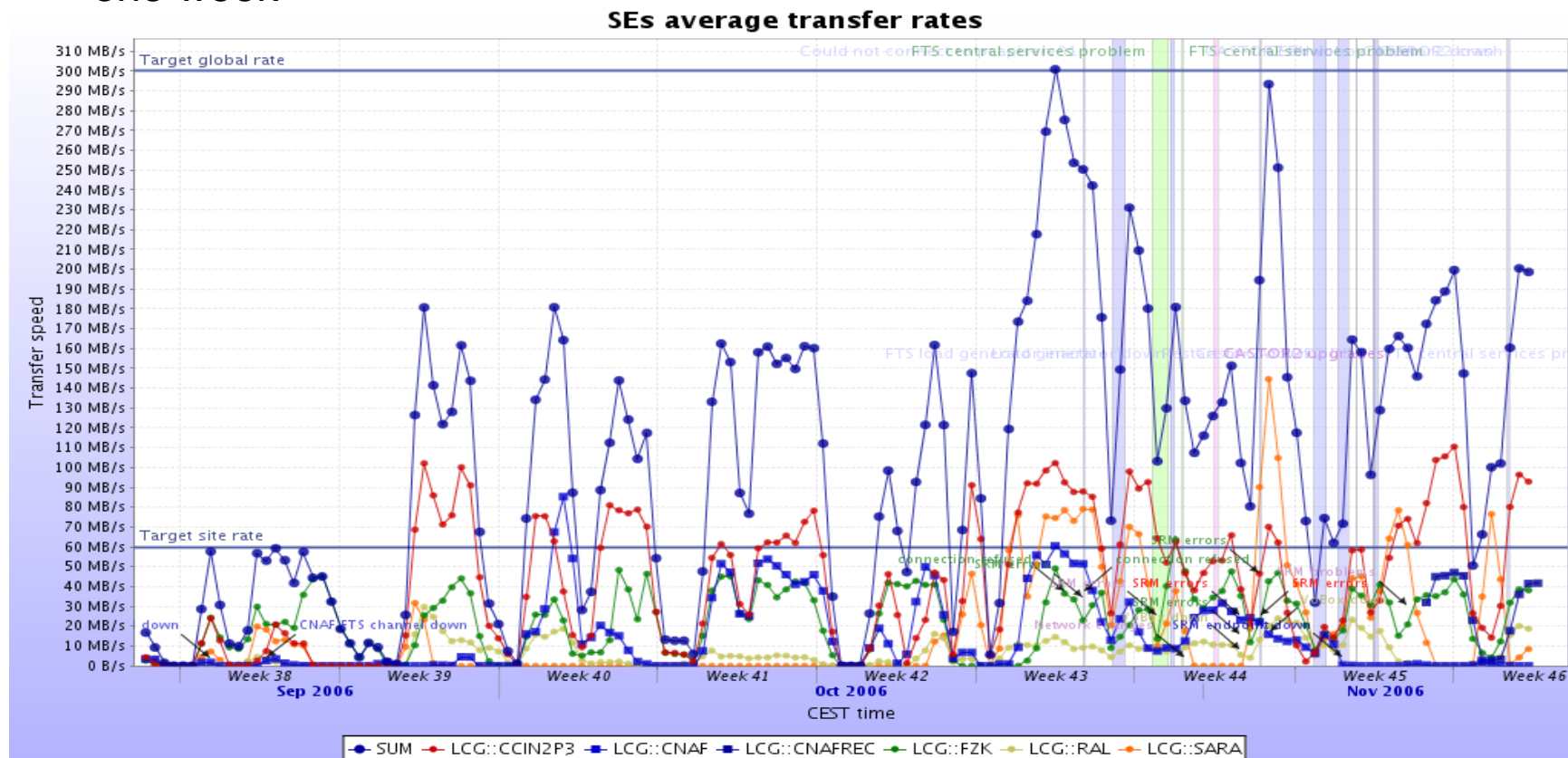
# Production data volume in the past 5 months

**Migration to Castor2**

0.5 PB stored in CASTOR2 since

the beginning of PDC'06 (April 2006)

# Files and rates in the past 5 months



**SE Traffic**

Average rate - 20 MB/sec

**Written files to SEs**

Total 2Mio files since beginning of PDC'06

~8000 files/day

# Data transfers with FTS

- CASTOR2@CERN is the data source
- Transfers to CNAF (Italy), GridKA (Germany), CCIN2P3 (France), RAL (UK), SARA (The Netherlands)
- Goal of the exercise - sustained throughput of 300MB/sec out of CERN for one week

# Issues in Offline use of CASTOR2

- See issues in slides 7-9
  - File recall from CASTOR2 - long waiting times for files to be staged from tape
    - Later requests are executed faster than early ones
  - Looping requests for a file recall (not properly retried in CASTOR2) can block the entire instance 'castoralice'
- Most of these issues are addressed by adding checks and retries in the Offline stager package
- In general the use of xrootd buffer 'dampens' the effects of CASTOR2 temporary problems, however this is not our long-term solution
  - We are about to exercise the new version of CASTOR2 with xrootd support (no buffers)
- FTS transfers: interplay between FTS software stack and remote storage
  - Rate out of CERN - no problem to attain 300MB/sec from CASTOR
  - Stability issues - we estimate that problems in CASTOR2 account only for ~7% of the failed transfers

# Expert support

- Methods of problem reporting
  - E-mails to the CASTOR2 list of experts (daily and whenever a problem occurs)
  - Regular meetings with CASTOR2 experts in the framework of ADCVII and PDC'06 challenges
  - Helpdesk tickets by direct CASTOR2 users (very little use compared to the data challenges)
- ALICE is satisfied with the level of expert support from CASTOR
  - All our queries are addressed in a timely manner
  - Wherever possible, solutions are offered and implemented immediately
  - Bug reports are evaluated and taken into consideration
- We are however worried about the sustainability of the current support in a 24/7 mode of operation - the problems often appear to be non-trivial and require intervention by very few (top level) experts

# Conclusions

- ALICE is using CASTOR2@CERN in a regime exactly as (DAQ) and very close (Offline) to the LHC startup 'production' mode

- The identified CASTOR2 issues reported in this presentation need to be addressed on a system level
  - Currently majority of the workarounds are built in the client software (DAQ/Offline frameworks)
  - Offline critically dependent on the integration of xrootd in CASTOR2 - progress is rather slow
  - All issues have been reported to the CASTOR2 team, however we are unable to judge at this point how quickly these will be resolved

- We are satisfied with the expert support provided by the CASTOR2 team
  - Reporting lines are clear and the feedback is quick
  - We are worried about the sustainability of present support structure