# Tier0 Status

## Tony Cass
(With thanks to Miguel Coelho dos Santos & Alex Iribarren)

## LCG-LHCC Mini Review, 1st July 2008

- Resource Ramp-up

- CASTOR Performance and Metrics

- Power Issues and Progress

# Agenda

- **Resource Ramp-up**
- CASTOR Performance and Metrics
- Power Issues and Progress

# Elonex Issues

- ## Resource Ramp-up
  - Tier0 purchasing affected by two issues
    - Elonex bankruptcy
    - Disk server problems under heavy load
- CASTOR Performance and Metrics
- Power Issues and Progress

# Ramp-up: Problems

- **Elonex Bankruptcy**
  - One disk server order rapidly switched to alternative suppliers.
  - Second disk server order plus CPU order switched to alternative suppliers after FC in March.

- **Disk server load issues**
  - Problems brought to light by improvements to hardware burn-in procedure.
    - Load to provoke issues significantly exceeds normal load on disk servers.
    - Previous generation servers also show the problems with extremely high load.
    - New capacity now released to deployment (and many servers have run well for some time with no issues.

# Ramp-up: Current state

- **CPU**
  - 100% of pledge delivered in early May, i.e. with one month delay

- **Disk**
  - 52% of pledge delivered to experiments in early May.
  - Balance of pledge is at CERN and will deployed progressively in coming weeks.
  - Delay, but minimal impact on CCRC exercise

- **Tape**
  - 100% of pledge available well before April

- **2009 Procurements**
  - On schedule: tenders for September FC adjudication opened; tenders for December adjudication to be sent out shortly.

# Agenda

- Resource Ramp-up

- **CASTOR Performance and Metrics**

- Power Issues and Progress

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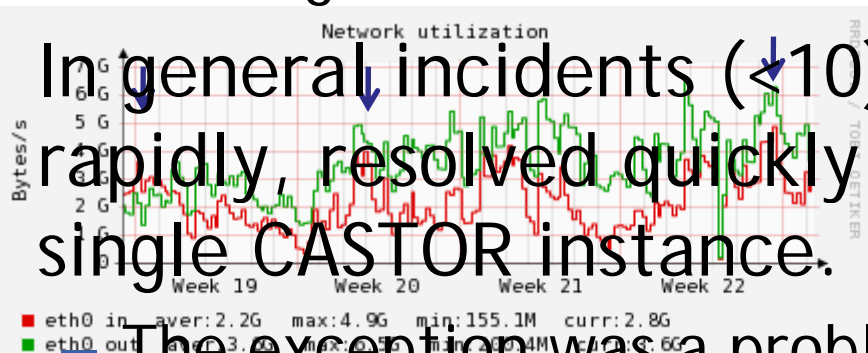**www.cern.ch/it**

- Resource Ramp-up

- **CASTOR Performance and Metrics**
  - CASTOR Service
  - SRM Interface
  - CASTOR metrics

- Power Issues and Progress
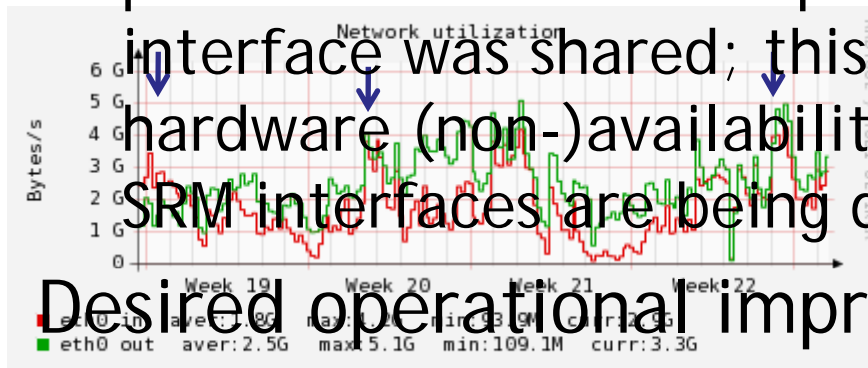
- CASTOR ran well throughout the May CCRC although end-user load seemed low…

- In general incidents (<10) were detected rapidly, resolved quickly and only affected a single CASTOR instance.
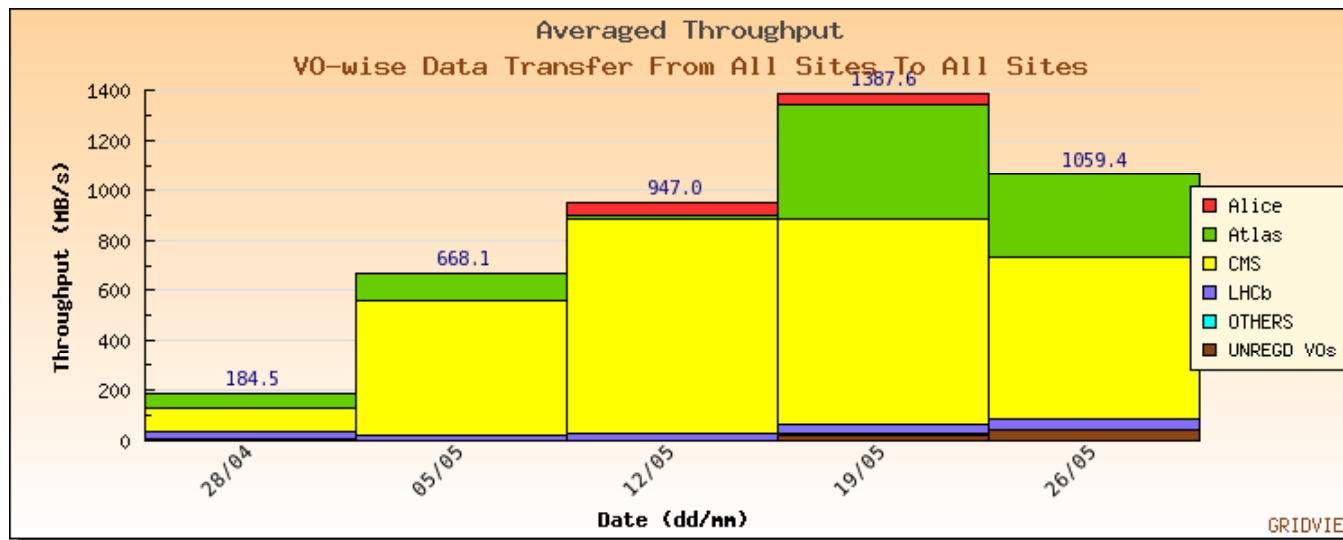
  - The exception was a problem on the CASTOR public service which impacted ATLAS as the SRM interface was shared; this configuration was due to hardware (non-)availability; the planned dedicated SRM interfaces are being deployed.

- Desired operational improvements (notably tape request prioritisation) are being deployed with promising initial results.

Overall disk cache end-user:

Network utilization

eth0 in  aver:2.2G  max:4.9G  min:155.1M  curr:2.8G
eth0 out

Three problems with the garbage collection mechanism for CMS contributed to large peaks of internal traffic…

CASTORCMS disk cache throughput

Network utilization

eth0 in
eth0 out  aver:2.5G  max:5.1G  min:109.1M  curr:3.3G

It showed the I/O capacity of the CMS setup (peaks of 9GB/s (in+out))

- When it works it works well
- A large volume of data was transferred
- The average rate was high



- Reliability is still an issue
- ~10 incidents with impact ranging from service degradation to complete unavailability

# SRM Interface --- II

- May 5 – redundant SRM back-ends lock each other in database [ALL VOs]
- May 13th – lack of space on SRM DB [LHCb]
- May 13th – DB "extreme locking" / DB deadlocks [ALL VOs]
- May 9th, May 14th, May 19th – SRM 'stuck' / no threads to handle requests [ATLAS]
- May 21st, May 24th – slow stager backend causes SRM stuck / DB overload [ All VOs]
- May 30th – get Timeouts due to slowness on Castor backend [ATLAS, LHCb]
- 3 times in May – problematic use of soft pinning caused GC problems [CMS]
- June 6th – patch update crashed backend servers [ATLAS, ALICE, CMS]

- ## To be improved:
  - Better resiliency to problems
  - More service decoupling
  - Some bugs need to be fixed
  - Better testing needs to be done
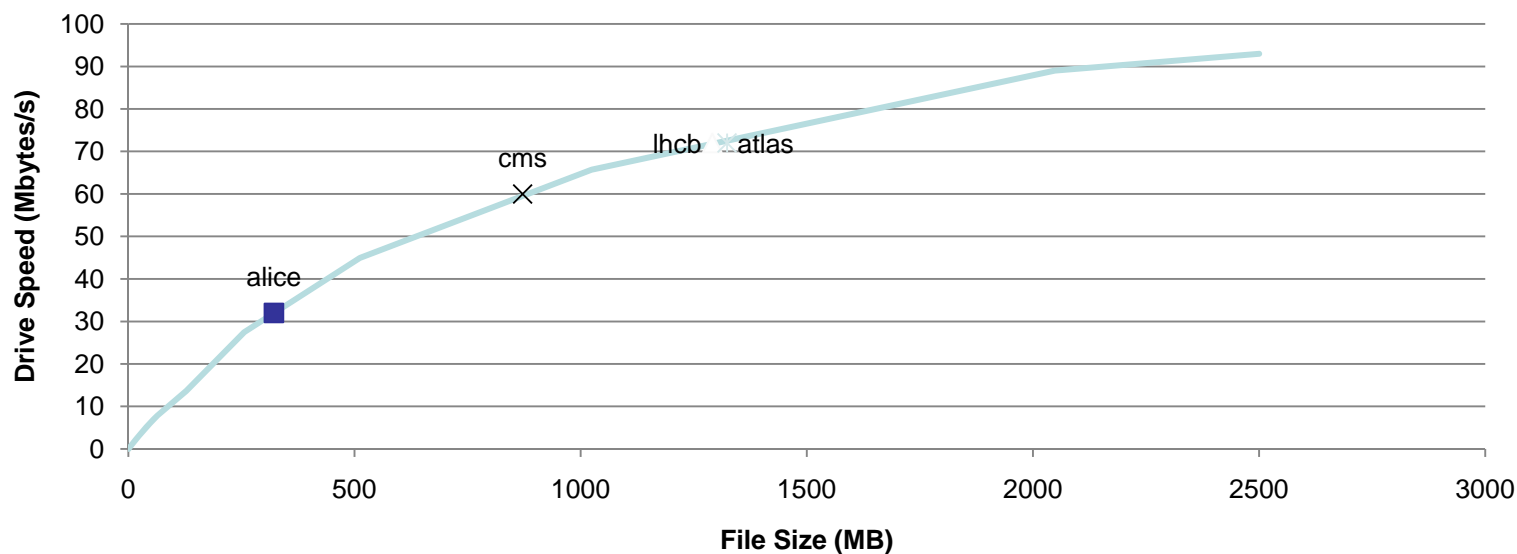
# SRM Interface --- III

- Separate out LHC VOs from shared instance **Done**
- Migrate all SRM databases to Oracle RAC **Done**
  (done for ATLAS)
- Upgrade to SRM 2.7 and deploy on SLC4*
  **In (very early) test**
  - Redundant backends
  - Uses CASTOR 2.1.7 API which allows deployment of redundant stager daemons
  - Deploy fixes for identified bugs
- Configure SRM DLF to send logs to **In test**
  appropriate stager DLF*
  - Improve our debugging response time
- Continue improving service monitoring

\* Required for "time to turl" metric; could be delivered in ~1 month.

- Metric implementation continues as collaboration between developers and operations teams.
- Improved instrumentation rolled out during CCRC in May
  - so few measurements for LHCb or ALICE
  - no upload to Lemon; excel plots only
- New Lemon sensor for CASTOR will be deployed in the near future to deliver automatic generation of metric plots
  - (this version centralises all daemon monitoring, so needs the 2.1.7-10 CASTOR release.)
- The following slides show a selection of metric plots which cover performance and issues during the May CCRC.
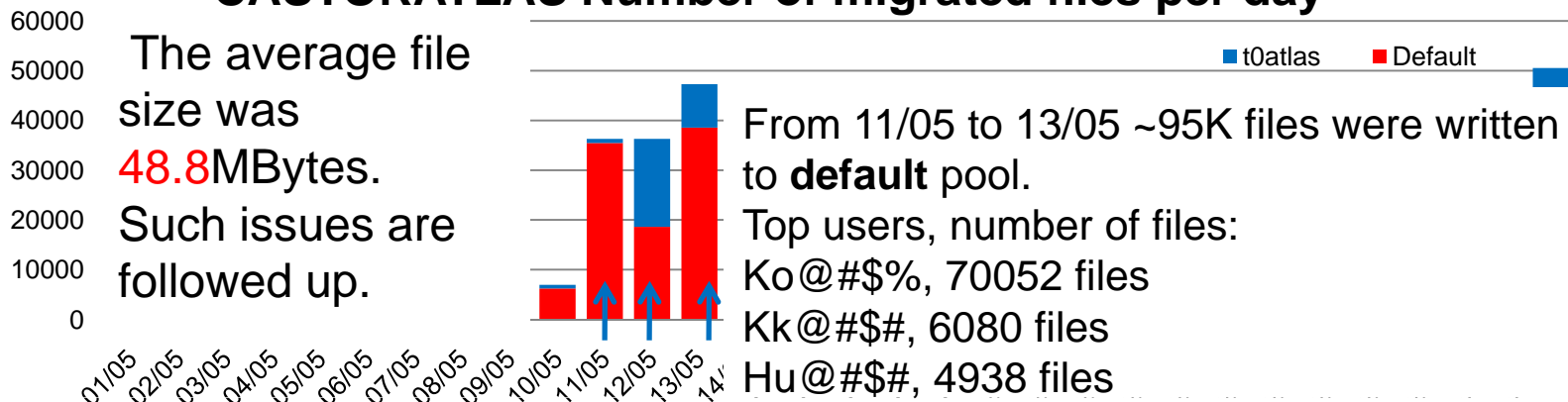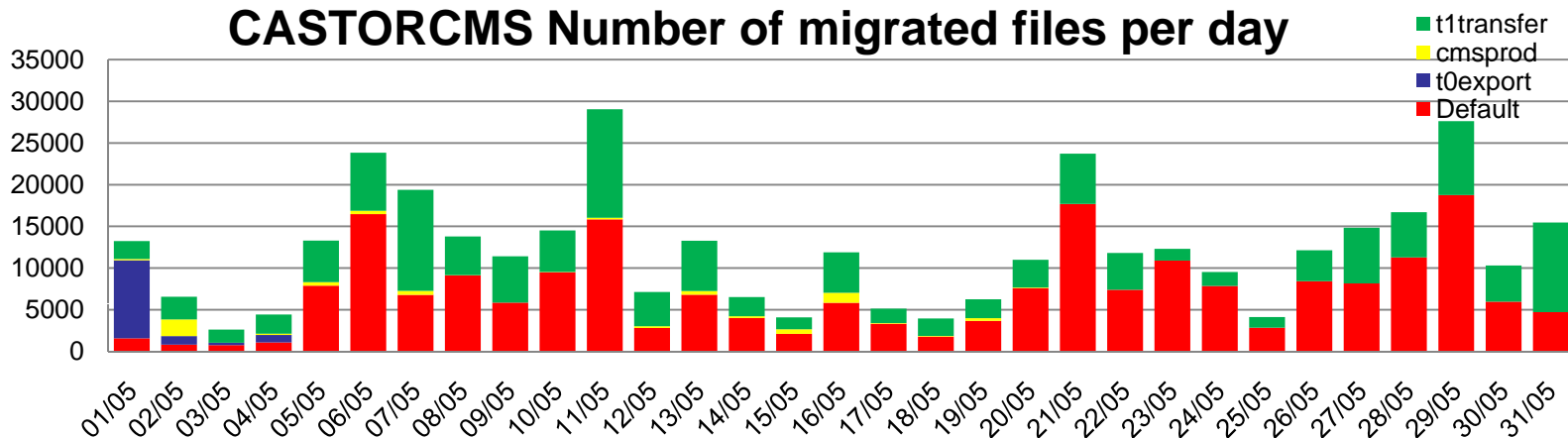
# File size and performance

## Typical Drive Performance



| Date | Alice | Atlas | CMS | LHCb |
|---|---|---|---|---|
| CCRC May '08 | **322 MB** | **1291 MB** | 872 MB | **1327 MB** |
| March '08 | 143 MB | 230 MB | **1490 MB** | 865 MB |
| CCRC Feb '08 | 340 MB | 320 MB | **1470 MB** | 550 MB |
| Jan '08 | 200 MB | 250 MB | **2000 MB** | 200 MB |

## CASTORATLAS Number of migrated files per day

The average file size was **48.8**MBytes. Such issues are followed up.

From 11/05 to 13/05 ~95K files were written to **default** pool.
Top users, number of files:
Ko@#$%, 70052 files
Kk@#$#, 6080 files
Hu@#$#, 4938 files

Legend: t0atlas, Default

## CASTORCMS Number of migrated files per day

Legend: t1transfer, cmsprod, t0export, Default

CASTORATLAS AVERAGE FILESIZE MIGRATED

t0atlas    default    t0atlas
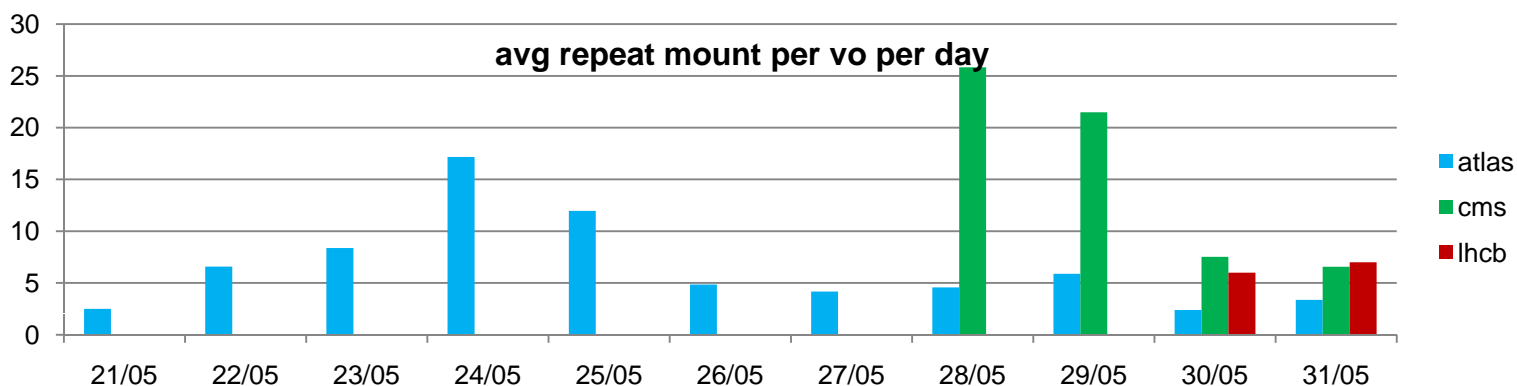
CASTORATLAS DATA RATE TO TAPE

CASTORATLAS FILES MIGRATED

**T0ATLAS is working well.**    Already mentioned issue on default pool...

# Repeat Tape Mounts

- Repeated (read) tape mounts per VO

**max repeat mount per vo per day**

atlas / cms / lhcb (bar chart, values 0–180 scale, dates 21/05 – 31/05)

**avg repeat mount per vo per day**

atlas / cms / lhcb (bar chart, values 0–30 scale, dates 21/05 – 31/05)

High value of repeat mounts for reading...

CASTORATLAS Get/Put latency

Tier-0 exercise

Power Cut!

CASTORCMS Get/Put latency

Tier1 Data Import

GC Problem
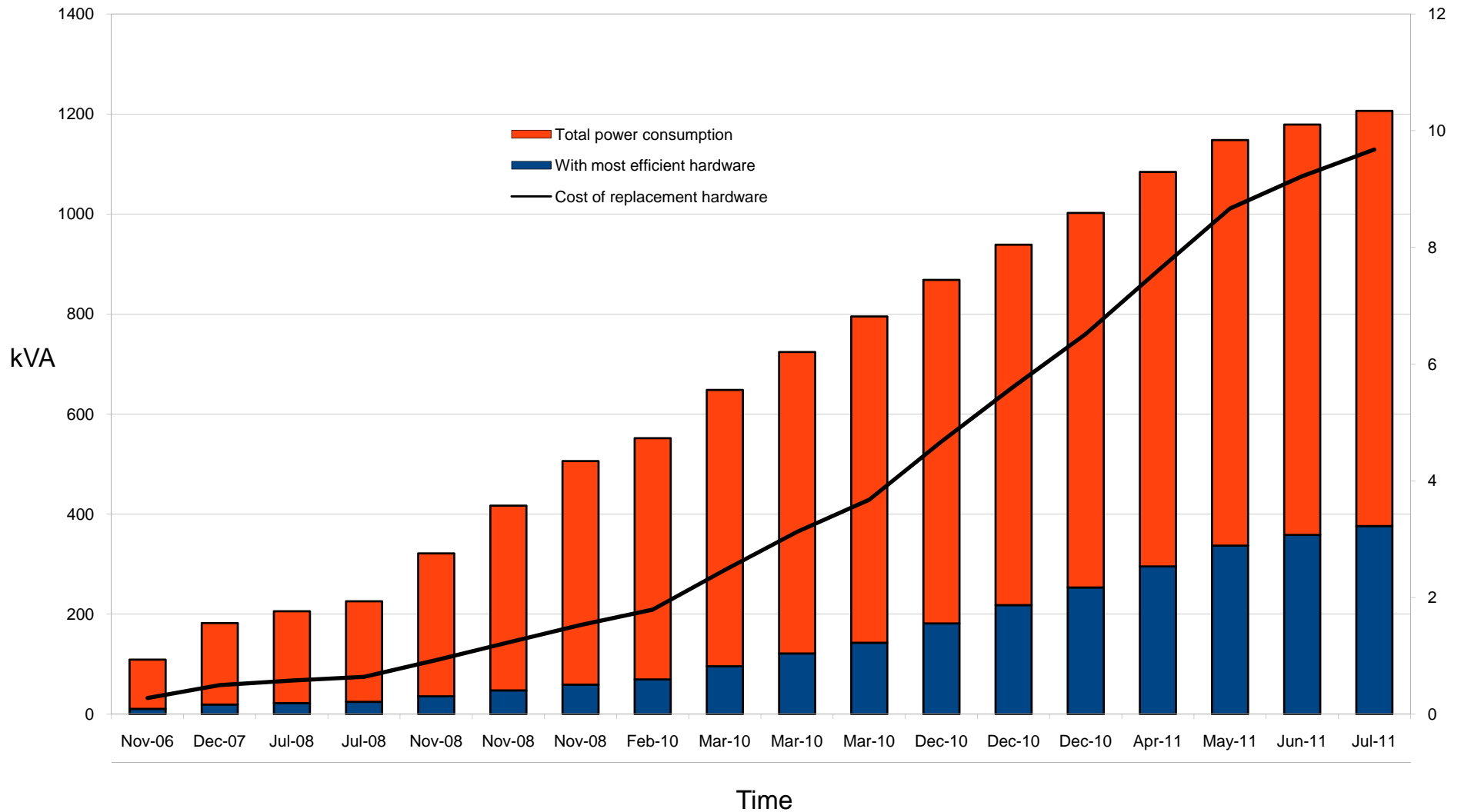
# Agenda

- Resource Ramp-up
- CASTOR Performance and Metrics
- **Power Issues and Progress**

- Resource Ramp-up

- CASTOR Performance and Metrics

- **Power Issues and Progress**

  – B513 status

  – New Computer Centre Planning

  – Covering the Gap

## CPU Server Power Consumption



Legend:
- Total power consumption
- With most efficient hardware
- Cost of replacement hardware

Y-axis (left): kVA — 0, 200, 400, 600, 800, 1000, 1200, 1400
Y-axis (right): 0, 2, 4, 6, 8, 10, 12

X-axis (Time): Nov-06, Dec-07, Jul-08, Jul-08, Nov-08, Nov-08, Nov-08, Feb-10, Mar-10, Mar-10, Mar-10, Dec-10, Dec-10, Dec-10, Apr-11, May-11, Jun-11, Jul-11

# New Computer Centre Planning

- In-house design and construction of new centre not possible (TS effort focussed on Linac4).
- No desire to tender for turn-key design and construction
  - Lowest cost bidder wins…
- Four phase process developed:
  1. Request (many) conceptual designs
  2. Commission 3-4 companies submitting conceptual designs to develop an outline design
  3. In-house, turn a selected outline design into plans and documents enabling
  4. Single tender for overall construction.
- "Call for proposals" for the conceptual designs sent out (deadline July 18th); process could lead to negotiation of construction contract end 2009.
  - Estimate subsequent detailed design phase of ~6 months and construction phase of ~18 months
  - New centre available for equipment installation in Jan 2012

# Covering the Gap

- B513 OK until end 2010 +
  New Centre from Jan 2012 ==>
  - Need to cover 2011 installations
  - Plus 1H2012 installations in case of construction delays.
- Tier1 centres asked for possible spare capacity in this window.
  - Oslo could be a possibility: Completing 2MW facility end-2009, but only need 1MW initially.
    - Discussions on modalities (hardware requirements, operation model, …) to start soon.
- Reviewing co-lo options within ~1hr of CERN.
  - No spare capacity at present,
  - Options possible on 2011 timeframe, but still at ~2kW/m$^2$, so likely very expensive.