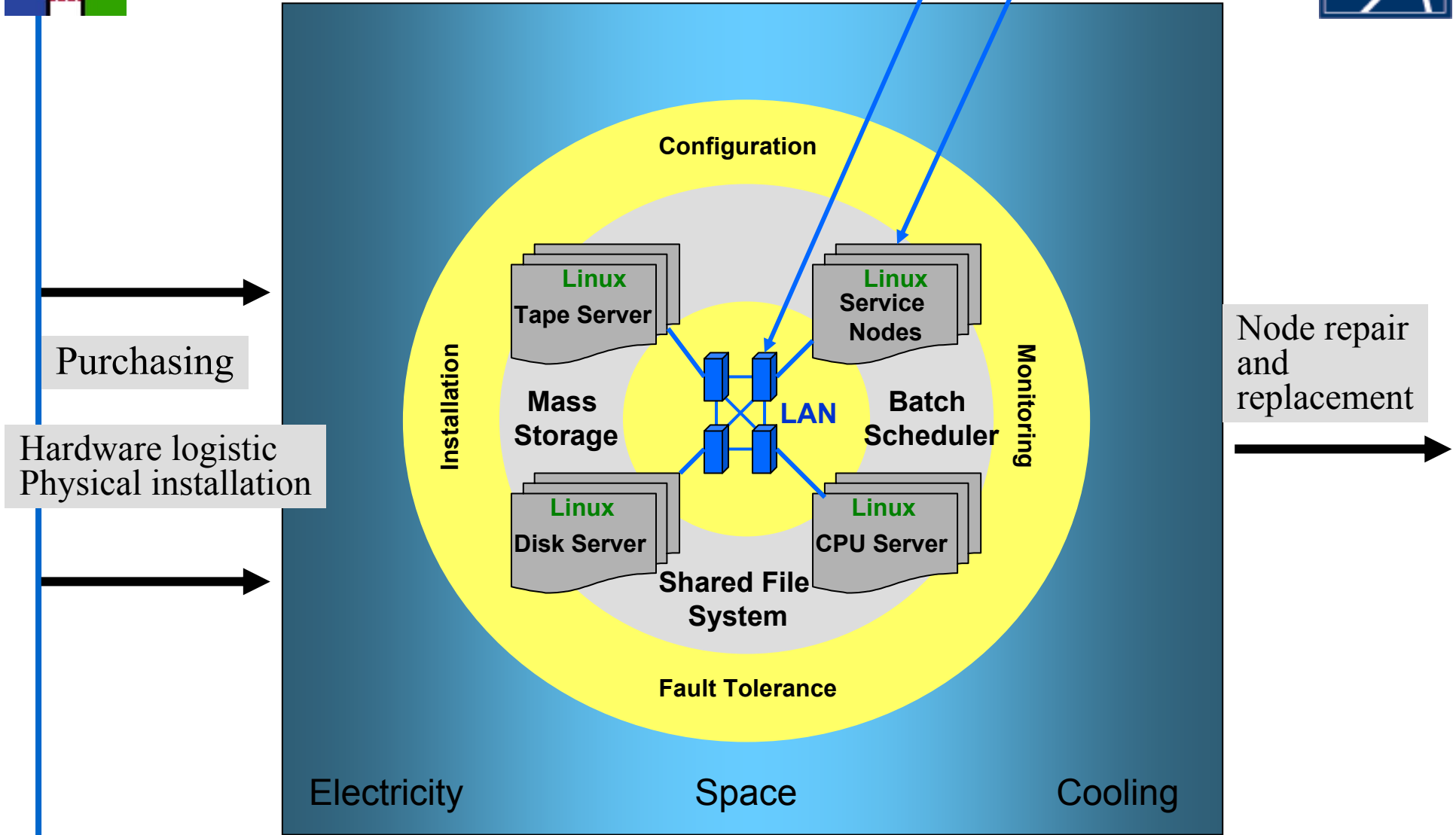# LHCC Review

# LCG Fabric Area

**Bernd Panzer-Steindel,  Fabric Area Manager**

# Material purchase procedures

**Long discussions inside IT and with SPL about
the best future purchasing procedures**

**new proposal to be submitted to finance committee in December:**

▪**for CPU and disk components**

▪**covers offline computing and physics data acquisition (online)**

▪**no 750 KCHF ceiling per tender**

▪**speed up of the process (e.g. no need to wait for a finance committee meeting)**

▪**effective already for 2005**

# Electricity and cooling

Upgrade of the electrical and cooling power to 2.5 MW



Installation of new transformer and their electrical connection to the existing infrastructure.  All milestones met within 1-2 weeks.
~900 KW available until mid 2006, currently running at ~550 KW

Cooling upgrade on track,  discussion about financial issues between IT and TS

# Space



Refurbishment of the left side of the Computer Center has started

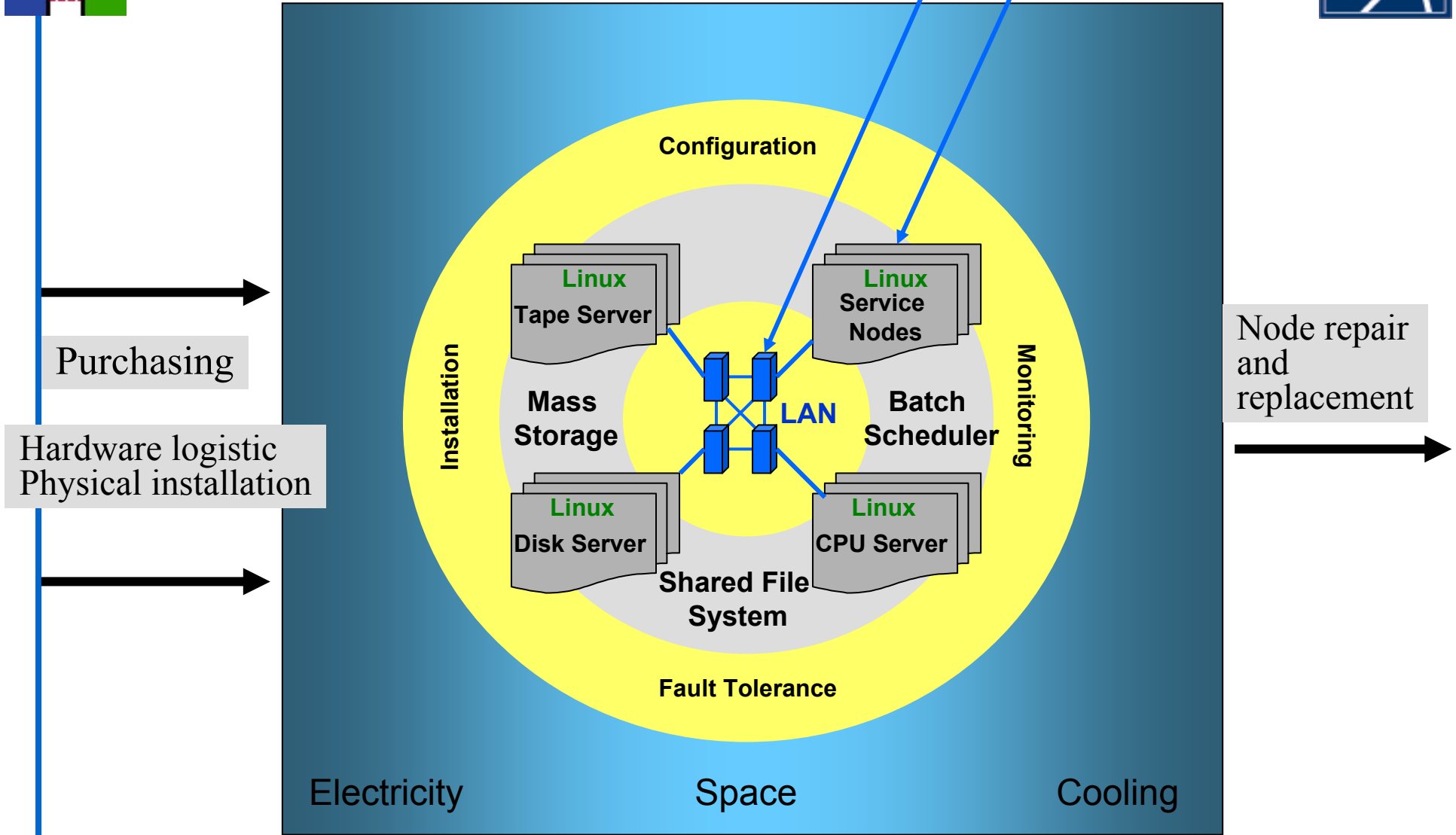New structure on the already refurbished right side

**During the period May-August more than 800 nodes were moved with minor service interruptions (AFS, NICE, MAIL, WEB, CASTOR server, etc.)**
**Very man-power intensive work , tedious and complicated scheduling.**

# Space



Before and after the move

WAN

Configuration

Installation

Monitoring

Linux Tape Server

Linux Service Nodes

Mass Storage

LAN

Batch Scheduler

Linux Disk Server

Linux CPU Server

Shared File System

Fault Tolerance

Electricity          Space          Cooling

Purchasing

Hardware logistic Physical installation

Node repair and replacement
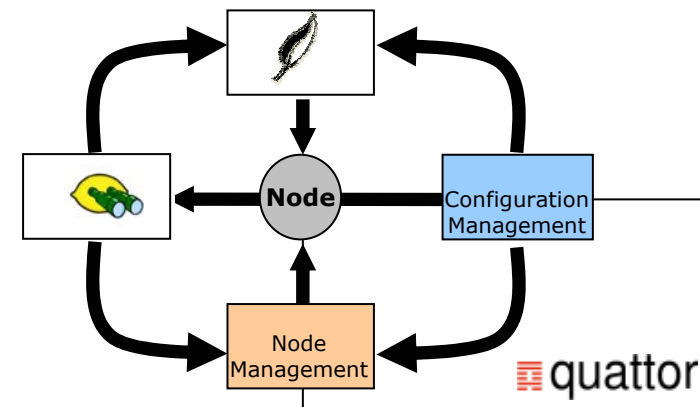
# Fabric Management with ELFms

ELFms stands for 'Extremely Large Fabric management system'

Subsystems:

- ◆ **quattor** : configuration, installation and management of nodes
- ◆ : system / service monitoring
- ◆ : hardware / state management



- ◆ ELFms manages and controls most of the nodes in the CERN CC
  - ~2100 nodes out of ~ 2700
  - Multiple functionality and cluster size (batch nodes, disk servers, tape servers, DB, web, …)
  - Heterogeneous hardware (CPU, memory, HD size,..)
  - Supported OS: Linux (RH7, RHES2.1, Scientific Linux 3 – IA32&IA64) and Solaris (9)

# Quattor

Quattor takes care of the *configuration, installation* and *management* of fabric nodes

➔ A **Configuration Database** holds the 'desired state' of all fabric elements

- Node setup (CPU, HD, memory, software RPMs/PKGs, network, system services, location, audit info…)

- Cluster (name and type, batch system, load balancing info…)

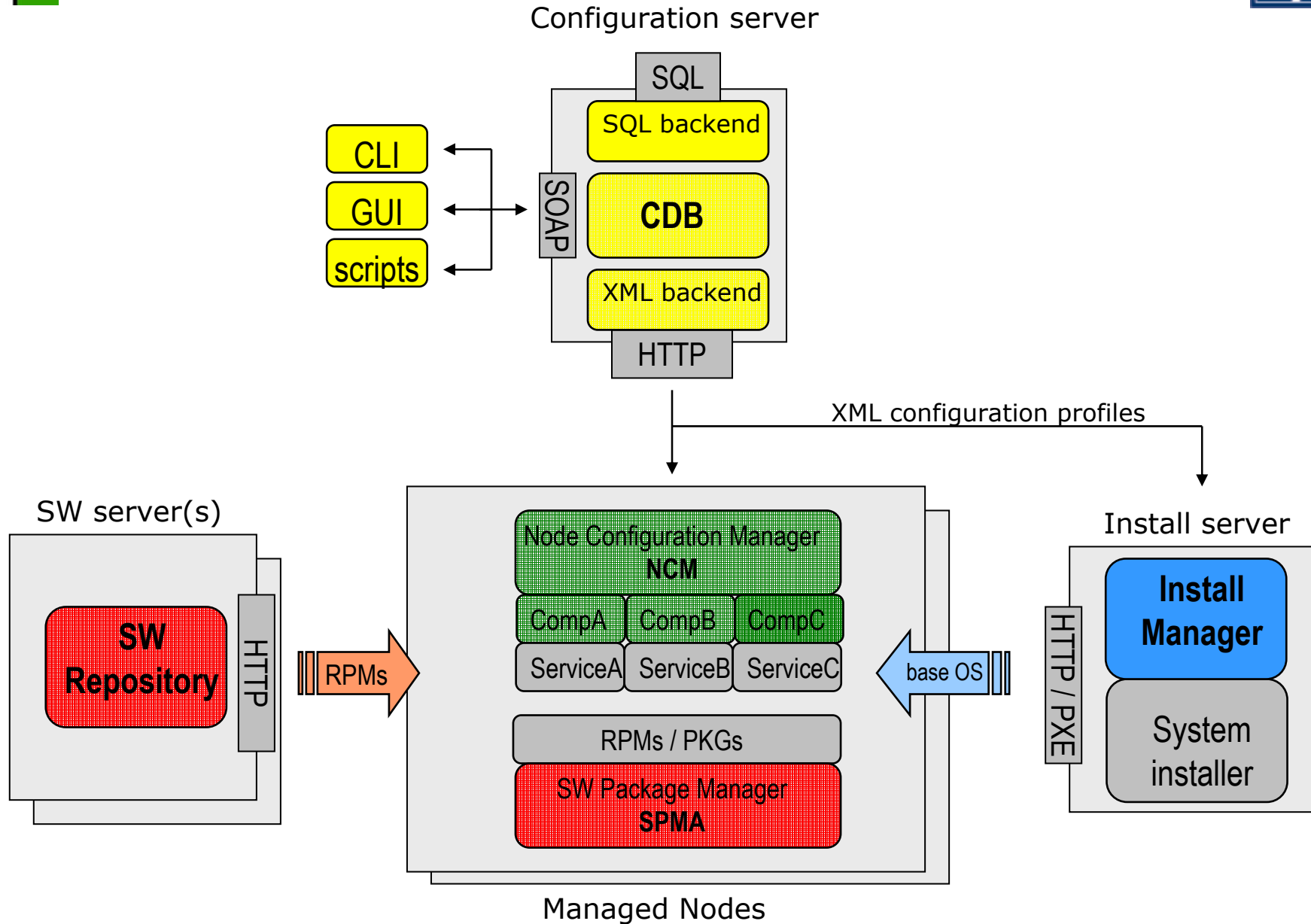- Defined in templates arranged in hierarchies – common properties set only once

➔ Autonomous management agents running on the node for

- **Base installation**

- **Service (re-)configuration**

- **Software installation and management**

- Quattor was initially developed in the scope of EU DataGrid. Development and maintenance now coordinated by CERN/IT

# Architecture



Configuration server

SQL

SQL backend

CDB

XML backend

HTTP

CLI
GUI
scripts

SOAP

XML configuration profiles

SW server(s)

SW Repository

HTTP

RPMs

Managed Nodes

Node Configuration Manager
NCM

CompA  CompB  CompC

ServiceA  ServiceB  ServiceC

RPMs / PKGs

SW Package Manager
SPMA

base OS

Install server

Install Manager

System installer

HTTP / PXE

# Quattor Deployment

- Quattor in complete control of Linux boxes (~ 2100 nodes, to grow to ~ 8000 in 2006-8)
  - Replacement of legacy tools (SUE and ASIS) at CERN during 2003

- CDB holding information of > 95% of systems in CERN-CC

- Over 90 NCM configuration components developed
  - From basic system configuration to Grid services setup… (including desktops)

- SPMA used for managing all software
  - ~ 2 weekly security and functional updates (including kernel upgrades)
  - Eg. KDE security upgrade (~300MB per node) and LSF client upgrade (v4 to v5) in 15 mins, without service interruption
  - Handles (occasional) downgrades as well

- Developments ongoing:
  - Fine-grained ACL protection to templates
  - Deployment of HTTPS instead of HTTP (usage of host certificates)
  - XML configuration profile generation speedup (eg. parallel generation)

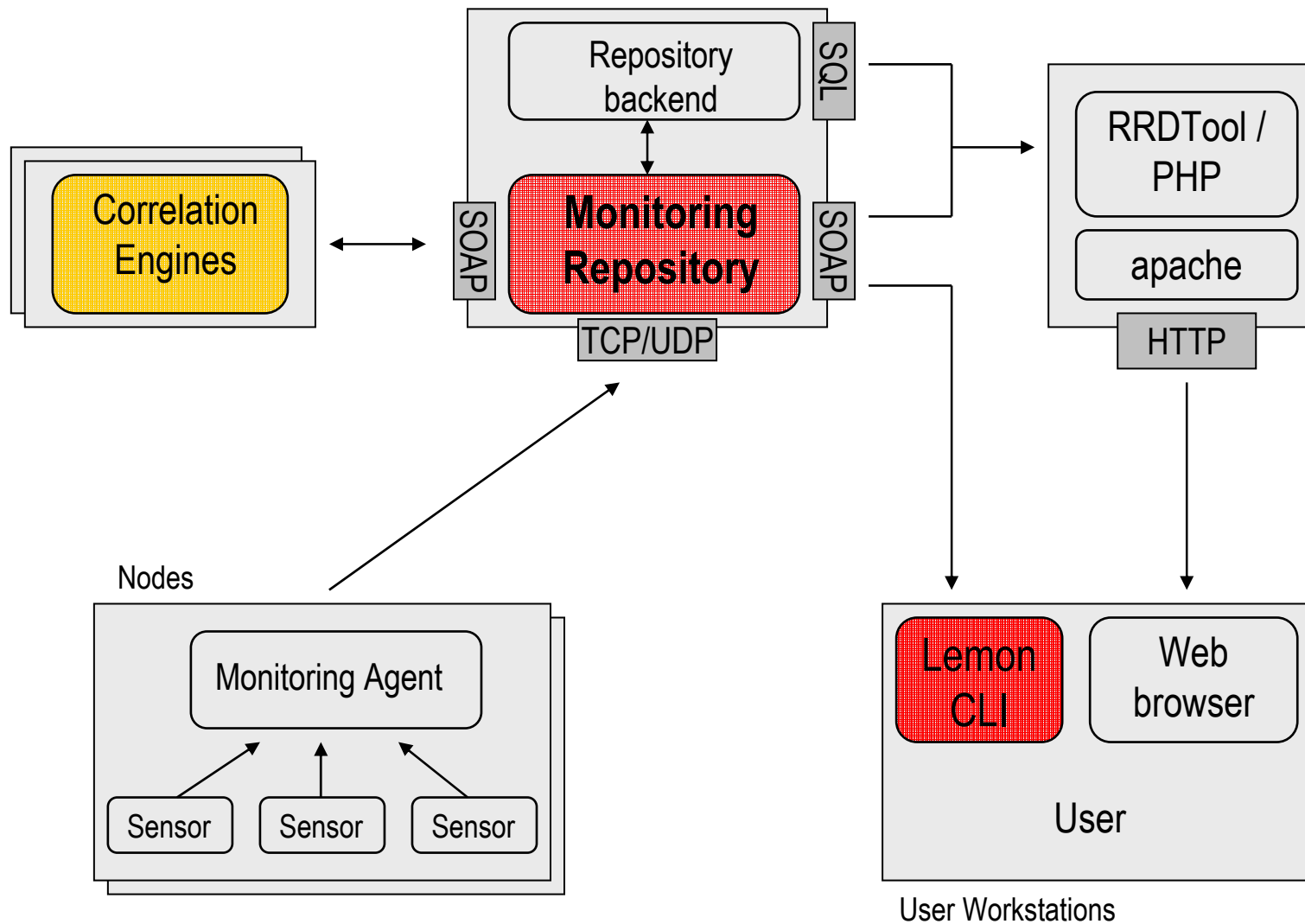- Proxy architecture for enhanced scalability …

# Quattor @ LCG/EGEE

- EGEE and LCG have chosen quattor for managing their integration testbeds

- Community effort to use quattor for fully automated LCG-2 configuration for all services
  - Aim is to provide a complete porting of LCFG configuration components
  - Most service configurations (WN, CE, UI, ..) already available
  - Minimal intrusiveness into site specific environments

- More and more sites (IN2P3, NIKHEF, UAM Madrid..) and projects (GridPP) discussing or adopting quattor as basic fabric management framework...

- ... leading to improved core software robustness and completeness
  - Identified and removed site dependencies and assumptions
  - Documentation, installation guides, bug tracking, release cycles

# Lemon – LHC Era MONitoring

# Deployment and Enhancements

- Smooth production running of Monitoring Agent and Oracle-based repository at CERN-CC
  - 150 metrics sampled every 30s -> 1d; ~ 1 GB of data / day on ~ 1800 nodes
  - No aging-out of data but archiving on MSS (CASTOR)

- Usage outside CERN-CC, collaborations
  - GridICE, CMS-Online (DAQ nodes)
  - BARC India (collaboration on QoS)
  - Interface with MonaLisa being discussed

- Hardened and enhanced EDG software
  - Rich sensor set (from general to service specific eg. IPMI/SMART for disk/tape..)

- Re-engineered Correlation and Fault Recovery
  - PERL-plugin based correlations engine for derived metrics (eg. average of LXPLUS users, load average & total active LXBATCH nodes)
  - Light-weight local self-healing module (eg. /tmp cleanup, restart daemons)

- Developing redundancy layer for Repository (Oracle Streams)

- Status and performance visualization pages …

# LEAF - LHC Era Automated Fabric

- ◆ LEAF is a collection of workflows for *high level* node hardware and state management, on top of Quattor and LEMON:

- ◆ HMS (Hardware Management System):

  - Track systems through all *physical* steps in lifecycle eg. installation, moves, vendor calls, retirement

  - Automatically requests installs, retires etc. to technicians

  - GUI to locate equipment physically

  - HMS implementation is CERN specific, but concepts and design should be generic

- ◆ SMS (State Management System):

  - Automated handling (and tracking of) high-level configuration steps
    - Reconfigure and reboot all LXPLUS nodes for new kernel and/or physical move
    - Drain and reconfig nodes for diagnosis / repair operations

  - Issues all necessary (re)configuration commands via Quattor

  - extensible framework – plug-ins for site-specific operations possible

# LEAF Deployment

◆ HMS in full production for all nodes in CC

  ▪ HMS heavily used during CC node migration (~ 1500 nodes)

◆ SMS in production for all quattor managed nodes


◆ Next steps:

  ▪ More automation, and handling of other HW types for HMS

  ▪ More service specific SMS clients (eg. tape & disk servers)

◆ Developing 'asset management' GUI

  ▪ Multiple select, drag&drop nodes to automatically initiate HMS moves and SMS operations

  ▪ Interface to LEMON GUI

# Summary

◆ ELFms is deployed in production at CERN

- Stabilized results from 3-year developments within EDG and LCG
- Established technology - from Prototype to Production
- Consistent full-lifecycle management and high automation level
- Providing real added-on value for day-to-day operations

◆ Quattor and LEMON are generic software

- Other projects and sites getting involved

◆ Site-specific workflows and "glue scripts" can be put on top for smooth integration with existing fabric environments

- LEAF HMS and SMS

◆ More information:
http://cern.ch/elfms

WAN

Configuration

Installation

Monitoring

Linux
Tape Server

Linux
Service
Nodes

Mass
Storage

LAN

Batch
Scheduler

Linux
Disk Server

Linux
CPU Server

Shared File
System

Fault Tolerance

Electricity

Space

Cooling

Purchasing

Hardware logistic
Physical installation

Node repair
and
replacement

# CPU Server

- The lifetime of this equipment is now from experience about 3 years
  → keep the equipment in production as long as 'useful' (stability, size of memory and local disk

- The cost contribution of processors to a node is only 30 %

- We still have price penalties for 1U, 2U and blade servers (between 10% and 100%)

- The technology trend moves away from GHz to multi-core processors to cope with the increasing power-envelope, this has major consequences for the needed memory size on a node because the memory requirement per job of the experiments is rising (towards 2 GB
  analysts see problem with re-programming applications (multithreading)
  → one main processor with multiple special cores for video, audio processing

- The power consumption worries have not yet been solved (lot's of announcements but e.g. the new Prescott runs at up to 130 W)
  Pentium M is a factor 2 more expensive per SI2000 unit

| Year | 2004 | 2005 | 2006 | 2007 | 2008 | 2009 | 2010 |
|------|------|------|------|------|------|------|------|
| CHF/SI2000 | 1.89 | 1.25 | 0.84 | 0.55 | 0.37 | 0.24 | 0.18 |

(more details here)

# CPU server expansion 2005

- **400 new nodes (dual 2.8 GHz, 2 GB memory) currently being installed**
- **will have than about 2000 nodes installed**
- **acceptance problem, too frequent crashes in test suites problem identified : RH 7.3 + access to memory > 1 GB**
- **outlook for next year: just node replacements, no bulk capacity upgrade**

**CPU resource requests for 2005**

# CPU server efficiency

- Will start at the end of November an activity in the area application performance

- Representatives from the 4 experiments and IT

- To evaluate the effects on the performance of
  1. different architectures (INTEL, AMD, PowerPC)
  2. different compilers (gcc, INTEL, IBM)
  3. compiler options

- Total Cost of Ownership in mind

- Influence on purchasing and farm architecture

# Disk server

**Today:**
**~400 disk server with 6000 disks and 450 TB of disk space (mirrored) installed**

**Issues :**

the lifetime of this equipment is now from experience 3 years
the MTBF figures in production are much lower than advertised
(usage pattern…)

the cost trends for the space are promising, faster than expected

while size and sequential speed are improving, is the random access
performance not changing (worry for analysis, but also multiple stream
productions)

| Year | 2004 | 2005 | 2006 | 2007 | 2008 | 2009 | 2010 |
|---|---|---|---|---|---|---|---|
| Disk Size [GB] | 200 | 330 | 540 | 900 | 1500 | 2400 | 4000 |

| Year | 2004 | 2005 | 2006 | 2007 | 2008 | 2009 | 2010 |
|---|---|---|---|---|---|---|---|
| CHF/GByte | 8.94 | 5.59 | 3.49 | 2.18 | 1.36 | 0.85 | 0.53 |

(more details here)

disk size is for the best price/performance units
i.e. today one can buy 400 GB disks, but the optimum isin the area 200-250 GB

# Disk server

problem in the first half of 2004
disk server replacement procedure  for 64 nodes took place
(bad bunch of disks, cable and cage problems)
reduced considerably the error rate

currently 150 TB being installed

we will try to buy ~ 500 TB disk space in 2005
- need more experience with much more disk space
- tuning of the new Castor system
- getting the load off the tape system
- test the new purchasing procedures

# GRID acccess

- 6 node DNS load balanced GridFTP service, coupled to the CASTOR disk pools of the experiments

- 80 % of the nodes in Lxbatch have the Grid software installed (using Quattor) (limits come from the available local disk space)

- Tedious IP renumbering during the year of nearly all nodes, to cope with the requirement for 'outgoing' connectivity from the current GRID software. Heavy involvement of the network and sysadmin teams.

- A set of Lxgate nodes dedicated to an experiment for central control, bookkeeping, 'proxy'

- Close and very good collaboration between the fabric teams and the Grid deployment teams

# Tape servers, drives, robots

**Today :**

**10 STK Silos with a total capacity of ~ 10 PB (200 GB cassettes)**

**50 9940B drives   fibre channel connected to Linux PCs on GE**

**Reaching 50000 tape mounts per week, close to the limit of the internal robot arm speed**

**will get before the end of the year IBM robot with 8 * 3592 drives and 8 STK LTO-2 drives for extensive tests**

**Boundary conditions for the choices of the next tape system for LHC running:**

- **Only three choices (linear technology) :   IBM,  STK, LTO Consortium**

- **The technology changes about every 5 years, with 2 generations within the 5 years (double density and double speed for version B, same cartridge)**

- **The expected lifetime of a drive type is about 4 years, thus copying of data has to start at the beginning of the 4th year**

- **IBM and STK or not supporting each others drives in their silos**

# Tape servers, drives, robots

- **Drives should have about a year establishment in the market**

- **Would like to have the new system in full production in the middle of 2007, thus purchase and delivery by mid 2006**

- **We have already 10 powderhorn STK silos, which will not host IBM or LTO drives**

- **LTO-2 and IBm 3592 drives are now about one year on the market**

- **LTO-3 and IBM 3592B by the end of 2005/beginning 2006
STK new drive available by ~ mid 2005**

- **Today estimated costs (certainly 20% error on the numbers)**
  **→ bare tape media costs**
      **IBM ~ 0.8 CHF/GB,  STK ~ 0.6 CHF/GB, LTO2 ~ 0.4 CHF/GB**
  **→ drive costs IBM ~ 24 KCHF,  STK ~ 37 KCHF,  LTO2 ~ 15KCHF**

- **High speed drives (> 100 MB/s) need more effort on the network/disk server/file system setup to ensure high efficiency**

**large over-constraint 'phase-space together with the performance/access pattern requirements  →  focus work in 2005**

# Tape storage



Analyzing the Lxbatch
inefficiency trends,
wait time due to tape queues

**stager hits #files limit**

# Mass storage performance

**First set of parameters defining the access performance of an application to the mass Storage system**

number of running batch jobs
internal organization of jobs ( exp.)
(e.g. just request file before usage)
priority policies (between exp. and within exp.)
CASTOR scheduling implementation

speed of the robot
distribution of tapes in silos
(at the time of writing the data)

CASTOR
database performance

tape drive speed
tape drive efficiency

disk server
filesystem
OS + driver

CASTOR load balancing mechanism
monitoring
Fault Tolerance
disk server optimization

data layout on disk
exp. policy
access patterns (exp.)
performance overall
file size

bugs and features

# Example : File sizes

Average file size on disk

| | |
|---|---|
| ATLAS | 43 MB |
| ALICE | 27 MB |
| CMS | 67 MB |
| LHCb | 130 MB |
| COMPASS | 496 MB |
| NA48 | 93 MB |

large amounts < 10MB

Fig 4. Size distribution of files currently STAGED



Total contents of above histogram is : 146847.00

File sizes: minimum: 0.0, maximum: 2000.0 and average: 71.3 (MB)

# Analytical calculation of tape drive efficiencies



Tape Drive Efficiency (30 MB/s max performance)

3 files per mount
10 files per mount
50 files per mount
100 files per mount
1000 files per mount

efficiency [%]

file size [MB]

average # files per mount ~ 1.3
large # of batch jobs requesting
files, one-by-one

tape mount time  ~ 120 s
file overhead      ~  4.4 s

# Tape storage

combination of problems        example : small files + randomness of access


possible solutions :

- concatenation of files on application or MSS level
- extra layer of disk cache,  Vendor or 'home-made'
- hierarchy of fast and slow access tape drives
- very large amounts of disk space
- ……….


Currently quite some effort is put into the analysis of all the available monitoring information to understand much better the influence of the different parameters on the overall performance.

the goal is to be able to calculate the cost of data transfers from tape to the application

→ CHF per MB/s  for volume of X TB

WAN

Configuration

Installation

Linux
Tape Server

Linux
Service
Nodes

Mass
Storage

LAN

Batch
Scheduler

Monitoring

Linux
Disk Server

Linux
CPU Server

Shared File
System

Fault Tolerance

Electricity          Space          Cooling

Purchasing

Hardware logistic
Physical installation

Node repair
and
replacement

# Network    LAN

- **2 * 10 GE switches were integrated in the CERN network backbone in June/July**

- **Two  generations of 10 GB high end routers from Enterasys and 3 switches (24/48 GE ports + 2*10 GE ports) from different vendors are on test in the high throughput cluster**

- **Market survey for the high end routers and the switches for the distribution layer (10 GE to multiple 1 GE) is currently finishing. Tenders will be out in Jan/Feb 2005**

- **First part of the new backbone deployment in mid 2005**

**Tomorrow's  schematic network topology**

WAN      10 Gigabit Ethernet
10000 Mbit/s

Backbone

Multiple 10 Gigabit Ethernet
200 * 10000 Mbit/s

10 Gigabit Ethernet
10000 Mbit/s

Gigabit Ethernet
1000 Mbit/s

**CPU Server**

**Disk Server**

**Tape Server**

# Network    WAN

- **Service data challenges  have started**
  **data transfers between CERN and Tier 1 centers**
  **setting up the routing between the sides is not trivial**
  **and takes some time**

- **10 Itanium nodes dedicated as GridFTP server**
  **Local disks, SRM interface, CASTOR**

- **Tests already with FNAL, BNL, NIKHEF, FZK**
  **e.g.  250 MB/s for days , FNAL pulling data via**
  **GridFTP from local disks**

# High Througput Prototype (openlab + LCG prototype)

**4 * GE connections to the backbone**

**10GE WAN connection**

**12 Tape Server STK 9940B**

**24 Disk Server (P4, SATA disks, ~ 2TB disk space each)**

**4 *ENTERASYS N7 10 GE Switches 2 * Enterasys X-Series**

**2 * 50 Itanium 2 (dual 1.3/1.5 GHz, 2 GB mem)**

**36 Disk Server (dual P4, IDE disks, ~ 1TB disk space each)**

**10GE**

**10 GE per node**

**80 * IA32 CPU Server (dual 2.4 GHz P4, 1 GB mem.)**

**10 GE per node**

**1 GE per node**

**10GE**

**10GE**

**28 TB , IBM StorageTank**

**40 * IA32 CPU Server (dual 2.4 GHz P4, 1 GB mem.)**

**80 IA32 CPU Server (dual 2.8 GHz P4, 2 GB mem.)**

**12 Tape Server STK 9940B**

# Planned data challenges

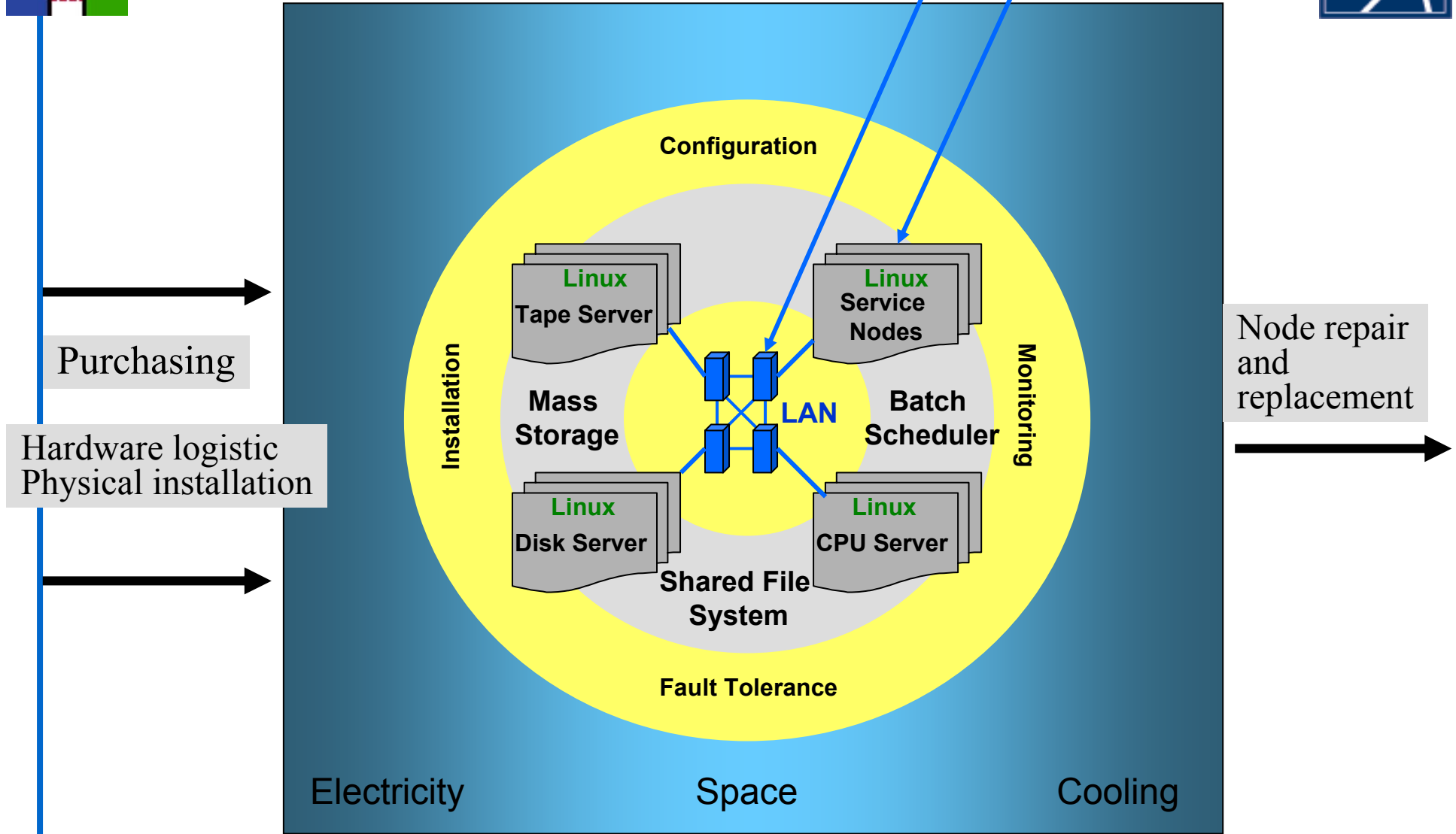- **Dec04 - Service Challenge 1 complete**
  **mass store-mass store, CERN+3 sites, 500 MB/sec between sites, 2 weeks sustained**

- **Mar05 - Service Challenge 2 complete**
  **reliable file transfer service, mass store-mass store, CERN+5 sites, 500 MB/sec between  sites, 1 month sustained**

- **Jul05 - Service Challenge 3 complete**
  **mock acquisition -  reconstruction - recording – distribution, CERN + 5 sites, 300 MB/sec., sustained 1 month**

- **Nov05 – ATLAS or CMS Tier-0/1 50% storage & distribution challenge complete 300 MB/sec, 5 Tier-1s (This is the experiment validation of Service Challenge 3)**

  Tier-0 data recording at 750 MB/sec
  → ALICE data storage challenge VII completed

- continuous data challenge mode in 2005
- use the high-throughput cluster for continues tests, expand the disk space
- start to use the new network backbone as soon as possible

**WAN**

Configuration

Installation

Linux
Tape Server

Linux
Service
Nodes

Mass
Storage

LAN

Batch
Scheduler

Monitoring

Linux
Disk Server

Linux
CPU Server

Shared File
System

Fault Tolerance

Electricity        Space        Cooling

Purchasing

Hardware logistic
Physical installation

Node repair
and
replacement

# Linux (I)

- **The official end of support for RedHat 7.3 was end of 2003**

- **Negotiations between CERN (HEP) and RedHat from October 2003 until February2004 (glutinous responses from RH) about licenses (RH strategy change in summer 2003)**

- **The 'price-breakthrough' came too late and was not competition with the chosen option : recompile the source code from RH (RH has to provide this due to GPL)**

- **First test versions of this CERN version were available at the end of February 2004**

- **The formal CERN Linux certification process (all experiments, AB, IT,..) started in March**

- **Collaboration with Fermi at Hepix in May 2004 on Scientific Linux (Fermi senior partner, reference repository) based on RedHat Enterprise version 3**

- **Community support for security patches of RH 7.3 deteriorated in Q2 2004 → started to buy patches from Progeny = no free CERN version of RH 7.3**

- **Hepix October 2004 : Scientific Linux is a success, many labs migrating to SL**

- **The SLC3 version is certified in November 2004**

# Linux (II)

**Strategy :**

1. Use Scientific Linux for the bulk installations, Farms and desktops

2. Buy licenses for the RedHat Enterprise version for special nodes (Oracle) ~100

3. Support contract with Redhat for 3rd level problems
   contract is in place since July 2004, ~50 calls opened, mixed experience
   review the status in Jan/Feb whether it is worthwhile the costs

4. We have regular contacts with RH to discuss further license and support issues.

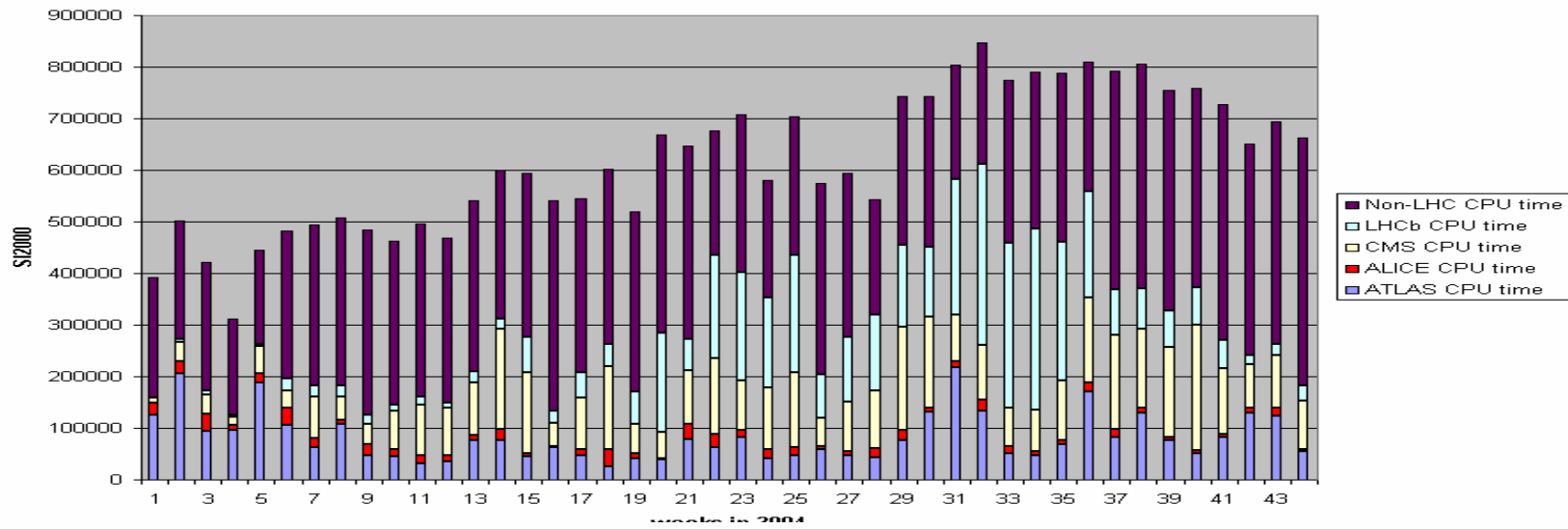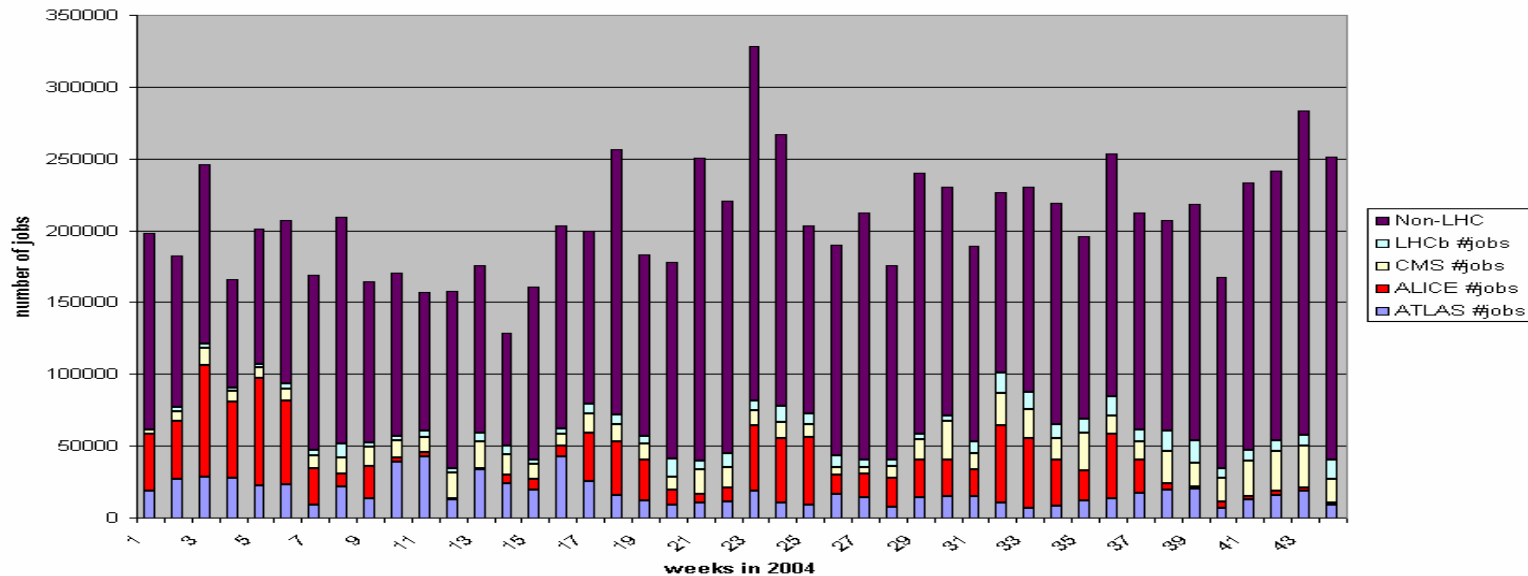The next RH version REL4 is in beta testing and needs some 'attention' during next year

# Batch Scheduler (I)

## CPU resources used in 2004



## Number of Lxbatch jobs per week

# Batch Scheduler (II)

- **Batch scheduler at CERN is LSF from Platform Computing**

- **Very good experience**

- **Complicated , but efficient and flexible fair-share configuration no scalability problems, good redundancy and fault tolerance**

- **Very good support line**

- **Site license and support contract, very cost effective**

- **Limited evaluation of other systems, experience from other sites in the last workshop (PBS, TORQUE+MAUI, Condor) evaluation of other systems was low priority, focus was on automation**

- **No argument for a change, will most likely stay with LSF for the next years**

- **Report in May, next Hepix includes a workshop on batch scheduler experience**

# File systems

Today's systems are AFS and CASTOR
(AFS : 27 server, 12 TB, 113 million files, availability ~ 99.86%
~ 40 MB/s I/O during day time, 660 million transactions/day)

Looked for global shared file system solutions :

- Tested and evaluated several possible file systems (together with CASPUR)
  (Storage Tank, Lustre, GFS, cXFS, StoreNext, ….)
  stability, fault tolerance, error recovery, scalability, SAN versus NAS, exporter,….

- Report at the end of the year

- No candidate for AFS replacement during the next 2-3 years

- Continue testing with Caspur (if interesting developments) from time to time

- Small investment into improving (performance, monitoring, scalability) of
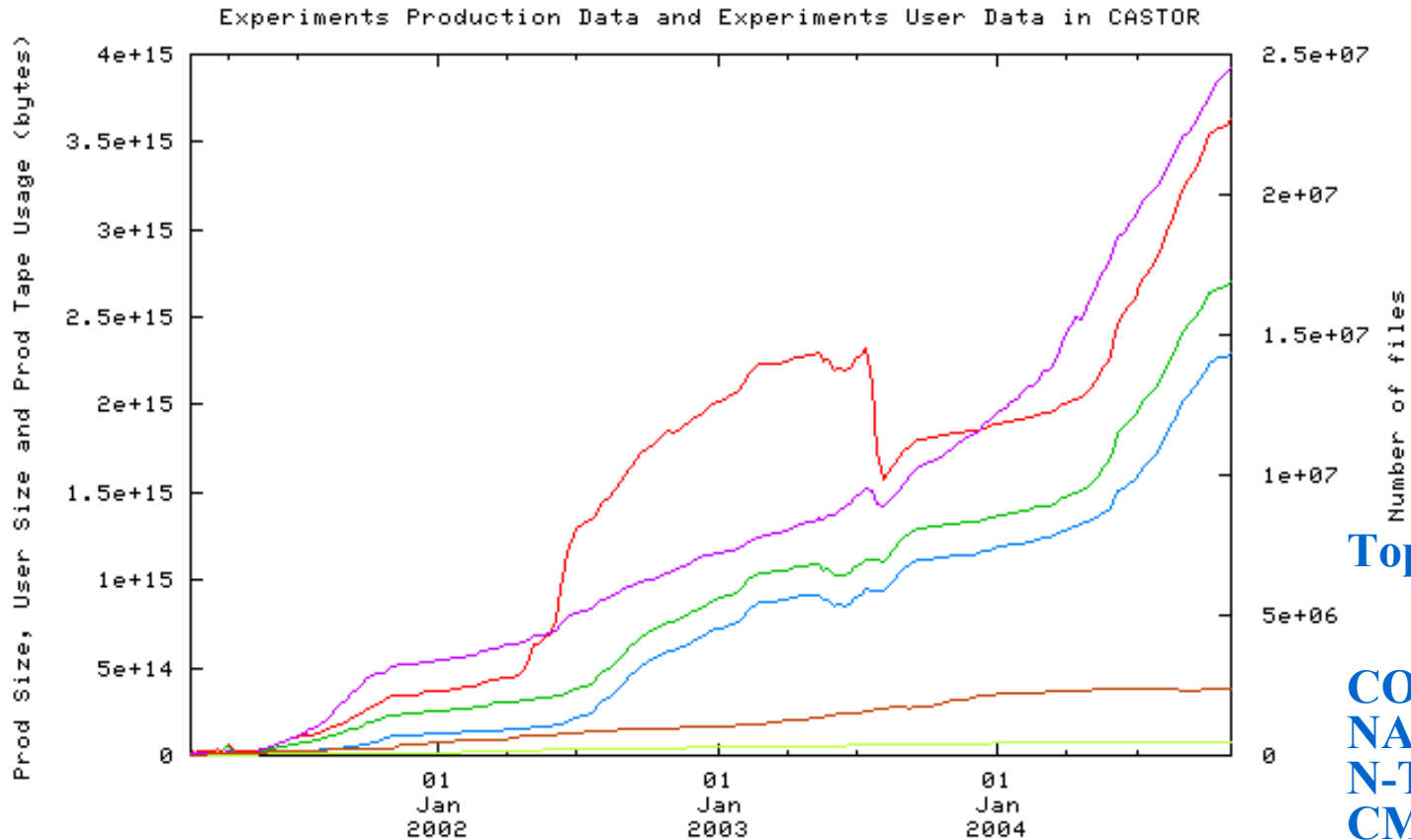  openAFS (collaboration with CASPUR and probably GSI)

# CASTOR status

- **Usage at CERN**
  - **~3.4 PB data**
  - **~26 million files**
  - **CDR running at up to 180 MB/s aggregate**

- **Operation**
  - **Repack in production (since 2003): >1PB of data repacked**
  - **Tape segments checksum calculation and verification is in production since March 2004**
  - **Sysreq/TMS definitely gone in July**
  - **VDQM prioritize tape write over read → no drive dedication for CDR needed since September**
  - **During 2004 some experiments hit stager catalogue limitation (~200k files) beyond which the stager response can be very slow**

- **Support at CERN**
  - **2nd and 3rd level separation works fine**
  - **4 FTE developer and 3 FTE operations**
  - **Increasing support for SRM and gridftp users**

- **Other sites**
  - **PIC and IHEP contribute to CASTOR development at CERN → liberate efforts for better CASTOR operational support to other sites**
  - **CNAF will soon contribute**
  - **RAL planning to evaluate CASTOR**

# CASTOR@CERN evolution



Experiments Production Data and Experiments User Data in CASTOR

Generated Nov 01, 2004 CASTOR (c) CERN/IT/ADC/CA

- TOTAL Prod Size
- TOTAL Prod Tape Usage
- TOTAL Prod Tape Usage Without Redwood Copy
- TOTAL Prod Nb Files
- TOTAL User Size
- TOTAL User Nb Files

**3.4 PB data
26 million files**

**Top 10 experiments**

|          | [TB] |
|----------|------|
| COMPASS  | 1066 |
| NA48     | 888  |
| N-Tof    | 242  |
| CMS      | 195  |
| LHCb     | 111  |
| NA45     | 89   |
| OPAL     | 85   |
| ATLAS    | 79   |
| HARP     | 53   |
| ALICE    | 47   |
| sum      | 2855 |

# New stager developments delay (I)

Several not foreseen but important extra activities :

The CASTOR development team has also the best knowledge of the internals of the current CASTOR system, thus are often involved in operational aspects as these have higher priority than developments.

1. The limits of the current system are seen now more frequently with the increased usage patterns of the experiments → urgent bug fixes or workarounds
   e.g. large number of small files (limit in the stager + tape technology limits)

2. Tape segments checksum calculation and verification deployed

3. Old service stopped :Sysreq/TMS

4. CDR priority scheme for writing tapes = better efficiency of drive usage

5. Bug fix in the repack procedure
   End November 2003 a bug was found in stager API during the certification of first production release of repack. The effect was that a fraction (~5%) of the repacked files got wrongly mapped in the CASTOR name server.
   Between December 2003 – May 2004 more or less one CASTOR developer working full time on finding and repairing incorrectly mapped CASTOR files
   A bit less than 50,000 files wrongly mapped out of >1 million
   Repair applied to the CASTOR name server the 26th of April 2004

# New stager developments delay (II)

**6.** SRM interoperability
- Drilling down the GSI (non-)interoperability details
- Holes in the SRM specs
- Time-zone difference (FNAL-CERN) does not favor efficient debugging of interoperability problems

**7.** Other grid activities: CASTOR as a disk pool manager without tape archive
- We provided a packaged solution for LCG
- But… support expectations pointed towards a development sidetrack
    - Castor is not well suited for such configurations
- Decided to drop all support for CASTOR disk-only configurations (Jan/Feb 2004) and focus on the CERN T0/T1 requirements
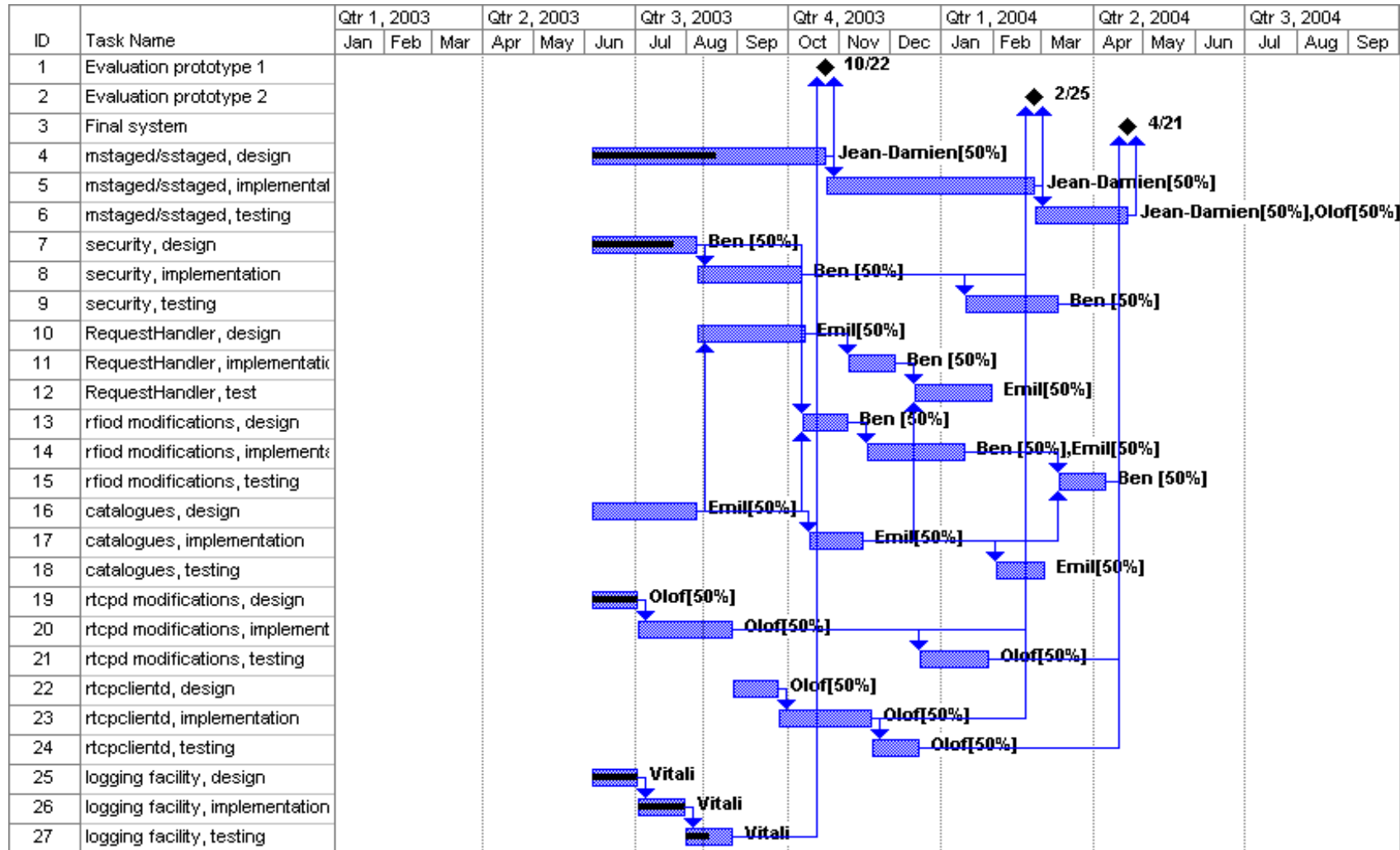
**8.** after the first prototype tests some small redesigns took place

To ease the heavy load on the CASTOR developers we were able to use man-power from our collaboration with PIC (Spain) and IHEP (Russia). These persons had already experience with CASTOR and were able to very quickly pick up some of the development tasks (there was no free time for any training of personal).
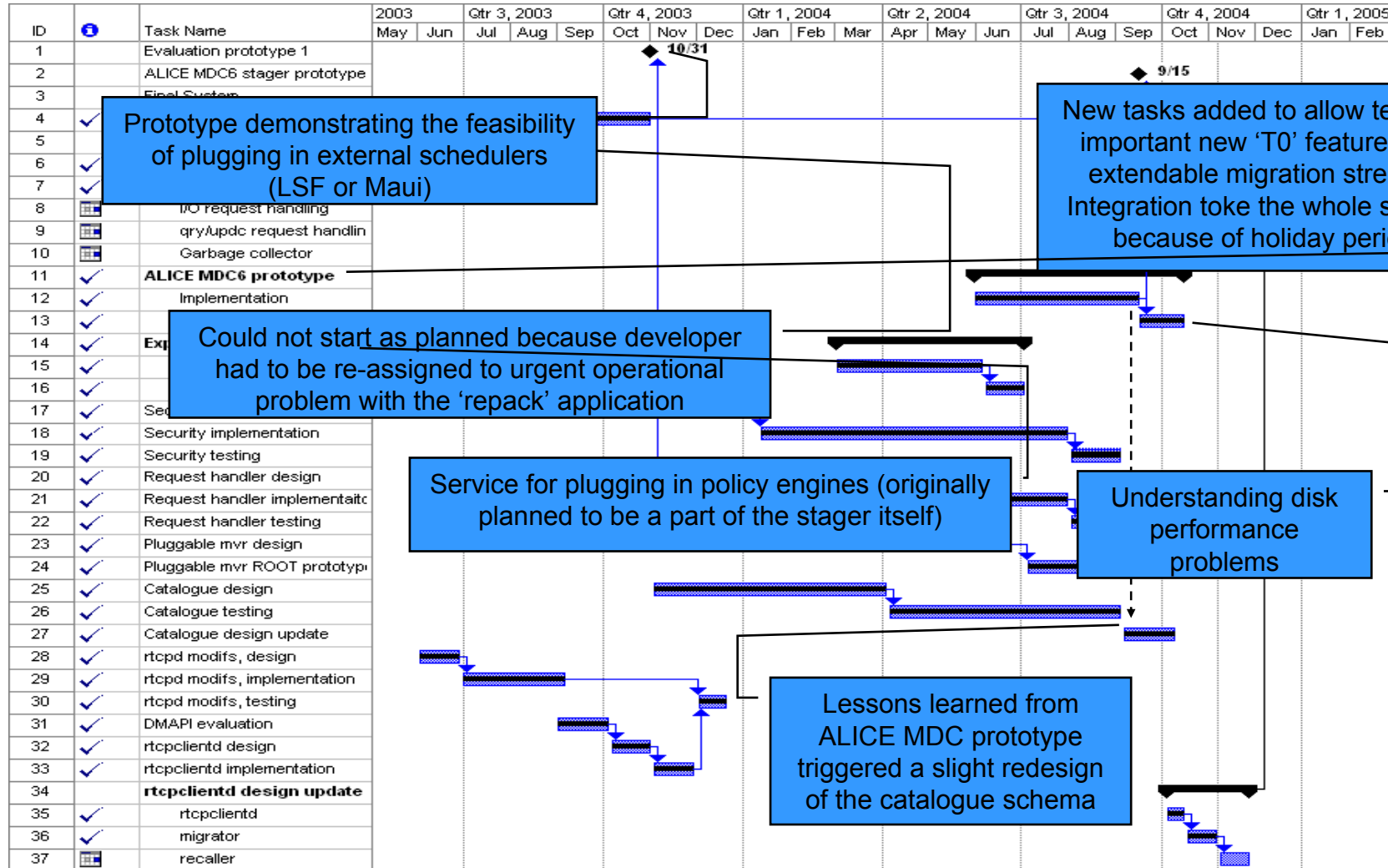
# New stager developments
# ALICE MDC-VI prototype

- ◆ **Because of the delays there was a risk to miss the ALICE MDC-VI milestone**

  - ▪ **New stager design addresses important Tier-0 issues:**

    - ⬥ **Dynamically extensible migration streams**

    - ⬥ **Just-in-time migration candidate selection based on file system load**

    - ⬥ **Scheduling and throttling of incoming streams**

  - ▪ **ALICE MDC-VI the ideal test environment. Could not afford to miss it…**

    - ⬥ **The features were ready but the central framework did not exist**

    - ⬥ **Decided to build a hybrid stager re-using a slimmed-down version of the current stgdaemon as central framework**

# New stager developments
# ALICE MDC-VI prototype



Application

ROOT TCastorFile with new stager API

Today's GC script

= old stager component

3rd party Policy Engine

stager_castor

stgdaemon

Recaller

Migrator

Request Handler

Request repository (Oracle or MySQL)

Resource Management Interface

Tape mover (RTCOPY) client daemon

mvr cntl

rootd (disk mover)

CASTOR tape archive components (VDQM, VMGR, RTCOPY)

LSF

Maui

Disk cache

file system load monitoring

# New stager developments
# Testing ALICE MDC-VI prototype

◆ **The prototype was very useful:**

- **Tuning of file-system selection policies**

- **The designed assignment of migration candidates to migration streams was not efficient enough →redesign of catalogue schema**

  - **Migration candidates initially assigned to all tape streams**

  - **The migration candidate is 'picked up' by the first stream that is ready to process it**

  - **Slow streams (e.g. bad tape or drive) will not block anything**

◆ **Also found that the disk servers used for our tests were not well tuned for competition between incoming and outgoing streams**

 **→ new procedures for the tuning of disk servers developed by the   Linux team**

# New stager developments
## Current status

- **Catalogue schema and state diagrams are ready**
  - **Code automatically generated**
  - **Only ORACLE supported for the moment**
  - **http://cern.ch/castor/DOCUMENTATION/STAGE/NEW/Architecture/**

- **The finalization of the remaining components is now running at full speed**
  - **Central request processing framework (the replacement of stgdaemon):**
    - **New stager API defined and published for feedback (http://cern.ch/castor/DOCUMENTATION/CODE/STAGE/NewAPI/index.html )**
    - **I/O (stagein/stageout) and query processors: implementation started. Ready in 3-4 weeks**
  - **Recaller**
    - **Implementation started. Ready 1 – 2 weeks**
  - **Garbage collector**
    - **Implementation not started. Estimated duration ~2 weeks**

- **Hopefully we will be able to replace the ALICE MDC6 prototype by the final system in early December**

- **will also start to test physics production type environment with large stager catalogue (millions of files) and tape recall frequency**

# New stager developments Deployment (cont)

- ◆ **Security issues**
  - ■ **All CASTOR services are technically prepared for strong authentication**
    - ◆ http://cern.ch/castor/DOCUMENTATION/CODE/SECURITY/CASTOR_Security_Implementation.pdf
    - ◆ **Kerberos-4, 5 and GSI supported**
  - ■ **CASTOR security plug-ins used by other projects (LCG, EGEE)**
  - ■ **A number of deployment issues remain:**
    - ◆ **Kerberos-5 infrastructure not yet in place**
    - ◆ **Batch job clients must have appropriate credentials**
    - ◆ **No solution yet for windows clients**
    - ◆ **Management of CASTOR service keys**
  - ■ **Propose to do first deployment without strong authentication and upgrade when all infrastructure issues are solved**
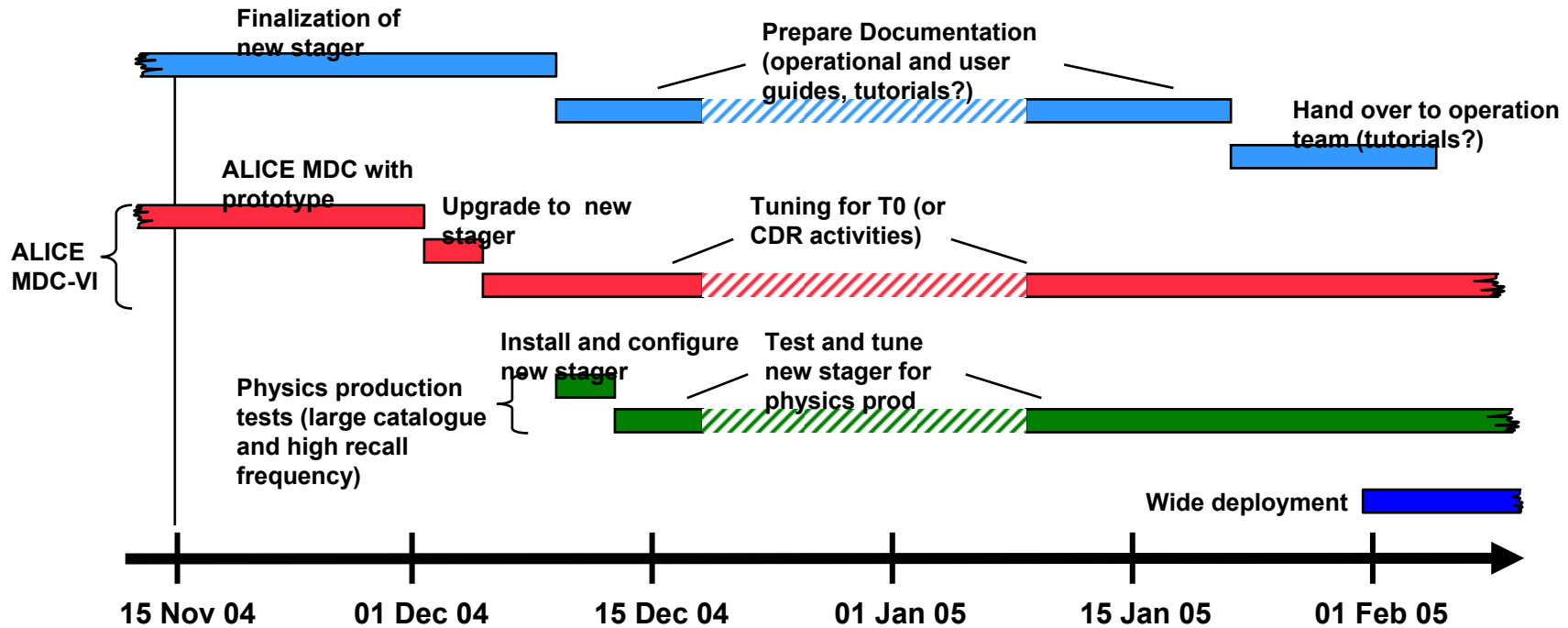
- ◆ **Packaging**
  - ■ **New packaging model envisaged:**
    - ◆ **One RPM for each CASTOR client and server**
      - ▪ **rfio**
      - ▪ **Stage**
      - ▪ **Nameserver**
      - ▪ **VMGR**
      - ▪ **…**
    - ◆ **One RPM for libraries, One 'devel' RPM (include files, man-pages)**

- ◆ **It will be possible to import disk servers from current to the new stager without having to re-stage the files**

**New system is deployed in the high-throughput cluster and heavily tested. One additional person has been added specifically for testing from IT using the ALICE MDC programs. Good performance but yet too many instabilities, debugging phase**

**ALICE MDC (goal 450 MB/s ) will be late, wait for the final Castor version stability**

# CERN T0 center

**first look at the costs of only the T0 part**
**of the CERN center (no analysis, limited reprocessing)**

**some basic assumptions**

- **data needs to be reconstructed in near real time**
- **one CDR processing and one re-processing per year of the raw data**
- **7 days disk buffer (load per disk is critical)**
- **sequential, fully organized, efficient access to tapes**

**CPU+Disk+Tape (from the table)  :   32.7  MCHF**
(share is ALICE:11.3 ATLAS:12.3, CMS:6.3 , LHCb:2.8)
**Plus**

| | |
|---|---|
| **Tape Infrastructure** | **4.5 MCHF** |
| **LAN bandwidth** | **7.4 MCHF** |
| **Sysadmin** | **2.6 MCHF** |
| **WAN** | **6.0  MCHF** |

-----------------------------------------------------------------

**Total T0 cost                                  53.2  MCHF**
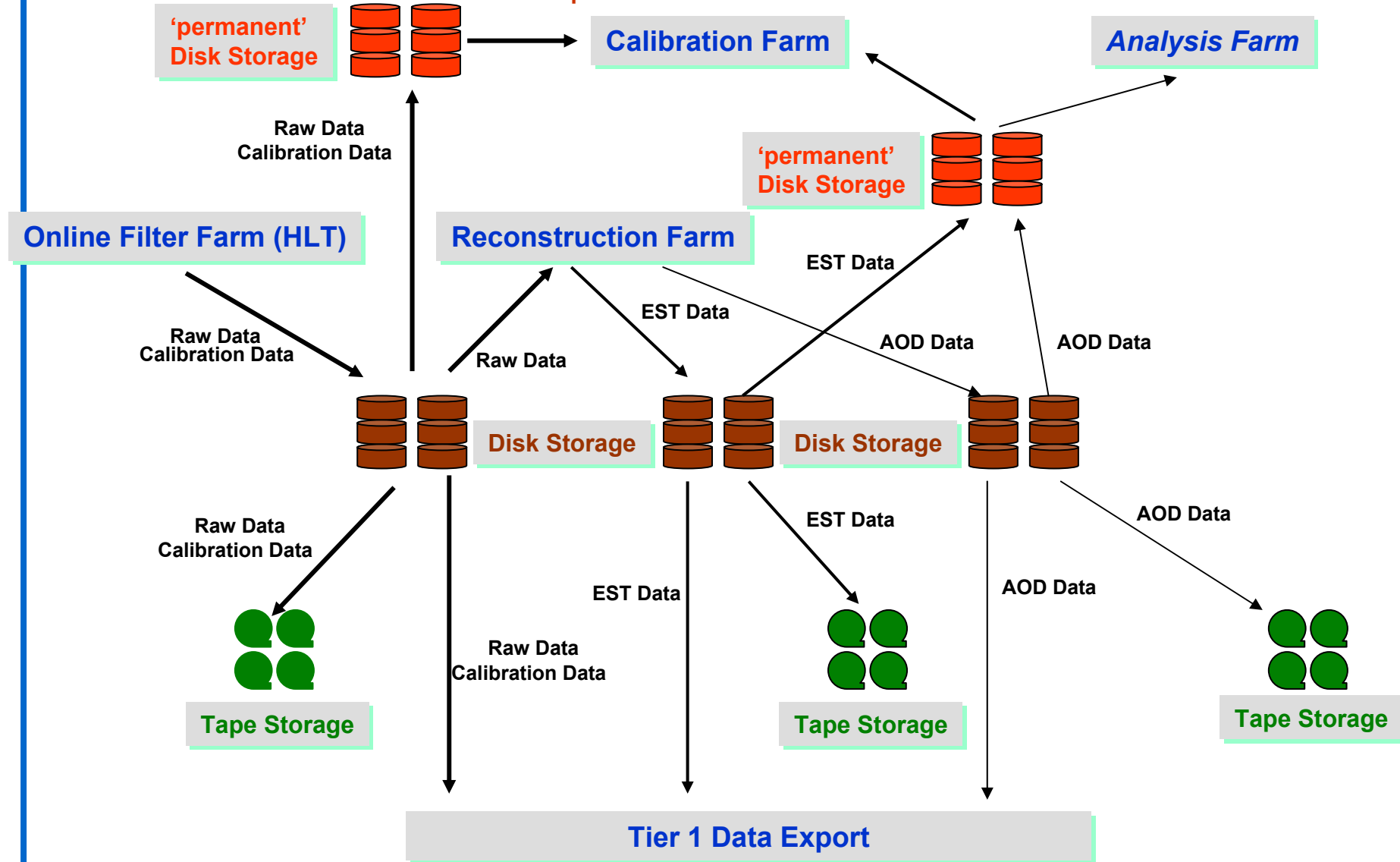
**of course very preliminary numbers and heavily dependent on the computing models**
**just to have an idea where we are….**                        **(more info here)**

# Dataflow local CERN Fabric 2007

Complex organization with high data rates (~10 GBytes/s) and ~100k streams in parallel

**'permanent' Disk Storage** → **Calibration Farm**

**Analysis Farm**

Raw Data
Calibration Data

**'permanent' Disk Storage**

**Online Filter Farm (HLT)**

**Reconstruction Farm**

EST Data

Raw Data
Calibration Data

Raw Data

EST Data

AOD Data

AOD Data

**Disk Storage**

**Disk Storage**

Raw Data
Calibration Data

EST Data

AOD Data

AOD Data

EST Data

Raw Data
Calibration Data

EST Data

AOD Data

**Tape Storage**

**Tape Storage**

**Tape Storage**

**Tier 1 Data Export**

# Complexity

## Hardware components

| | end 2004 | 2008 |
|---|---|---|
| CPU capacity [SI2000] | 2 Million | 20 million |
| Disk space [TB] | 450 | 4000 |
| # CPU server | 2000 | 4000 |
| # disks | 6000 | 8000 |
| # disk server | 400 | 800 |
| # tape drives | 50 | 200? |
| # tape cartridges | 50000 | 50000 |

(these are estimates for 2008, assuming CPU capacity and disk space are continue to grow as in the last 2 years, Moore's Law)

→ today we are less than a factor 2 in hardware complexity away from the system in 2008

# Summary

- **Major activity and success was the automation developments in the farms (ELFms)**

- **Space, cooling, electricity infrastructure on track**

- **no surprises in the CPU, disk server and network area**

- **Delays in the CASTOR area, pre-production system now under heavy tests**

- **Focus on Tape technology developments and market for 2005**

- **Tape system will be under heavy stress in 2005 (data challenges and their preparations)**