



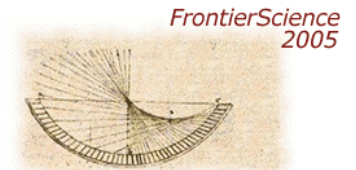
Running the Italian Tier-1 for CMS using Grid tools

D. Bonacorsi

(on behalf of INFN-CNAF Tier-1 staff and the CMS experiment)



New Frontiers in Subnuclear Physics 2005
Fourth International Conference on Frontier Science
September 12th-17th, 2005 - Milano Bicocca, Italy





Use of Grid in CMS



- CMS computing relies on a distributed infrastructure of Grid resources, services and toolkits

- ❑ building blocks provided by Worldwide LHC Computing Grid [WLCG]

- ❖ CMS builds application layers able to interface with few - at most - different Grid flavors (LCG-2, Grid-3, EGEE, NorduGrid, OSG)

→ see [M.Mazzucato, this conf, day 5]

- CMS "C-TDR" delivered (Jun05)

- ❑ CMS Computing Model (CERN/LHCC 2004-035) revised, in preparation for the first year of LHC running (2008)

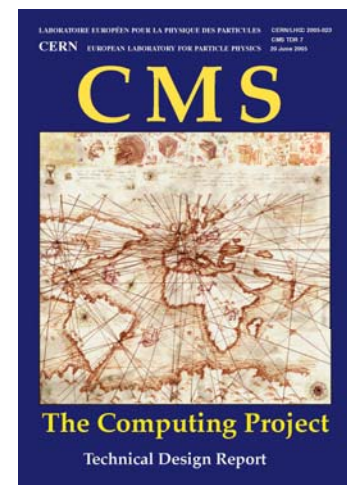
- ❖ not "blueprint", but "baseline" targets (+ development strategies)

- ❑ hierarchy of computing tiers using WLCG tools

- ❖ MONARC + CERN/LHCC 2001-004 + LCG MoU

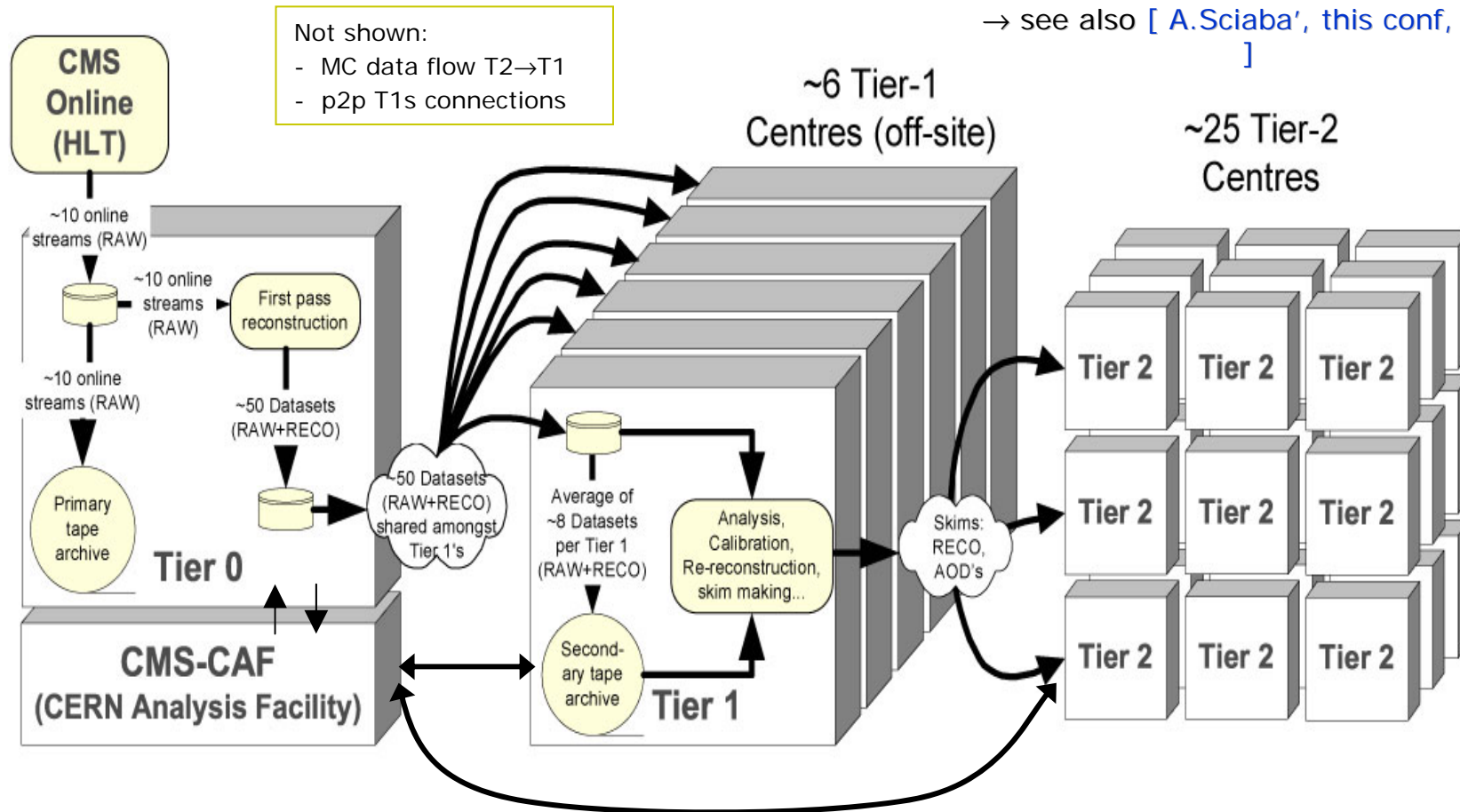
- ❖ focus on tiers role, functionality and responsibility

CERN-LHCC 2005-023





Tiered architecture



This talk: focus on Italian Tier-1, among the set of large CMS T1 sites

- ASCC (Taipei), CCIN2P3 (Lyon), FNAL (Chicago), GridKA (Karlsruhe), **INFN-CNAF (Bologna)**, PIC (Barcelona), RAL (Oxford)





CMS Workload and Data Management

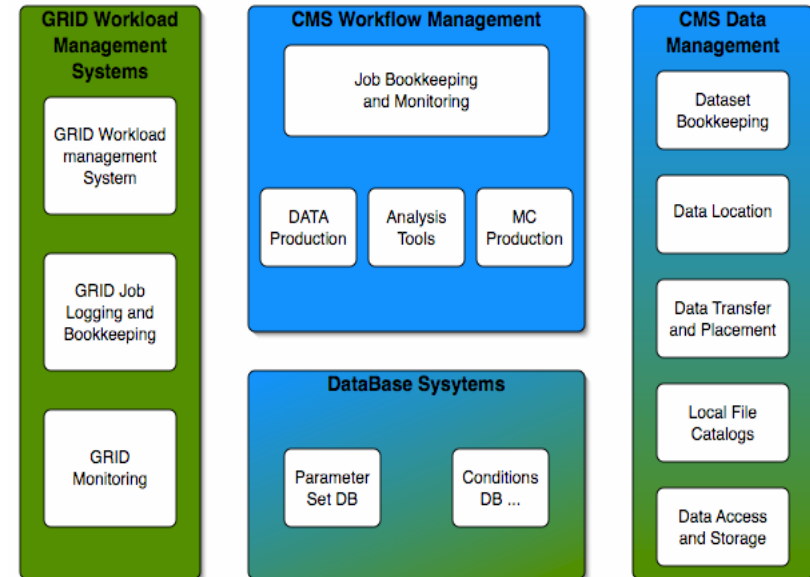
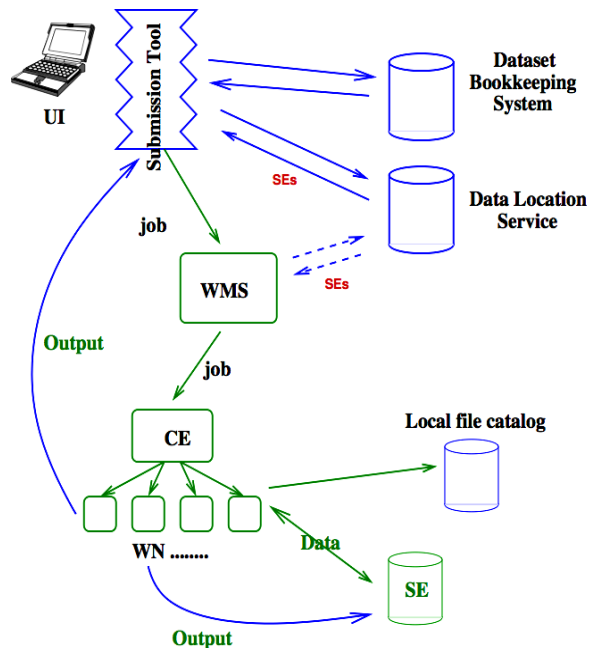


All this: "baseline"

WMS

not a specific flavor..

- ❑ relies on Grid workload management for:
 - ❖ properly configured environment for CMS jobs
 - ❖ good submission rate: O(1000) jobs/few-secs
 - ❖ processing reliability: 24/7, >95% job success rate
 - ❖ policy/priorities between CMS sub-groups
 - ❖ redundant monitoring of job submission/status



DMS

- ❑ Must address needs for a distributed workflow
- ❑ Data book-keeping system: "what data do exist?"
 - ❖ contains references to parameters, lumi, data quality info
- ❑ Data location service: "what is the data location?"
- ❑ Data transfer and placement
 - ❖ relying ('not trusting') on underlying transfer systems
- ❑ Site local services
 - ❖ data storage and access + local file catalogues

FrontierScience 2005





Required functionalities in a CMS T1



- Scheduled data-processing operations:
 - ❑ later-pass reco, AOD extraction, skimming, reprocessing, ...
- Data archiving:
 - ❑ custody of raw+reco and subsequently produced data
- Disk storage management:
 - ❑ fast cache to mass storage archives, buffer for data transfer, ...
- Data distribution:
 - ❑ data import/export from/to any CMS Tier (0/1/2/N)
- Analysis/User support
 - ❑ grant proficient data access (via local, CMS and WLCG services)
- Possible use as a co-located Tier-2
 - ❑ access to local long-term storage and local batch facilities but no interference with T1 commitments towards CMS





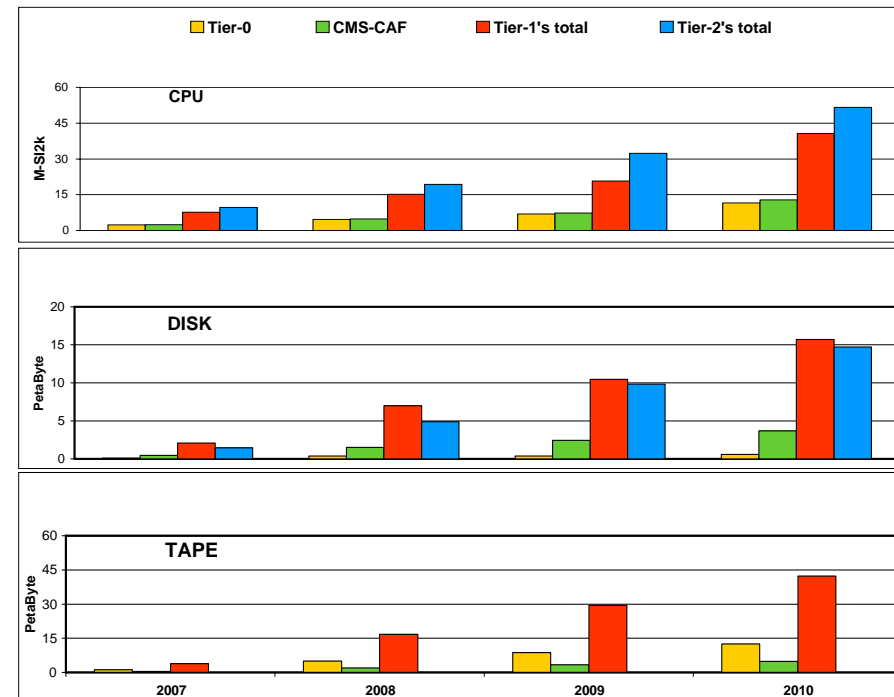
CMS T1 Services & Resources



⇒ such functionalities imply:

➤ Services:

- Disk/MSS management
- WLCG-enabled farming
- Site security
- Accounting
- Database services
- User support



➤ Resource requirements (nominal for an average Tier-1 in 2008):

- ✓ WAN: > 10 Gb/s
- ✓ CPU: 2.5 MSI2k (scheduled reprocessing : analysis = 2 : 1)
- ✓ Disk: 1.2 PB (~85% for analysis data serving)
- ✓ MSS: 2.8 PB (losses ~tens of GB per PB stored)





Outcome of CMS Data Challenge 2004



- Validate the CMS computing model on a sufficient number of Tier-0/1/2's
 - large scale test of the computing/analysis model
- ❑ CMS Pre-Challenge Production (PCP)
- ❑ CMS Data Challenge (DC04)
 - ❖ Reconstruction and analysis on CMS data sustained over 2 months at the 5% of the LHC rate at full lumi

Outcome:

- reconstruction/data-transfer/analysis may run at 25 Hz
- automatic registration and distribution of data, key role of the TMDB
 - ❑ was the embrional PhEDEx (see later)
- support a (reasonable) variety of different data transfer tools and set-up
 - ❑ Tier-1's: different performances, related to operational choices
 - ❖ SRB, LCG Replica Manager and SRM investigated; INFN T1 good performance using LCG-2 tools
- register all data and metadata (POOL) to a world-readable catalogue
 - ❑ RLS: sufficient as a global file catalogue, inadequate as a global metadata catalogue
- analyze the reconstructed data at the Tier-1's as data arrive
 - ❑ INFN T1: LCG-2 components: dedicated bdII+RB; UIs, CEs+WNs at CNAF and PIC
 - ❑ real-time analysis at Tier-2's was demonstrated to be possible
 - ❖ ~15k jobs submitted, time window reco-done/analysis-start: ~20 mins
- optimize file size (i.e. increase $\langle \# \text{events} \rangle / \langle \# \text{files} \rangle$)
 - ❑ more efficient use of bandwidth, reduce overhead of commands, address MSS scalability issues





The INFN-CNAF Tier-1



- Located at INFN-CNAF centre, in Bologna (Italy)
 - ❑ computing facility for INFN HNEP community
 - ❖ one of the main nodes of GARR network

- Multi-experiment Tier-1
 - ❑ 13 expts supported
 - ❖ LHC and others, e.g. AMS, Argo, BaBar, CDF, Magic, Virgo, ...
 - ❑ towards a dynamic share of access to resources to all involved expts

- CNAF is a special Italian site from a Grid perspective
 - ❑ participating to LCG, EGEE, INFN-GRID projects
 - ❑ support to R&D activities
 - ❖ develop Grid prototypes/components, testing Grid interfaces, ...
 - ❑ ... “traditional” access to resources may be granted as well





Tier-1 resources and services



- computing power
 - ❑ mainly on global farm: ~2250 CPU slots available (+ servers)
 - ❖ biproc boxes [320 @0.8-2.4 GHz, 350 @3 GHz], ht activated
 - ❑ nb queues > nb exps
 - ❑ fair share on a monthly time window, hard limit at 1000 jobs/exp

- storage
 - ❑ on-line data access (**disks**)
 - ❖ IDE, SCSI, FC; 4 NAS systems [~60 TB], 2 SAN systems [~225 TB]
 - ❑ custodial task on MSS (**tapes** in Castor HSM system)
 - ❖ Stk L180 lib - overall ~18 TB
 - ❖ Stk 5500 lib - 6 LTO-2 [~ 240 TB] + 2 9940b [~ 136 TB] (more to be installed)

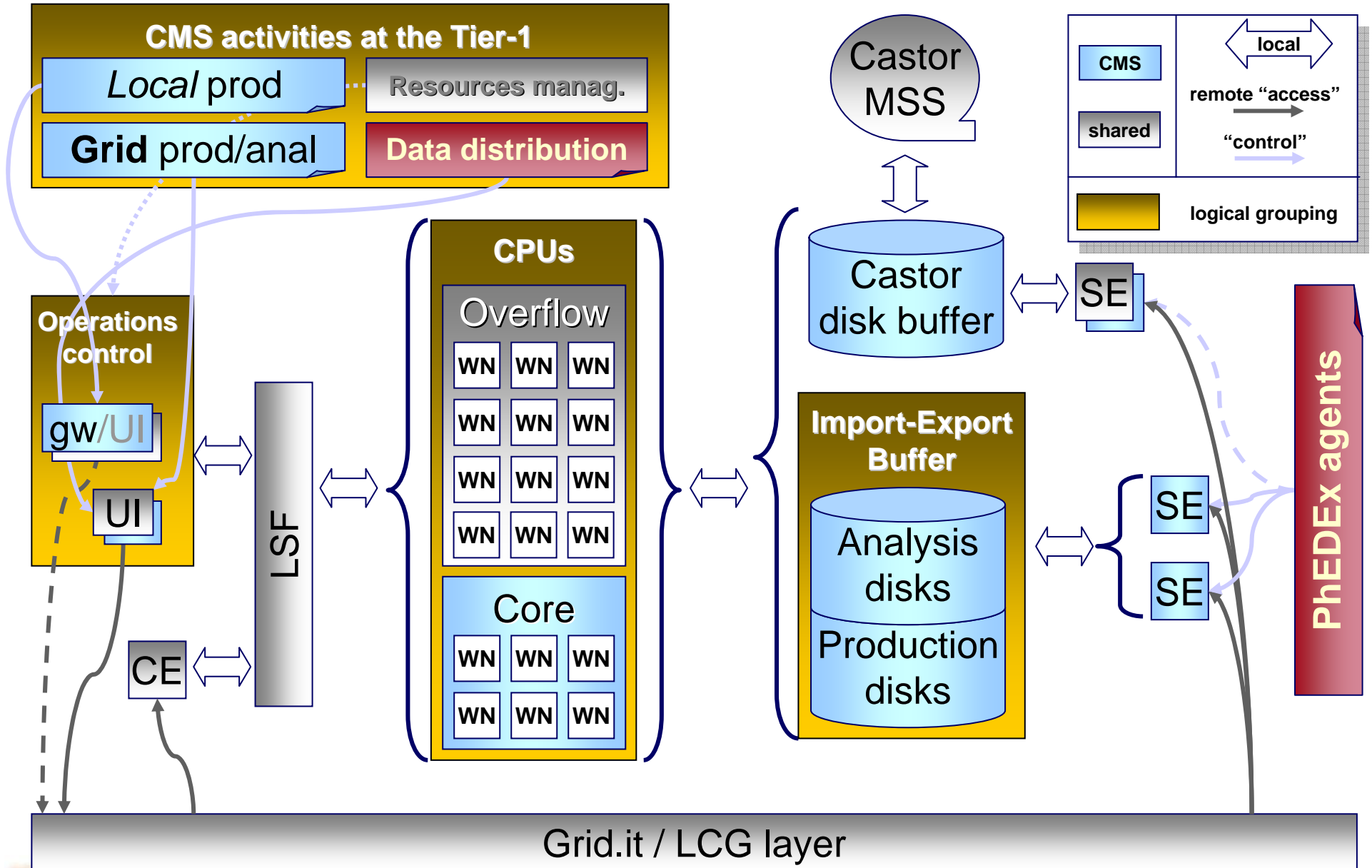
- networking
 - ❑ T1 LAN
 - ❖ rack FE switches with 2xGbps uplinks to core switch (ds → via GE to core)
 - ❖ upgrade foreseen → rack Gb switches [Q4 2005]
 - ❑ 1 Gbps T1 link to WAN (+ dedicated links for LCG Service Challenges)
 - ❖ upgrade foreseen to 10 Gbps [Q3 2005]

- more:
 - ❑ *low-level*: infrastructure, sys-admin, db services, ...
 - ❑ *high-level*: support to exp-specific activities, coordination with tiers, ...





Current CMS set-up at the Tier-1





The PhEDEx project



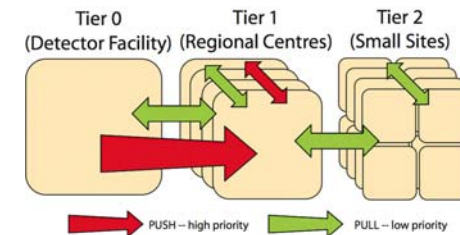
<http://cern.ch/cms-project-phedex>



→ see poster at this conf

Physics Experiment Data Export

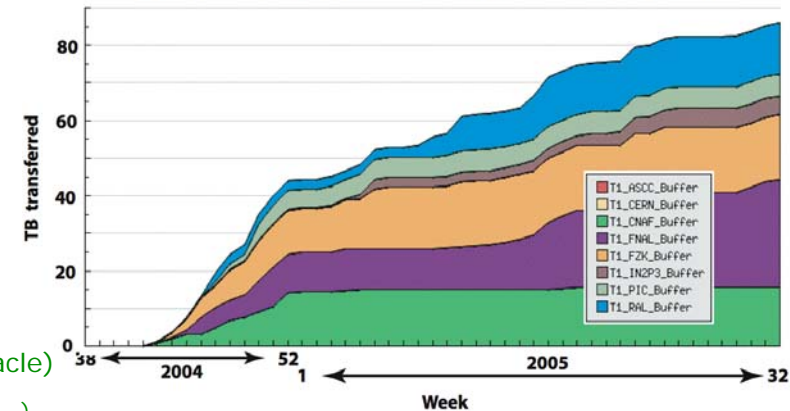
large-scale, reliable, scalable dataset replication



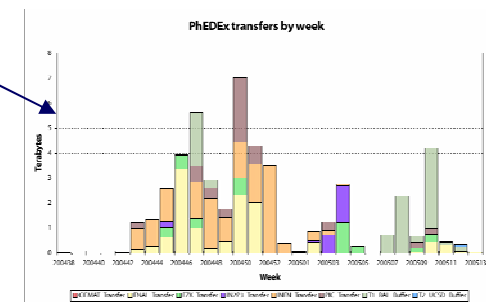
currently in use by CMS as official data distribution management system

- Primary concerns and achievements:
 - ❑ address HEP use-cases via a push-pull negotiation
 - ❑ replication, data safety, tape migration/stage
 - ❑ layered architecture and agnosticism
 - ❑ reliability and replication robustness
- Deployment
 - ❑ in production use for over a year for CMS
 - ❑ Transfer Management DB (TMDB) as blackboard (Oracle)
 - ❑ topology with T0, 7 T1, 16 smaller site (INFN: T1 + 4 T2 + ...)
 - ❑ handles ~110 TB, replicated on average twice
 - ❑ SRM provides a std interface to storage systems
- Replication performances
 - ❑ Day-to-day production service + high-throughput SC3 test transfers
 - ❖ exhibits sustained rates of 1 TB/day/T1 (100 TB in 6 days in SC3)

PhEDEx Tier 1 production transfers by week



Important role of INFN Tier-1 (see orange color)



Interfacing with LCG LFC under evaluation, and with gLite FTS in progress





Storage issues for CMS at INFN T1



→ see P.P.Ricci and D.Bonacorsi contributions at ACAT05

disks

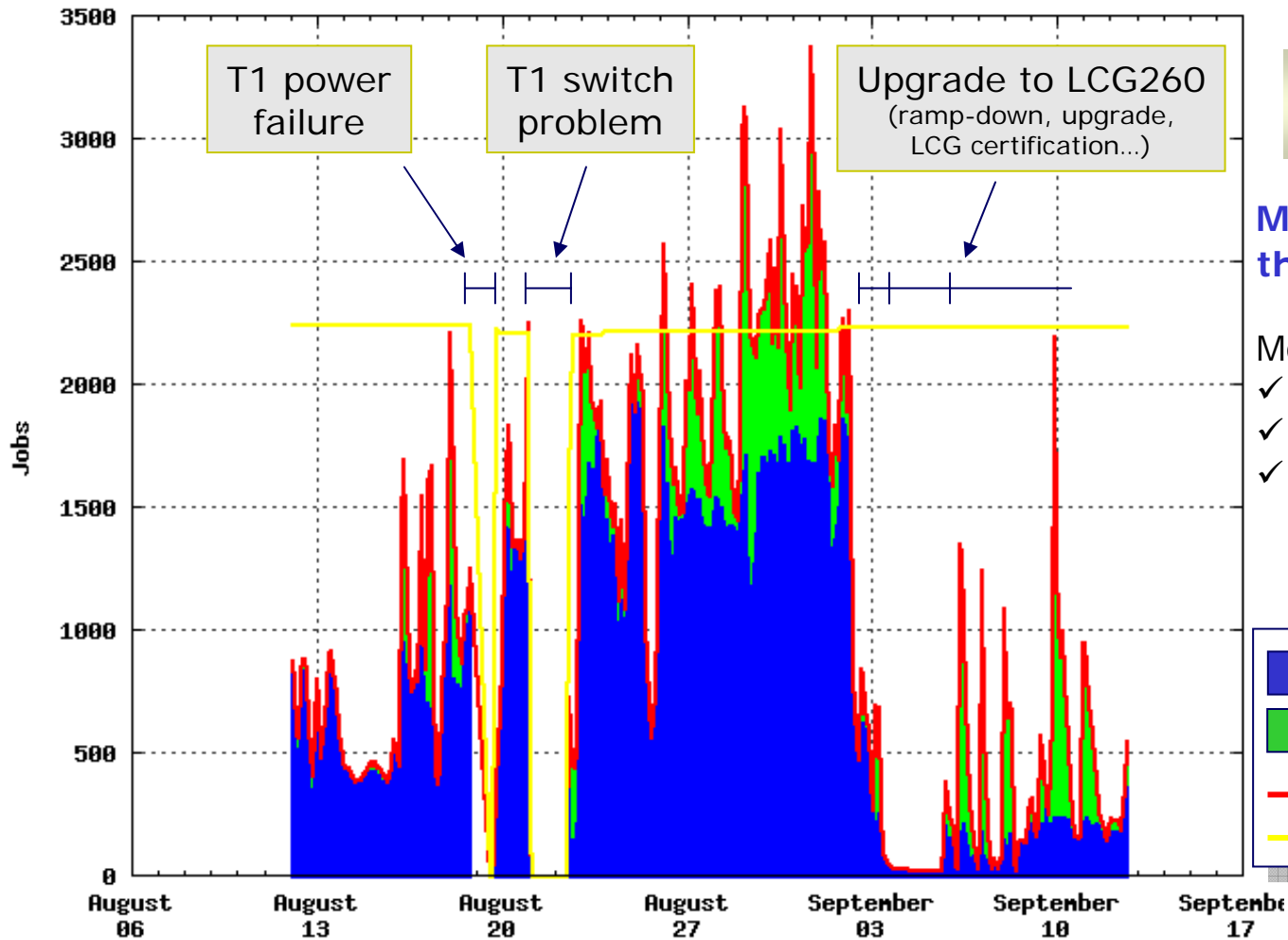
- driven by requirements of LHC data processing at the Tier-1
 - ❑ i.e. simultaneous access to ~PBs of data from ~1000 nodes at high rate
- main focus is on robust, load-balanced, redundant solutions to grant proficient and stable data access to distributed users
 - ❑ a *SAN approach with a parallel filesystem on-top* looks promising

tapes

- CMS DC04 helped to focus some problems:
 - ❑ LTO-2 drives not efficiently used by expts in production at T1
 - ❖ performance degradation increases as file size decreases
 - ❖ hangs on locate/fskip after ~100 not-sequential reading
 - ❖ not-full tapes are labelled 'RDONLY' after 50-100 GB written only
 - ❑ CASTOR performances increase with sequential pre-stage reading
- solutions?
 - ❑ CASTOR-2? 9940b drives? + file-merging, import buffer, ...

- a DC is exp-specific, its conclusions may be wider
- a problem may be *better* addressed if conceived as "shared" one in a shared T1
- key role of the experiment expertise at Tier-1





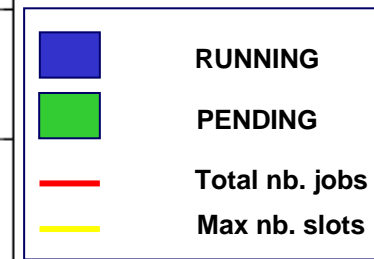
T1 farm occupancy

May be capable to address the needs of experiments

More focused effort needed on

- ✓ reliability and 24/7 support
- ✓ scheduled interventions
- ✓ troubleshooting

[all jobs at the T1]



➤ Underwent some migration seasons:

- ❑ OS: RH → **SLC v.3.0.4**
- ❑ mw: → **LCG v.2.6.0**
 - ❖ ~90% WN migrated, running LCG cert
- ❑ install WNs/servers: → **Quattor**
- ❑ batch scheduler: → **LSF v.6.1** | + LCG interf.

CMS share w.r.t total farm occupancy

month	Jun05		Jul05	
type of jobs	all jobs	Grid jobs	all jobs	Grid jobs
Jobs done [%]	24.4	43.6	14.8	27.6
CPU time [%]	5.1	8.3	5.0	7.9
Total time [%]	25.6	33.1	18.8	24.3





CMS analysis via CRAB



→ see [F.Fanzago et al, poster at this conf]

... "controlled" and "fake" (DC04) vs. "unpredictable" and "real" (now)

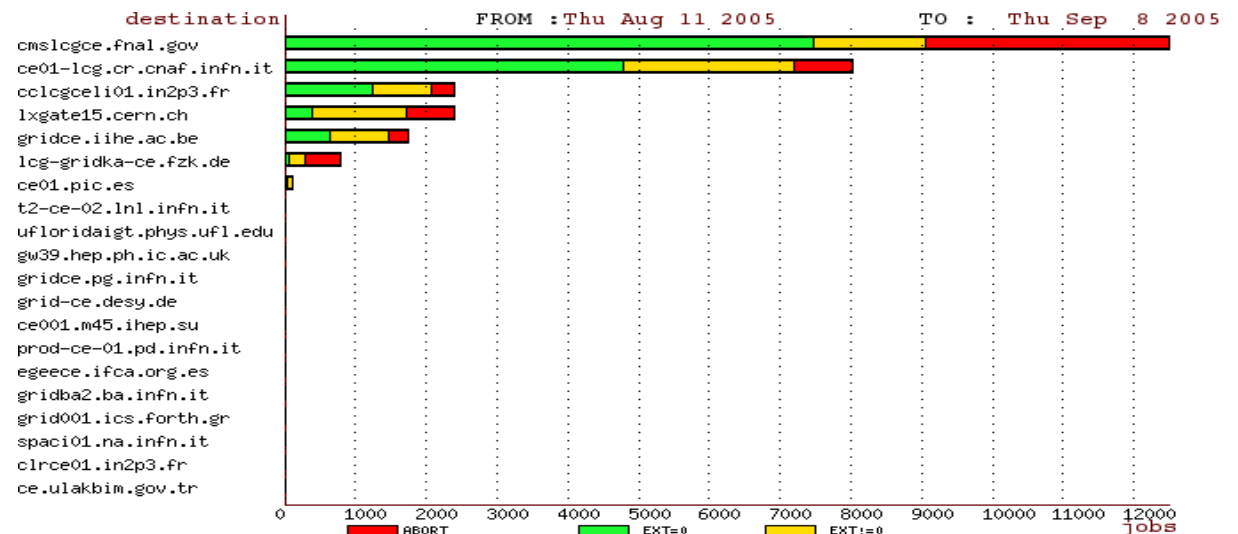
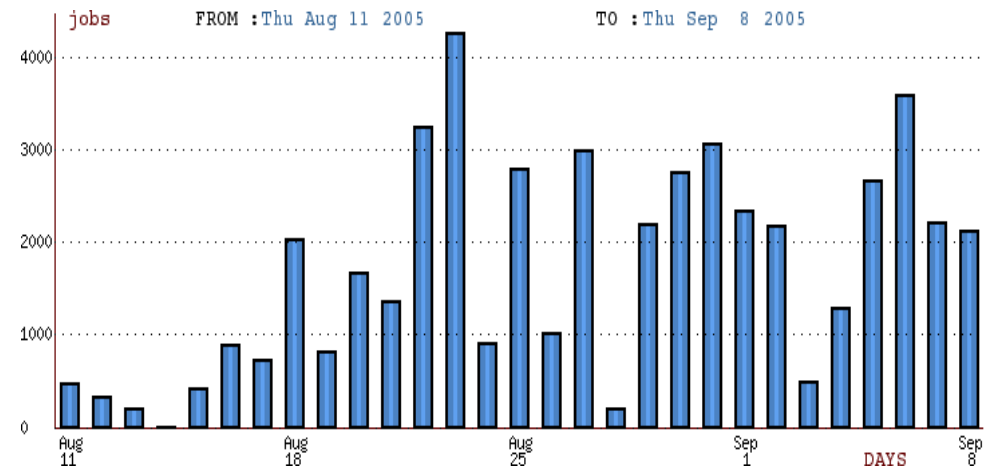
- T1 provides one full LCG site + 2 dedicated RBs/bdII + support to analysts

Data samples (~80M evts) for the P-TDR distributed at T1s

Grid allows end-to-end analysis:
send jobs to prelocated data

CRAB (CMS Remote Analysis Builder): a tool for job preparation, submission, monitoring

- important INFN contribution
- ✓ 10s users, 100s jobs/day





LCG Service Challenge 3 ('SC3')

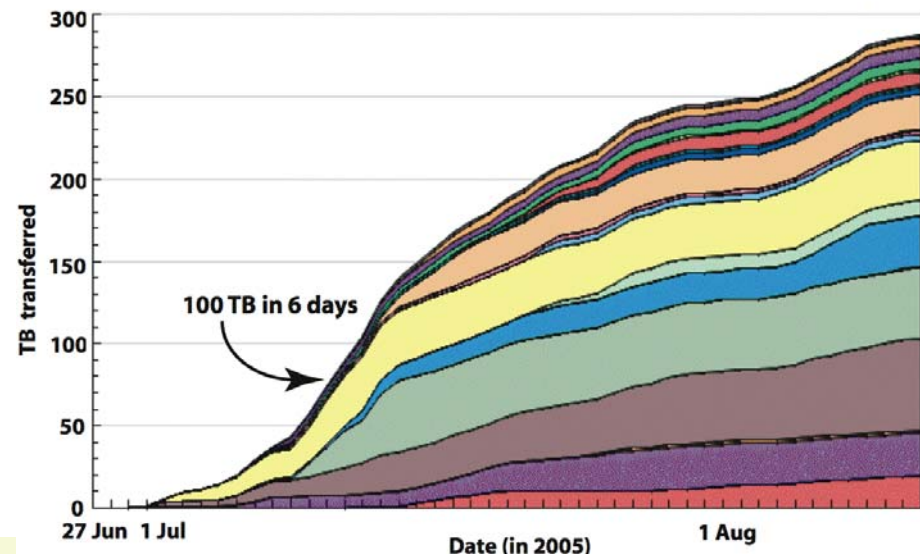


- data transfer and data serving in *real* use-cases
 - ❑ review existing infrastructure/tools (but analysis: SC4) and give a boost
 - ❑ "test and evaluate", and push crucial components to prod-quality
- Two phases:
 1. Jul05: SC3 "throughput" phase
 - ❖ Tiers simultaneous import/export, MSS involved: move real files, store on real hw
 2. >Sep05: SC3 "service" phase
 - ❖ small scale replica of the overall system
 - modest throughput, main focus on testing a quite complete environment
 - space for experiment-specific tests and inputs

INFN T1 participates to SC3,
also with CMS applications

- ❑ phase-1 post-mortem in progress
- ❑ phase-2 starting..

PhEDEx transfers during LCG SC3 by day





Summary



- C-TDR has served to converge on architecture for a baseline system
 - ❑ basic Grid infrastructure and services in places
 - ❑ relevant issue in Grid use is **reliability**

- key role of T1s in this scenario
 - ❑ focus on synergy among R&D activities and stability in providing services
 - ❑ experiment people expertise at the site is crucial

- ... now: technical design of missing components
 - ❑ use Grid services as much as possible for both WMS and DMS
 - ❑ bring up such a system over next 6-9 months

Past/forthcoming scheduled computing challenges:

- ✓ CMS DC04
- ✓ SC3 part 1 (Jul05): 'throughput'
 - SC3 part 2 (Sep05-Dec05): most services (but analysis)
 - SC4 (from Apr06): all services
 - CSA (summer 2006): full test of CMS computing system

