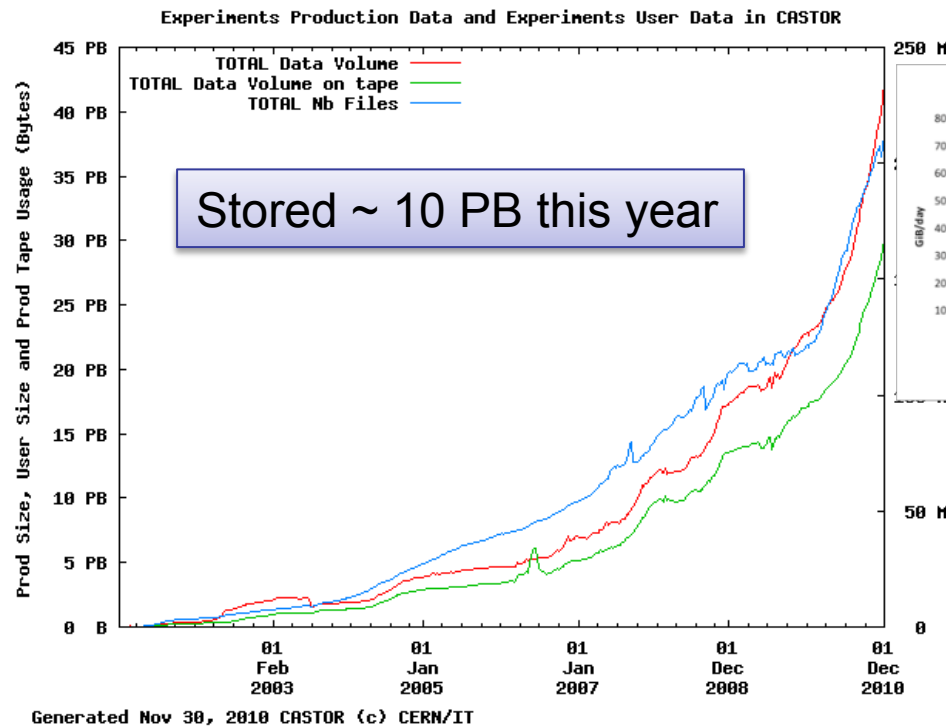# L'evoluzione dei modelli di calcolo a LHC
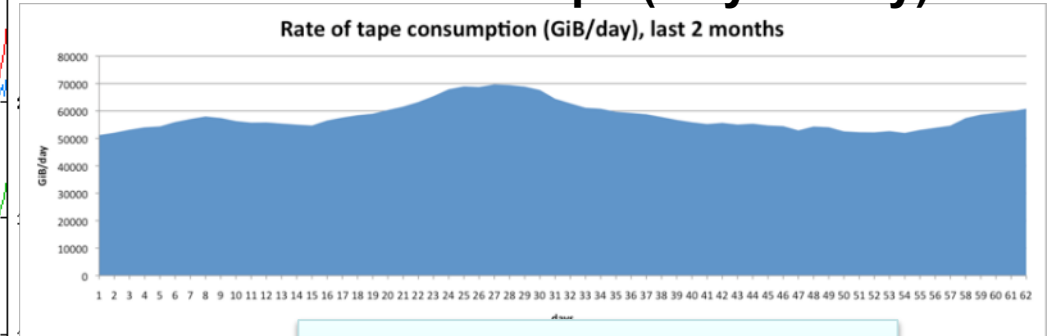
Concezio Bozzi, INFN Ferrara
Workshop CCR su stato e prospettive del calcolo scientifico
Legnaro, 17 febbraio 2011

*gratefully acknowledging*
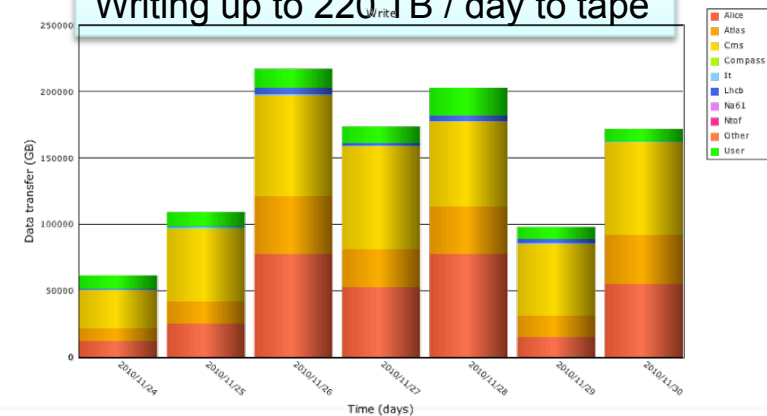*I. Bird, H. Newman, I. Fisk and R. Jones*
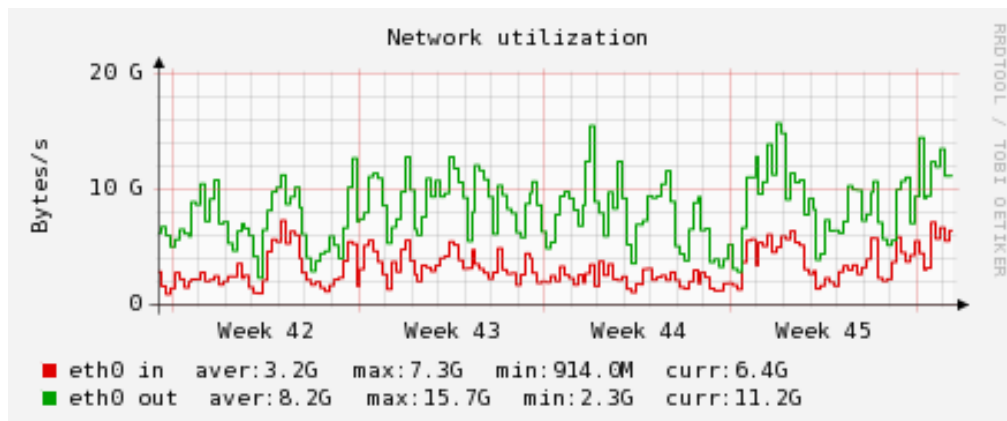
# Data did not fall on the floor

Experiments Production Data and Experiments User Data in CASTOR

Stored ~ 10 PB this year

Generated Nov 30, 2010 CASTOR (c) CERN/IT

## Data written to tape (Gbytes/day)

Rate of tape consumption (GiB/day), last 2 months

Writing up to 220 TB / day to tape

## Disk Servers (Gbytes/s)

Network utilization

eth0 in   aver:3.2G   max:7.3G   min:914.0M   curr:6.4G
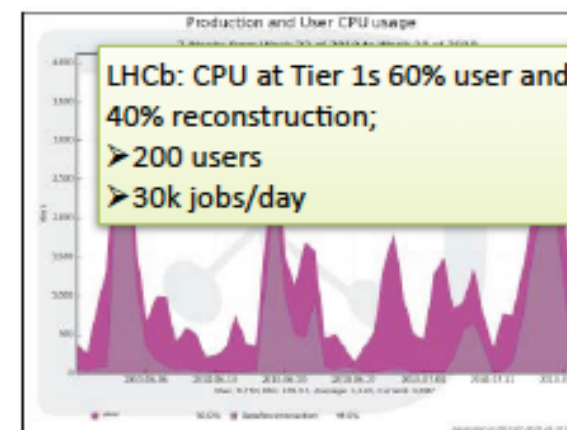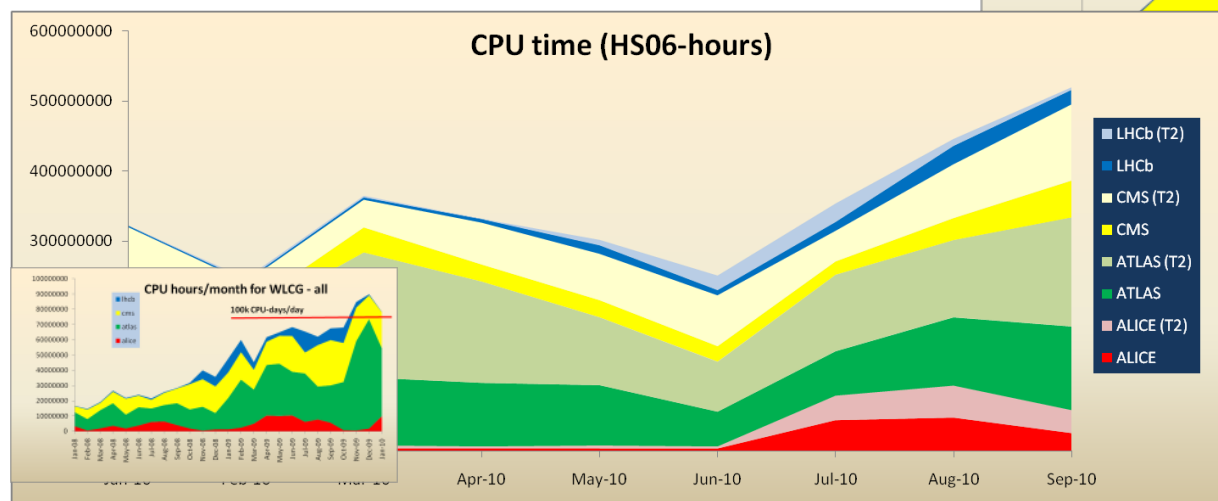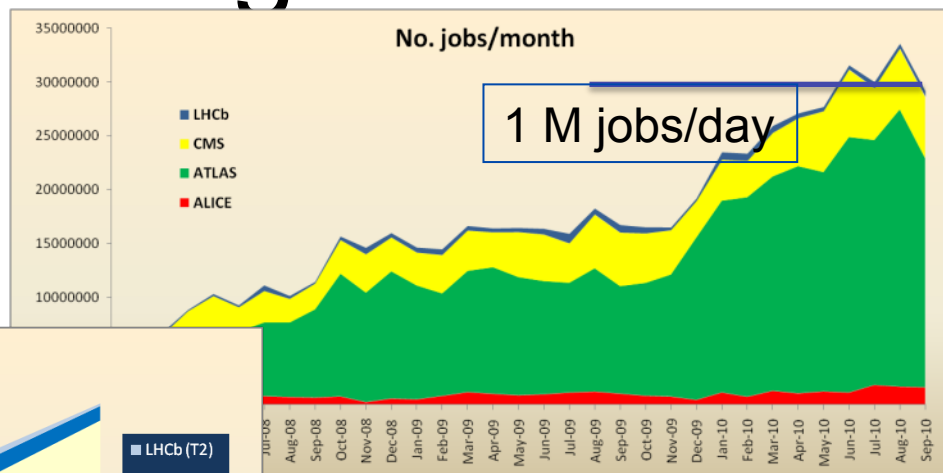eth0 out  aver:8.2G   max:15.7G  min:2.3G     curr:11.2G

## Tier 0 storage:

• Accepts data at average of 2.6 GB/s; peaks > 7 GB/s
• Serves data at average of 7 GB/s; peaks > 18 GB/s
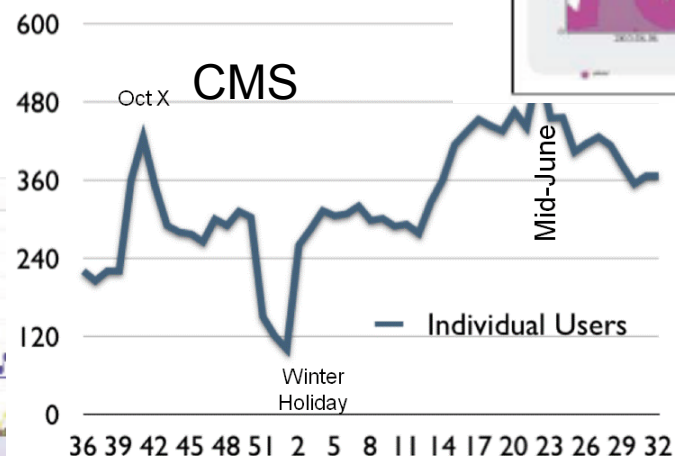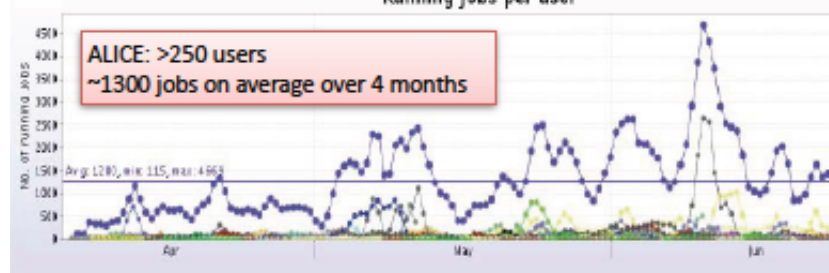• **CERN Tier 0 moves ~ 1 PB data per day**

2

# WLCG Usage

- **Use remains consistently high**
  - 1 M jobs/day; >>100k CPU-days/day
  - Actually much more inside pilot jobs

**No. jobs/month**

1 M jobs/day

Legend: LHCb, CMS, ATLAS, ALICE

**CPU time (HS06-hours)**

Legend: LHCb (T2), LHCb, CMS (T2), CMS, ATLAS (T2), ATLAS, ALICE (T2), ALICE

CPU hours/month for WLCG - all
100k CPU-days/day

As well as LHC data, large simulation productions ongoing

ALICE: ~200 users, 5-10% of Grid resources

Running jobs per user

ALICE: >250 users
~1300 jobs on average over 4 months

Production and User CPU usage

LHCb: CPU at Tier 1s 60% user and 40% reconstruction;
➤ 200 users
➤ 30k jobs/day

CMS

Oct X
Winter Holiday
Mid-June
Individual Users

36 39 42 45 48 51 2 5 8 11 14 17 20 23 26 29 32

Large numbers of analysis users
CMS ~800,
ATLAS ~1000,
LHCb/ALICE ~200

3

# CPU – July



**CPU Delivered - July 2010 - by Tier**

- CERN
- RAL
- KIT
- BNL
- NL Tier-1
- CC-IN2P3
- CNAF
- FNAL
- PIC
- NDGF
- TRIUMF

- Significant use of Tier 2s for analysis
  - frequently-expressed concern that too much analysis would be done at CERN is not reflected

**Tier 2 CPU Delivered - July 2010 - by Country**

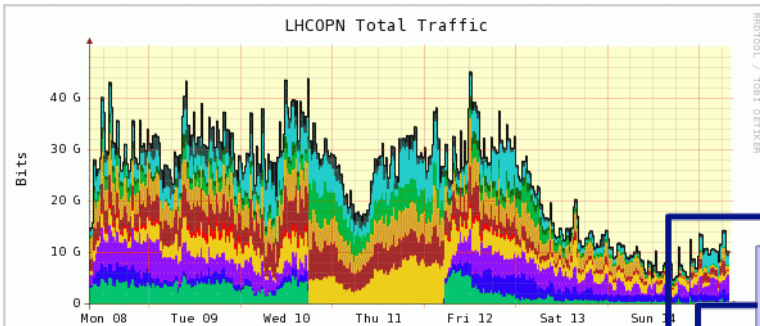| | |
|---|---|
| USA | France |
| Germany | UK |
| Italy | Russian Fed. |
| Spain | Czech Rep. |
| Japan | Poland |
| Sweden | Switzerland |
| Romania | China |
| Canada | Portugal |
| Hungary | Norway |
| Taipei | Republic of Korea |
| Austria | Israel |
| Estonia | Turkey |
| Finland | Australia |
| Belgium | Brazil |
| India | Pakistan |
| Slovenia | |

- **Tier 0 capacity underused in general**
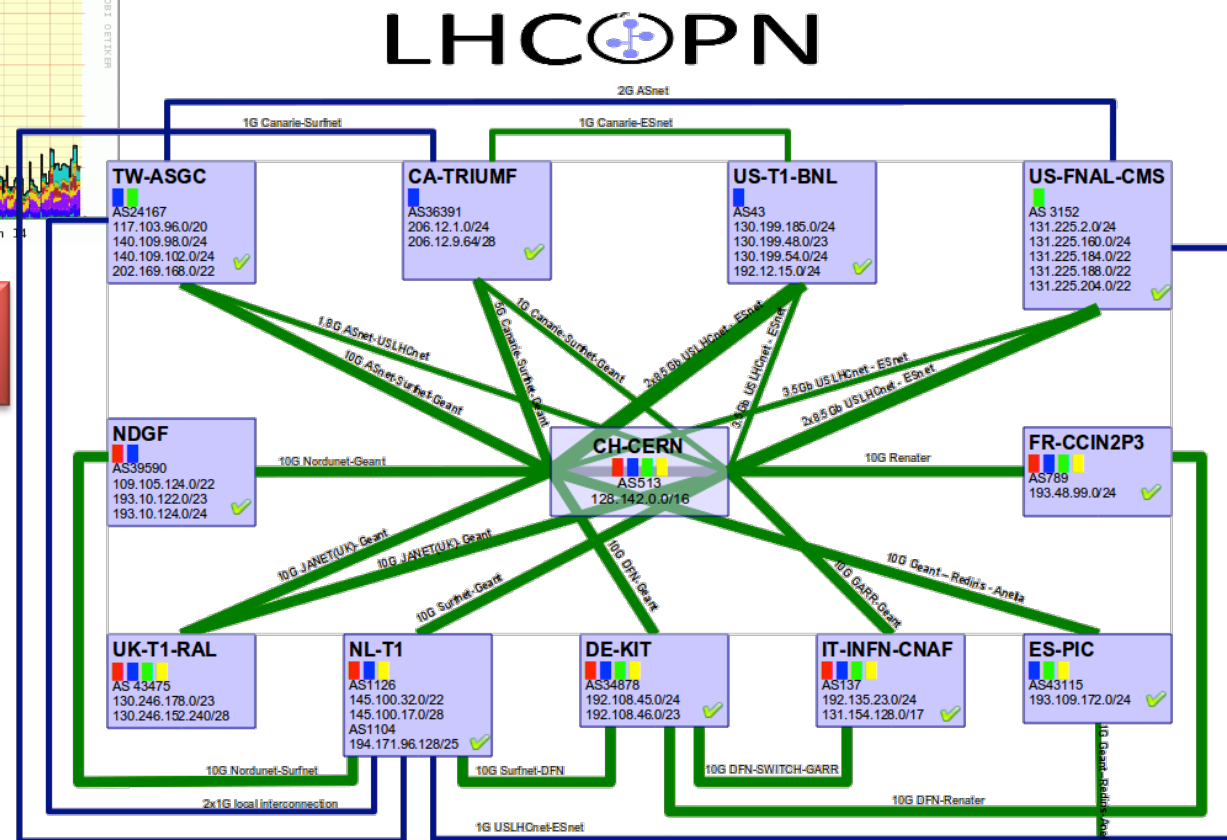  - But this is expected to change as luminosity increases

# Data transfer

- Data transfer capability today able to manage much higher bandwidths than expected/feared/planned

**LHCOPN Total Traffic**

Fibre cut during STEP'09:
Redundancy meant no interruption

**Data transfer:**
- SW: gridftp, FTS (interacts with endpoints, recovery), experiment layer

- HW: light paths, routing, coupling to storage

- Operational: monitoring

## LHCOPN

2G ASnet

1G Canarie-Surfnet          1G Canarie-ESnet

**TW-ASGC**
AS24167
117.103.96.0/20
140.109.98.0/24
140.109.102.0/24
202.169.168.0/22

**CA-TRIUMF**
AS36391
206.12.1.0/24
206.12.9.64/28

**US-T1-BNL**
AS43
130.199.185.0/24
130.199.48.0/23
130.199.54.0/24
192.12.15.0/24

**US-FNAL-CMS**
AS 3152
131.225.2.0/24
131.225.160.0/24
131.225.184.0/22
131.225.188.0/22
131.225.204.0/22

**NDGF**
AS39590
109.105.124.0/22
193.10.122.0/23
193.10.124.0/24

**CH-CERN**
AS513
128.142.0.0/16

**FR-CCIN2P3**
AS789
193.48.99.0/24

**UK-T1-RAL**
AS 43475
130.246.178.0/23
130.246.152.240/28

**NL-T1**
AS1126
145.100.32.0/22
145.100.17.0/28
AS1104
194.171.96.128/25

**DE-KIT**
AS34878
192.108.45.0/24
192.108.46.0/23

**IT-INFN-CNAF**
AS137
192.135.23.0/24
131.154.128.0/17

**ES-PIC**
AS43115
193.109.172.0/24

10G Nordunet-Surfnet          10G Surfnet-DFN          10G DFN-SWITCH-GARR          10G DFN-Renater

2x1G local interconnection

1G USLHCnet-ESnet

T0-T1 and T1-T1 traffic
T1-T1 traffic only
Not deployed yet
(thick) >=10Gbps
(thin) <10Gbps

= Alice    = Atlas
= CMS    = LHCb
= internet backup available
p2p prefix: 192.16.166.0/24
edoardo.martelli@cern.ch 20100916

& the academic/research networks for Tier1/2!
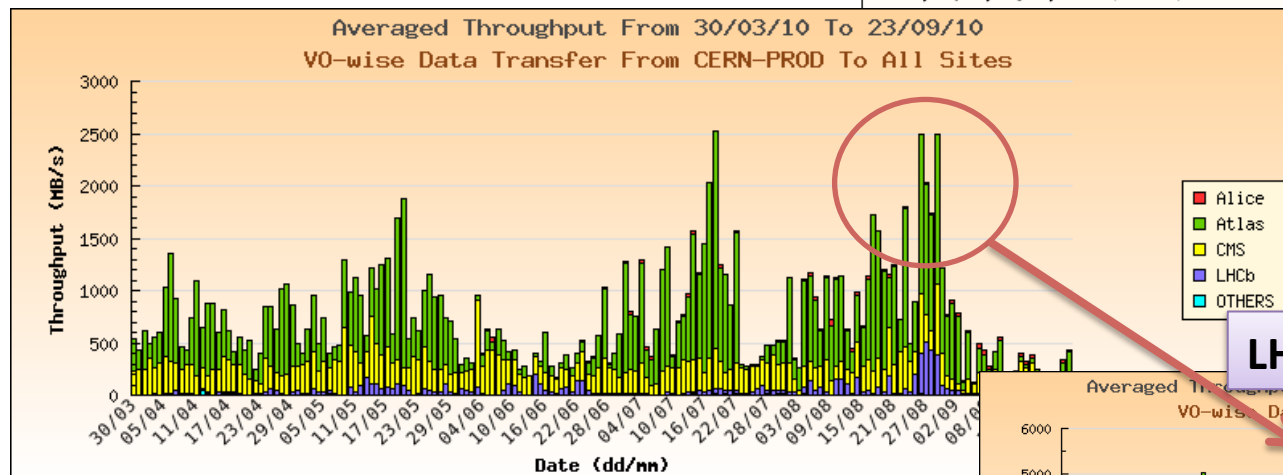
5
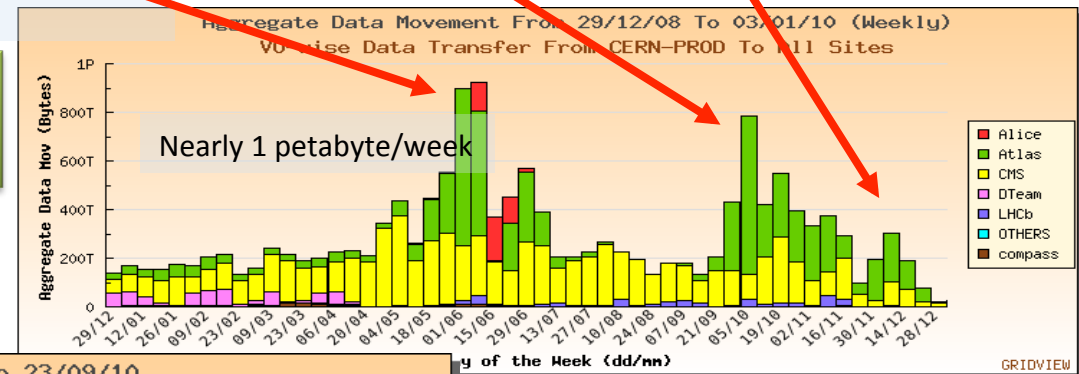
# Data transfers

Final readiness test (STEP'09)

Preparation for LHC startup

LHC physics data

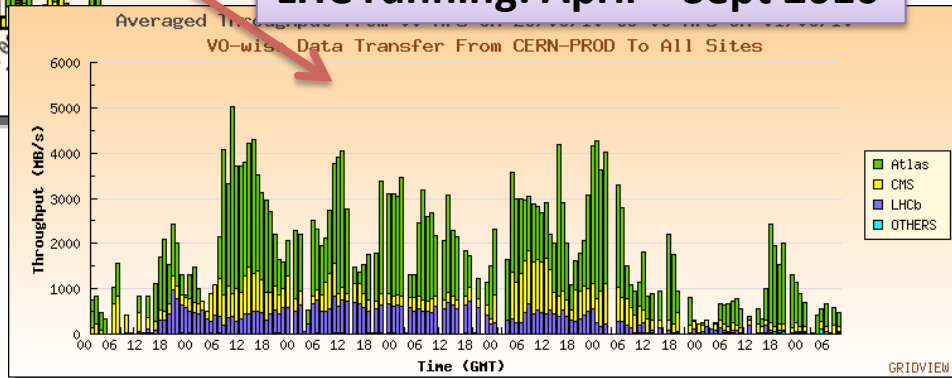2009: STEP09 + preparation for data

Aggregate Data Movement From 29/12/08 To 03/01/10 (Weekly)
VO-wise Data Transfer From CERN-PROD To All Sites

Nearly 1 petabyte/week

Alice
Atlas
CMS
DTeam
LHCb
OTHERS
compass

GRIDVIEW

Averaged Throughput From 30/03/10 To 23/09/10
VO-wise Data Transfer From CERN-PROD To All Sites

Alice
Atlas
CMS
LHCb
OTHERS

**LHC running: April – Sept 2010**

Averaged Throughput From 00 HRS on 28/07/10 To 00 HRS on 01/08/10
VO-wise Data Transfer From CERN-PROD To All Sites

Atlas
CMS
LHCb
OTHERS

GRIDVIEW

LHCOPN TOTAL Traffic Flow 1 (Out-bound)

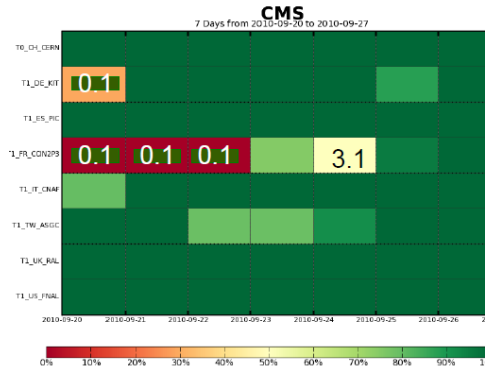Traffic on OPN up to 70 Gb/s!
- ATLAS early reprocessing campaigns

6

# Reliabilities



**Site Reliability: CERN + Tier 1s**

Legend: Average — Average - 8 best sites — Target
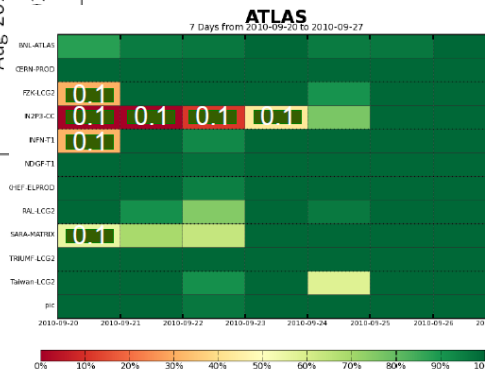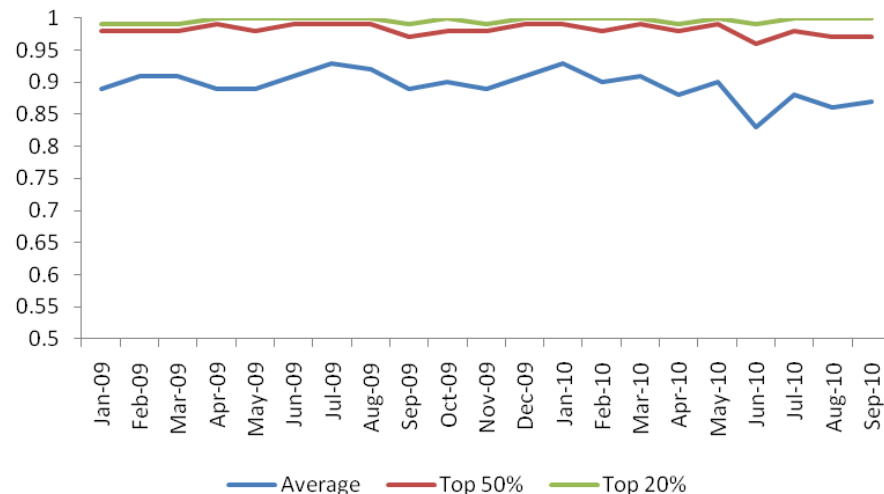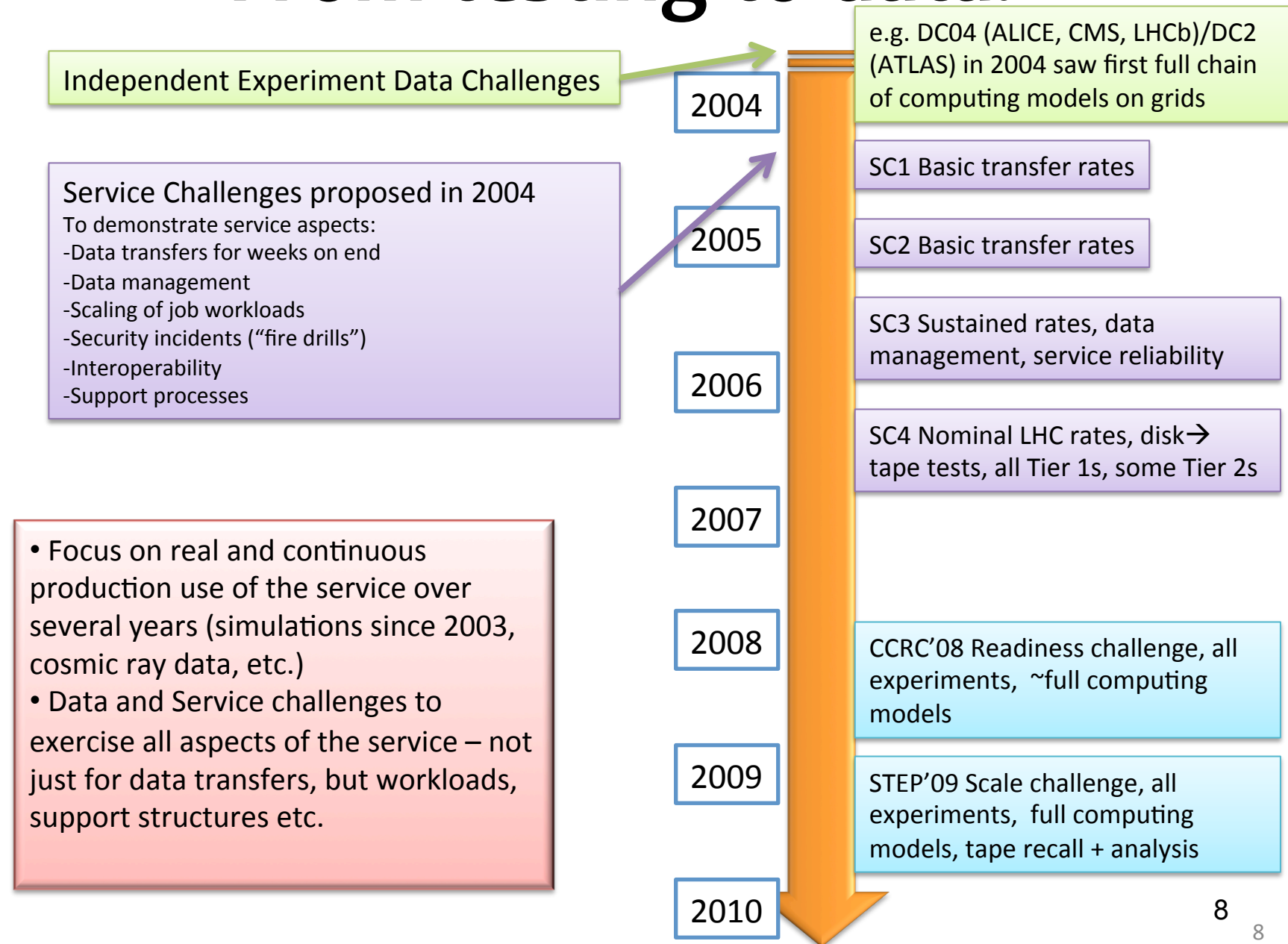
Experiment-measured site availabilities:

Includes down times during security patching;
At times ~50% of resources were unavailable.
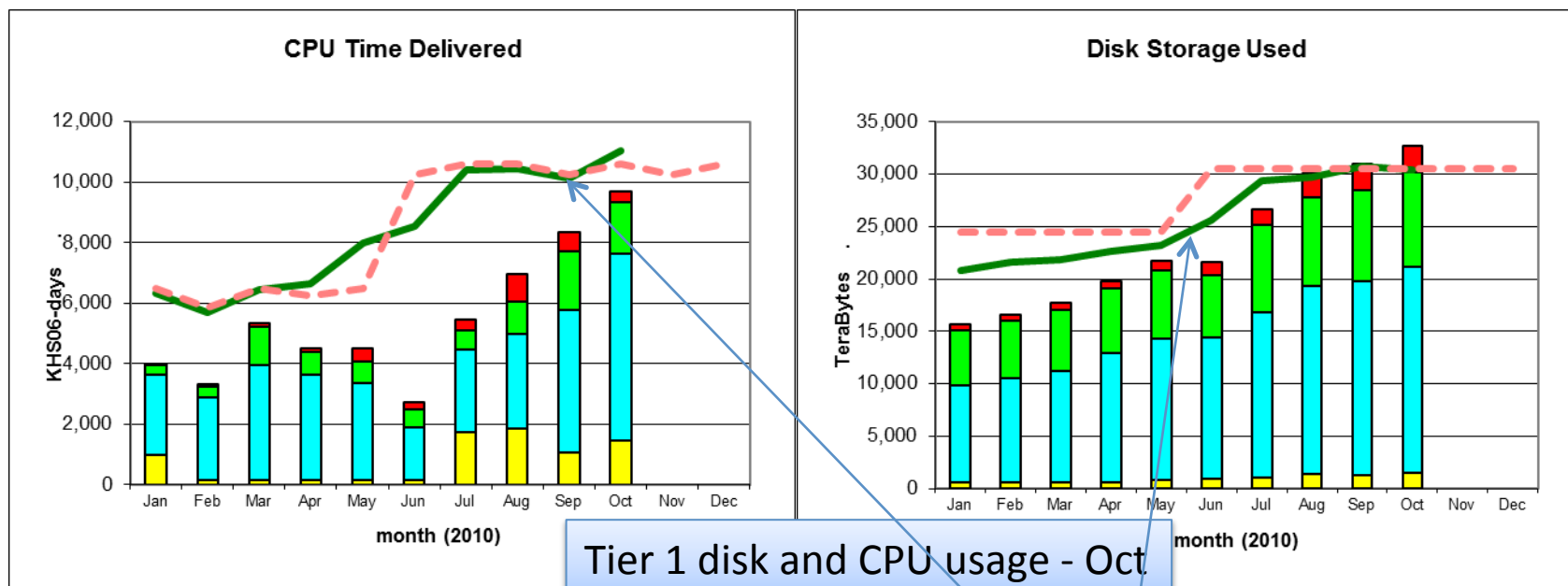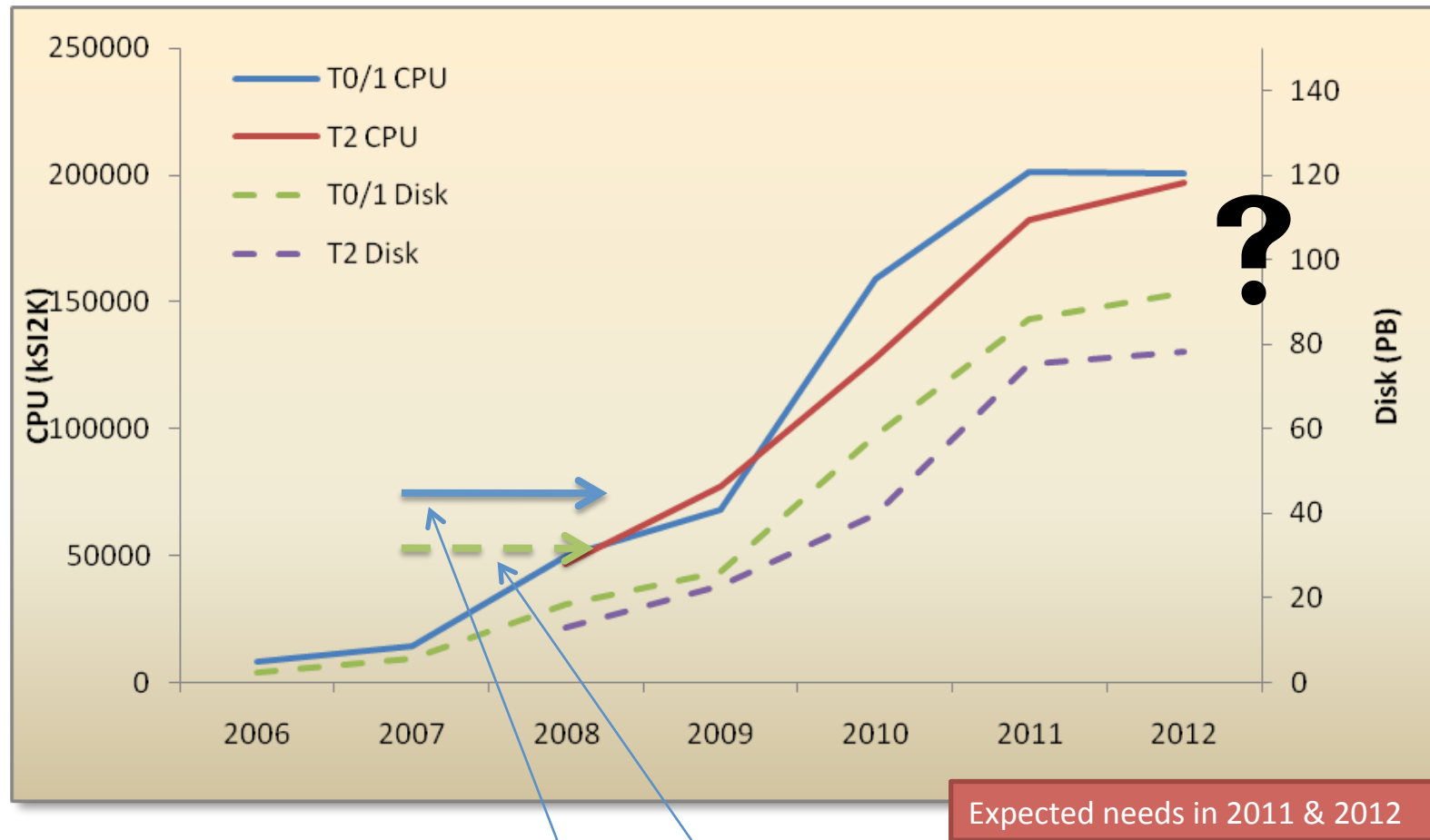
**Tier 2 Reliabilities**

Legend: Average — Top 50% — Top 20%

ATLAS
7 Days from 2010-09-20 to 2010-09-27

ALICE
7 Days from 2010-09-20 to 2010-09-27

CMS
7 Days from 2010-09-20 to 2010-09-27

LHCb
7 Days from 2010-09-20 to 2010-09-27

# From testing to data:

Independent Experiment Data Challenges

e.g. DC04 (ALICE, CMS, LHCb)/DC2 (ATLAS) in 2004 saw first full chain of computing models on grids

**2004**

Service Challenges proposed in 2004
To demonstrate service aspects:
-Data transfers for weeks on end
-Data management
-Scaling of job workloads
-Security incidents ("fire drills")
-Interoperability
-Support processes

SC1 Basic transfer rates

**2005**

SC2 Basic transfer rates

SC3 Sustained rates, data management, service reliability

**2006**

SC4 Nominal LHC rates, disk→ tape tests, all Tier 1s, some Tier 2s

**2007**

• Focus on real and continuous production use of the service over several years (simulations since 2003, cosmic ray data, etc.)
• Data and Service challenges to exercise all aspects of the service – not just for data transfers, but workloads, support structures etc.

**2008**

CCRC'08 Readiness challenge, all experiments, ~full computing models

**2009**

STEP'09 Scale challenge, all experiments, full computing models, tape recall + analysis

**2010**

8

# Resource usage

- Now Tier 1s and Tier 2s start to be fully occupied; as planned with reprocessing, analysis, and simulation loads



Tier 1 disk and CPU usage - Oct

| Tier 1 use - Oct | CPU use/pledge | Disk use/pledge |
|---|---|---|
| ALICE | 1.04 | 0.25 |
| ATLAS | 0.94 | 0.89 |
| CMS | 0.54 | 0.74 |
| LHCb | 0.27 | 0.79 |
| **Overall** | **0.78** | **0.75** |

NB: Assumed effic factors
0.85 for CPU
0.70 for disk

9

# Resource Evolution (no run in 2012)



Expected needs in 2011 & 2012

Need foreseen @ TDR for T0+1 CPU and Disk for 1st nominal year

NB. In 2005 only 10% of 2008 requirement was available. The ramp-up has been enormous!

# Elements of a computing model

- Basic parameters
  - How many events, how many event types
  - Event size, event types
  - Processing times

- Data distribution
  - Filtering, skimming, slimming
  - How many copies in Tier1/Tier2 ensembles

- Data processing
  - "Scheduled" activities: how many processes in a year? How long is a reprocessing cycle? How many versions on disk?
  - "Chaotic" activities: how many analysis groups/users? How frequently do they access data? How much time for a full pass?

# Experiment models have evolved

- Models all ~based on the MONARC tiered model of 10 years ago
- Several significant variations, however

# The Monarc rationale

- The MONARC computing model of 2000 relied heavily on data placement
- Jobs were sent to datasets already resident on sites
- Multiple copies of the data would be hosted on the distributed infrastructure
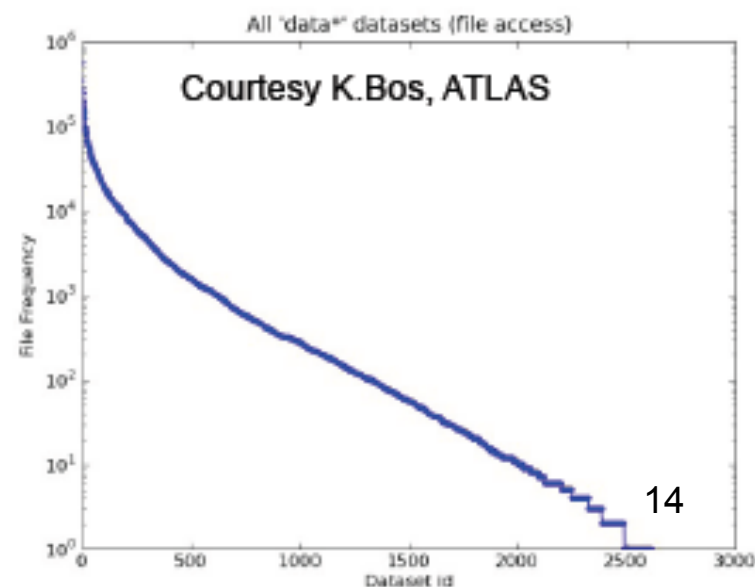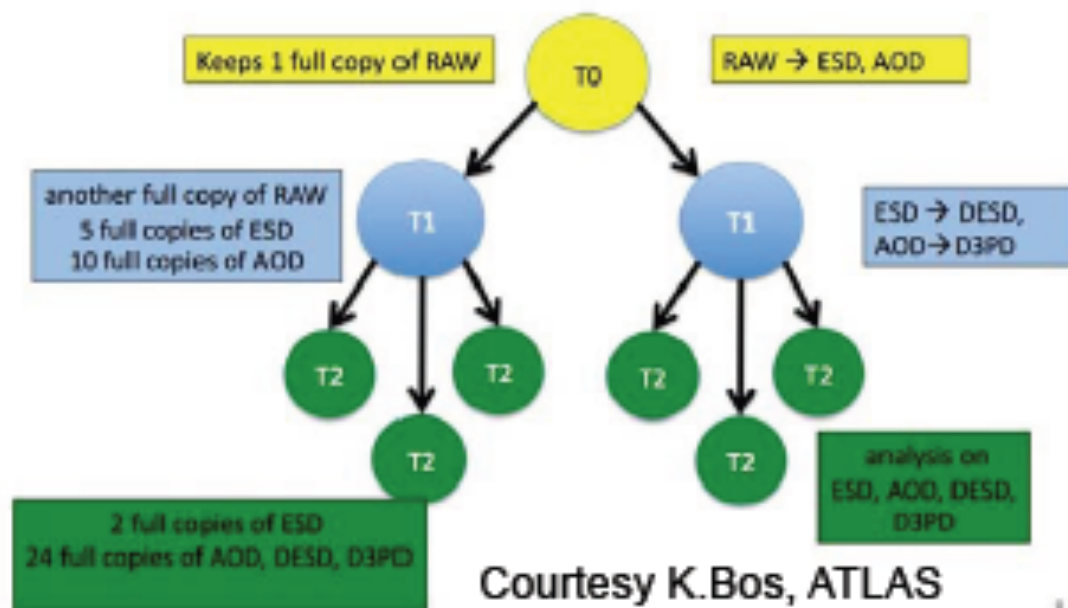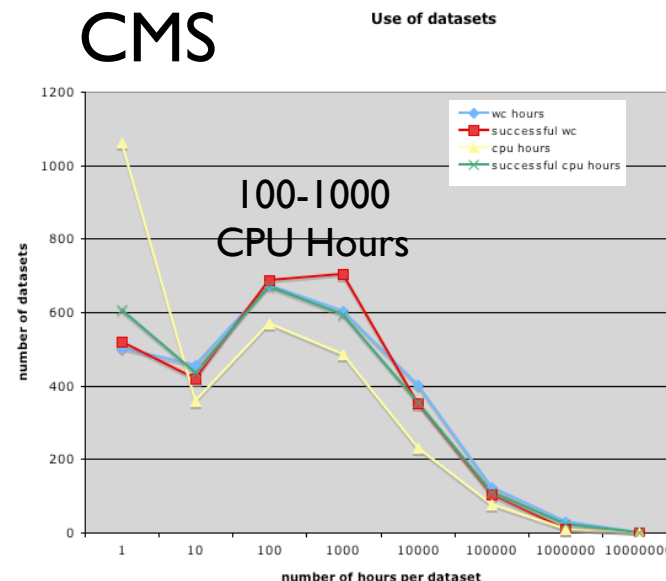- General concern that the network would be insufficient or unreliable

- As we have just seen, this is no longer the case nowadays
- Look at ways to make more efficient use of the resources



13

Richard Mount

## Reuse of PD2P Datasets - Month 2
### (July 15 - August 14)

| Value | Count |
|-------|-------|
| 0 | 1274 |
| 10 | 72 |
| 100 | 45 |
| 1000 | 66 |
| 10000 | 103 |
| 100000 | 15 |

24 full copies of AOD, DESD, D3PD

Courtesy K.Bos, ATLAS

15

0000

5

14

5

# Evolution of data placement

- Move towards caching of data rather than strict planned placement
- Download the data when required
  - Selects popular datasets automatically
  - When datasets no longer used will be replaced in the caches
- Data sources can be any (Tier 0, 1, 2)
- Can still do some level of intelligent pre-placement
- Understanding a distributed system built on unreliable and asynchronous components means
  - Accepting that catalogues may be not fully updated
  - Data may not be where you thought it was
  - Thus must allow remote access to data (either by caching on demand and/or by remote file access)
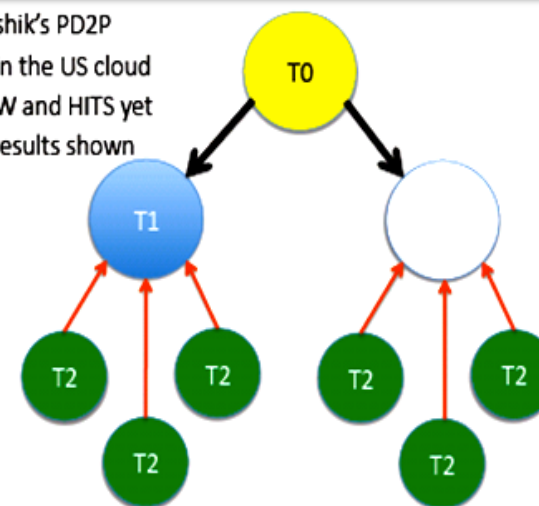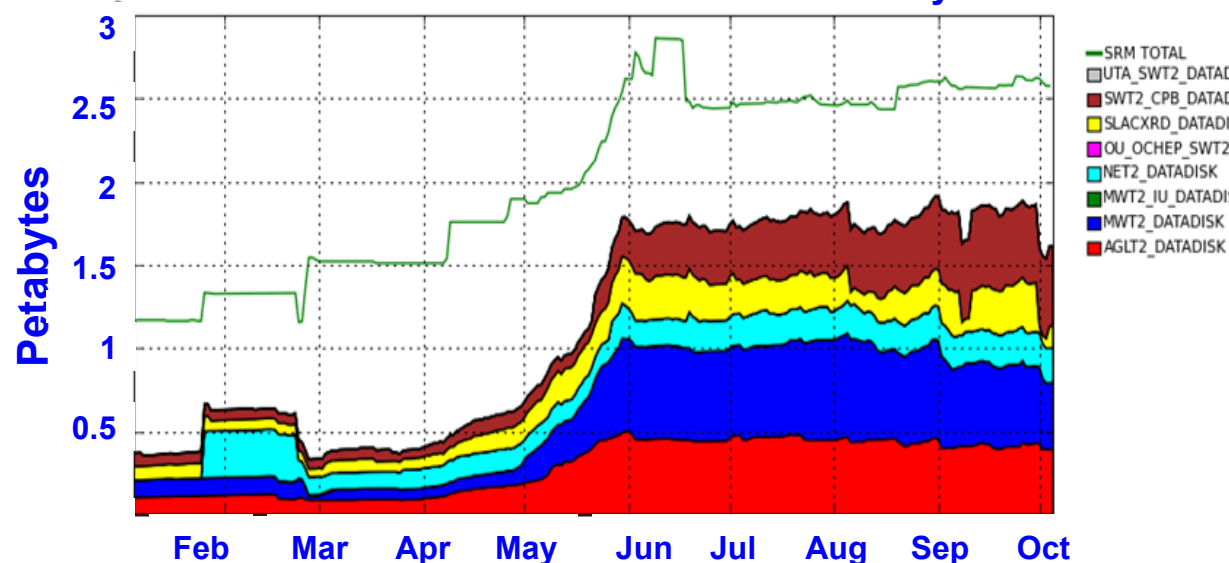
# Pull Model in Atlas BNL Cloud

- **PD2P: Atlas implementation of the pull model**
  - **Tier1 used as repository (Tier0-Tier1: Push)**
  - **Dynamic data placement at Tier2s**
  - **Dataset is subscribed to a Tier 2 if no other copies are available (except at a Tier 1), as soon as any user needs it**
- **Deployed in the US (BNL) cloud in June**

## Data Pull Model I

- This is Kaushik's PD2P
- Runs now in the US cloud
- Not for RAW and HITS yet
- Intersting results shown



**Cumulative evolution of DATADISK by site**



Legend:
- SRM TOTAL
- UTA_SWT2_DATAD
- SWT2_CPB_DATAD
- SLACXRD_DATADI
- OU_OCHEP_SWT2
- NET2_DATADISK
- MWT2_IU_DATADIS
- MWT2_DATADISK
- AGLT2_DATADISK

**Kaushik De, Atlas Week Oct 2010**

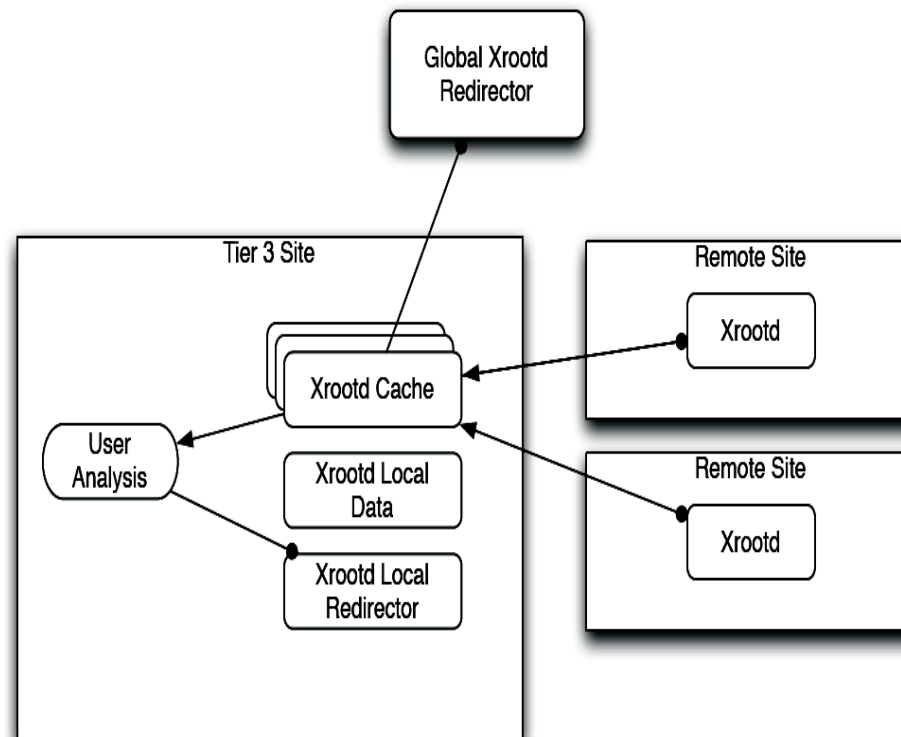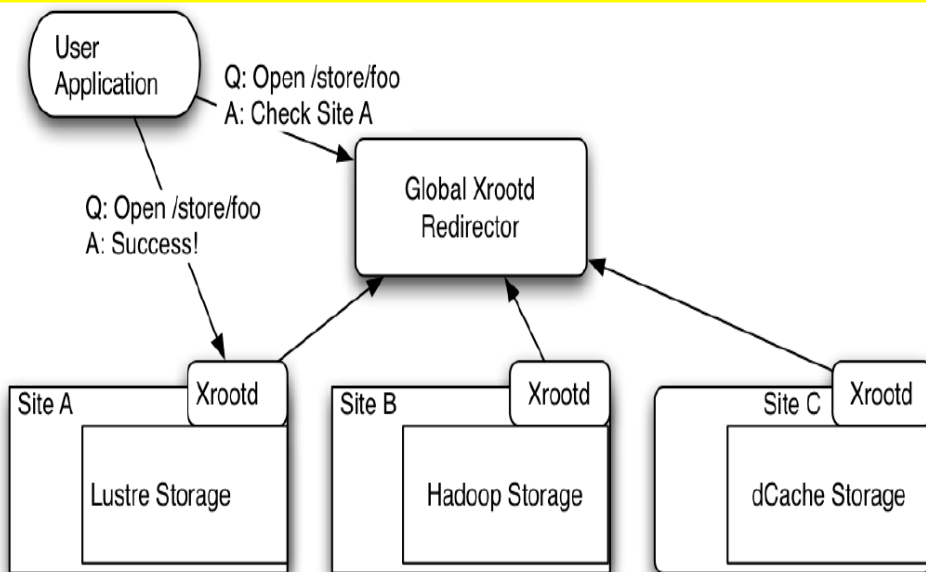**Before: Exponential rise from right after LHC start**

**Much slower rise in disk utilization since July**

16

# Remote Data Access and Local Processing with Xrootd (CMS)

- **Useful for smaller sites with less (or even no) data storage**

- **Only selected objects are read (with object read-ahead).** No transfer of entire data sets

- **CMS demonstrator: Omaha diskless Tier3, served data from Caltech and Nebraska (Xrootd)**





**Strategic Decisions:**
Remote Access
Vs
Data Transfers

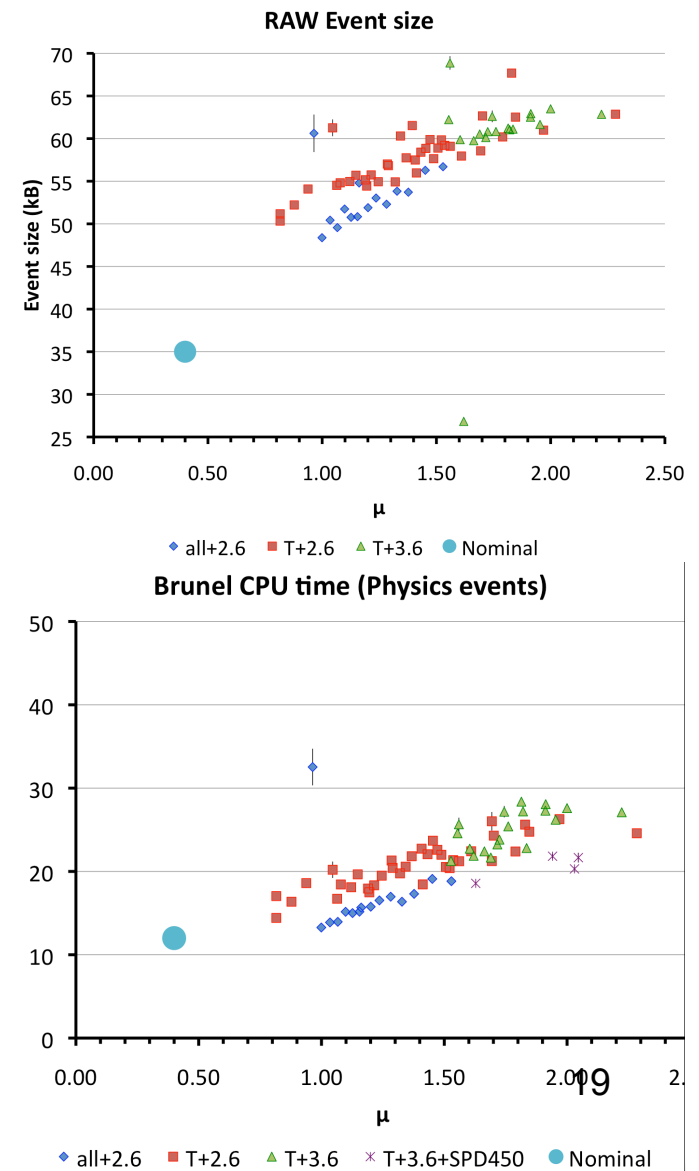Brian Bockelman, September 2010

17

# Implications for networks

- Hierarchy of Tier 0, 1, 2 no longer so important
- Tier 1 and Tier 2 may become more equivalent for the network
- Traffic could flow more between countries as well as within (already the case for CMS)
- Network bandwidth (rather than disk) will need to scale more with users and data volumes
- Data placement will be driven by demand for analysis and not pre-placement
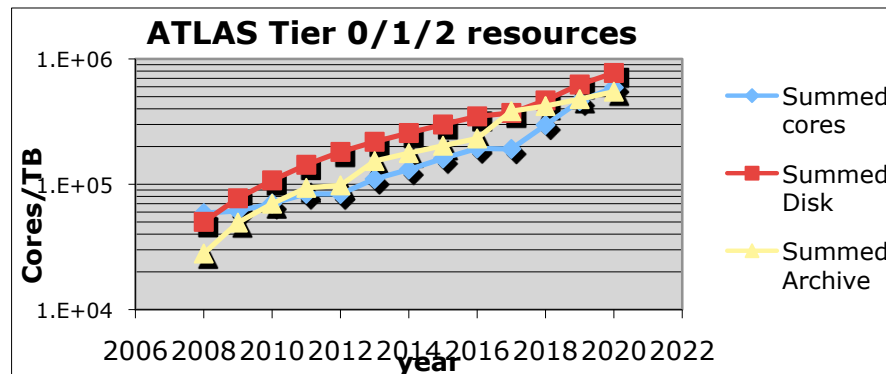
# Processing challenges

Global Grid usage



19

# Future Challenges

- We assume we can use growth in CPU
  - But this implies changing architectures
  - And handle the data throughput
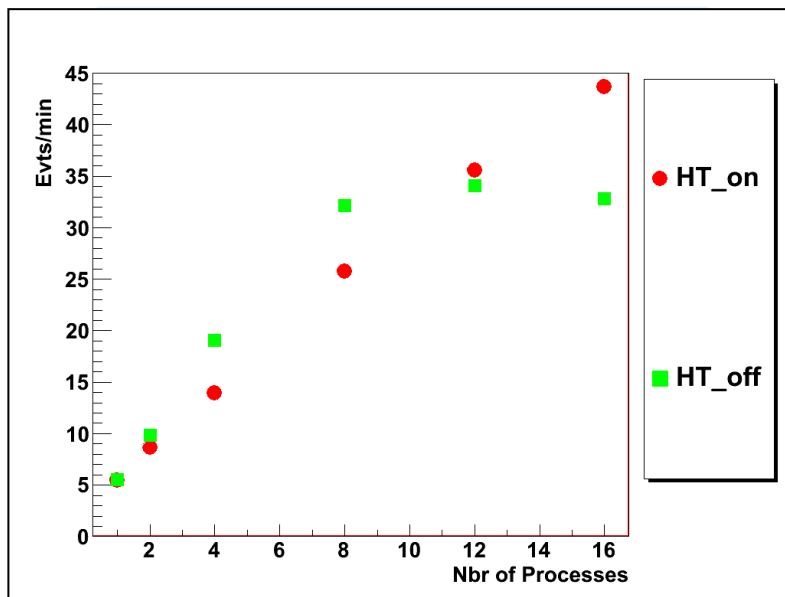


**ATLAS Tier 0/1/2 resources**

- Experiments already working to deal with multi cores
  - Many cores and GPGPUs are down the line
- We need to use them or be very clear why we cannot

# Parallelism

- Generally work smarter!
  - E.g. AthenaMP Event level parallelism
    - Share common memory between parent and daughter processes to allow many on a single node
    - Some speed-up using event loop parallelism
  - Also share common pages between processes with KSM
    - Real gains in memory use, but some slow-down
  - Cache as much as you can (e.g. pile-up events)
  - Also Non-Uniform Memory Access, simultaneous multi-threading
- Issues: hard to monitor performance in parallel jobs



- Other approaches
  - Job level parallelism (e.g parallel Gaudi) & hyperthreading
  - CMS working on this sort of 'workflow' parallelism
  - Pinning of processes to cores or hyperthreads with Affinty

21

# Virtualisation and "clouds"

- .... Another hype / marketing / diversion ???

- Yes, but
  - Virtualisation is already helping in several areas
    - Breaking the dependency nightmare
    - Improving system management, provision of services on demand
    - Potential to help use resources more effectively and efficiently (many of us have power/cooling limitations)
    - Use of remote computer centres
  - Cloud technology
    - Let's not forget why we have and need a "grid"; much of this cannot be provided by today's "cloud" offerings
      - Collaboration (VO's), worldwide AAI and trust, dispersed resources (hw and people),
    - Although we should be able to make use of commercial clouds transparently

# What about Grid middleware?

## The *Basic* Baseline Services – from the TDR (2005)

- **Storage Element**
  - **Castor, dCache** [SRM is too complex]
  - **Storm added in 2007**
  - **SRM 2.2 – deployed in production – Dec 2007**

- **Basic transfer t** [OK, but why not HTTP?]

- **File Tran** [OK for some use cases]

  [OK, but must sync with storage / No need for distributed catalogue]

- **LCG File**

- **LCG data mgt tools - lcg-utils**

- **"Posix" I/O –** ✔
  - **Grid File Access Library (GFAL)**

- **Synchronised databases T0↔T1s**
  - **3D project** [Frontier/Squid for many use cases]

- **Informa** [LDAP → messaging? Static vs dynamic info]
  - **BDII**

- **Compute E** [Still have LCG-CE, not yet replaced; MUPJs!]
  - **Globus/**
  - **web ser**
  - **Support for multi-user pilot jobs**

  [Actual LHC use cases much simpler / Pilot frameworks may supercede it]

- **W**
  - **WMS, LB**

- **VO Management System (VOMS), MyProxy** ✔

- **VO Boxes** → Virtual machine

- **Application** → CVMFS or Squid

- **Job Monitoring Tools** → MSG, Nagios, etc

- **APEL etc.** ✔

# What about grid middleware?

- Clearly a thinner layer today than originally imagined
  - And the actual usage is far simpler
- Experiment layer is deeper ... And different from one to the other
- Experiments had to work hard to (mostly) hide the grid details from users
- Pilot jobs are (almost) ubiquitous in all experiments
- Simplification of some services is possible and helps long term maintenance and support
- The current grid infrastructure can sit transparently over virtualised (cloud) services
  - And provide a potential path for evolutionary change

# Automation, monitoring and testing

- Operations are still too effort-intensive
  - increase automation
- Monitoring is essential to keep system going and understand its usage patterns
  - More to be done for storage systems
  - Tendency to have too much!
  - Keep distinct views for experiments, sites, and managers
- Lots of testing results in outstanding availability and reliability
  - Revealed many configuration problems (e.g. ATLAS Hammercloud)

# Conclusioni

- Il sistema di calcolo distribuito degli esperimenti a LHC ha funzionato molto bene in questo primo periodo di presa dati
- Le risorse a disposizione degli esperimenti erano "comode"
  - Che succederà quando LHC arriverà a regime?
- I modelli di calcolo si stanno evolvendo allo scopo di ottimizzare l'utilizzo delle risorse sfruttando gli "asset" consolidati
  - Bisogna capire bene le implicazioni sulla rete
- Occorre rimanere al passo con le tecnologie di punta…
  - Cambiamenti di architettura per many-core? GPU?
  - Virtualizzazione?
  - Cloud computing?
- …continuando a garantire il buon funzionamento di quanto è stato fatto finora
  - Automatizzare, testare, monitorare