

CERN

European Organization for Nuclear Research
Organisation Européenne pour la Recherche Nucléaire

The LHC Computing Challenge

Tony Cass
Leader, Fabric Infrastructure & Operations Group
Information Technology Department

15th November 2007



Outline

- Introduction to CERN and Experiments
- LHC Computing
- Challenges
- Summary/Conclusion



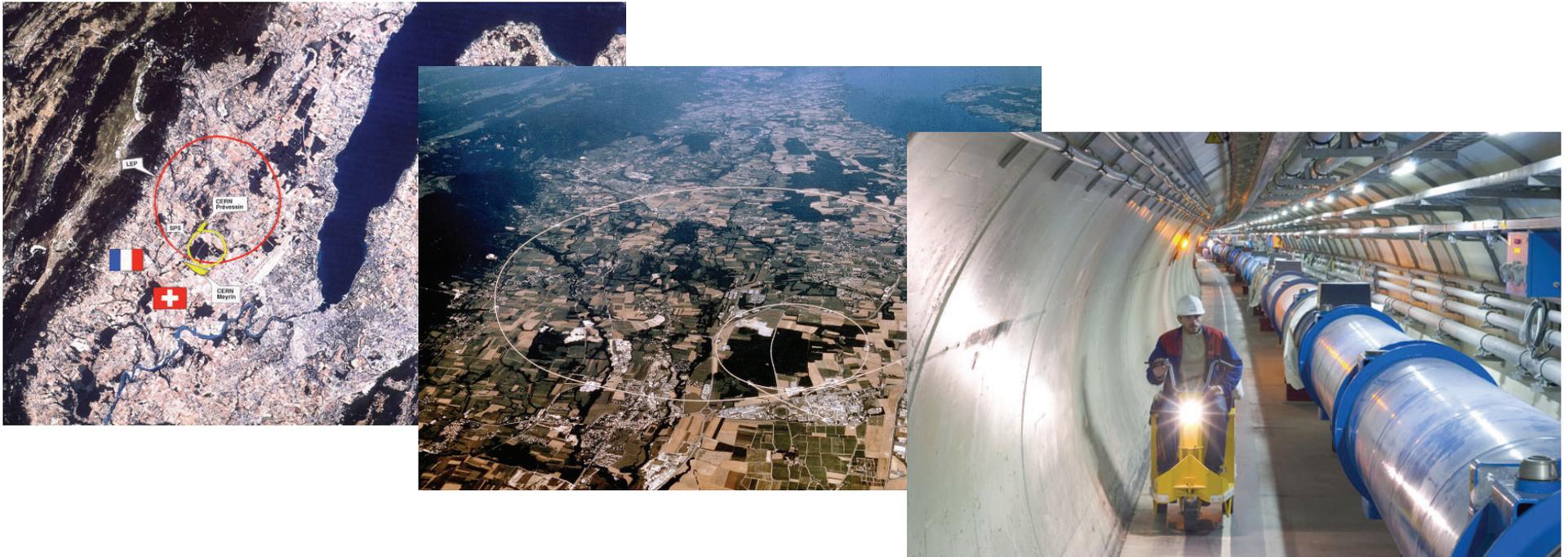
Outline

- Introduction to CERN and Experiments
- LHC Computing
- Challenges
- Summary/Conclusion



CERN

The fastest racetrack on the planet...

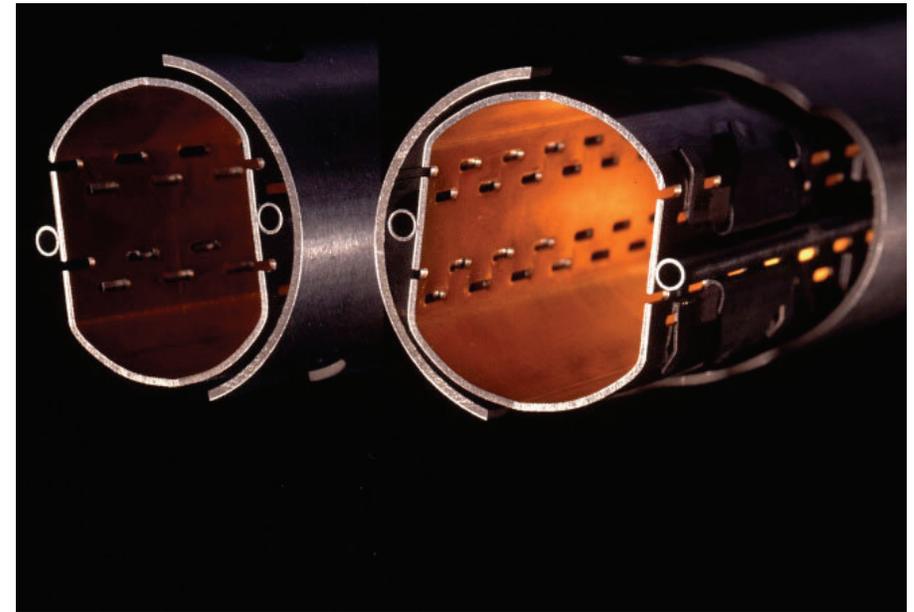
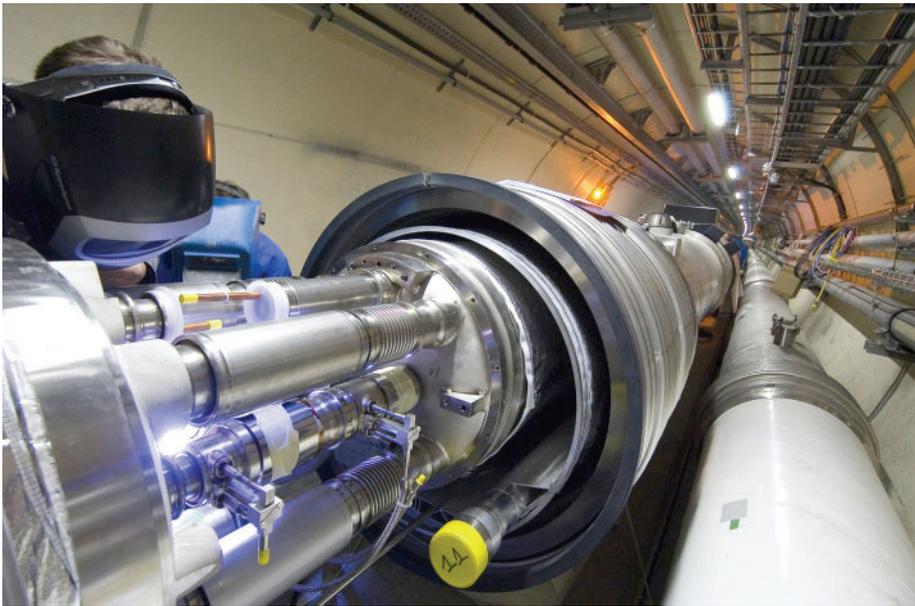


Trillions of protons will race around the 27km ring in opposite directions over 11,000 times a second, travelling at 99.999999991 per cent the speed of light.



CERN

The emptiest space in the solar system...



To accelerate protons to almost the speed of light requires a vacuum as empty as interplanetary space. There is 10 times more atmosphere on the moon than there will be in the LHC.



CERN

One of the coldest places in the universe...

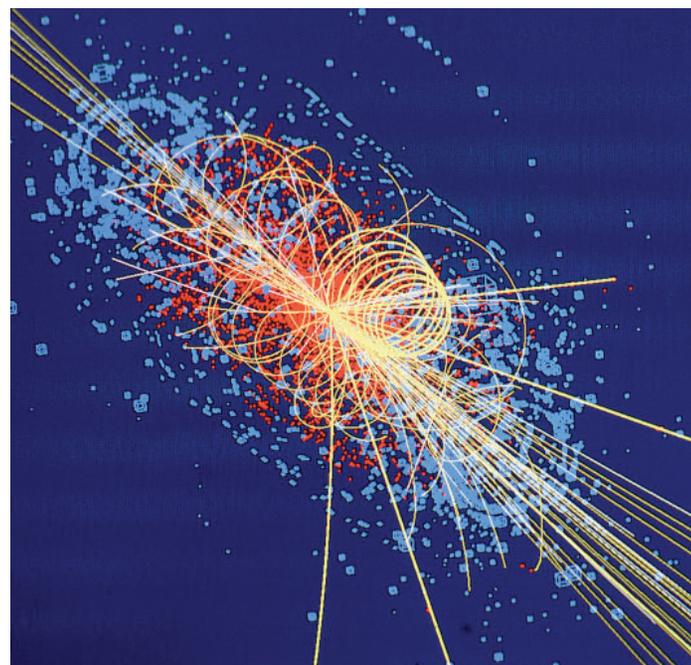
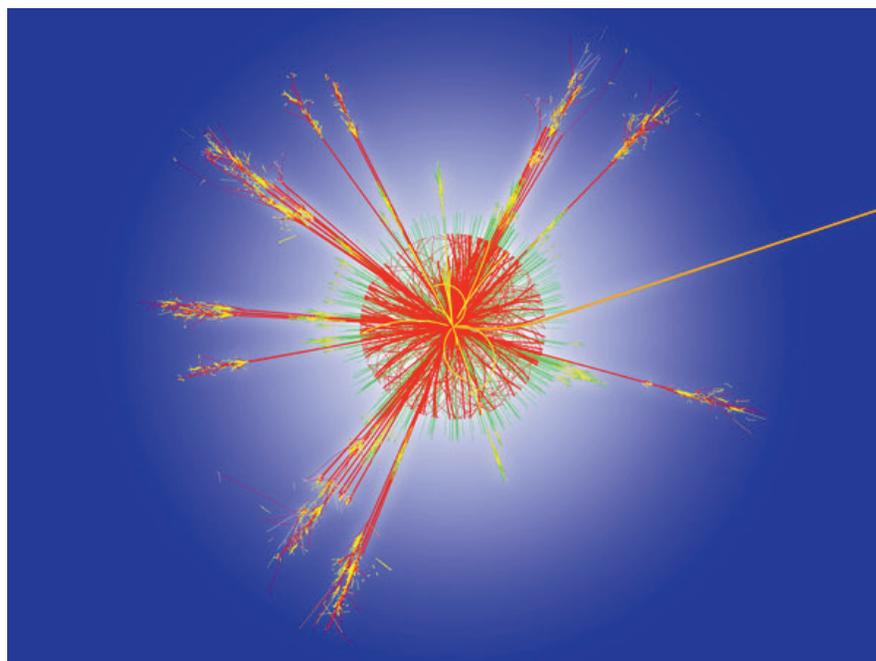


With an operating temperature of about -271 degrees Celsius, just 1.9 degrees above absolute zero, the LHC is colder than outer space.



CERN

The hottest spots in the galaxy...

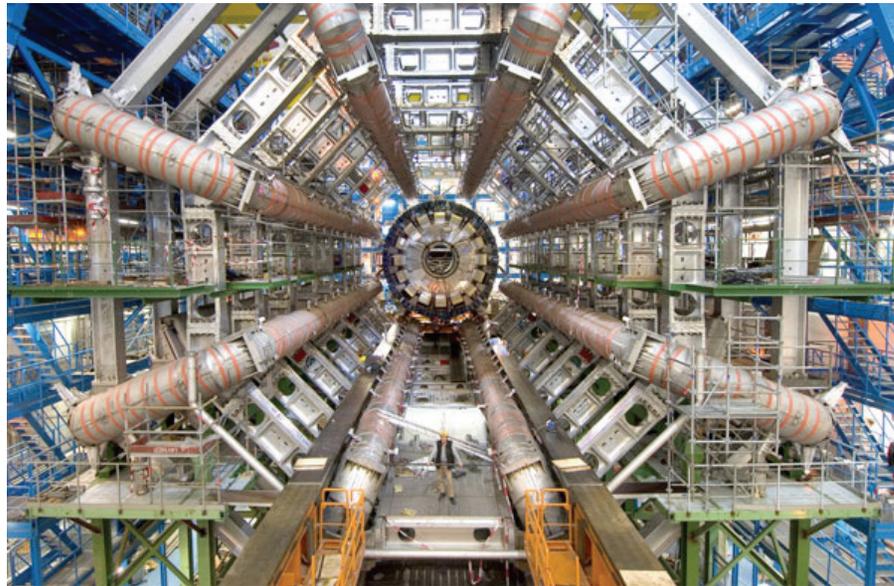


When two beams of protons collide, they will generate temperatures 1000 million times hotter than the heart of the sun, but in a minuscule space.



CERN

The biggest most sophisticated detectors ever built...



To sample and record the debris from up to 600 million proton collisions per second, scientists are building gargantuan devices that measure particles with micron precision.



CERN

The most extensive computer system in the world...



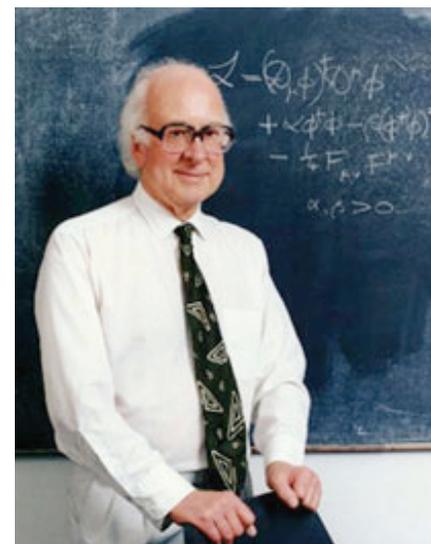
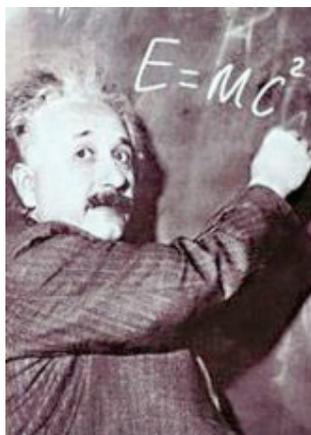
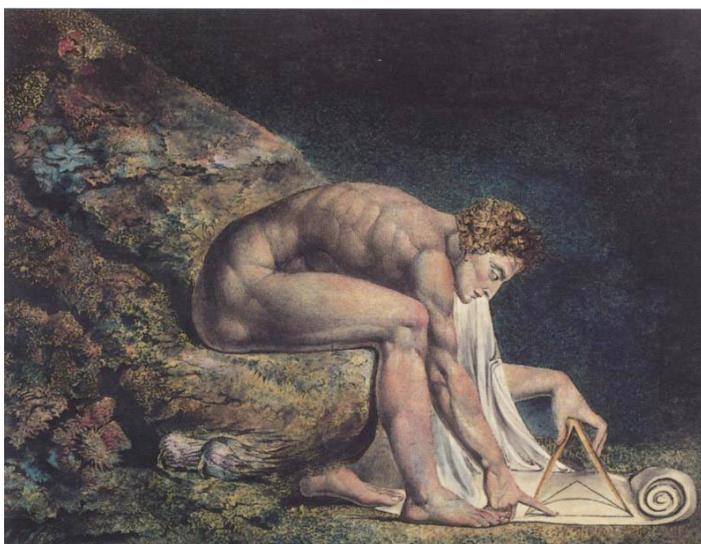
To analyse the data, tens of thousands of computers around the world are being harnessed in the Grid. The laboratory that gave the world the web, is now taking distributed computing a big step further.



CERN

Why?

To push back the frontiers of knowledge...

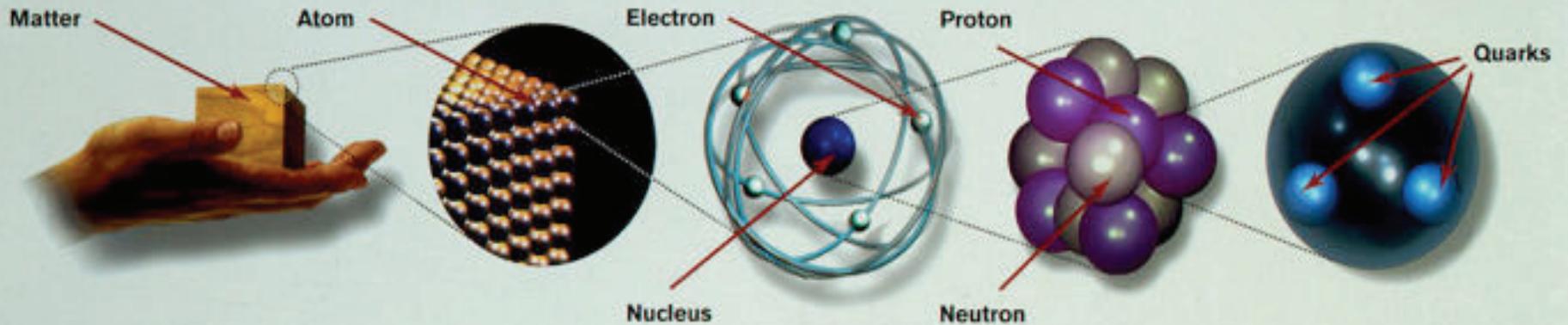


Newton's unfinished business... what is mass?

Science's little embarrassment... what is 96% of the Universe made of?

Nature's favouritism... why is there no more antimatter?

The secrets of the Big Bang... what was matter like within the first second of the Universe's life?



Matter particles

All ordinary particles belong to this group

These particles existed just after the Big Bang. Now they are found only in cosmic rays and accelerators

LEPTONS		
FIRST FAMILY	Electron Responsible for electricity and chemical reactions; it has a charge of -1	Electron neutrino Particle with no electric charge, and possibly no mass; billions fly through your body every second
SECOND FAMILY	Muon A heavier relative of the electron; it lives for two-millionths of a second	Muon neutrino Created along with muons when some particles decay
THIRD FAMILY	Tau Heavier still; it is extremely unstable. It was discovered in 1975	Tau neutrino not yet discovered but believed to exist

QUARKS		
Up Has an electric charge of plus two-thirds; protons contain two, neutrons contain one		Down Has an electric charge of minus one-third; protons contain one, neutrons contain two
Charm A heavier relative of the up; found in 1974		Strange A heavier relative of the down; found in 1964
Top Heavier still		Bottom Heavier still; measuring bottom quarks is an important test of electroweak theory

Force particles

These particles transmit the four fundamental forces of nature although gravitons have so far not been discovered

Gluons
Carriers of the **strong force** between quarks

Felt by: quarks

The explosive release of nuclear energy is the result of the **strong force**

Photons
Particles that make up light; they carry the **electromagnetic force**

Felt by: quarks and charged leptons

Electricity, magnetism and chemistry are all the results of **electro-magnetic force**

Intermediate vector bosons
Carriers of the **weak force**

Felt by: quarks and leptons

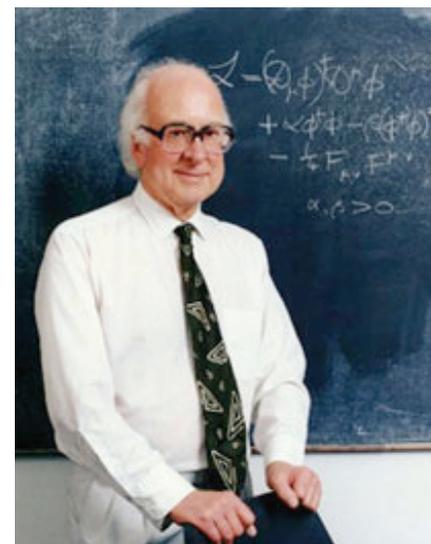
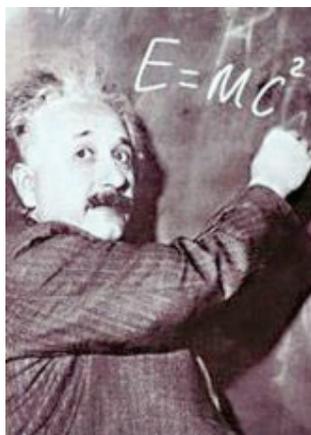
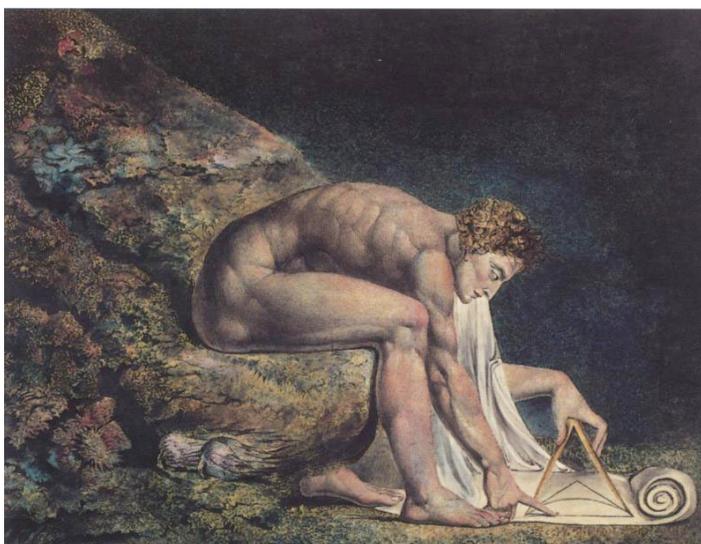
Some forms of radio-activity are the result of the **weak force**

Gravitons
Carriers of **gravity**

Felt by: all particles with mass

All the weight we experience is the result of the **gravitational force**

To push back the frontiers of knowledge...

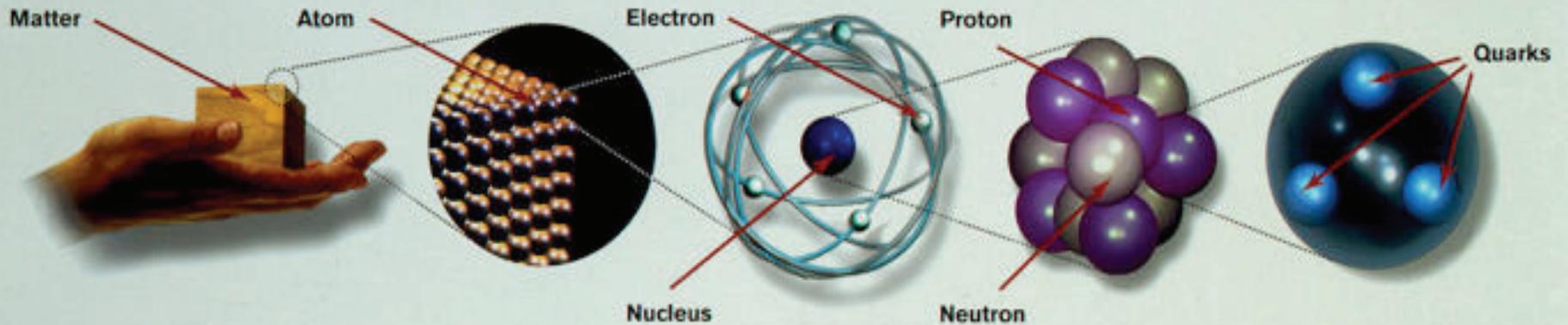


Newton's unfinished business... what is mass?

Science's little embarrassment... what is 96% of the Universe made of?

Nature's favouritism... why is there no more antimatter?

The secrets of the Big Bang... what was matter like within the first second of the Universe's life?



Matter particles
All ordinary particles belong to this group

These particles existed just after the Big Bang. Now they are found only in cosmic rays and accelerators

LEPTONS		
FIRST FAMILY	Electron Responsible for electricity and chemical reactions; it has a charge of -1	Electron neutrino Particle with no electric charge, and possibly no mass; billions fly through your body every second
SECOND FAMILY	Muon A heavier relative of the electron; it lives for two-millionths of a second	Muon neutrino Created along with muons when some particles decay
THIRD FAMILY	Tau Heavier still; it is extremely unstable. It was discovered in 1975	Tau neutrino Not yet discovered but believed to exist

QUARKS		
Up Has an electric charge of plus two-thirds; protons contain two, neutrons contain one	Down Has an electric charge of minus one-third; protons contain one, neutrons contain two	
Charm A heavier relative of the up; found in 1974	Strange A heavier relative of the down; found in 1964	
Top Heavier still	Bottom Heavier still; measuring bottom quarks is an important test of electroweak theory	

Force particles
These particles transmit the four fundamental forces of nature although gravitons have so far not been discovered

Gluons
Carriers of the **strong force** between quarks

Felt by: quarks

The explosive release of nuclear energy is the result of the **strong force**

Photons
Particles that make up light; they carry the **electromagnetic force**

Felt by: quarks and charged leptons

Electricity, magnetism and chemistry are all the results of **electro-magnetic force**

Intermediate vector bosons
Carriers of the **weak force**

Felt by: quarks and leptons

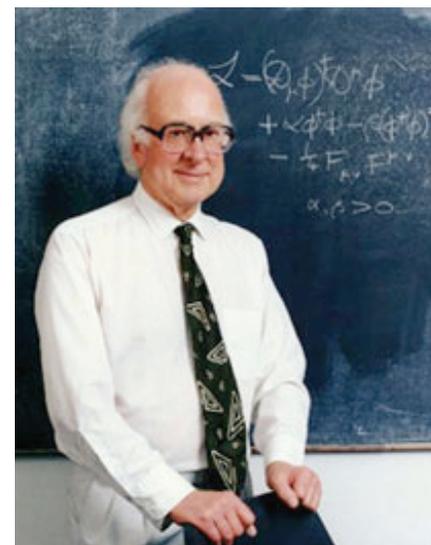
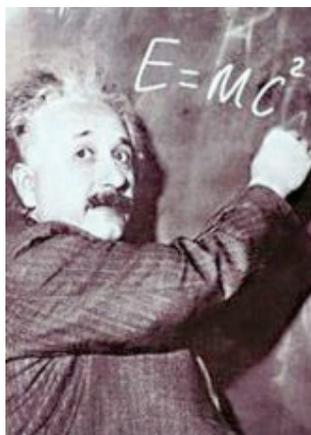
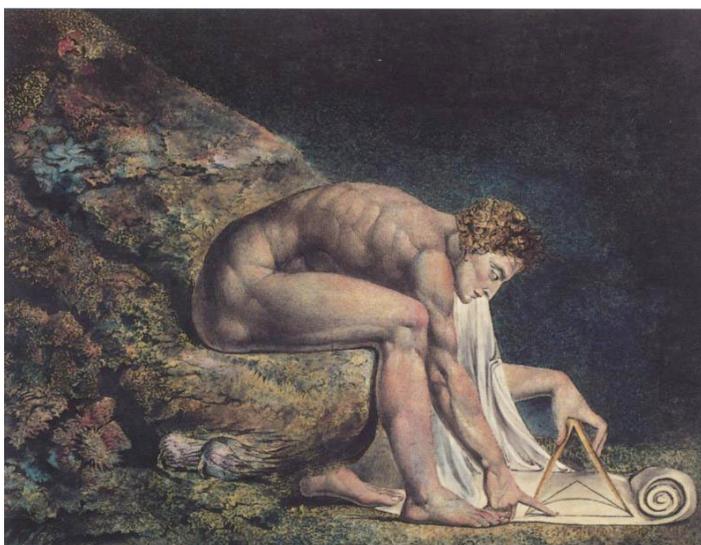
Some forms of radio-activity are the result of the **weak force**

Gravitons
Carriers of **gravity**

Felt by: all particles with mass

All the weight we experience is the result of the **gravitational force**

To push back the frontiers of knowledge...



Newton's unfinished business... what is mass?

Science's little embarrassment... what is 96% of the Universe made of?

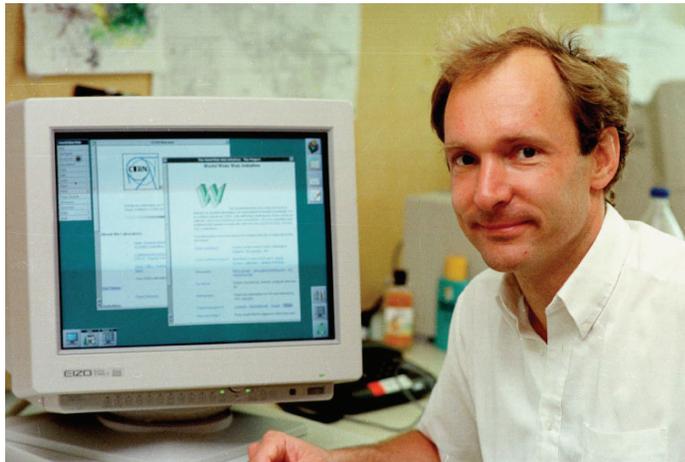
Nature's favouritism... why is there no more antimatter?

The secrets of the Big Bang... what was matter like within the first second of the Universe's life?



CERN

To develop new technologies...



Information technology - the Web and the Grid

Medicine - diagnosis and therapy

Security - scanning technologies for harbours and airports

Vacuum - new techniques for flat screen displays or solar energy devices



CERN

To unite people from different countries and cultures...



20 Member states

38 Countries with cooperation agreements

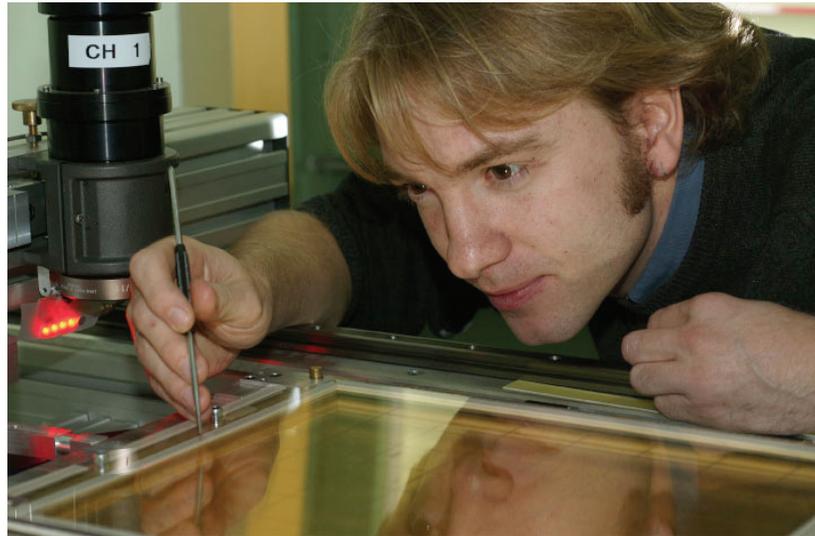
111 Nationalities

10000 People



CERN

To train the scientists and engineers of tomorrow...



From mini-Einstein workshops for five to sixes, through to professional schools in physics, accelerator science and IT, CERN plays a valuable role in building enthusiasm for science and providing formal training..



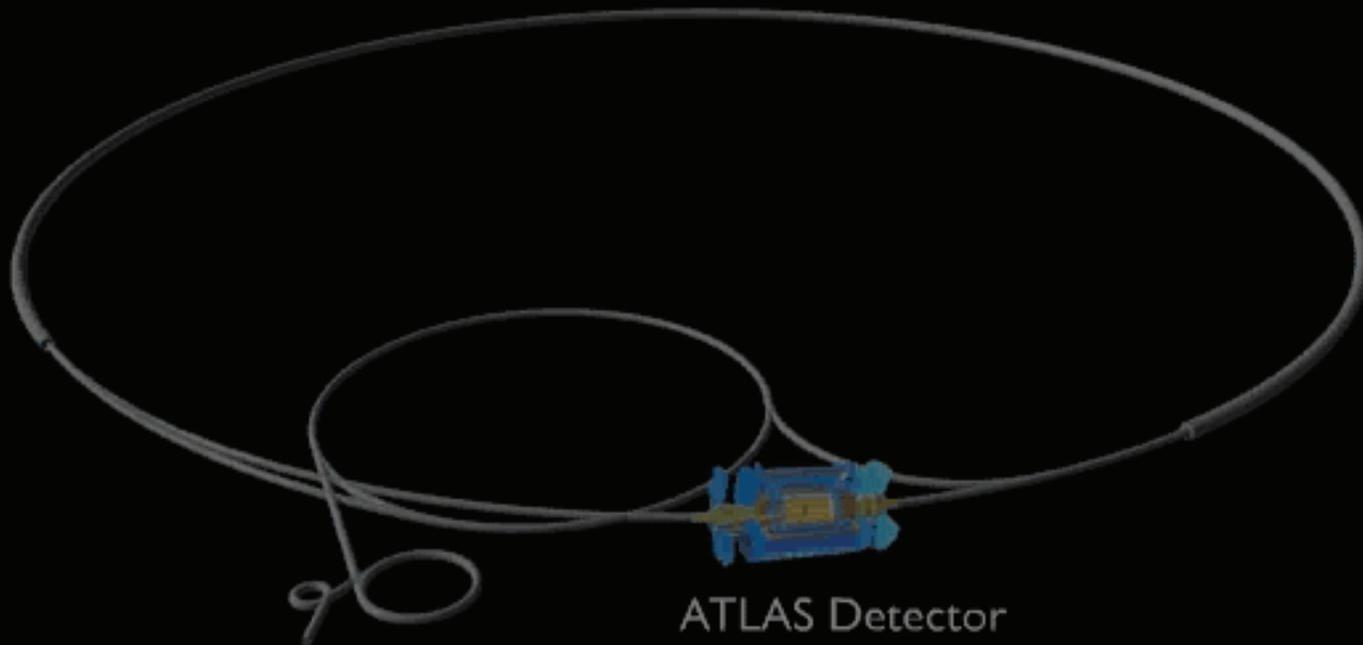
“Compact” Detectors!





PLAY ▶

Large Hadron Collider



ATLAS Detector



Outline

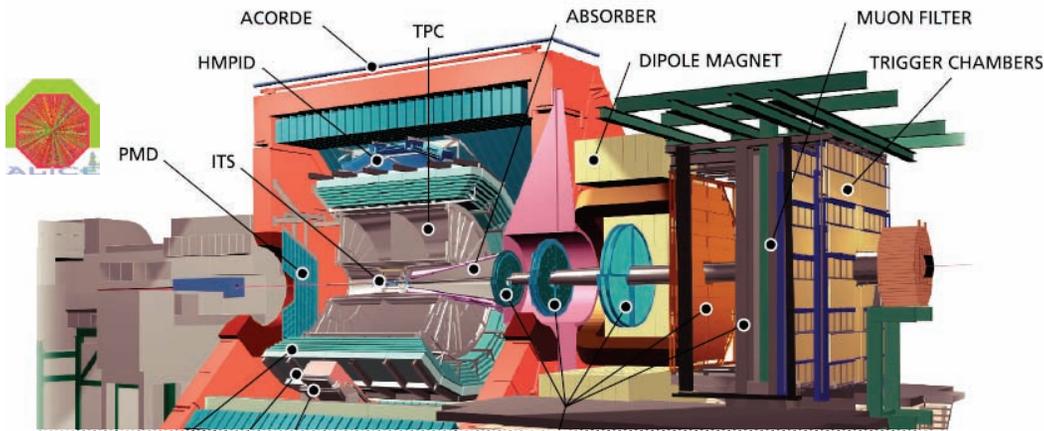
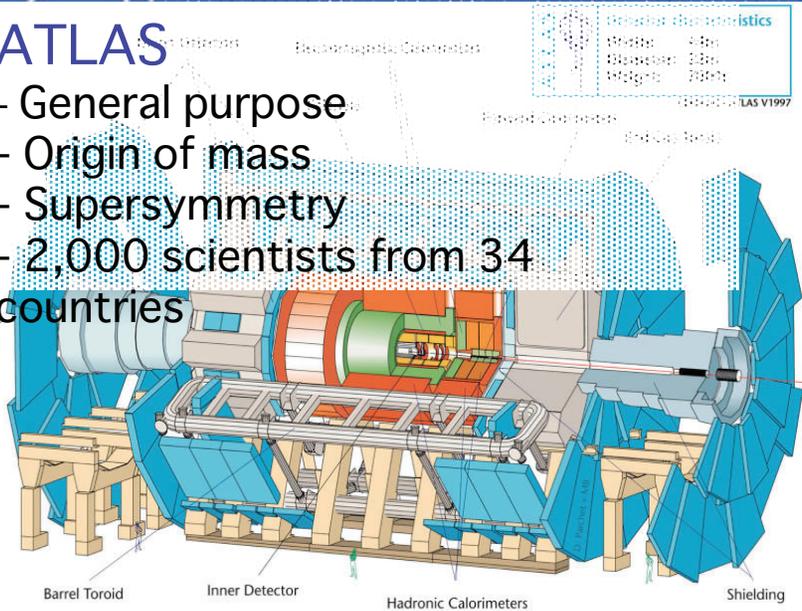
- Introduction to CERN and Experiments
- **LHC Computing**
- Challenges
- Summary/Conclusion



The Four LHC Experiments...

ATLAS

- General purpose
- Origin of mass
- Supersymmetry
- 2,000 scientists from 34 countries

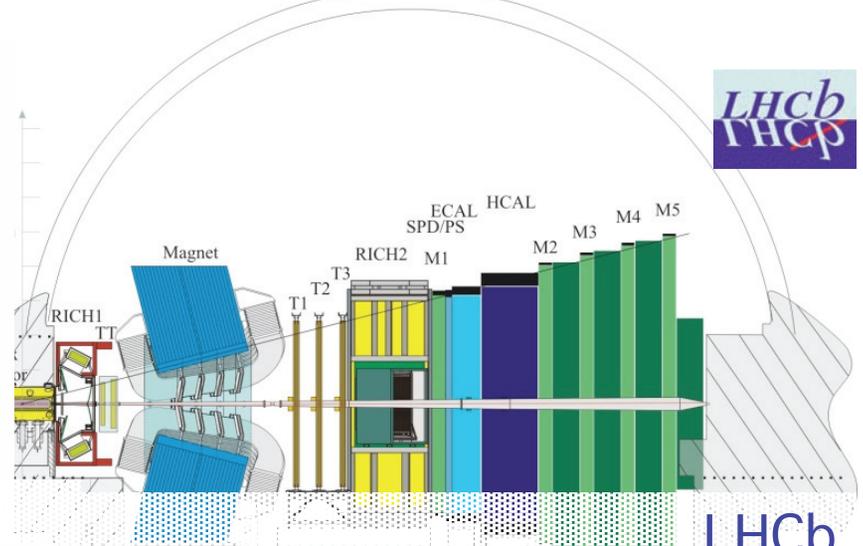
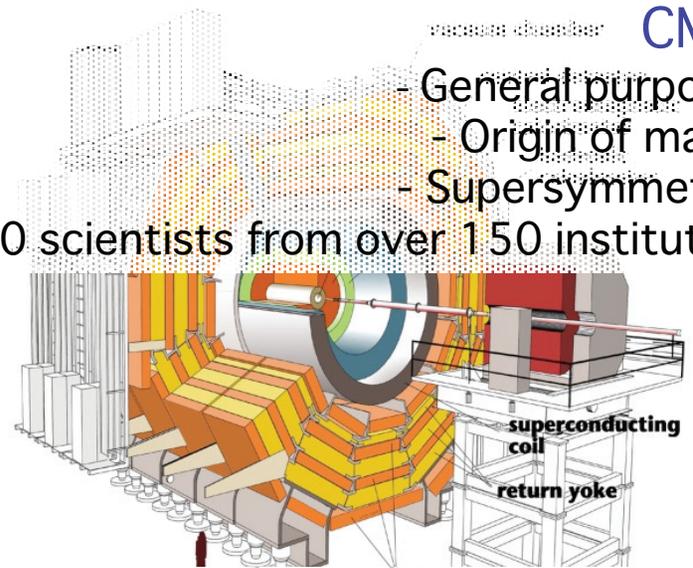


ALICE

- heavy ion collisions, to create quark-gluon plasmas
- 50,000 particles in each collision

CMS

- General purpose
- Origin of mass
- Supersymmetry
- 1,800 scientists from over 150 institutes



LHCb

- to study the differences between matter and antimatter



... generate lots of data ...



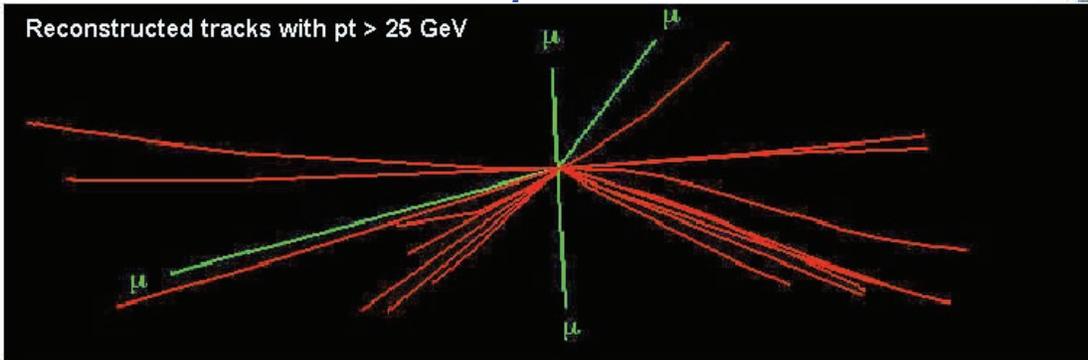
The accelerator generates 40 million particle collisions (events) every second at the centre of each of the four experiments' detectors





... generate lots of data ...

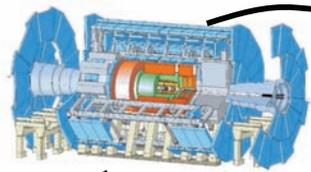
reduced by online computers to a few hundred “good” events per second.



Which are recorded on disk and magnetic tape at 100-1,000 MegaBytes/sec → **~15 PetaBytes per year**

for all four experiments

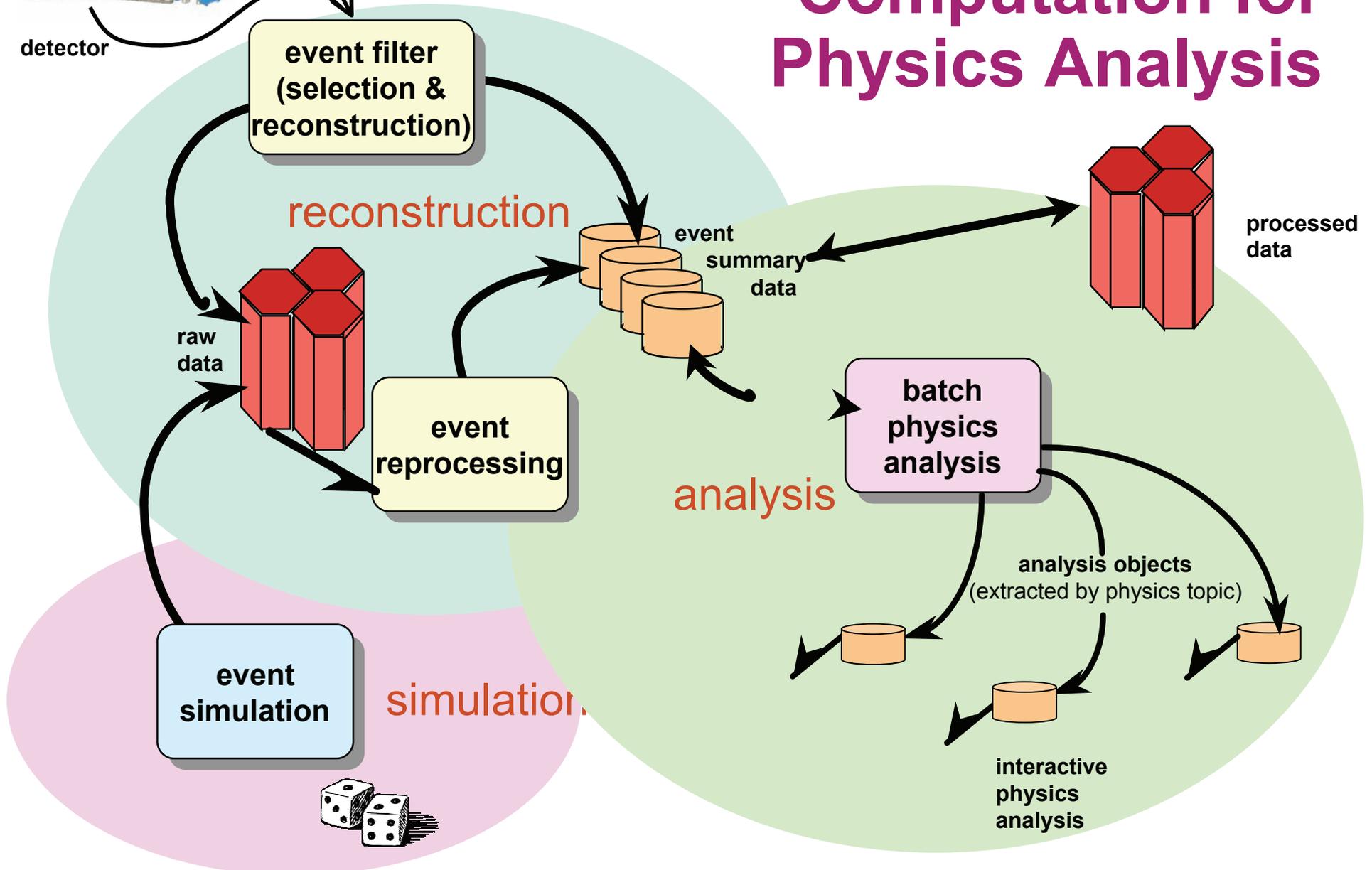




detector



Data Handling and Computation for Physics Analysis





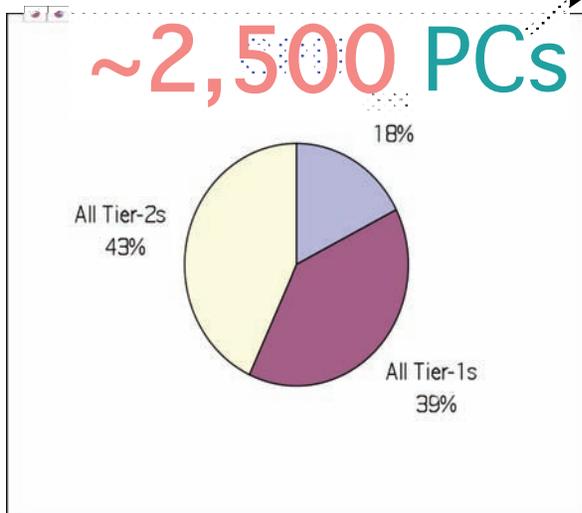
... leading to a high box count

Summary of Computing Resource Requirements

All experiments - 2008

From LCG TDR - June 2005

	CERN	All Tier-1s	All Tier-2s	Total
CPU (MSPECint2000s)	25	56	61	142
Disk (PetaBytes)	7	31	19	57
Tape (PetaBytes)	18	35		53



Another ~1,500 boxes



Outline

- Introduction to CERN and Experiments
- LHC Computing
- **Challenges**
- Summary/Conclusion



Outline

- Introduction to CERN and Experiments
- LHC Computing
- **Challenges**
 - _ Capacity Provision
 - _ Box Management
 - _ Data Management and Distribution
 - _ What' s Going On?
- Summary/Conclusion



Outline

- Introduction to CERN and Experiments
- LHC Computing
- **Challenges**
 - _ Capacity Provision
 - _ Box Management
 - _ Data Management and Distribution
 - _ What' s Going On?
- Summary/Conclusion



Computing Service Hierarchy

Tier-0 – the accelerator centre

- § Data acquisition & initial processing
- § Long-term data curation
- § Distribution of data Tier-1 centres



Canada – Triumf (Vancouver)
France – IN2P3 (Lyon)
Germany – Forschungszentrum Karlsruhe
Italy – CNAF (Bologna)
Netherlands – NIKHEF/SARA (Amsterdam)
Nordic countries – distributed Tier-1
Spain – PIC (Barcelona)
Taiwan – Academia Sinica (Taipei)
UK – CLRC (Oxford)
US – FermiLab (Illinois)
– Brookhaven (NY)

Tier-1 – “online” to the data acquisition process high availability

- § Managed Mass Storage
- § Data-heavy analysis
- § National, regional support

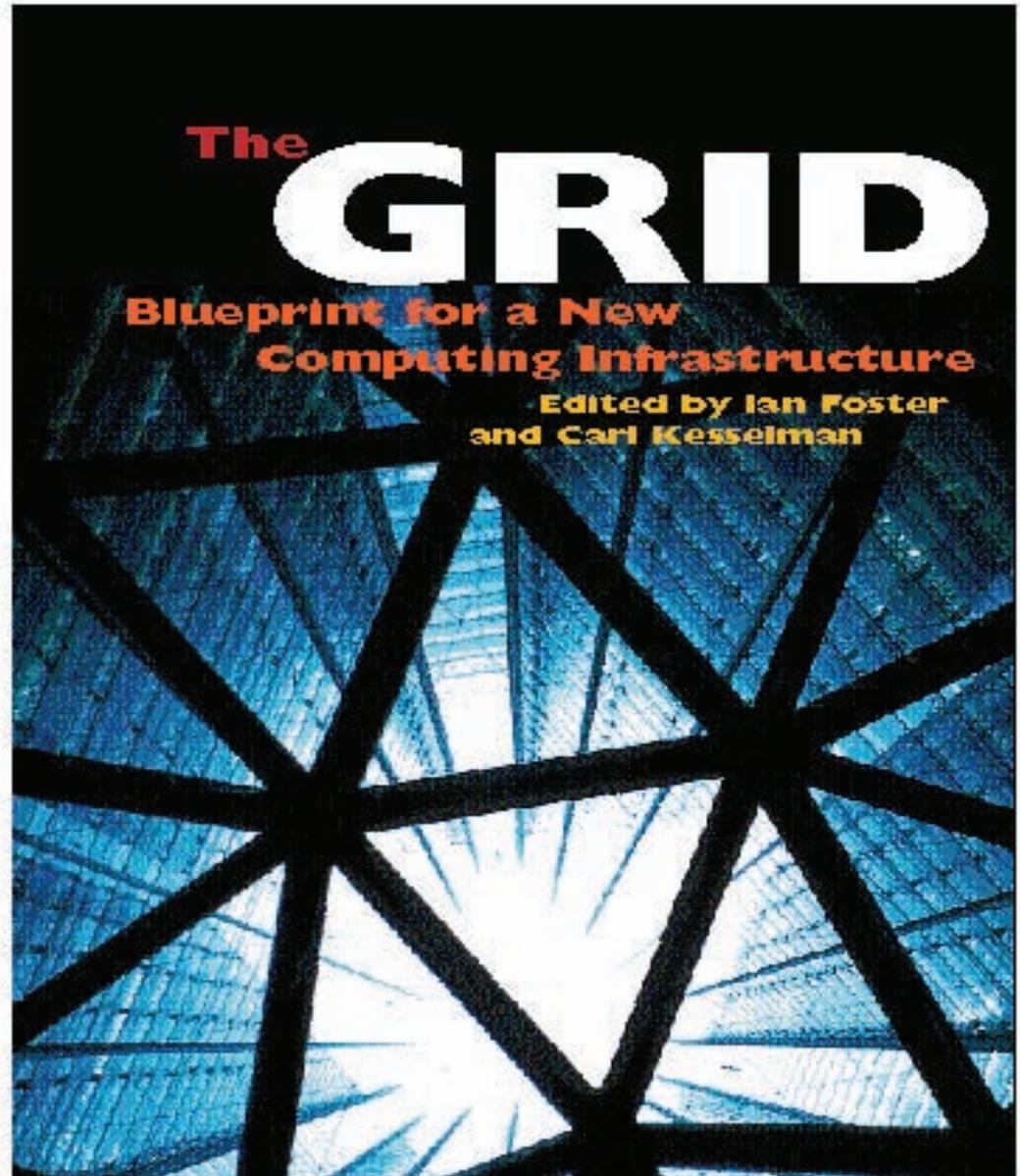
Tier-2 – ~100 centres in ~40 countries

- § Simulation
- § End-user analysis – batch and interactive



The Grid

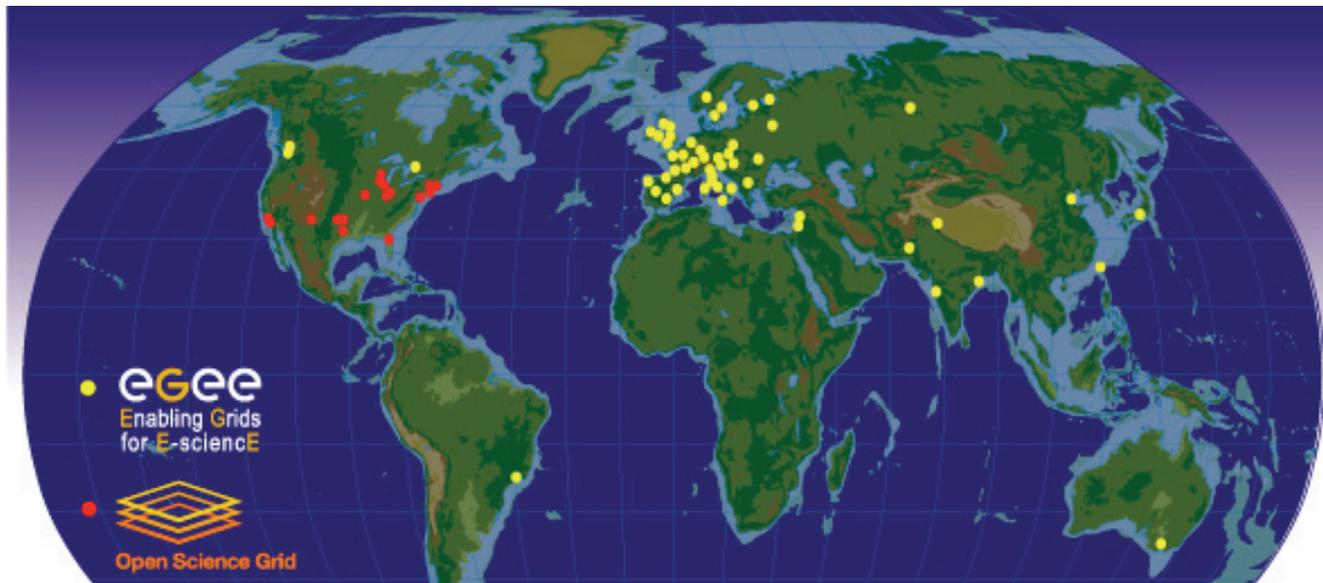
- Timely Technology!
- **Deploy** to meet LHC computing needs.
- Challenges for the **Worldwide LHC Computing Grid Project** due to
 - _ worldwide nature
 - competing middleware...
 - _ newness of technology
 - competing middleware...
 - _ scale





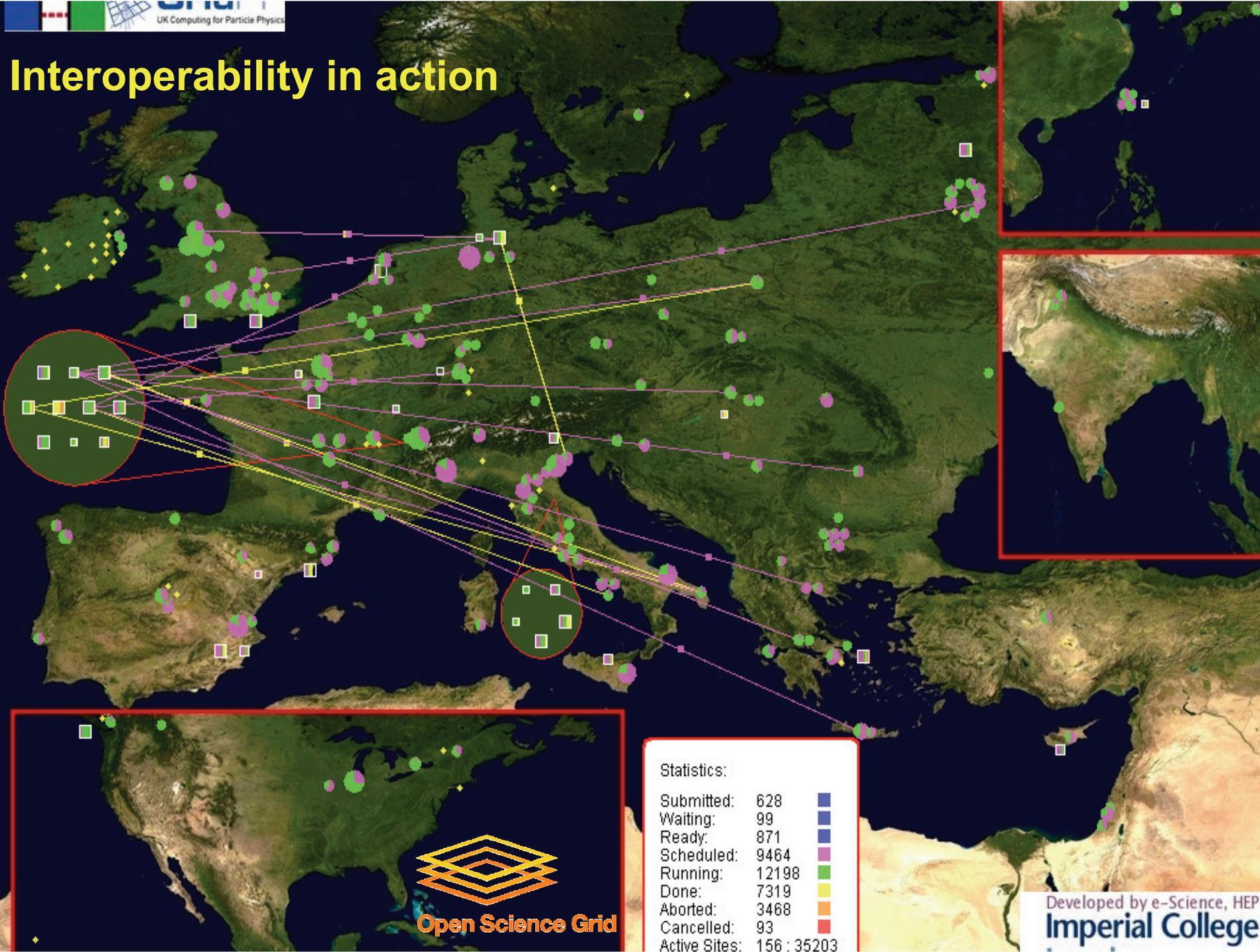
Interoperation between Grid Infrastructures

- WLCG service relies on three Grid infrastructures: EGEE, OSG and NorduGrid
- Interoperability required (and achieved) for
 - _ users (job submission)
 - _ administration (identity, monitoring, accounting, ...)



A map of the worldwide LCG infrastructure operated by EGEE and OSG.

Interoperability in action





- TECHNOLOGY
 - Infrastructure
 - Network
 - Security
 - Client
 - Server
 - Mobile
 - Operating Systems
 - Data Center
 - Applications
 - Development
 - Architecture
- LEADERSHIP

CIO Resource Alerts
GET NOTIFIED!

NEWSLETTERS

CIO.com updates, insights and advice on technology, management and your career.

- Advice and Opinion
- CIO Consumer IT
- CIO Leader
- CIO Enterprise
- CIO Insider

More Newsletters | Edit Profile

enter e-mail **SIGN-UP**

Feature

Seven Wonders of the IT World

The fastest supercomputer. The most intriguing data center. The constantly changing core at the heart of Linux. Take a tour of the most impressive and most unusual marvels of the IT world.

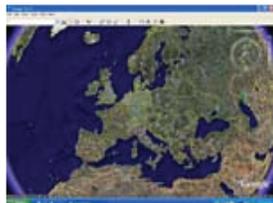
Leave a comment (65)

By C.G. Lynch

PAGE 4

World's largest scientific grid computing project: The E-science II (EGEE-II) project

Launched: September 2006, for use by scientists around the world.



A Google Earth view of European sites hooked into the EGEE grid computing project

Helps power: Large-scale scientific research projects in fields from geology to chemistry—for example, will analyze data from CERN's Large Hadron Collider, a particle accelerator being built to help investigate details around the Big Bang and related physics questions.

Amount of work it does: 98,000 jobs a day, more than 1 million per month.

Juggling ability: Runs about 30,000 jobs concurrently, on average.

RELATED SOLUTIONS

How-To »

- How to Lay Dead Technologies to Rest
- Ghostly Gear: Technology Tools for Paranormal Investigations

Research & Analysis »

- Lax Laptop Security Can Be Dangerous...and Expensive
- Reduce Information Technology Complexity, Costs with Consolidation

Advice & Opinion »

- Join the Conversation!
- Ask, Answer and Interact with the CIO Community

News »

- Nokia Makes the Ferrari of Phones...Literally
- Holiday Gift Guide 2007: Best Technology Bling

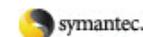
VIDEO

Business Innovation Video Series



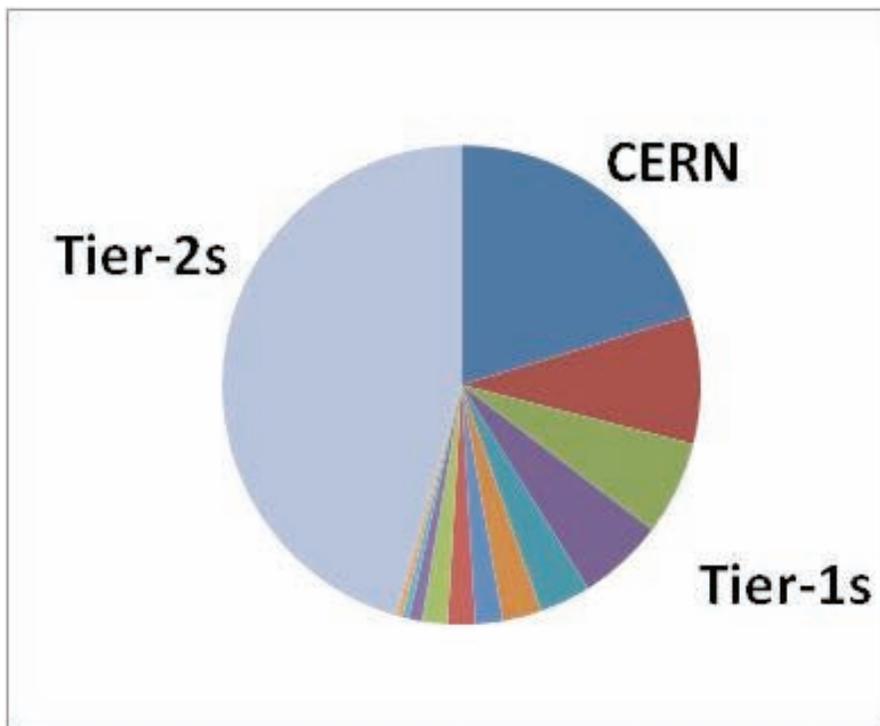
IT Leaders are discussing how IT is becoming part of the innovation cycle.

Watch the videos »

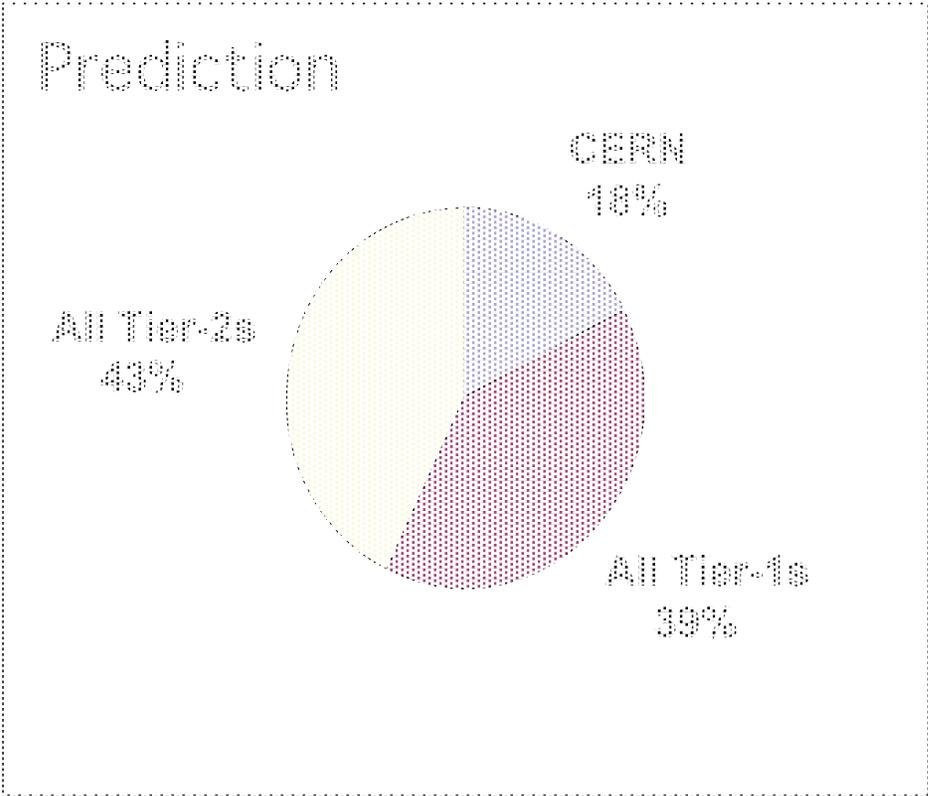




Accounted CPU Usage



80 Tier-2s 45%
11 Tier-1s 35%
CERN 20%



530M SI2K-days/month (CPU)

9 PB disk at CERN + Tier-1s



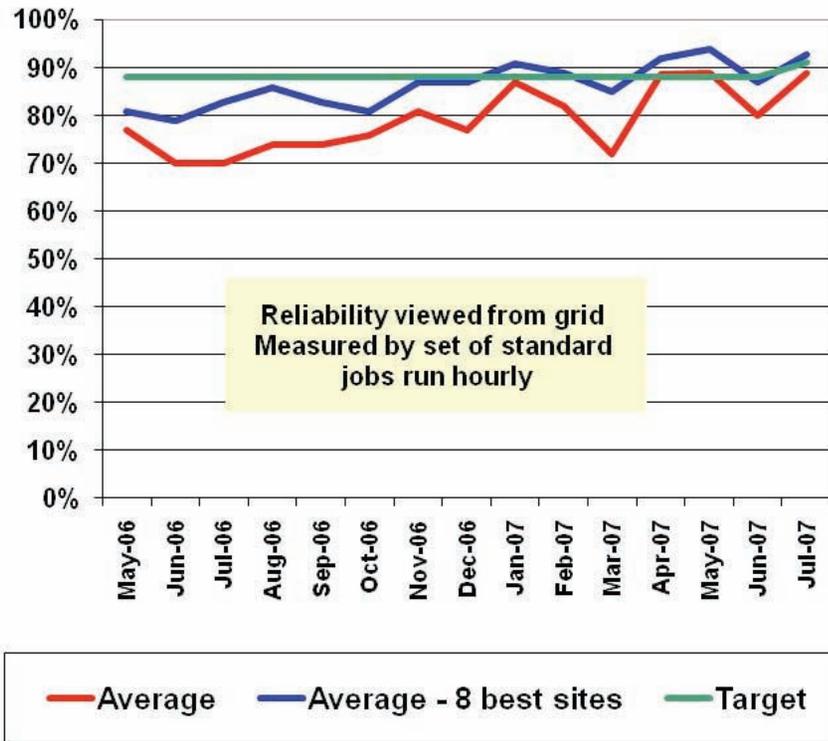
Remaining Challenges

- Creating a working Grid service across multiple infrastructure is clearly a success, but challenges remain
 - _ Reliability

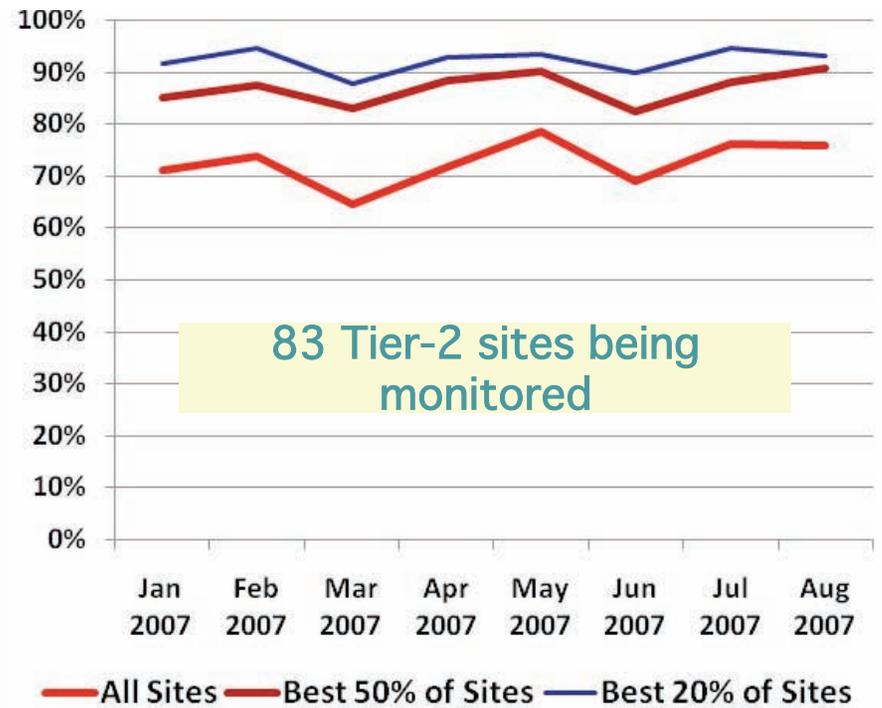


Reliability

Site Reliability CERN + Tier-1s



Site Reliability Tier-2 Sites



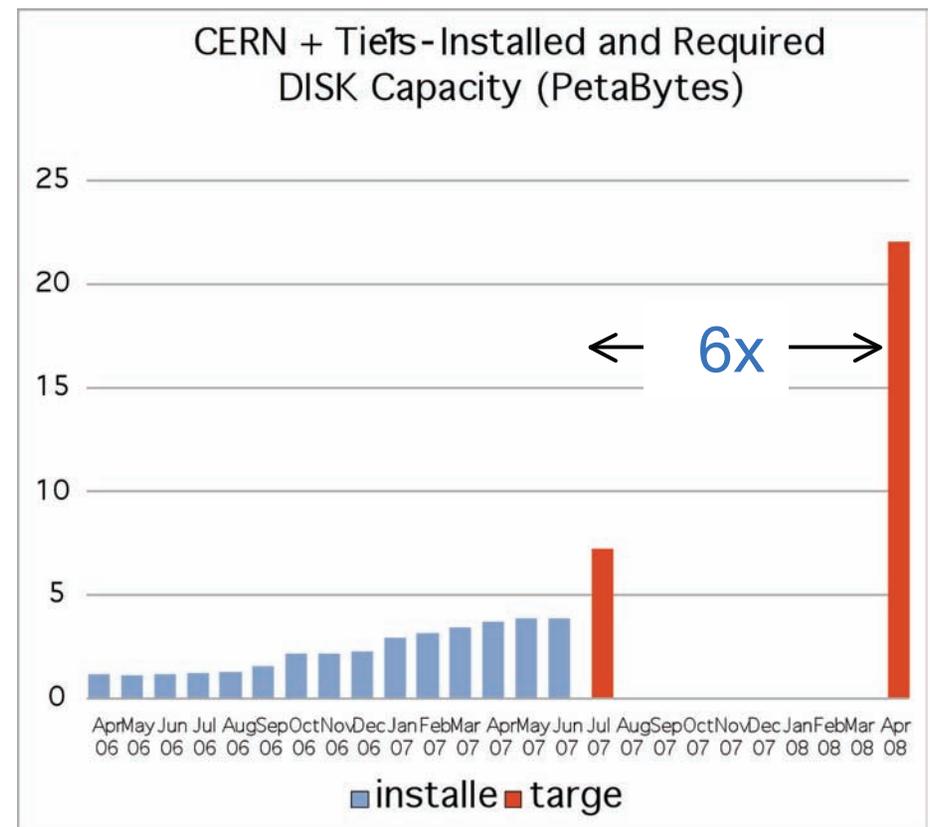
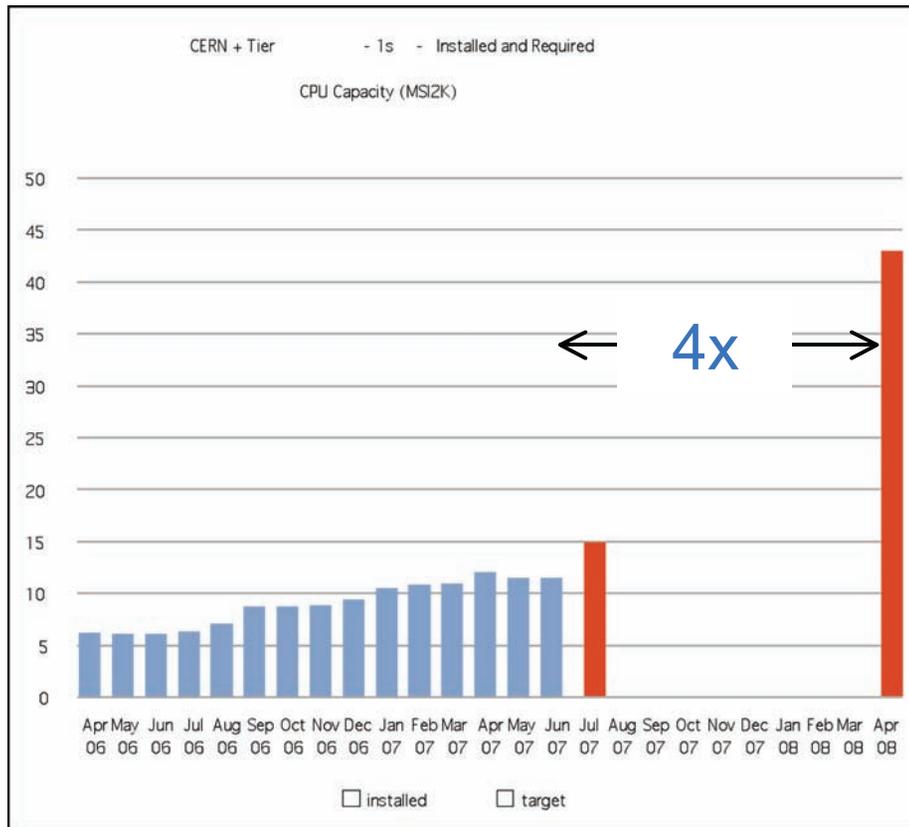


Remaining Challenges

- Creating a working Grid service across multiple infrastructure is clearly a success, but challenges remain
 - _ Reliability
 - _ Ramp-up



Steep ramp-up still needed before first physics run



Evolution of installed capacity from April 06 to June 07
Target capacity from MoU pledges for 2007 (due July07)
and 2008 (due April 08)



Remaining Challenges

- Creating a working Grid service across multiple infrastructure is clearly a success, but challenges remain
 - _ Reliability
 - _ Ramp-up
 - _ Collaboration
 - From computer centre empires to a federation
 - consensus rather than control
 - _ ...



Outline

- Introduction to CERN and Experiments
- LHC Computing
- **Challenges**
 - _ Capacity Provision
 - _ **Box Management**
 - _ Data Management and Distribution
 - _ What' s Going On?
- Summary/Conclusion

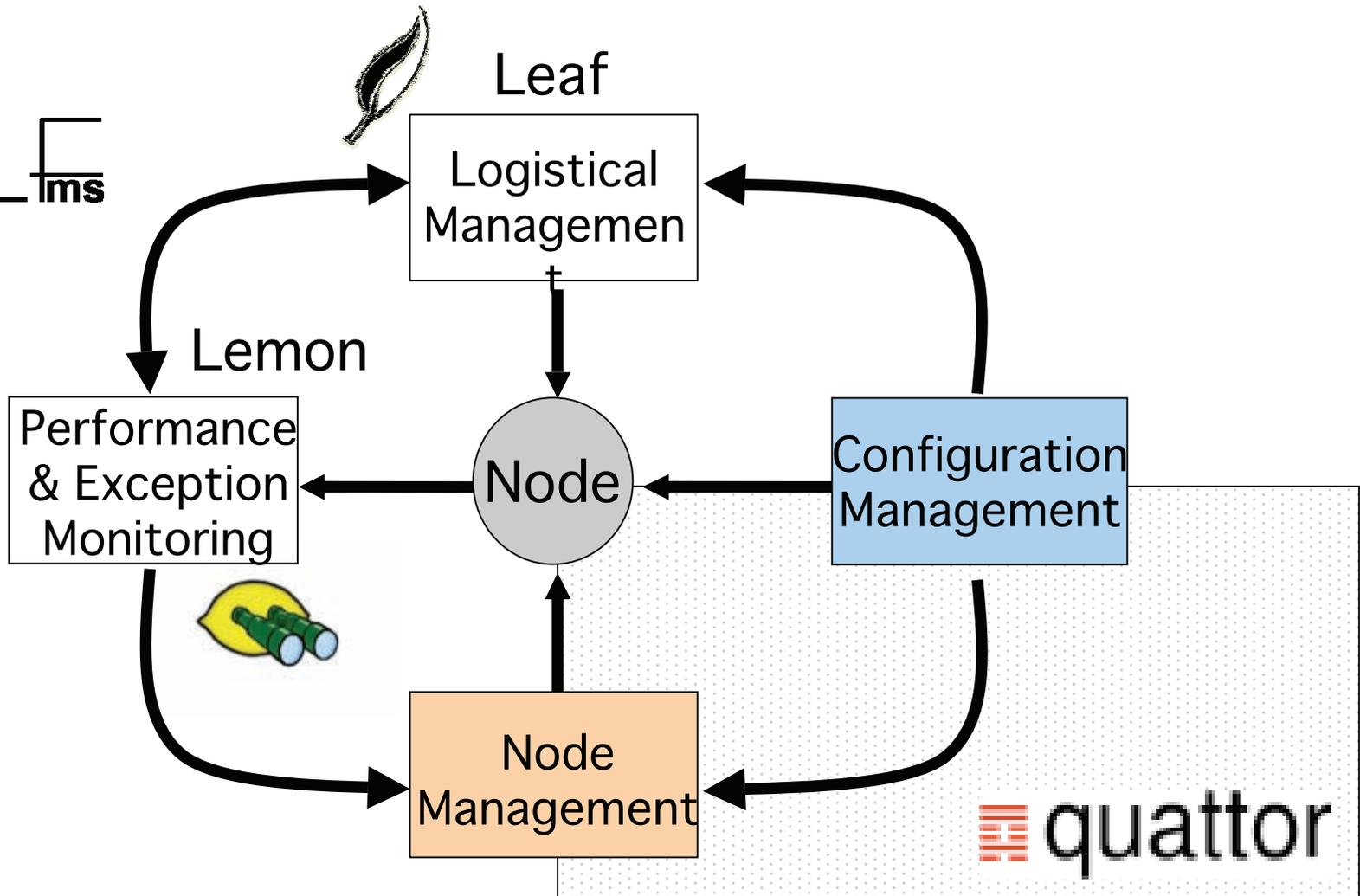
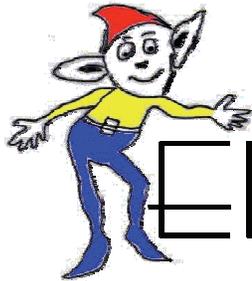


Outline

- Introduction to CERN and Experiments
- LHC Computing
- **Challenges**
 - _ Capacity Provision
 - _ **Box Management**
 - Installation & Configuration
 - Monitoring
 - Workflow
 - _ Data Management and Distribution
 - _ What's Going On?
- Summary/Conclusion



ELFms Vision



Toolkit developed by CERN in collaboration with many HEP sites and as part of the European DataGrid Project.
See <http://cern.ch/ELFms>



Why a bespoke system?

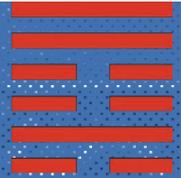
- Commercial Management Suites
 - _ (Full) Linux support rare (5+ years ago...)
 - _ Much work needed to deal with specialist HEP applications; insufficient reduction in staff costs to justify license fees.
- Open Source Systems
 - _ Many packages with interesting features, but none featuring all of items considered essential
 - Declarative, hierarchical configuration specification permitting validation, integrated software distribution and configuration management, separation of configuration data and code, feedback loop to avoid configuration drift, ability to update running systems

See EDG/WP4 report on current technology (http://cern.ch/hep-proj-grid-fabric/Tools/DataGrid-04-TED-0101-3_0.pdf) or "Framework for Managing Grid-enabled Large Scale Computing Fabrics" (<http://cern.ch/quattor/documentation/poznanski-phd.pdf>) for reviews of various packages.



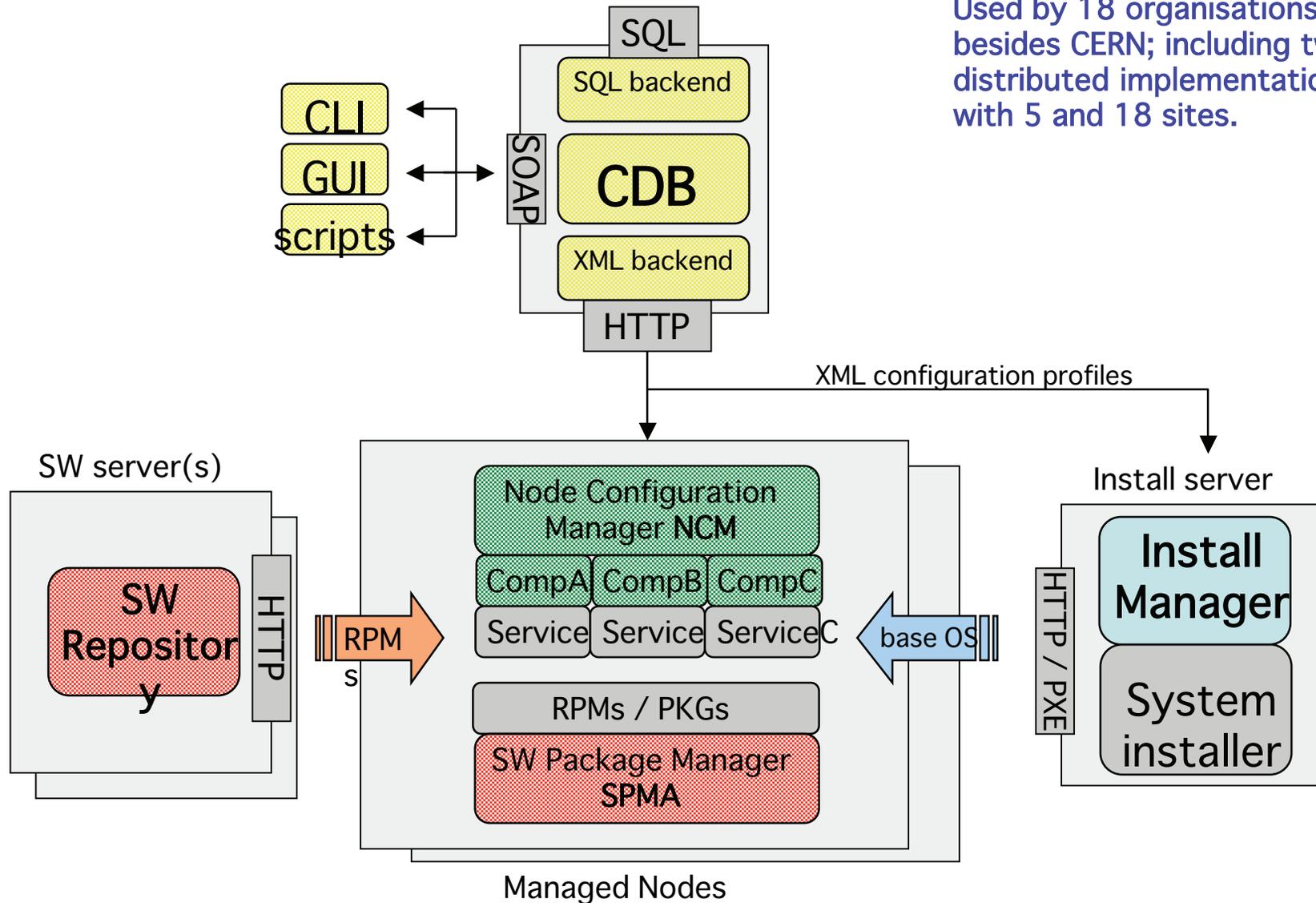
Outline

- Introduction to CERN and Experiments
- LHC Computing
- **Challenges**
 - _ Capacity Provision
 - _ **Box Management**
 - Installation & Configuration
 - Monitoring
 - Workflow
 - _ Data Management and Distribution
 - _ What's Going On?
- Summary/Conclusion



Configuration server

Used by 18 organisations besides CERN; including two distributed implementations with 5 and 18 sites.





Configuration Hierarchy

http://lemonweb.cern.ch/lemon-status/tpl_view.php?profile=pro_type_lxbatch_slc3

Tony Cass

Lemon Monitoring Web Pages - C...

Home Documentation Alarms PCFinder Metrics Error trending Help

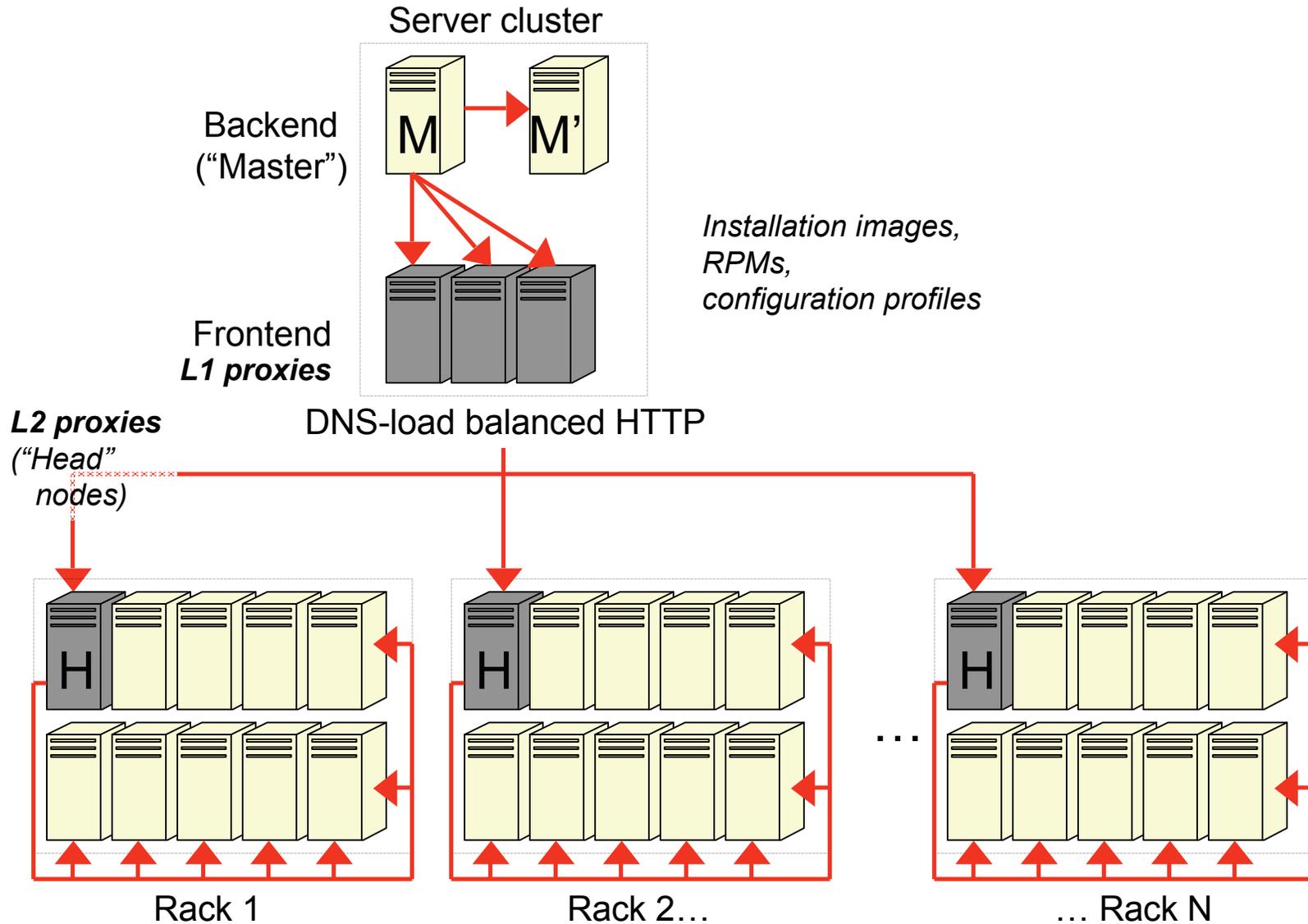
Template info: pro_type_lxbatch_slc3 17 Nov 2006 Fri 18:20:22

CERN Monitoring

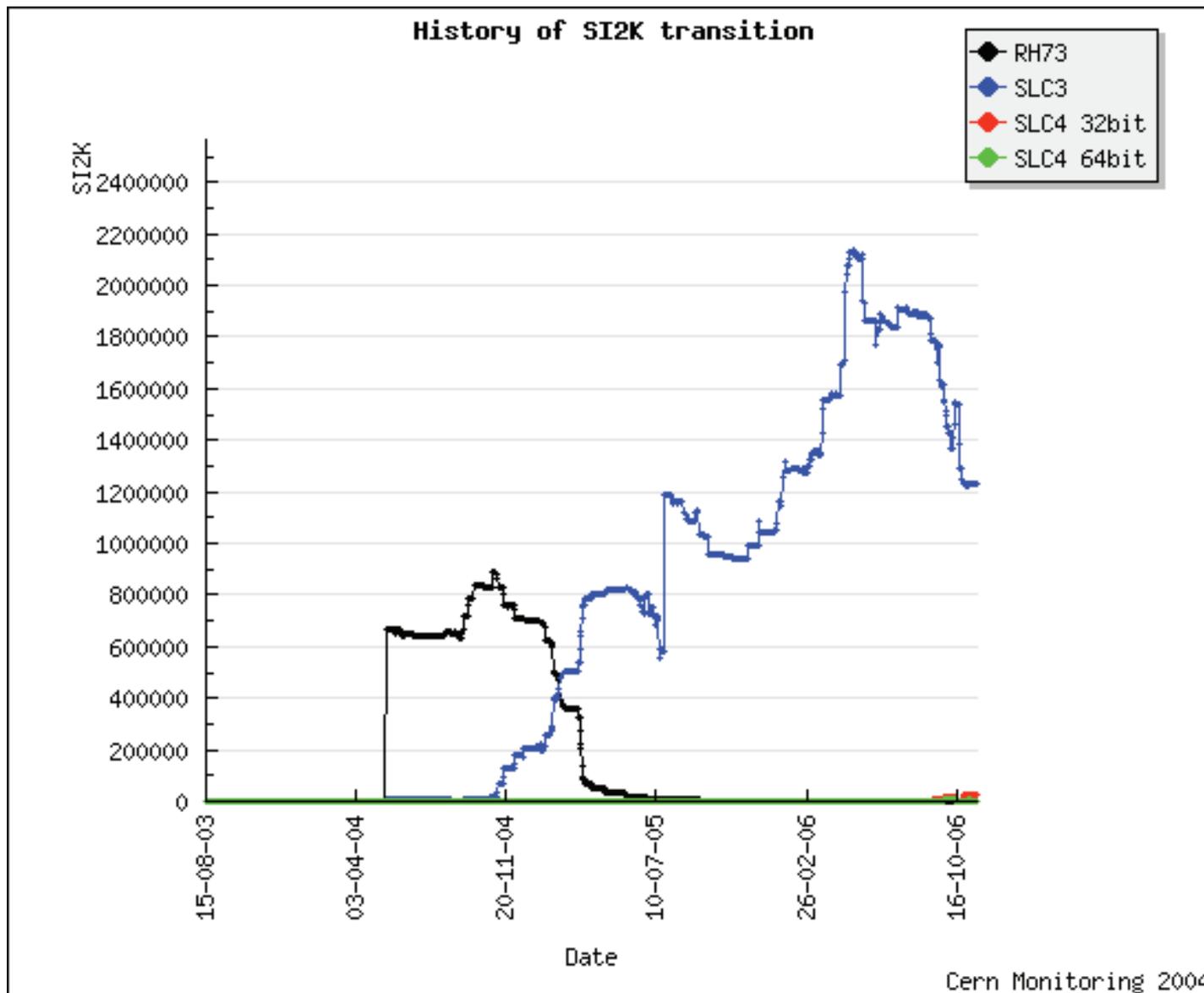
```
#####  
#  
# template pro_type_lxbatch_slc3  
#  
# RESPONSIBLE: Thorsten Kleinwort  
#  
#####  
template pro_type_lxbatch_slc3;  
  
include pro_software_components_slc3;  
include pro_system_lxbatch;  
include pro_os_slc3;  
  
"/system/cluster/tplname" = "pro_type_lxbatch_slc3";  
  
#  
# Yaim for gLite 3.0  
#  
include pro_software_components_lcg_yaim_3_0;  
  
"/software/components/yaim/active" = true;  
"/software/components/yaim/nodetype/glite-WN" = true;  
"/software/components/yaim/configure" = true; # Do automatically configure YAIM  
"/system/cluster/subname" = "public";  
"/system/accounting/name" = "share";  
  
#  
# SPMA proxy configuration  
#  
# use head node as proxy server  
"/software/components/spma/headnode"=true;  
# active SPMA proxy  
"/software/components/spma/proxy"="yes";  
  
"/software/components/lsfclient/lstype" = lstype() ;
```



Scalable s/w distribution...



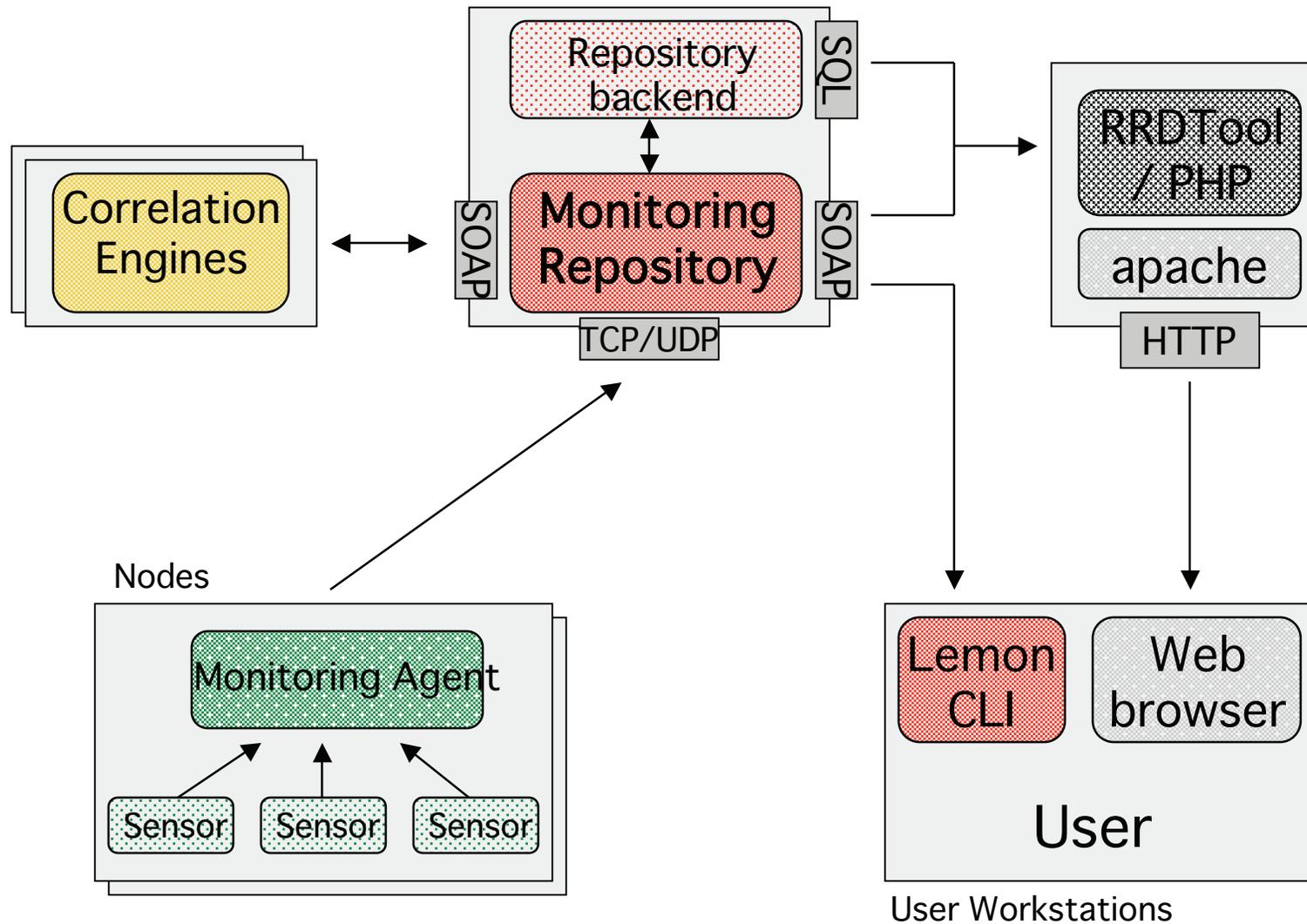
... in practice!





Outline

- Introduction to CERN and Experiments
- LHC Computing
- **Challenges**
 - _ Capacity Provision
 - _ **Box Management**
 - Installation & Configuration
 - Monitoring
 - Workflow
 - _ Data Management and Distribution
 - _ What's Going On?
- Summary/Conclusion





What is monitored

- Node based Lemon sensors cover all the usual system parameters and more
 - _ system load, file system usage, network traffic, daemon count, software version...
 - _ SMART monitoring for disks
 - _ Oracle monitoring
 - number of logons, cursors, logical and physical I/O, user commits, index usage, parse statistics, ...
 - _ AFS client monitoring
 - _ ...
- It is also possible to provide “non-node” sensors. At CERN these allow integration of
 - _ information from the building management system
 - Power demand, UPS status, temperature, ...
 - _ and full feedback is possible (although not implemented): e.g. system shutdown on power failure
 - _ high level mass-storage and batch system details
 - Queue lengths, file lifetime on disk, ...
 - _ hardware reliability data



Monitoring displays

Browser address bar: http://lemonweb.cern.ch/lemon-status/pdu_units.php?time=3&pdu_units=EOD53*43&detailed=&auto

Page Title: Lemon Monitoring Web Pages - P...

CERN Monitoring

kW - last year

■ kW aver: 68.48 max: 76.95 min: 39.52 curr: 72.23

kW - last year

■ kW1	aver: 22.85	max: 25.82	min: 12.44	curr: 23.37
■ kW2	aver: 19.53	max: 22.50	min: 11.18	curr: 19.82
■ kW3	aver: 26.10	max: 29.40	min: 15.89	curr: 29.03

kVA - last year

■ kVA aver: 71.99 max: 80.32 min: 43.27 curr: 75.52

kVA - last year

■ kVA1	aver: 23.91	max: 26.81	min: 13.57	curr: 24.32
■ kVA2	aver: 20.57	max: 23.58	min: 12.33	curr: 20.72
■ kVA3	aver: 27.52	max: 30.83	min: 17.36	curr: 30.48

A - last year

■ A aver: 40.59 max: 45.50 min: 30.52 curr: 45.50

A - last year

■ A1	aver: 105.12	max: 117.76	min: 59.53	curr: 106.71
■ A2	aver: 90.34	max: 103.60	min: 54.14	curr: 90.66
■ A3	aver: 120.59	max: 135.31	min: 75.93	curr: 133.83

Last

Search host: Virtual Clusters | Clusters | Racks | HMM models | Databases | Virtual Organizations | Power Distribution



Dynamic cluster definition

- As LHC goes into operation, the quality of the data will be a key specification. The way in which we monitor the system will be crucial. Lemmonfly” - e.g. ATLAS

LEMON Monitoring

Home | Documentation | Alarms | PCFinder | Metrics | Error trending | Help

Vo info: vo_atlas 19 Nov 2006 Sun 12:23:26

Virtual Organization Information

# of hosts (down):	325 (14)
operating system(s):	2.4.21-47.EL.cernsmp, 2.6.9-42.0.2.EL.1.cernsmp, ▶
# of CPUs (down):	656 (26)
average up time:	52 days, 23h:03m (boots per host)
hosts down:	lxb1326, lxb0471 ▶
exceptions:	ATD_WRONG, SPMA_ERROR, ▶
ITCM history	
Select from hosts:	None ▼
Metric Distributions	Correlations

Load Percentages

44.9%

0-0.5	44.9%
0.5-1.0	10.8%
1.0-2.0	12.3%
> 2.0	27.7%
down	4.3%

CPU utilization - last day

User CPU aver:	390.50m	max: 9567.51m	min: 110.89m	curr: 121.03m
System CPU aver:	1178.57m	max: 8015.01m	min: 667.16m	curr: 1509.13m
Nice CPU aver:	47280.61m	max: 50101.48m	min: 41057.70m	curr: 42895.94m
Idle CPU aver:	50504.46m	max: 56486.85m	min: 39921.95m	curr: 54742.38m
IO wait CPU aver:	525.01m	max: 4687.60m	min: 224.10m	curr: 607.26m
IRQ CPU aver:	19.06m	max: 93.60m	min: 14.88m	curr: 17.78m
Soft IRQ CPU aver:	102.70m	max: 152.51m	min: 84.84m	curr: 107.53m

Network utilization - last day

eth0 in	7.94M	max: 29.03M	min: 11.24M	curr: 11.64M
eth0 out	3.43M	max: 11.48M	min: 1.05M	curr: 2.48M

Search host: Virtual Clusters ▶ Clusters ▶ Racks ▶ HW models ▶ Databases ▶ Virtual Organizations ▶ Power Distribution ▶

Last modified by mirsi (CERN IT/FIO-FS), May 08 2006 11:59:58. PHP version: 5.1.6

W3C HTML 4.01 ✓ W3C CSS ✓

Done lemonweb.cern.ch AdBlock



Outline

- Introduction to CERN and Experiments
- LHC Computing
- **Challenges**
 - _ Capacity Provision
 - _ **Box Management**
 - Installation & Configuration
 - Monitoring
 - **Workflow**
 - _ Data Management and Distribution
 - _ What's Going On?
- Summary/Conclusion



LHC Era Automated Fabric

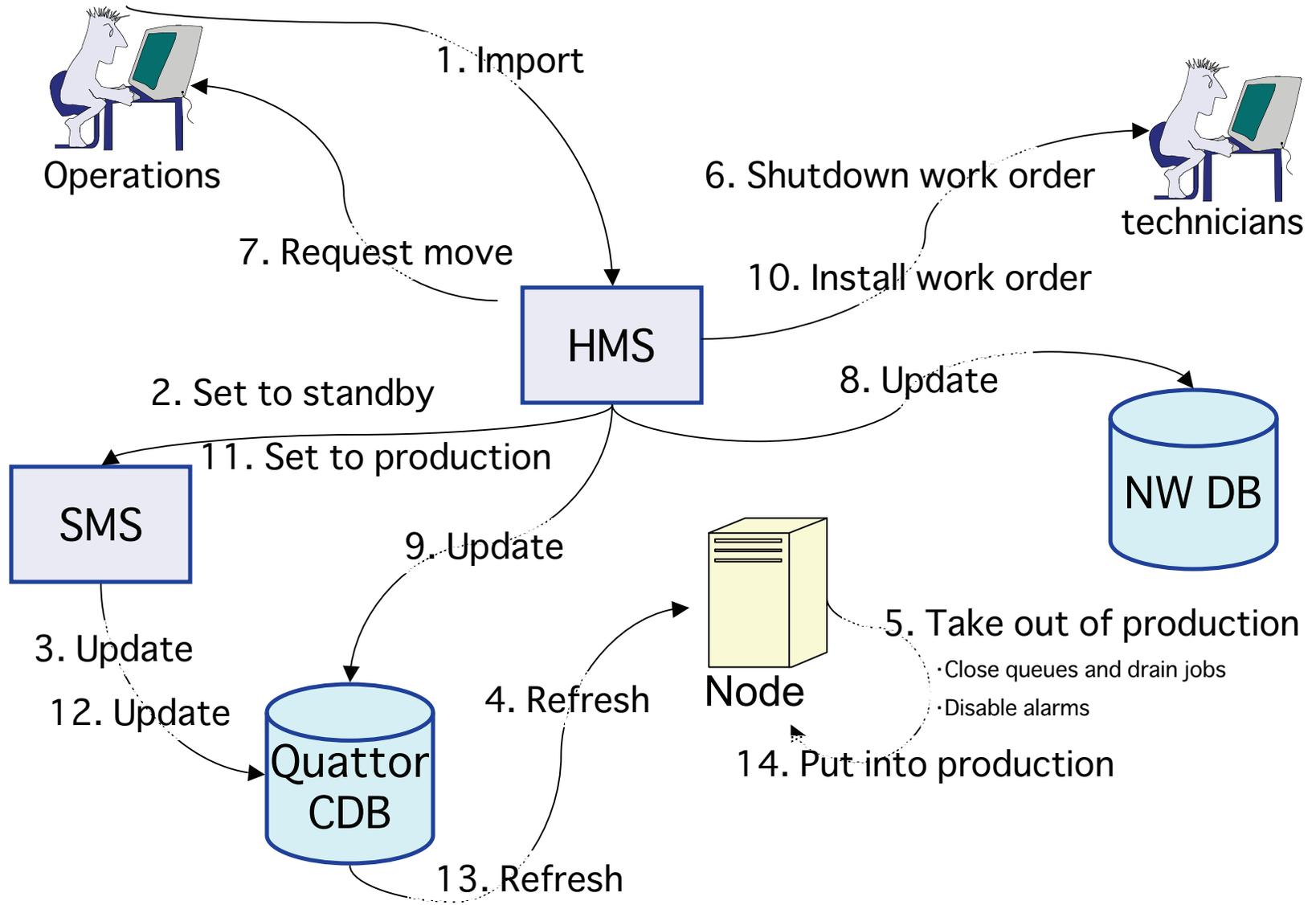
LEAF is a collection of workflows for *high level* node hardware and state management, on top of Quattor and LEMON:

- **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



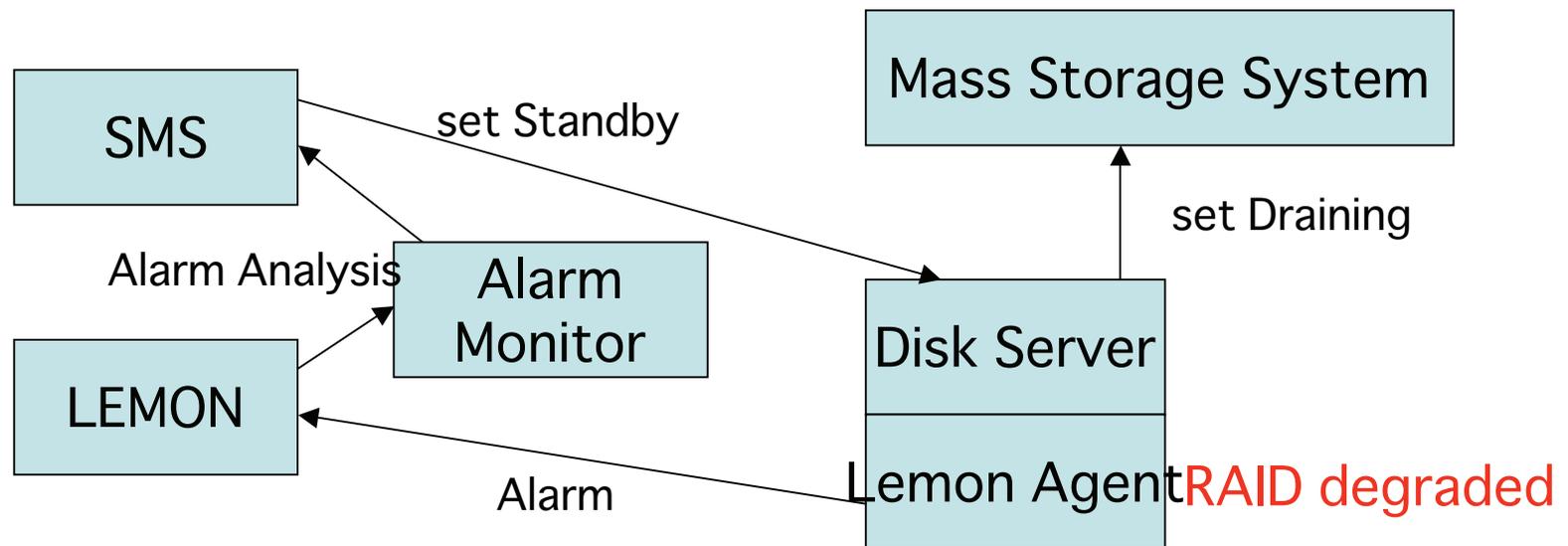


LEAF workflow example



Integration in Action

- Simple
 - _ Operator alarms masked according to system state
- Complex
 - _ Disk and RAID failures detected on disk storage nodes lead automatically to a reconfiguration of the mass storage system:



Draining: no new connections allowed; existing data transfers continue.



Box Mgmt Successes &

- At CERN, the ELFms toolkit has allowed us to cope with a significant increase in box count with reduced staffing levels.
- We have confidence the software will scale further
 - _ although changes needed (e.g. to cope with virtualisation).
- Large scale farm operation, though, remains a challenge!
 - _ ramp-up, purchasing, h/w failures, ...
 - (even if we are not at the Google scale)

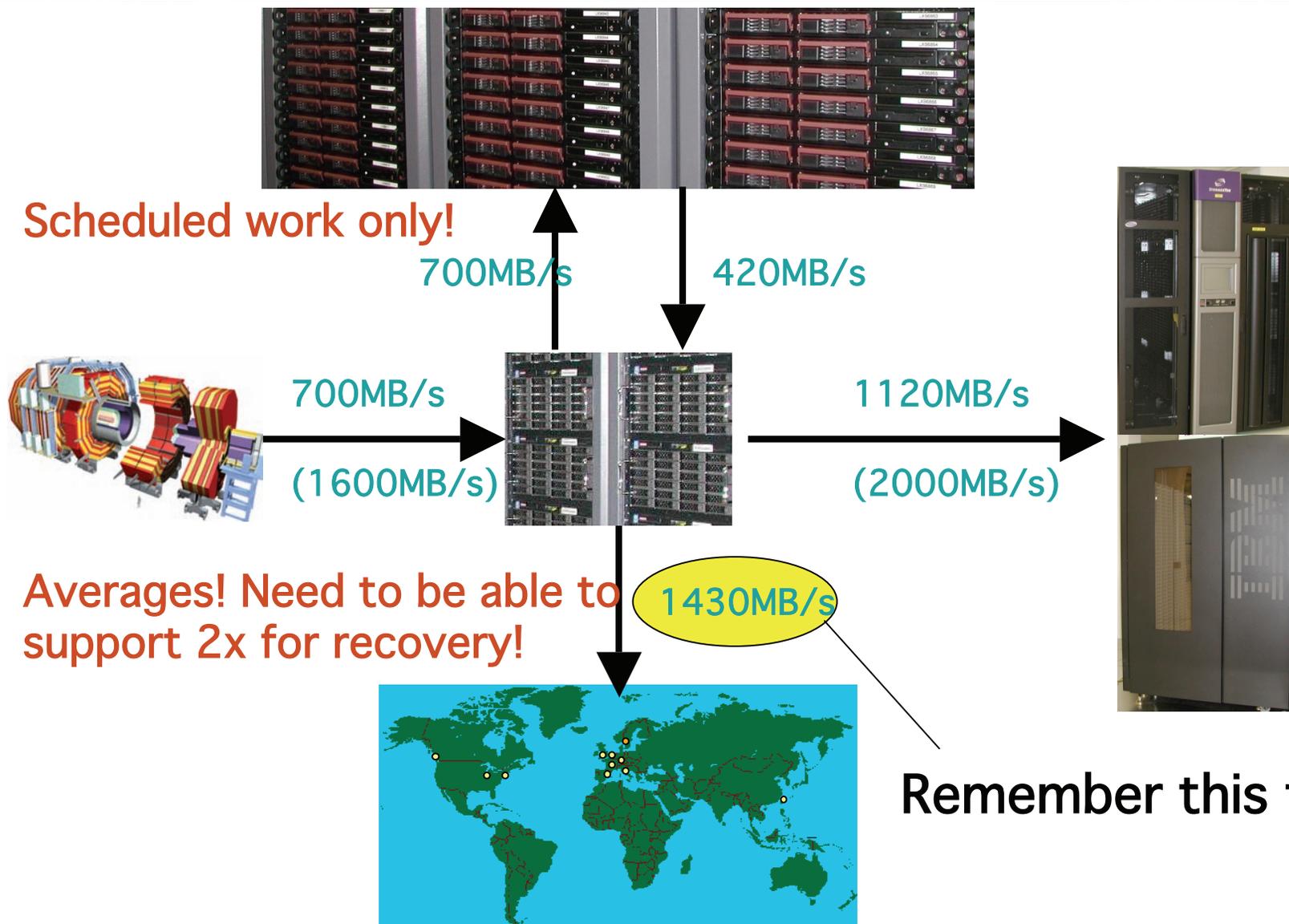


Outline

- Introduction to CERN and Experiments
- LHC Computing
- **Challenges**
 - _ Capacity Provision
 - _ Box Management
 - _ **Data Management and Distribution**
 - _ What's Going On?
- Summary/Conclusion



Dataflows and rates





Volumes & Rates

- 15PB/year. Peak rate to tape >2GB/s
 - _ 3 full SL8500 robots/year
- Requirement in first 5 years to reread all past data between runs
 - _ 60PB in 4 months: 6GB/s
- Can run drives at sustained 80MB/s
 - _ 75 drives flat out merely for controlled access
- Data Volume has interesting impact on choice of technology
 - _ Media use is advantageous: high-end technology (3592, T10K) favoured over LTO.



Access Patterns

- Multiple use cases...
 - _ Sustained transfer to remote site
 - WAN visibility; I/O intensive
 - _ Rapid transfer of data set to CPU node
 - LAN access; I/O intensive
 - _ Long running analysis access to data on server
 - LAN access, low I/O, long duration
- ...all have different footprint on disk servers



Mass Storage Systems @

- Commercial Mass Storage Systems have been evaluated at CERN, but with little success. Key systems evaluated were
 - _ Lachmann/Legent OSM
 - Still in use at DESY, but interest at CERN much reduced due to lack of long-term support (DESY provide their own support)
 - _ IBM's HPSS
 - In use at SLAC, BNL (US labs) and IN2P3 (French Computer Centre)
 - Experience at CERN showed random access to files (a major use case) was poor; addressing this required additional software and disk buffers
 - _ At the time, HPSS also required a DCE infrastructure and had limited O/S & hardware support.
 - _ IEEE "vision" of companies providing pluggable components of an overall system didn't work out in practice; we ended up with single vendors providing all the components...
- ... and so CERN developments became more and more capable leading to Castor: **CERN Advanced Storage System**



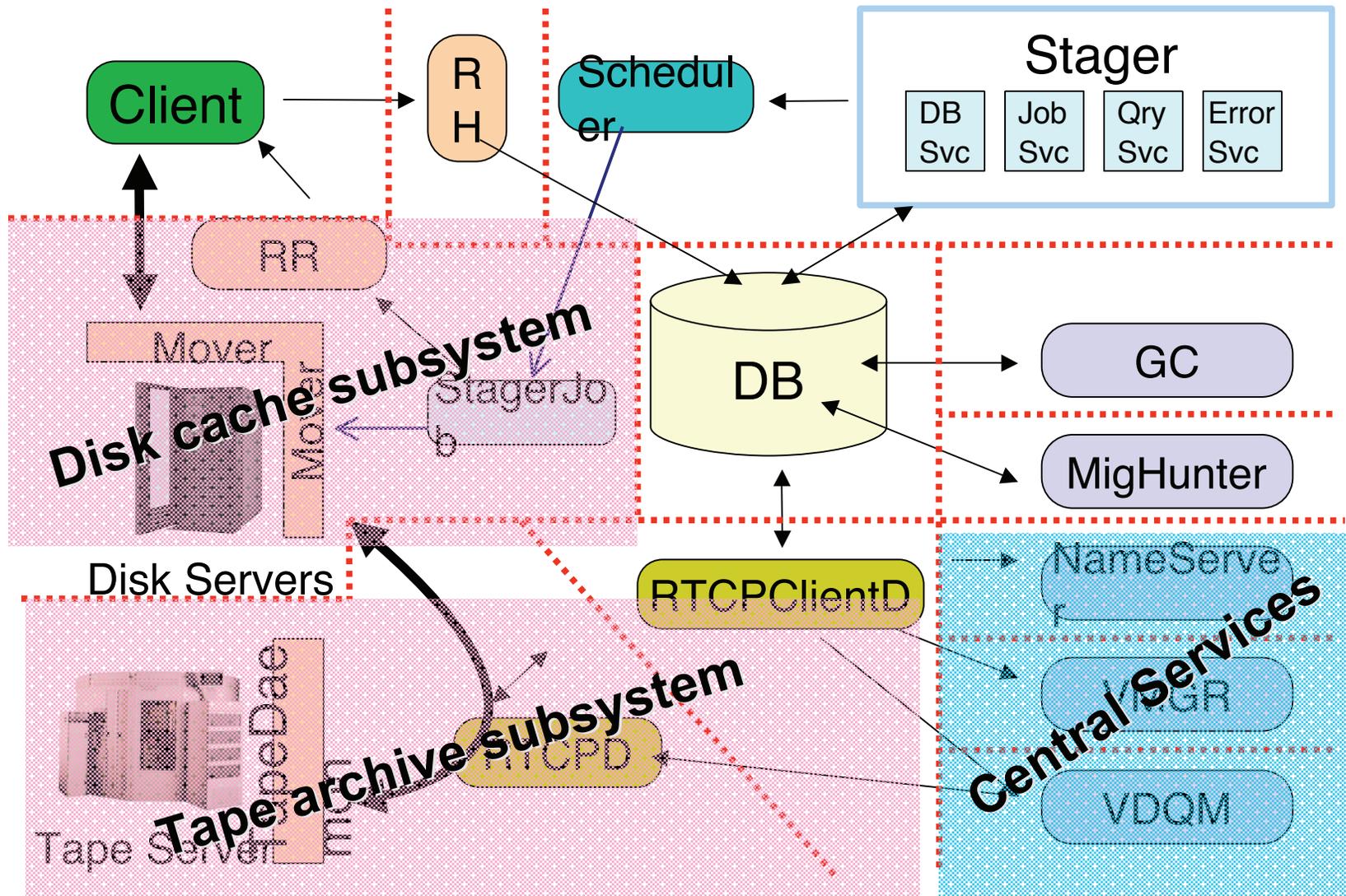
Key Castor Features



- Database Centric
 - _ Stateless agents; can restart easily on error
 - _ No direct connection from users to critical services
- Scheduled Access to I/O
 - _ No overloading of disk servers
 - Per-server limit set according to type of transfer
 - _ servers can support many random access style accesses, but only a few sustained data transfers
 - _ I/O requests can be scheduled according to priority
 - Fair shares access to I/O just as for CPU
 - Prioritise requests from privileged users



Castor Architecture



Detailed view



Castor Performance



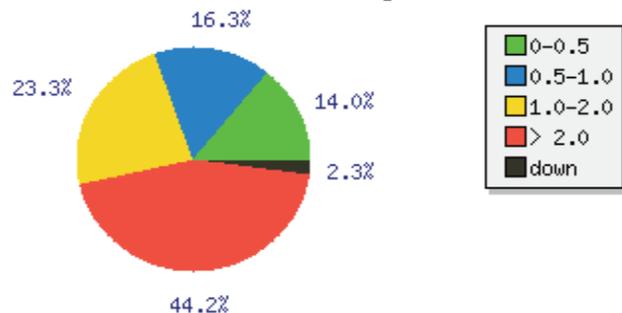
Cluster info: c2sc4 subcluster wan

19 Apr 2006 Wed 16:14:27

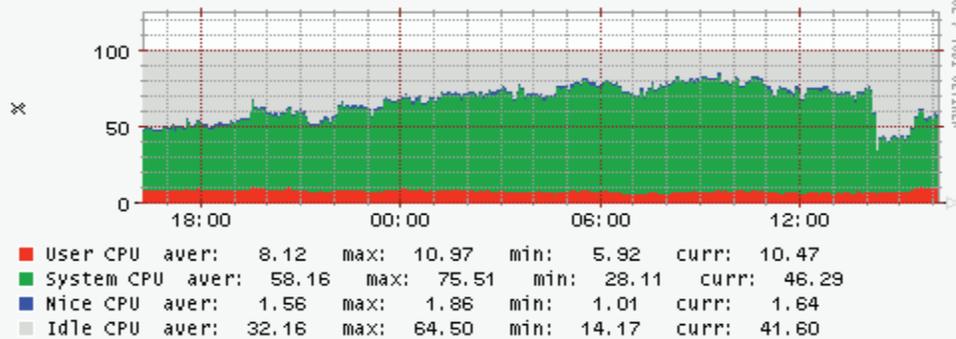
Cluster Information

# of hosts (down):	43 (1)
operating system(s):	2.6.9-34.EL.cernsmp
# of CPUs (down):	51 (2)
average up time:	20 days, 18h:50m (boots per host)
hosts down:	lxfkra3004
exceptions:	RPC_STATD_WRONG, MIRROR_BROKEN,
ITCM history	View template
Select from hosts:	<input type="text" value="None"/>
Metric Distributions	Correlations

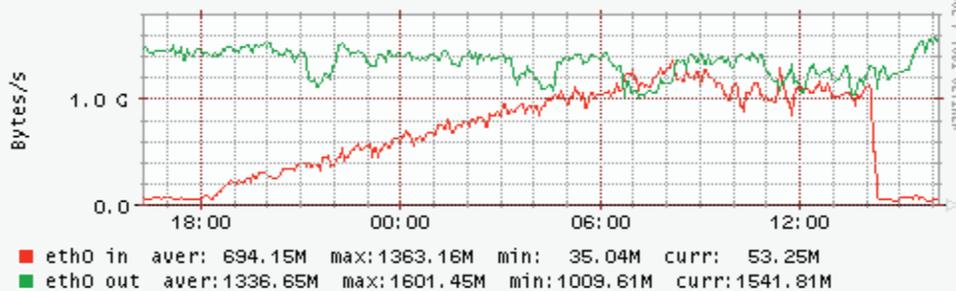
Load Percentages



CPU utilization - last day



Network utilization - last day



Last

Sustained transfer from disk of 1.2GB/s as data import ramps



Castor Performance



Cluster info: ITDC

03 Jan 2006 Tue 10:19:12

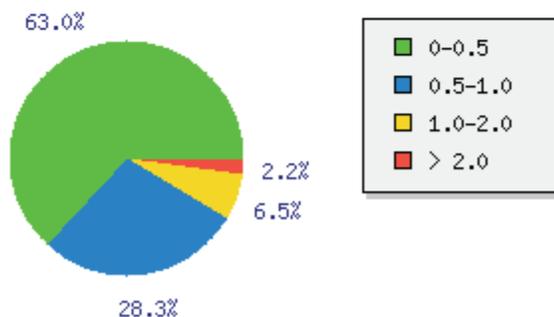
Cluster Information

# of hosts (down):	46 (0)
operating system(s):	2.4.21-37.EL.cernsmp
# of CPUs (down):	62 (0)
average up time:	47 days, 14h:14m (boots per host)
hosts down:	none
exceptions:	FILESYSTEM_ERROR
ITCM history	View template
Select from hosts:	<input type="text" value="None"/>

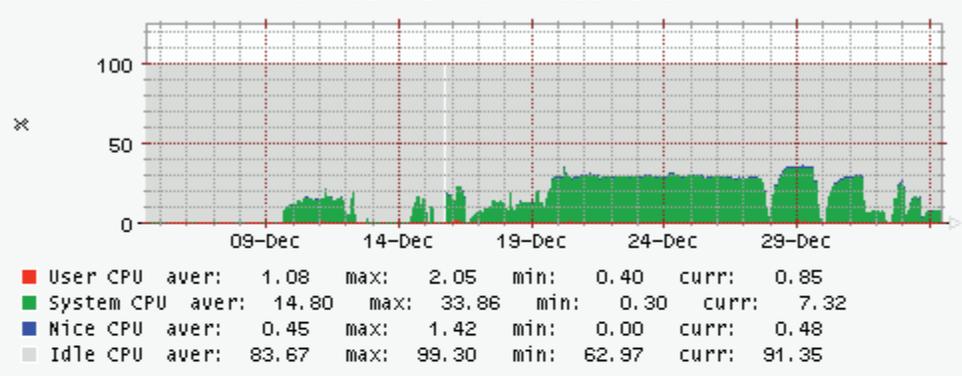
Metric Distributions

Correlations

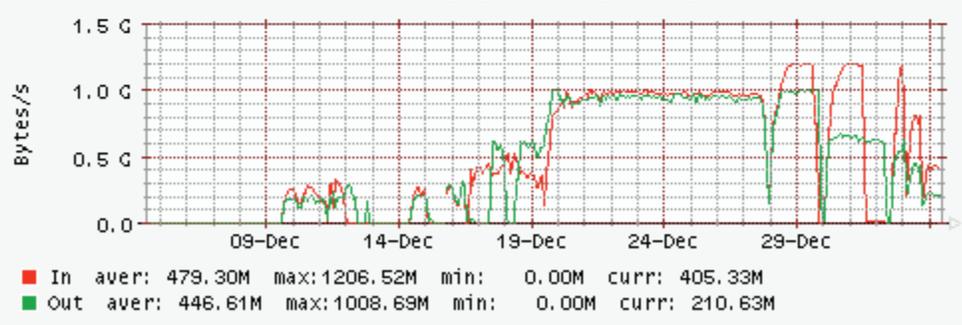
Load Percentages



CPU utilization - last month



Network utilization - last month



Last

Sustained transfer of incoming data to tape at 1GB/s
Note the dates! Failed hardware was left down.



Castor Performance



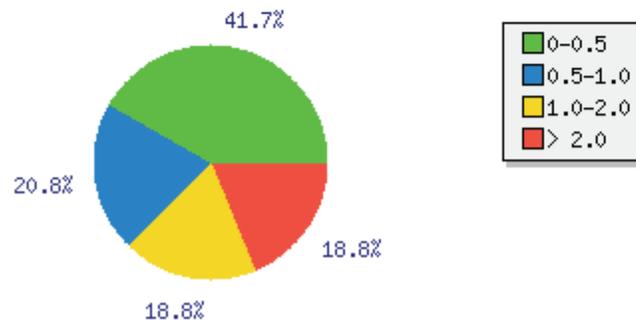
Cluster info: castor2 subcluster ITDC

06 Mar 2006 Mon 08:15:45

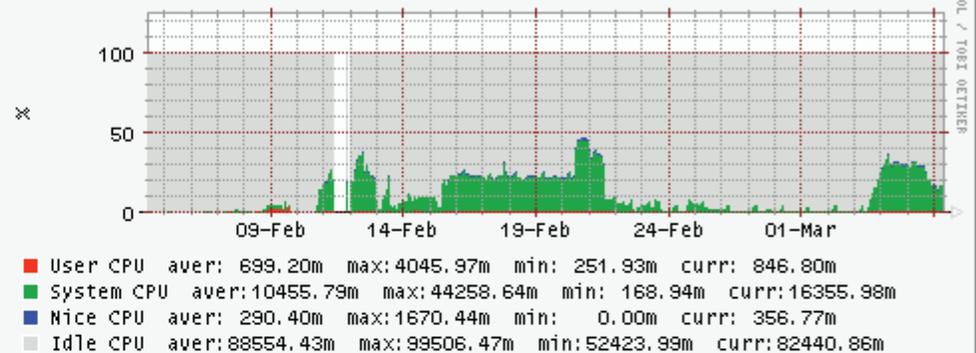
Cluster Information

# of hosts (down):	48 (0)
operating system(s):	2.4.21-37.EL.cernsmp, 2.4.21-37.0.1.EL.cernsmp
# of CPUs (down):	88 (0)
average up time:	38 days, 7h:43m (boots per host)
hosts down:	none
ITCM history	View template
Select from hosts:	<input type="button" value="None"/> <input type="button" value="v"/>
Metric Distributions	Correlations

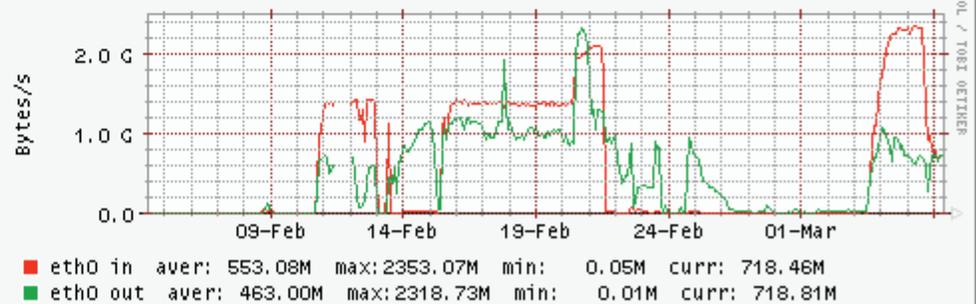
Load Percentages



CPU utilization - last month



Network utilization - last month

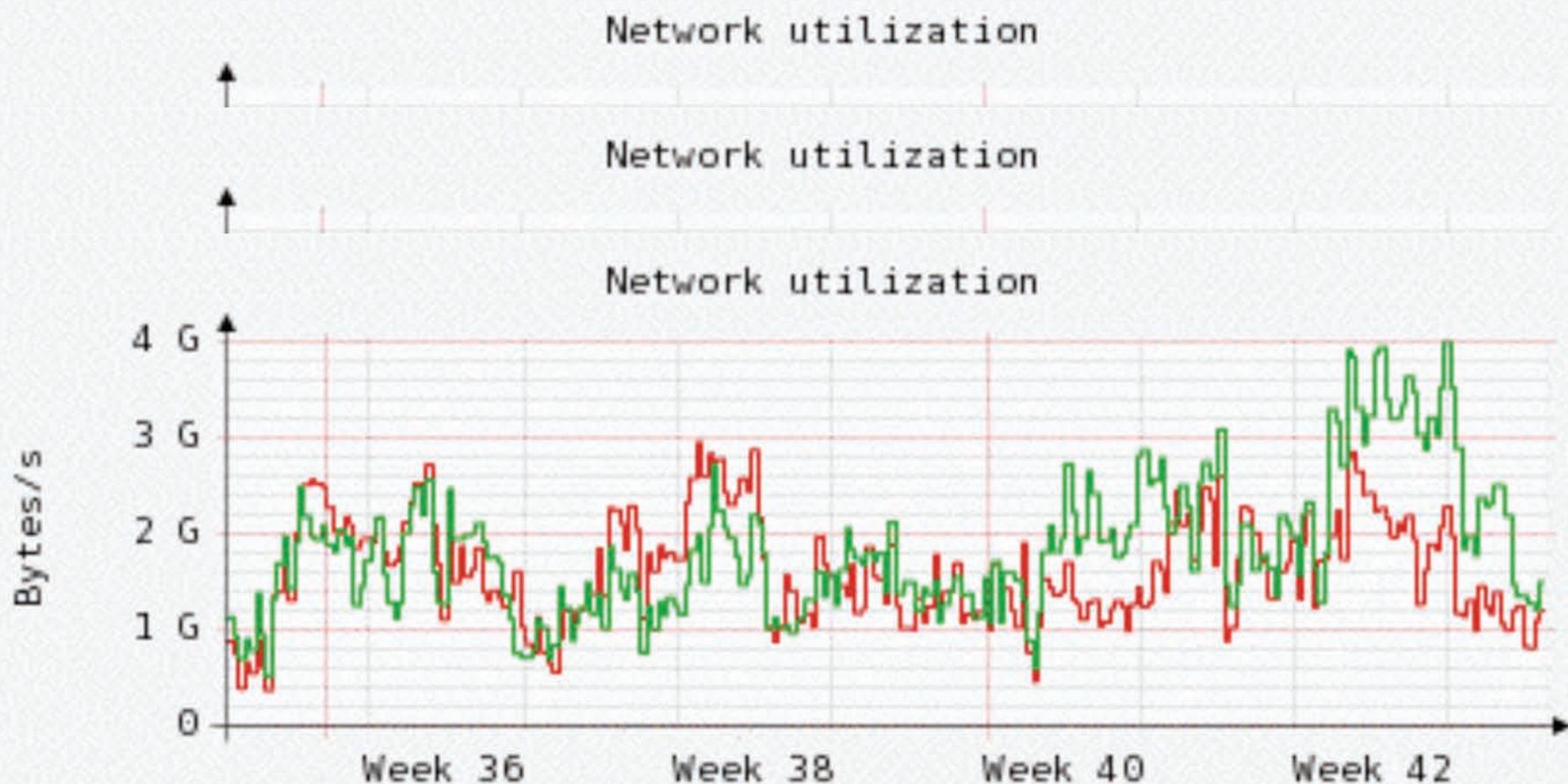


Last

Peak transfer of incoming data to tape at over 2GB/s



Castor Performance



RRDTOI RRDTOI RRDTOOL / TOBI OETIKER

■ eth0 in	aver:1.6G	max:2.9G	min:373.6M	curr:1.2G
■ eth0 out	aver:1.8G	max:4.0G	min:511.8M	curr:1.5G
■ eth1 in	aver:0.0	max:0.0	min:0.0	curr:0.0
■ eth1 out	aver:0.0	max:0.0	min:0.0	curr:0.0



More Data Challenges

- Long data lifetime
- Disk capacity vs I/O rates
- File sizes
- Multiple Mass Storage Systems
- Organised Data Export

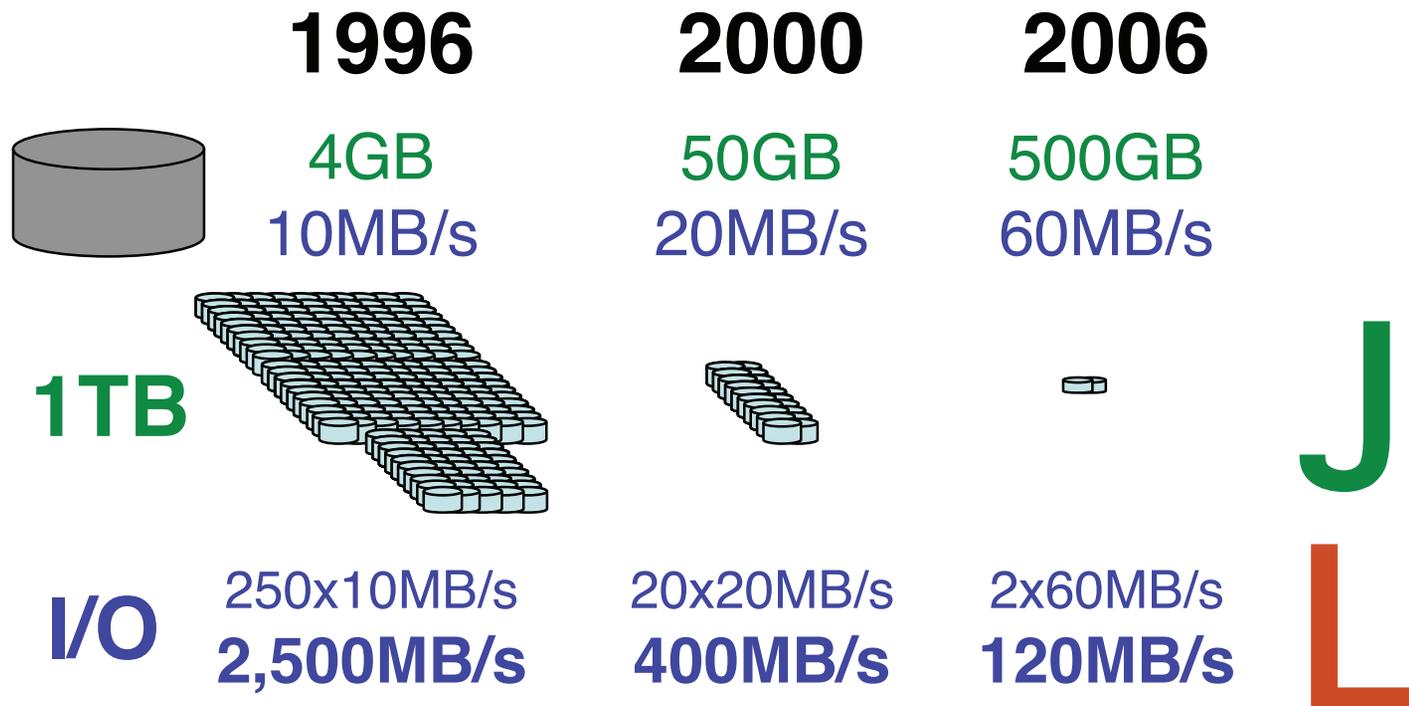


Long lifetime

- LEP, CERN's last accelerator, started in 1989 and shutdown 10 years later.
 - _ First data recorded to IBM 3480s; at least 4 different technologies used over the period.
 - _ **All data ever taken**, right back to 1989, was reprocessed and reanalysed in 2001/2.
- LHC starts in 2007 and will run until at least 2020.
 - _ What technologies will be in use in 2022 for the final LHC reprocessing and reanalysis?
- Data repacking required every 2-3 years.
 - _ Time consuming
 - _ Data integrity must be maintained



Disk capacity & I/O rates



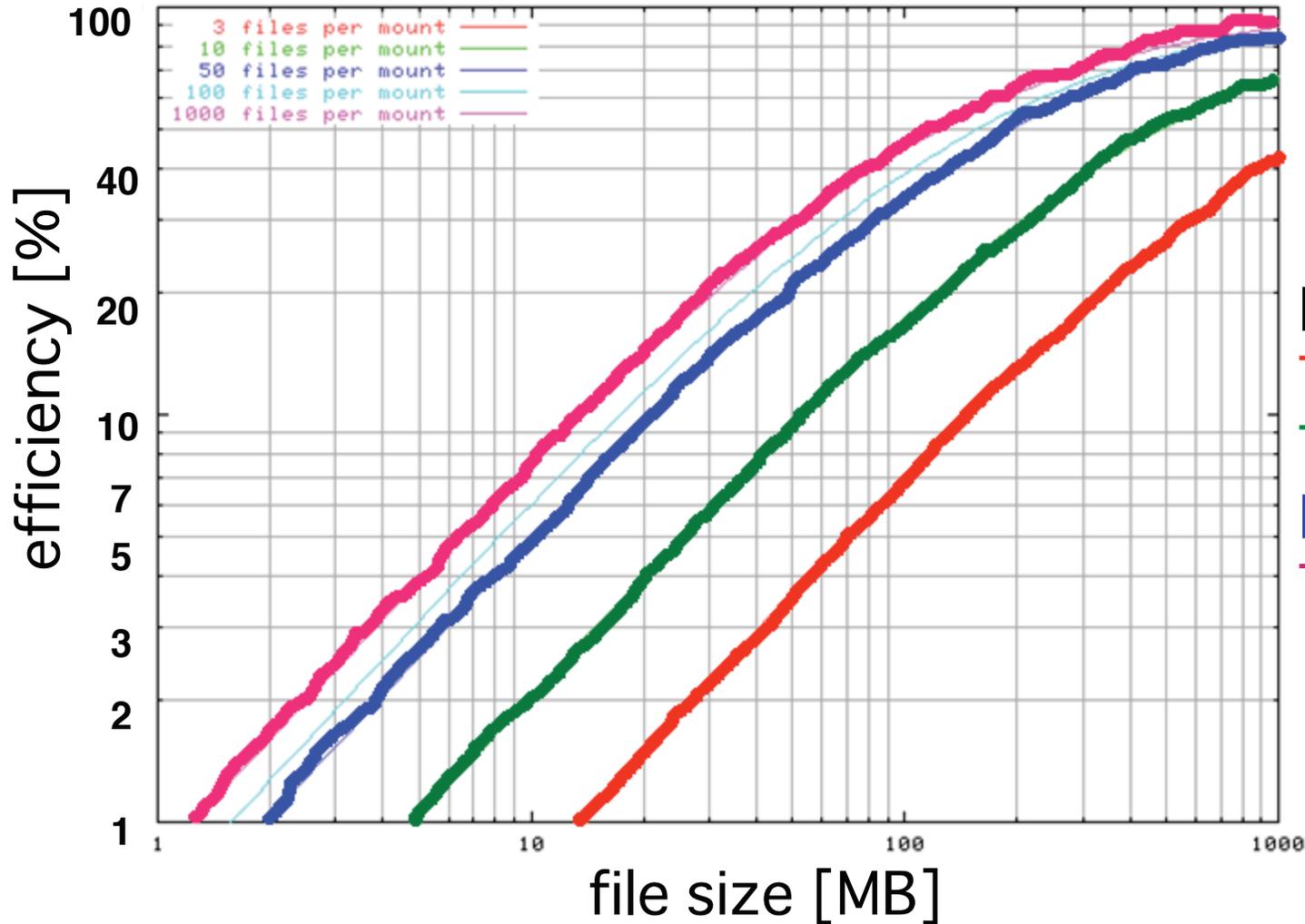
CERN now purchases two different storage server models: capacity oriented and throughput oriented.

- fragmentation increases management complexity
- (purchase overhead also increased...)



(File) Size Matters!

Tape Drive Efficiency (30 MB/s max performance)



Files per mount
Three
Ten
Fifty
Thousand

tape mount time ~120s
file overhead ~5s (400MB)



Multiple Mass Storage

- CASTOR is not the sole MASS for LHC
 - _ Fermilab's dCache is used at many sites; DPM, a disk-only storage manager is also common.
- Users, of course, don't want to know...
 - _ ... and experiment code needs to run at many sites
 - _ SRM, the Storage Resource Manager, provides a common interface layer to the various mass storage systems
 - See <http://sdm.lbl.gov/srm-wg/>
 - _ These multiple and independent implementations of the interface **all talk to each other.**
 - _ Key element to have successful...



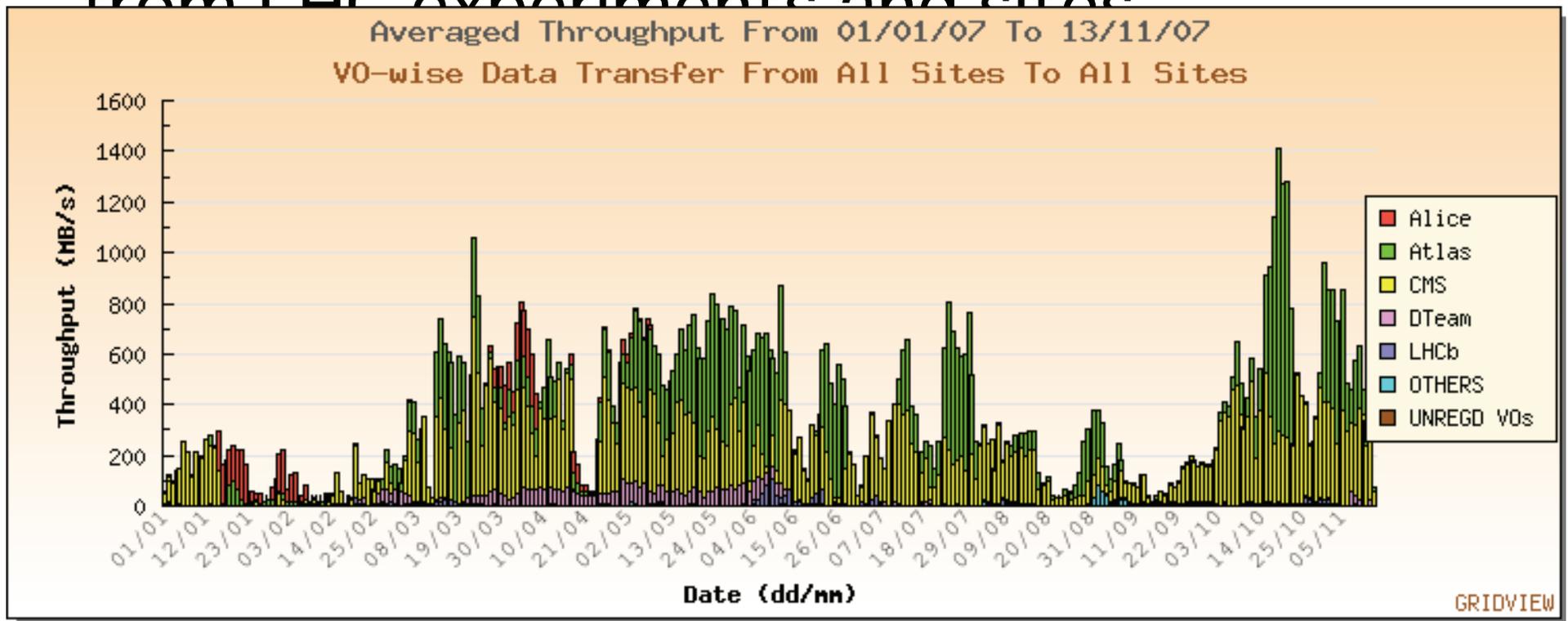
Organised Data Export

- LHC experiments need to ship data between sites
 - _ Raw data export, Analysis Data updates, Monte Carlo data import, ...
- This is complicated at our scale
 - _ with petabytes of data transferred, a 0.1% failure rate can't be easily rescued or followed up manually
 - _ Sites policies (e.g. fraction of resources allocated to a given VO) must be respected
 - _ ...



The File Transfer Service

- Developed as part of the EGEE data management activity to meet requirements from LHC experiments and sites



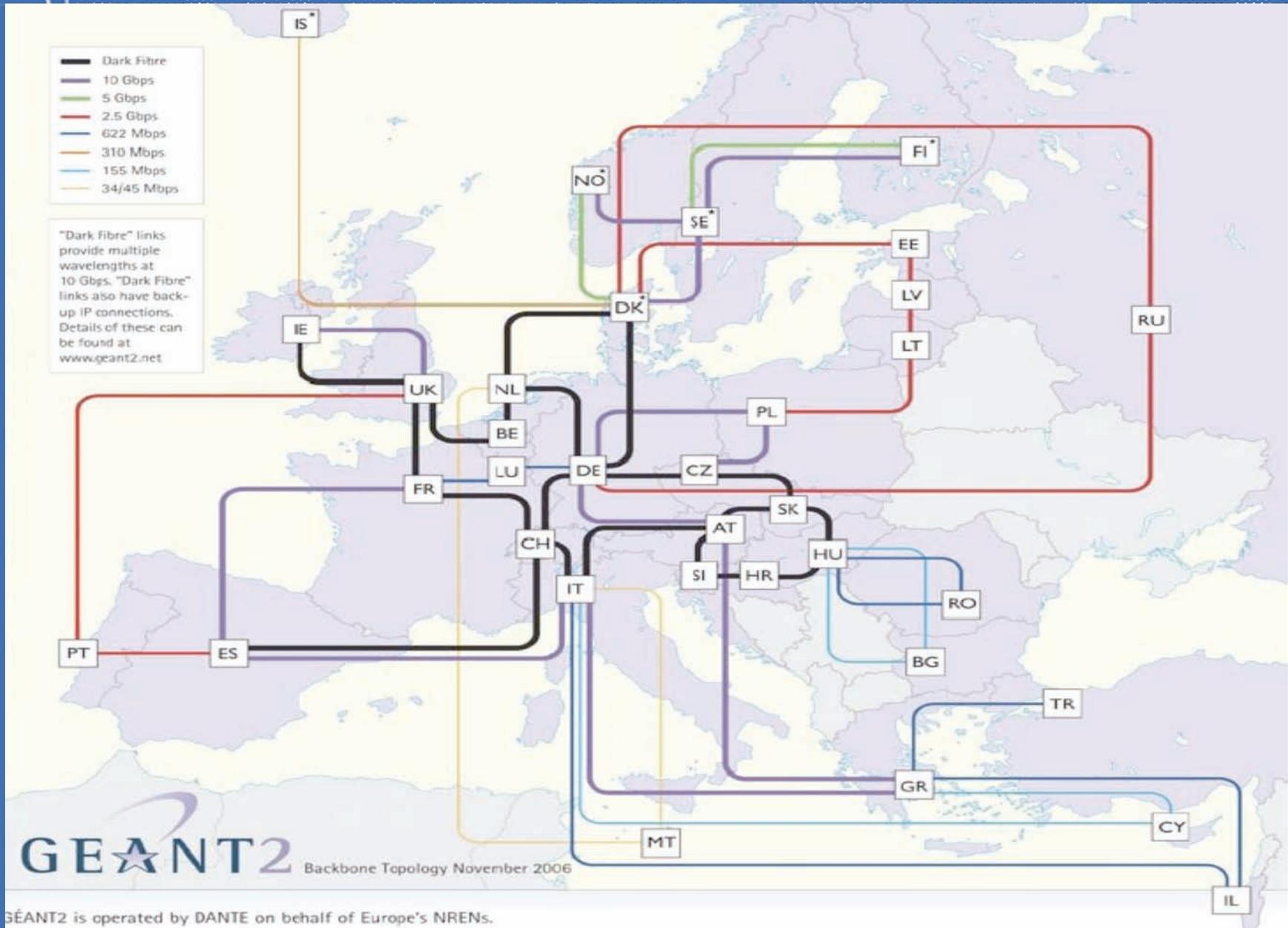
• Prevent network overload

– (But no management on the network level)

- Prevent storage overload



... and we needed a network





Data Successes & Challenges

- Successes:
 - _ We have an advanced Mass Storage System at CERN able to meet the demanding requirements for Data Acquisition and export.
 - _ Large scale data transfers between sites are becoming routine.
- But yet to demonstrate
 - _ exports for multiple experiments simultaneously
 - _ operations for large scale user analysis
 - most work so far has been controlled “production”



Outline

- Introduction to CERN and Experiments
- LHC Computing
- **Challenges**
 - _ Capacity Provision
 - _ Box Management
 - _ Data Management and Distribution
 - _ **What' s Going On?**
- Summary/Conclusion



A Complex Overall Service

- Site managers understand systems (we hope!).
- But do they understand the service?
 - _ and do the users?
 - _ and what about cross site issues?
 - Are things working?
 - If not, just where is the problem?
 - _ how many different software components, systems and network service providers are involved in a data transfer site X to site Y?



User Status Views @ CERN

Home
Dependencies
Admin
Documentation
Help

IT Services
choose another view: ▼
4 Oct 2006 Wed 23:34:05

<div style="border: 1px solid #ccc; padding: 2px; margin-bottom: 5px;"> Administrative applications f 100% available (more) </div>	<div style="border: 1px solid #ccc; padding: 2px; margin-bottom: 5px;"> Services for physics * 82% available (more) </div>	<div style="border: 1px solid #ccc; padding: 2px; margin-bottom: 5px;"> For developers and engineers {} 75% available (more) </div>
<div style="border: 1px solid #ccc; padding: 2px; margin-bottom: 5px;"> Windows Services W 100% available (more) </div>	<div style="border: 1px solid #ccc; padding: 2px; margin-bottom: 5px;"> Infrastructure services I 99% available (more) </div>	<div style="border: 1px solid #ccc; padding: 2px; margin-bottom: 5px;"> Databases D 99% available (more) </div>
<div style="border: 1px solid #ccc; padding: 2px; margin-bottom: 5px;"> Mail and Web Services M 97% available (more) </div>	<div style="border: 1px solid #ccc; padding: 2px; margin-bottom: 5px;"> Worldwide LHC Computing Grid W 68% degraded (more) </div>	<div style="border: 1px solid #ccc; padding: 2px; margin-bottom: 5px;"> Networking N 100% available (more) </div>
<div style="border: 1px solid #ccc; padding: 2px; margin-bottom: 5px;"> Documents and collaborative services D 100% available (more) </div>		

SLS by CERN IT/FIO

Service performance for 16 Nov 2006

Key Performance Indicators:

Number of services: ●

Number of valid updates: ●

All OK: Number of valid updates (180) was higher than minimum expected/ target level (150)

Services for physics

availability:
[\(more\)](#)

percentage: 82%

status: **available**

this service consist of:

- Castor2
- Batch service
- LXPLUS
- LXBUILD
- LXGATE
- CERN LXGATE Facility (LXGATE)**
Availability: 92%, **affected**
2 out of 24 tested nodes are not reachable.
- AFS
- Kerberos
- CDB
- Lemon
- Linux software

More on availability in the last week

Grid Tier-1 sites
availability - last 7 days

Statistics:

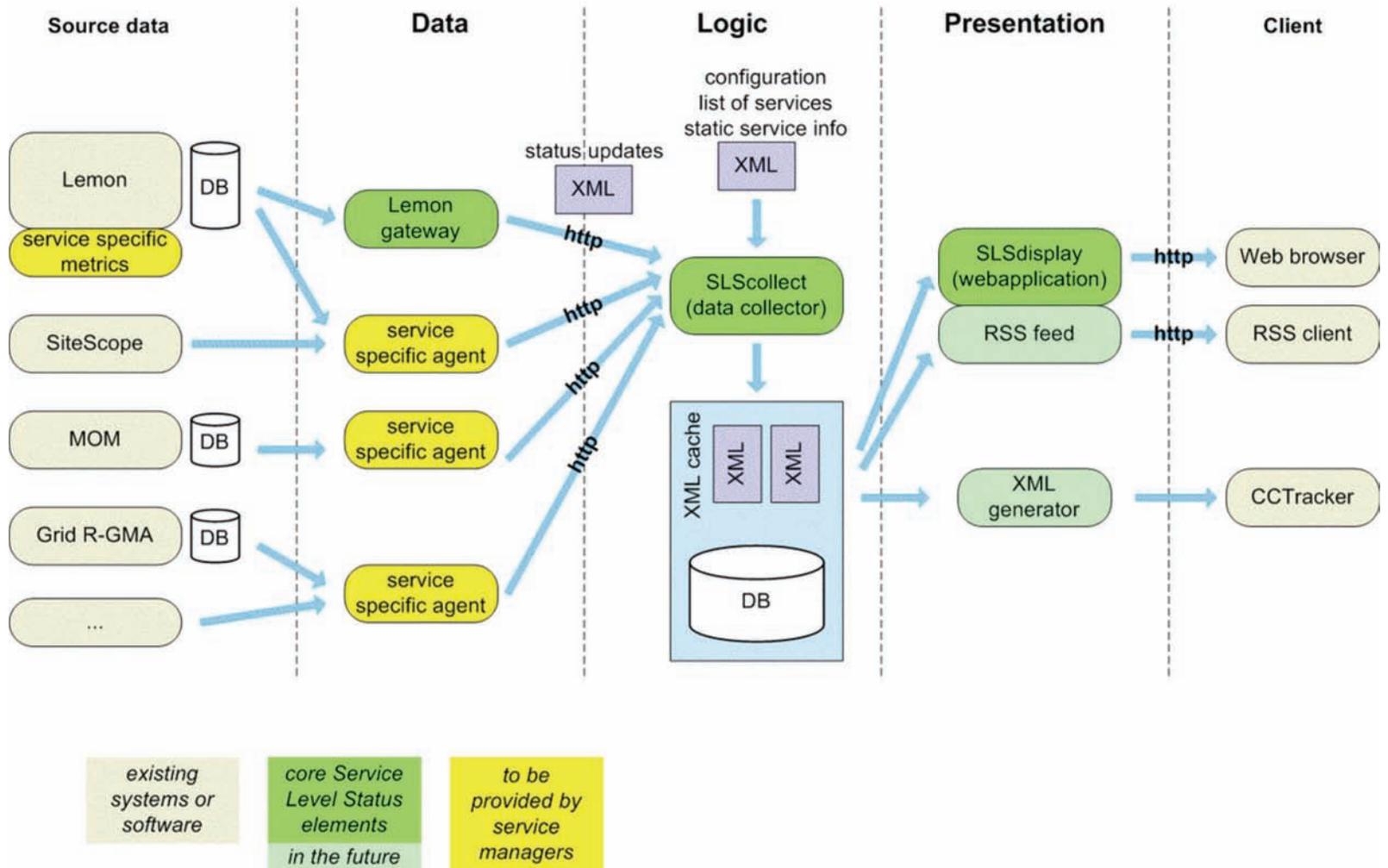
last week avg: **47% affected**

max: **70% affected** last at 07 Oct 2006, 13:00

min: **11% available** last at 05 Oct 2006, 22:00

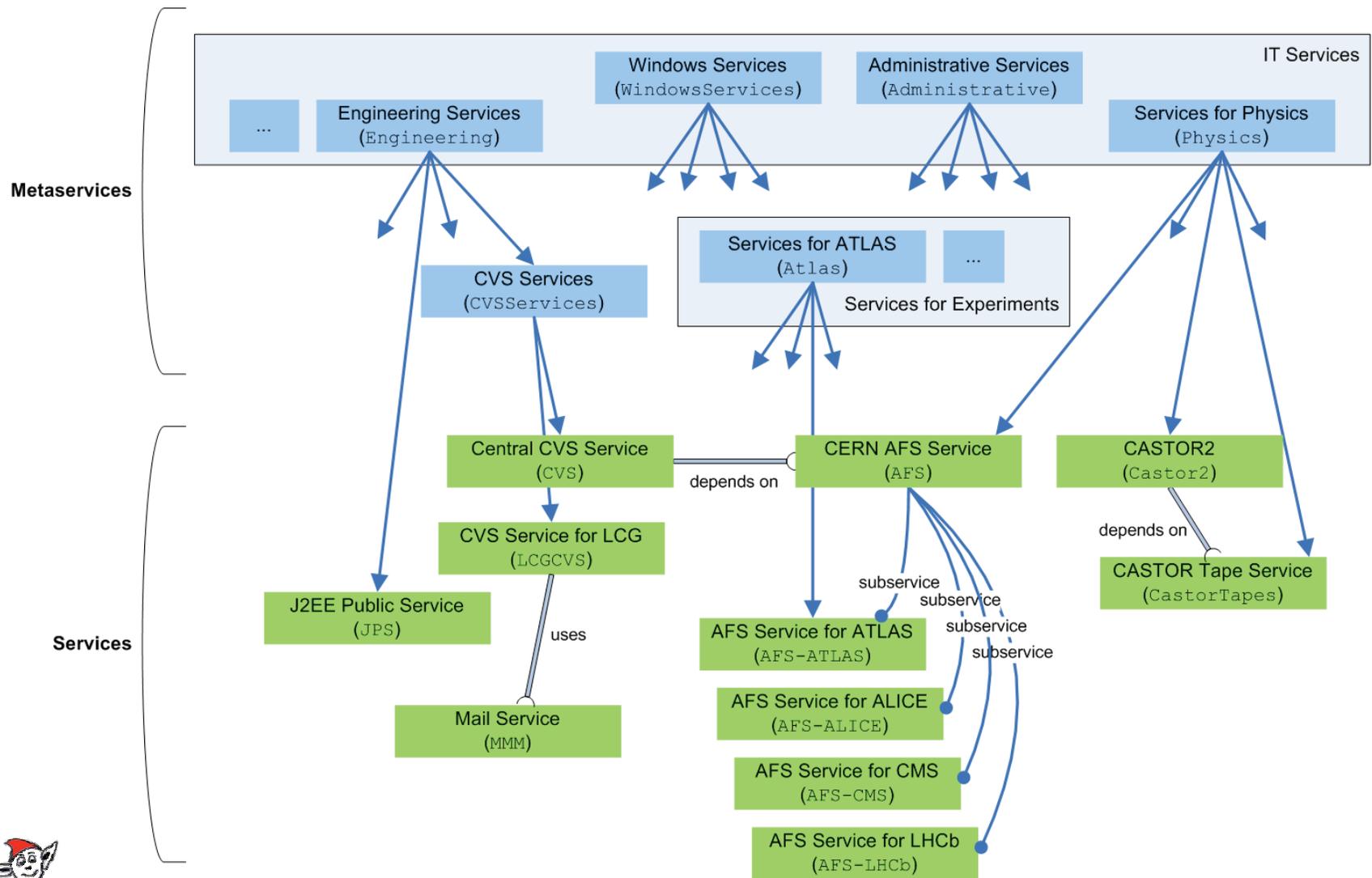


SLS Architecture





SLS Service Hierarchy





SLS Service Hierarchy

Service Level Status overview

Home Dependencies Admin Documentation Help

Services used by LHCb people 19 Nov 2006 Sun 12:58:09

Services used by LHCb

availability: (more)

percentage: 84%

status: **available**

this service consist of:

- Services for physics
- LXPLUS
- Batch service
- AFS
- AFS Service for LHCb
- LXBUILD (LHCb)
- Mail and Web Services
- Indico
- Twiki
- EDH system
- Electronic Document Handling system (EDH)**
Availability: **100%, available**
23 SiteScope test(s) out of 23 succeeded
- Castor2LHCb
- CASTOR Tape Service

availability in the last 24 hours (more):

Additional information

full name: **Services used by LHCb people**

short name: Services used by LHCb

vo: LHCb

email: Helpdesk@cern.ch

web site: <http://lhcb.cern.ch>

Availability update

last update: 12:51:10, 19 Nov 2006
(7 minutes ago)

refreshed every: 2 minutes

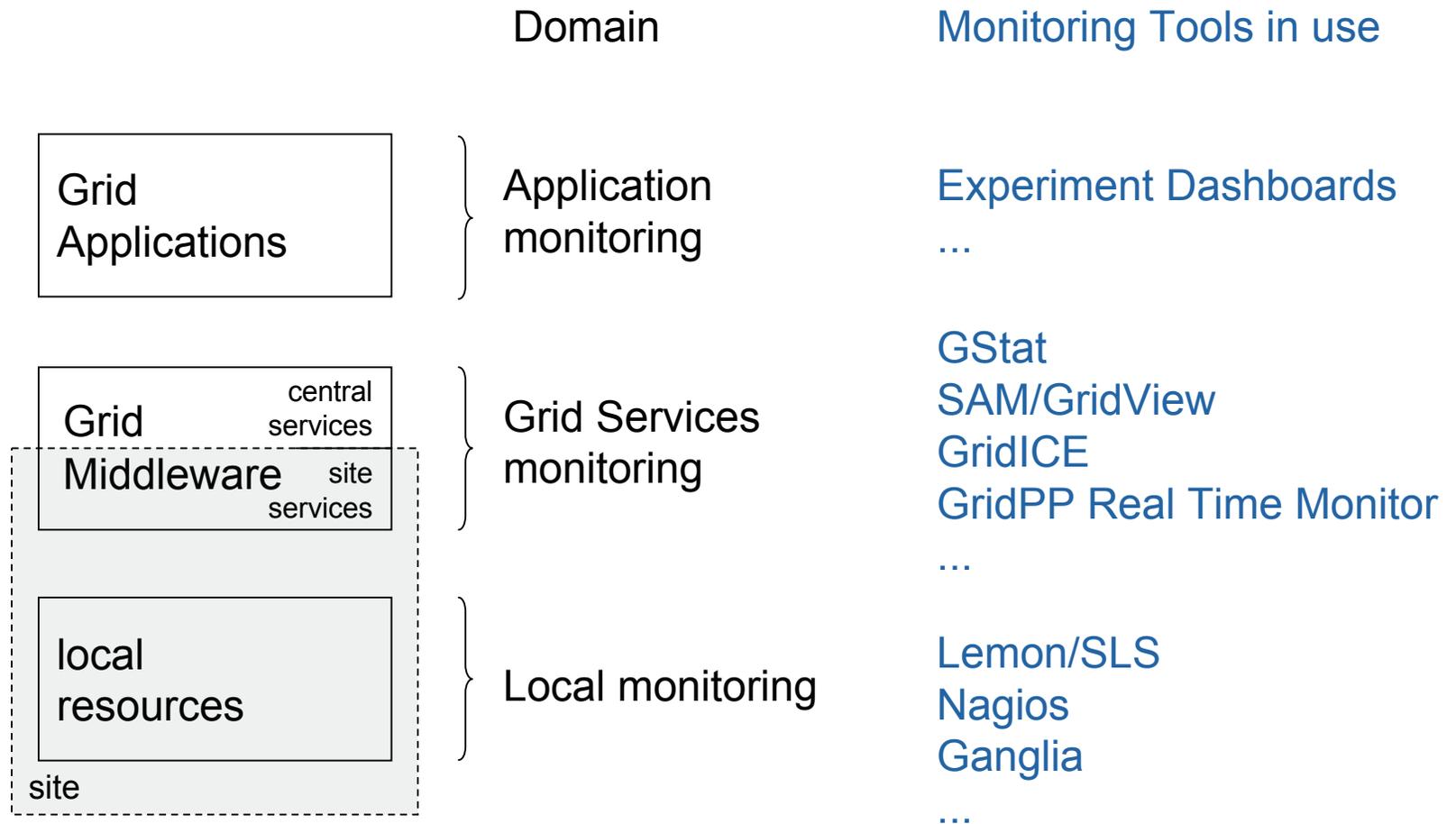
expires after: 60 minutes

Admin

[admin tools](#)

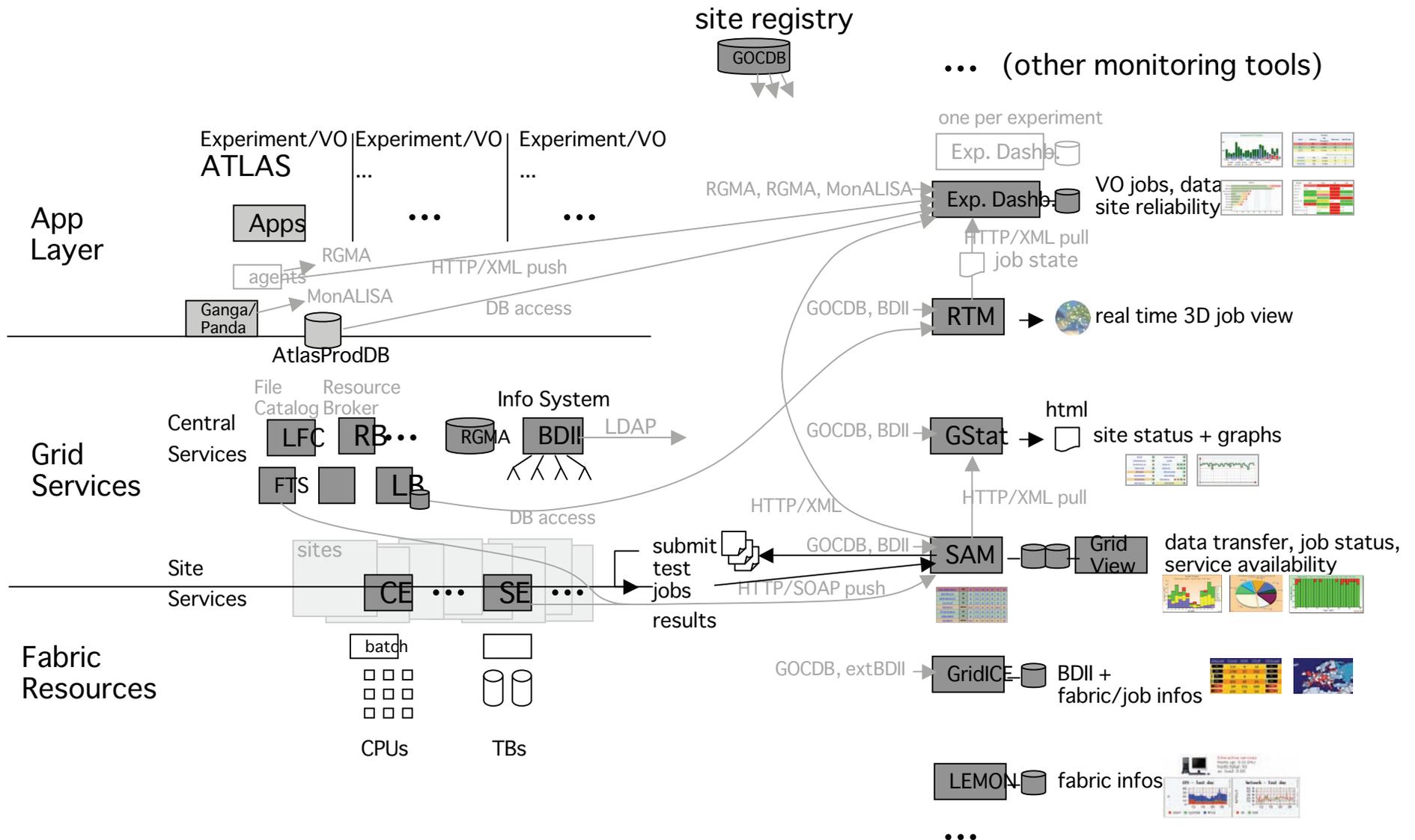


WLCG Grid Monitoring





Grid Monitoring Landscape





New visualizations for the

- Grid monitoring data is complex!
 - _ And there are many sites...
- Current tools visualize data by sorted tables, bar charts, etc.
- Difficult to present an easy to understand **top-level view** which provides
 - _ quick, action oriented oversight and insight
 - _ help understand job failures and availability patterns

Can new visualizations help?



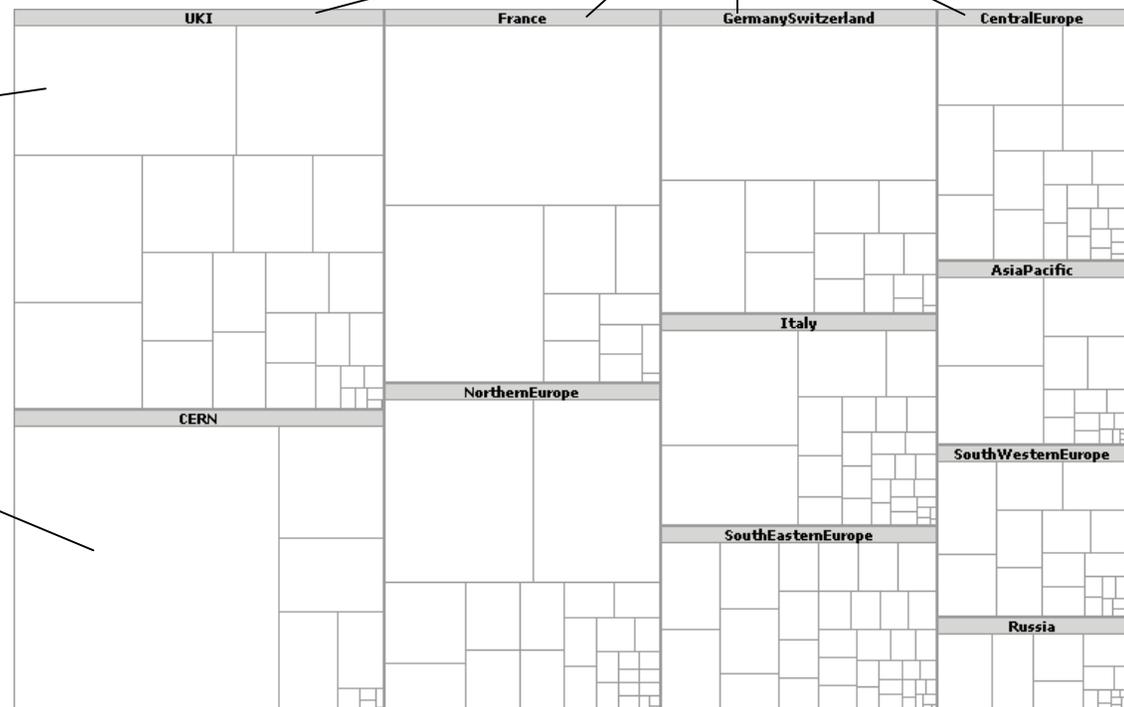
GridMap Visualization

- Idea
 - _ visualize the Grid by using *Treemaps*
 - (Grid + Treemap = *GridMap*)

- Example *GridMap*

site

regions



Size of rectangle is e.g.

- size of site (#CPUs)
- #running jobs
- ...



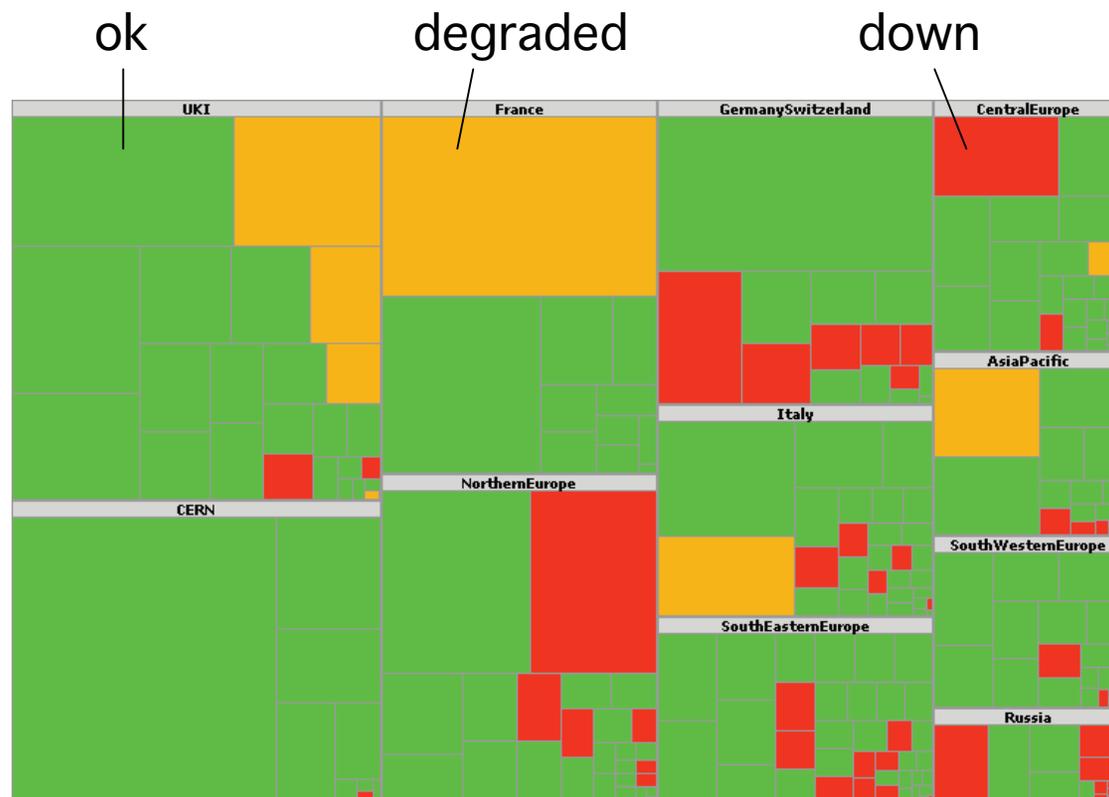
GridMap Visualization

- Idea
 - _ visualize the Grid by using *Treemaps*
 - (Grid + Treemap = *GridMap*)

- Example *GridMap*

Colour of rectangle is e.g.

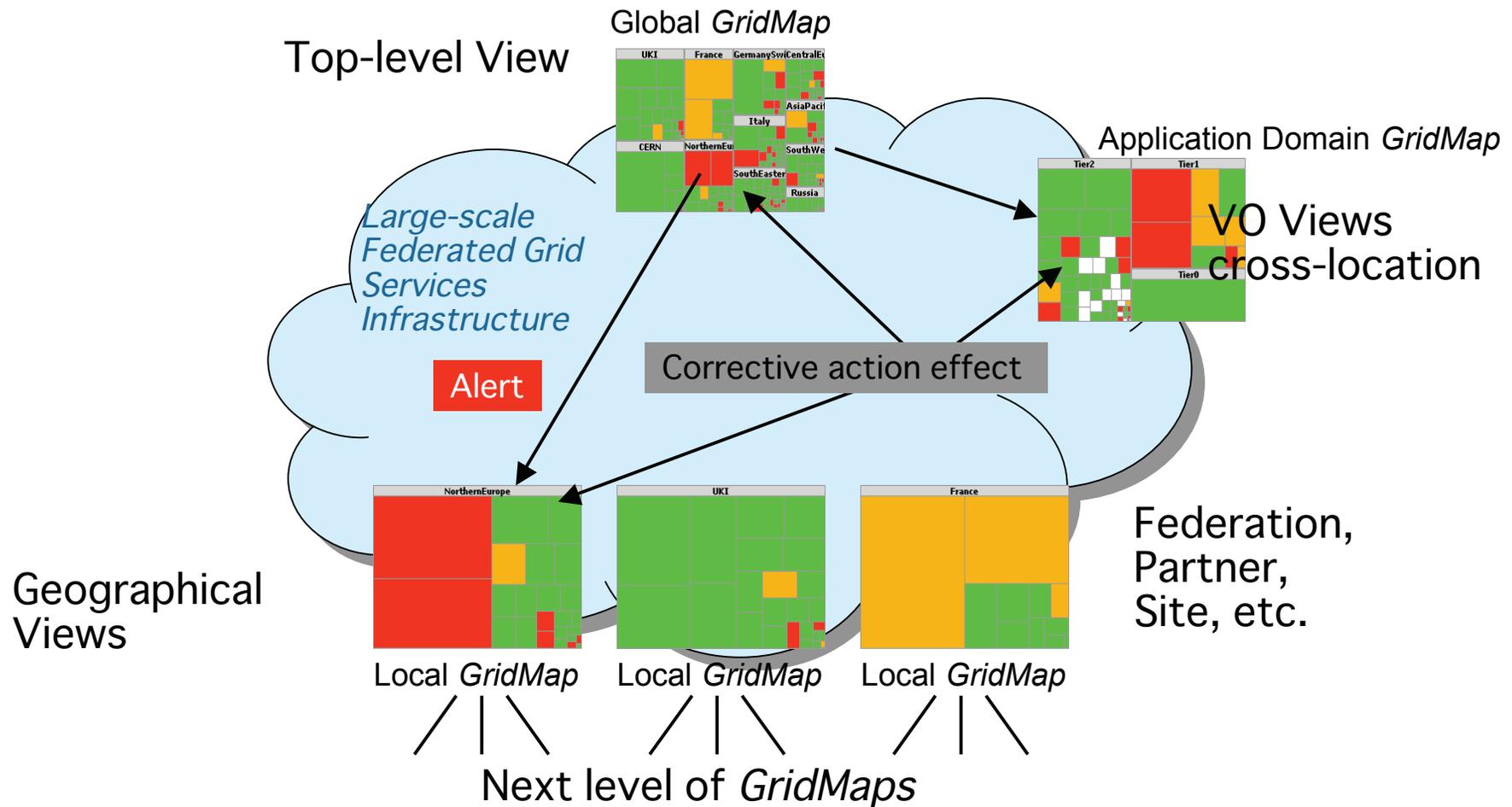
- SAM status of site / service
- Availability of site / service
- ...





Multiple Views

- *GridMaps* can be used for *top-level*, *geographical* and *VO* views

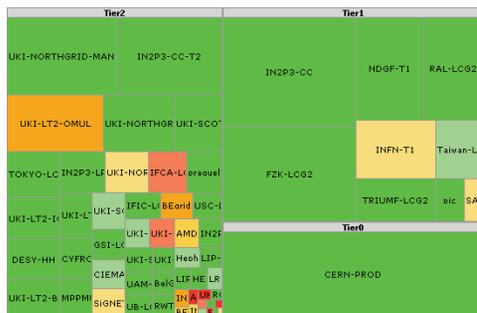




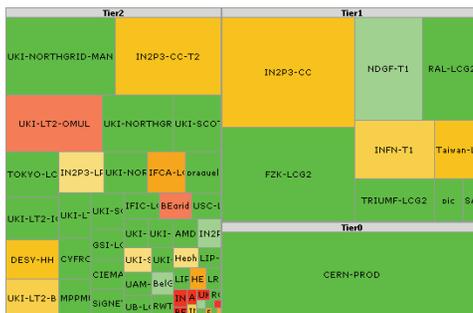
Trends

- Trends can be understood by looking at a sequence of *GridMaps*

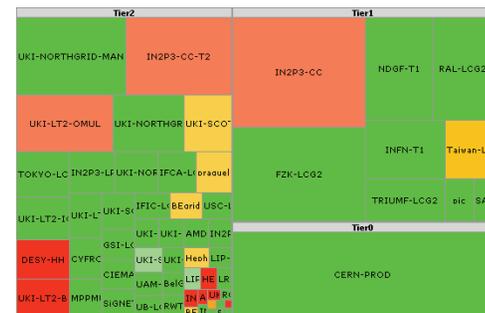
Site Availability over time:



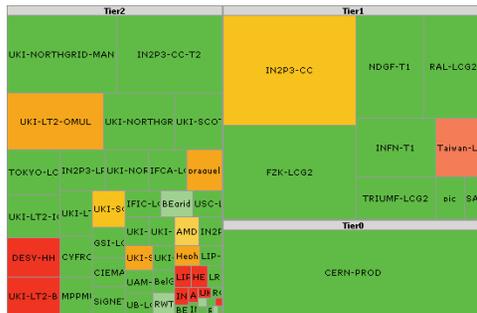
20 Sep 2007



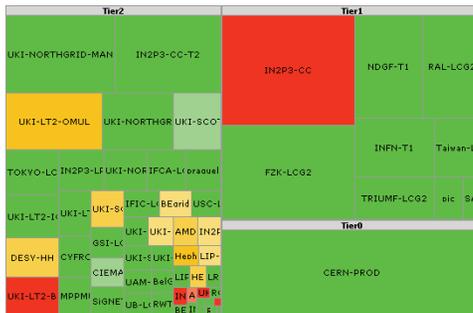
21 Sep 2007



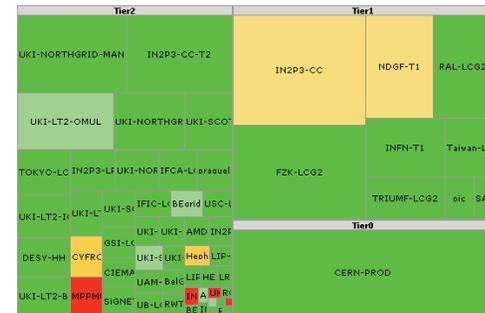
22 Sep 2007



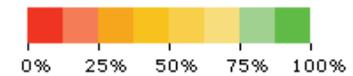
23 Sep 2007



24 Sep 2007



25 Sep 2007



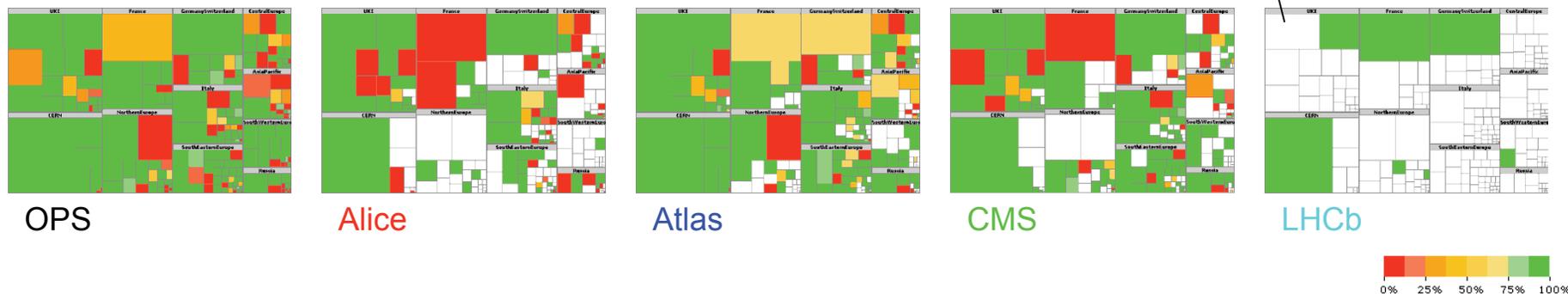


More Views

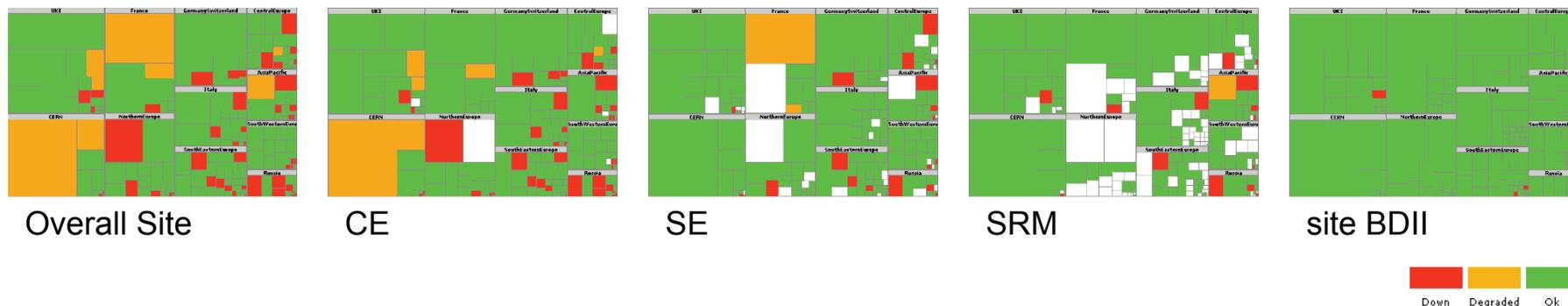
- Correlations of metrics can be discovered by switching between different views

Site Availability from different VO perspectives:

sites without colour do not support the VO

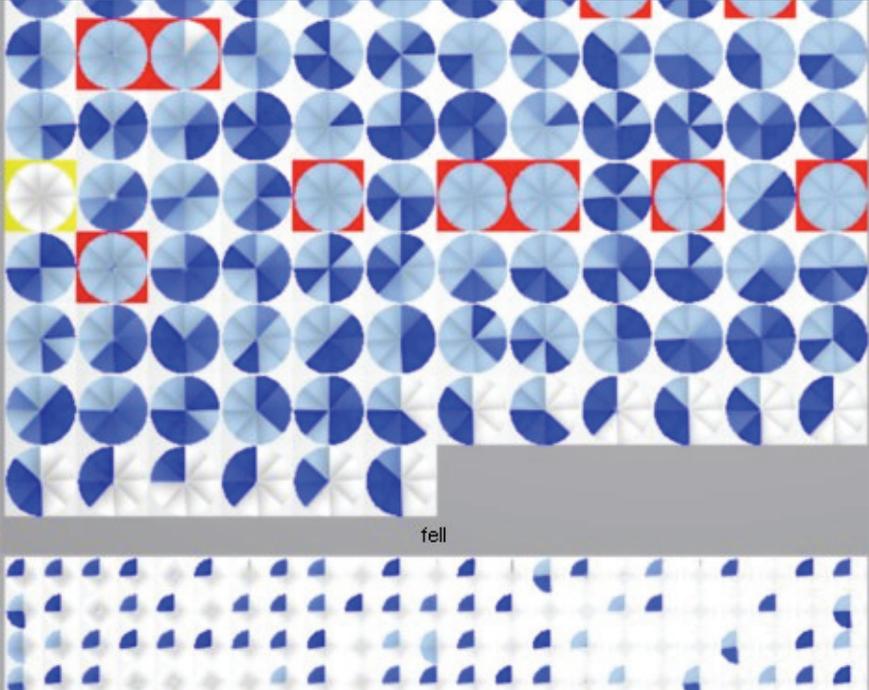
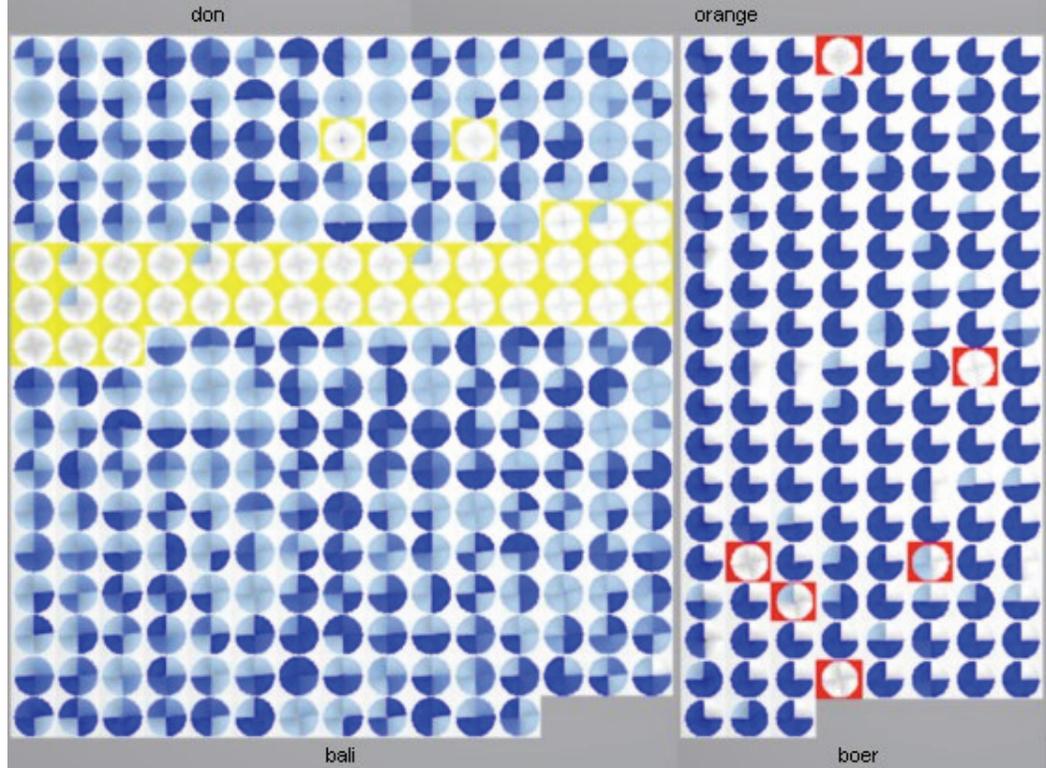
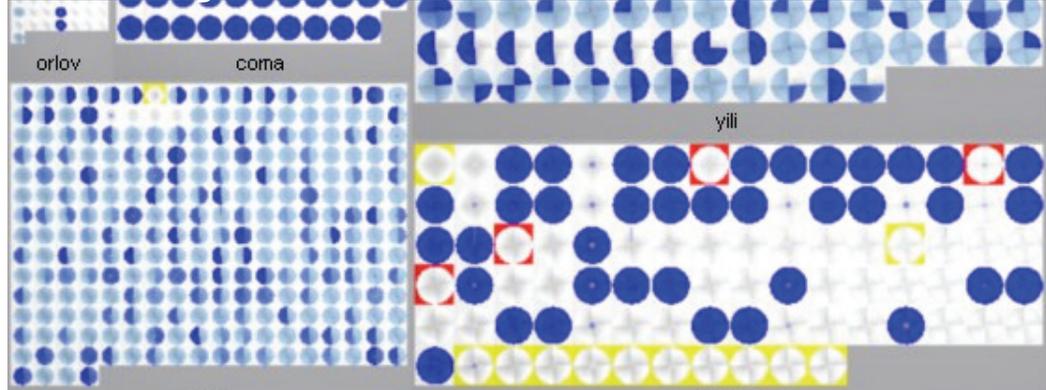


Status of different Site Services:



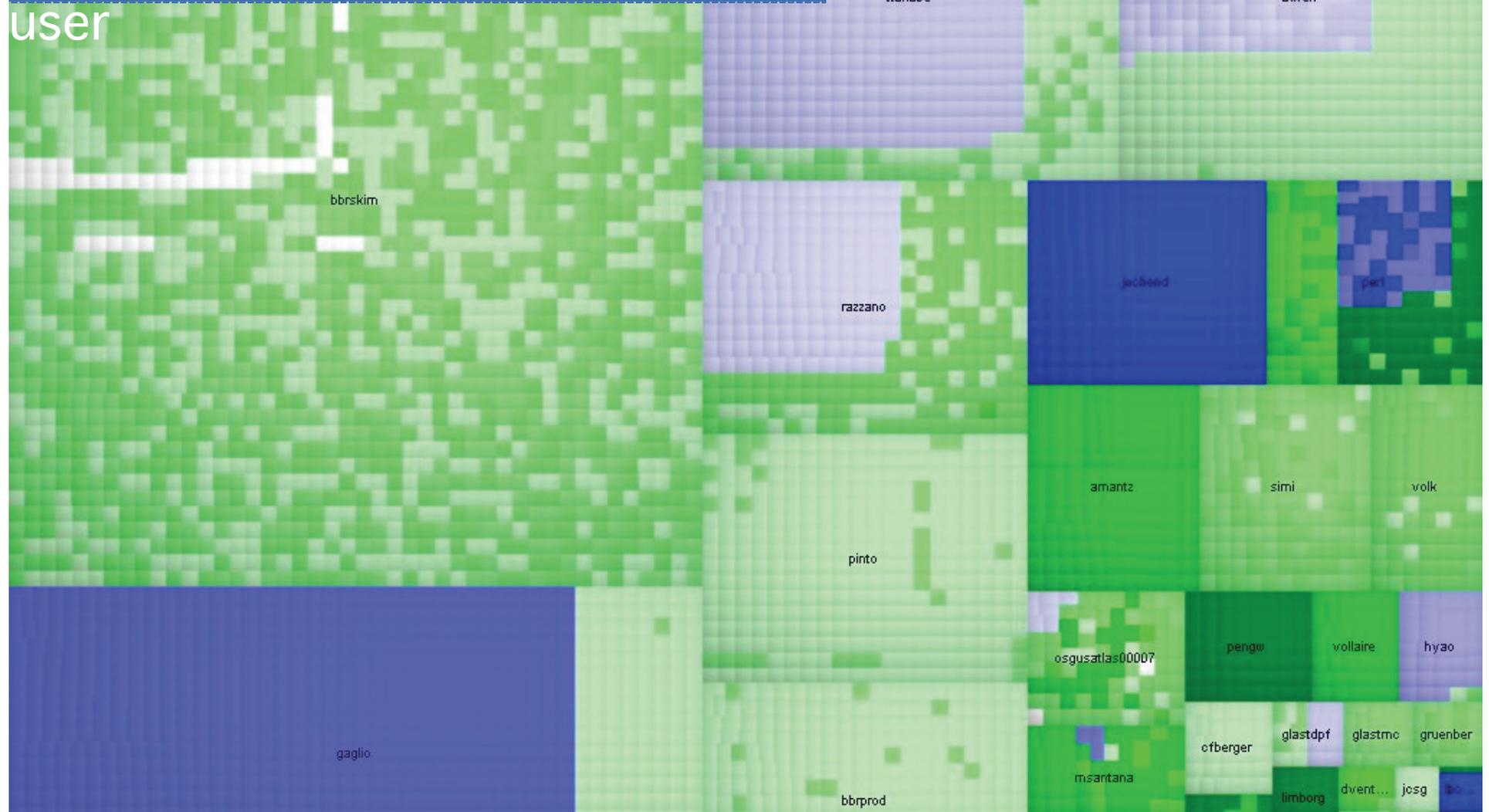


Similar idea being introduced at SLAC to aid understanding of local system behaviour



Status of batch system
 Size of node
 Size of cluster
 Job distribution
 Efficiency of job
 Status of node

Tree Map of Users
identifies large users
correlates queue/run-time with
user

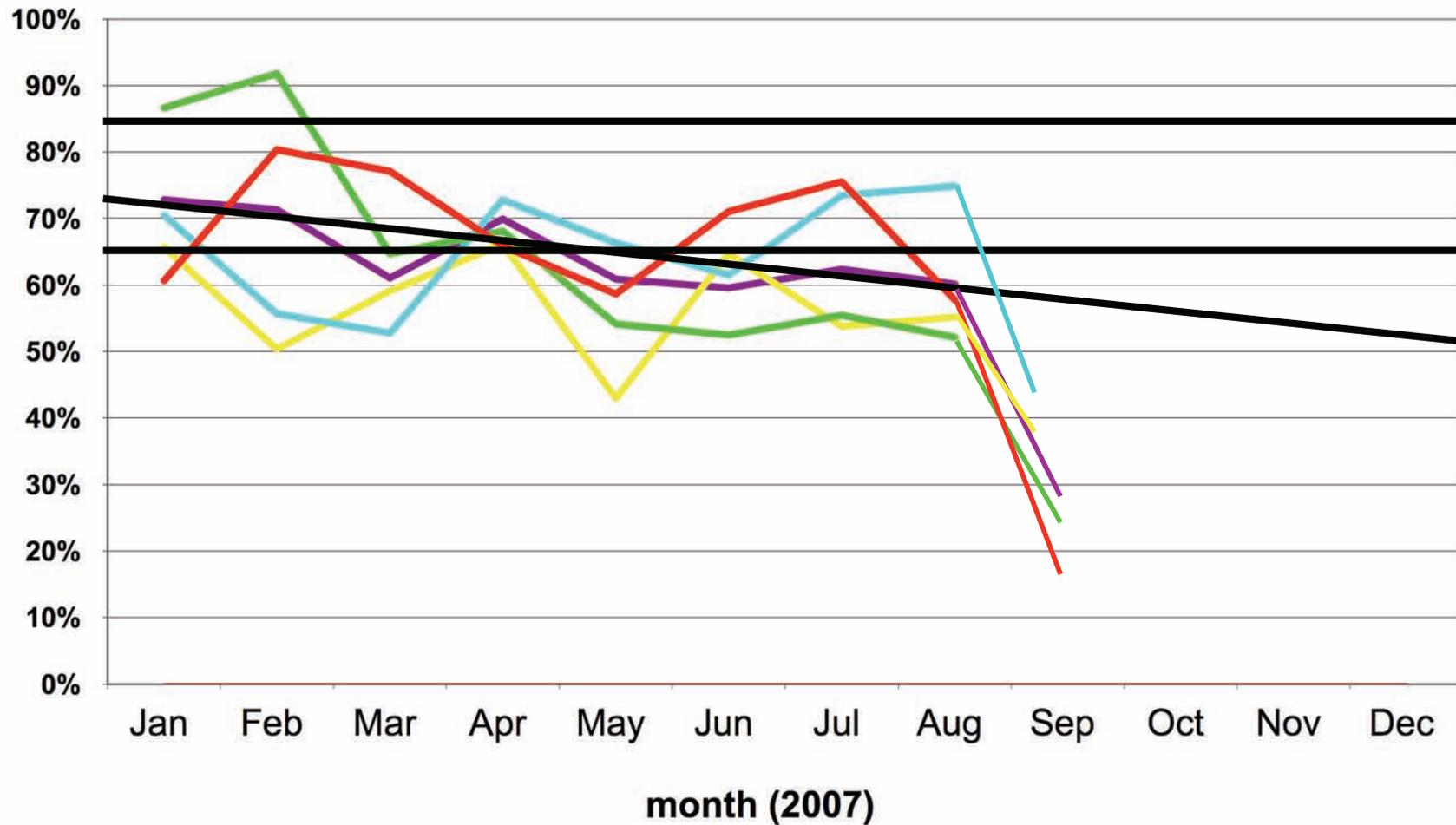


Uses Magnaview, a commercial data exploration tool (www.magnaview.nl)
See also



So, how to solve this?

Ratio of CPU : Wall_clock Times





Outline

- Introduction to CERN and Experiments
- LHC Computing
- Challenges
 - _ Capacity Provision
 - _ Box Management
 - _ Data Management and Distribution
 - _ What' s Going On?
- **Summary/Conclusion**



Summary

- Immense Challenges & Complexity
 - _ Data rates, developing software, lack of standards, worldwide collaboration, ...
- Considerable Progress in last ~5-6 years
 - _ WLCG service exists
 - _ Petabytes of data transferred
 - _ ...
- But real data is nearly here...
 - _ Will the system cope with chaotic analysis?
 - _ Will we understand the system enough to identify problems—and fix underlying causes?
 - _ Major “Dress Rehearsals” in Feb & May 2008
 - last chance to shake system down before operation
- Answer(s) at LISA ’ 08?

Thank You!