



The LHC Computing Grid

The Computing Environment for LHC Data Analysis

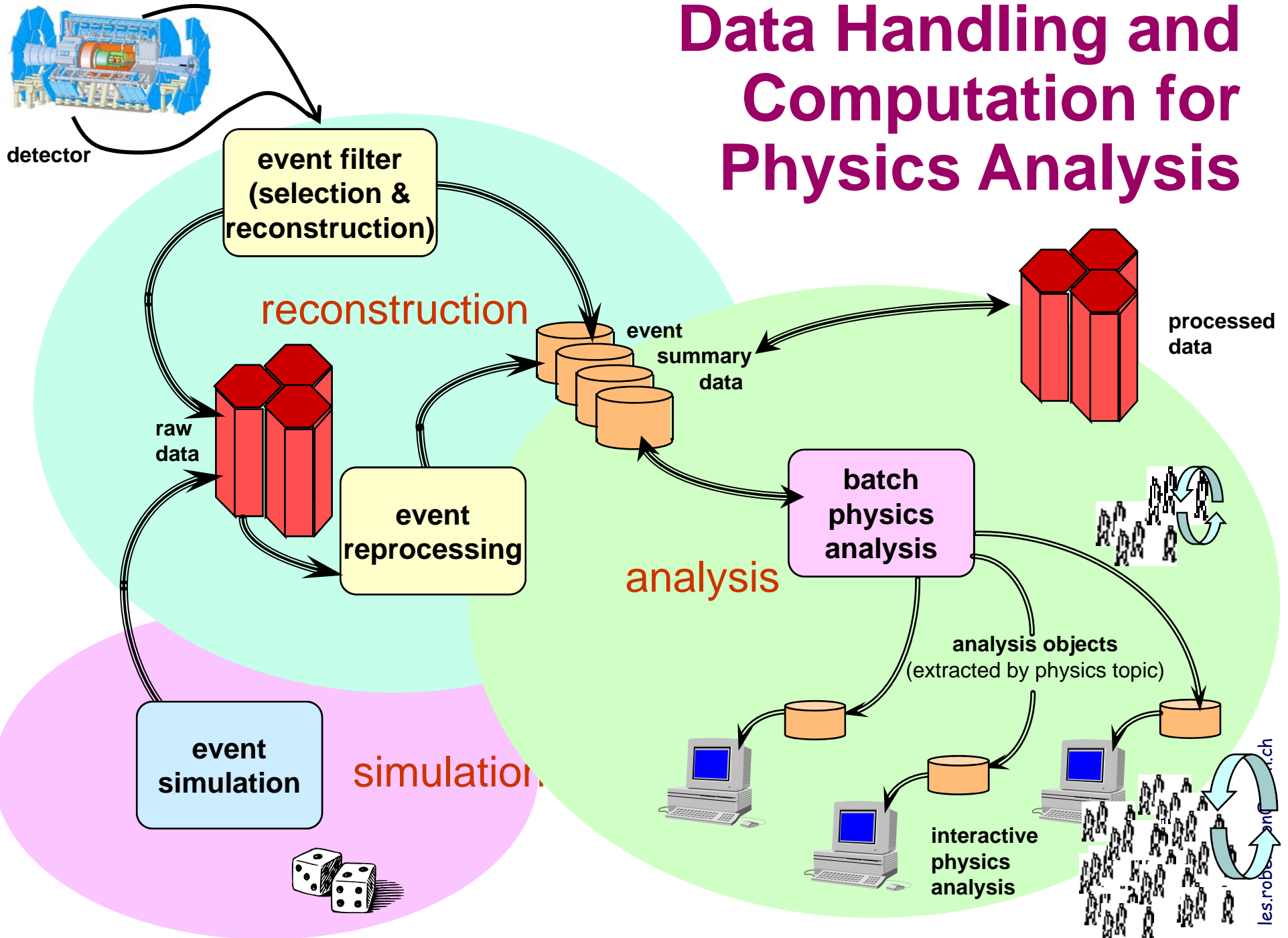
CHIPP Phoenix Cluster
Inauguration

Manno, Switzerland
30 May 2008

Les Robertson
IT Department - CERN
CH-1211 Genève 23



Data Handling and Computation for Physics Analysis





The Computing System for LHC DATA ANALYSIS

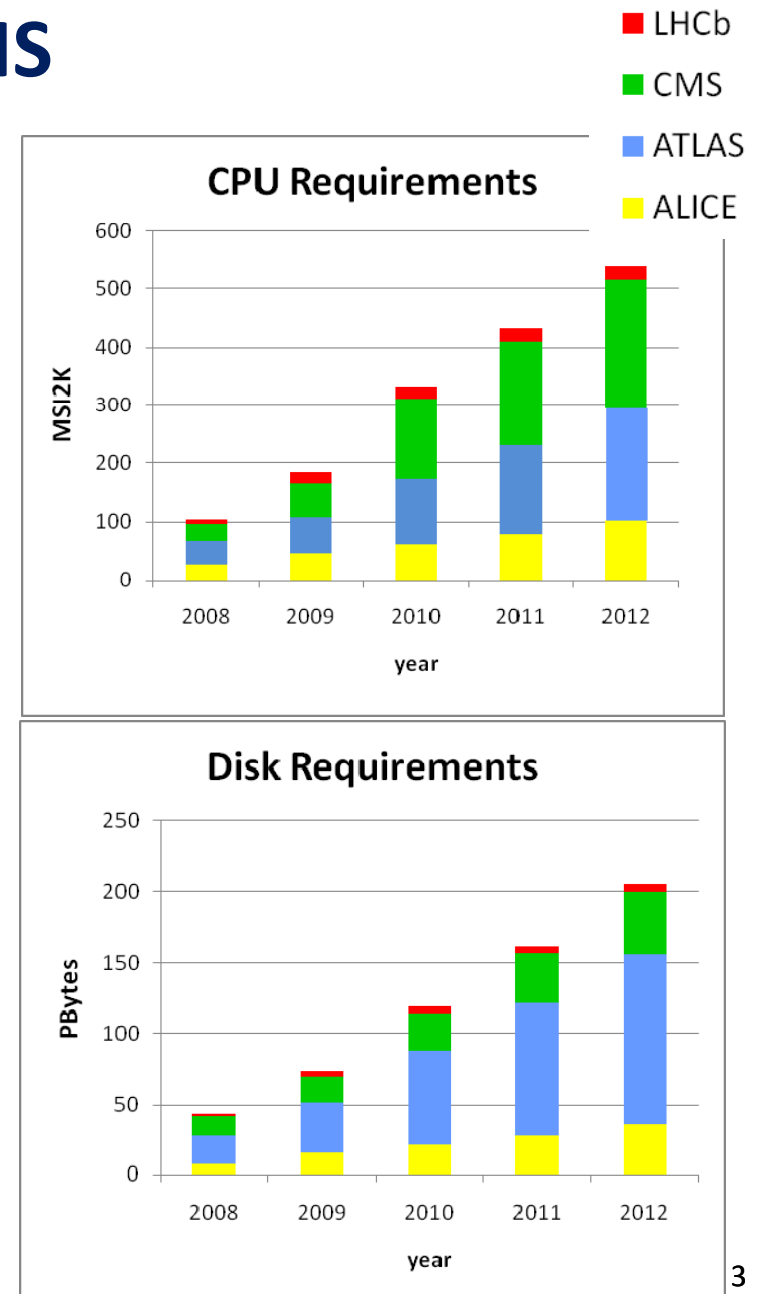
Capacity and Evolution

Computing requirements for all four experiments in first full year (2009)

- ~70 PetaBytes disk
- ~100K processor cores (2009)
- used by > 5,000 scientists & engineers

Growth driven by new data, accelerator enhancements and improving efficiency, new analysis techniques, expanding physics scope,

- disk storage growing at ~40 PB/year
- BUT – evolution of access patterns unclear**



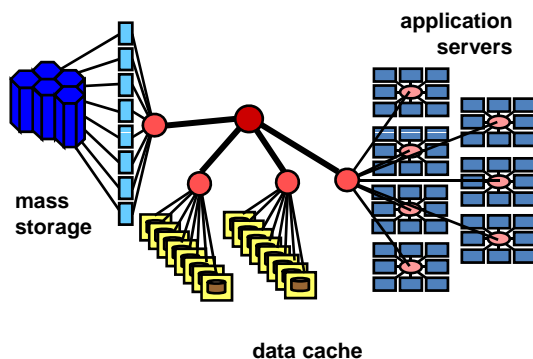


The key characteristics of experimental HEP data analysis that dictate the design of the computing system

- independent events
→ **easy parallelism**
- codes have
 - modest memory needs (~2GB)
 - modest floating point content→ **perform well on PCs**

BUT --

- enormous data collections
→ **PetaBytes of new data every year**
- **shared by very large user collaborations**, many different groups, independent approaches to analysis
→ **unpredictable data access patterns**



- a simple distributed architecture developed ~1990 enabled experimental HEP to migrate from supercomputers and mainframes to clusters
- with the flexibility for easy evolution to new technologies
- and benefit from the mass market driven growth in the performance and capacity of PCs, disks, and local area networking



Why did we decide on a geographically distributed computing system for LHC?

- CERN's budget for physics computing was insufficient
- Easy parallelism, use of simple PCs, availability of high bandwidth international networking make it *possible* to extend the distributed architecture to the wide area

AND

- The ~5,000 LHC collaborators are distributed across institutes all around the world with access to local computing facilities, ...
... and funding agencies prefer to spend at home if they can
- Mitigates the risks inherent in the computing being controlled at CERN, subject to the lab's funding priorities and with access and usage policies set by central groups within the experiments

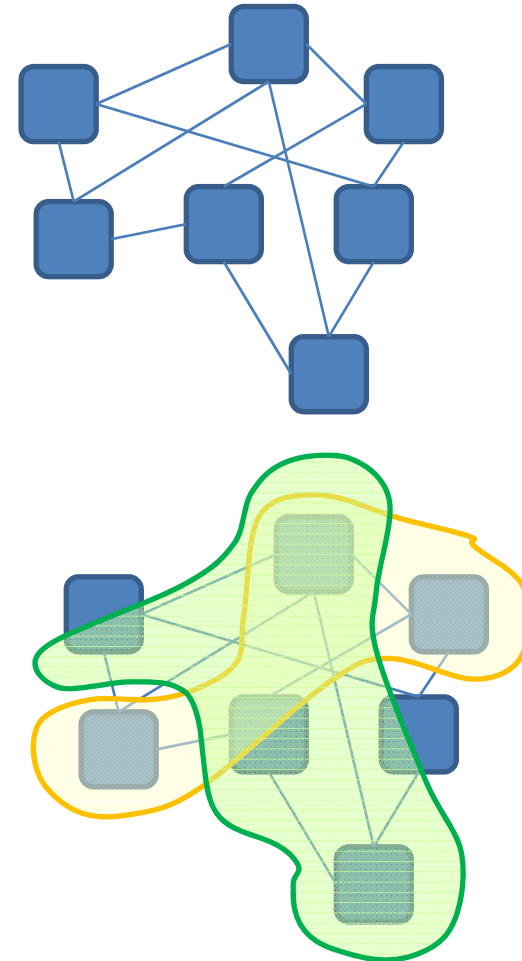
ALSO

- Active participation in the LHC computing service gives the institute (not just the physicist) a continuing and key role in the data analysis
-- which is where the physics discovery happens
- Encourages novel approaches to analysis
... and to the provision of computing resources



What do we mean by a Computing Grid**?

- Collaborating computing centres
- Interconnected with good networking
- Interfaces and protocols that enable the centres to advertise their resources and exchange data and work units
- Layers of software that hide all the complexity from the user
- So the end-user does not need to know where his data sits and where his jobs run
- The Grid does not itself impose a hierarchy or centralisation of services
- Application groups define *Virtual Organisations* that map users to subsets of the resources attached to the Grid



** There are many different variations on the term Grid – this is the HEP definition



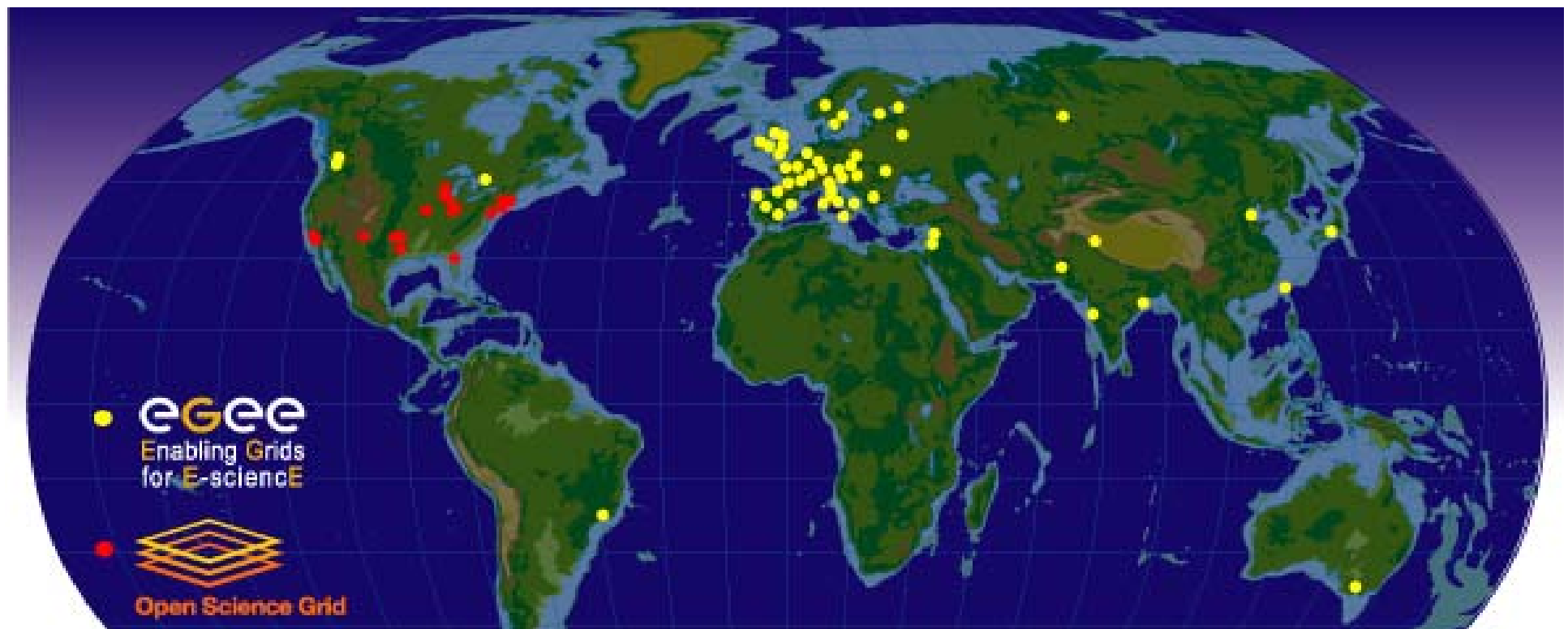
- The advantage for the *computer centre* is that the basic services can be provided in a standard way for different application groups
 - e.g. user authentication, job submission, storage access, data transfer...
 - ATLAS, CMS, LHCb, DZERO,, BioMed, Fusion,
- The advantage for the *application group* is that it can integrate resources from different centres and view them as a single service without having to support all of the software layers, negotiate the installation of special software, register users on each site, etc.
- But they have the *flexibility* to pick and choose - replace software layers with their own products, decide which services are provided at which sites,



LCG depends on two major science grid infrastructures

EGEE - Enabling Grids for E-Science (*with EU funding*)

OSG - US Open Science Grid (*with DoE and NSF funding*)





The Middleware for the Baseline Services needed for the LHC Experiments

- Information system
- Security framework
- Storage Element
- SRM interface to Mass Storage
dCache, DPM, CASTOR, STORM
- Basic data transfer tools -
Gridftp, srmCopy.
- Reliable file transfer service -
FTS
- Catalogue services -
LFC, Globus RLS
- Catalogue and data management tools - lcg-utils
- Compute element -
Globus/Condor-G based CE,
Cream (web services)
- Reliable messaging service
- Virtual Organisation Management Services
- Database distribution services -
ORACLE streams, SQUID
- POSIX-I/O interfaces to storage
- Workload Management -
EGEE Resource Broker,
VO-specific schedulers
- Job monitoring tools
- Grid monitoring tools
- Application software installation
- GUIs for analysis, production
GANGA, CRAB, PANDA, ..

For LCG, grid interoperability is required at the level of the baseline service
→ same software or standard interfaces or compatible functionality



The Middleware for the Baseline Services needed for the LHC Experiments

- Information system
- Security framework
- Storage Element
- SRM interface to Mass Storage -
dCache, DPM, STORM
- Basic data transfer -
Gridftp, sftp
- Reliable file transfer -
FTS
- Catalogue services -
LFC, Globus RLS
- Catalogue and data management tools - lcg-utils
- Compute element -
Globus/Condor-G based CE,
Crab (web services)
- Messaging service
- Organisation Management
services
- Distribution services -
XCache streams, SQUID
- UNIX-I/O interfaces to storage
- Workload Management -
EGEE Resource Broker,
VO-specific schedulers
- Job monitoring tools
- Grid monitoring tools
- Application software installation
- GUIs for analysis, production
GANGA, CRAB, PANDA, ..

GRID - a simple idea -
Not so simple to
implement it!

For LCG, grid interoperability is required at the level of the baseline service
→ same software or standard interfaces or compatible functionality



The LHC Computing Grid - A Collaboration of 4 Experiments + ~130 Computer Centres

Tier-0 - the accelerator centre

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



11 Tier-1 Centres - "online" to the data acquisition process → high availability

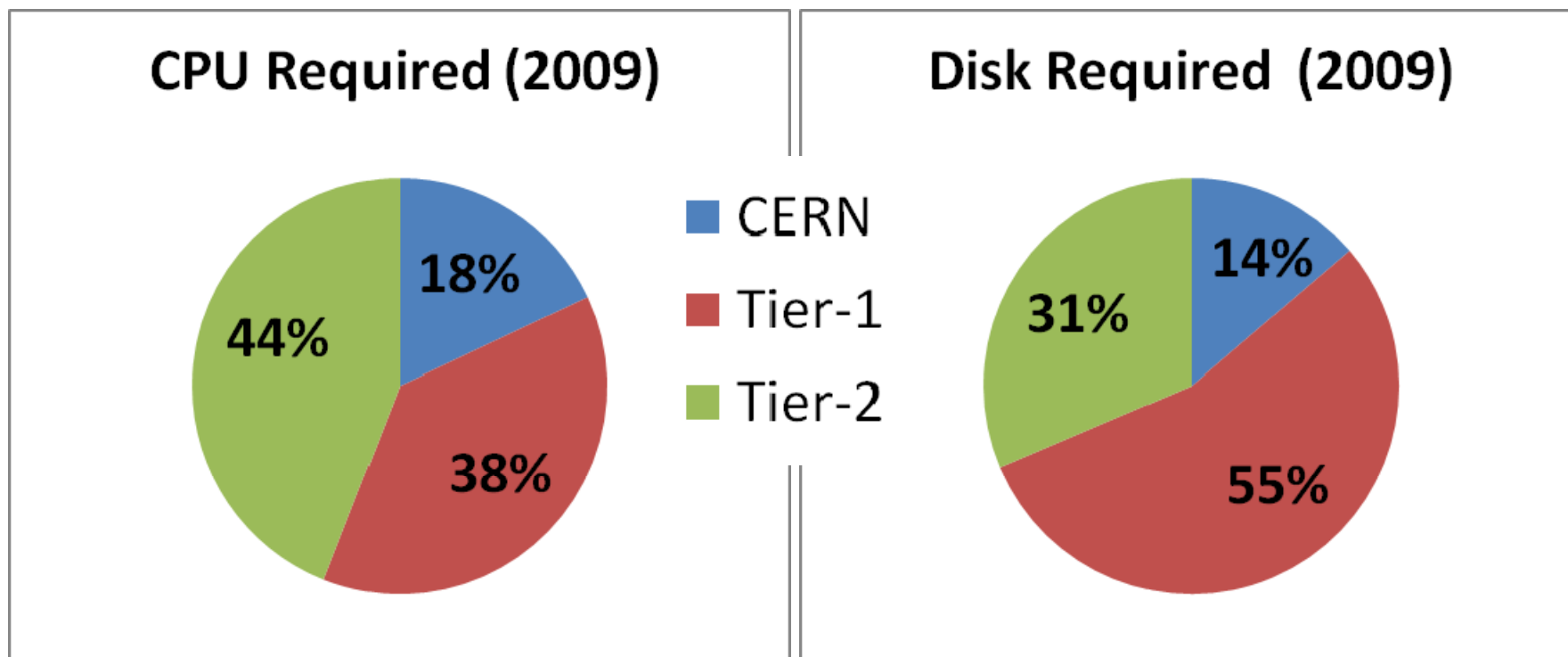
- Managed Mass Storage -
→ grid-enabled data service
- Data-heavy analysis
- National, regional support

Tier-2 - 120 Centres in 60 Federations in 35 countries

- **End-user (physicist, research group) analysis** -
where the discoveries are made
- Simulation



Distribution of Resources across Tiers



- 2009 – the first full year of data taking
- <20% at CERN – less than half of that at the Tier-2s
→ the distributed system must work from Day 1

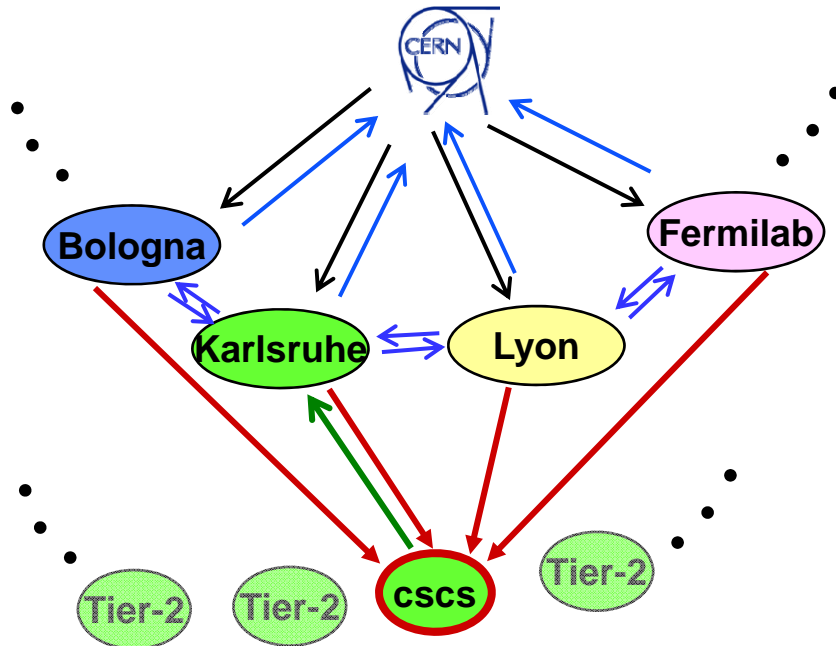


Middleware & Software

- From many sources: Globus, Condor, EGEE gLite, High Energy Physics common tools and packages, experiments, open-source projects, proprietary packages, ...
- Two fundamental middleware packages are integrated, tested and distributed by the infrastructure grids
 - gLite by EGEE
 - built on the Virtual Data Toolkit by OSG
- The mass storage systems are crucial but complicated components of the LCG service - HEP developments
- And a thick layer of software is maintained and distributed by the experiments (data management, resource scheduling, analysis GUIs, ..)

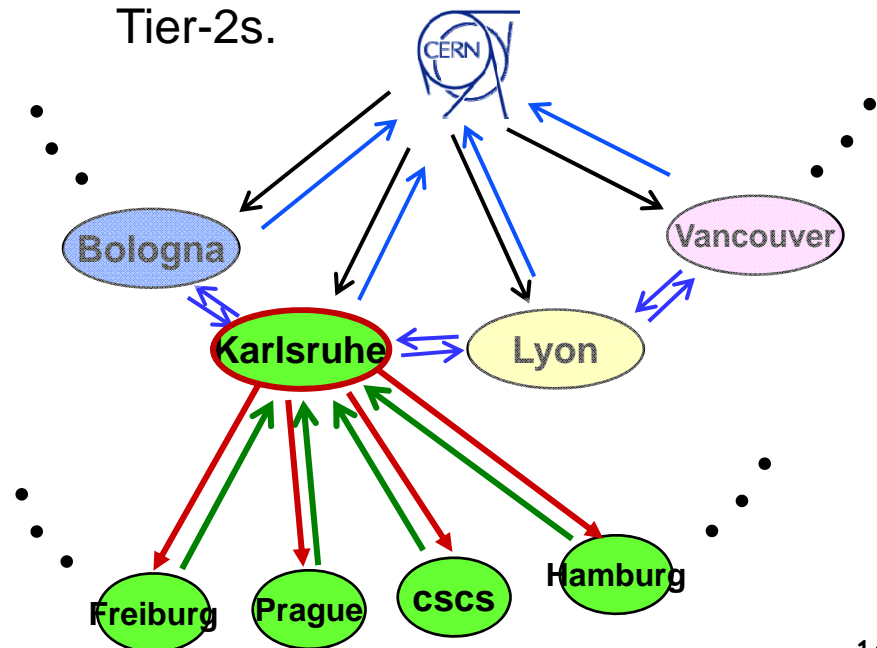


Experiment computing models define specific data flows between CERN, Tier-1s and Tier-2s

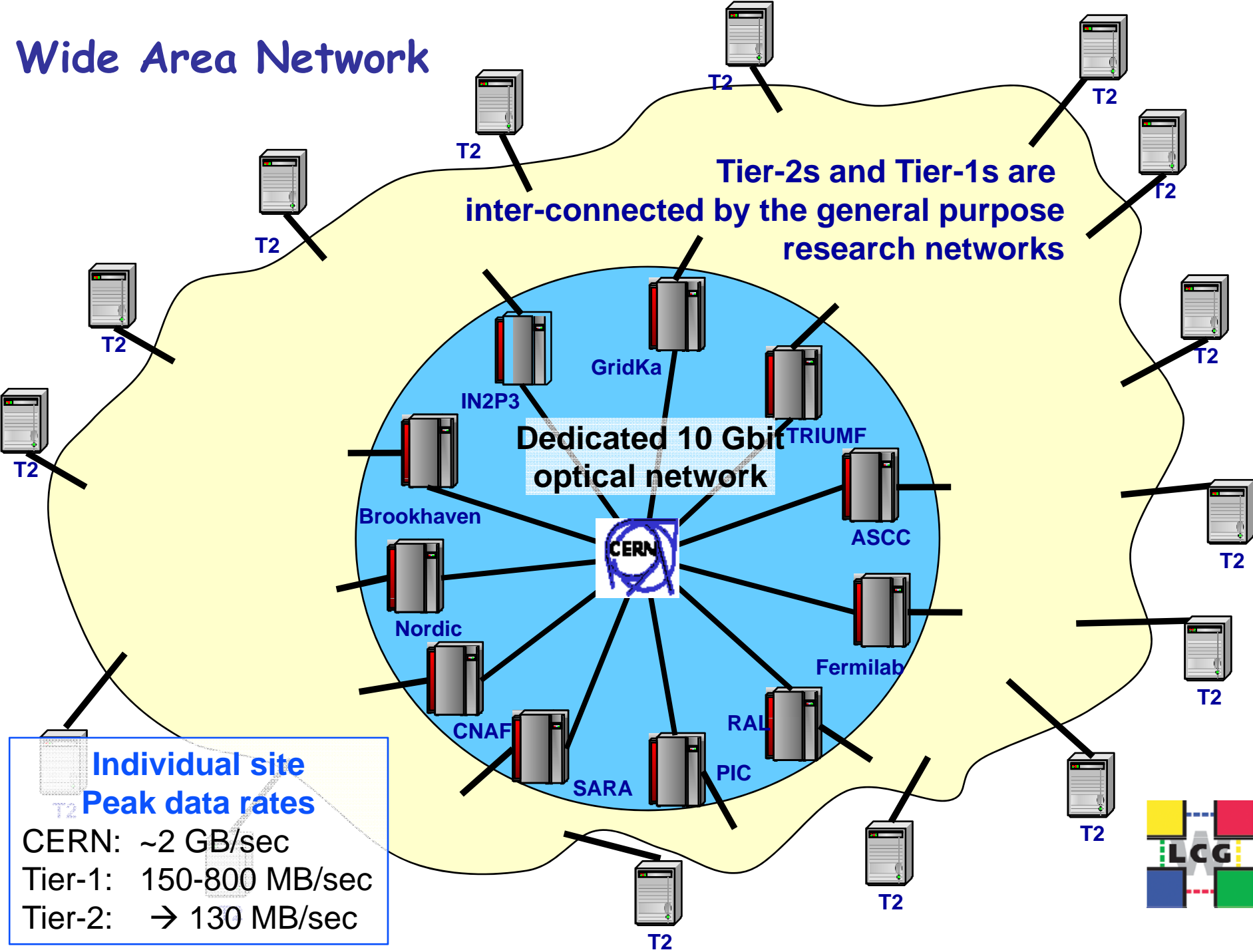


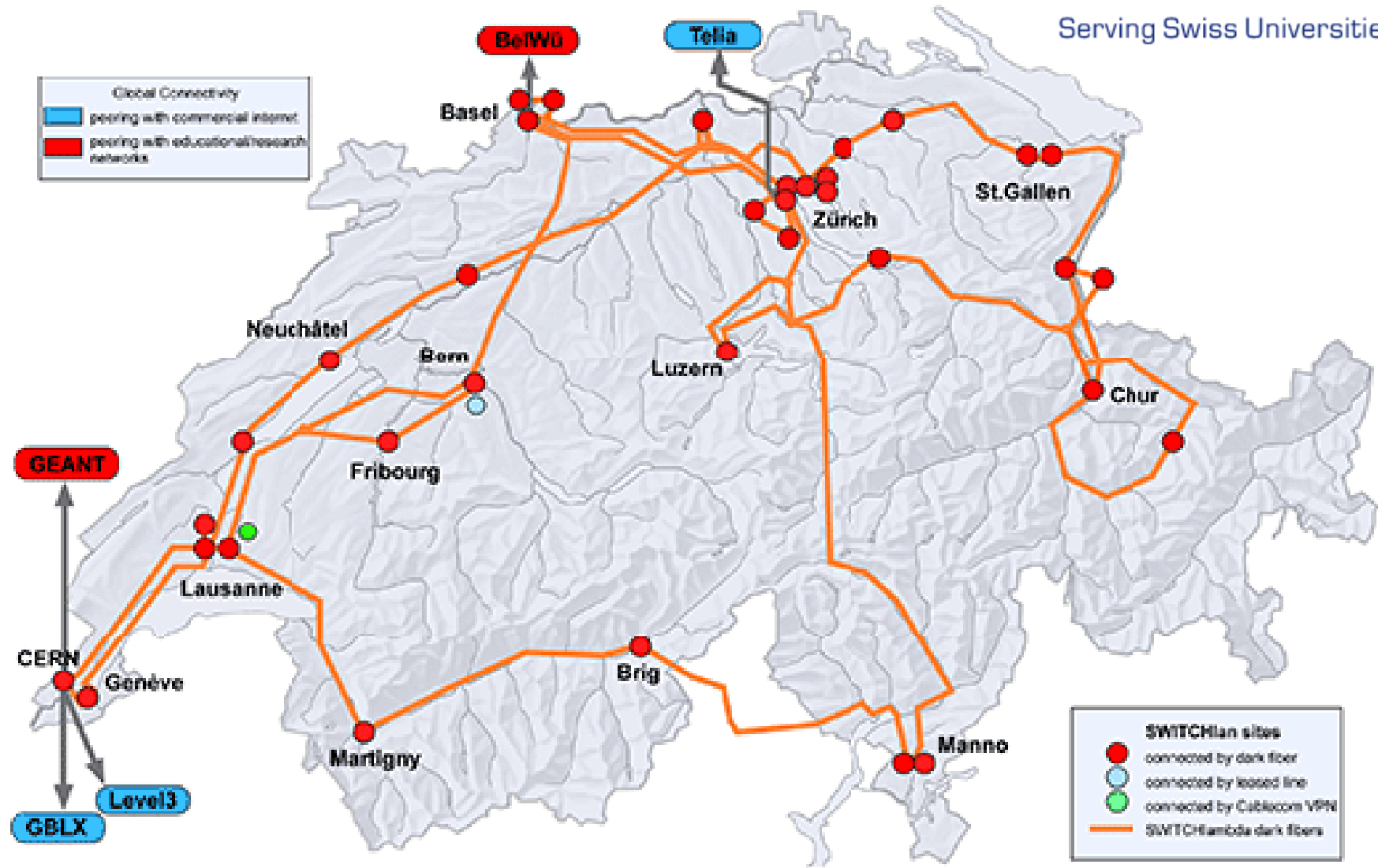
CMS – Tier-2s send simulated data to a specific Tier-1, but obtain data for analysis from any of the 7 Tier-1s:
Taipei, Bologna, Fermilab, Karlsruhe, Lyon, Barcelona, Rutherford Lab

ATLAS – Each Tier-1 acts as the data repository for a “cloud” of associated Tier-2s.



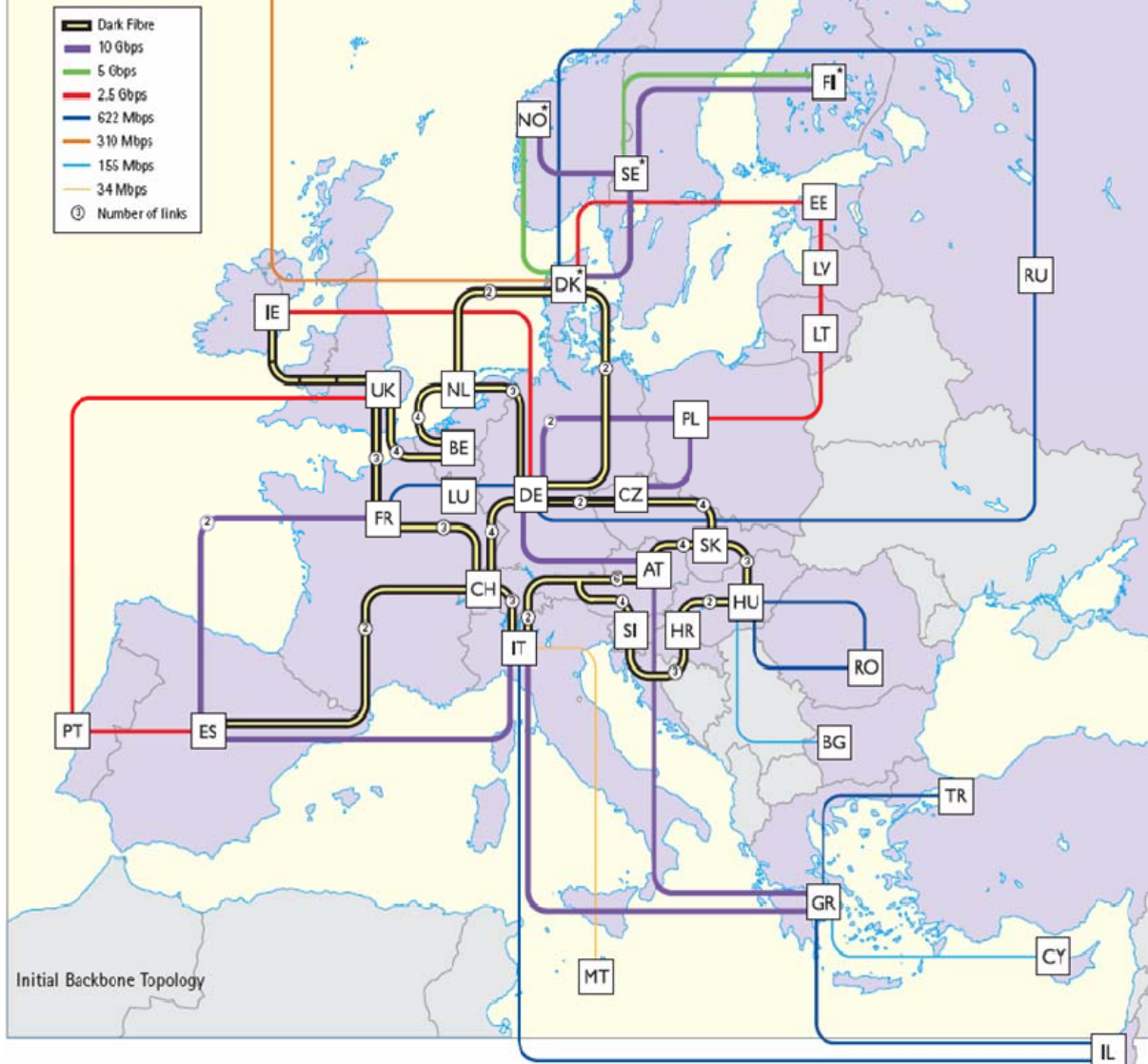
Wide Area Network







European Research Network Backbone



GEANT2 is operated by DANTE on behalf of Europe's NRENs.

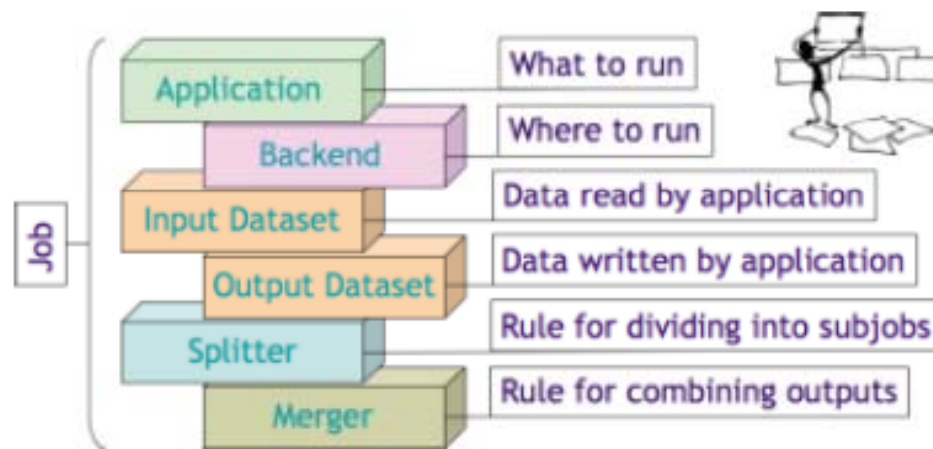
SWITCH

Serving Swiss Universities



Topographical map: ©2004 swissmap

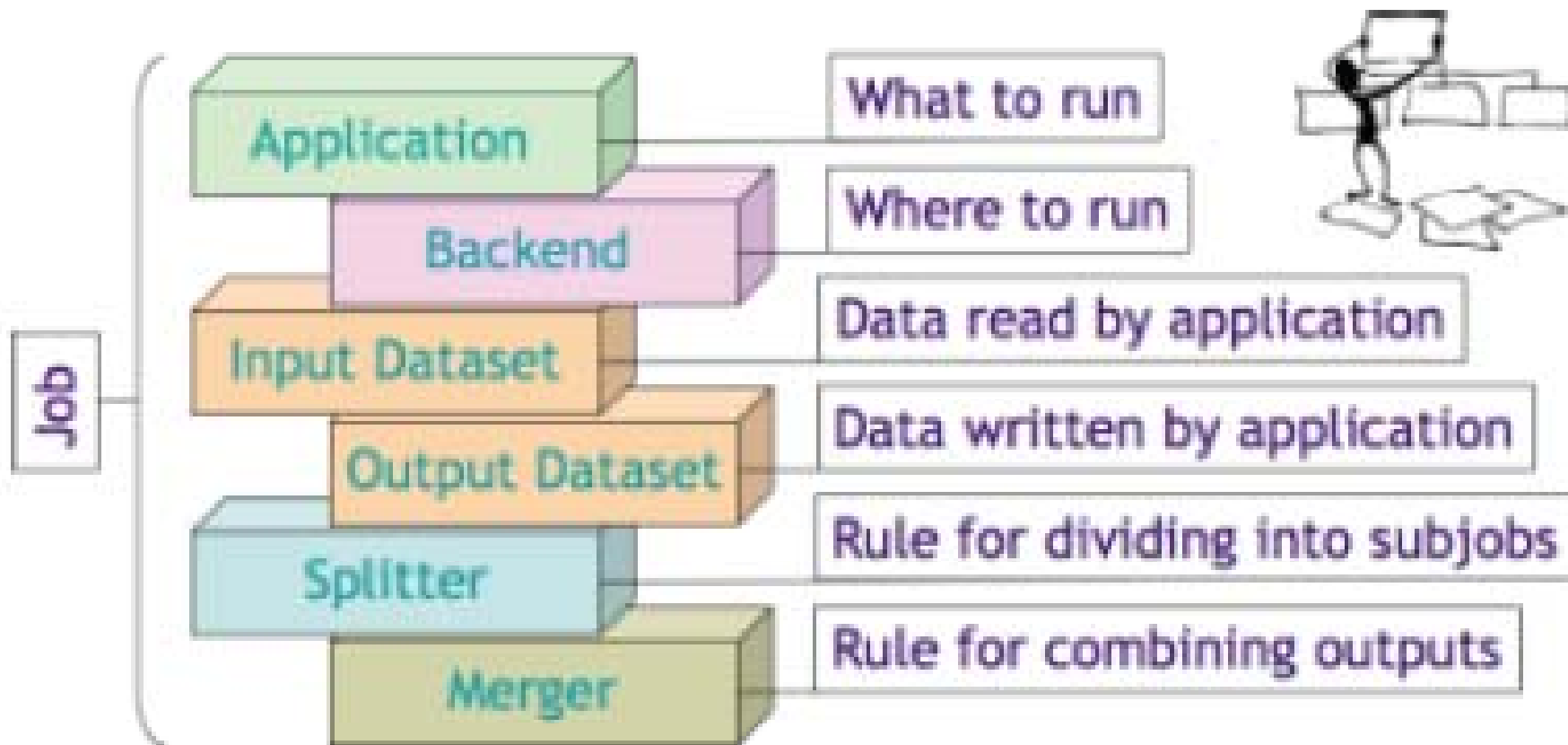
- Grid user interface: **GANGA**
 - ▶ Single tool for all Grid-based work (including analysis and small Monte Carlo productions)
 - ▶ Trivial switching between Grid running and local execution (for testing purposes)
 - ▶ Grid backends include LCG (EGEE), NorduGrid, OSG (Panda)
 - ▶ Interface via a command-line or a GUI



www.cern.ch/ganga



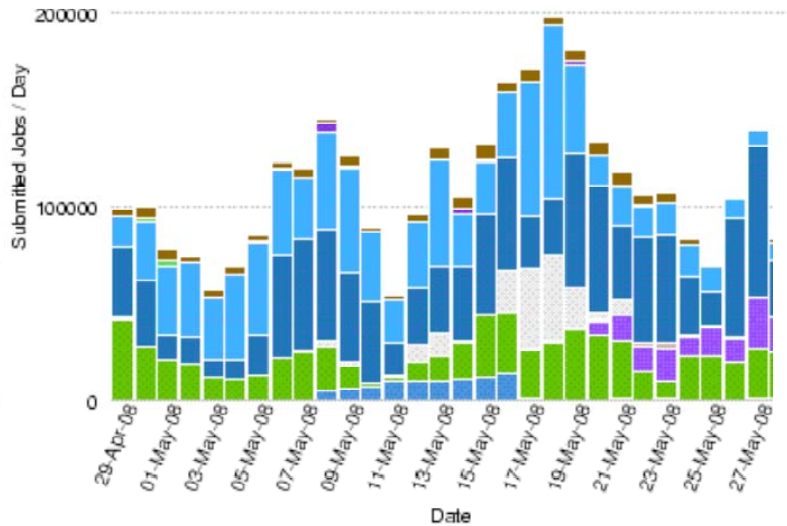
- Grid user interface: **GANGA**
 - Single tool for all Grid based work (including analysis and small Monte Carlo)





Experiment Dashboards provide tools for monitoring and debugging

CMS jobs May 2008 sorted by activities
Up to 200K jobs per day - 35% end-user analysis



ATLAS data transfer status - 28 May 2008 Throughput ~1100 MB/s



Data: Tier 0 Data: Production Jobs: Production Jobs: Analysis Panda: Production

Overview Dataset Info Page Help User Guide Feedback

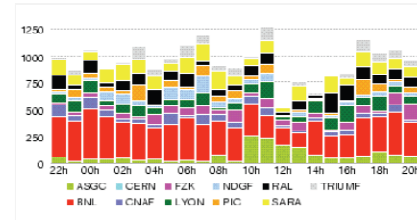
OVERVIEW

- OVERVIEW Activity
- Activity in Last Hour
- Activity in Last 4 Hours
- Activity in Last 24 Hours
- Activity in Last 7 Days
- Activity in Last 30 Days
- Activity in ...

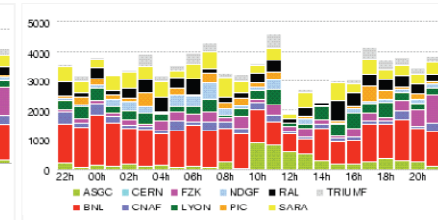
Cloud Activity

- ASGC Cloud
- BNL Cloud
- CERN Cloud
- CNAF Cloud
- FZK Cloud
- LYON Cloud
- NDGF Cloud
- PIC Cloud
- RAL Cloud
- SARA Cloud
- TRIUMF Cloud

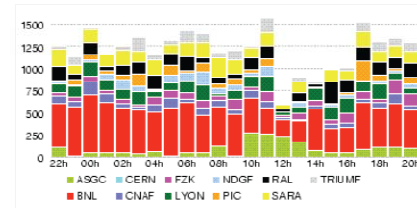
Throughput (MB/s)



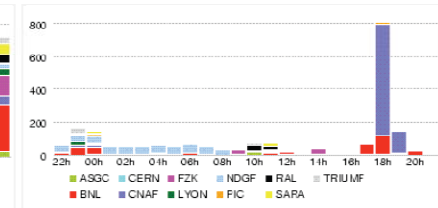
Data Transferred (gbytes)



Completed File Transfers

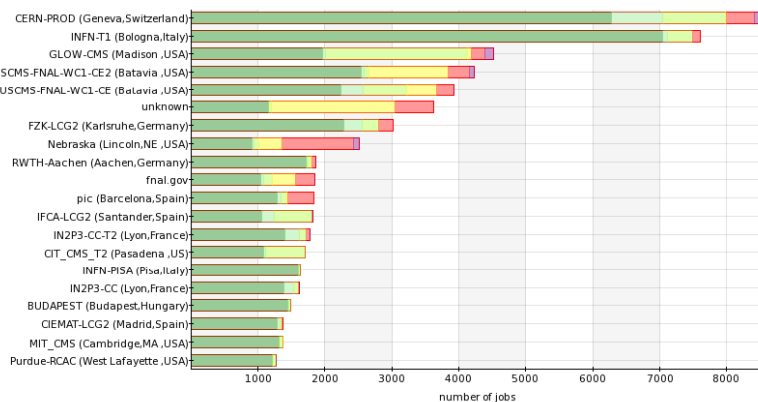


Total Number Transfer Errors



Activity Summary (Last 24 Hours)
Click on the cloud name to view list of sites

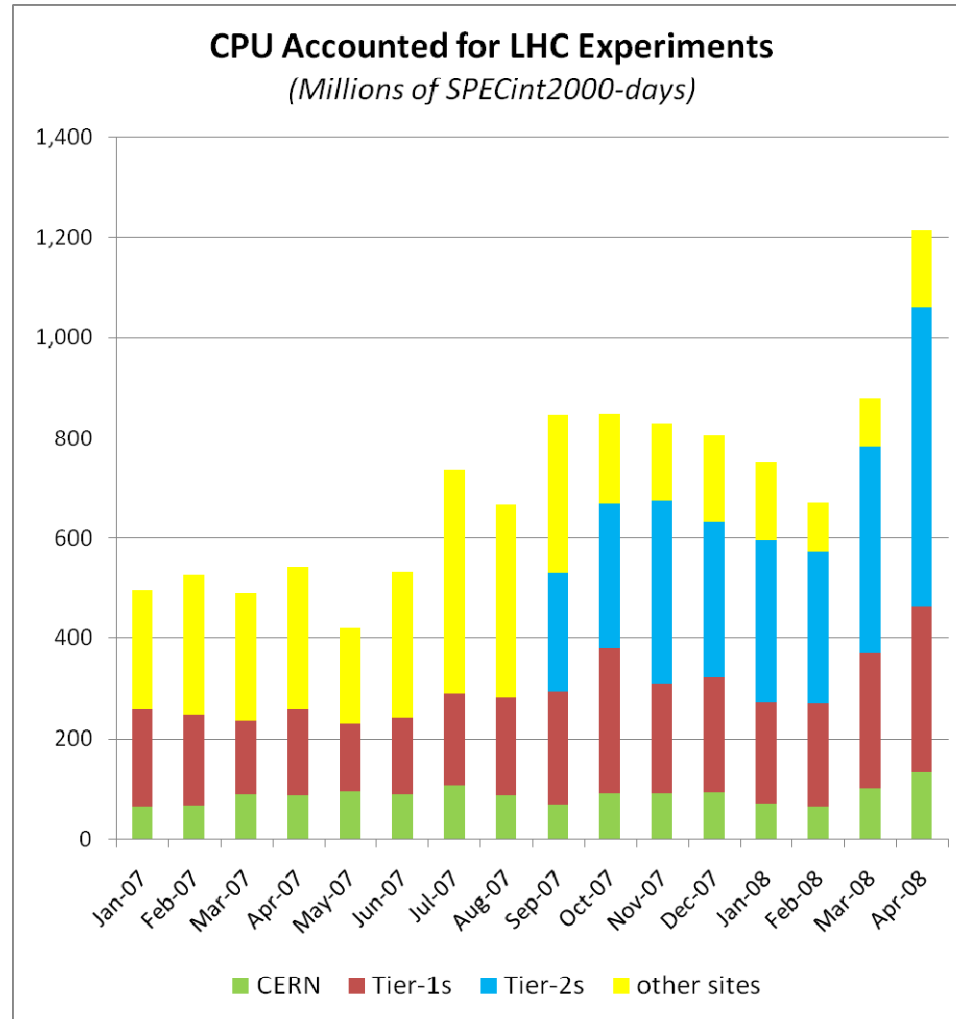
[suggest](#) **job status by site – user can drill down to find details of the errors**



submitted app-succeeded app-failed app-unknown pending running aborted cancelled



Evolution of LHC grid cpu usage



150% growth since Jan 07

More than half from Tier-2s

~800K core-days

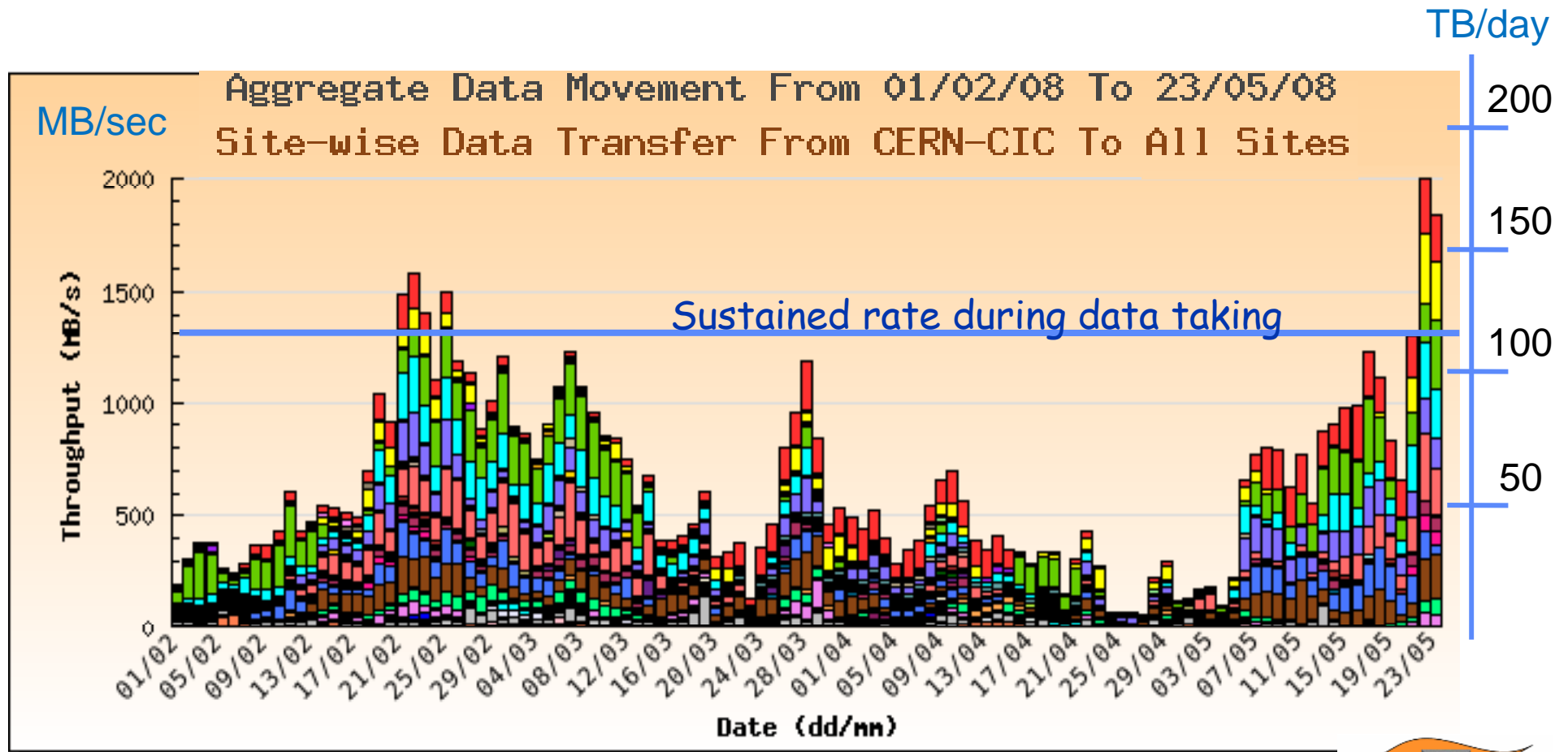
~25K cores at 100% utilisation

~55% of committed capacity

Experiment services are still in test mode - awaiting the real data



Data Transfer CERN→Tier-1s





CMS Data Volume TBytes/day - all sites to all sites



PhEDEx - CMS Data Transfers

[Info](#) [Activity](#) [Data](#) [Requests](#) [Components](#) [Reports](#)

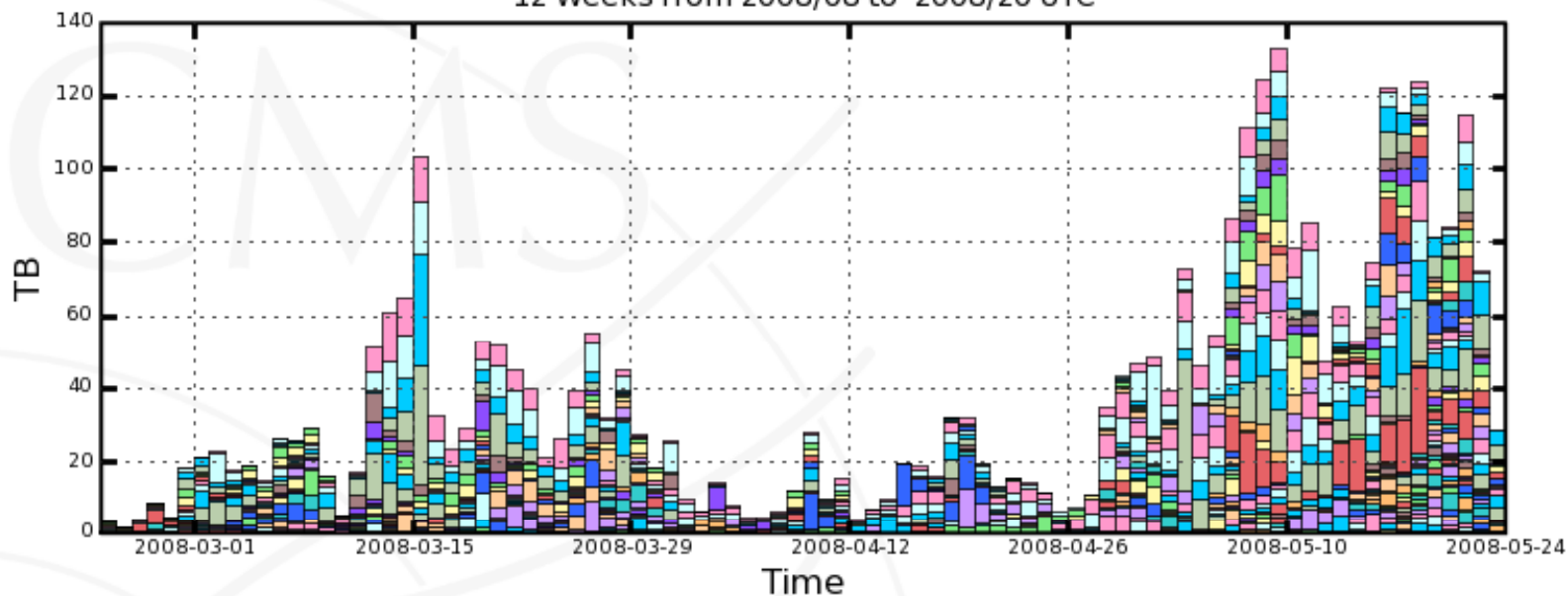
[Rate](#) [Rate Plots](#) [Queue Plots](#) [Quality Plots](#) [Routing](#) [Transfer Details](#) [Deletions](#) [Recent Errors](#)

Graph by filter source destination show MSS nodes

Period up to

CMS PhEDEx - Transfer Volume

12 Weeks from 2008/08 to 2008/20 UTC





LHC Computing Grid Status Summary

- The “final” pre-startup testing is now going on
→ all experiments exercise simultaneously their full computing chains - from the data acquisition system to end-user analysis at the Tier-2s - at the full 2008 scale
- No show-stoppers since before the tests began in February - day-to-day issues being rapidly resolved
- Most performance and scaling metrics have been achieved

BUT this is research

- The actual load, access patterns, user behaviour is unpredictable - depends on how physicists react to what they find in the real data
- We can look forward to an exciting time when the beam starts!!



Summary

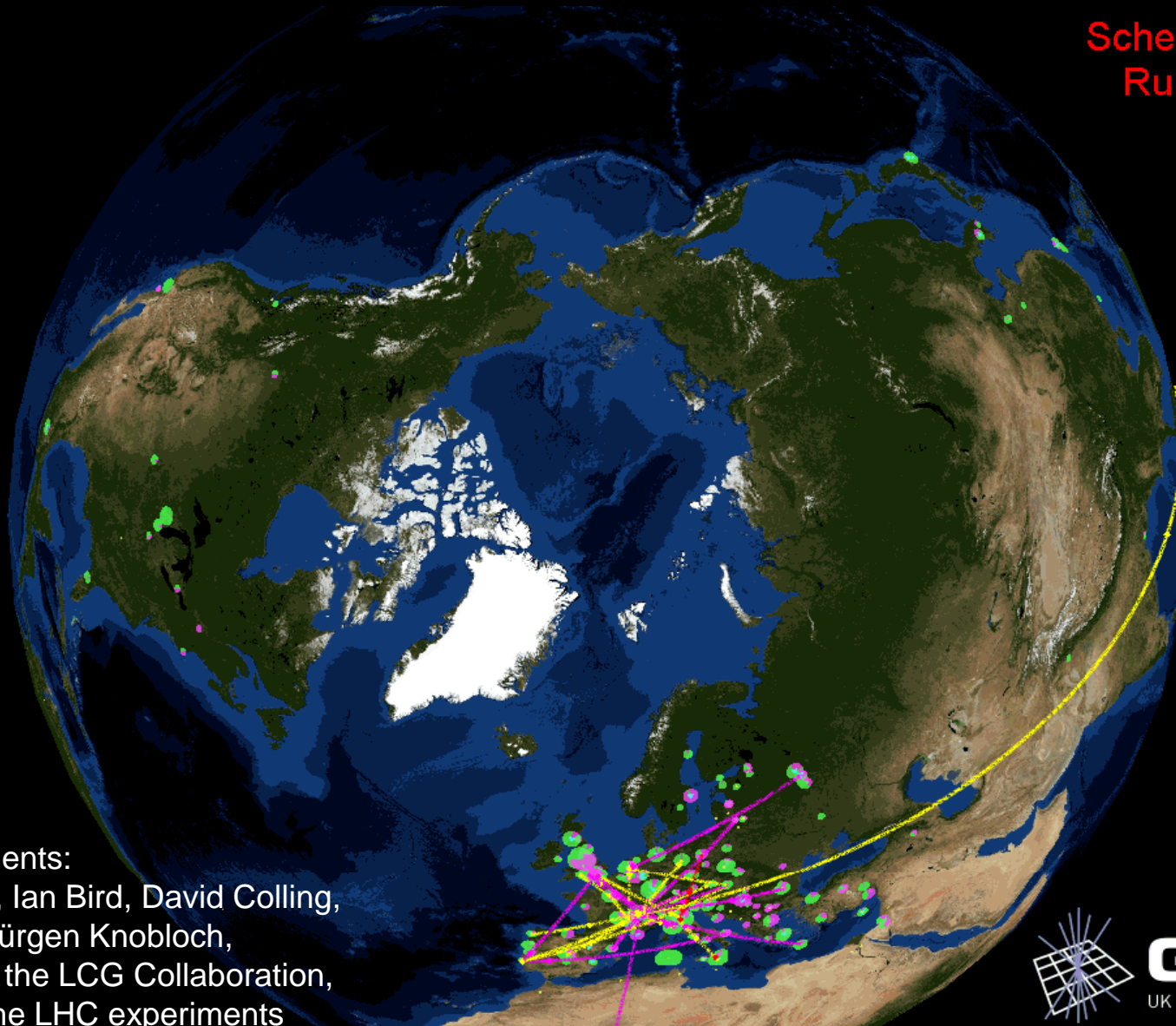
- **Grids are all about sharing**
 - they are a means whereby groups distributed around the world can pool their computing resources
 - large centres and small centres can all contribute
 - users everywhere can get equal access to data and computation
 - without having to spend all of their time seeking out the resources
 - **Grids also allow the flexibility to place the computing facilities in the most effective and efficient places -**
 - exploiting funding wherever it is provide,
 - piggy-backing on existing computing centres,
 - or exploiting cheap and renewable energy sources
 - **The LHC provides a pilot application -**
 - with massive computing requirements, world-wide collaborations
 - that is already demonstrating that grids can deliver in production
- and the scientific success of LHC will depend on the grid from day 1**



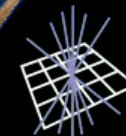
Computing at the Terra-Scale

eGEE
Enabling Grids
for E-science

Scheduled = 21539
Running = 25374



Acknowledgements:
Julia Andreeva, Ian Bird, David Colling,
David Foster, Jürgen Knobloch,
Fairouz Malek, the LCG Collaboration,
EGEE, OSG, the LHC experiments



GridPP
UK Computing for Particle Physics