

LHC Computing Grid today

Did it work ?

Sept. 9th 2011, Günter Quast

Institut für Experimentelle Kernphysik

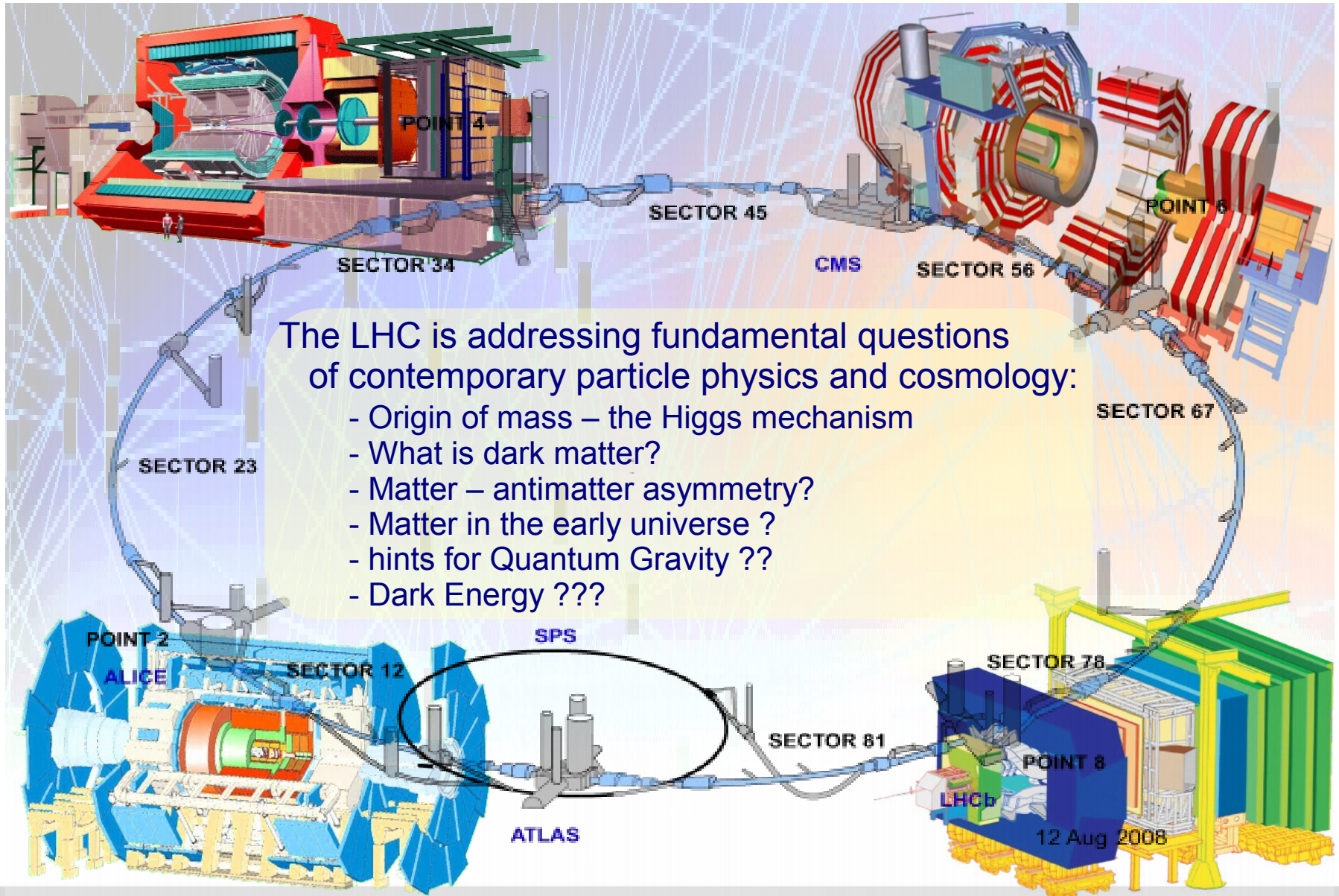


GridKa
School

9th International

GridKa School 2011

Large Hadron Collider and Experiments



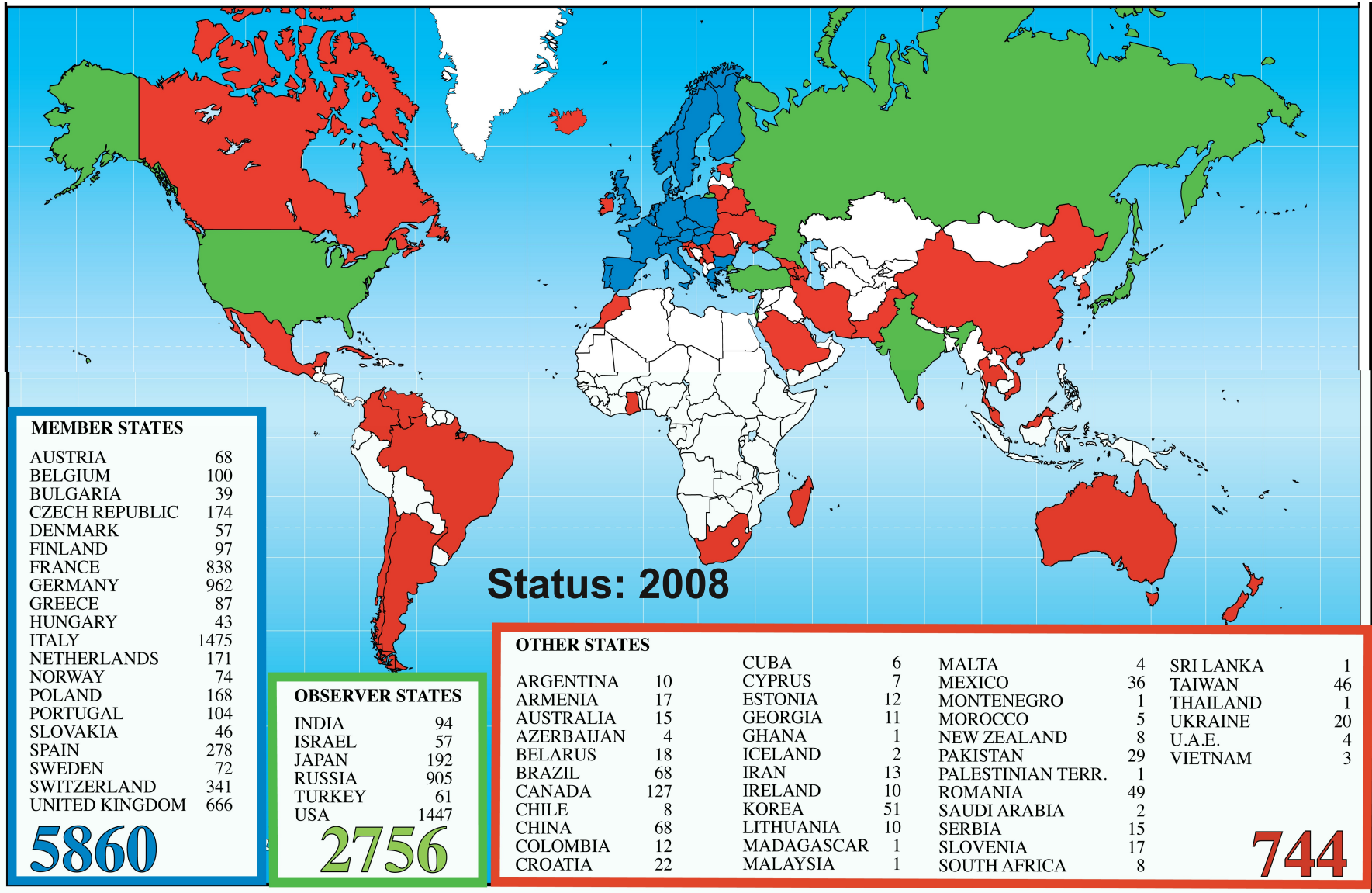
The LHC is addressing fundamental questions of contemporary particle physics and cosmology:

- Origin of mass – the Higgs mechanism
- What is dark matter?
- Matter – antimatter asymmetry?
- Matter in the early universe ?
- hints for Quantum Gravity ??
- Dark Energy ???

Particle Physics is international teamwork

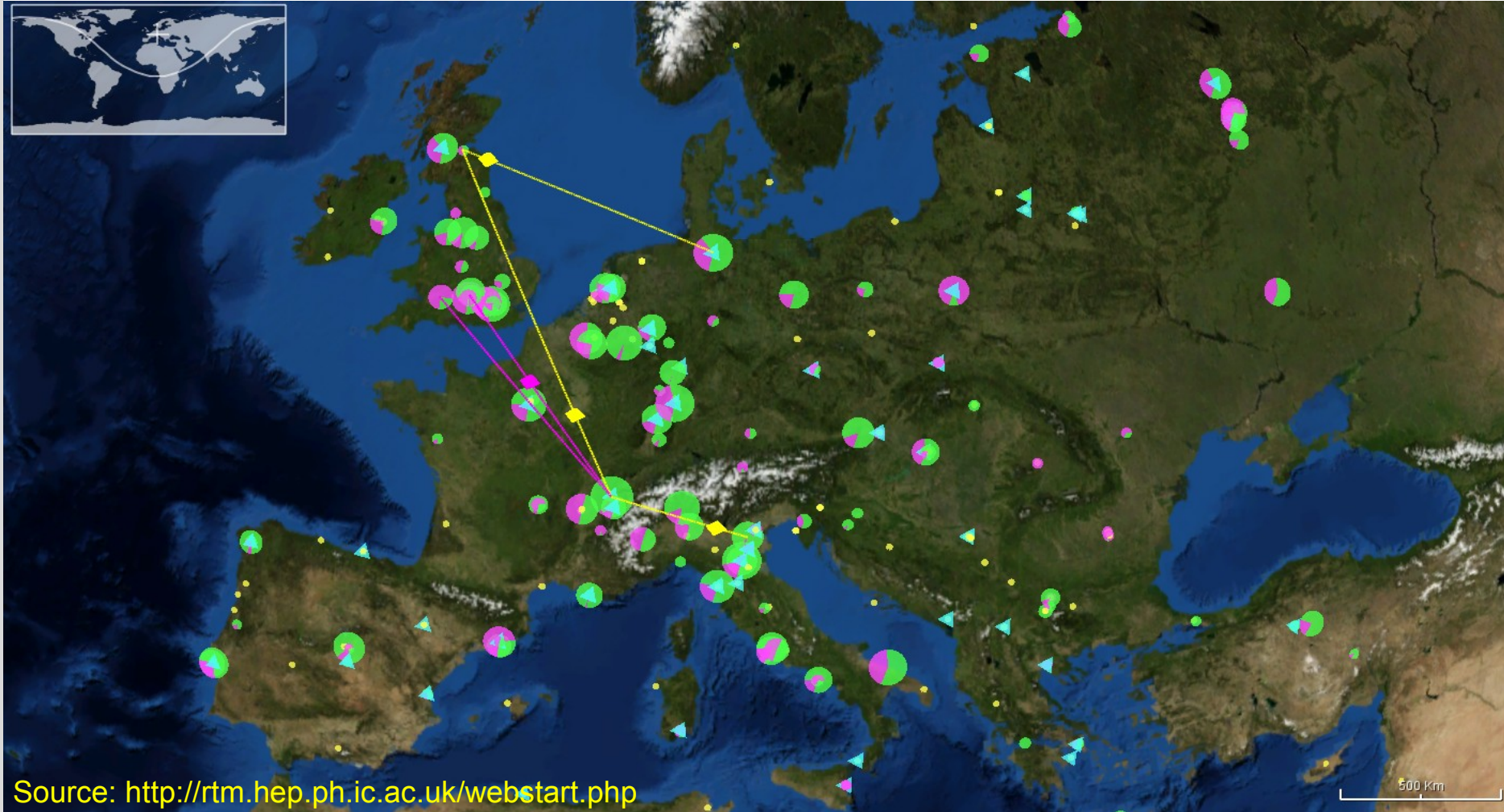


Working @Cern: 290 institutes from Europe, ~ 6349 Users
318 institutes elsewhere, ~ 3766 Users



- Given the international and collaborative nature of HEP
Computing must be distributed
 - harvest intellectual contributions from all partners,
also funding issues
- Early studies in 1999 (Monarc Study group) suggested a hierarchical approach, following the typical data reduction schemes usually adopted in data analysis in high energy physics
- Grid paradigm came at the right time and was adopted by LHC physicists as the base line for distributed computing
- Major contributions by physicists to developments in Grid computing
- Other HEP communities also benefit and contributed

WLHC Computing Grid in action



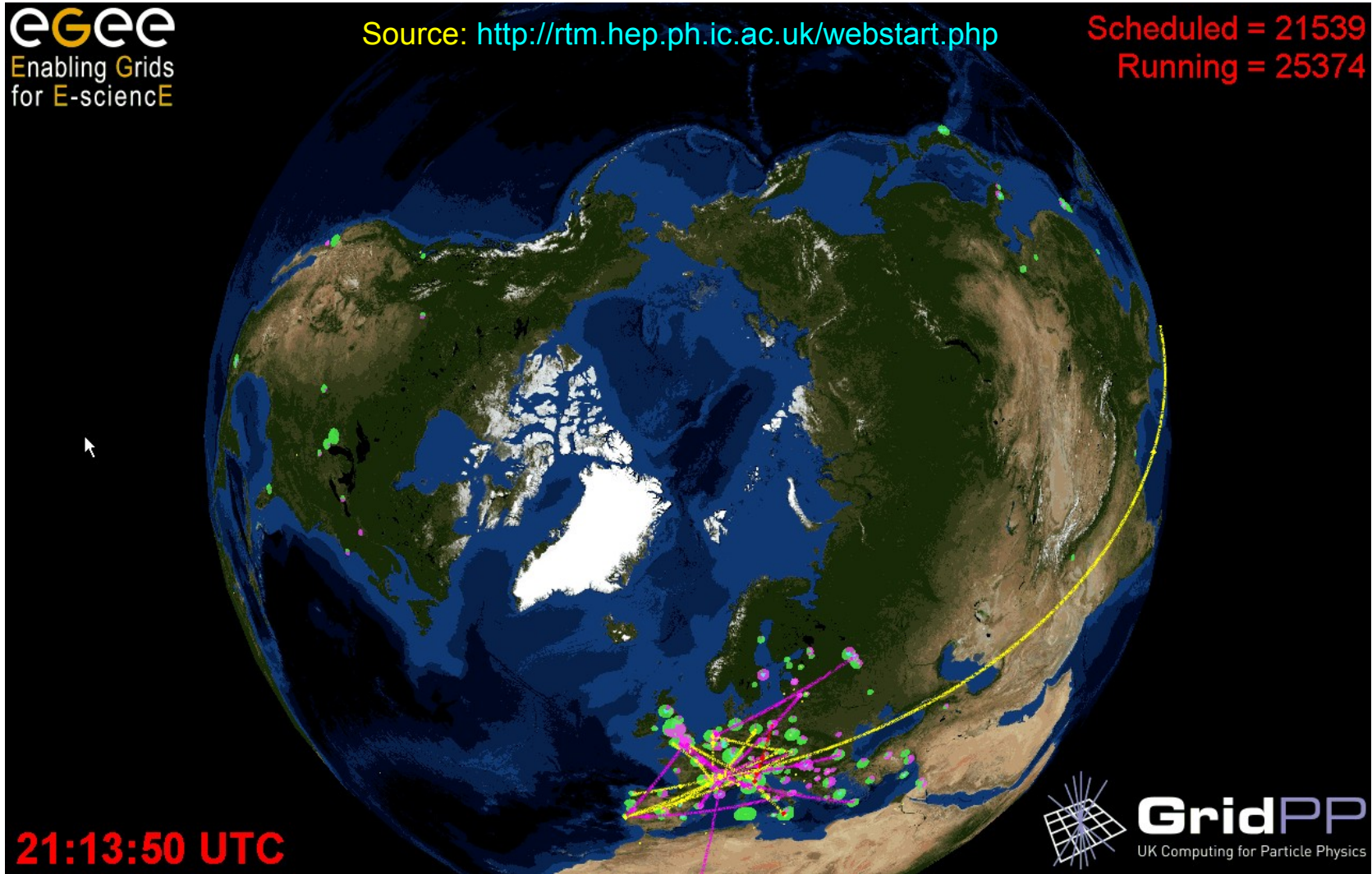
A truly international, world-spanning Grid for LHC data processing, simulation and analysis

WLHC Computing Grid in action

eGEE
Enabling Grids
for E-science

Source: <http://rtm.hep.ph.ic.ac.uk/webstart.php>

Scheduled = 21539
Running = 25374

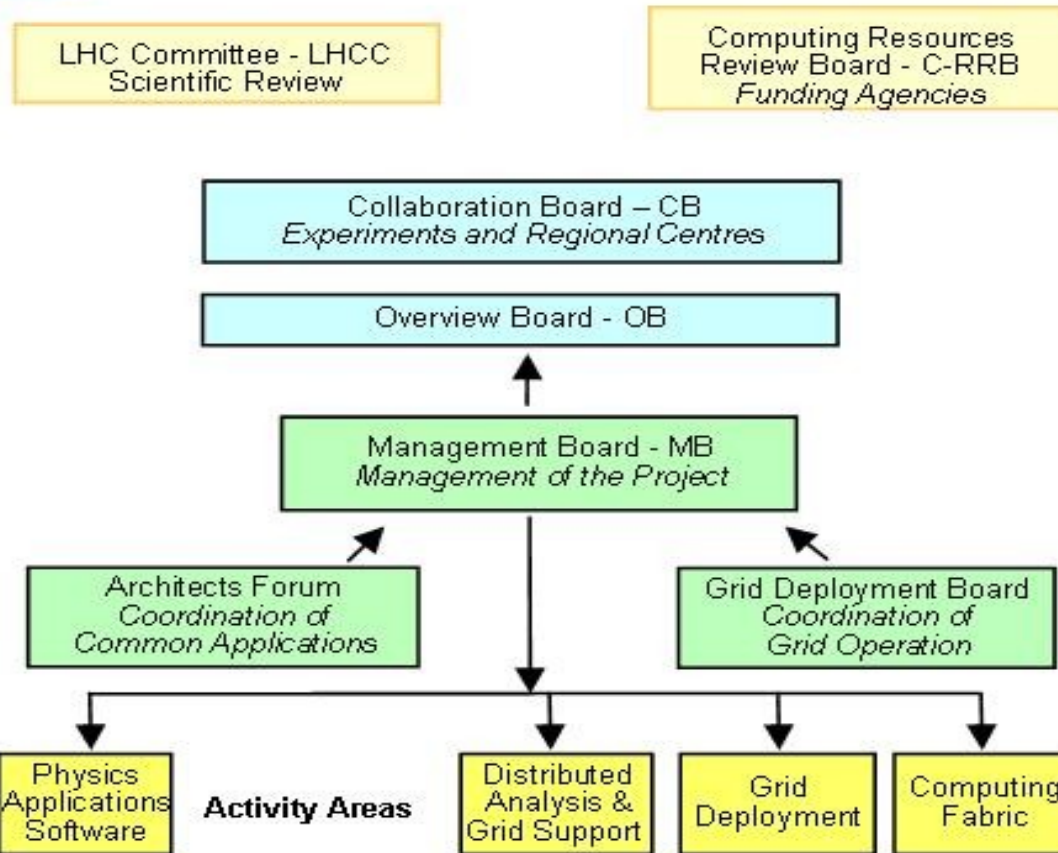


A truly international, world-spanning Grid for LHC data processing, simulation and analysis

Organisation of the World-wide LHC computing Grid



Worldwide LCG Organisation



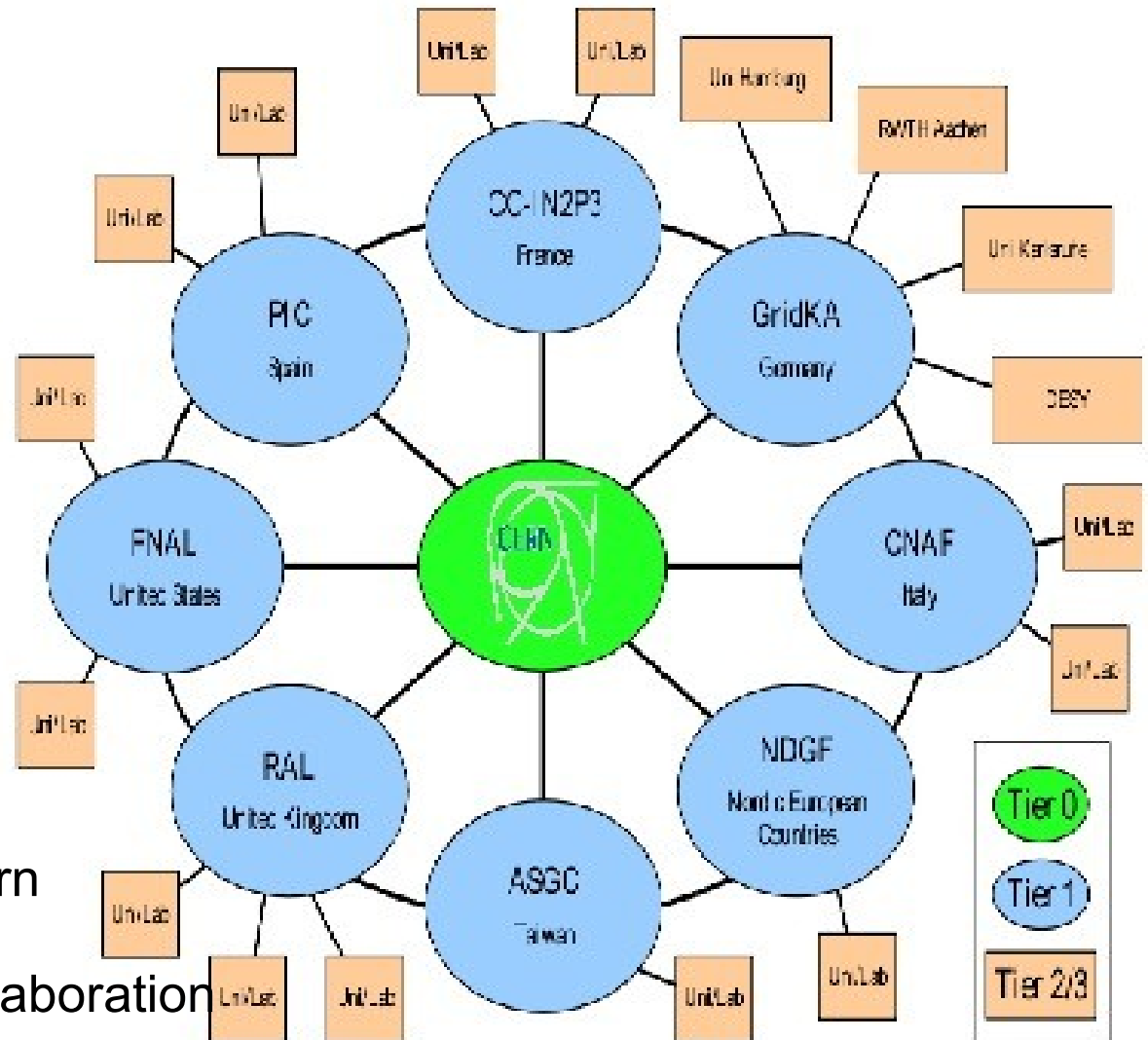
Grids can't work without an organisational structure representing all parties involved

Structure of the LHC Grid

A grid with hierarcies and different tasks at different levels

In addition, it is a **“Grid of Grids”** with interoperability between different middlewares:

- EGEE middleware in most of Europe
- Open Science Grid in USA
- NorduGrid in Northern Europe
- Alien by the Alice collaboration




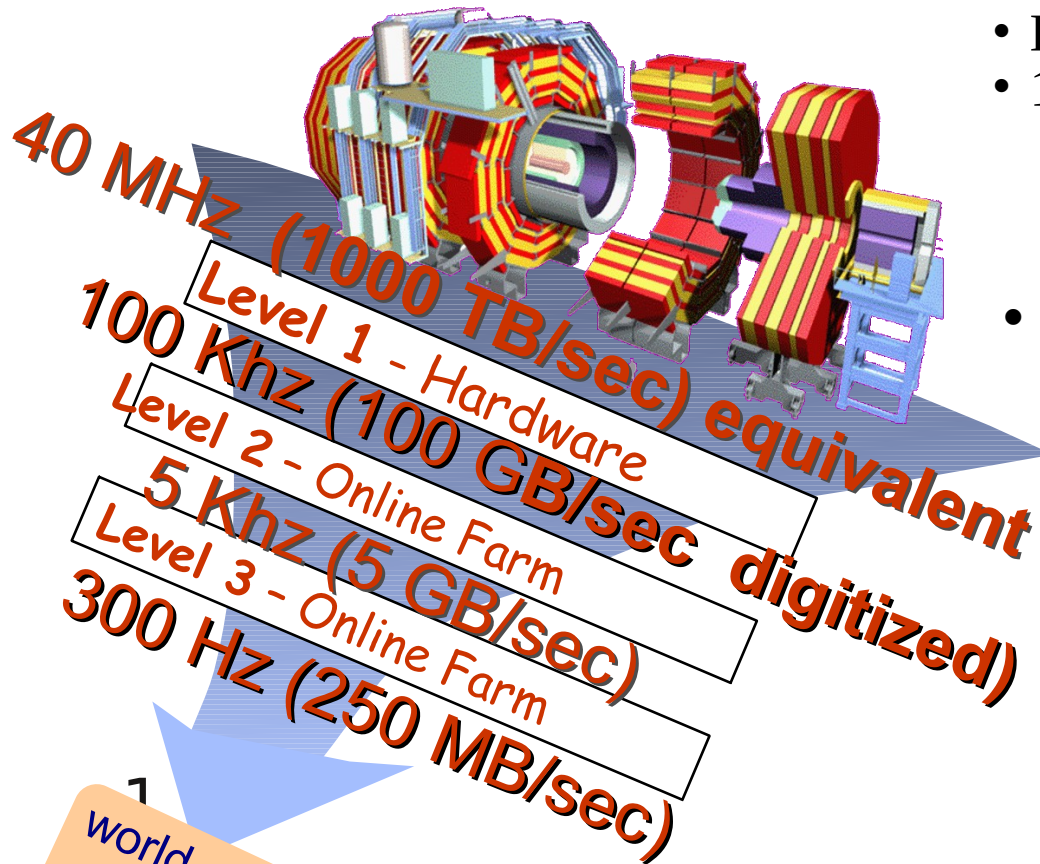
LHC Experiments are huge data sources

- ~ 100 Millionen detector cells
- LHC collision rate: 40 MHz
- 10-12 bit/cell

→ **~1000 Tbyte/s raw data**

- Zero-suppression and Trigger reduce this to „only“ some 100 Mbyte/s

i.e. 1  /sec

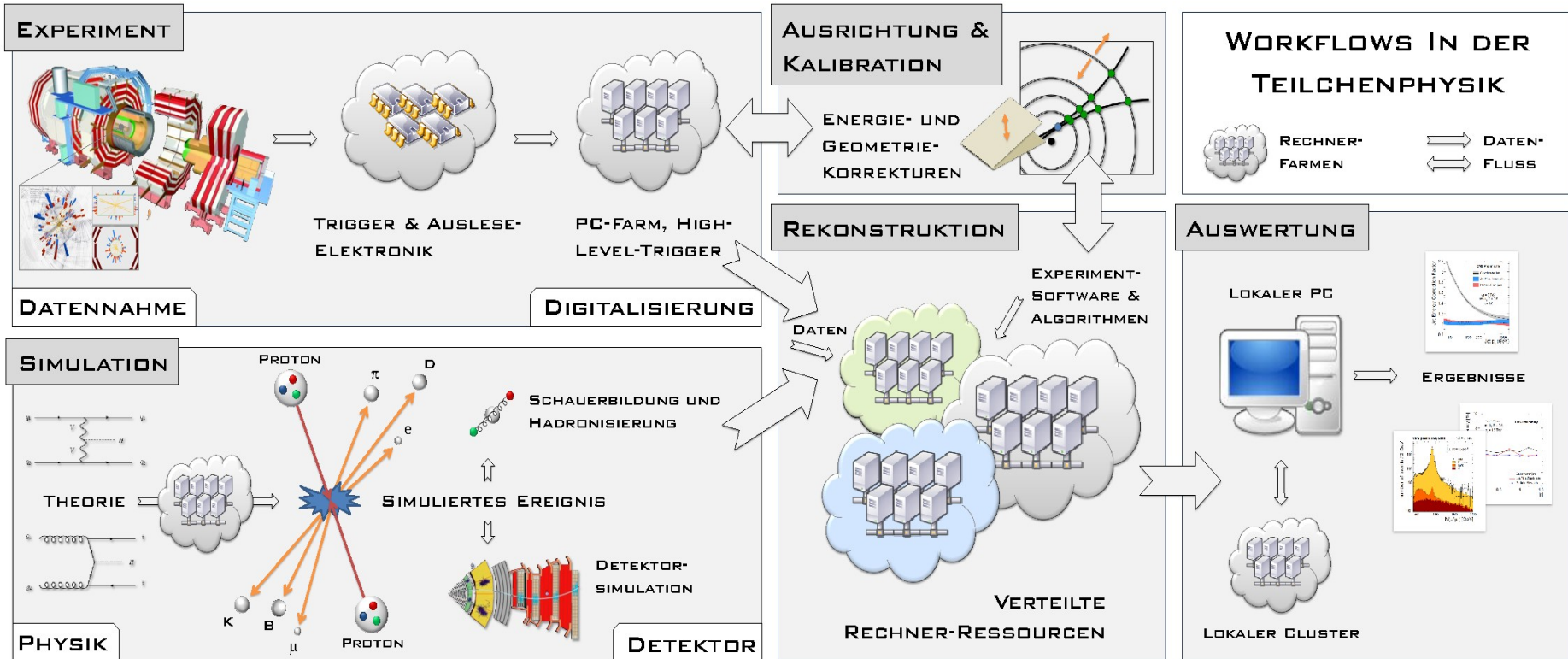


1
world-wide
physics
community

Grid computing
ideal for the job !

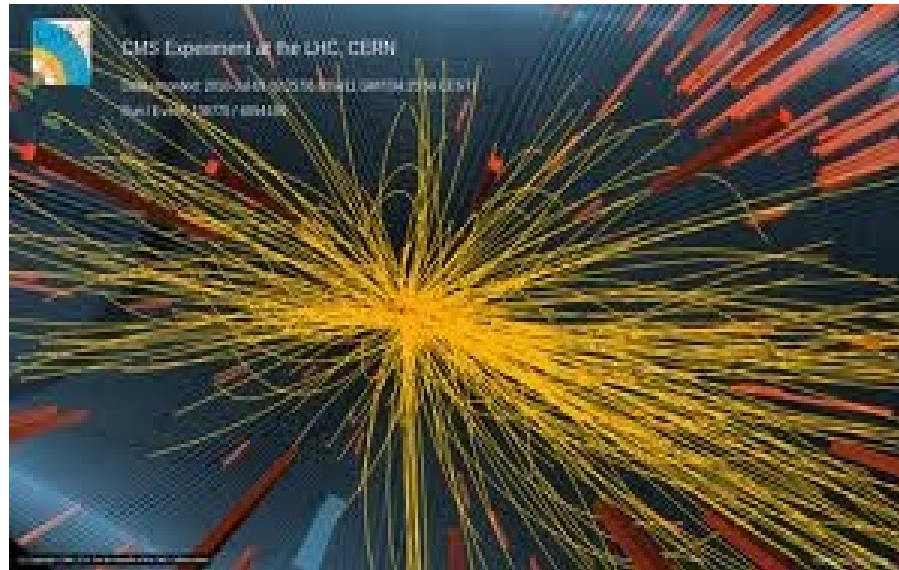


What Particle Physicists do on the Grid



- CPU-intensive simulation of particle physics reactions and detector response
- processing (=reconstruction) of large data volumes
- I/O-intensive filtering and distribution of data
- transfer to local clusters and workstations for final physics interpretation

Why Grid is well suited for HEP



Experimental HEP codes - key characteristics:

- modest memory requirement (~2GB) & modest floating point
 - **perform well on PCs**
- independent events
 - **easy parallelism**
- large data collections (TB → PB)
- shared by very large collaborations

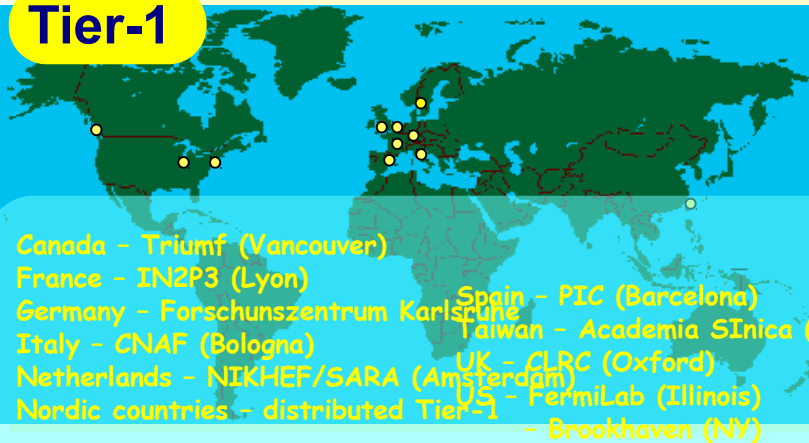


Tier-0 the accelerator centre

- Data acquisition & initial processing
- Long-term data curation
- Distribution of data to T1/T2



Tier-1



11 Tier-1 Centres

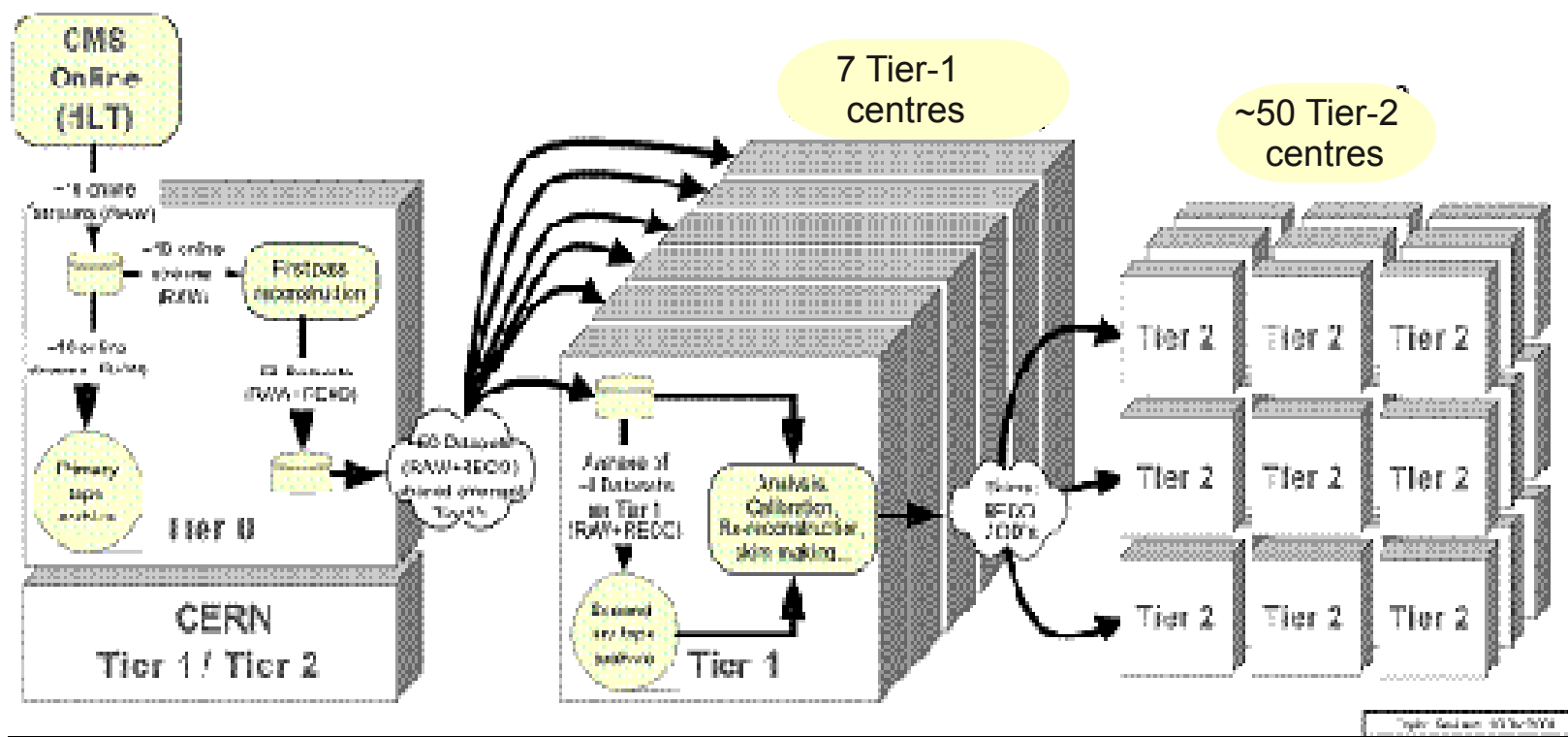
- “online” to the data acquisition process
→ high availability
- Managed Mass Storage
→ grid-enabled data service
- Data-intensive analysis
- National, regional support

Tier-2 150 Centres in 60 Federations in 35 countries

- **End-user (physicist, research group) analysis & collaboration with T3**
(= institute resources) – **where the discoveries are made**
- Monte Carlo Simulation

Tier-3 several 100 grid-enabled PC clusters @ institutes

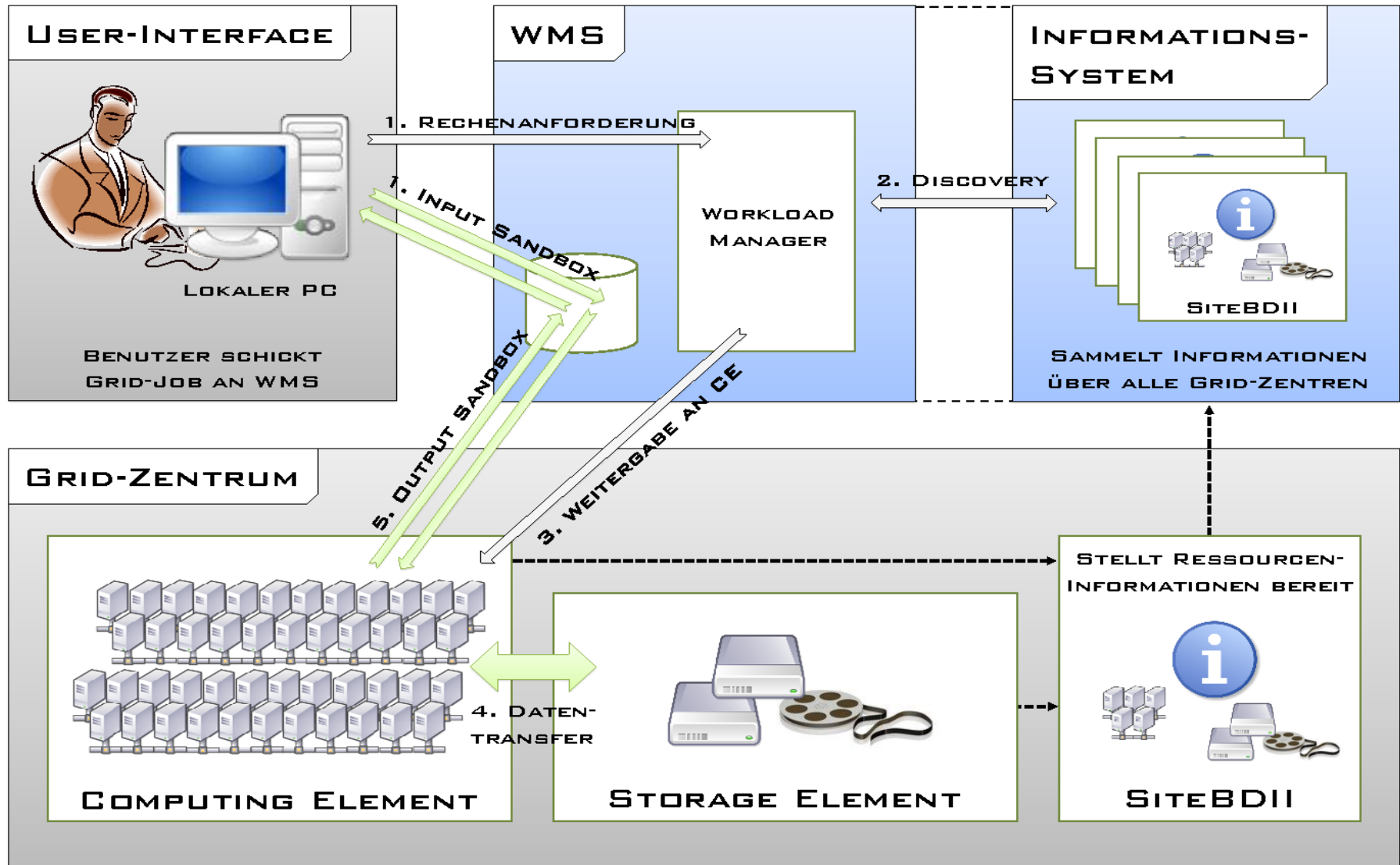
Example: CMS computing model



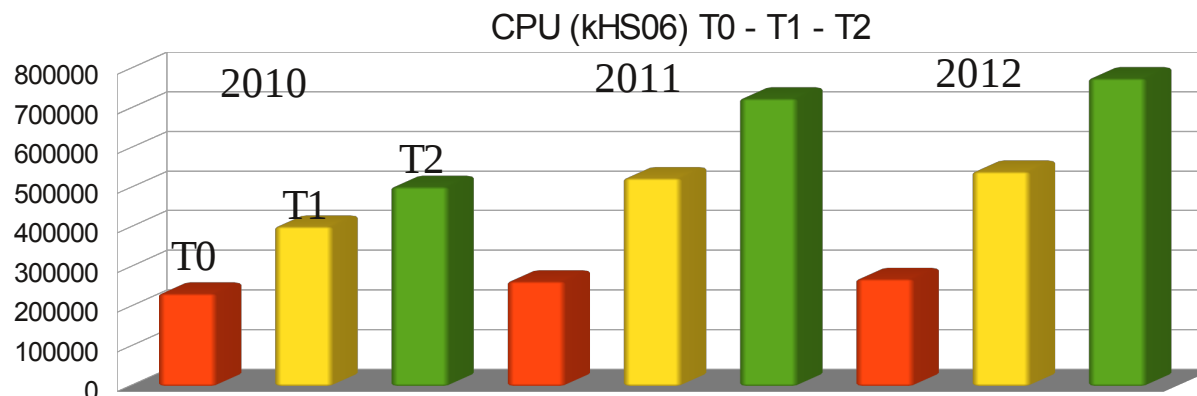
LHC-Experiments

- typically share big Tier-1s, take responsibility for experiment-specific services
- have a large number of Tier2s, usually supporting only one experiment
- have an even larger number of Tier-3s without any obligations towards WLCG

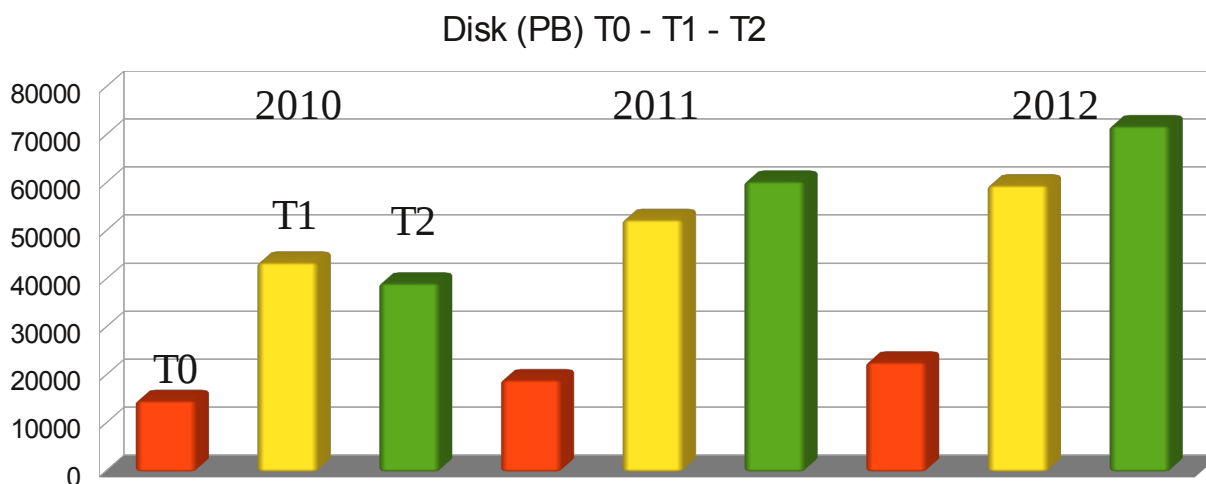
Typical workflows on the Grid



How big is WLCG today ?



Total CPU 2011:
 1500 kHS06
 approx. equiv.
 150'000 CPU cores

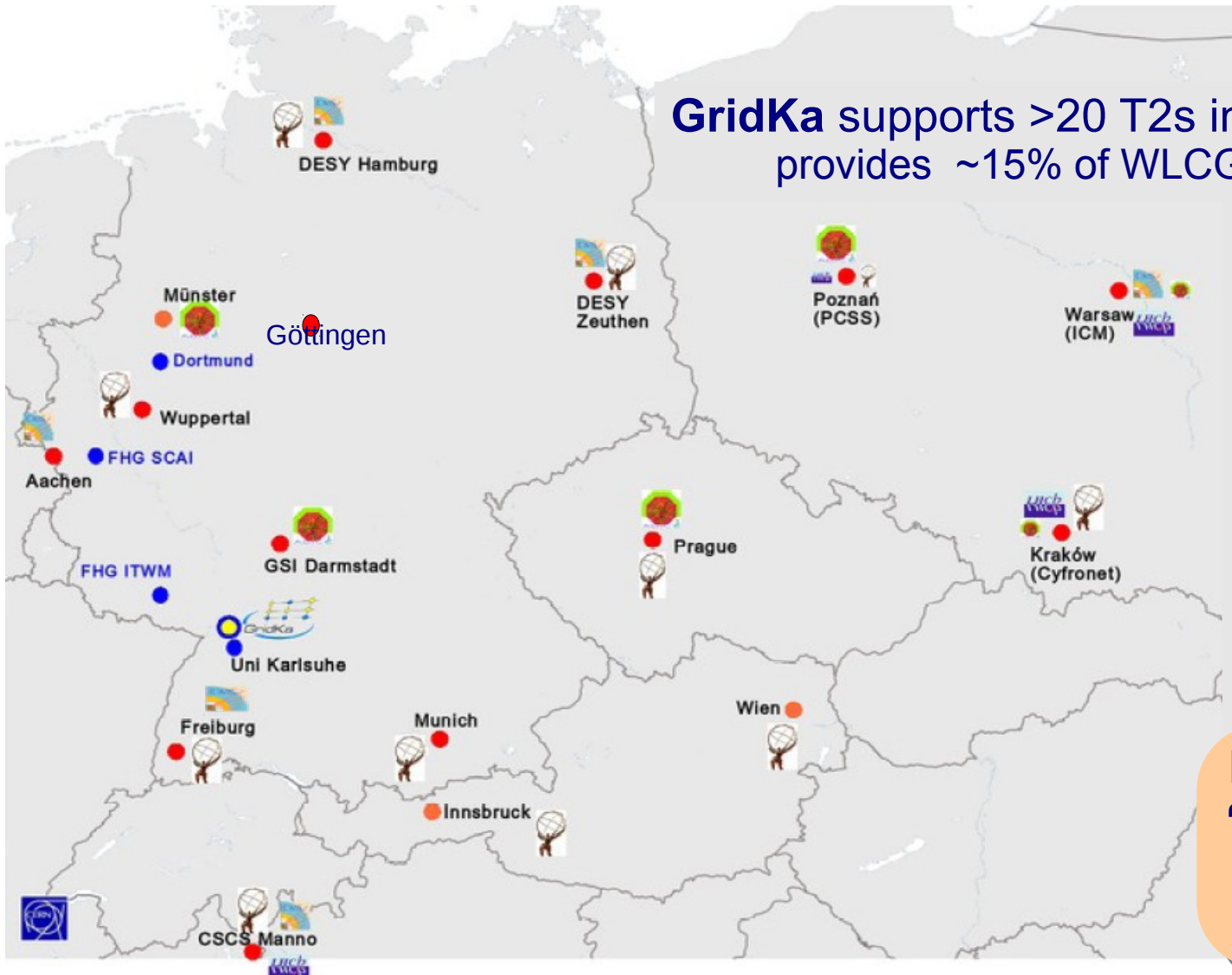


Total Disk 2011:
 130 PB,
 same amount as
Tape Storage 2011:
 130 PB
 (40 PB at CERN,
 none @ T2s)

2012 numbers still being negotiated !

The largest Science Grid in the World

A closer look to the surroundings of GridKa



GridKa supports >20 T2s in 6 countries,
provides ~15% of WLCG T1 resources

Alice T2
sites in
Russia



**Most complex
“T2 cloud” of
any T1 in
WLCG**

**After almost 2 years of experience
with LHC operation:**

How well did it work ?

Almost 2 years of experience - Did it work ?

Up to the users to give feedback:

D. Charlton, ATLAS, EPS HEP 2011

Computing Grid Delivers Physics

Data preparation:

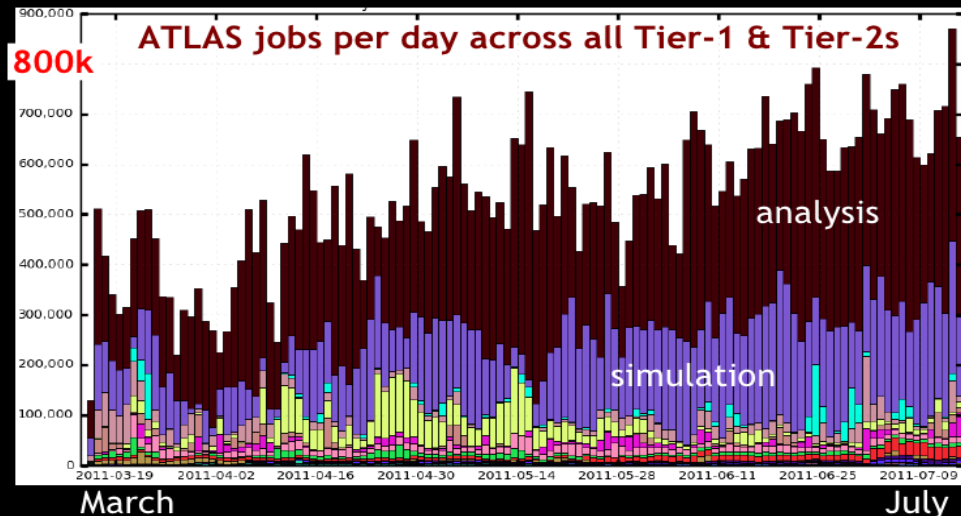
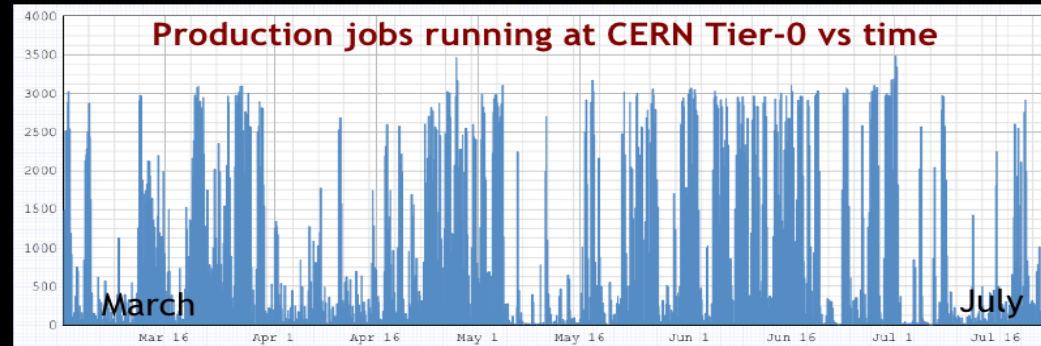
- First-pass reco. at Tier-0 within ~2 days
- Calibration/DQ good for physics analysis
- Data analysable on Grid within ~1 week

Tier-1 and Tier-2's process ~ $\frac{2}{3}$ M jobs per day

- simulation
- re-reconstruction (campaigns)
- group production (ntuples...)
- physics analysis

The high quality computing system allows us to show results on data taken until the end of June

Payback for the years of investment and hard work

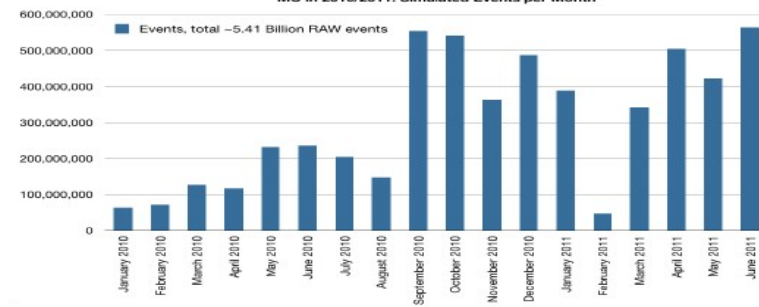
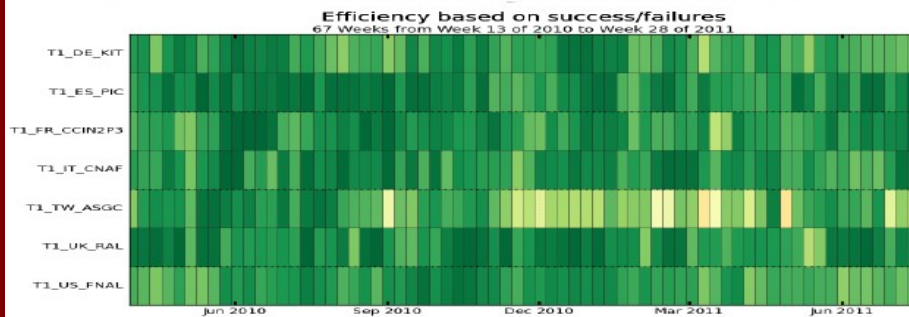
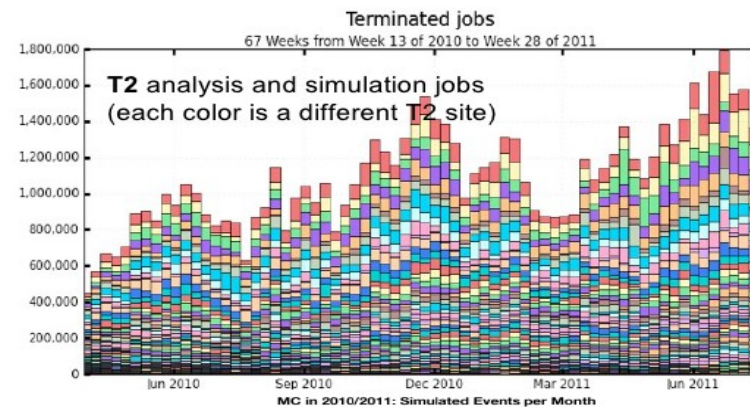
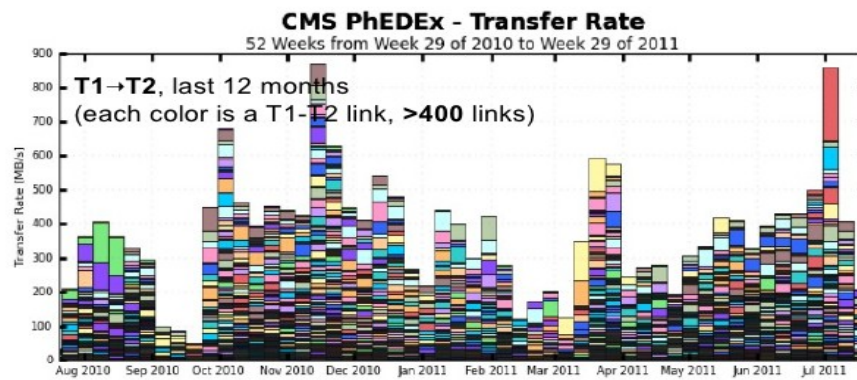


G. Tonelli, CMS, EPS HEP 2011



Offline and Computing running smoothly

- Smooth **Tier-0** operation, keeping up with the data taking. Increase in **Tier-1** utilization, for reprocessing and skimming jobs; High usage of **Tier-2** for analysis, **>400 (800)** individual users per week (month). More than 5.4 Billions MC events.



G. Tonelli, CERN/INFN/UNIPI

HEP_2011_GRENOBLE

July 25 2011

36

Did it work ?

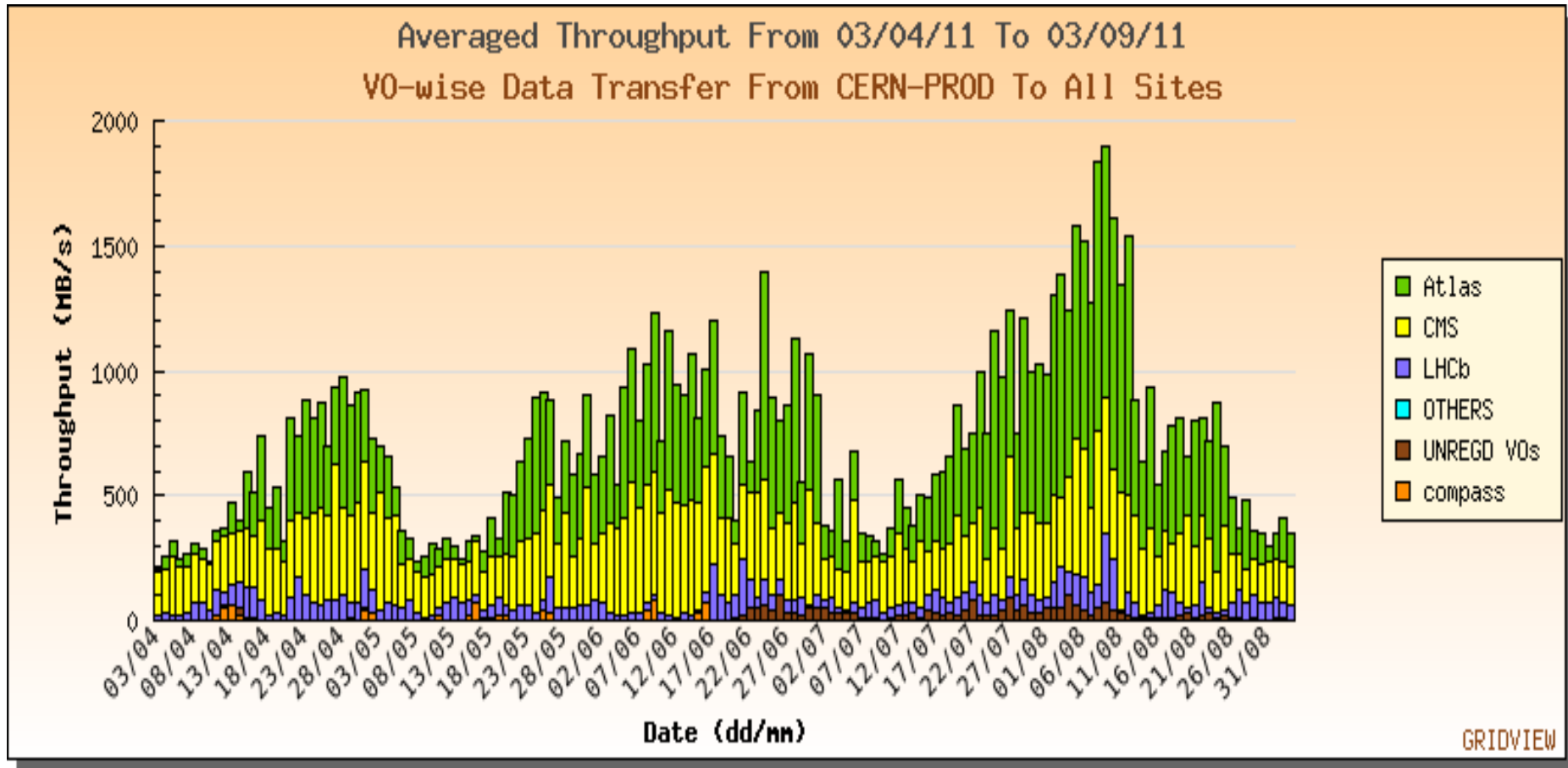
Obviously it did!

- Grid infrastructure for the LHC performed extremely well
- physics results from freshly recorded data
- but: effort for running computing infrastructure is high!

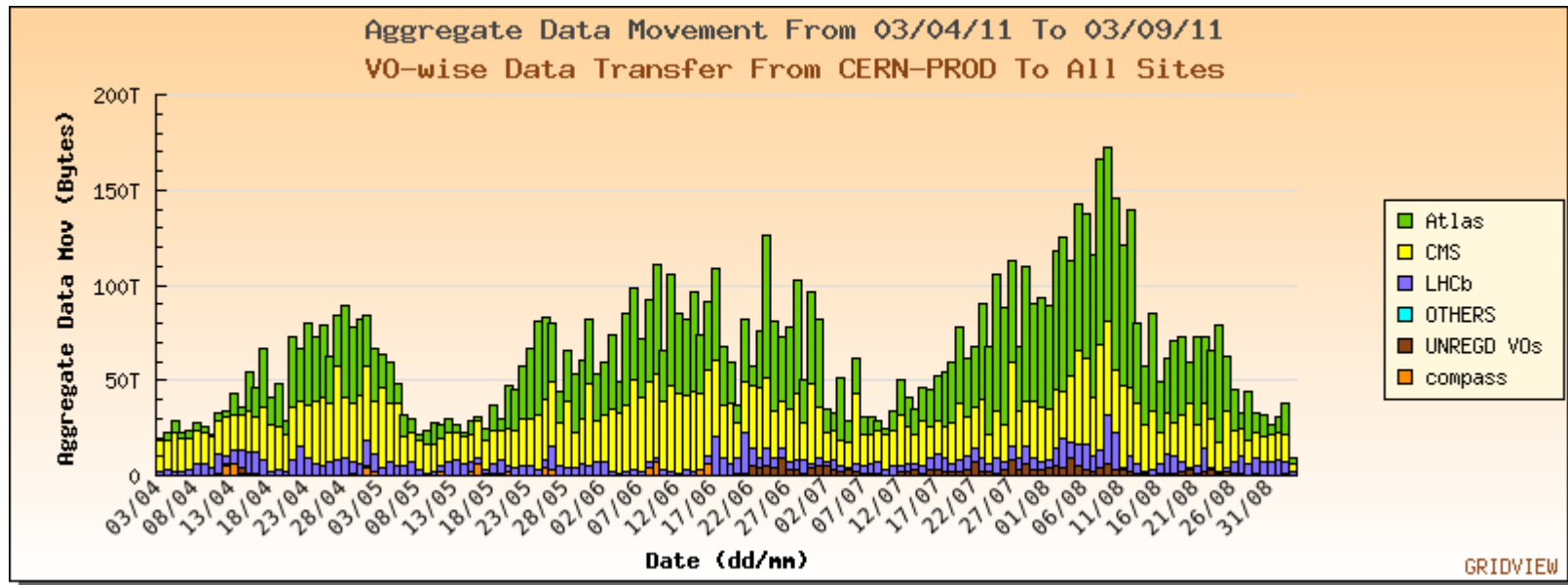
There are challenges ahead!

Let`s have a look in detail ...

Data Export Rate from CERN

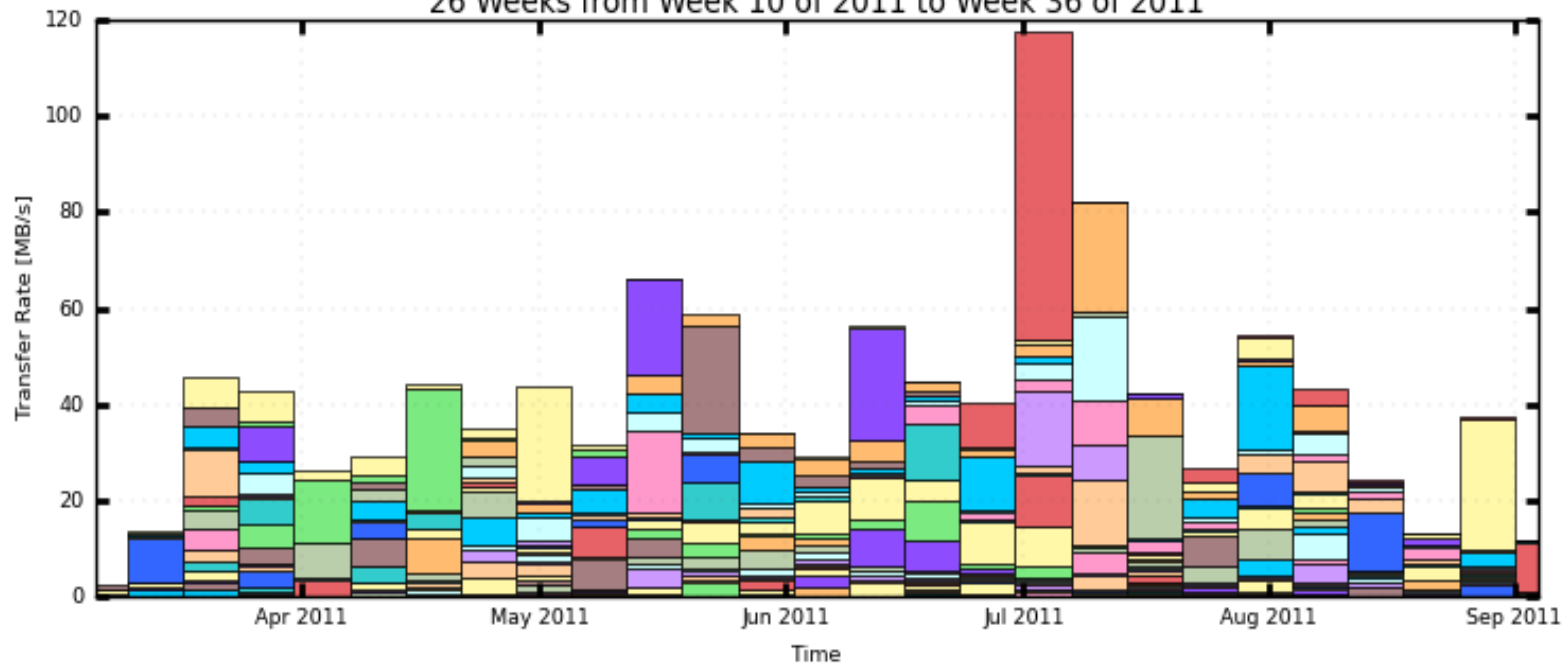


Daily data export from CERN



CMS PhEDEx - Transfer Rate

26 Weeks from Week 10 of 2011 to Week 36 of 2011



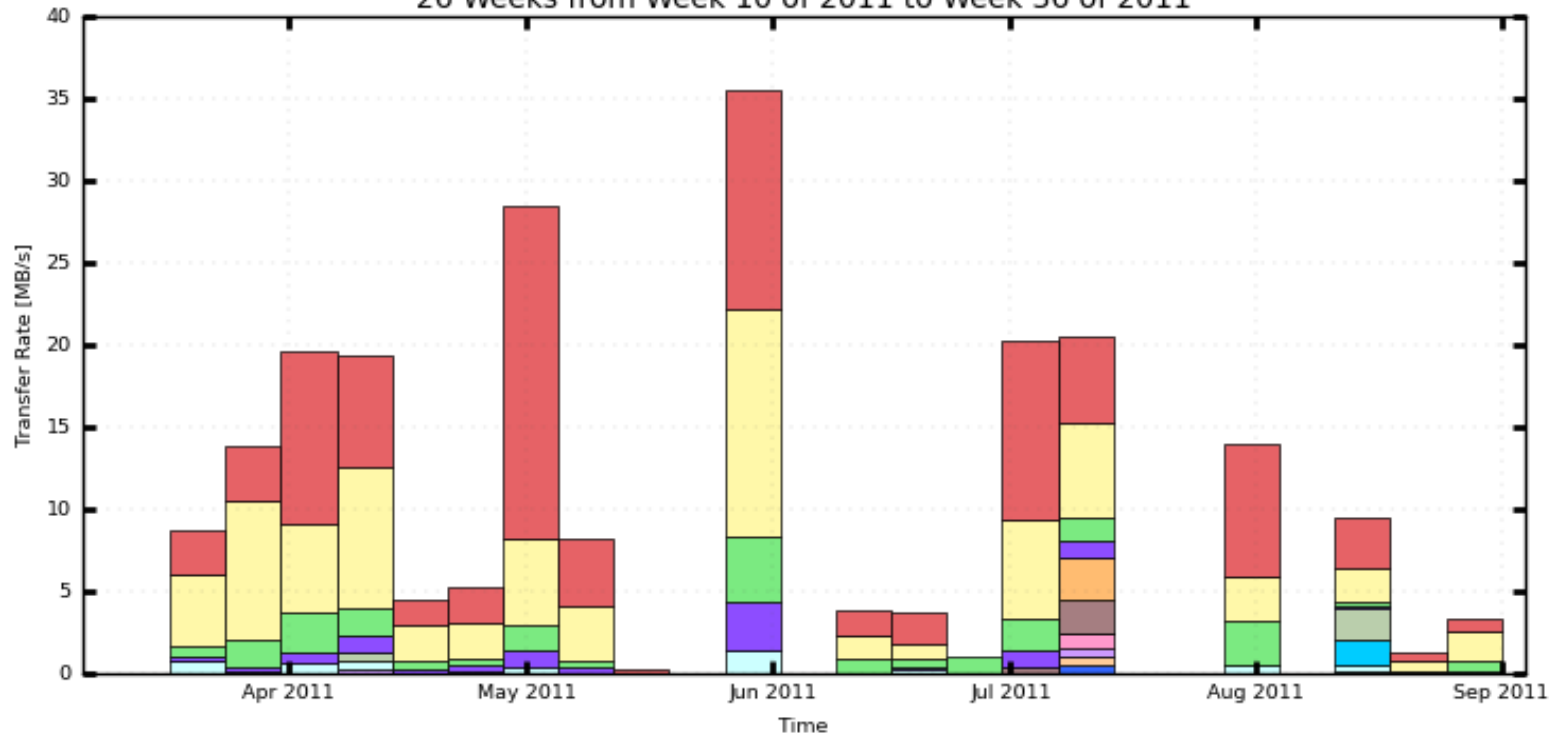
- | | | | | |
|---------------------|----------------|------------------|------------------|---------------------|
| T2_CH_CERN | T2_US_Florida | T2_US_Vanderbilt | T2_IT_Pisa | T2_US_Caltech |
| T2_CH_CAF | T2_HU_Budapest | T2_DE_DESY | T2_US_MIT | T2_UK_London_IC |
| T2_US_UCSD | T2_TW_Taiwan | T2_US_Wisconsin | T2_US_Purdue | T2_IN_TIFR |
| T2_UK_SGrid_RALPP | T2_FR_GRIF_LL | T2_RU_JINR | T2_CH_CSCS | T2_FR_IPHC |
| T2_CN_Beijing | T2_RU_RRC_KI | T2_RU_INR | T2_ES_IFCA | T2_UK_London_Brunel |
| T2_FR_GRIF_IRFU | T2_AT_Vienna | T2_US_Nebraska | T2_BE_UCL | T2_IT_Rome |
| T2_IT_Legnaro | T2_IT_Bari | T2_DE_RWTH | T2_BR_SPRACE | T2_ES_CIEMAT |
| T2_EE_Estonia | T2_RU_SINP | T2_BR_UERJ | T2_BE_IHE | T2_TR_METU |
| T2_RU_IHEP | T2_PL_Warsaw | T2_UA_KIPT | T2_PT_LIP_Lisbon | T2_KR_KNU |
| T2_UK_SGrid_Bristol | T2_FR_CCIN2P3 | T2_PT_NCG_Lisbon | T2_RU_ITEP | |

Maximum: 117.55 MB/s, Minimum: 2.32 MB/s, Average: 40.60 MB/s, Current: 11.49 MB/s

Data Import T2->T1 (MC): example CMS@KIT

CMS PhEDEx - Transfer Rate

26 Weeks from Week 10 of 2011 to Week 36 of 2011

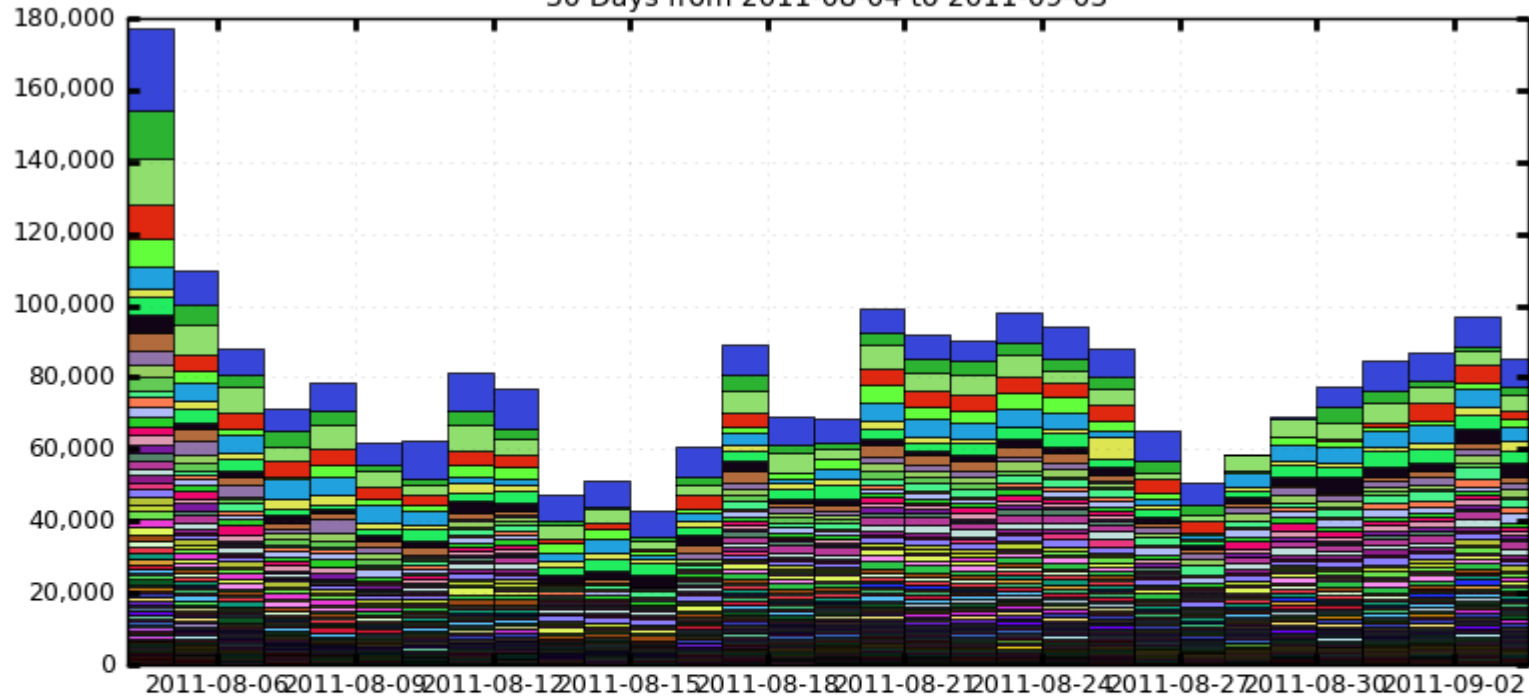


Maximum: 35.51 MB/s, Minimum: 0.00 MB/s, Average: 8.19 MB/s, Current: 0.03 MB/s

Processed Jobs ATLAS

Running jobs

30 Days from 2011-08-04 to 2011-09-03

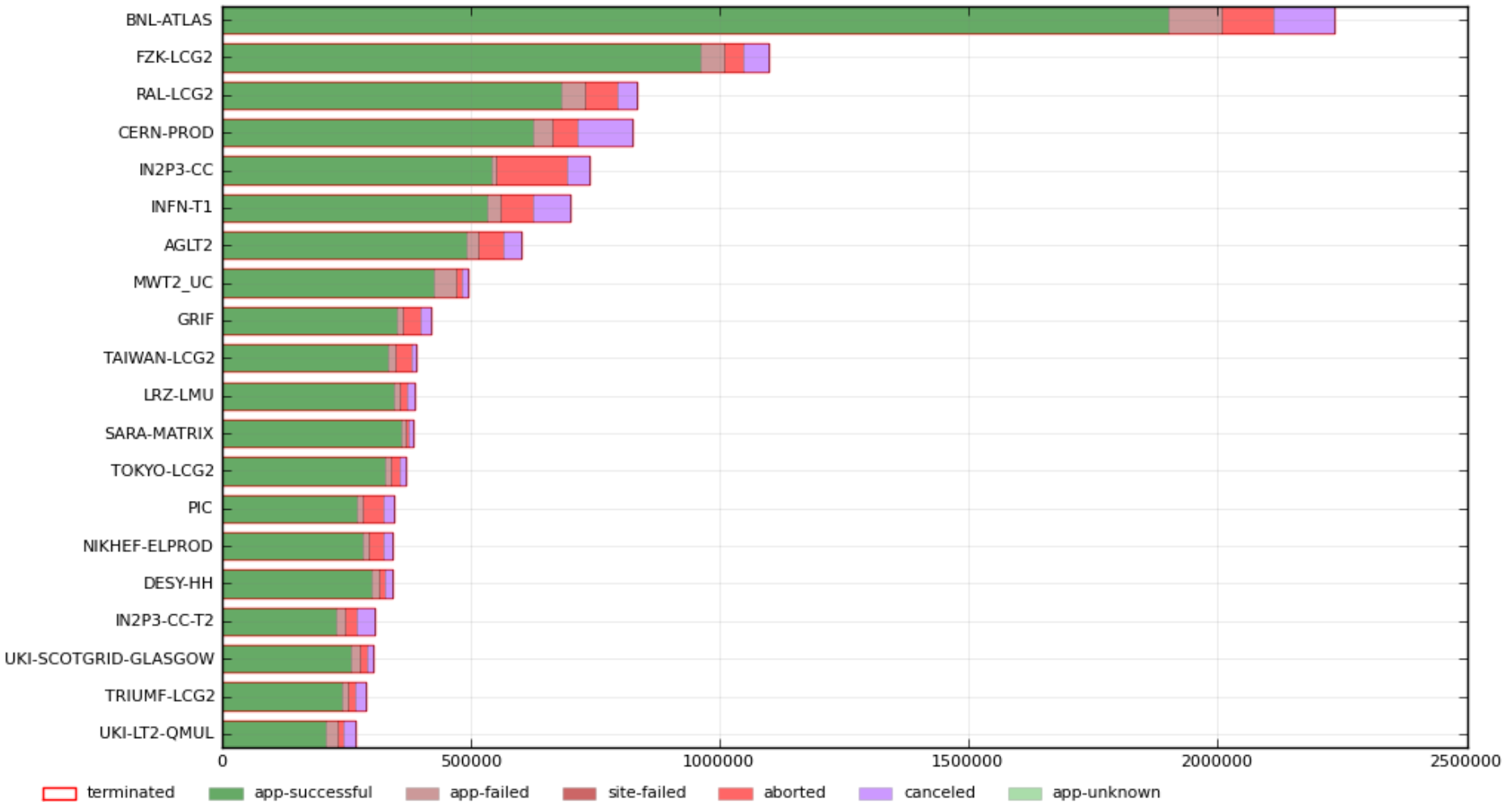


- | | | | | |
|-------------|----------------------|-------------------|---------------------|---------------------|
| BNL-ATLAS | MWT2_UC | FZK-LCG2 | AGLT2 | INFN-T1 |
| IN2P3-CC | TAIWAN-LCG2 | RAL-LCG2 | CERN-PROD | GRIF |
| NDGF-T1 | LRZ-LMU | CYFRONET-LCG2 | UKI-NORTHGRID-MAN-H | IN2P3-CC-T2 |
| DESY-HH | SIGNET | UKI-LT2-QMUL | TOKYO-LCG2 | SE-SNIC-T2 |
| SWT2_CPB | UKI-SCOTGRID-GLASGOW | NIKHEF-ELPROD | INFN-NAPOLI-ATLAS | IFAE |
| SARA-MATRIX | MPPMU | UKI-LT2-RHUL | IFIC-LCG2 | GOEGRID |
| PIC | UNI-DORTMUND | WT2 | TRIUMF-LCG2 | WUPPERTALPROD |
| PRAGUELCG2 | UKI-NORTHGRID-LANCS | HEP-SCINET-T2 | IN2P3-LPC | UTA_SWT2 |
| CSCS-LCG2 | UNI-FREIBURG | RO-07-NIPNE | INFN-MILANO-ATLASC | UKI-SOUTHGRID-RALPP |
| DESY-ZN | IN2P3-LAPP | UKI-SCOTGRID-ECDF | OU_OCHEP_SWT2 | ... plus 54 more |

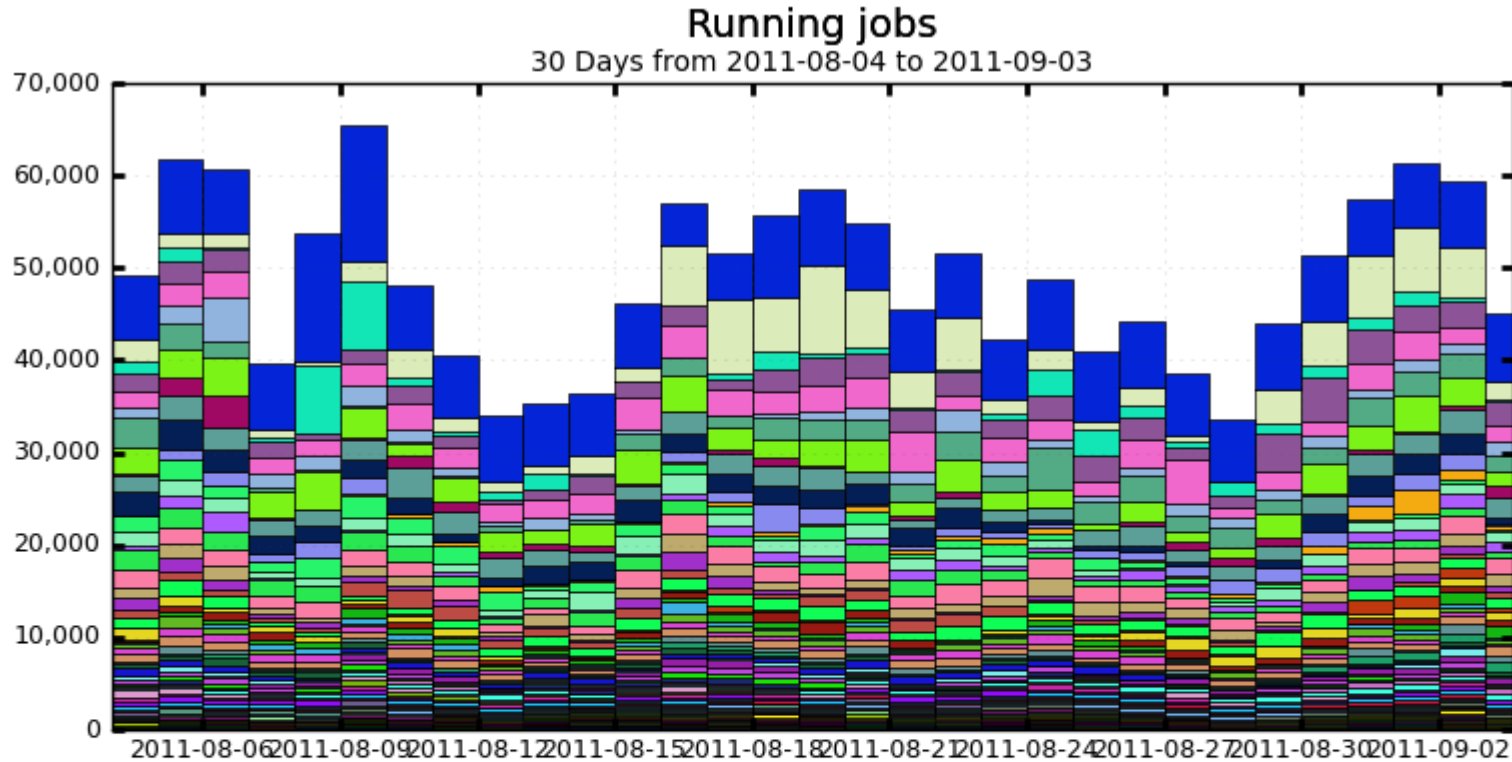
Maximum: 177,292 , Minimum: 43,131 , Average: 79,788 , Current: 85,226

Processed jobs success rates (ATLAS)

Terminated Jobs per site



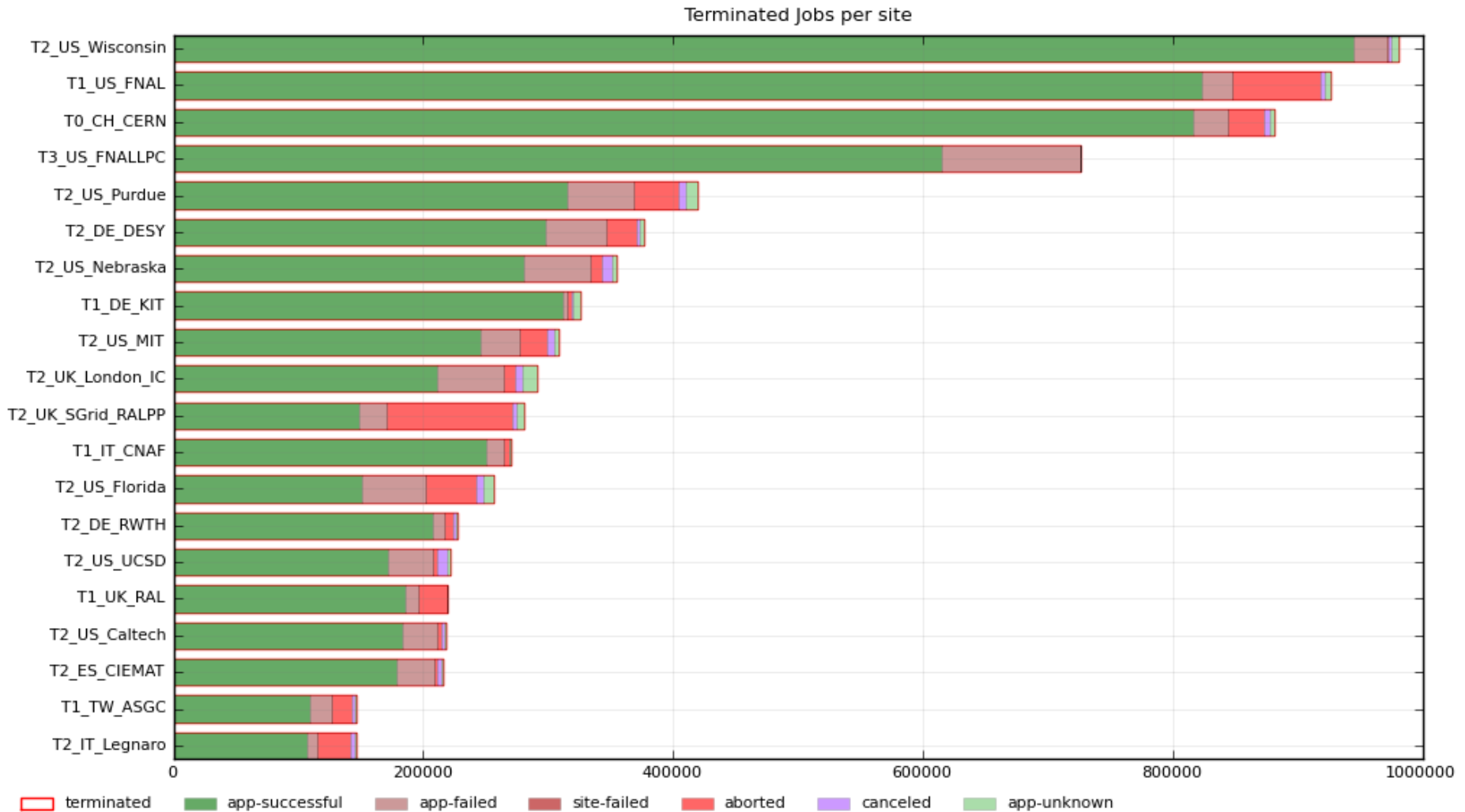
Processed Jobs - CMS



- | | | | | |
|---------------------|-----------------------|---------------------|--------------------|-------------------|
| ■ T1_US_FNAL | ■ T2_US_Purdue | ■ T1_UK_RAL | ■ T2_US_MIT | ■ T3_US_FNALLPC |
| ■ T1_IT_CNAF | ■ T2_US_Florida | ■ T2_US_Wisconsin | ■ T1_FR_CIN2P3 | ■ T2_DE_DESY |
| ■ T2_US_Nebraska | ■ T1_DE_KIT | ■ T2_FR_GRIF_IRFU | ■ T0_CH_CERN | ■ T2_US_Caltech |
| ■ T2_FR_CIN2P3 | ■ T2_DE_RWTH | ■ T2_UK_London_IC | ■ T2_IT_Pisa | ■ T2_ES_IFCA |
| ■ T2_UK_SGrid_RALPP | ■ T2_US_UCSD | ■ T3_UK_London_QMUL | ■ T1_TW_ASGC | ■ T1_ES_PIC |
| ■ T2_FR_GRIF_LLQ | ■ T2_RU_JINR | ■ T2_CH_CSCS | ■ T2_BE_IHHE | ■ T2_ES_CIEMAT |
| ■ T2_PT_NCG_Lisbon | ■ T2_FR_IPHC | ■ T2_BR_SPRACE | ■ T2_IT_Legnaro | ■ T3_US_NotreDame |
| ■ T2_EE_Estonia | ■ T2_UK_London_Brunel | ■ T3_US_Colorado | ■ T3_FR_IPNL | ■ T2_CN_Beijing |
| ■ T2_BE_UCL | ■ T2_FI_HIP | ■ T3_US_TTU | ■ T2_US_Vanderbilt | ■ T2_IT_Bari |
| ■ T2_IT_Rome | ■ T2_KR_KNU | ■ T3_IT_Trieste | ■ T2_TR_METU | ... plus 37 more |

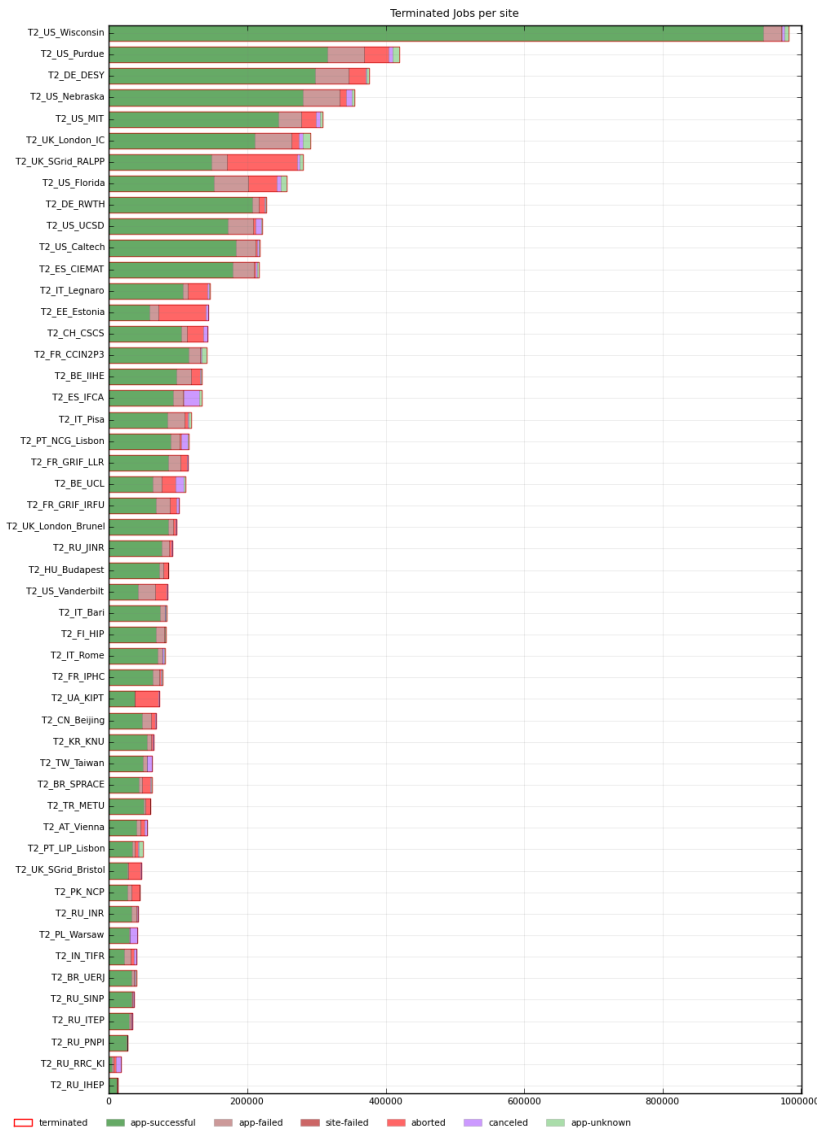
Maximum: 65,452 , Minimum: 33,647 , Average: 48,783 , Current: 45,065

Processed Jobs success rates (CMS)



Job success rates @ T2s

(CMS T2s shown here
ATLAS looks similar)



T2s see a mixture of
- MC production jobs
- user analysis and skimming jobs

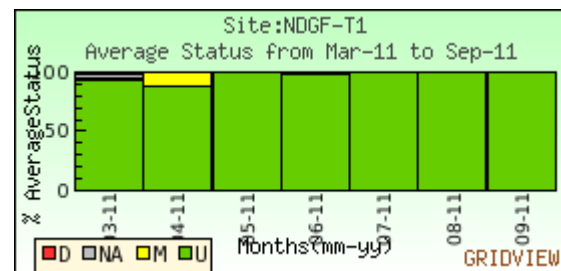
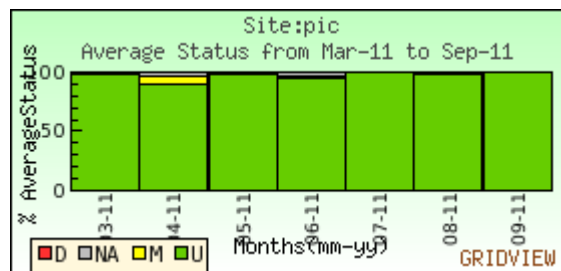
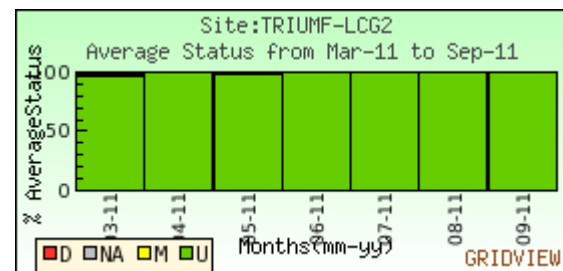
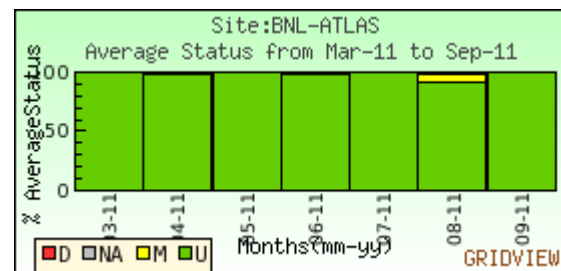
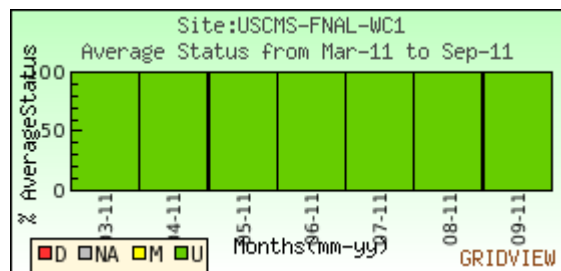
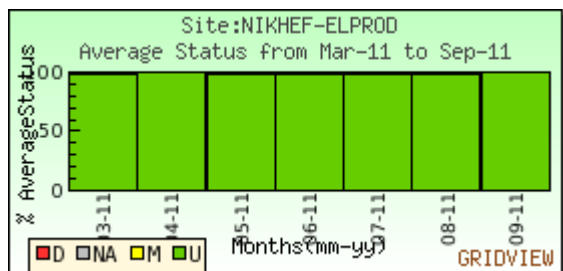
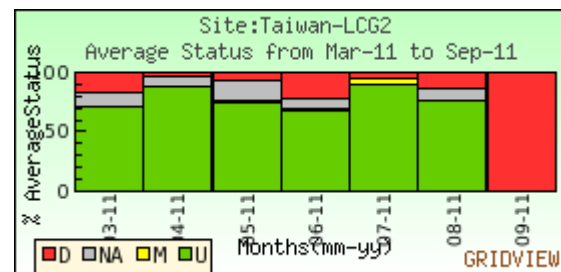
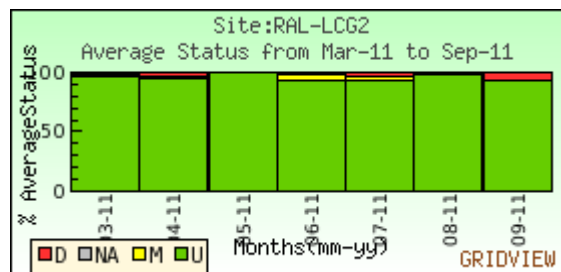
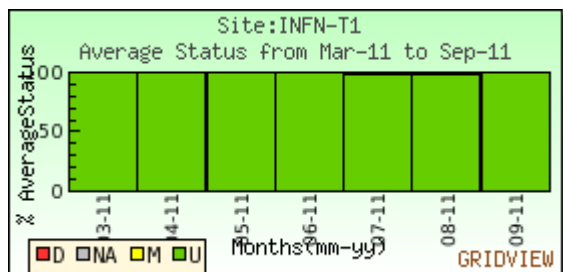
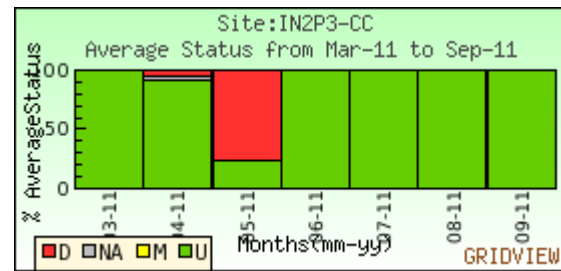
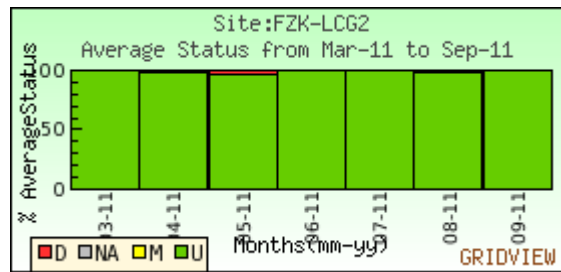
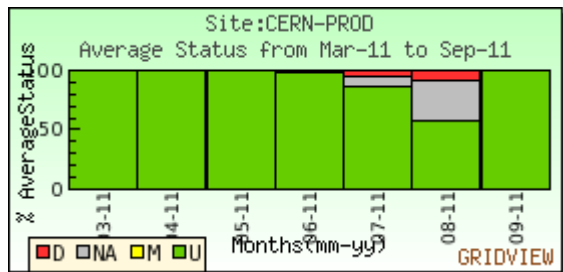
Job success rates not in all cases
above 90%

90% success rate would be considered
very low for a classical computer centre

This must improve ...

Not easy to disentangle failures of
the system from “user errors”

Site Reliability (examples)



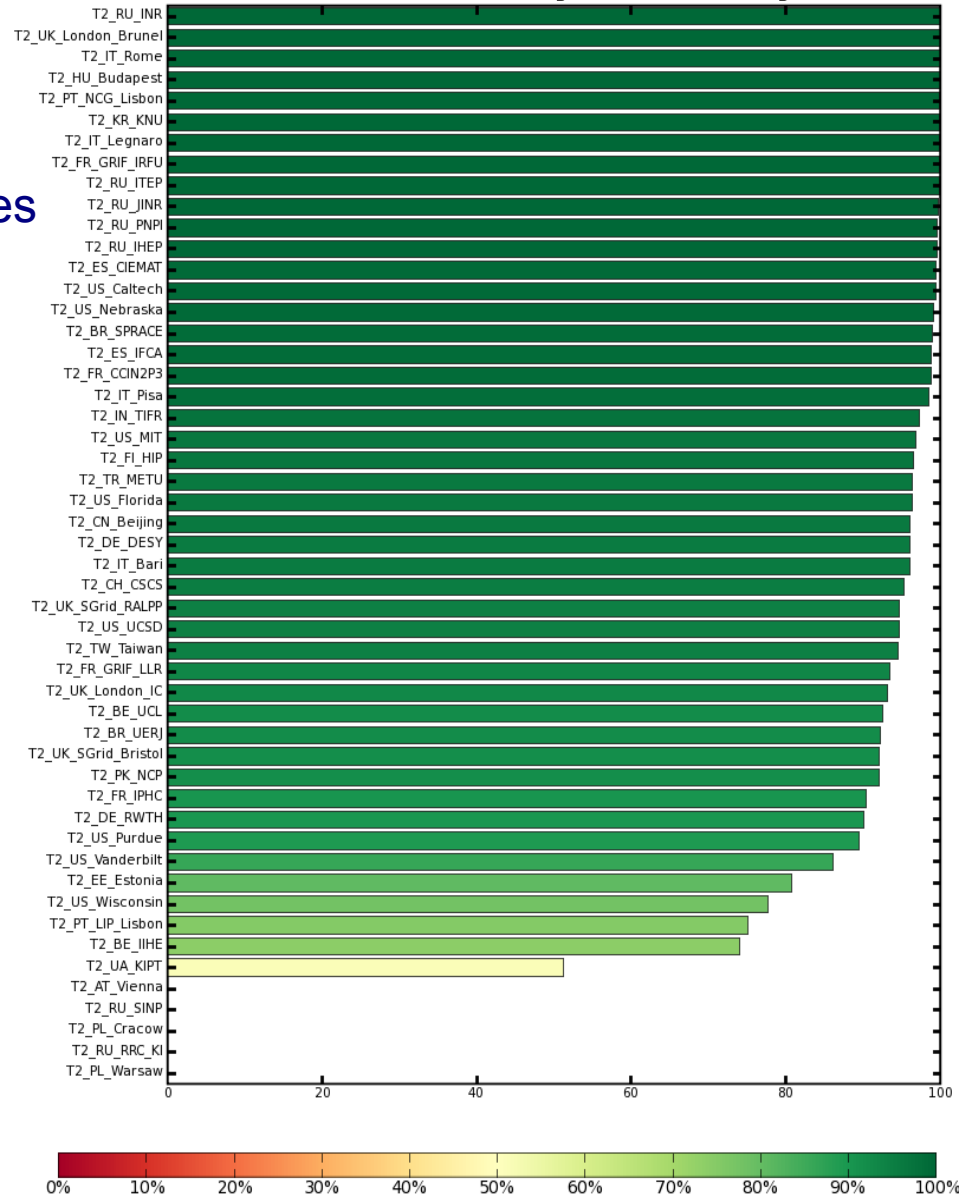
T2 performance

Site Availability, last 31 days

“Availability” of CMS T2 sites

There are sites with performance issues !

Typically, less well performing sites are very small !



T2 Performance – ATLAS example

Site Availability using WLCG_NAGIOS

30 Days from 2011-08-04 to 2011-09-03



T2 Performance Atlas (cont'd)

This is an example of a time-resolved measurement of site availability

Message similar as previously.

This kind of graphs helps the site responsables to monitor their sites and act on problems.

Provided centrally by WLCG and experiments !



Again: Does it work ?

YES !

- Routinely running ~150'000 jobs simultaneously on the Grid
- Shipping over 100 TB/day to T1 centres
- data distribution to T2 works well
- some T2s have performance issues
- very little is known about T3 usage and success rates - responsibility of the institutes
- plenty of resources available at LHC start-up, now approaching “resource limited operation”
- Users have adapted to the “GridWorld” -
Grid is routinely used as a huge batch system, output is transferred home

but ...

Does it work?

Message:

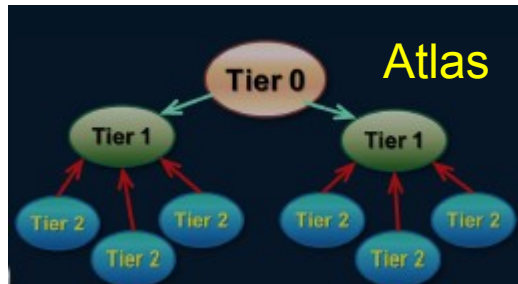
**it worked better than expected by many,
but running such a complex computing infrastructure
as the WLCG is tedious (and expensive!)**

Reliability and cost of operation can be improved by

- simplified and more robust middleware
- redundancy of services and sites,
requires dynamic placement of data and investment in network bandwidth
- automated monitoring and triggering of actions
- use of commercially supported approaches to distributed computing:
 - private clouds are particularly important for shared resources at universities
 - eventually off-load simple tasks (simulation, statistics calculations) to commercial clouds

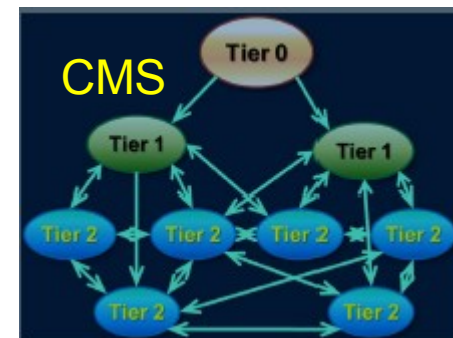
Many of the new developments were addressed at this School

Let`s have a look at some future developments ...



ATLAS and CMS
computing models differ slightly

CMS already more “distributed”

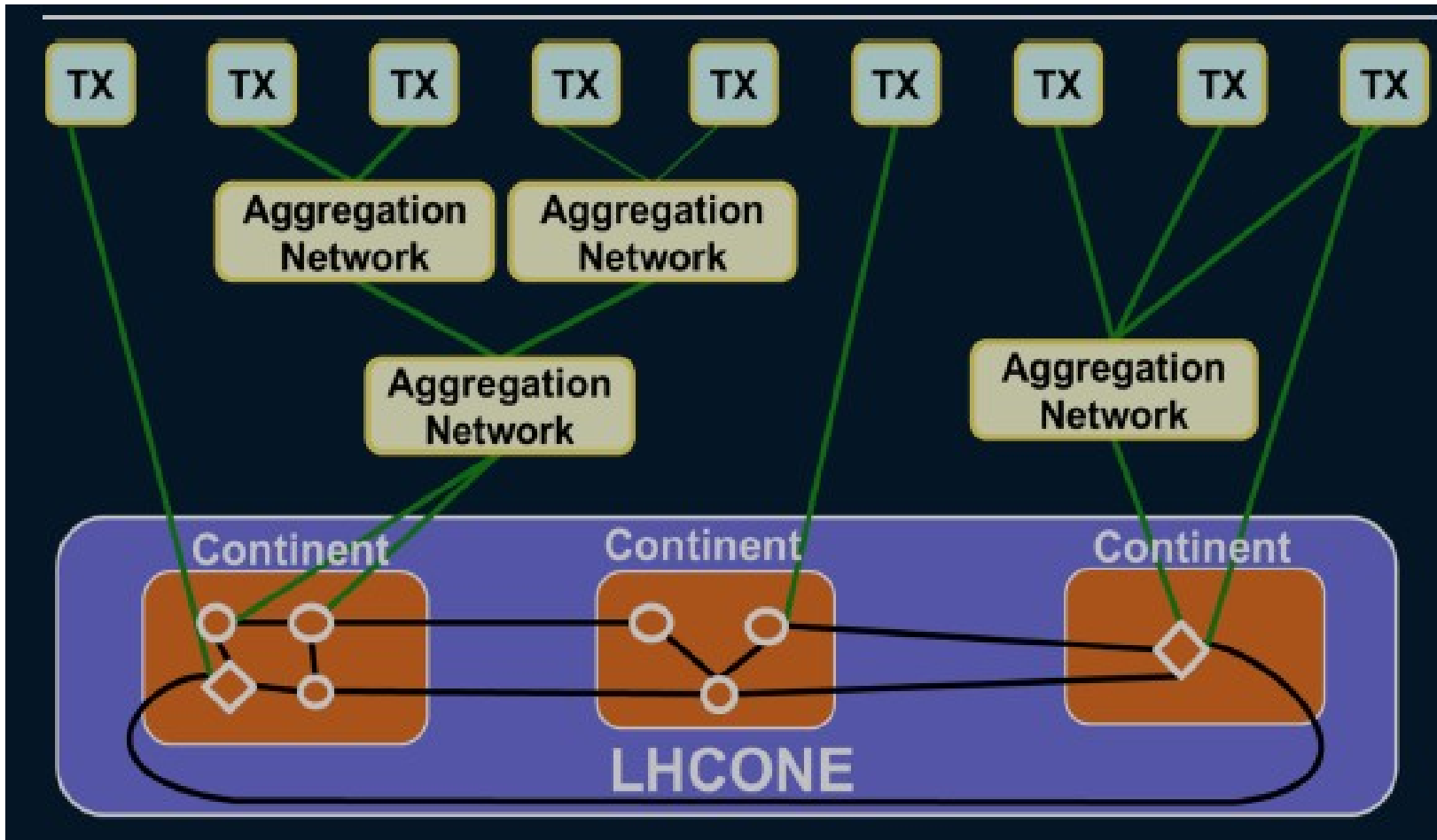


Aim of LHCONE project is

better trans-regional networking for data analysis,
complementary to **LHCOPN** network connecting LHC T1s

- **flat(er) hierarchy:** any site has access to any other site's data
- **dynamic data caching:** pull data “on demand”
- **remote data access:** jobs may use data remotely

by interconnecting open exchange points between regional networks



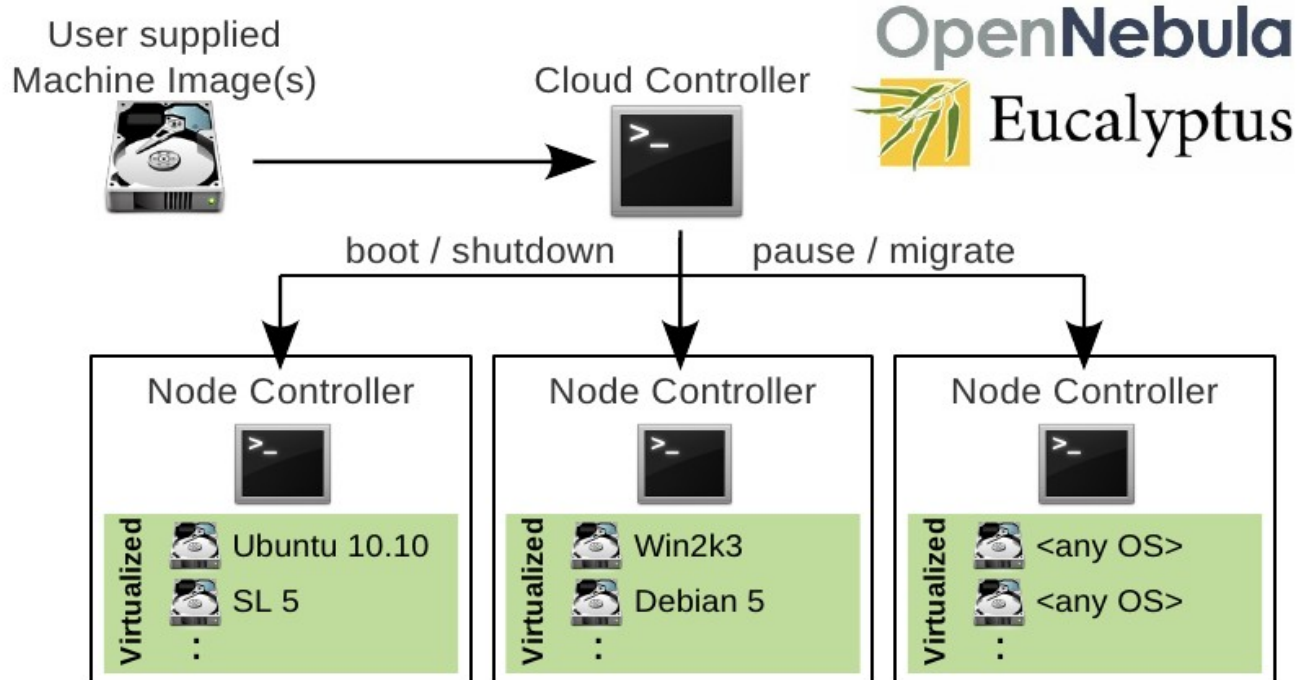
Schematic layout of LHCONE network infrastructure

A dedicated HEP network infrastructure – what is the cost ?

Virtualisation

You have heard a lot about clouds and virtualisation at this school in a nutshell:

- Clouds offer “Infrastructure as a Service”
- easy provision of resources “on demand”
even by including (private) cloud resources as a classical batch queue
(e.g. ROCED project developed at EKP, KIT)
- independent of local hardware and operating system
(Scientific Linux 5 for Grid middleware and experiment software)

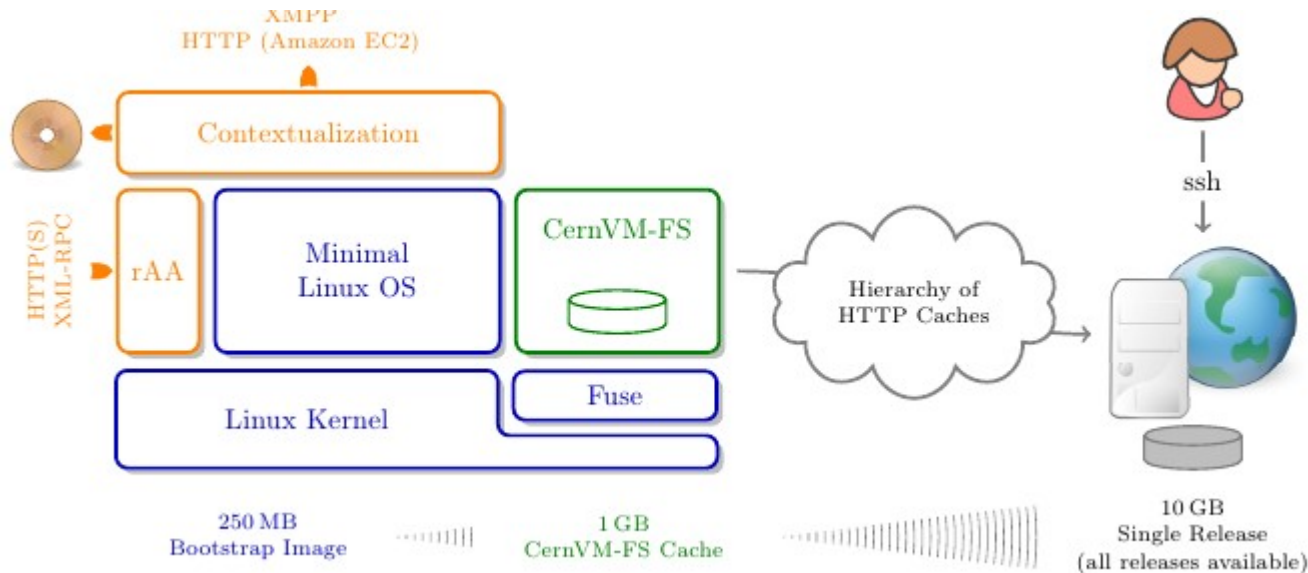


CernVM & CernVM-FS

CernVM is a virtual machine (“Virtual Software Appliance”) based on Scientific Linux with CERN software environment, runs on    



CernVM-FS is a client-server file system based on http and implemented as a user-space file system optimized for read-only access to software repositories with a performant caching mechanism. Allows a CernVM instance to efficiently access software installed remotely.



A recap of this School

This week, you have heard about many of the new developments:

● J. Templon, *Grid and Cloud*

“Grids need Clouds to prosper, Clouds need Grids to scale”

● O. Synge, *Virtualisation*

● P. Millar, *Data Storage*

● C. Witzig, *European Grid Projects*

● T. Beckers, *Storage Architectures for Petaflops Computing*

● S. Reißer, *Grid User Support*

● N. Abdennadher, *Combining Grid, Cloud and Volunteer Computing*

● U. Schwickerath, *Cloud Computing*

● S. Maffioletti, *ARC for developers*

● T. Metsch, B. Schott, *Sustainable DCI Operations*

● A. Aeschlimann, *Grid and Cloud Security*

and a number of Hands-on workshops in parallel sessions.

**HEP Grid runs fine in its initial version,
but virtualisation and “clouds” offer new possibilities
for resource increase, efficiency, cost effectiveness , operation and reliability.**

It's up to you, the participants of this school, to shape the future !

Thanks to all speakers, session teams and organizers and also from my side.