

From raw data to Pbytes on disk

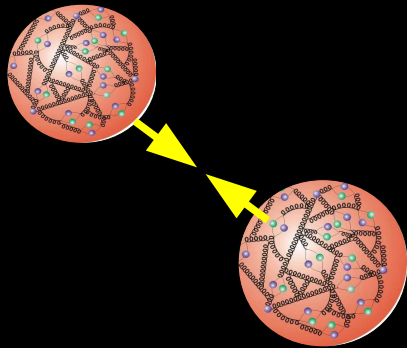
The world wide LHC Computing Grid

Günter Quast

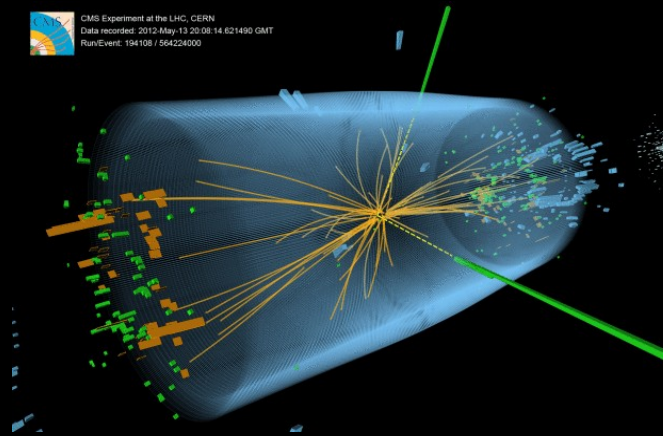
HAP Workshop Bad Liebenzell, Dark Universe

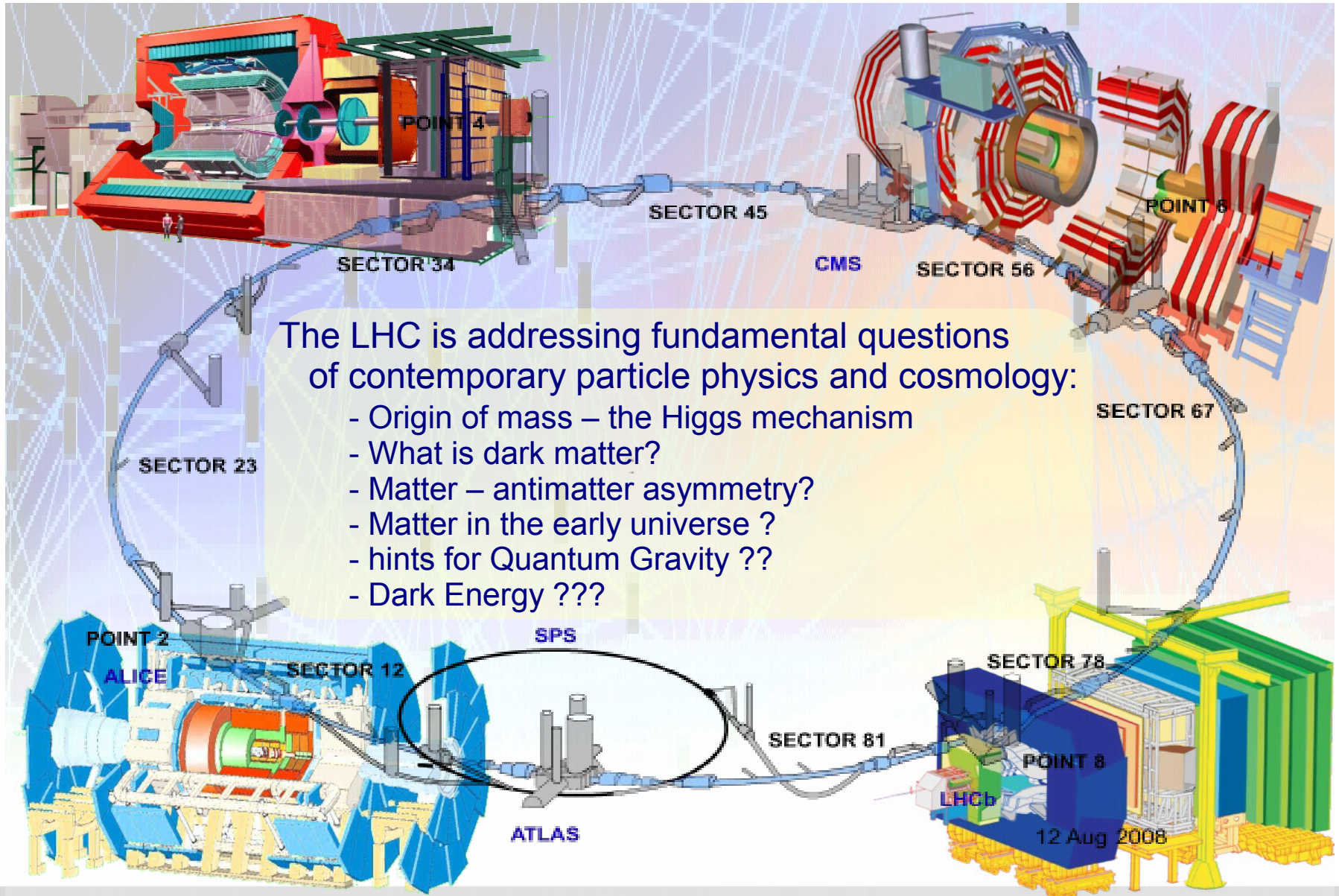
Nov. 22nd 2012

Institut für Experimentelle Kernphysik



CMS Experiment of the LHC, CERN
Data recorded: 2012-May-13 20:08:14.621490 GMT
Run/Event: 194108 / 564224000





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

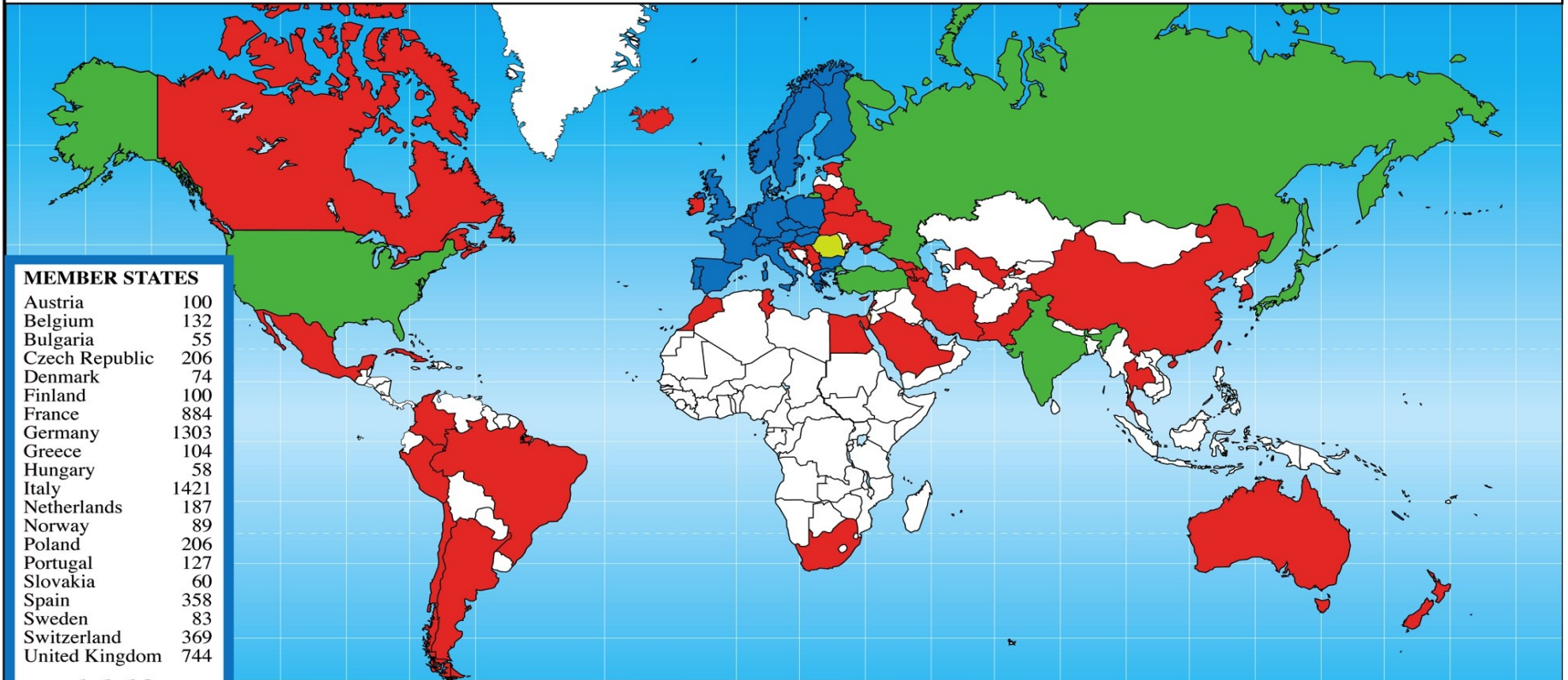
- Origin of mass – the Higgs mechanism
- What is dark matter?
- Matter – antimatter asymmetry?
- Matter in the early universe ?
- hints for Quantum Gravity ??
- Dark Energy ???

Particle Physics is international teamwork



Working @Cern: ~300 institutes from Europe, ~ 7000 Users
~300 institutes elsewhere, ~ 4000 Users

Distribution of All CERN Users by Nation of Institute on 9 January 2012



MEMBER STATES

Austria	100
Belgium	132
Bulgaria	55
Czech Republic	206
Denmark	74
Finland	100
France	884
Germany	1303
Greece	104
Hungary	58
Italy	1421
Netherlands	187
Norway	89
Poland	206
Portugal	127
Slovakia	60
Spain	358
Sweden	83
Switzerland	369
United Kingdom	744

6660

OBSERVERS

India	115
Japan	225
Russia	856
Turkey	77
USA	1708

2981

CANDIDATE FOR ACCESSION

Romania	75
---------	----

ASSOCIATE MEMBER IN THE PRE-STAGE TO MEMBERSHIP

Israel	62
--------	----

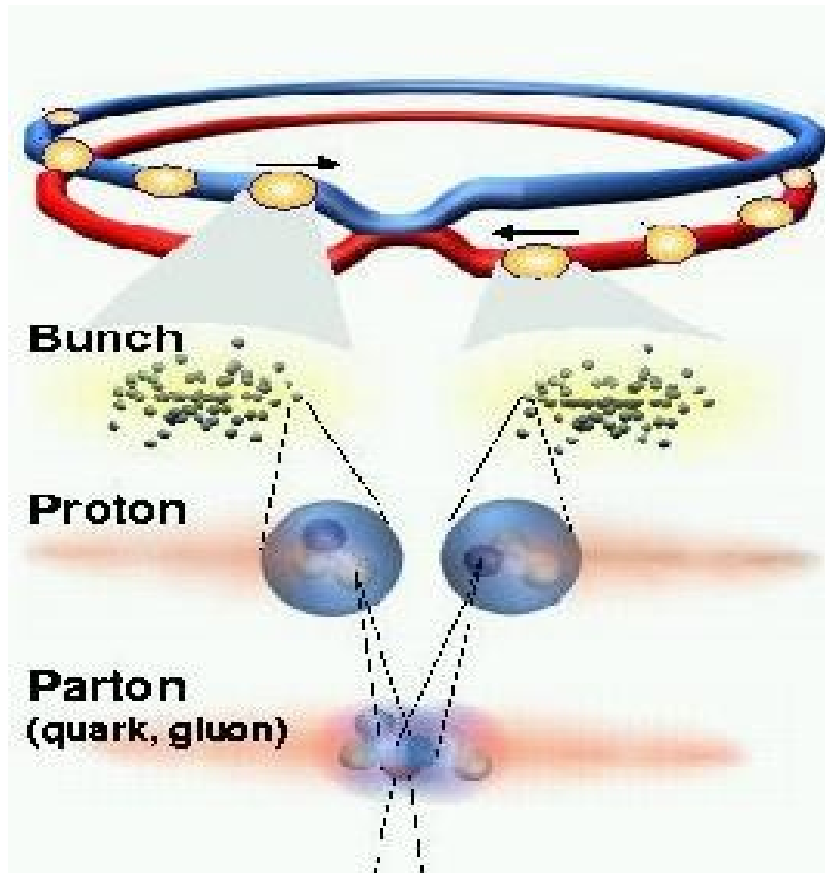
OTHERS

Argentina	18	China	95	Iran	14	Pakistan	19	Ukraine	21
Armenia	12	China (Taipei)	67	Ireland	10	Peru	2	Uzbekistan	1
Australia	24	Colombia	10	Korea	89	Qatar	1		
Azerbaijan	1	Croatia	17	Lebanon	1	Saudi Arabia	3		
Belarus	22	Cuba	4	Lithuania	12	Serbia	26		
Brazil	93	Cyprus	9	Malta	1	Slovenia	37		
Canada	167	Egypt	7	Mexico	43	South Africa	21		
Chile	4	Estonia	18	Montenegro	1	Thailand	5		
		Georgia	10	Morocco	5	T.F.Y.R.O.M.	2		
		Iceland	3	New Zealand	11	Tunisia	1		

907

LHC Parameter

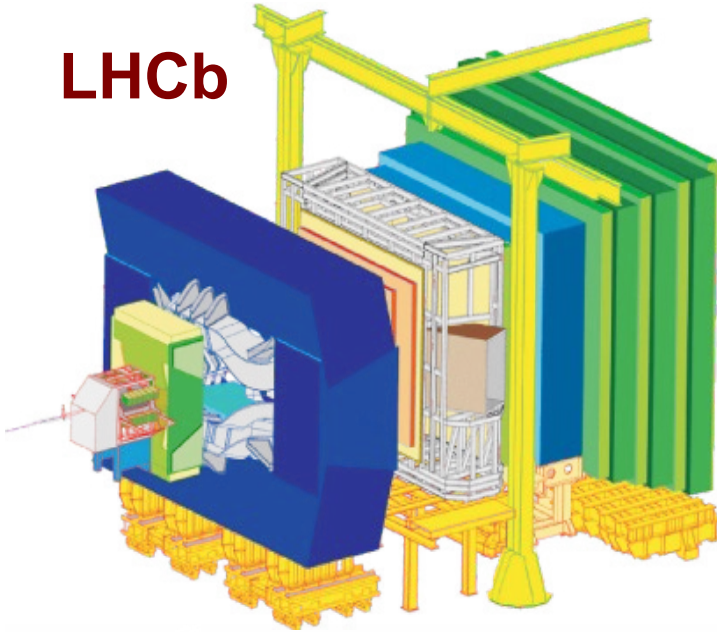
Luminosity: $\sim 10^{34} \text{ cm}^{-2}\text{s}^{-1}$ design, already exceeded by now



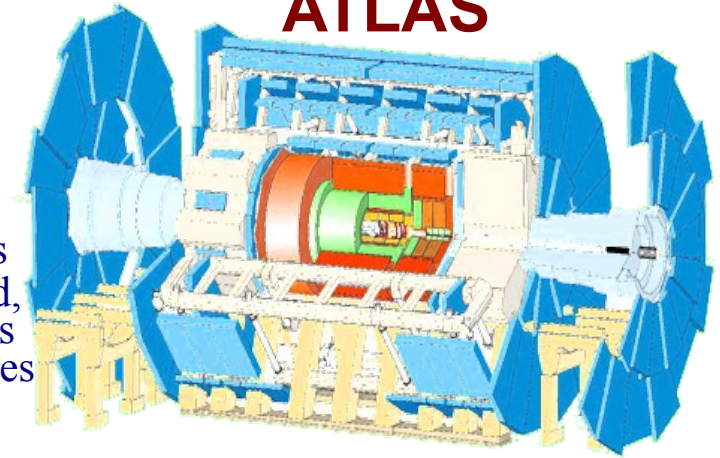
- max. 2835x2835 Proton-Proton bunches
- 10^{11} protons/bunch
- proton energie: 7 TeV
- collision rate: 40 Mhz
- up to 10^9 pp-interactions/sec
- 4 detektors
 - ATLAS** mult-purpouse, pp
 - CMS** multi-purpouse, pp
 - LHCb** asymmetric for b physics
 - ALICE** heavy ion physics
(Pb \leftrightarrow Pb and p \leftrightarrow Pb)

Data sources

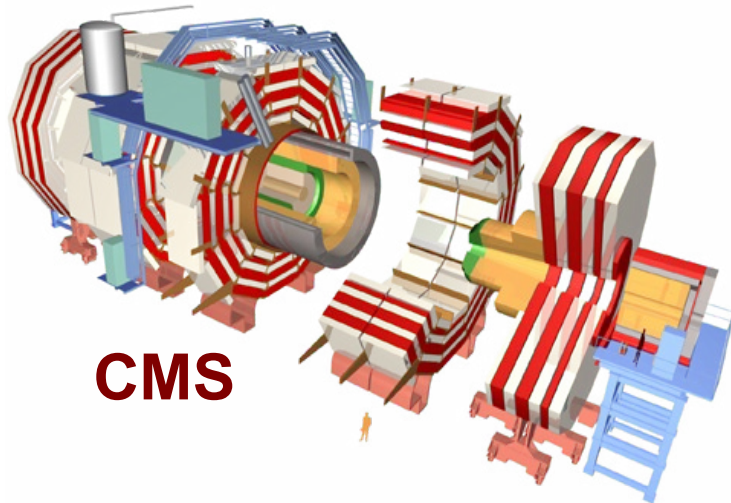
LHCb



ATLAS

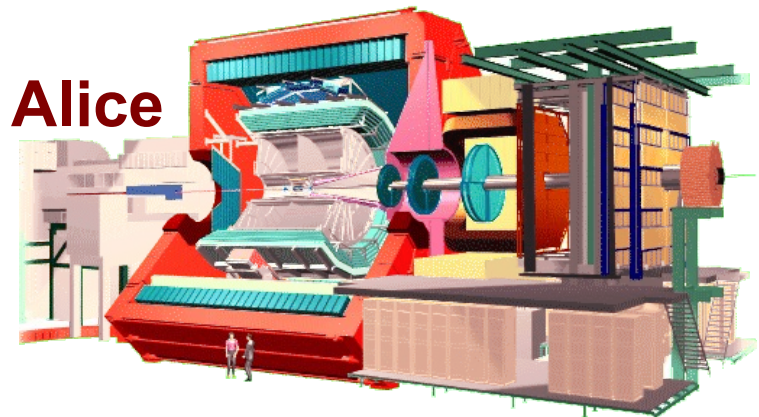


- Built by institutes all over the world, ~10'000 physicists from ~70 countries
- Each detector has more than 100 million sensors
- 40 million pp-collisions recorded per second
- detectors specialized on different questions



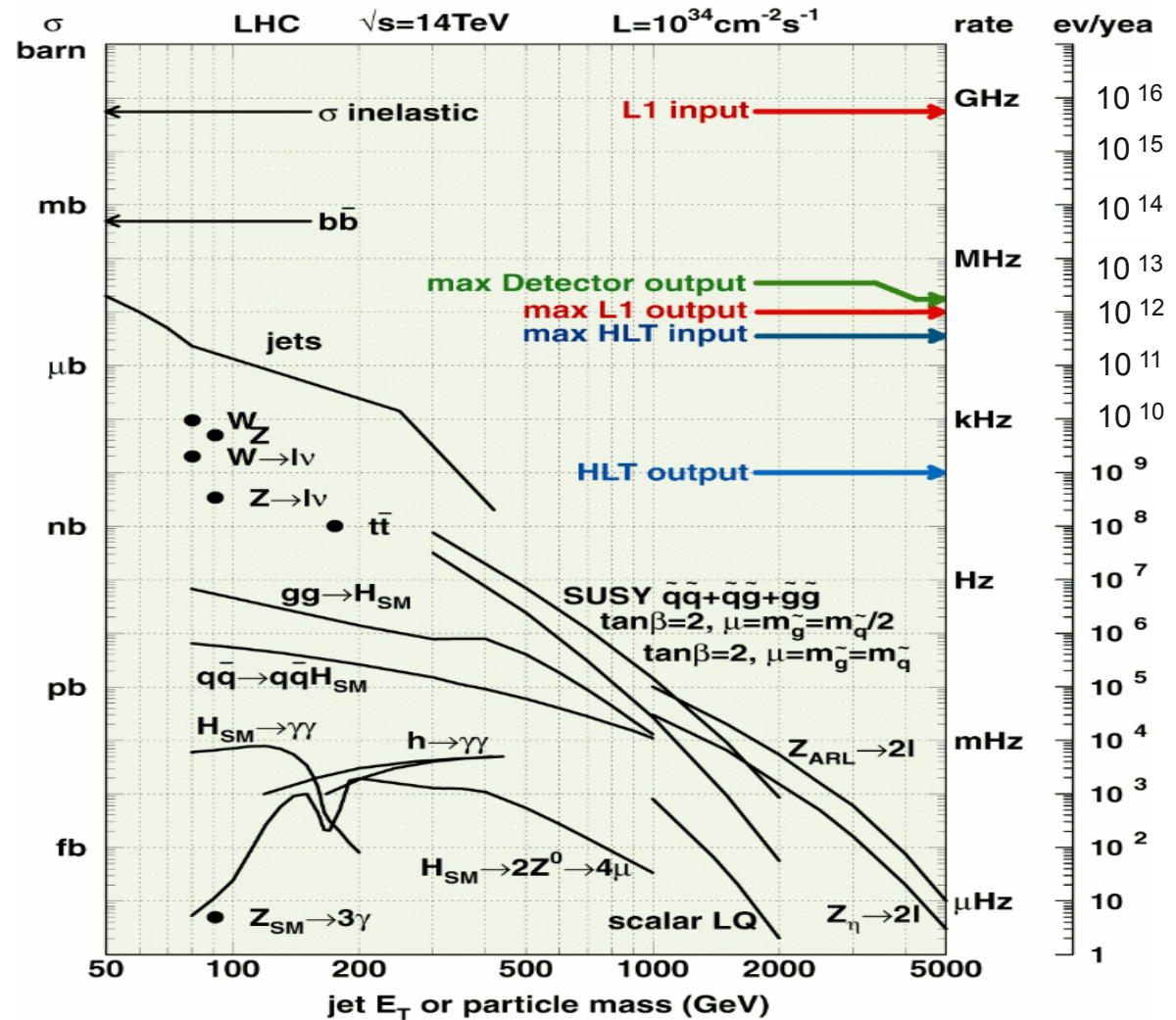
CMS

Alice



The Challenge

Rates of
 “interesting physics”
 ~ 10^{10} times smaller than
 inelastic pp cross section



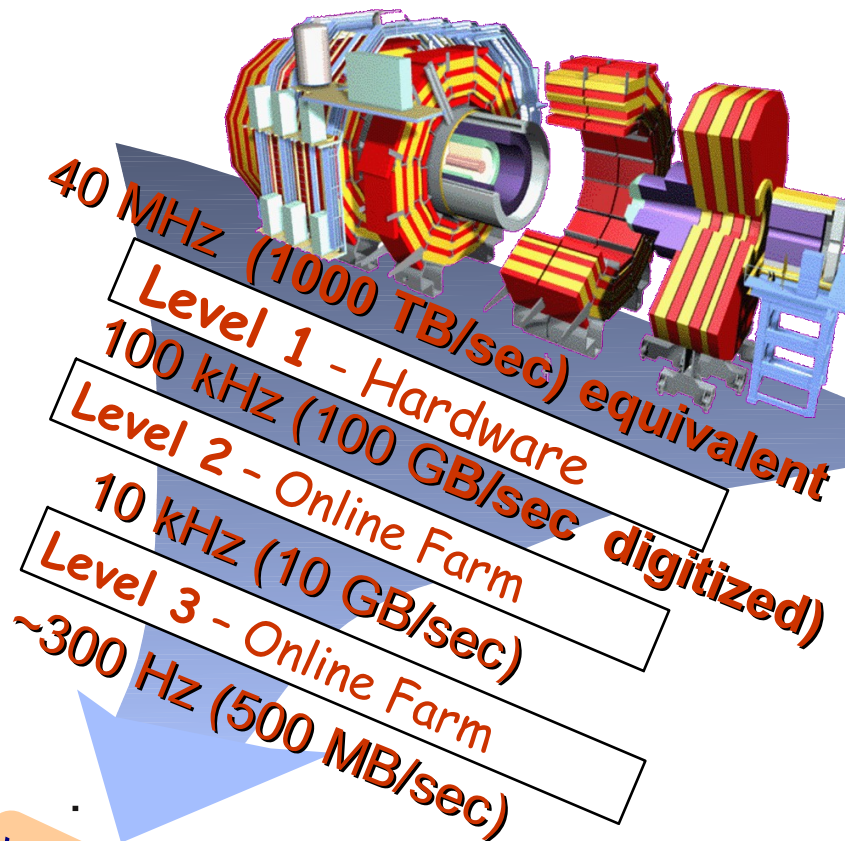
Data Sources – example CMS

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

→ **~1000 Tbyte/s raw data**

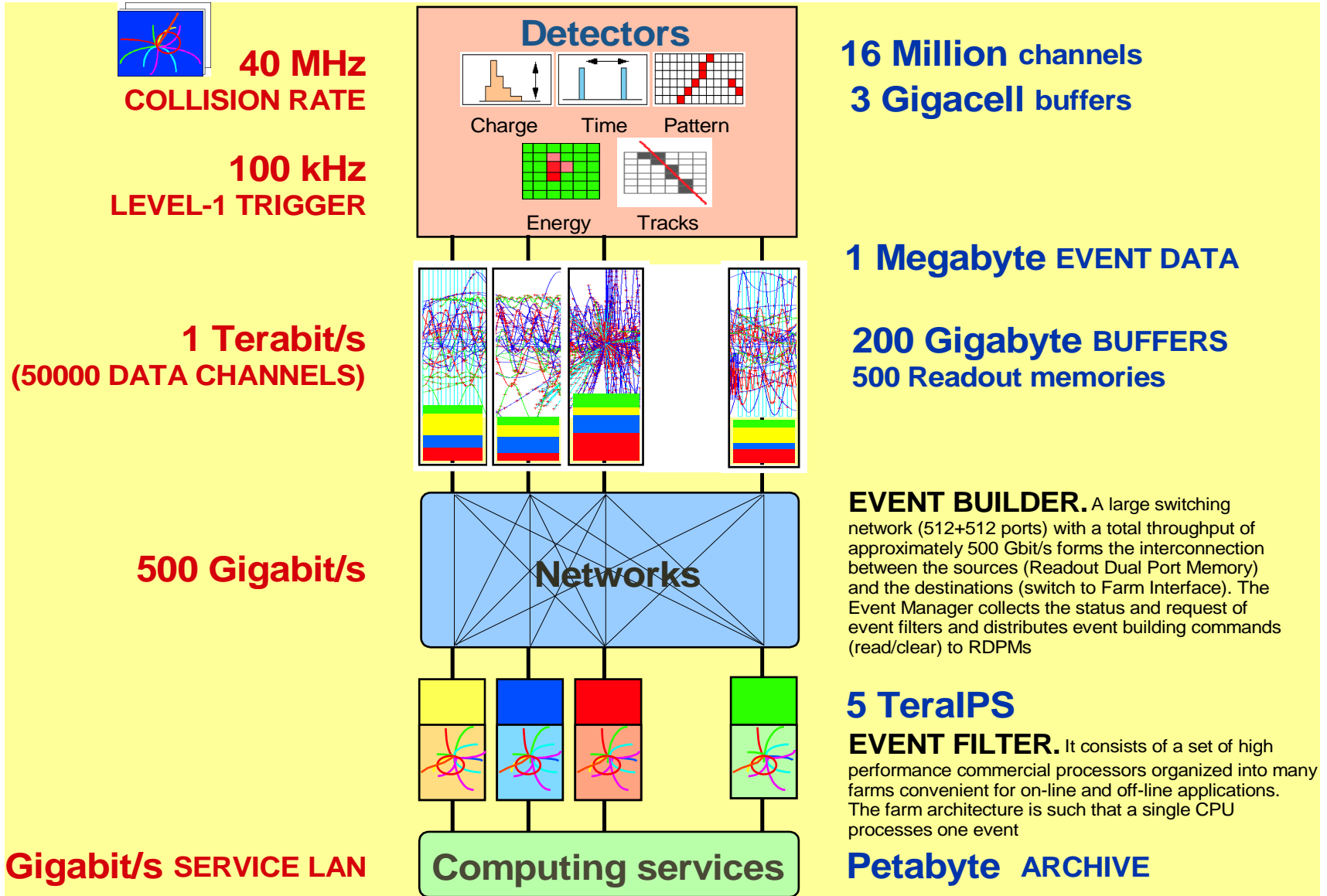
- Zero-suppression and Trigger reduce this to „only“ some 100 Mbyte/s

i.e. 1  /sec

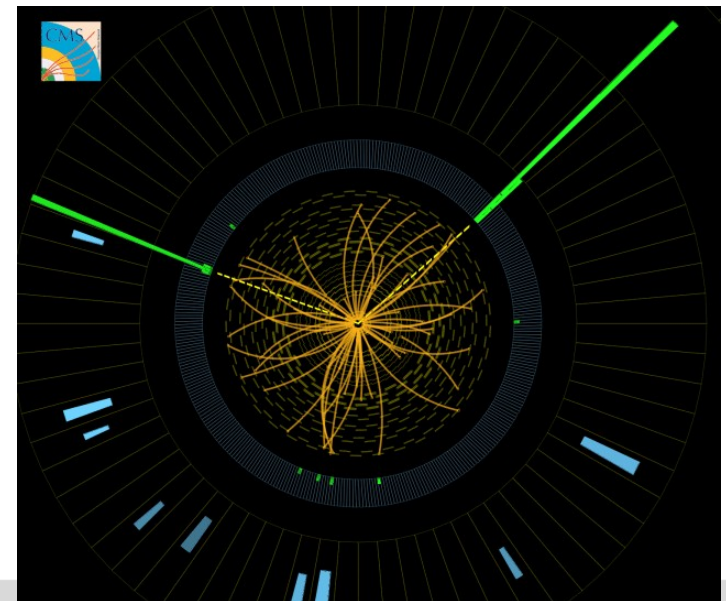
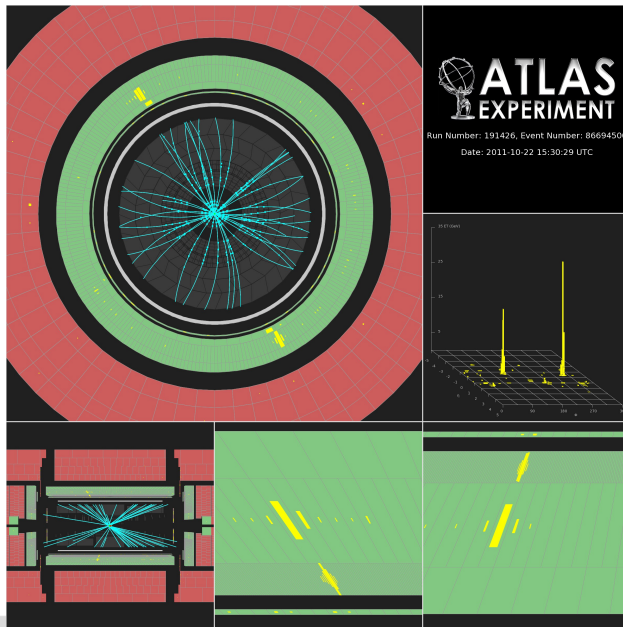
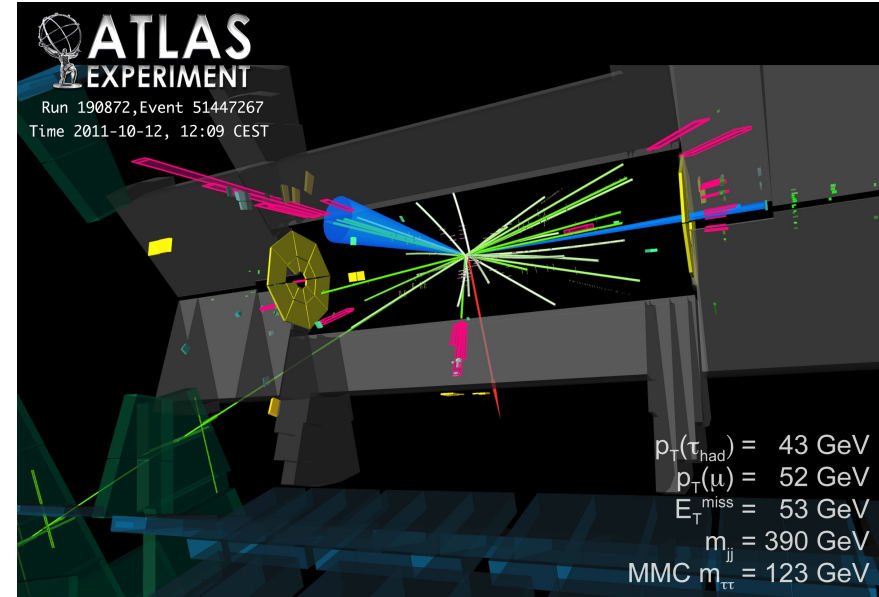
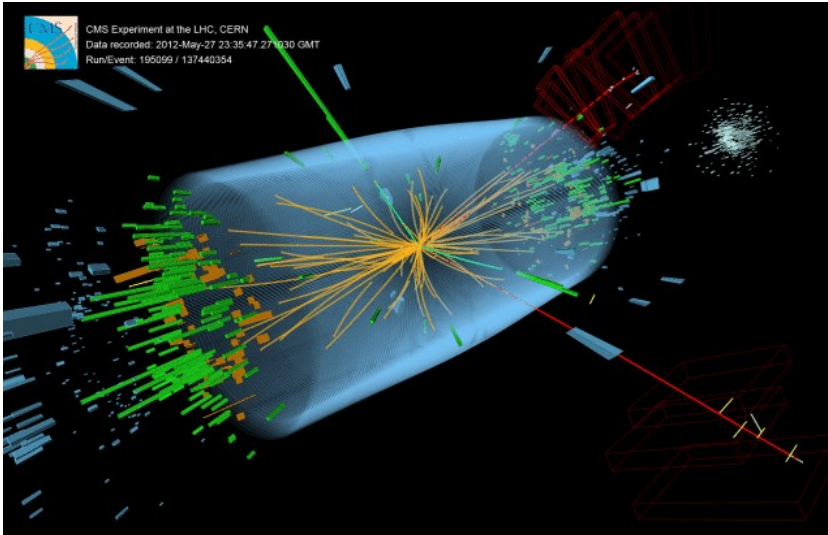


Distributed Computing
is ideal to do this job !

world-wide
physics
community

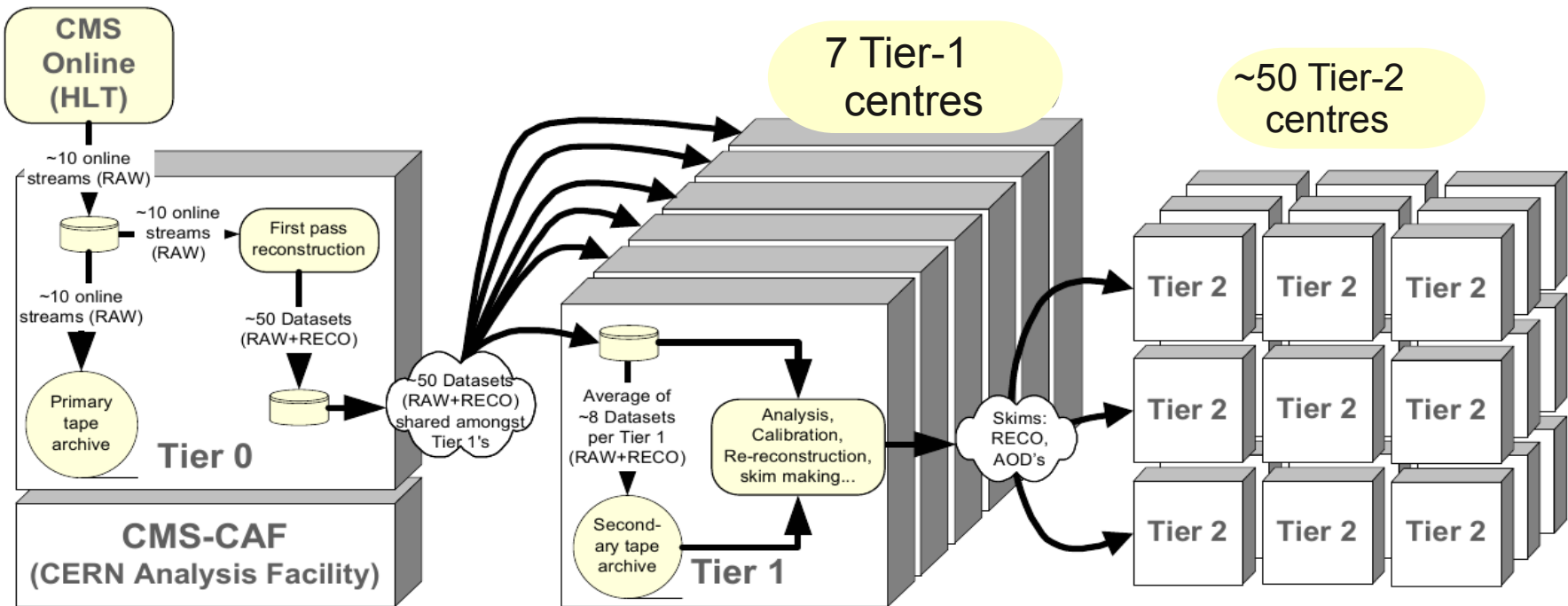


The output: reconstructed events



Computing Models – a hierarchical structure

example: C MS



CMS Computing TDR, 2005

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

The Grid Paradigm ...

... came at the right time to provide a model for distributed computing in HEP

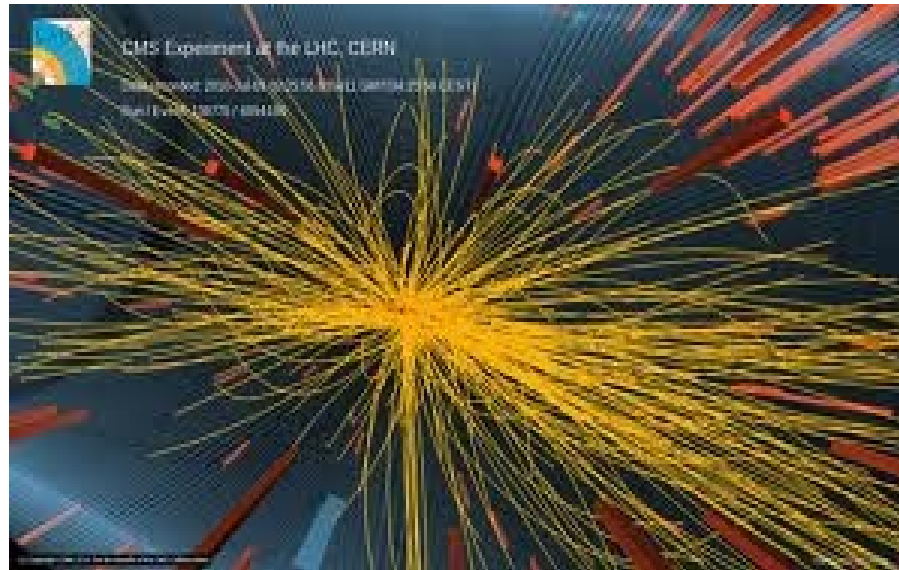
“A computational grid is a hardware and software infrastructure that provides dependable, consistent, pervasive, and inexpensive access to high-end computational capabilities.”

C. Kesselman, I. Foster, 1998

- *Coordinates resources that are not subject to central control ...*
- *... using standard, open, general-purpose protocols and interfaces ...*
- *... to deliver nontrivial quality of services*

I. Foster, 2002

Why Grid is well suited for HEP



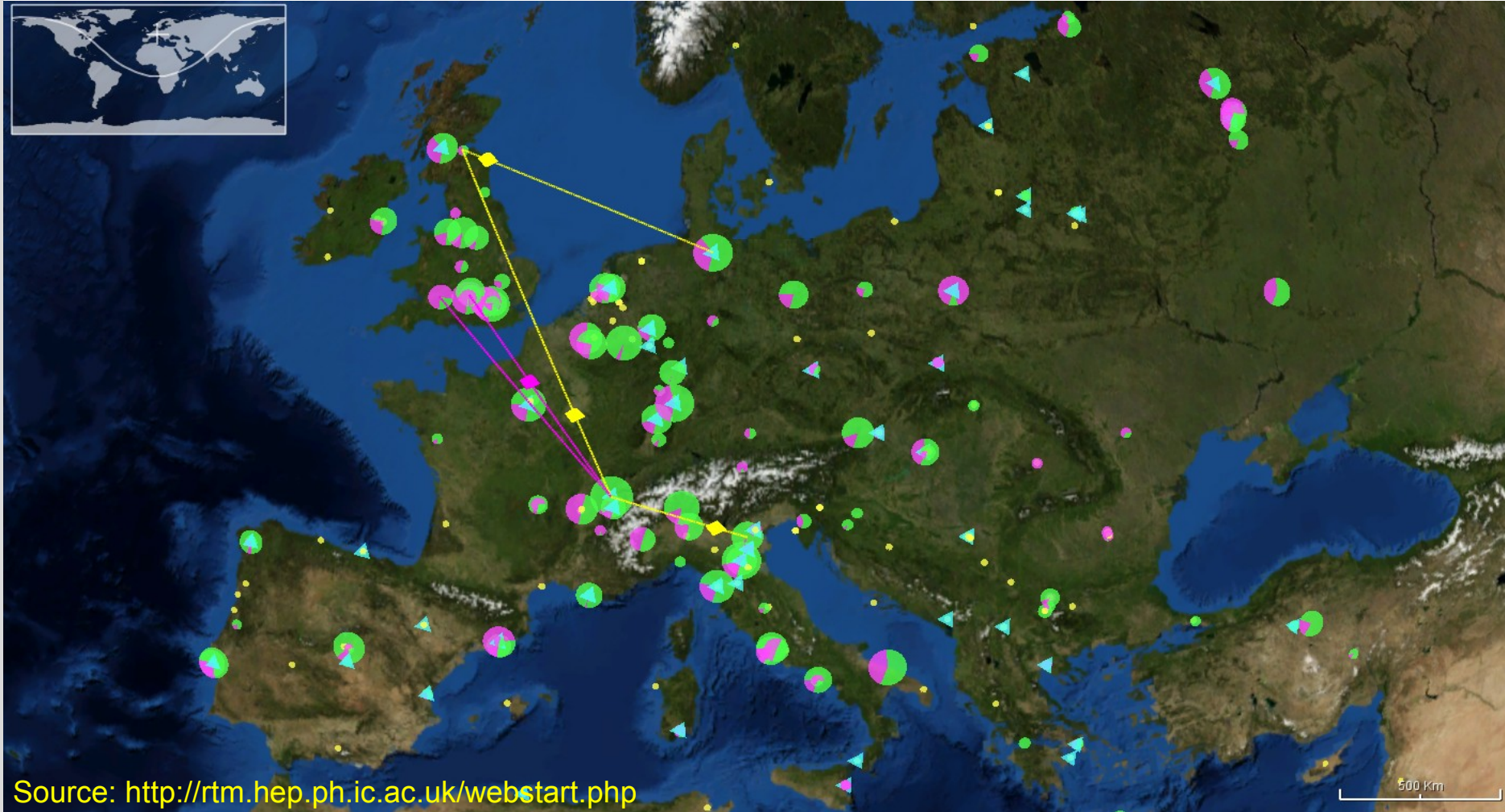
Experimental HEP codes - key characteristics:

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



- 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 hierarcical approach, following the typical data reduction schemes usually adopted in data analysis in high enery 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

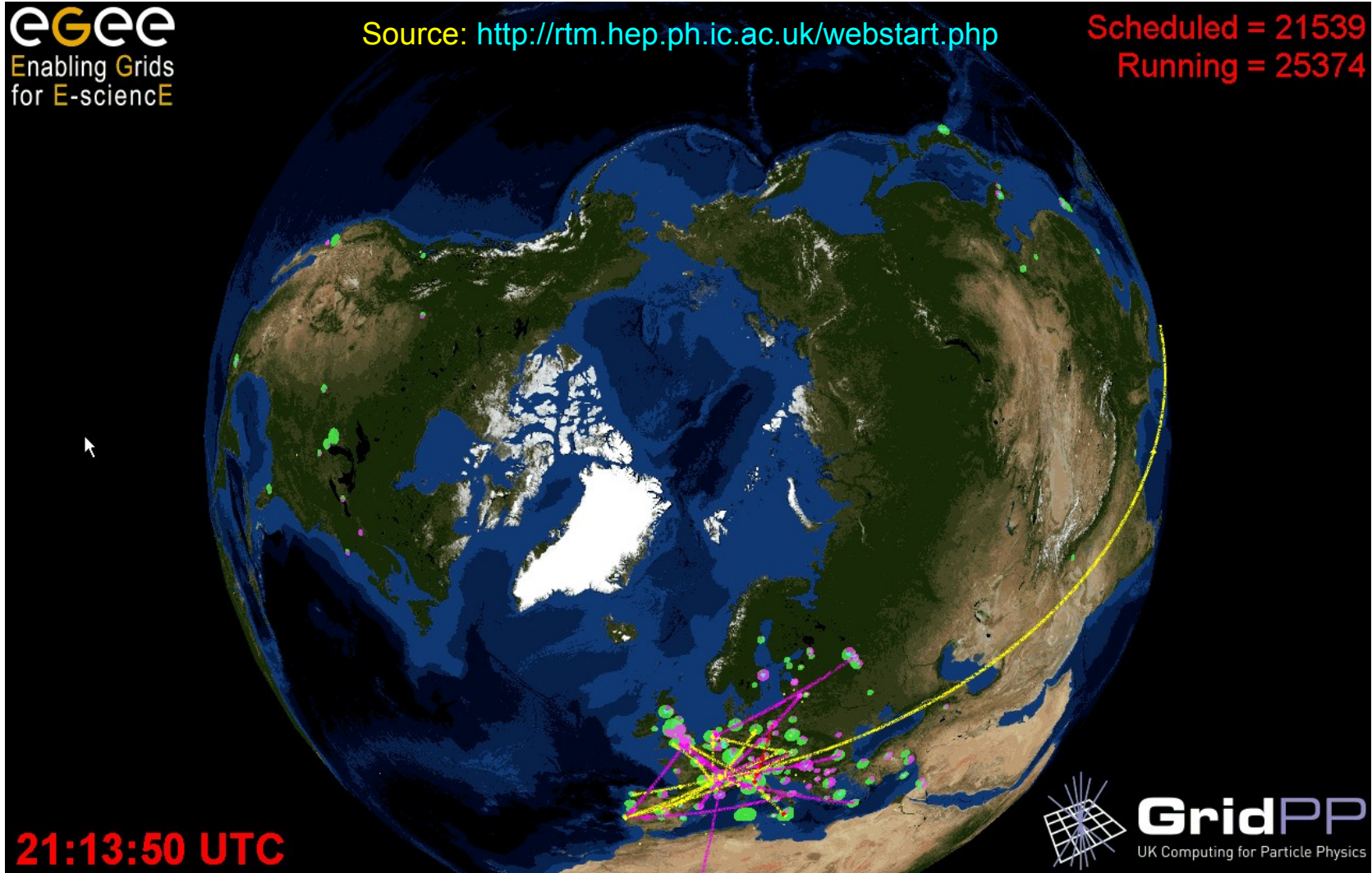


A truly international, world-spanning Grid for LHC data processing, simulation and analysis

eGEE
Enabling Grids
for E-science

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

Scheduled = 21539
Running = 25374



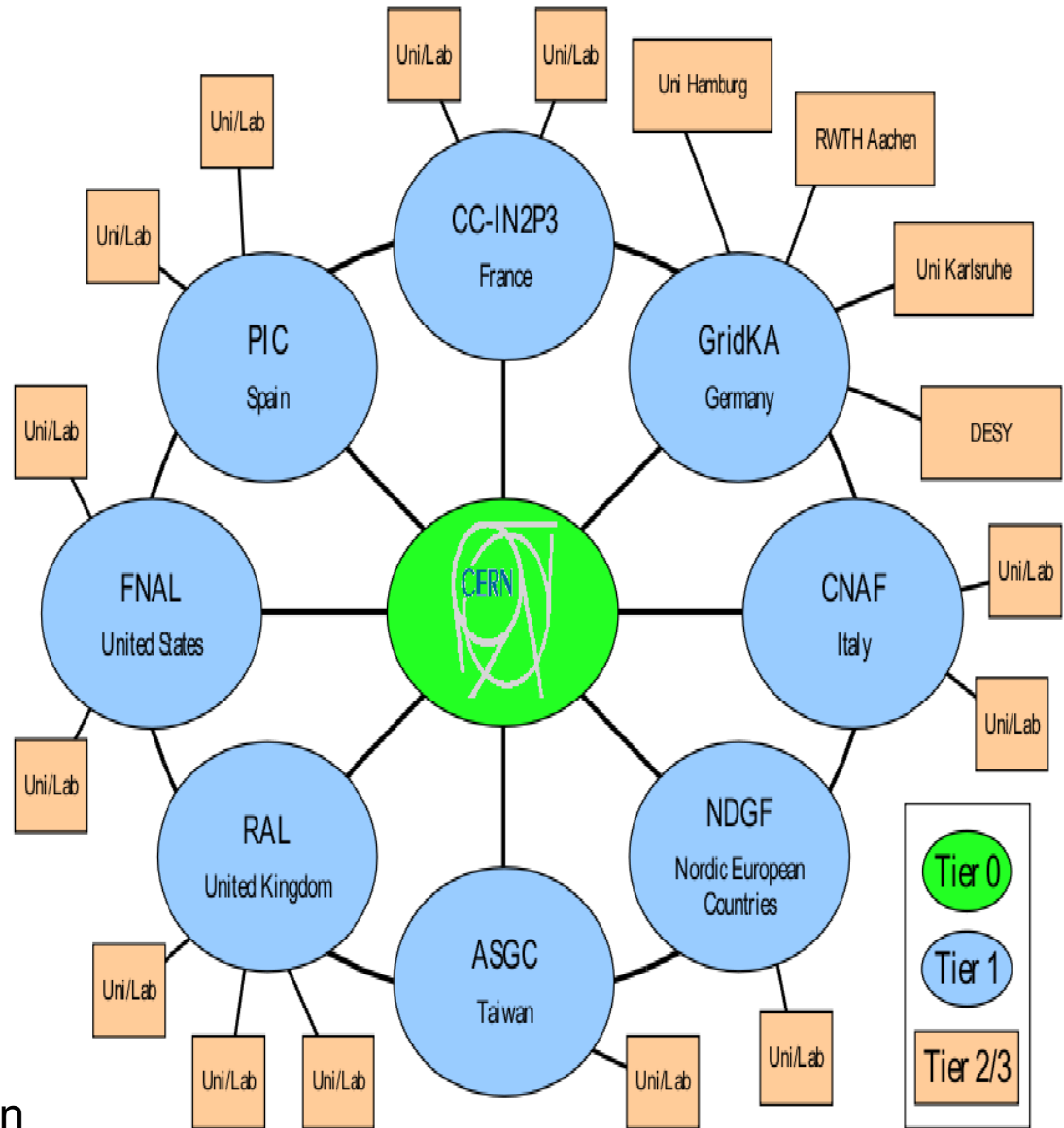
A truly international, world-spanning Grid for LHC data processing, simulation and analysis

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:

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

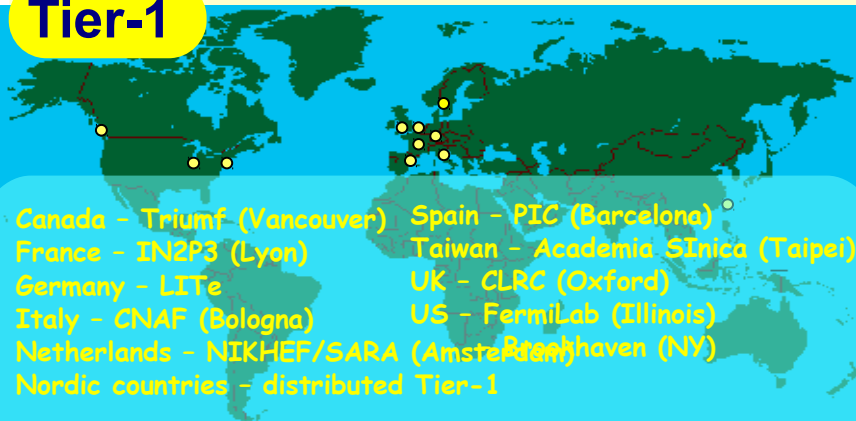


Tier-0 the accelerator centre

- 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

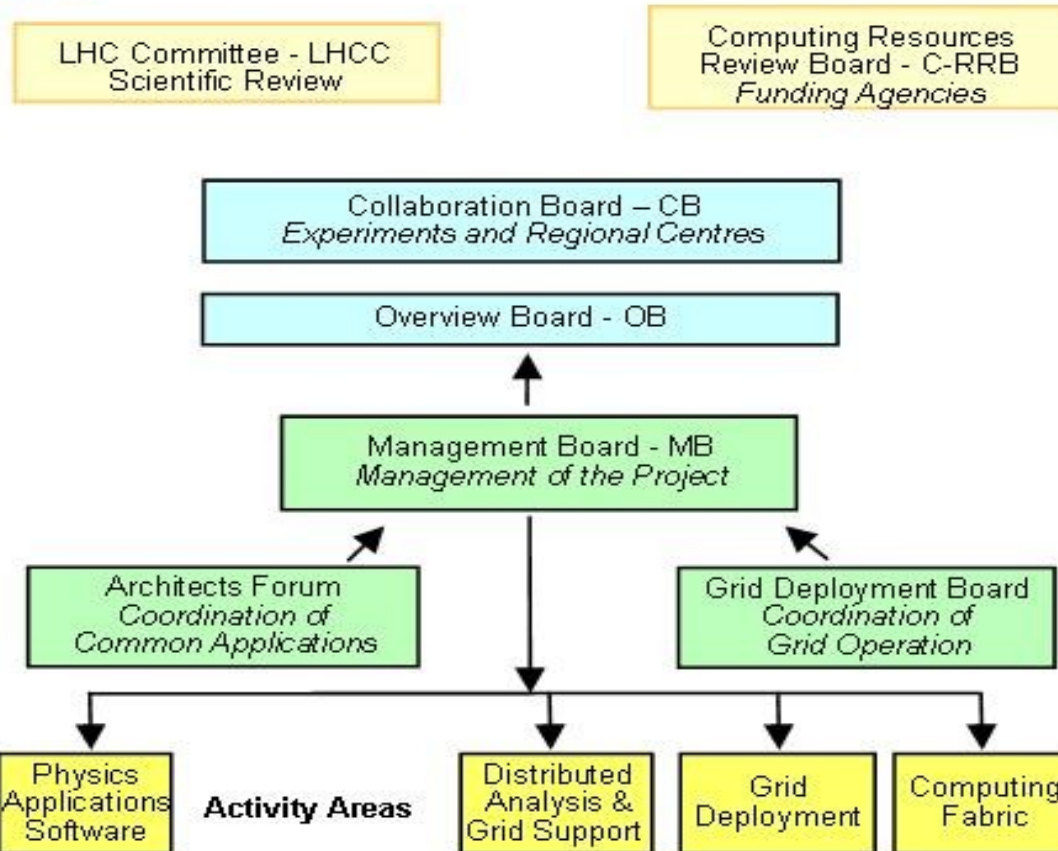
- End-user (physicist, research group) analysis & collaboration with T3
(= institute resources) – **where the discoveries are made**
- Monte Carlo Simulation

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

Organisation of the World-wide LHC computing Grid

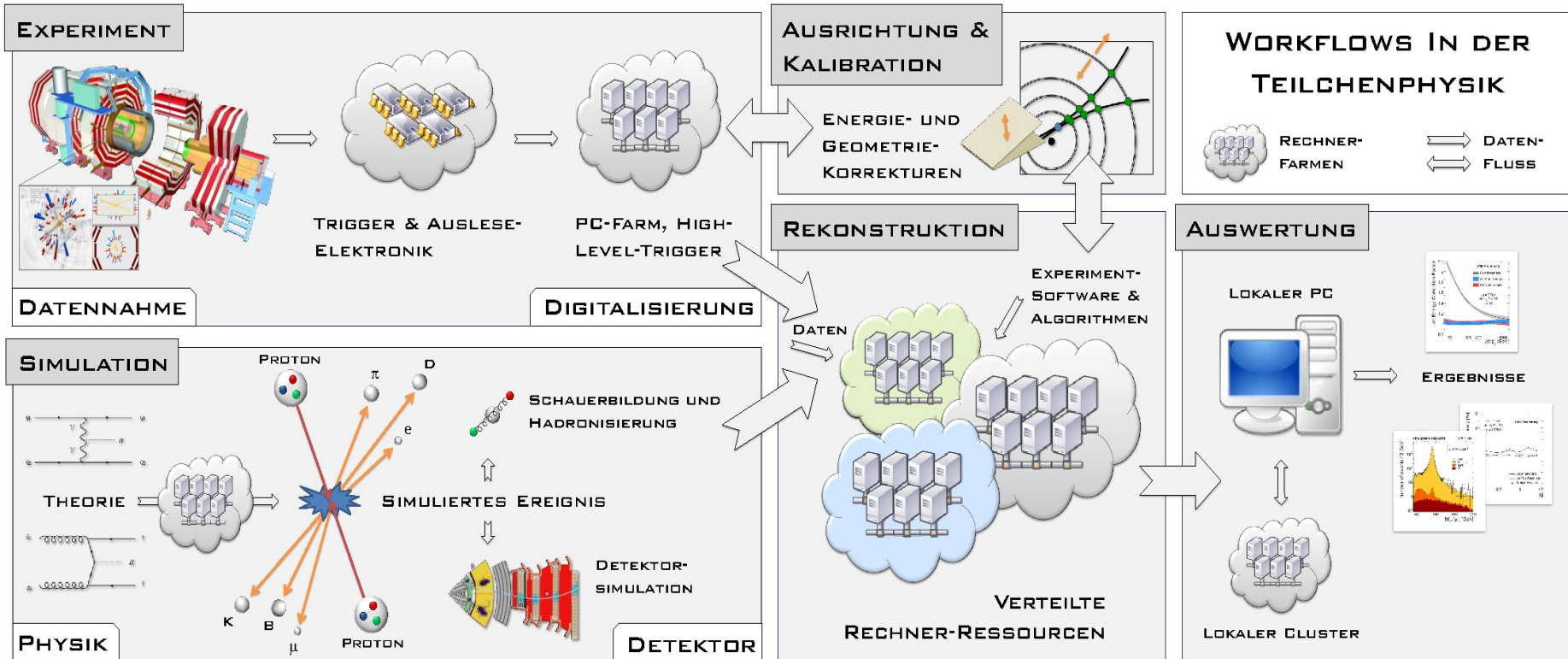


Worldwide LCG Organisation



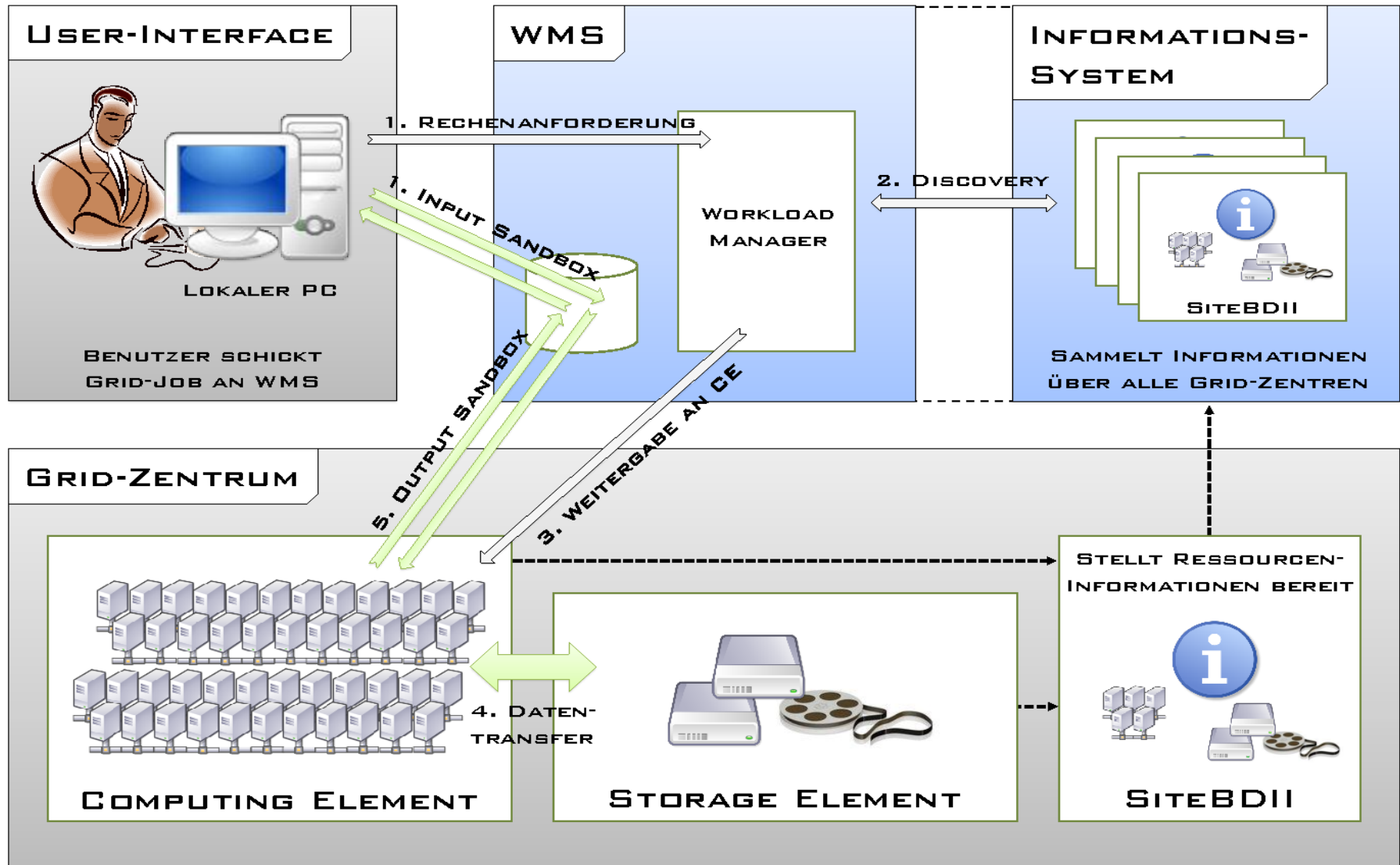
Grids can't work without an organisational structure representing all parties involved

What Particle Physicists do on the Grid

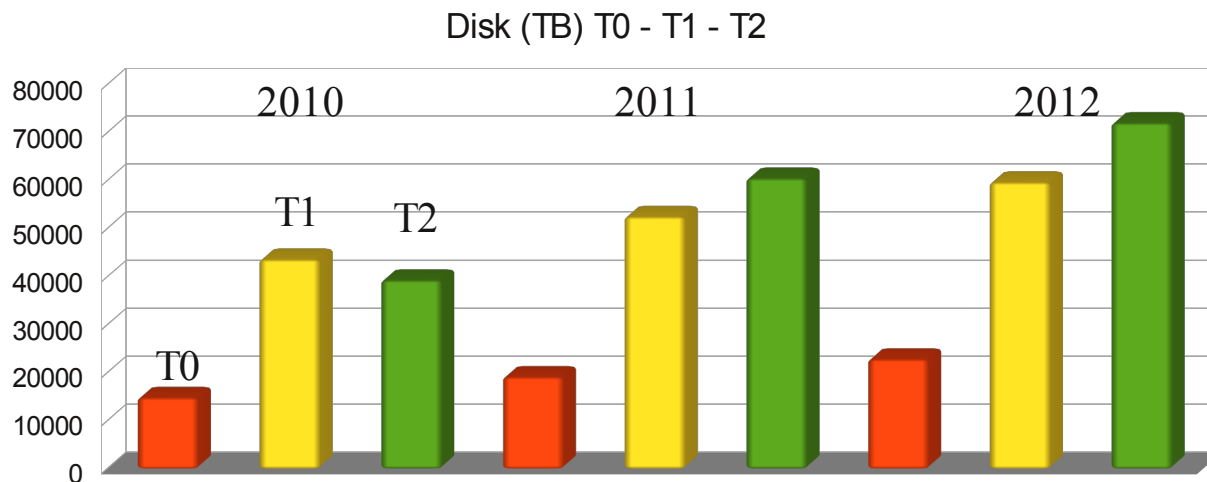
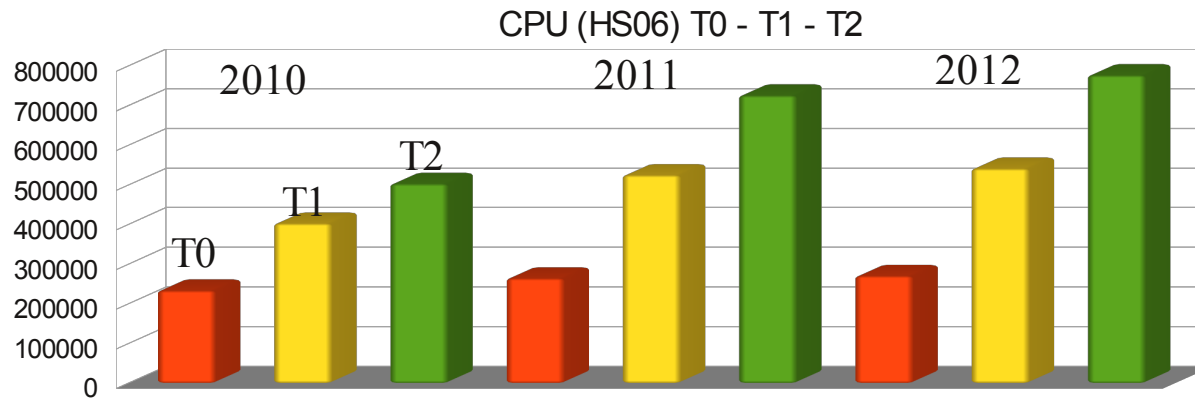


- **CPU-intensive simulation** of particle physics reactions and detector response
- **processing** (=reconstruction) of **large data volumes**
- **I/O-intensive filtering and distribution** of data
- **transfer** to local clusters and workstations for final physics interpretation

Typical workflows on the Grid



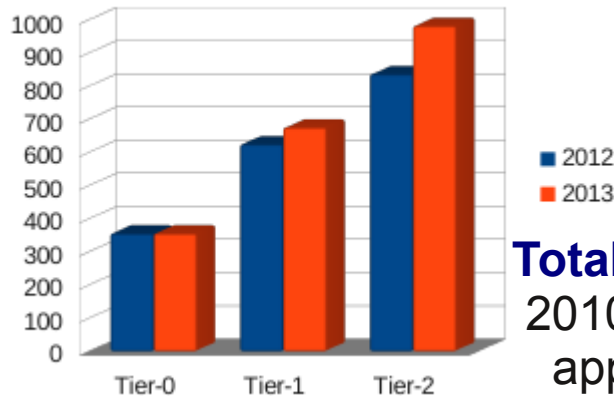
How big is WLCG ?



The largest Science Grid in the World

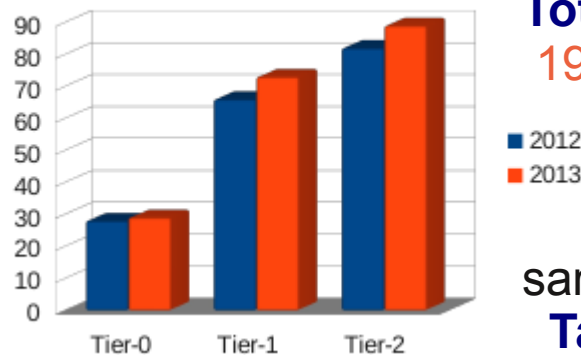
How big is WLCG today and compared with others ?

CPU kHS06



Total CPU 2013:
2010 kHS06
approx. equiv.
200'000 CPU cores

Disk Storage (PB)



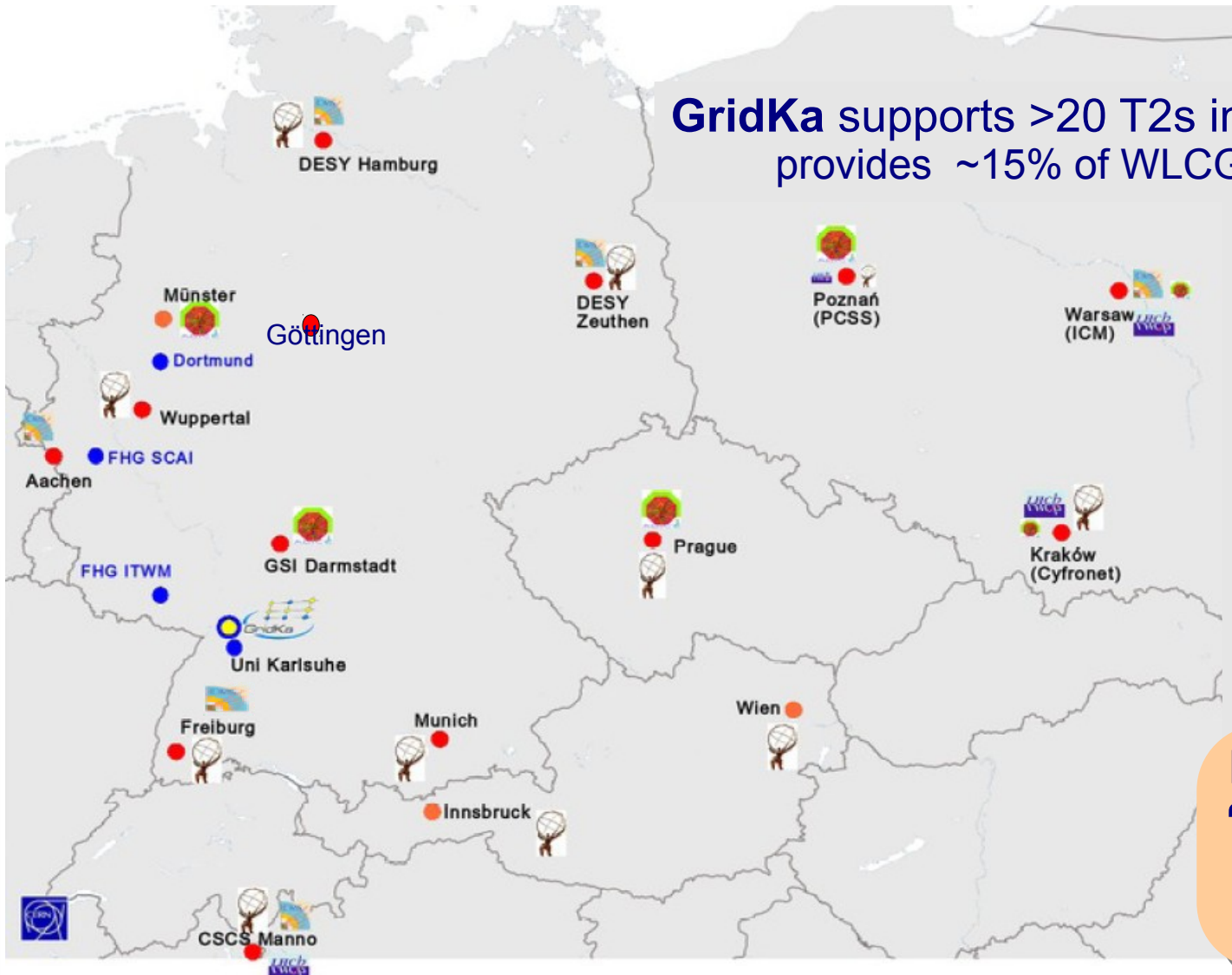
Total Disk 2013:
190 PB,

same amount as
Tape Storage
190 PB

Other “big players” ?

- new HPC Cluster in Munich (SuperMuc, #64 on Top 500 list)
155 k Cores, 10 PB disk
- 1&1 Rechenzentrum Karlsruhe
30'000 Server, ~2 PB/Monat I/O
- Amazon Cloud:
~ 500 PB
- Facebook
~30 PB
- Google Data Center
power 200 MW (WLCG: ~6MW)

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 three years of experience
with LHC operation:**

How well did it work ?

After 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:

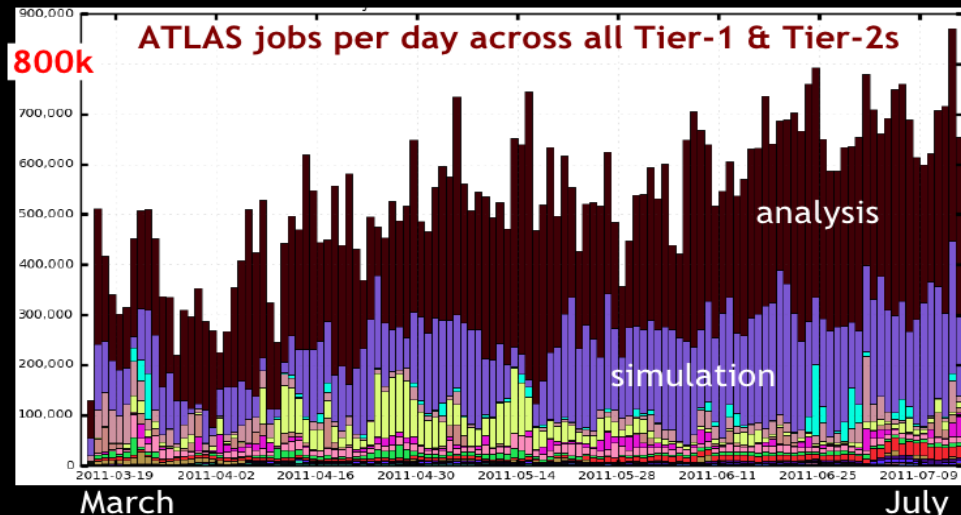
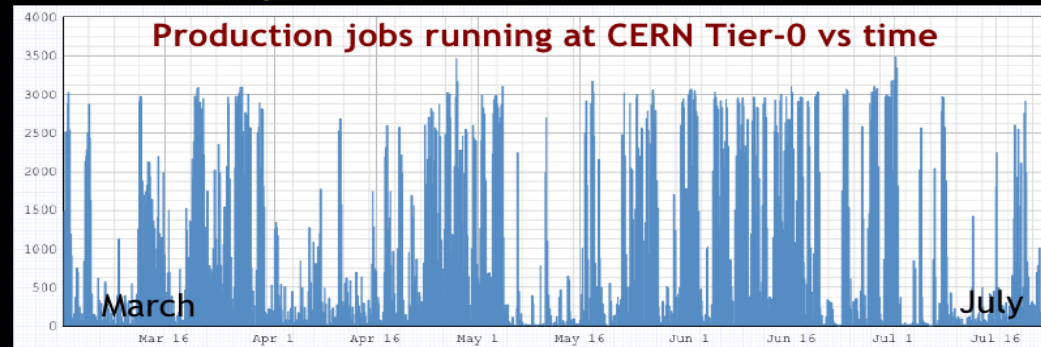
- First-pass reco. at Tier-0 within ~2 days
- Calibration/DQ good for physics analysis
- Data analysable on Grid within ~1 week

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

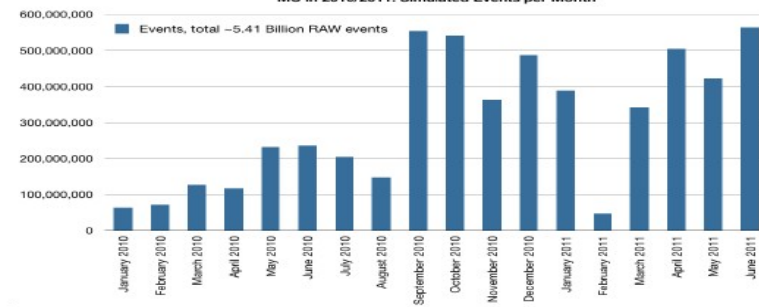
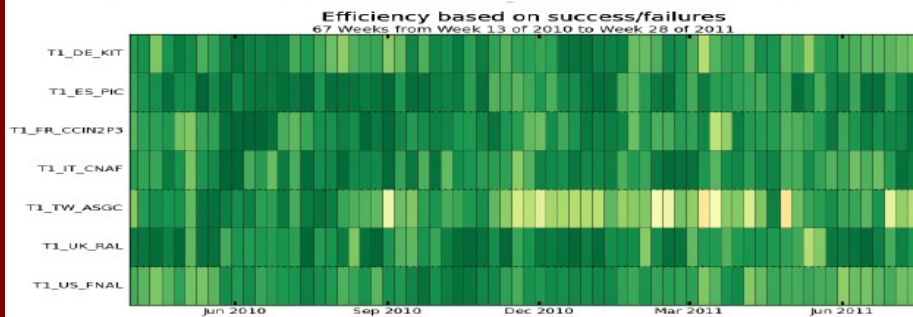
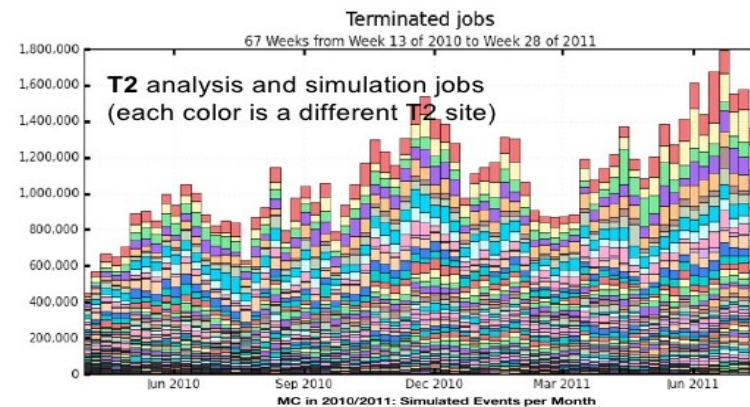
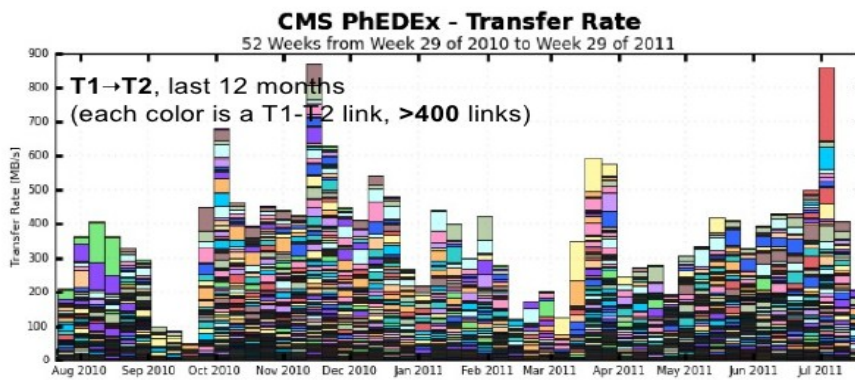


G. Tonelli, CMS, EPS HEP 2011



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

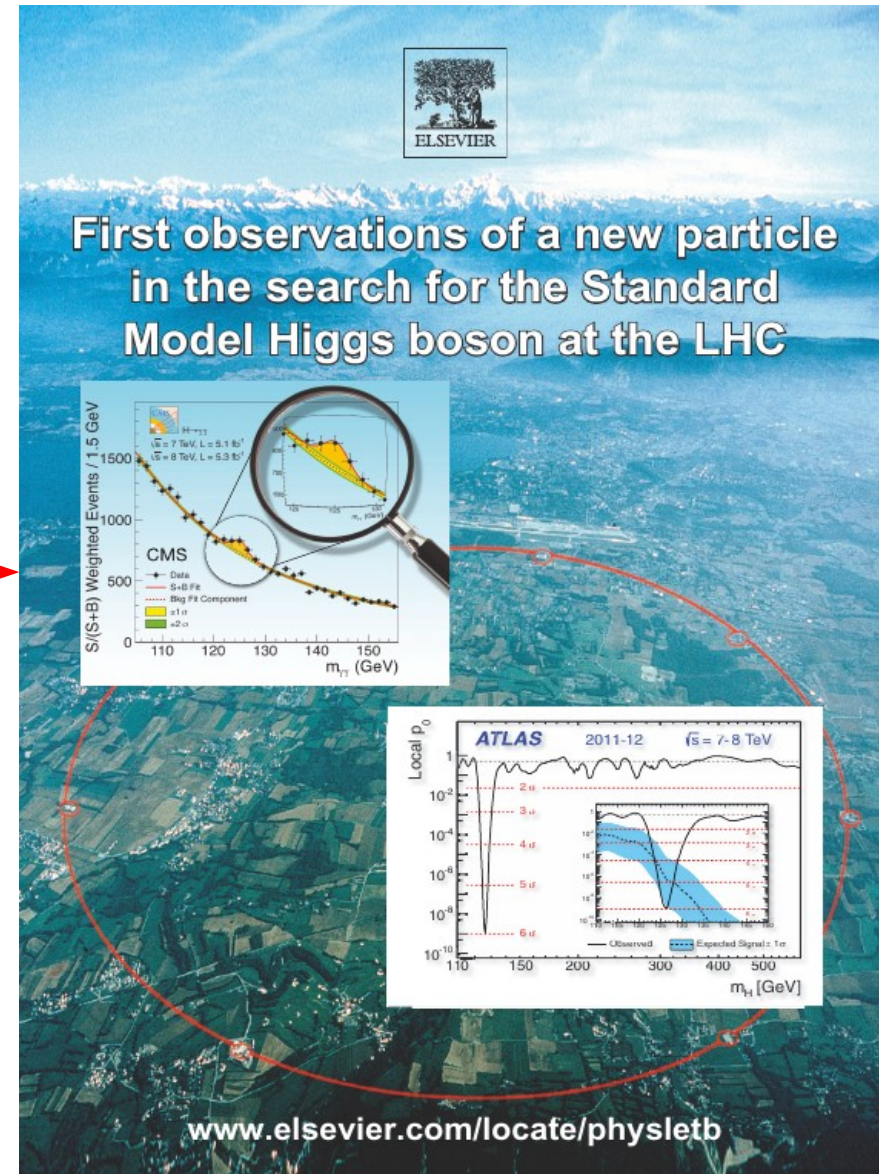
36

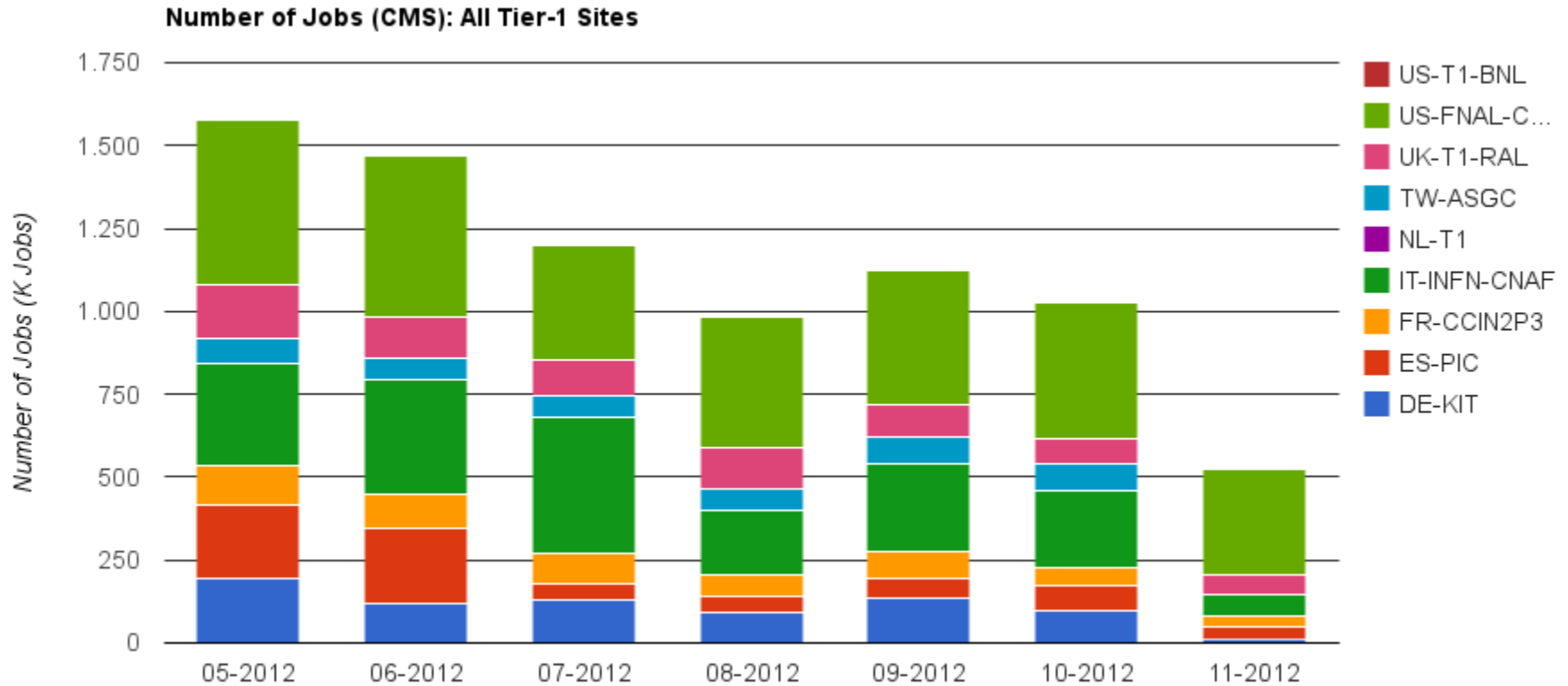
Did it work ?

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!

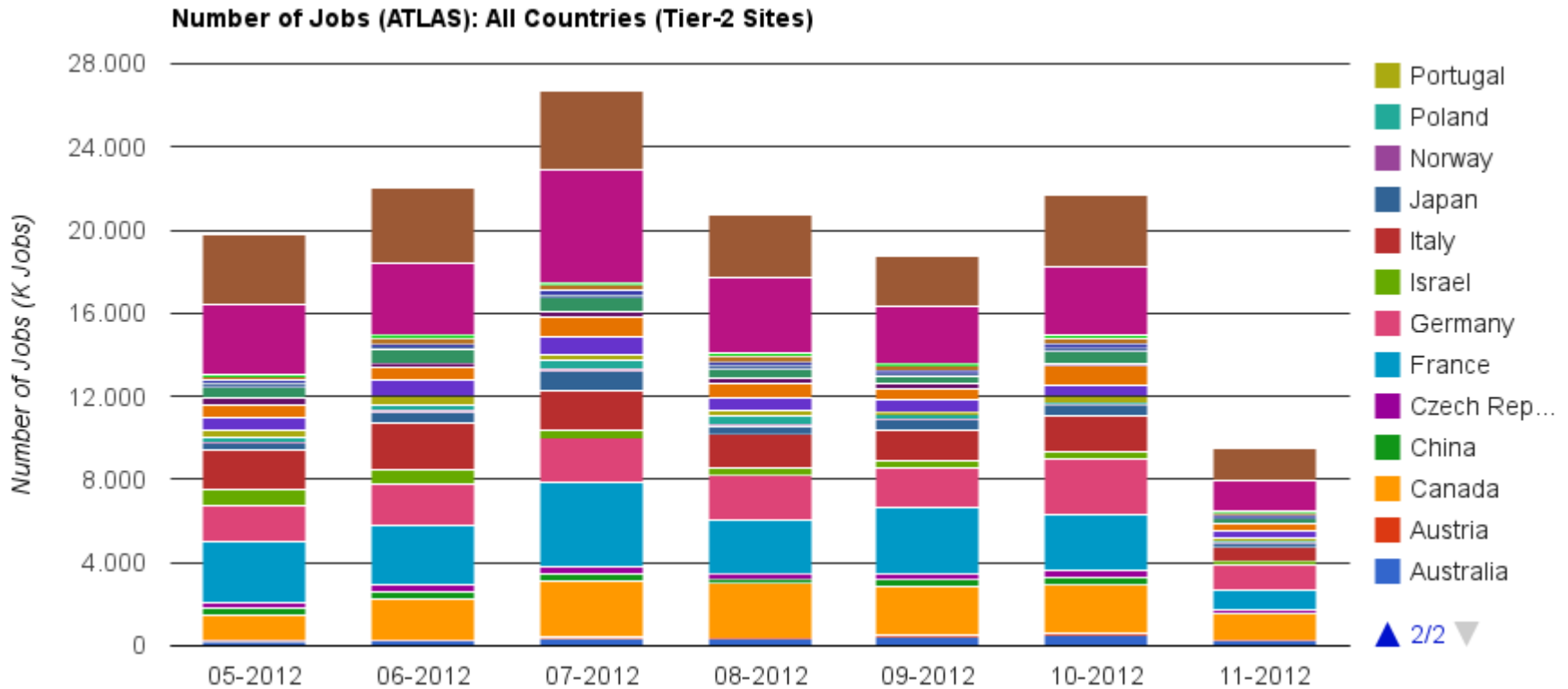
Let`s have a look in detail ...





CMS: ~1 million jobs/month

ATLAS: ~5 million (shorter) jobs/month on Tier1s

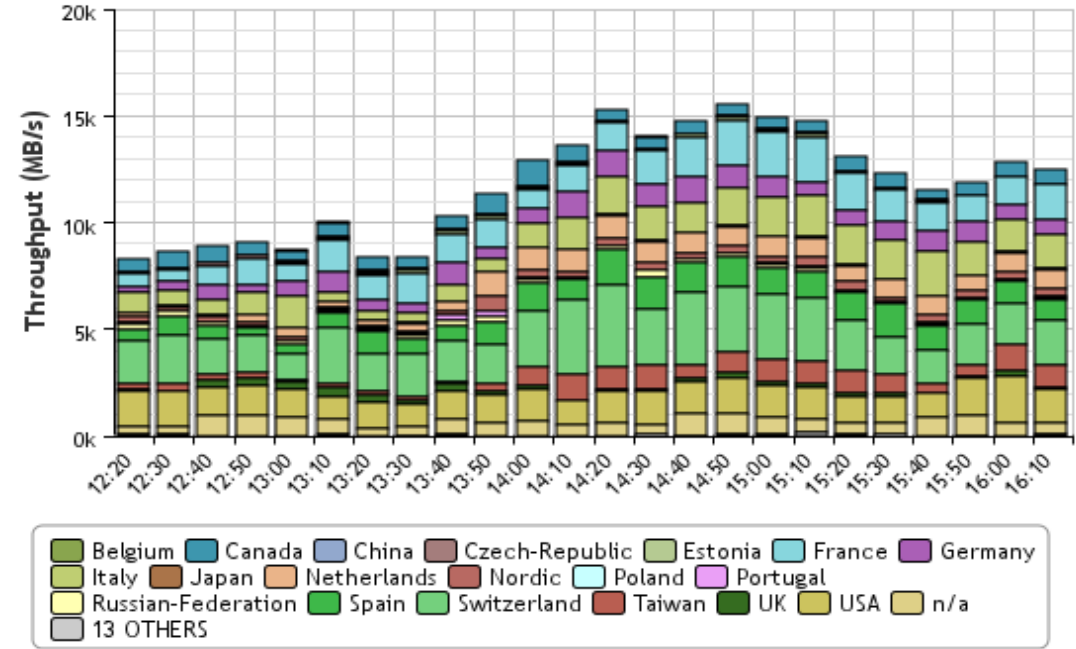
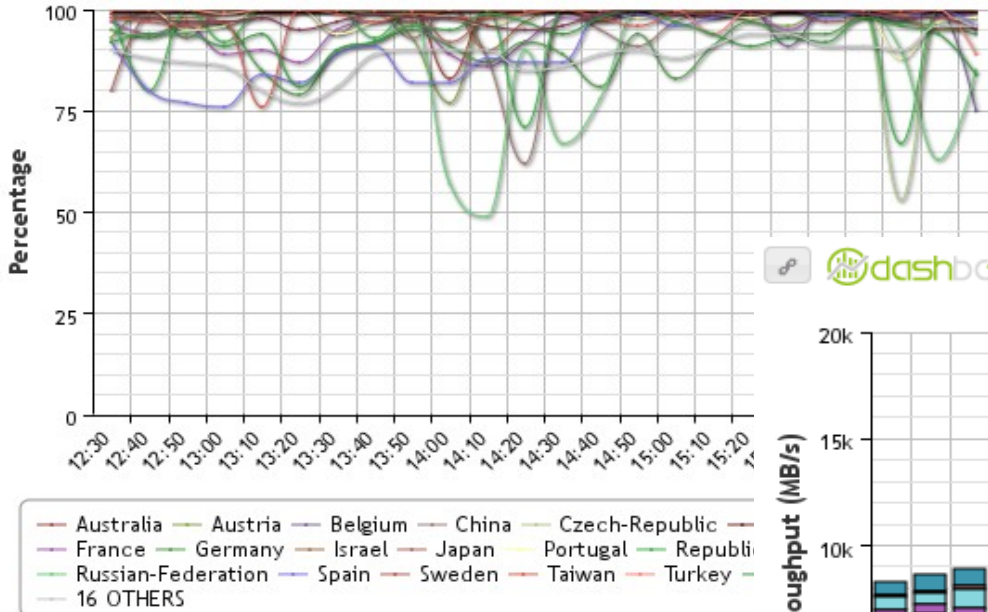


ATLAS: ~20 million jobs/month

CMS: ~ 7 million jobs/month

on Tier2s

Data rates



Routinely run data transfers across the world at a transfer volume of ~50 TBytes/h

Again: Does it work ?

YES !

- Routinely running ~1 million jobs per day on the Grid
- Shipping over 1 PB/day of data around
- data distribution to T2 works well
- very little is known about T3 usage and success rates - responsibility of the institutes
- plenty of resources were available at LHC start-up, now reached “resource limited operation”
- Users have adapted to the “GridWorld” - Grid is routinely used as a huge batch system, output is transferred home

but ...

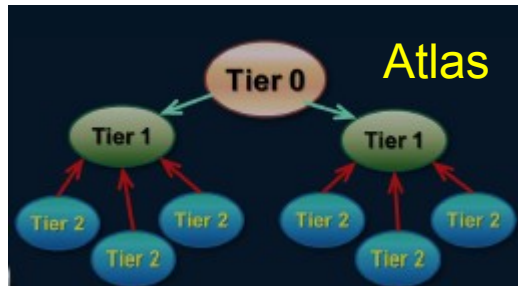
Message:

**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

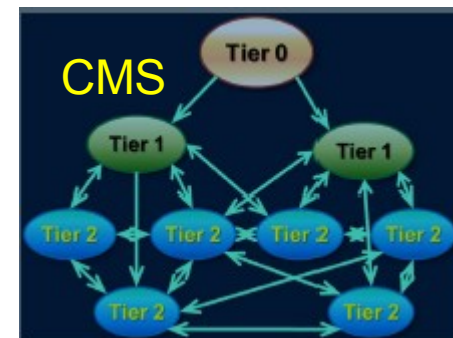
- simplified and more robust middleware
- redundancy of services and sites,
requires dynamic placement of data and investment in network bandwidth
- automated monitoring and triggering of actions
- use of commercially supported approaches to distributed computing:
 - private clouds are particularly important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

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



ATLAS and CMS
computing models differ slightly

CMS already more “distributed”

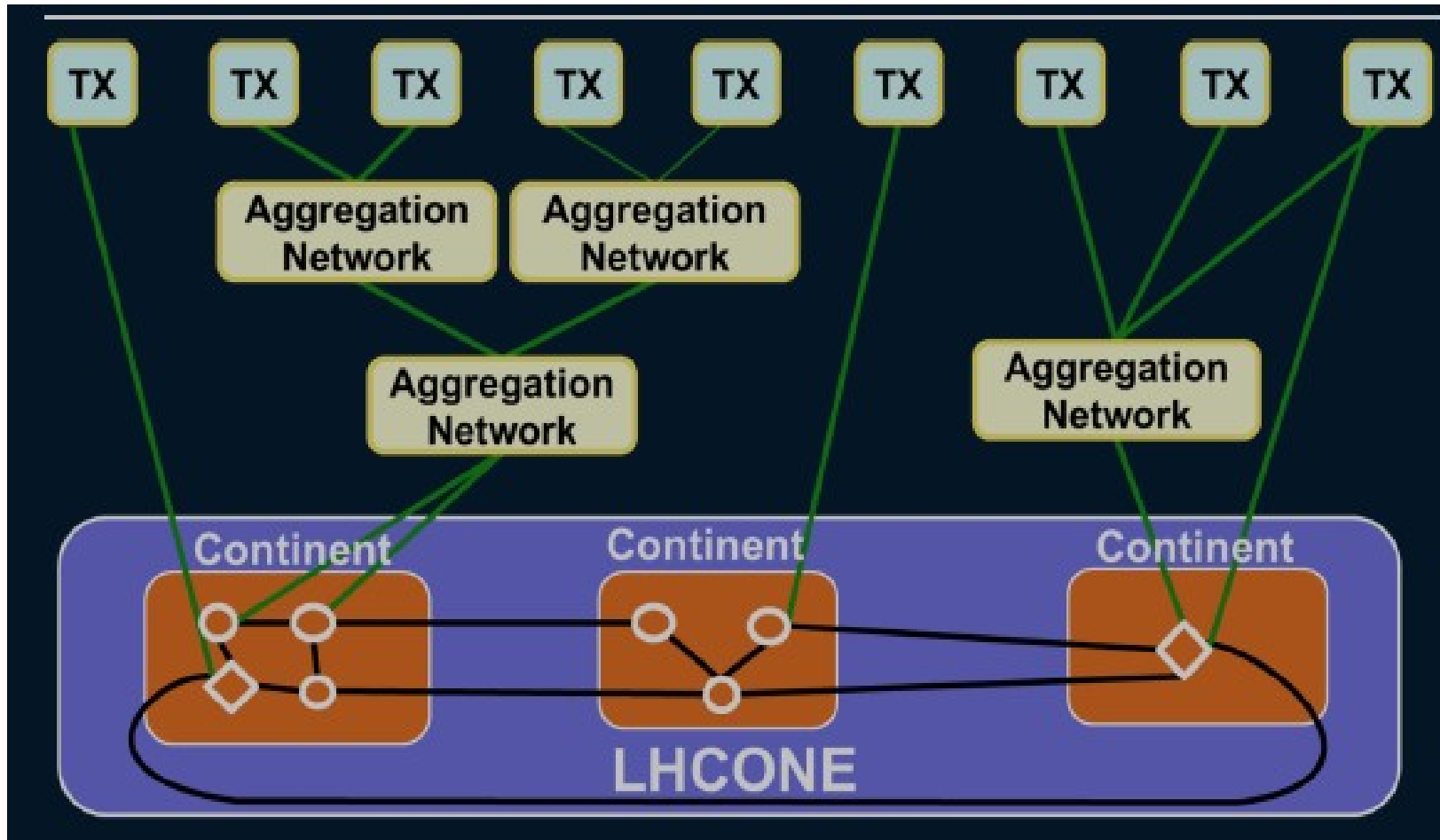


Aim of LHC ONE project is

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



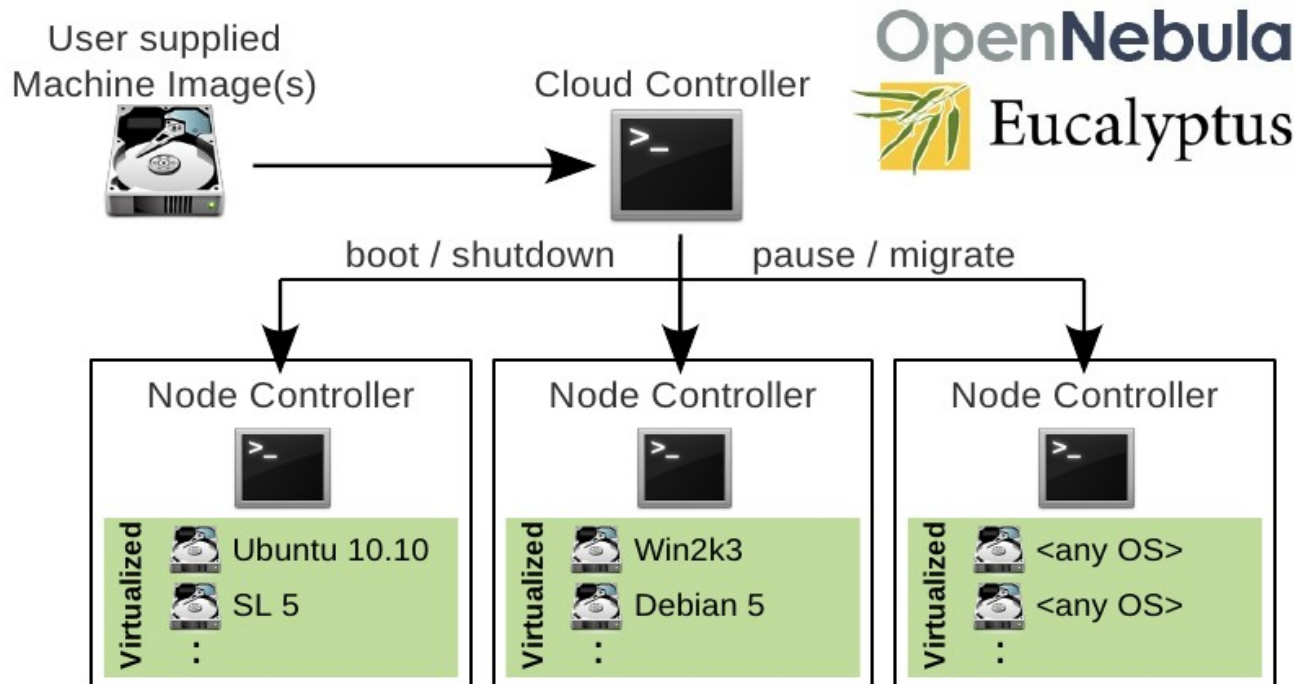
Schematic layout of LHCONE network infrastructure

A dedicated HEP network infrastructure – what is the cost ?

Virtualisation

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)

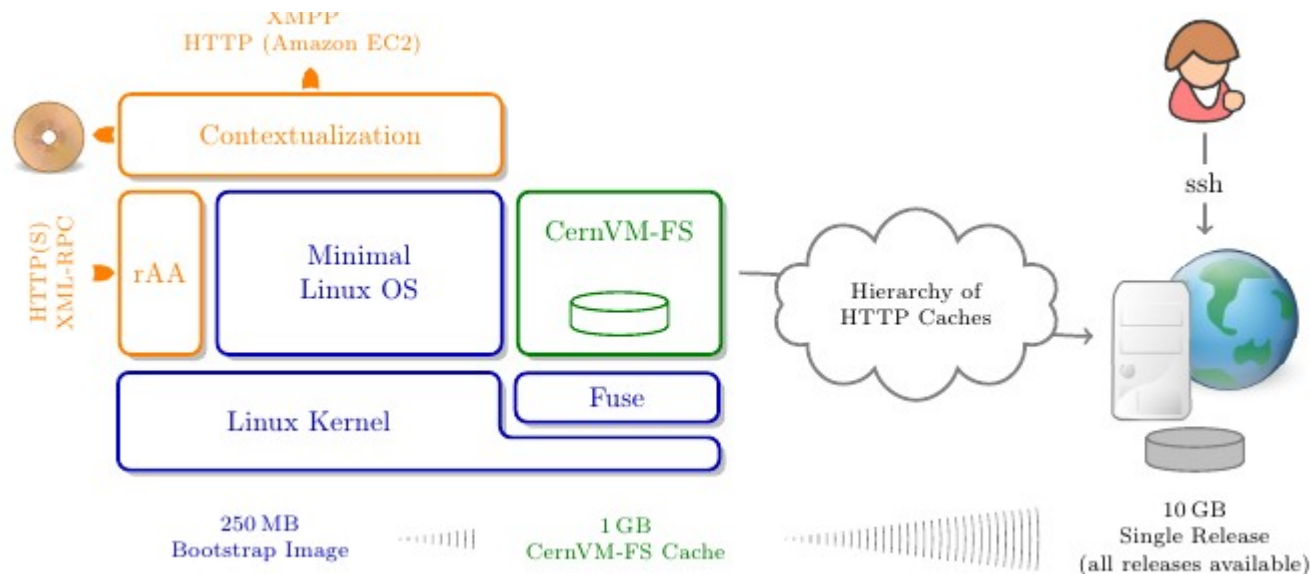


CernVM & CernVM-FS

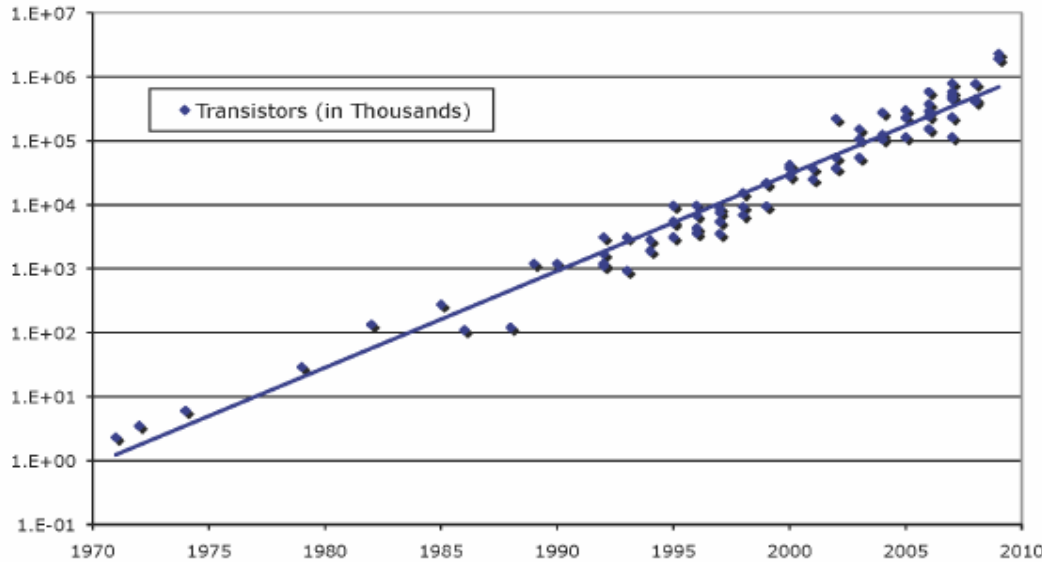
CernVM is a virtual machine (“Virtual Software Appliance”) based on Scientific Linux with CERN software environment, runs on    



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.



Most important challenge: Hardware evolves

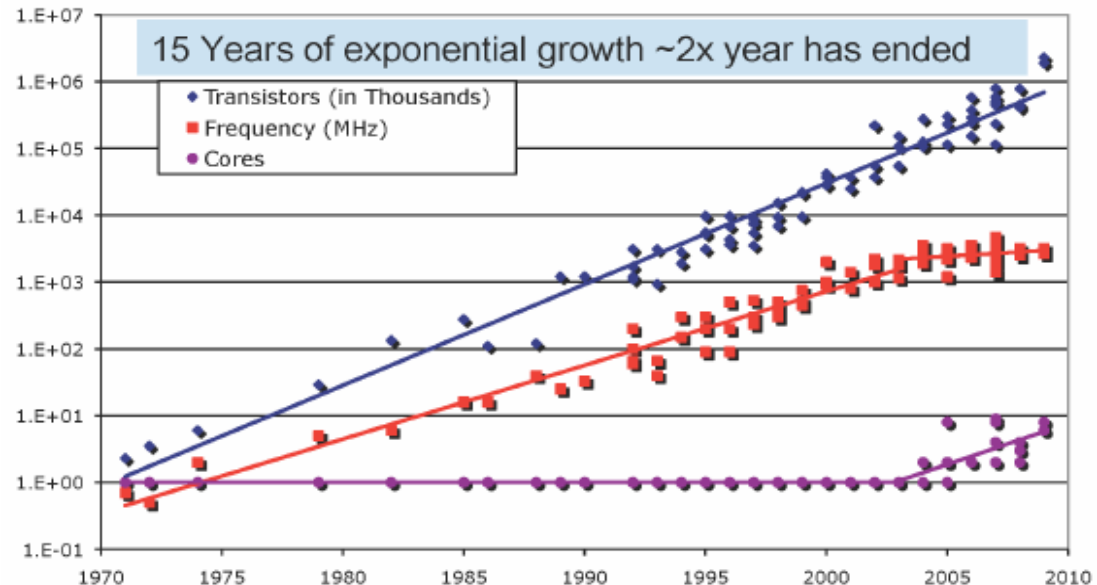


of transistors per chip keeps growing exponentially

but clock rates saturate

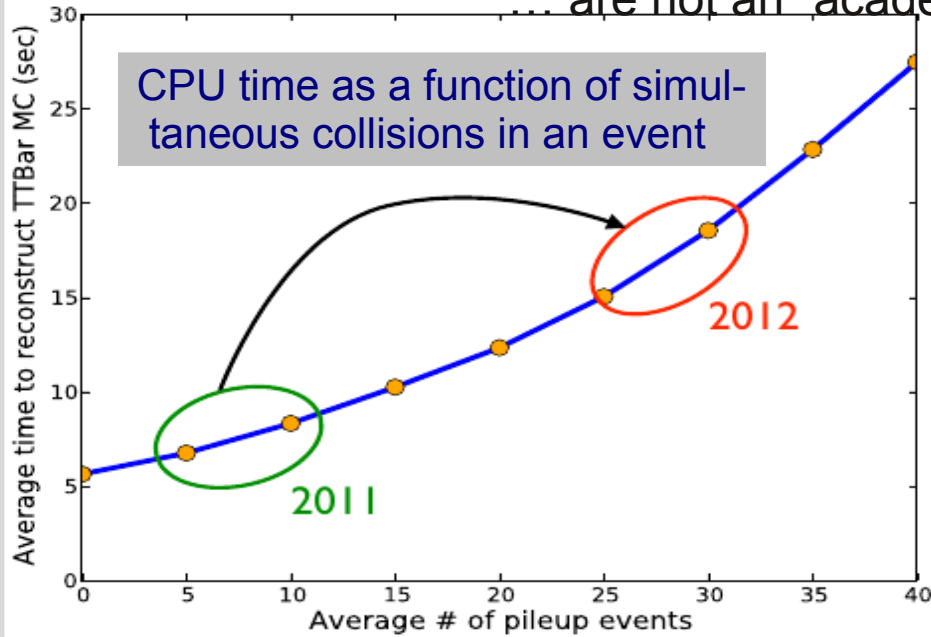
**instead:
more CPU cores on a chip !**

→ **affects the way we
(have to) write software**



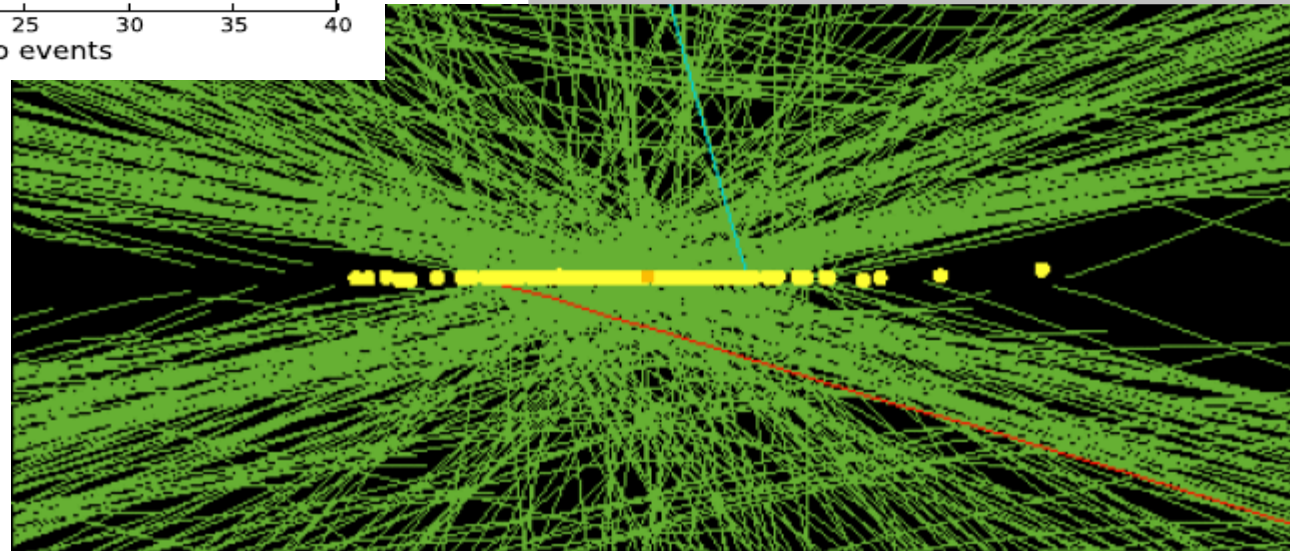
Software challenges ...

... are not an “academic” problem, already hit us:



Such a “**high pile-up**” environment will become “normal” with rising LHC luminosity

Event recorded by CMS with 78 simultaneously colliding protons in a bunch



Need (to learn) new ways of writing code

- exploit vector units on CPU cores: “data oriented programming”
(to replace “object oriented programming in time-critical algorithms”
exploit auto-vectorisation in new gcc to use
Single Instruction Multiple Data (MMX, SSE(2), AVX,...) features of modern CPUs
- memory / core remains roughly constant
- multi-threading and multi-cores
“one job per core”, i.e. “trivial parallelism” does not work any more
requires clever workflows and thread control
- use of Graphical Processing Units (GPUs) for time-critical parts
development environments get more user-friendly, will have to learn this!

Hint: *there will be a course on Parallel Software Development at the “1st KSETA topical days” in February 2013 by Benedikt Hegner (CERN) et al.
(from whom I've stolen some of the material just shown)*

- **existing LHC Grid provides excellent performance** to particle physics
- situation was rather comfortable in the past:
“sit and wait” for better hardware
this will change: have to adapt to new hardware environments
??? Will PhD students be able to master **parallel programming & physics & detector service tasks** **???**
→ need new, parallel-capable frame works
- **resource demands of experiments will rise:**
 - detector upgrades with more readout channels
 - higher trigger rates
 - higher luminosity and more crowded events**will count on an evolutionary ansatz, as we did so far**
- **new, open and standard ways for distributing computing emerged:**
various variants of **cloud computing** to be integrated in computing models
- **network bandwidth will continue rising,**
decouples storage from CPU,
enables clever, transregional data management