

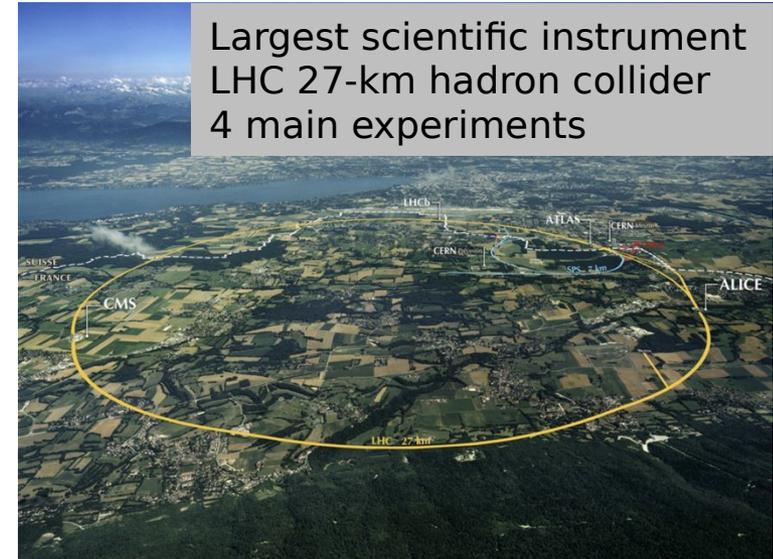


**Is the Big Data Experience at  
CERN useful for other  
applications?**

# CERN



- **Founded in 1954 (Science for Peace)**
  - Research
    - develop, build, run unique 'frontier' facilities for European High Energy physics
  - Provide an environment for training physicists and engineers
  - Facilitate and actively pursue Technology Transfer
  - Foster international collaboration
- **20 member states**
  - Austria, Belgium, Bulgaria, Czech republic, Denmark, Finland, France, Germany, Greece, Hungary, Italy, Netherlands, Norway, Poland, Portugal, Slovak republic, Spain, Sweden, Switzerland, United Kingdom
  - Candidate for accessions: Romania
- **Observer states**
  - European Commission, India, Japan, Russian Federation, Turkey, UNESCO and US
- **Pre-stage to membership**
  - Israel and Serbia
- **People**
  - ~ 2,300 staff
  - ~ 1,000 other paid personnel
  - 11,000 users



Largest scientific instrument  
LHC 27-km hadron collider  
4 main experiments

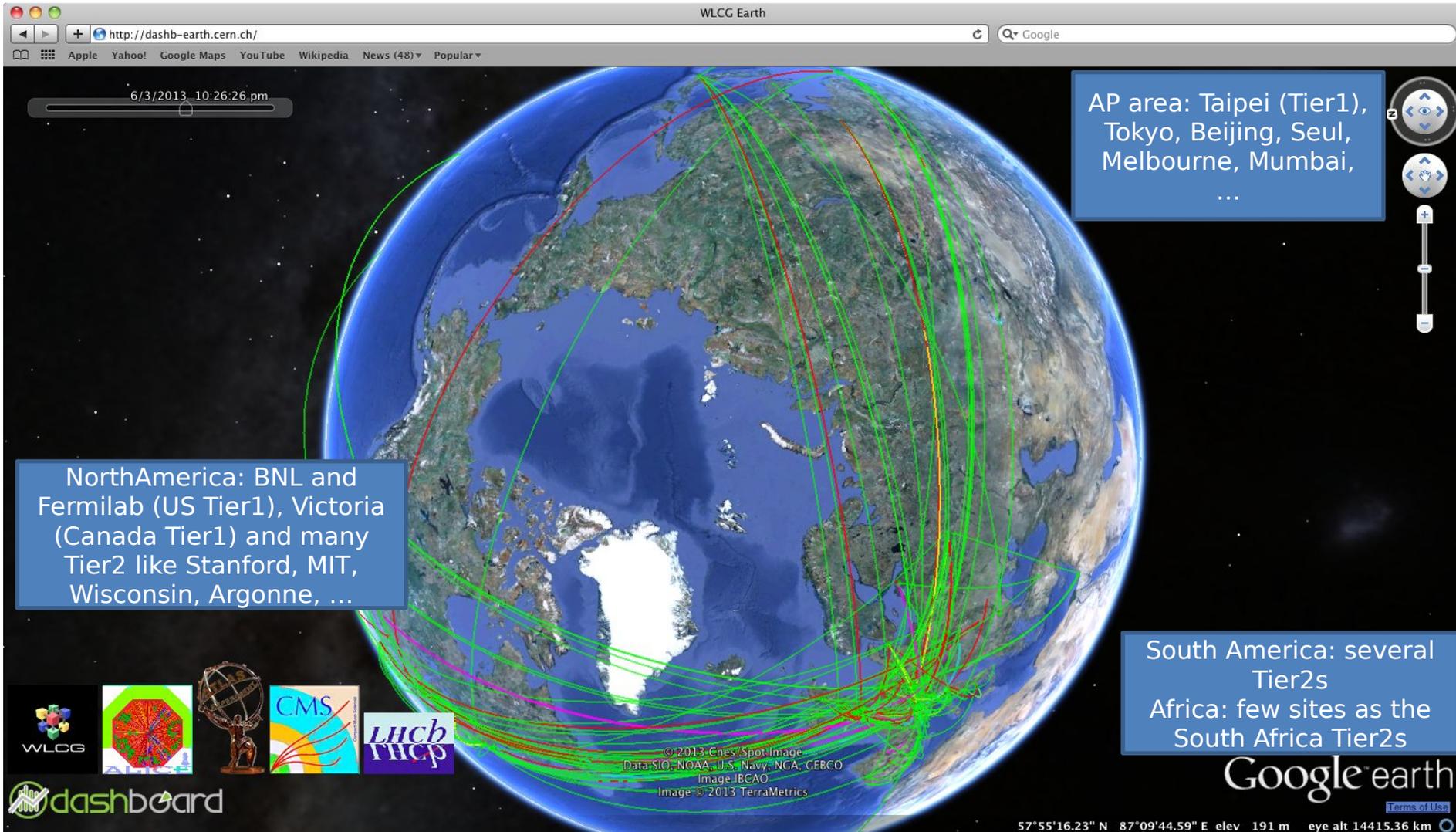


ATLAS  
~ 100 million electronic channels  
~ 3 000 physicists

# WLCG: Worldwide LHC Computing Grid collaboration)



The WLCG collaboration was set up in 2002 for the LHC physics programme at CERN. It links up national and international grid infrastructures (200+ sites). CERN provides ~20% of the resources (CPU and storage)



CERN: 90,000 cores and 100 PB data; growing to cope with the evolution of LHC computing needs

# CERN Computer Centre

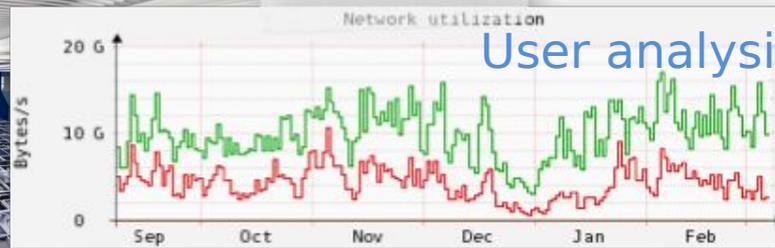
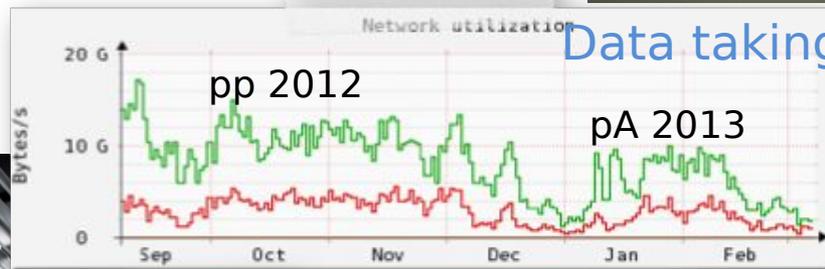


## CERN computer centre:

- Built in the 70s on the CERN site (Meyrin-Geneva)
- ~3000 m<sup>2</sup> (3 main machine rooms)
- 3.5 MW for equipment
- Est. PUE ~ 1.6

## New extension:

- Located at Wigner (Budapest)
- ~1000 m<sup>2</sup>
- 2.7 MW for equipment
- 2x100Gb links (21 and 24 ms RTT)



Number of 10GB NICs	2,794
Number of 1GB NICs	18,444
Number of cores	90,408
Number of disks	76,644
Number of memory modules	64,116
Number of processors	17,520
Number of servers	10,233
Total disk space (TiB)	116,952.07

# (Innovation in) Computing in/from High-Energy Physics



- Demanding science → Demanding computing

Power usage



- Innovation
  - Web invention
  - Grid computing (LHC Computing Grid)



# ITU conference (2006)

The problem:  
Assign frequencies for  
digital radio and  
television (international  
treaty)

Critical point:  
Need on dependability:  
verify (iteratively)  
the compatibility  
between radio stations

Solution:  
Use the EGEE grid + a  
system used for CERN LHC  
to increase the reliability  
of the Grid

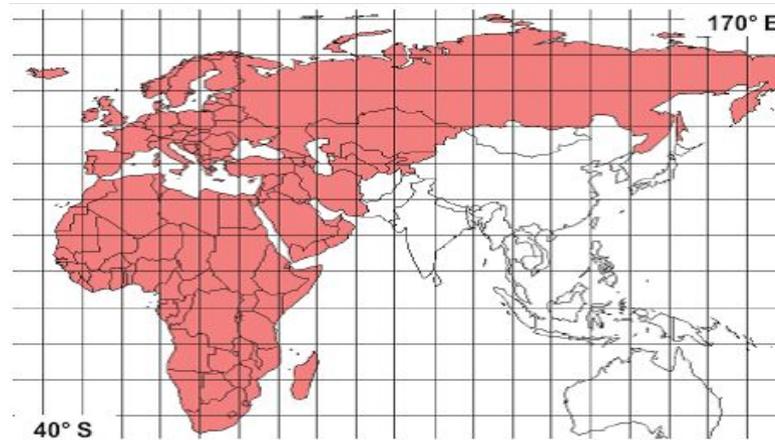
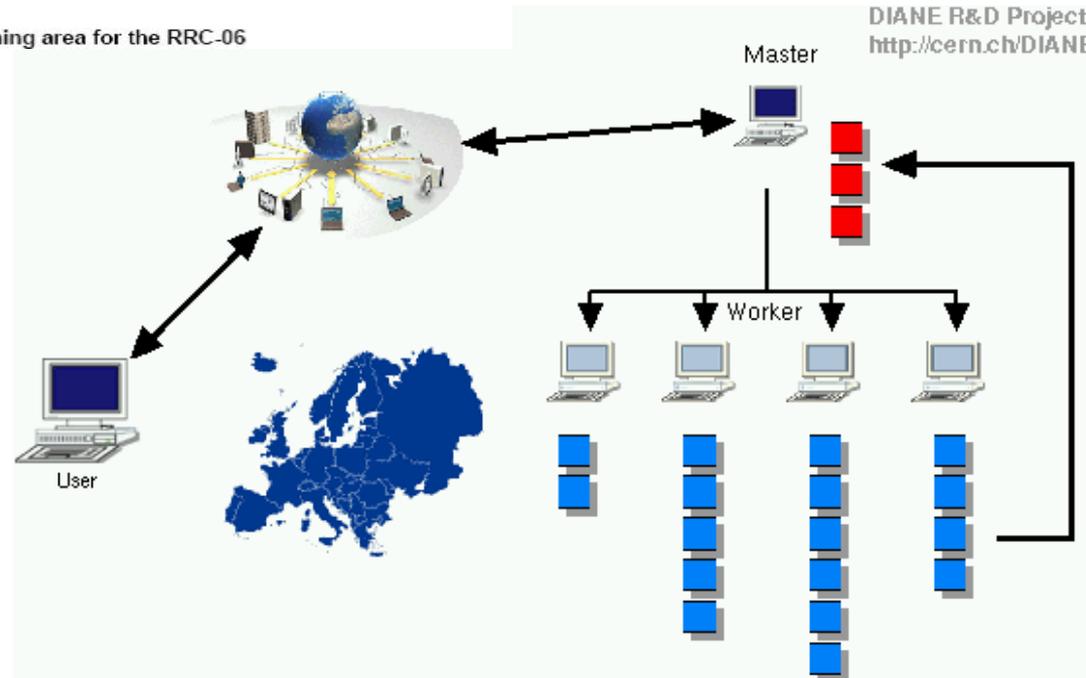
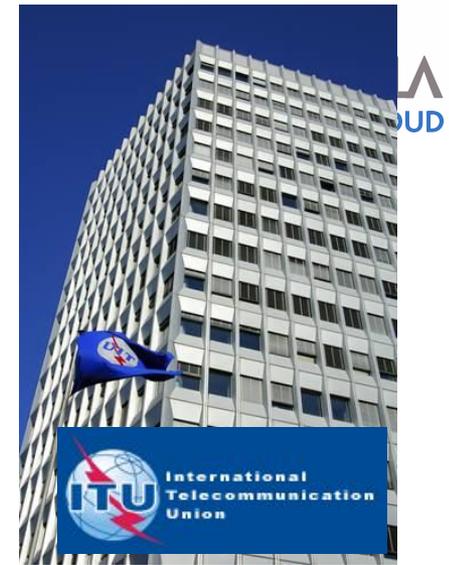
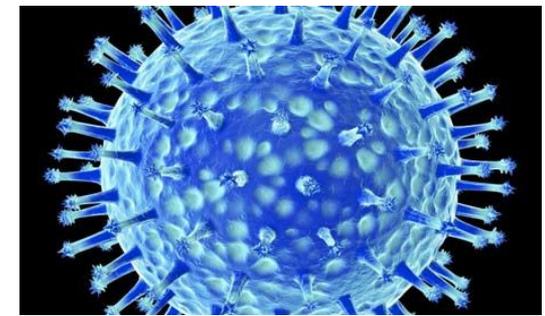


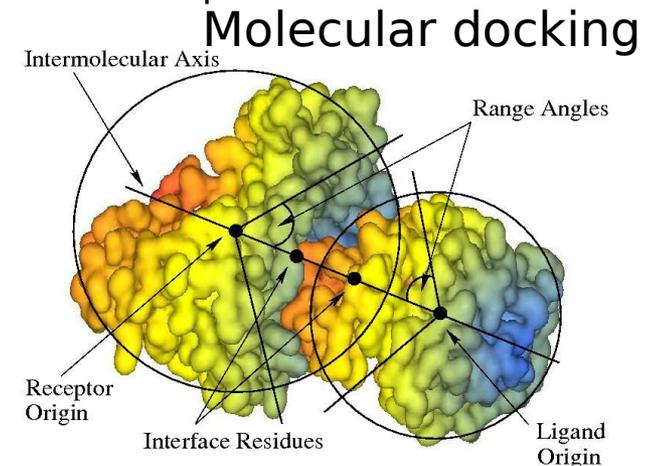
Figure 1  
The extent of the planning area for the RRC-06



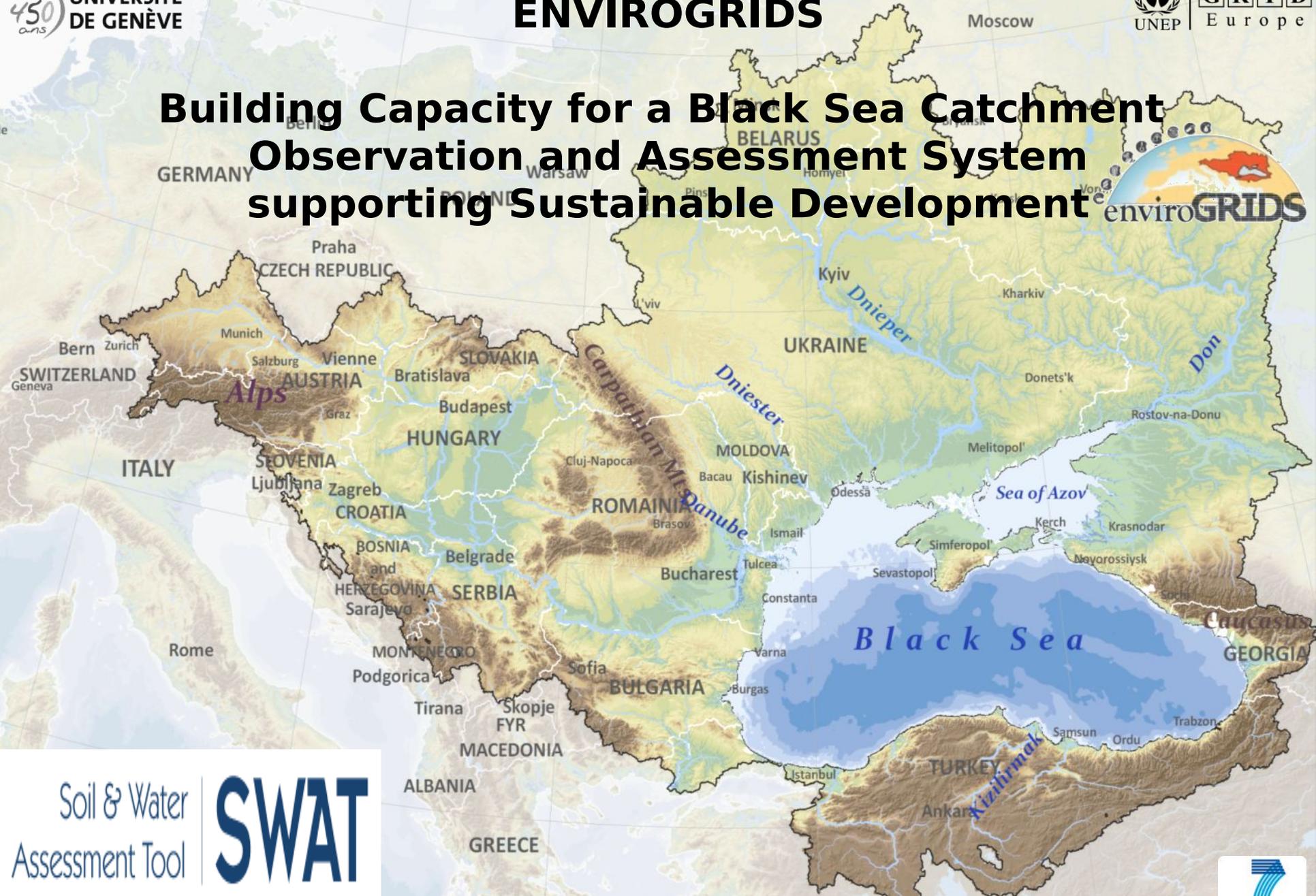
# Bird Flu



- Basic idea:
  - Compute how a given chemical interacts with a protein (e.g. belonging to a virus)
  - High affinity means the chemical is a potential drug against the virus
- In silico (i.e. use your PC):
  - Scan millions of chemicals ( $\sim 10^3$  s per chemical-protein pair)
    - With 1,000 PCs, 1 docking per second
  - Good candidate given to biologist (verification longer -and more expensive- than in silico docking)
  - In practice, you enrich the initial sample saving time
    - (and money)
    - Essential to fight to pandemic (H5N1) or to fight neglected diseases (like Malaria)
- WISDOM collaboration
  - Malaria
  - H5N1 (Bird Flu)



## Building Capacity for a Black Sea Catchment Observation and Assessment System supporting Sustainable Development



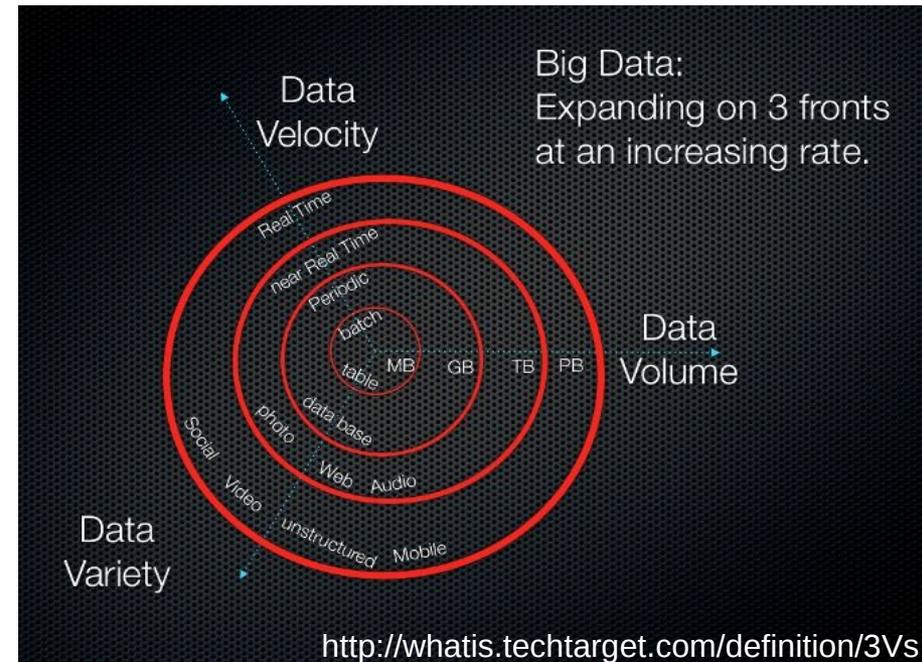
Soil & Water Assessment Tool **SWAT**

Courtesy Dr A. Lehmann (UniGenève and UNEP)

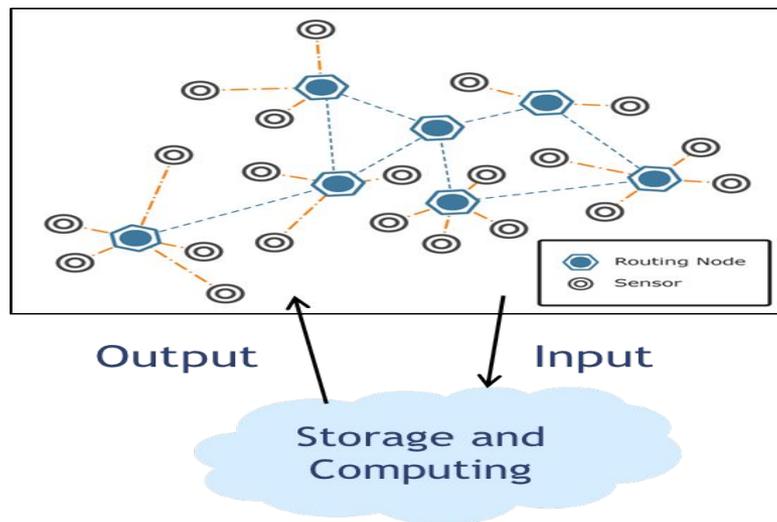
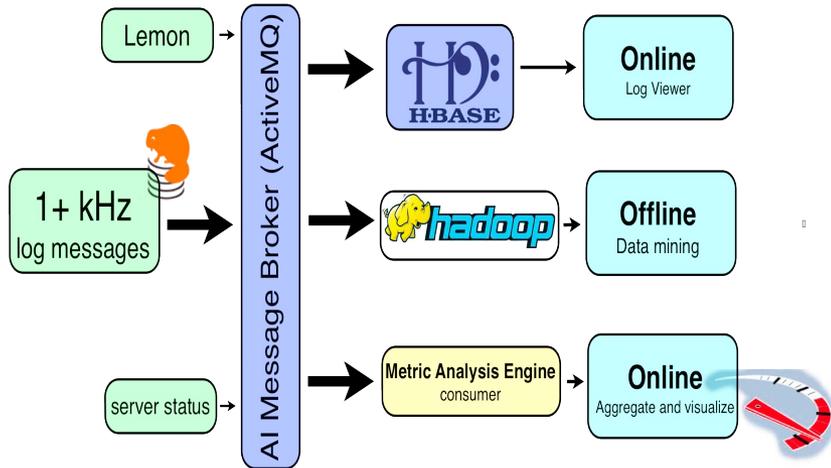
# Big Data is not only in HEP



- Big data means different things to different people...
  - Volume Variety Velocity
  - HEP “Volume” is **big**; Velocity ranges from real-time to batch; Variety is relatively **low**
    - Although the ‘queries’ we run on our data are really complex (compared to a standard Google search)
    - And each user often creates new custom data (also distributed!) from large distributed input (e.g. event selection)
      - Cfr. with Netflix delivering a movie files to their users
      - 
      - **Cfr. Smart Cities/Traffic analysis**
- Volume of LHC is not (yet) what a
  - single area can ‘produce’
  - $o(100)$  of sensors giving traffic info
  - in relative frequent samples (10 s)
  - → 100 GB/year
- Variety?
  - Much higher!
  - Federating data from different admin
  - (e.g. highways, phone companies,
  - police dept, city administration,
  - sensor embedded in cars, ...)



# How can we model it?



- **Control system of for data management services**

- O(1000) 'sensor'
- 100 GB/day

- **We need:**

- Long-term storage ("batch")
- Auditing, trend analysis, ...
- Troubleshooting ("fast")
- Optimisation and system evolution

- **How it is implemented?**

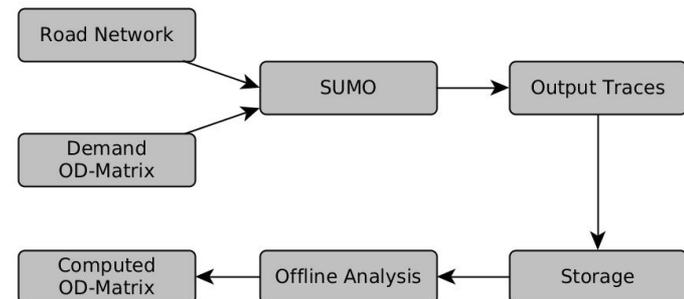
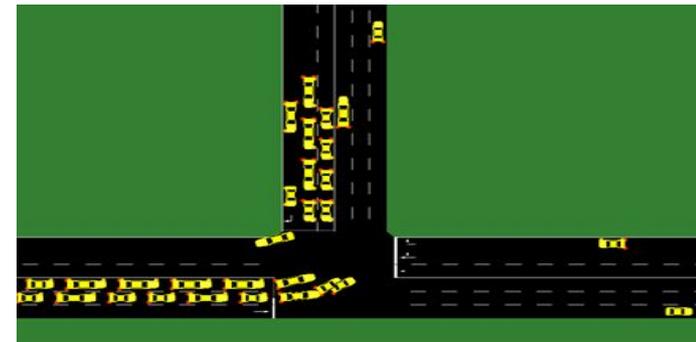
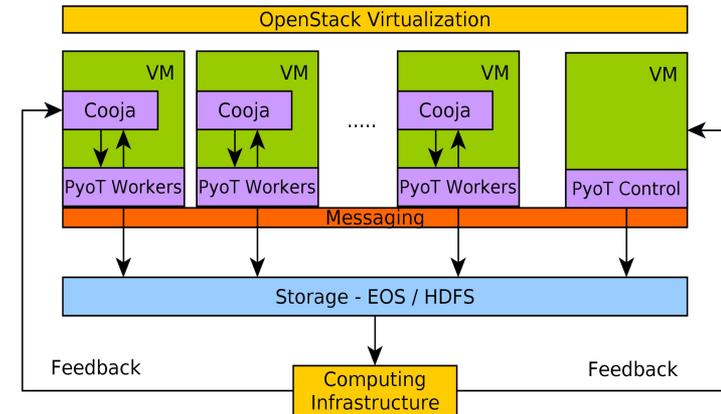
- Overlay of the existing infrastructure
- Sensors + message bus
- Flexibility + 'Real time'
- Publishers decide what to publish
- Multiple consumers
- One consumer is the long-term aggregation

- **Can it be reused as a prototype?**

# Exercise



- National Inter-University Consortium of Telecommunications (CNIT) and Scuola Superiore S. Anna Pisa (SSSA):
  - P. Pagano and A. Azzara` (internship at CERN)
  - CNIT/SSSA is active in a number of projects notably: URBELOG (URBan Electronic LOGistics) with Telecom Italia prime contractor, Turin as pilot site); inter-modal Intelligent Transport System (Port of Livorno)
- CERN IT DSS
  - CERN OpenLab internship
- Full simulation
  - Wireless sensors (realistic data collection including network losses)
  - Traffic flow (simple mesh using SUMO)
  - Data aggregated on a storage backend



# Our contribution



- Scalable data platforms for data analysis

## HADOOP (HDFS) use

- Open-source implementation of frameworks for reliable, scalable, distributed computing and data storage
- Used by us for Monitoring
- Evaluation of the product

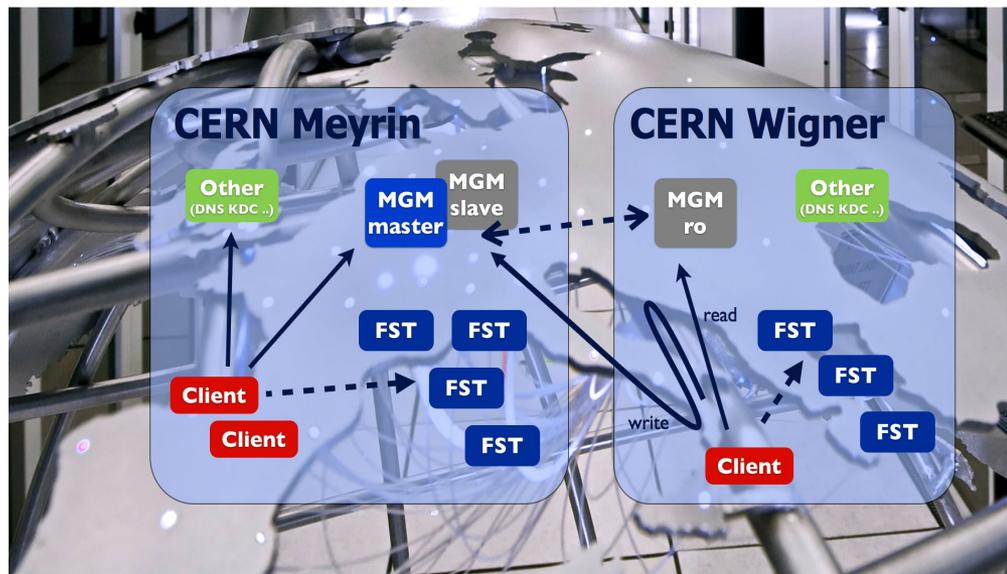
EOS

## – EOS use

- Developed at CERN; used for LHC physics analysis
- Concurrent usage
- Scalable operations
- Scalable performances

## – Common Key features

- Data aggregation
- Reference data
- Multiple users/usages
- Short-term monitoring and offline trend analysis
- Part of our new services:





# CERN openlab in a nutshell

A science – industry partnership to drive R&D and innovation with over a decade of success

- Evaluate state-of-the-art technologies in a challenging environment and improve them
- Test in a research environment today what will be used in many business sectors tomorrow
- Train next generation of engineers/employees
- Disseminate results and outreach to new audiences



## PARTNERS



## ASSOCIATE



# A European cloud computing partnership: big science teams up with big business



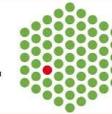
## Strategic Plan

- ▢ Establish multi-tenant, multi-provider cloud infrastructure
- ▢ Identify and adopt policies for trust, security and privacy
- ▢ Create governance structure
- ▢ Define funding schemes



To support the computing capacity needs for the ATLAS experiment

EMBL



Setting up a new service to simplify analysis of large genomes, for a deeper insight into evolution and biodiversity



To create an Earth Observation platform, focusing on earthquake and volcano research



# Q&A



## Big Data kan in sport het verschil maken tussen winst en verlies

Door **SANDRA BOUTEN** | 10 OKTOBER, 2013 13:25 |



Xavier Cortada (with the participation of physicist Pete Markowitz), "The search of the Higgs boson:  $H \rightarrow ZZ$ ", digital art, 2013.

