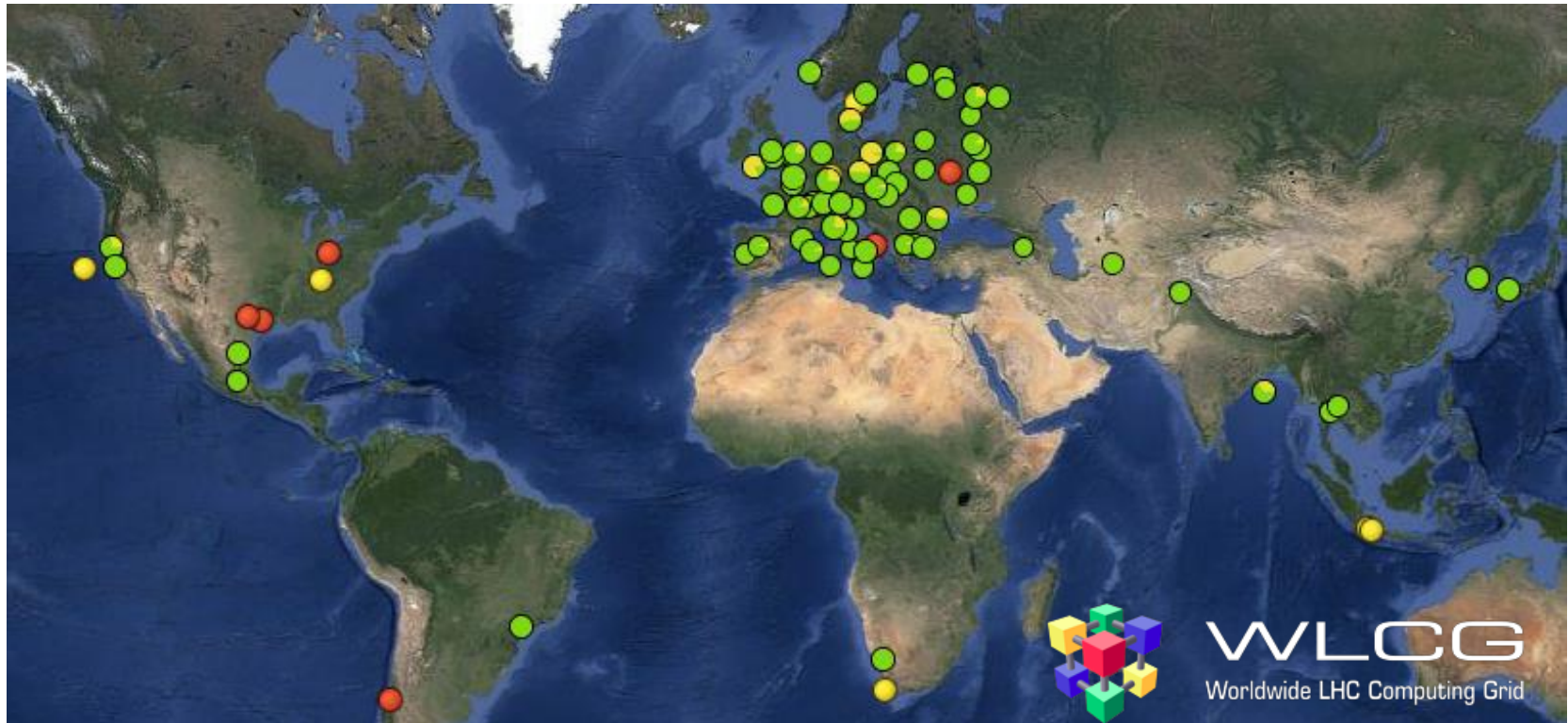A Large Ion Collider Experiment
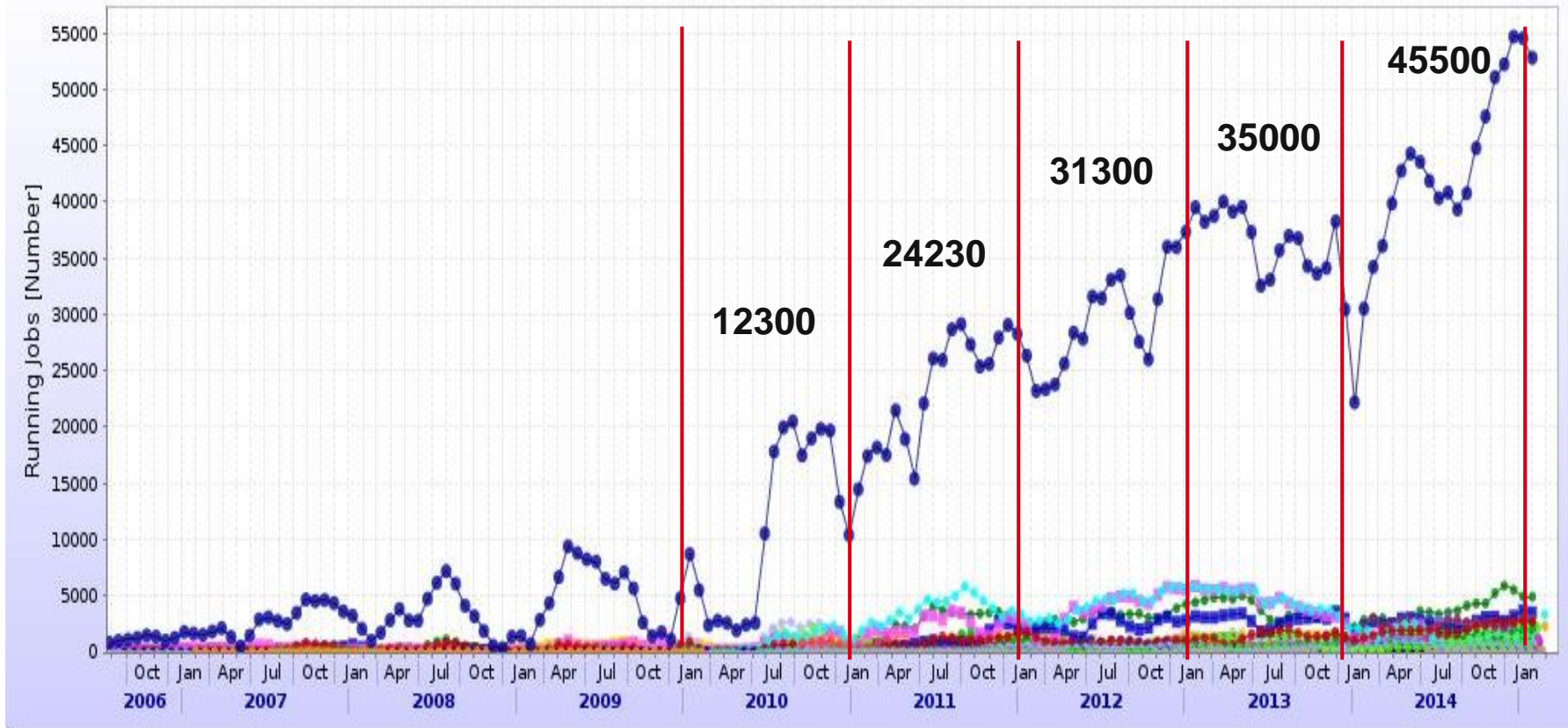
# ALICE Computing in Run 2+

P. Buncic

# Computing Model in Run 1&2



- Distributed grid computing model on top of WLCG services
  - Faithful implementation of the original grid ideas
  - Any job can run anywhere and access data from any place
  - Scheduling takes care of optimizations and brings "computing to data"
  - Every file and its replicas are accounted in the file catalogue

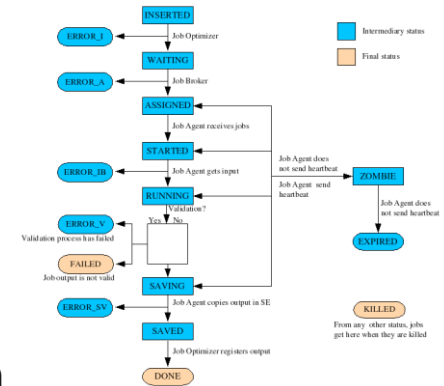# Grid evolution from ALICE perspective



**Year on year increase**    +97%    +30%    +12%    +30%

# **Distributed computing**

- ALICE has used Grid for the distributed production of Monte Carlo data, reconstruction and analysis
  - Elaborated and well debugged job state machine
  - Used for execution of all workflows in ALICE
  - File Catalog accounting for over 10^9 files
- So far more than 280 million jobs have been successfully run worldwide from the central Task Queue (TQ), resulting in the currently active volume of 50 PB of data
  - System scaled from 500 to 80000 concurrently running jobs
- The job handling model based on pilot jobs and late binding to the real job has proven to be the solution for the large scale computing problems of the LHC experiments
  - Central Task Queue governing priorities

- This concept will be retained for Run 3

# Workflows & Efficiencies
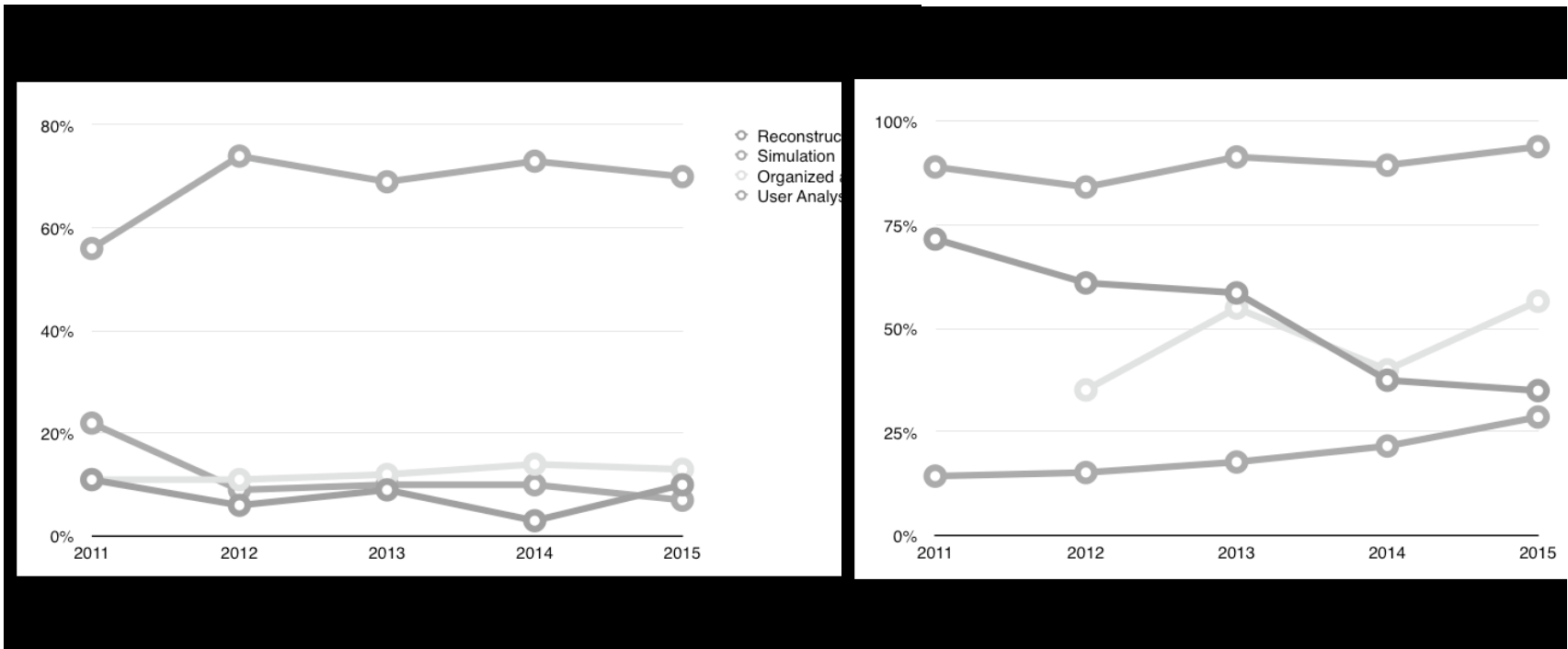


- All types of jobs (workflows) are run on the Grid
  - But not all of them are equally efficient
  - We run mostly simulation (over 70% of CPU) which happens to be also very efficient (over 90% CPU/Wall time)
  - Analysis is the least efficient but it takes only 20% of overall CPU budget
- We cannot afford to be inefficient in Run 3!

# ALICE Upgrade

- Event rate will grow by a factor 100 in Run 3

- Grid capacity will continue to grow within the constraints of a flat budget
    - Available resources must be used with maximum efficiency

- We have to reduce data volume early and by a large factor
    - Requires considerable computing capacity @ P2

- Use h/w accelerators to speed up the computation and reduce the cost
    - Requires new and flexible s/w framework that can adapt itself to different environments

# O2: Big, heterogeneous facility

+ 463 FPGAs
  - Detector readout and fast cluster finder
+ 100'000 CPU cores
  - To compress 1.1 TB/s data stream by overall factor 14
+ 5000 GPUs
  - To speed up the reconstruction
  - 3 CPU[1] + 1 GPU[2] = 28 CPUs
+ 50 PB of disk
  - To buy us an extra time and allow more precise calibration

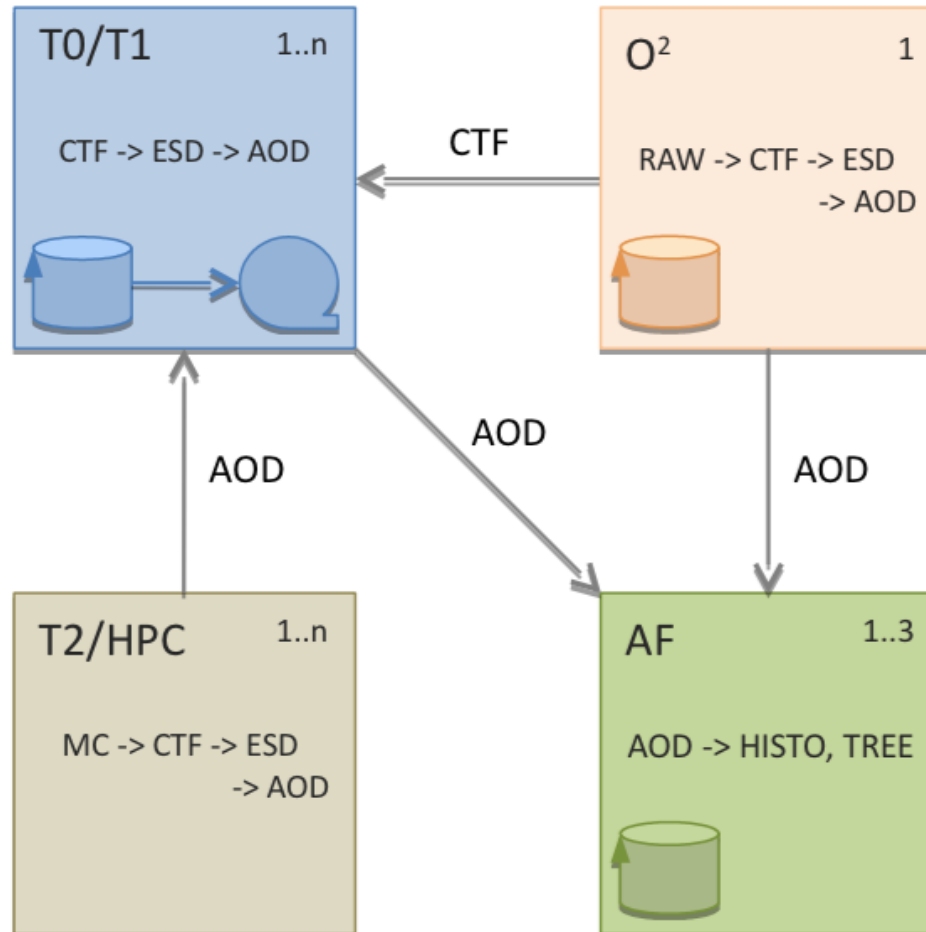------------------------------------------------------------------

= Considerable computing capacity that will be used for Online and Offline tasks
  - ✧ Identical s/w should work in Online and Offline environments

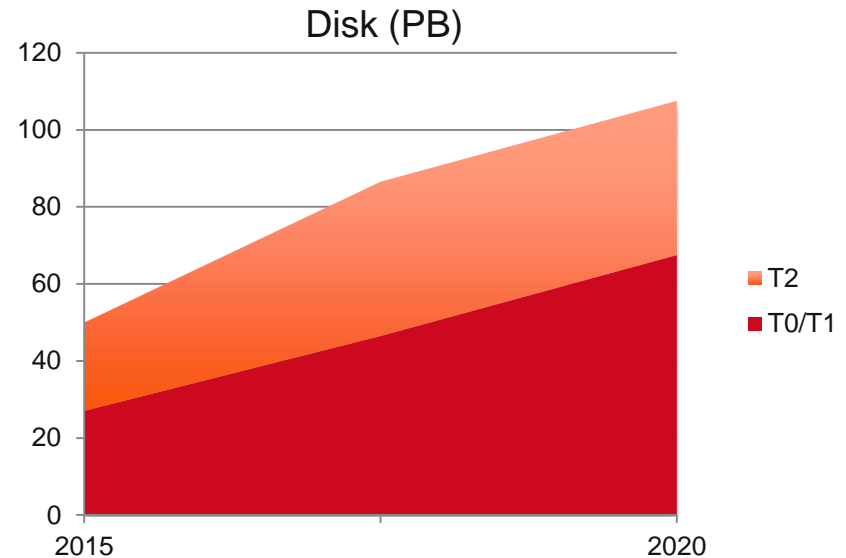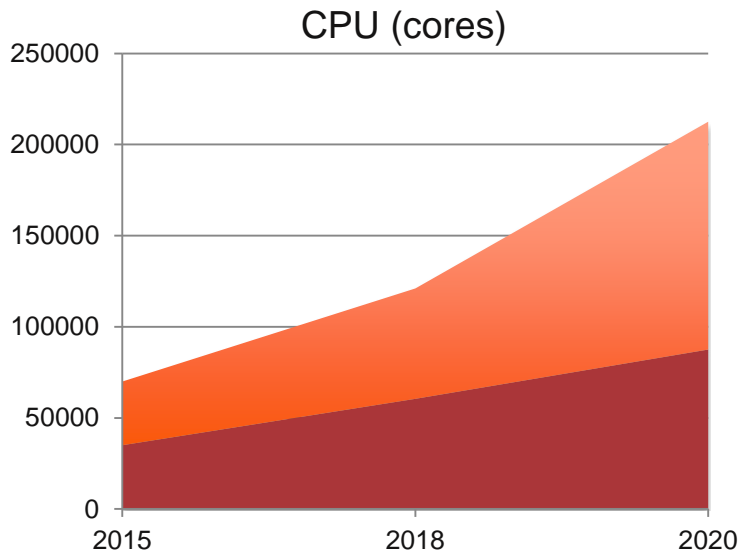[1] Intel Sandy Bridge, 2GHz
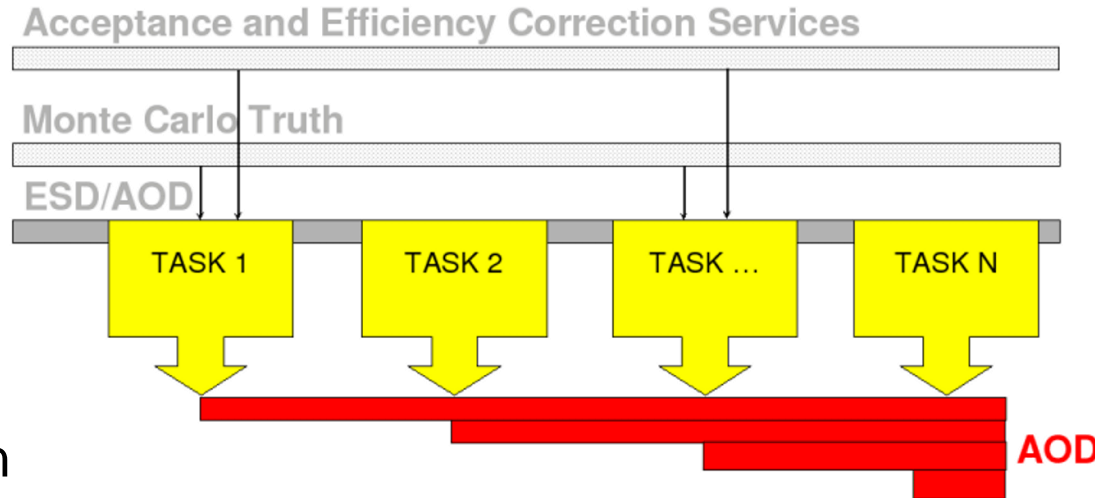[2] AMD S9000

# Specializing Grid Tiers

# Grid vs O2



- Expecting Grid resources (CPU, storage) to grow at 20% per year rate (x.2.5 over 5 years)
  - Large number of disk will be used by Run 1 and Run 2 data
- Since T2s will be used almost exclusively for simulation jobs (no input) and resulting AODs will be exported to T1s/AFs, we expect to significantly lower needs for storage on T2s
  - Use available funding to buy more CPUs
- We expect to use 20Gbps share of network connectivity between CERN and T1s
- Given the network capacity and storage available at various Tiers as well the resources at P2 we expect 2/3 raw data processing to be done at P2 and 1/3 at T1s

# Analysis Facilities

Acceptance and Efficiency Correction Services

Monte Carlo Truth

ESD/AOD

TASK 1   TASK 2   TASK ...   TASK N

AOD

- Motivation
  - Analysis is the least efficient of all workloads that we run on the Grid
  - I/O bound in spite of attempts to make it more efficient by using the analysis trains
  - Increased data volume will only magnify the problem
- Solution
  - Collect AODs on a few dedicated sites that are capable of locally processing quickly large data volume
  - Typically (a fraction of) HPC facility (20-30'000 cores) and 5-10 PB of disk on very performant file system
  - Run organized analysis on local data like we do today on the Grid

# Complexity management



- Virtually joining together the sites based on proximity (latency) and network capacity into Regional Data Clouds
- Each cloud/region provides reliable data management and sufficient processing capability
  - Dealing with handful of clouds/regions instead of the individual sites

# Summary

- With x100 increase of event rate at 1TB/s  data storage and management becomes very serious problem for new computing model
- We aim to minimize amount of data moving between Tiers and carry out most of processing on local datasets
- Analysis Facilities are new component in our Computing Model that are meant to improve the analysis efficiency
- In general, our strategy is to avoid problems and have some ideas how to tackle our computing problems but  there is a lot of room for improvements
    - Data management
    - Monitoring & Control
    - Automating (Grid) Operations
    - New s/w framework (multi process, message based, 0MQ, GPUs, vectorization, performance tuning… )