# Support for multiple virtual organizations in the Romanian LCG Federation

**M. Dulea, S. Constantinescu, M. Ciubancan**

**Department of Computational Physics and Information Technologies**

**'Horia Hulubei' National Institute for R&D in Physics**

**and Nuclear Engineering (IFIN-HH),**

**Magurele, Romania**

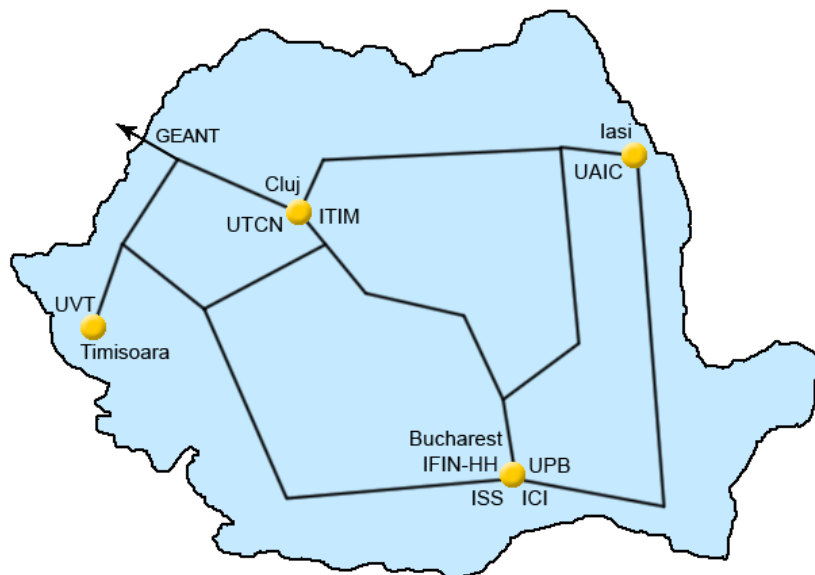# ROMANIAN GRID INFRASTRUCTURE

**Resource centres**:

- 10 active sites hosted by 6 institutions
- Total number of cores: 6200
- Total disk capacity: 1.8 PB

**Certification authority:**

- RomanianGRID CA, operated by ROSA

**Network infrastructure** provided by

- NREN: RoEduNet (backbone min. 10 Gbps)



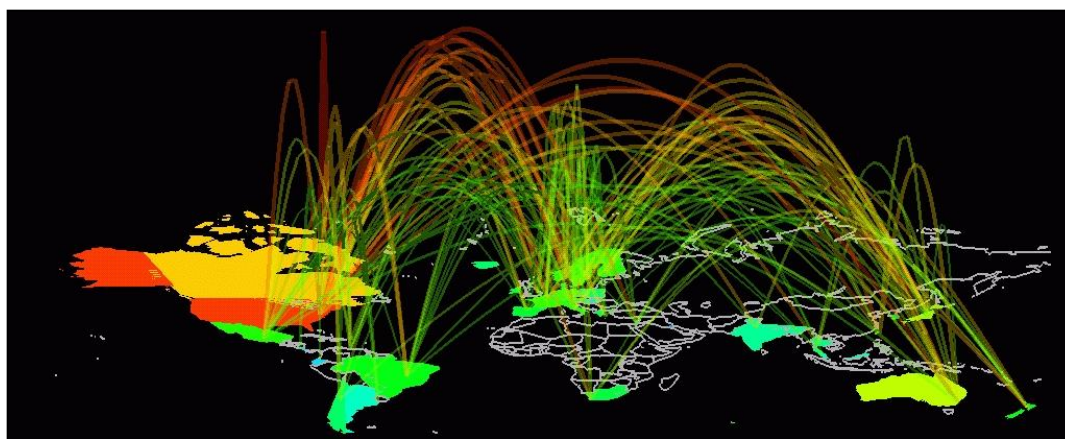**Research communities (VOs) supported** (last 12 m.):

| VO | Collaboration supported | # of sites | % of total |
|---|---|---|---|
| `alice` | ALICE experiment - LHC | 3 | 43.3% |
| `atlas` | ATLAS experiment - LHC | 4 | 50.6% |
| `lhcb` | LHCb experiment - LHC | 3 | 4.7% |
| `envirogrids .vo.eu- egee.org` | FP7 enviroGRIDS project | 1 | 0.6% |
| `gridifin` | RO physics & related areas | 1 | 0.2% |
| `see` | Multidisciplinary, SEE | 2 | 0.3% |
| `seegrid` | FP6 SEE-GRID project | 1 | 0.2% |
| `hone` | H1 experiment - DESY | 1 | 0.1% |

**More than 98% of the national grid production is dedicated to the community of elementary particle physicists involved in LHC research.**

The 5 institutions that contribute to the grid support of the 3 LHC experiments, within the **Worldwide LHC Computing Grid** collaboration (WLCG), are partners in the **Romanian Tier-2 Federation** (**RO-LCG**).

# RO-LCG WITHIN THE WLCG COLLABORATION

## WHY LHC COMPUTING GRID ?

The discovery of the Higgs-like particle "*has only been possible because of the extraordinary achievements of the experiments, infrastructure, and the grid computing*" (Rolf Heuer, CERN, 4.07.2012)





## THE ROLE OF RO-LCG

As a Tier-2 centre, the main tasks of RO-LCG consist of:

❖ providing disk storage resources and computing power for the Monte Carlo simulations and data analysis needed by the experiments;

❖ fulfill the resource pledge and the service level agreement (SLA) stipulated in the WLCG Memorandum of Understanding (MoU) (minimal annual availability of 95%, and minimal connectivity of 10 Gbps with the GEANT Education and Research network)

**Resources provided**:

- 7 grid centres
- 4800 cores
- ~ 1.8 Petabytes storage capacity

# RO-LCG TOPOLOGY

**Resource centres**:

- IFIN-HH:

    NIHAM (`alice`)

    RO-02-NIPNE (`atlas`)

    RO-07-NIPNE (`alice`, `atlas`, `lhcb`)

    RO-11-NIPNE (`lhcb`)

- ISS = Institute of Space Science:

    RO-13-ISS (`alice`)

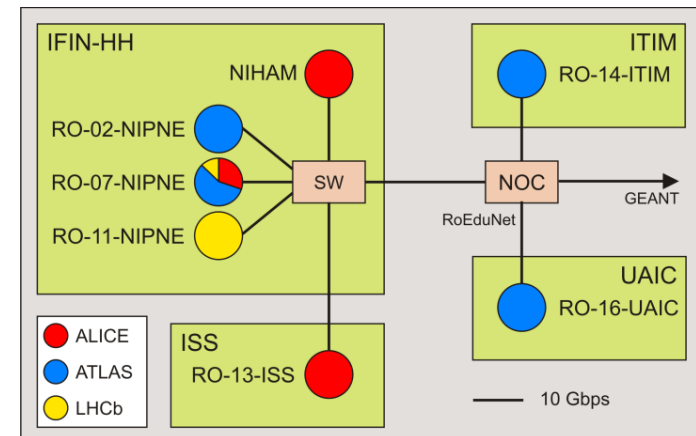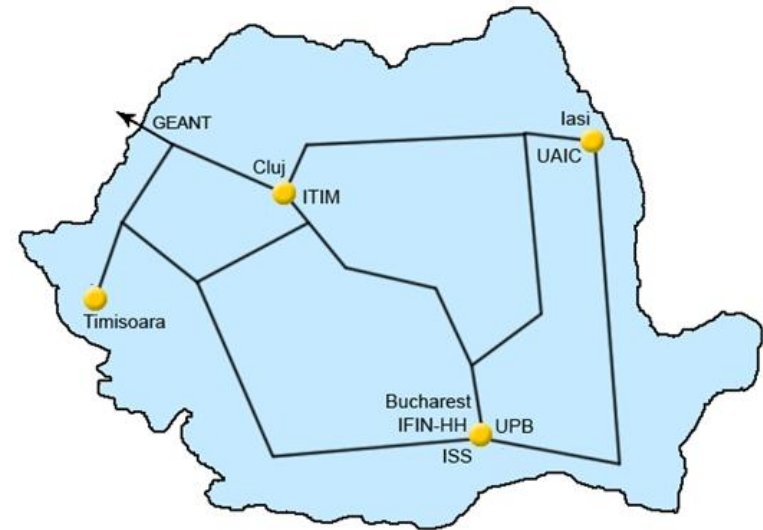- ITIM = Natl. Inst. for R&D in Isotopic and Molecular Technologies from Cluj-Napoca:

    RO-14-ITIM (`atlas`)

- UAIC = 'Alexandru Ioan Cuza' University of Iasi

    RO-16-UAIC (`atlas`)

NOC = RoEduNet's National Operation Centre

*Bucharest Operation Centre not represented (located between SW and NOC)*

## BASIC CONFIGURATION

Most of the sites share the following configuration

- Linux platform (mostly Scientific Linux),

- LCG middleware (migration from gLite 3.2 to recent EMI(-2) releases),

- CREAM-CE Computing Element (CE) servers

- PBS/TORQUE queuing system

- MAUI job scheduler

- DPM/SRM Storage Elements (SE).

Each site is dedicated to a single LHC experiment, with the exception of RO-07-NIPNE,

administrated by DFCTI, which supports three LHC VOs.

They will be presented by the other speakers.

# NETWORK TOPOLOGY @ DFCTI

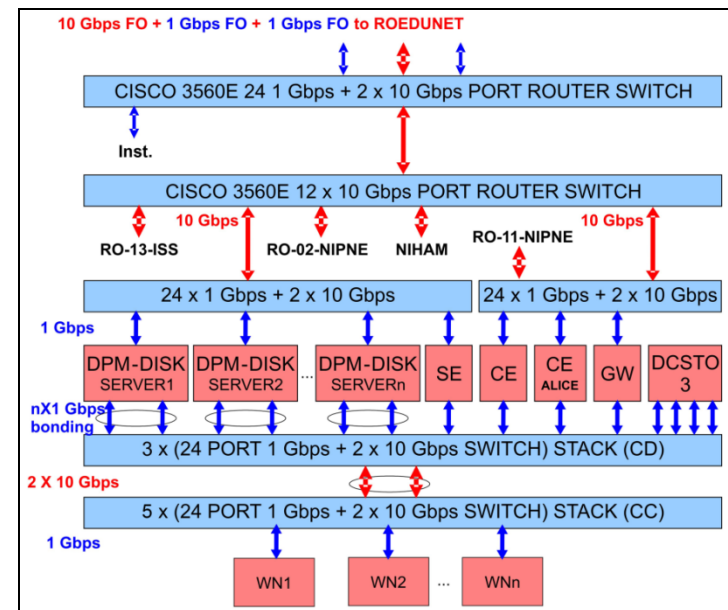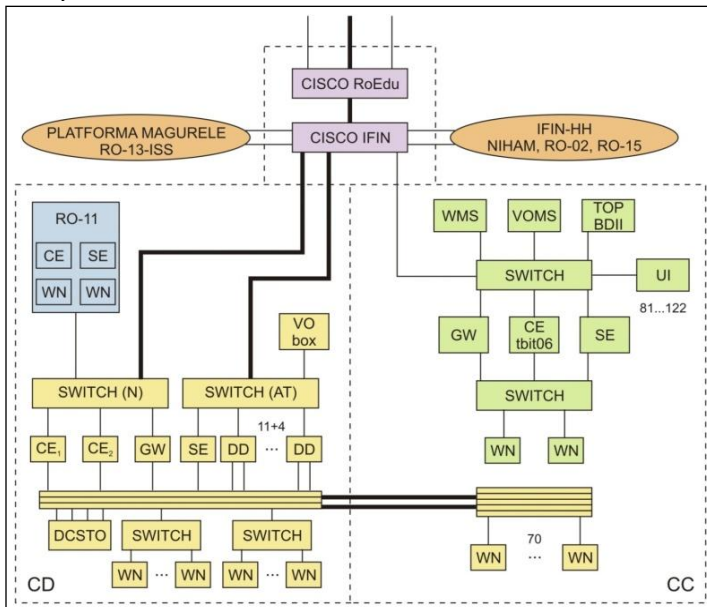DFCTI hosts 3 grid sites: RO-07-NIPNE ( brown); RO-11-NIPNE (blue); IFIN GRID (green, part of GriNFiC)

## RO-07

Job management: ensured by two CREAM-CE computing elements

Storage management: by a DPM/SRM head node that manages 15 DPM disk servers

Jobs are processed by more than 100 workernodes (WNs) with 4-32 cores each, and 2GB RAM per core. The WNs configuration also includes 4 GB virtual memory per core, minimum 400 GB disk scratch space (800 GB on 32-core nodes), and 1-2 Gbps network link.

The submission of alice jobs required the implementation of a gLite VO-box server as an user interface, that submits job agents to the CE and provides proxy management. To increase the availability of the description of the local resources and their status, a top BDII server was installed within the IFIN GRID cluster, which shares the same network.

"**Computing Centre"** (APC InfraStruXure™ )

**"Data Centre"**

**RO-07-NIPNE**:
1456 processing cores, 416 TB storage capacity
(alice production, atlas, lhcb, gridifin, ifops)

**RO-11-NIPNE**:
120 processing cores, 40 TB storage capacity
(lhcb)

## LCG EVOLUTION

### LHC

- experiments increase the trigger rates
- increase of LHC luminosity

This is expected to generate larger amounts of experimental data -> will require:

- a major upgrade of the storage and processing resources
- a significantly grow the throughput both across LCG and within the resource centres

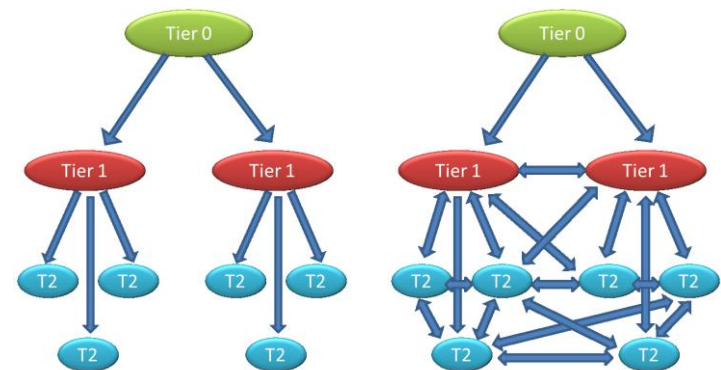### IT ADVANCES (networking, hardware & software)

### IMPACT ON LCG

The initial computing models are adapted to the current processing requirements of the experiments

Examples:

ATLAS - evolution from the hierarchical model

to the mesh model, in order to improve network

performance (ex.: LHCONE for T2s)

    - LAN throughput upgrade for analysis jobs

LHCb - starts running reprocessing jobs on Tier-2
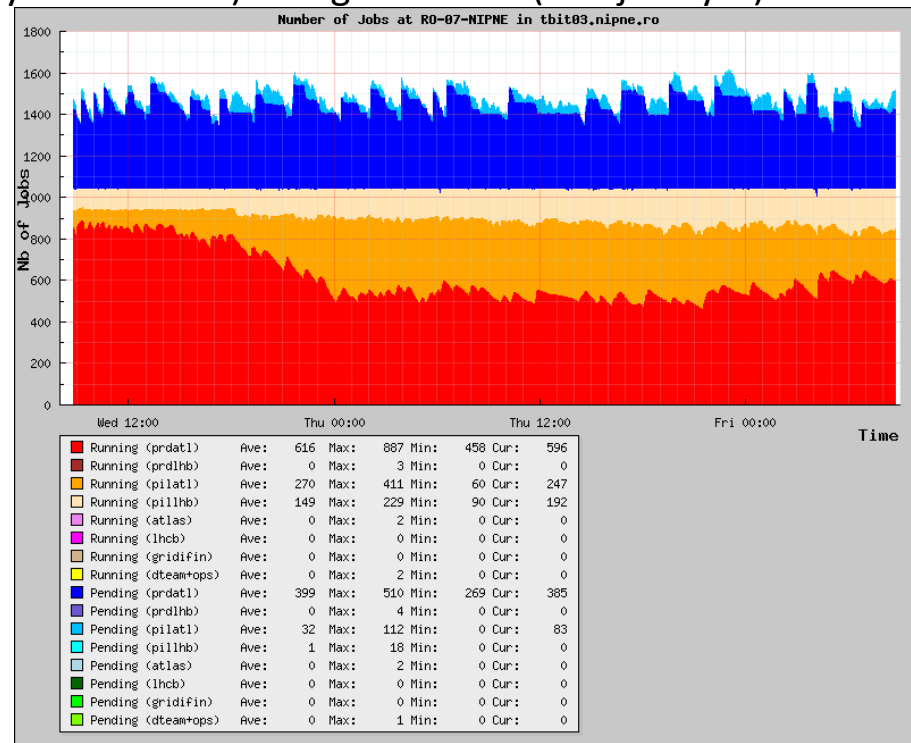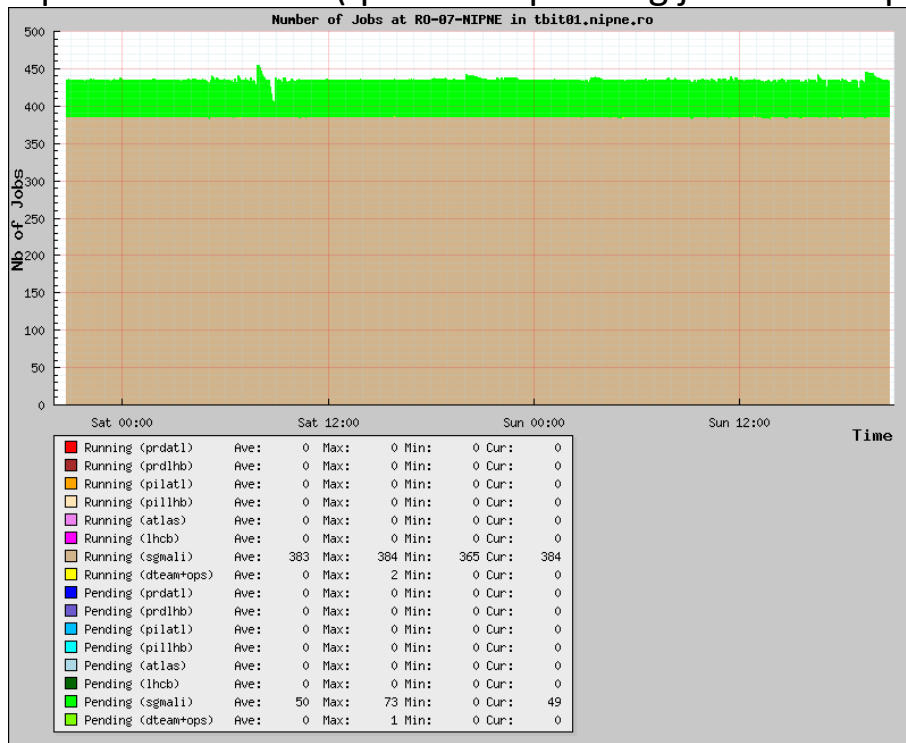
centres

The efficiency of resource usage in multiple-VO sites that run in intensive production depends on: job scheduling strategy, resource allocation policy, data throughput capacity, CPU efficiency per job type, etc. These factors were adapted to the specificities of the LCG job workflow, RO-07 being optimized for the simultaneous support of the alice, atlas, lhcb, and gridifin VOs.

Balance should be reached between each VO's requests (minimum queuing time and successful run of as much as possible jobs) and centre's requirements (optimal use of the infrastructure, fulfillment of the pledged SLA and resource levels for each experiment).

For alice, atlas, and lhcb (*after starting reprocessing on T2s*): the incoming job streams (run requests) are quasi-continuous (queues of pending jobs not empty - see below). For gridifin not (few jobs yet)

**=> Resources can be booked to 'busy' VOs**

The resource allocation and queuing policy was set, and the job scheduling rules were configured according to the analysis above.

One CE with its own job queue and associated WNs is dedicated to the alice production jobs, upon experiment's request. The remaining WNs are allotted to a second CE with queues dedicated to atlas, lhcb and gridifin VOs, with basic first-come-first-served (FCFS) scheduling strategy.

Queuing time limits can be fixed through MAUI to convenient values for CPU and wall times, respectively. Usually, these are longer than necessary for LCG jobs, which come with their own timers. Maximum limits for the number of running and queued LCG jobs were not imposed, as these are controlled locally by the available resources and centrally by the condition of not having more queued jobs than half of the number of running jobs on each site.

The number of queued gridifin jobs is limited by agreement with users.

To favor the running of atlas analysis jobs, which are shorter but more abundant, being most important for the ATLAS physics, an additional rule was imposed over the FCFS algorithm, that gives higher priority to all such jobs at regular intervals of time (e.g. 5 mins.).

# IMPROVING NETWORK PERFORMANCE

Optimal **external communication** (in terms of bandwidth, throughput and availability) is necessary for operations like: a) receiving JDL input data and asynchronous transfers of input files for data analysis (4-5 GB); b) sending job logs plus results files from SE to the users (4-10 GB), and sending the results of the simulations from SE to the Tier-1 centres for archiving and processing; c) grid monitoring activities like the SAM tests.

Poor bandwidth and/or throughput can hinder the simultaneous transfer of multiple result files, slowing down considerably the transfer rate of individual files. This can eventually lead to the abort of the transfers because the service quality conditions are not met, or to the blocking of job submission towards the site until the number of simultaneous transfers drops below a given limit (which is atlas policy). The situation described in this last example can lead to idle WNs, if the number of running jobs on site becomes smaller than the number of available CPUs.

The capacity of the **internal (local) network** is important as well, especially for the sites that perform data analysis, which requires the transfer of large input files from SE to WNs. The recent increase of the number of *atlas* analysis jobs that need to be concurrently processed led to a higher bandwidth consumption which can create bottlenecks in the local network, and the abortion of the jobs that reach the time limit.

To prevent this situation, measures were taken to upgrade the bandwidth in the RO-07-NIPNE cluster through switch cascading,  building the stack configuration depicted in Fig. 2. This allows to preserve the scalability of the cluster, by increasing the bandwidth available for data transfer at a constant rate whenever the storage capacity is upgraded.

## EXTERNAL NETWORK PERFORMANCE

In mesh topology, a new category of grid sites is defined: Direct T2s (**T2Ds**), that
- are primary hosts for datasets (analysis) and for group analysis
- get and send data from different clouds
- participate in cross cloud production

To became a T2D, sites have to improve the network transfers with T1s.

A site is T2D qualified if all the transfers of large files betwen the candidate site and 9/12 T1s are higher than a threshold value for some periods of time.

Standard measurements of transfer capacity are performed with perfSonar, which is a tool analysing the point-to-point throughput of a link.

To test the transfer capacity between the RO-LCG sites and T1s, and to detect points of failure, several perfSonar servers have been deployed :
    2 in DFCTI's datacenter (1 for bandwidth 1 for latency)
    2 servers in Bucharest RoEduNet Operations Centre
    2 servers in the RoEduNet NOC.

In the next 2 slides:

    measurements of throughput for transfers between IFIN-HH / RoEduNet: and T1 sites.

## PERFSONAR TESTS

# NETWORK SUPPORT

First figure: the total data transfer performance of the IFIN-HH – GEANT link during the concurrent transfer of multiple files from NIHAM's SE to various alice centres, reaching a global transfer rate of 8.5 Gbps on the external 10 Gbps interface.

Second figure: a total data traffic of 7.4 Gbps, which was obtained between the SE of RO-07-NIPNE and its WNs on the occasion of the concurrent transfer of 150 initial data files for the atlas analysis

# MONITORING SERVICE AVAILABILITY

ifops + nagios

->

Web interface developed for keeping track of the history of the results of ifops SAM tests



<-
Monitoring desk, automatically displaying the status of the grid services of 6 RO-LCG sites every 5 mins. (last 24-hours' window)

# CONCLUSIONS

❑ During the last six years since its founding, RO-LCG has evolved towards providing increasing computational resources and better quality services to the LHC community. Besides the fulfillment of the common Grid requirements, the necessity of concurrently supporting three LCG virtual organizations in intensive production regime required the design and implementation of specific technical solutions adapted to the local hardware and software environment. This was especially true in the case of the RO-07-NIPNE, which currently supports four VOs. Specific measures were taken regarding the resource allocation, the job scheduling, the data communication and management. A system of global monitoring of the Grid services availability was implemented for all the RO-LCG sites. A significant contribution to the performance management of the Grid system came from the NGI-independent set of tools implemented for monitoring the data transfer, storage efficacy, and service availability in the resource centres; this avoids possible inconvenients related to the end of the EGI support.

❑ Further work should focus on improving site and network efficiency. RO-LCG centres will join LHCONE, but measures are still to be enforced for improving the network connectivity, in order to become a T2D site.

## ACKNOWLEDGEMENTS

# THANK YOU
# FOR YOUR ATTENTION !